

FORMAL RULES

1. Algebraic structures

1.1. Typical structures

A vector space has been introduced as a 4-tuple $\mathcal{V} = (V, S, +, \cdot)$ with specific behavior of the vector addition and scaling operations. Arithmetic operations between scalars were implicitly assumed to be similar to those of the real numbers, but also must be specified to obtain a complete definition of a vector space. Algebra defines various structures that specify the behavior operations with objects. Knowledge of these structures is useful not only in linear algebra, but also in other mathematical approaches to data analysis such as topology or geometry.

Groups. A group is a 2-tuple $\mathcal{G} = (G, +)$ containing a set G and an operation $+$ with properties from Table 2. If $\forall a, b \in G, a + b = b + a$, the group is said to be commutative. Besides the familiar example of integers under addition $(\mathbb{Z}, +)$, symmetry groups that specify spatial or functional relations are of particular interest. The rotations by $0, \frac{\pi}{2}, \pi, \frac{3\pi}{2}$ or vertices of a square form a group.

Addition rules	
$a + b \in G$	Closure
$a + (b + c) = (a + b) + c$	Associativity
$0 + a = a$	Identity element
$a + (-a) = 0$	Inverse element

Table 1. Group $\mathcal{G} = (G, +)$ properties, for $\forall a, b, c \in G$

Rings. A ring is a 3-tuple $\mathcal{R} = (R, +, \cdot)$ containing a set R and two operations $+, \cdot$ with properties from Table 1. As is often the case, a ring is more complex structure built up from simpler algebraic structures. With respect to addition a ring has the properties of a commutative group. Only associativity and existence of an identity element is imposed for multiplication. Matrix addition and multiplication has the structure of ring $(\mathbb{R}^{m \times m}, +, \cdot)$.

Addition rules	
$(R, +)$ is a commutative (Abelian) group	
Multiplication rules	
$a \cdot b \in R$	Closure
$(a \cdot b) \cdot c = a \cdot (b \cdot c)$	Associativity
$a \cdot 1 = 1 \cdot a = a$	Identity element
Distributivity	
$a \cdot (b + c) = (a \cdot b) + (a \cdot c)$	on the left
$(a + b) \cdot c = (a \cdot c) + (b \cdot c)$	on the right

Table 2. Ring $\mathcal{R} = (R, +, \cdot)$ properties, for $\forall a, b, c \in R$.

Fields. A ring is a 3-tuple $\mathcal{F} = (F, +, \cdot)$ containing a set F and two operations $+, \cdot$, each with properties of a commutative group, but with special behavior for the inverse of the null element. The multiplicative inverse is denoted as a^{-1} . Scalars S in the definition of a vector space must satisfy the properties of a field. Since the operations are often understood from context a field might be referred to as the full 3-tuple, or, more concisely just through the set of elements as in the definition of a vector space.

Addition rules	
$(F, +)$ is a commutative (Abelian) group	
Multiplication rules	
(F, \cdot) is a commutative group except that 0^{-1} does not exist	
Distributivity	
$a \cdot (b + c) = (a \cdot b) + (a \cdot c)$	

Table 3. Field $\mathcal{F} = (F, +, \cdot)$ properties, for $\forall a, b, c \in F$.

Using the above definitions, a vector space $\mathcal{V} = (V, S, +, \cdot)$ can be described as a commutative group $(V, +)$ combined with a field S that satisfies the scaling properties $au \in V, a(u + v) = au + bv, (a + b)u = au + bu, a(bu) = (ab)u, 1u = u$, for $\forall a, b \in S, \forall u, v \in V$.

1.2. Vector subspaces

A central interest in data science is to seek simple description of complex objects. A typical situation is that many instances of some object of interest are initially given as an m -tuple $v \in \mathbb{R}^m$ with large m . Assuming that addition and scaling of such objects can cogently be defined, a vector space is obtained, say over the field of reals with an Euclidean

distance, E_m . Examples include for instance recordings of medical data (electroencephalograms, electrocardiograms), sound recordings, or images, for which m can easily reach in to the millions. A natural question to ask is whether all the m real numbers are actually needed to describe the observed objects, or perhaps there is some intrinsic description that requires a much smaller number of descriptive parameters, that still preserves the useful idea of linear combination. The mathematical transcription of this idea is a vector subspace.

DEFINITION. (VECTOR SUBSPACE). $\mathcal{U} = (U, S, +, \cdot)$, $U \neq \emptyset$, is a *vector subspace* of vector space $\mathcal{V} = (V, S, +, \cdot)$ over the same field of scalars S , denoted by $\mathcal{U} \leq \mathcal{V}$, if $U \subseteq V$ and $\forall a, b \in S$, $\forall \mathbf{u}, \mathbf{v} \in U$, the linear combination $a\mathbf{u} + b\mathbf{v} \in U$.

The above states a vector subspace must be closed under linear combination, and have the same vector addition and scaling operations as the enclosing vector space. The simplest vector subspace of a vector space is the null subspace that only contains the null element, $U = \{\mathbf{0}\}$. In fact any subspace must contain the null element $\mathbf{0}$, or otherwise closure would not be verified for the particular linear combination $\mathbf{u} + (-\mathbf{u}) = \mathbf{0}$. If $U \subset V$, then \mathcal{U} is said to be a *proper subspace* of \mathcal{V} , denoted by $\mathcal{U} < \mathcal{V}$.

Setting $n - m$ components equal to zero in the real space \mathcal{R}_m defines a proper subspace whose elements can be placed into a one-to-one correspondence with the vectors within \mathcal{R}_m . For example, setting component m of $\mathbf{x} \in \mathbb{R}^m$ equal to zero gives $\mathbf{x} = [x_1 \ x_2 \ \dots \ x_{m-1} \ 0]^T$ that while not a member of \mathbb{R}^{m-1} , is in a one-to-one relation with $\mathbf{x}' = [x_1 \ x_2 \ \dots \ x_{m-1}]^T \in \mathbb{R}^{m-1}$. Dropping the last component of $\mathbf{y} \in \mathbb{R}^m$, $\mathbf{y} = [y_1 \ y_2 \ \dots \ y_{m-1} \ y_m]^T$ gives vector $\mathbf{y}' = [y_1 \ y_2 \ \dots \ y_{m-1}] \in \mathbb{R}^{m-1}$, but this is no longer a one-to-one correspondence since for some given \mathbf{y}' , the last component y_m could take any value.

```
octave] m=3; x=[1; 2; 0]; xp=x(1:2); disp(xp)
1
2

octave] y=[1; 2; 3]; yp=y(1:2); disp(yp)
1
2

octave]
```

Vector subspaces arise in decomposition of a vector space. The converse, composition of vector spaces $\mathcal{U} = (U, S, +, \cdot)$ $\mathcal{V} = (V, S, +, \cdot)$ is also defined in terms of linear combination. A vector $\mathbf{x} \in \mathbb{R}^3$ can be obtained as the linear combination

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} x_1 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ x_2 \\ x_3 \end{bmatrix},$$

but also as

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 - a \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ a \\ x_3 \end{bmatrix},$$

for some arbitrary $a \in \mathbb{R}$. In the first case, \mathbf{x} is obtained as a unique linear combination of a vector from the set $U = \{[x_1 \ 0 \ 0]^T | x_1 \in \mathbb{R}\}$ with a vector from $V = \{[0 \ x_2 \ x_3]^T | x_2, x_3 \in \mathbb{R}\}$. In the second case, there is an infinity of linear combinations of a vector from V with another from $W = \{[x_1 \ x_2 \ 0]^T | x_1, x_2 \in \mathbb{R}\}$ to the vector \mathbf{x} . This is captured by a pair of definitions to describe vector space composition.

DEFINITION. Given two vector subspaces $\mathcal{U} = (U, S, +, \cdot)$, $\mathcal{V} = (V, S, +, \cdot)$ of the space $\mathcal{W} = (W, S, +, \cdot)$, the *sum* is the vector space $\mathcal{U} + \mathcal{V} = (U + V, S, +, \cdot)$, where the sum of the two sets of vectors U, V is $U + V = \{\mathbf{u} + \mathbf{v} | \mathbf{u} \in U, \mathbf{v} \in V\}$.

DEFINITION. Given two vector subspaces $\mathcal{U} = (U, S, +, \cdot)$, $\mathcal{V} = (V, S, +, \cdot)$ of the space $\mathcal{W} = (W, S, +, \cdot)$, the *direct sum* is the vector space $\mathcal{U} \oplus \mathcal{V} = (U \oplus V, S, +, \cdot)$, where the direct sum of the two sets of vectors U, V is $U \oplus V = \{\mathbf{u} + \mathbf{v} \mid \exists! \mathbf{u} \in U, \exists! \mathbf{v} \in V\}$. (unique decomposition)

Since the same scalar field, vector addition, and scaling is used, it is more convenient to refer to vector space sums simply by the sum of the vector sets $U + V$, or $U \oplus V$, instead of specifying the full tuple for each space. This shall be adopted henceforth to simplify the notation.

```
octave] u=[1; 0; 0]; v=[0; 2; 3]; vp=[0; 1; 3]; w=[1; 1; 0]; disp([u+v
vp+w])
```

```
1 1
2 2
3 3
```

```
octave]
```

In the previous example, the essential difference between the two ways to express $\mathbf{x} \in \mathbb{R}^3$ is that $U \cap V = \{\mathbf{0}\}$, but $V \cap W = \{[0 \ a \ 0]^T \mid a \in \mathbb{R}\} \neq \{\mathbf{0}\}$, and in general if the zero vector is the only common element of two vector spaces then the sum of the vector spaces becomes a direct sum. In practice, the most important procedure to construct direct sums or check when an intersection of two vector subspaces reduces to the zero vector is through an inner product.

DEFINITION. Two vector subspaces U, V of the real vector space \mathbb{R}^m are *orthogonal*, denoted as $U \perp V$ if $\mathbf{u}^T \mathbf{v} = 0$ for any $\mathbf{u} \in U, \mathbf{v} \in V$.

DEFINITION. Two vector subspaces U, V of $U + V$ are *orthogonal complements*, denoted $U = V^\perp$, $V = U^\perp$ if they are orthogonal subspaces, $U \perp V$, and $U \cap V = \{\mathbf{0}\}$, i.e., the null vector is the only common element of both subspaces.

```
octave] disp([u' * v vp' * w])
```

```
0 1
```

```
octave]
```

The above concept of orthogonality can be extended to other vector subspaces, such as spaces of functions. It can also be extended to other choices of an inner product, in which case the term conjugate vector spaces is sometimes used.

The concepts of sum and direct sum of vector spaces used linear combinations of the form $\mathbf{u} + \mathbf{v}$. This notion can be extended to arbitrary linear combinations.

DEFINITION. In vector space $\mathcal{V} = (V, S, +, \cdot)$, the *span* of vectors $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n \in V$, is the set of vectors reachable by linear combination

$$\text{span}\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n\} = \{\mathbf{b} \in V \mid \exists x_1, \dots, x_n \in S \text{ such that } \mathbf{b} = x_1 \mathbf{a}_1 + \dots + x_n \mathbf{a}_n\}.$$

Note that for real vector spaces a member of the span of the vectors $\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n\}$ is the vector \mathbf{b} obtained from the matrix vector multiplication

$$\mathbf{b} = \mathbf{A}\mathbf{x} = [\mathbf{a}_1 \ \mathbf{a}_2 \ \dots \ \mathbf{a}_n] \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}.$$

From the above, the span is a subset of the co-domain of the linear mapping $f(x) = Ax$.

2. Vector subspaces of a linear mapping

The wide-ranging utility of linear algebra essentially results a complete characterization of the behavior of a linear mapping between vector spaces $f: U \rightarrow V, f(au + bv) = af(u) + bf(v)$. For some given linear mapping the questions that arise are:

1. Can any vector within V be obtained by evaluation of f ?
2. Is there a single way that a vector within V can be obtained by evaluation of f ?

Linear mappings between real vector spaces $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$, have been seen to be completely specified by a matrix $A \in \mathbb{R}^{m \times n}$. It is common to frame the above questions about the behavior of the linear mapping $f(x) = Ax$ through sets associated with the matrix A . To frame an answer to the first question, a set of reachable vectors is first defined.

DEFINITION. The *column space* (or *range*) of matrix $A \in \mathbb{R}^{m \times n}$ is the set of vectors reachable by linear combination of the matrix column vectors

$$C(A) = \text{range}(A) = \{b \in \mathbb{R}^m \mid \exists x \in \mathbb{R}^n \text{ such that } b = Ax\}.$$

By definition, the column space is included in the co-domain of the function $f(x) = Ax$, $C(A) \subseteq \mathbb{R}^m$, and is readily seen to be a vector subspace of \mathbb{R}^m . The question that arises is whether the column space is the entire co-domain $C(A) = \mathbb{R}^m$ that would signify that any vector can be reached by linear combination. If this is not the case then the column space would be a proper subset, $C(A) \subset \mathbb{R}^m$, and the question is to determine what part of the co-domain cannot be reached by linear combination of columns of A . Consider the orthogonal complement of $C(A)$ defined as the set vectors orthogonal to all of the column vectors of A , expressed through inner products as

$$a_1^T y = 0, a_2^T y = 0, \dots, a_n^T y = 0.$$

This can be expressed more concisely through the transpose operation

$$A = [a_1 \ a_2 \ \dots \ a_n], A^T y = \begin{bmatrix} a_1^T \\ a_2^T \\ \vdots \\ a_n^T \end{bmatrix} = \begin{bmatrix} a_1^T y \\ a_2^T y \\ \vdots \\ a_n^T y \end{bmatrix},$$

and leads to the definition of a set of vectors for which $A^T y = \mathbf{0}$

DEFINITION. The *left null space* (or *cokernel*) of a matrix $A \in \mathbb{R}^{m \times n}$ is the set

$$N(A^T) = \text{null}(A^T) = \{y \in \mathbb{R}^m \mid A^T y = \mathbf{0}\}.$$

Note that the left null space is also a vector subspace of the co-domain of $f(x) = Ax$, $N(A^T) \subseteq \mathbb{R}^m$. The above definitions suggest that both the matrix and its transpose play a role in characterizing the behavior of the linear mapping $f = Ax$, so analogous sets are define for the transpose A^T .

DEFINITION. The *row space* (or *corange*) of a matrix $A \in \mathbb{R}^{m \times n}$ is the set

$$R(A) = C(A^T) = \text{range}(A^T) = \{c \in \mathbb{R}^n \mid \exists y \in \mathbb{R}^m c = A^T y\} \subseteq \mathbb{R}^n$$

DEFINITION. The *null space* of a matrix $A \in \mathbb{R}^{m \times n}$ is the set

$$N(A) = \text{null}(A) = \{x \in \mathbb{R}^n | Ax = \mathbf{0}\} \subseteq \mathbb{R}^n$$

Examples. Consider a linear mapping between real spaces $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$, defined by $y = f(x) = Ax = [y_1 \dots y_n]^T$, with $A \in \mathbb{R}^{m \times n}$.

1. For $n = 1, m = 3$,

$$A = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, A^T = [1 \ 0 \ 0],$$

the column space $C(A)$ is the y_1 -axis, and the left null space $N(A^T)$ is the y_2y_3 -plane. Vectors that span these spaces are returned by the Octave `orth` and `null` functions.

```
octave] A=[1; 0; 0];
         disp(orth(A));
         disp('-----');
         disp(null(A'))

-1
-0
-0
-----
 0  0
 1  0
 0  1

octave]
```

2. For $n = 2, m = 3$,

$$A = \begin{bmatrix} 1 & -1 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} = [a_1 \ a_2], A^T = \begin{bmatrix} 1 & 0 & 0 \\ -1 & 0 & 0 \end{bmatrix},$$

the columns of A are colinear, $a_2 = -a_1$, and the column space $C(A)$ is the y_1 -axis, and the left null space $N(A^T)$ is the y_2y_3 -plane, as before.

```
octave] A=[1 -1; 0 0; 0 0];
         disp(orth(A));
         disp('-----');
         disp(null(A'))

-1.00000
-0.00000
-0.00000
-----
 0  0
 1  0
 0  1

octave]
```

```
octave]
```

3. For $n = 2, m = 3$,

$$A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}, A^T = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix},$$

the column space $C(A)$ is the y_1y_2 -plane, and the left null space $N(A^T)$ is the y_3 -axis.

```
octave] A=[1 0; 0 1; 0 0];
         disp(orth(A));
         disp('-----');
         disp(null(A'))

-1  -0
-0  -1
-0  -0
-----
 0
 0
 1

octave]
```

4. For $n = 2, m = 3$,

$$A = \begin{bmatrix} 1 & 1 \\ 1 & -1 \\ 0 & 0 \end{bmatrix}, A^T = \begin{bmatrix} 1 & 1 & 0 \\ 1 & -1 & 0 \end{bmatrix},$$

the same $C(A), N(A^T)$ are obtained, albeit with a different set of spanning vectors returned by `orth`.

```
octave] A=[1 1; 1 -1; 0 0];
         disp(orth(A));
         disp('-----');
         disp(null(A'))

 0.70711  0.70711
 0.70711 -0.70711
-0.00000 -0.00000
-----
 0
 0
 1

octave]
```

```
octave]
```

5. For $n=3, m=3$,

$$A = \begin{bmatrix} 1 & 1 & 3 \\ 1 & -1 & -1 \\ 1 & 1 & 3 \end{bmatrix} = [a_1 \ a_2 \ a_3],$$

$$A^T = \begin{bmatrix} 1 & 1 & 1 \\ 1 & -1 & 1 \\ 3 & -1 & 3 \end{bmatrix} = \begin{bmatrix} a_1^T \\ a_2^T \\ a_3^T \end{bmatrix}, A^T y = \begin{bmatrix} a_1^T y \\ a_2^T y \\ a_3^T y \end{bmatrix}$$

since $a_3 = a_1 + 2a_2$, the orthogonality condition $A^T y = \mathbf{0}$ is satisfied by vectors of form $y = [a \ 0 \ -a]$, $a \in \mathbb{R}$.

```
octave] A=[1 1 3; 1 -1 -1; 1 1
          3]; disp(orth(A));
          disp('-----');
          disp(null(A'))
```

```
0.69157    0.14741
-0.20847   0.97803
0.69157    0.14741
-----
0.70711
0.00000
-0.70711
```

```
octave]
```

The above low dimensional examples are useful to gain initial insight into the significance of the spaces $C(A), N(A^T)$. Further appreciation can be gained by applying the same concepts to processing of images. A gray-scale image of size p_x by p_y pixels can be represented as a vector with $m = p_x p_y$ components, $b \in [0, 1]^m \subset \mathbb{R}^m$. Even for a small image with $p_x = p_y = 128 = 2^7$ pixels along each direction, the vector b would have $m = 2^{14}$ components. An image can be specified as a linear combination of the columns of the identity matrix

$$b = Ib = [e_1 \ e_2 \ \dots \ e_m] \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix},$$

with b_i the gray-level intensity in pixel i . Similar to the inclined plane example from §1, an alternative description as a linear combination of another set of vectors a_1, \dots, a_m might be more relevant. One choice of greater utility for image processing mimics the behavior of the set $\{1, \cos t, \cos 2t, \dots, \sin t, \sin 2t, \dots\}$ that extends the second example in §1, would be for $m=4$

$$A = [a_1 \ a_2 \ a_3 \ a_4] = \begin{bmatrix} 1 & 1 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 \end{bmatrix}.$$

DATA REDUNDANCY

1. Linear dependence

For the simple scalar mapping $f: \mathbb{R} \rightarrow \mathbb{R}, f(x) = ax$, the condition $f(x) = 0$ implies either that $a = 0$ or $x = 0$. Note that $a = 0$ can be understood as defining a zero mapping $f(x) = 0$. Linear mappings between vector spaces, $f: U \rightarrow V$, can exhibit different behavior, and the condition $f(x) = Ax = \mathbf{0}$, might be satisfied for both $x \neq \mathbf{0}$, and $A \neq \mathbf{0}$. Analogous to the scalar case, $A = \mathbf{0}$ can be understood as defining a zero mapping, $f(x) = \mathbf{0}$.

In vector space $\mathcal{V} = (V, S, +, \cdot)$, vectors $u, v \in V$ related by a scaling operation, $v = au, a \in S$, are said to be colinear, and are considered to contain redundant data. This can be restated as $v \in \text{span}\{u\}$, from which it results that $\text{span}\{u\} = \text{span}\{u, v\}$. Colinearity can be expressed only in terms of vector scaling, but other types of redundancy arise when also considering vector addition as expressed by the span of a vector set. Assuming that $v \notin \text{span}\{u\}$, then the strict inclusion relation $\text{span}\{u\} \subset \text{span}\{u, v\}$ holds. This strict inclusion expressed in terms of set concepts can be transcribed into an algebraic condition.

DEFINITION. The vectors $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n \in V$, are *linearly dependent* if there exist n scalars, $x_1, \dots, x_n \in S$, at least one of which is different from zero such that

$$x_1 \mathbf{a}_1 + \dots + x_n \mathbf{a}_n = \mathbf{0}.$$

Introducing a matrix representation of the vectors

$$\mathbf{A} = [\mathbf{a}_1 \ \mathbf{a}_2 \ \dots \ \mathbf{a}_n]; \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

allows restating linear dependence as the existence of a non-zero vector, $\exists \mathbf{x} \neq \mathbf{0}$, such that $\mathbf{A}\mathbf{x} = \mathbf{0}$. Linear dependence can also be written as $\mathbf{A}\mathbf{x} = \mathbf{0} \not\Rightarrow \mathbf{x} = \mathbf{0}$, or that one cannot deduce from the fact that the linear mapping $f(\mathbf{x}) = \mathbf{A}\mathbf{x}$ attains a zero value that the argument itself is zero. The converse of this statement would be that the only way to ensure $\mathbf{A}\mathbf{x} = \mathbf{0}$ is for $\mathbf{x} = \mathbf{0}$, or $\mathbf{A}\mathbf{x} = \mathbf{0} \Rightarrow \mathbf{x} = \mathbf{0}$, leading to the concept of linear independence.

DEFINITION. The vectors $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n \in V$, are *linearly independent* if the only n scalars, $x_1, \dots, x_n \in S$, that satisfy

$$x_1 \mathbf{a}_1 + \dots + x_n \mathbf{a}_n = \mathbf{0}, \tag{1}$$

are $x_1 = 0, x_2 = 0, \dots, x_n = 0$.

2. Basis and dimension

Vector spaces are closed under linear combination, and the span of a vector set $\mathcal{B} = \{\mathbf{a}_1, \mathbf{a}_2, \dots\}$ defines a vector subspace. If the entire set of vectors can be obtained by a spanning set, $V = \text{span } \mathcal{B}$, extending \mathcal{B} by an additional element $\mathcal{C} = \mathcal{B} \cup \{\mathbf{b}\}$ would be redundant since $\text{span } \mathcal{B} = \text{span } \mathcal{C}$. This is recognized by the concept of a basis, and also allows leads to a characterization of the size of a vector space by the cardinality of a basis set.

DEFINITION. A set of vectors $\mathbf{u}_1, \dots, \mathbf{u}_n \in V$ is a *basis* for vector space $\mathcal{V} = (V, S, +, \cdot)$ if

1. $\mathbf{u}_1, \dots, \mathbf{u}_n$ are linearly independent;
2. $\text{span}\{\mathbf{u}_1, \dots, \mathbf{u}_n\} = V$.

DEFINITION. The number of vectors $\mathbf{u}_1, \dots, \mathbf{u}_n \in V$ within a basis is the *dimension* of the vector space $\mathcal{V} = (V, S, +, \cdot)$.

3. Dimension of matrix spaces

The domain and co-domain of the linear mapping $f: U \rightarrow V, f(\mathbf{x}) = \mathbf{A}\mathbf{x}$, are decomposed by the spaces associated with the matrix \mathbf{A} . When $U = \mathbb{R}^n, V = \mathbb{R}^m$, the following vector subspaces associated with the matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ have been defined:

- $C(\mathbf{A})$ the column space of \mathbf{A}
- $C(\mathbf{A}^T)$ the row space of \mathbf{A}
- $N(\mathbf{A})$ the null space of \mathbf{A}
- $N(\mathbf{A}^T)$ the left null space of \mathbf{A} , or null space of \mathbf{A}^T

DEFINITION. The *rank* of a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ is the dimension of its column space and is equal to the dimension of its row space.

DEFINITION. The *nullity* of a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ is the dimension of its null space.

DATA INFORMATION

1. Partition of linear mapping domain and codomain

A partition of a set S has been introduced as a collection of subsets $P = \{S_i | S_i \subset P, S_i \neq \emptyset\}$ such that any given element $x \in S$ belongs to only one set in the partition. This is modified when applied to subspaces of a vector space, and a partition of a set of vectors is understood as a collection of subsets such that any vector except $\mathbf{0}$ belongs to only one member of the partition.

Linear mappings between vector spaces $f: U \rightarrow V$ can be represented by matrices A with columns that are images of the columns of a basis $\{\mathbf{u}_1, \mathbf{u}_2, \dots\}$ of U

$$A = [f(\mathbf{u}_1) \ f(\mathbf{u}_2) \ \dots].$$

Consider the case of real finite-dimensional domain and co-domain, $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$, in which case $A \in \mathbb{R}^{m \times n}$,

$$A = [f(\mathbf{e}_1) \ f(\mathbf{e}_2) \ \dots \ f(\mathbf{e}_n)] = [\mathbf{a}_1 \ \mathbf{a}_2 \ \dots \ \mathbf{a}_n].$$

The column space of A is a vector subspace of the codomain, $C(A) \leq \mathbb{R}^m$, but according to the definition of dimension if $n < m$ there remain non-zero vectors within the codomain that are outside the range of A ,

$$n < m \Rightarrow \exists \mathbf{v} \in \mathbb{R}^m, \mathbf{v} \neq \mathbf{0}, \mathbf{v} \notin C(A).$$

All of the non-zero vectors in $N(A^T)$, namely the set of vectors orthogonal to all columns in A fall into this category. The above considerations can be stated as

$$C(A) \leq \mathbb{R}^m, \ N(A^T) \leq \mathbb{R}^m, \ C(A) \perp N(A^T) \ C(A) + N(A^T) \leq \mathbb{R}^m.$$

The question that arises is whether there remain any non-zero vectors in the codomain that are not part of $C(A)$ or $N(A^T)$. The fundamental theorem of linear algebra states that there no such vectors, that $C(A)$ is the orthogonal complement of $N(A^T)$, and their direct sum covers the entire codomain $C(A) \oplus N(A^T) = \mathbb{R}^m$.

LEMMA 1. Let \mathcal{U}, \mathcal{V} , be subspaces of vector space \mathcal{W} . Then $\mathcal{W} = \mathcal{U} \oplus \mathcal{V}$ if and only if

- i. $\mathcal{W} = \mathcal{U} + \mathcal{V}$, and
- ii. $\mathcal{U} \cap \mathcal{V} = \{\mathbf{0}\}$.

Proof. $\mathcal{W} = \mathcal{U} \oplus \mathcal{V} \Rightarrow \mathcal{W} = \mathcal{U} + \mathcal{V}$ by definition of direct sum, sum of vector subspaces. To prove that $\mathcal{W} = \mathcal{U} \oplus \mathcal{V} \Rightarrow \mathcal{U} \cap \mathcal{V} = \{\mathbf{0}\}$, consider $\mathbf{w} \in \mathcal{U} \cap \mathcal{V}$. Since $\mathbf{w} \in \mathcal{U}$ and $\mathbf{w} \in \mathcal{V}$ write

$$\mathbf{w} = \mathbf{w} + \mathbf{0} \quad (\mathbf{w} \in \mathcal{U}, \mathbf{0} \in \mathcal{V}), \quad \mathbf{w} = \mathbf{0} + \mathbf{w} \quad (\mathbf{0} \in \mathcal{U}, \mathbf{w} \in \mathcal{V}),$$

and since expression $\mathbf{w} = \mathbf{u} + \mathbf{v}$ is unique, it results that $\mathbf{w} = \mathbf{0}$. Now assume (i),(ii) and establish an unique decomposition. Assume there might be two decompositions of $\mathbf{w} \in \mathcal{W}$, $\mathbf{w} = \mathbf{u}_1 + \mathbf{v}_1$, $\mathbf{w} = \mathbf{u}_2 + \mathbf{v}_2$, with $\mathbf{u}_1, \mathbf{u}_2 \in \mathcal{U}$, $\mathbf{v}_1, \mathbf{v}_2 \in \mathcal{V}$. Obtain $\mathbf{u}_1 + \mathbf{v}_1 = \mathbf{u}_2 + \mathbf{v}_2$, or $\mathbf{x} = \mathbf{u}_1 - \mathbf{u}_2 = \mathbf{v}_2 - \mathbf{v}_1$. Since $\mathbf{x} \in \mathcal{U}$ and $\mathbf{x} \in \mathcal{V}$ it results that $\mathbf{x} = \mathbf{0}$, and $\mathbf{u}_1 = \mathbf{u}_2$, $\mathbf{v}_1 = \mathbf{v}_2$, i.e., the decomposition is unique. \square

In the vector space $U + V$ the subspaces U, V are said to be orthogonal complements is $U \perp V$, and $U \cap V = \{\mathbf{0}\}$. When $U \leq \mathbb{R}^m$, the orthogonal complement of U is denoted as U^\perp , $U \oplus U^\perp = \mathbb{R}^m$.

THEOREM. Given the linear mapping associated with matrix $A \in \mathbb{R}^{m \times n}$ we have:

1. $C(A) \oplus N(A^T) = \mathbb{R}^m$, the direct sum of the column space and left null space is the codomain of the mapping
2. $C(A^T) \oplus N(A) = \mathbb{R}^n$, the direct sum of the row space and null space is the domain of the mapping

3. $C(\mathbf{A}) \perp N(\mathbf{A}^T)$ and $C(\mathbf{A}) \cap N(\mathbf{A}^T) = \{\mathbf{0}\}$, the column space is orthogonal to the left null space, and they are orthogonal complements of one another,

$$C(\mathbf{A}) = N(\mathbf{A}^T)^\perp, N(\mathbf{A}^T) = C(\mathbf{A})^\perp.$$

4. $C(\mathbf{A}^T) \perp N(\mathbf{A})$ and $C(\mathbf{A}^T) \cap N(\mathbf{A}) = \{\mathbf{0}\}$, the row space is orthogonal to the null space, and they are orthogonal complements of one another,

$$C(\mathbf{A}^T) = N(\mathbf{A})^\perp, N(\mathbf{A}) = C(\mathbf{A}^T)^\perp.$$

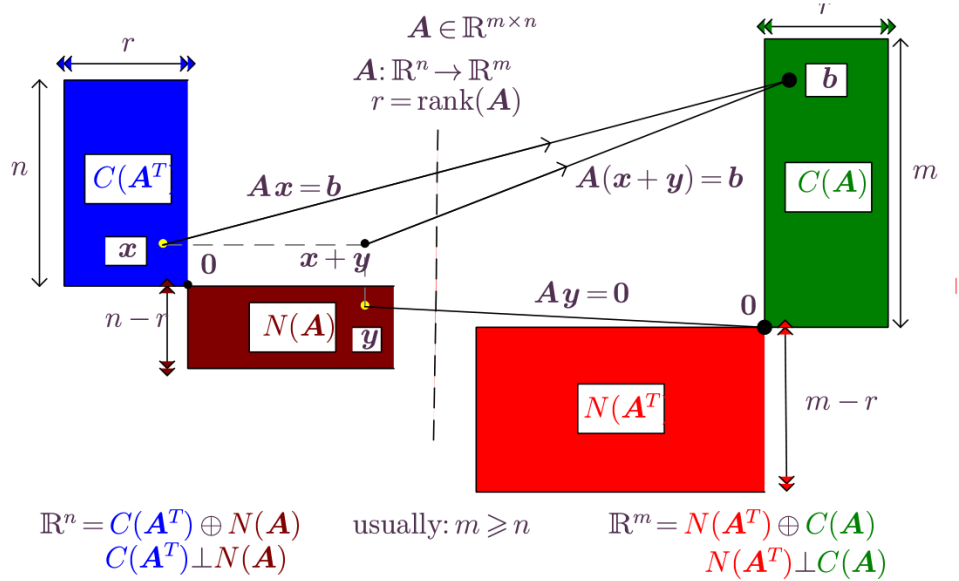


Figure 1. Graphical representation of the Fundamental Theorem of Linear Algebra, Gil Strang, *Amer. Math. Monthly* **100**, 848-855, 1993.

Consideration of equality between sets arises in proving the above theorem. A standard technique to show set equality $A = B$, is by double inclusion, $A \subseteq B \wedge B \subseteq A \Rightarrow A = B$. This is shown for the statements giving the decomposition of the codomain \mathbb{R}^m . A similar approach can be used to decomposition of \mathbb{R}^n .

- i. $C(\mathbf{A}) \perp N(\mathbf{A}^T)$ (column space is orthogonal to left null space).

Proof. Consider arbitrary $\mathbf{u} \in C(\mathbf{A}), \mathbf{v} \in N(\mathbf{A}^T)$. By definition of $C(\mathbf{A})$, $\exists \mathbf{x} \in \mathbb{R}^n$ such that $\mathbf{u} = \mathbf{A}\mathbf{x}$, and by definition of $N(\mathbf{A}^T)$, $\mathbf{A}^T \mathbf{v} = \mathbf{0}$. Compute $\mathbf{u}^T \mathbf{v} = (\mathbf{A}\mathbf{x})^T \mathbf{v} = \mathbf{x}^T \mathbf{A}^T \mathbf{v} = \mathbf{x}^T (\mathbf{A}^T \mathbf{v}) = \mathbf{x}^T \mathbf{0} = 0$, hence $\mathbf{u} \perp \mathbf{v}$ for arbitrary \mathbf{u} , \mathbf{v} , and $C(\mathbf{A}) \perp N(\mathbf{A}^T)$. □

- ii. $C(\mathbf{A}) \cap N(\mathbf{A}^T) = \{\mathbf{0}\}$ ($\mathbf{0}$ is the only vector both in $C(\mathbf{A})$ and $N(\mathbf{A}^T)$).

Proof. (By contradiction, *reductio ad absurdum*). Assume there might be $\mathbf{b} \in C(\mathbf{A})$ and $\mathbf{b} \in N(\mathbf{A}^T)$ and $\mathbf{b} \neq \mathbf{0}$. Since $\mathbf{b} \in C(\mathbf{A})$, $\exists \mathbf{x} \in \mathbb{R}^n$ such that $\mathbf{b} = \mathbf{A}\mathbf{x}$. Since $\mathbf{b} \in N(\mathbf{A}^T)$, $\mathbf{A}^T \mathbf{b} = \mathbf{A}^T (\mathbf{A}\mathbf{x}) = \mathbf{0}$. Note that $\mathbf{x} \neq \mathbf{0}$ since $\mathbf{x} = \mathbf{0} \Rightarrow \mathbf{b} = \mathbf{0}$, contradicting assumptions. Multiply equality $\mathbf{A}^T \mathbf{A}\mathbf{x} = \mathbf{0}$ on left by \mathbf{x}^T ,

$$\mathbf{x}^T \mathbf{A}^T \mathbf{A}\mathbf{x} = \mathbf{0} \Rightarrow (\mathbf{A}\mathbf{x})^T (\mathbf{A}\mathbf{x}) = \mathbf{b}^T \mathbf{b} = \|\mathbf{b}\|^2 = 0,$$

thereby obtaining $\mathbf{b} = \mathbf{0}$, using norm property 3. Contradiction. □

- iii. $C(\mathbf{A}) \oplus N(\mathbf{A}^T) = \mathbb{R}^m$

Proof. (iii) and (iv) have established that $C(\mathbf{A}), N(\mathbf{A}^T)$ are orthogonal complements

$$C(\mathbf{A}) = N(\mathbf{A}^T)^\perp, N(\mathbf{A}^T) = C(\mathbf{A})^\perp.$$

By Lemma 2 it results that $C(\mathbf{A}) \oplus N(\mathbf{A}^T) = \mathbb{R}^m$. □

The remainder of the FTLA is established by considering $\mathbf{B} = \mathbf{A}^T$, e.g., since it has been established in (v) that $C(\mathbf{B}) \oplus N(\mathbf{A}^T) = \mathbb{R}^n$, replacing $\mathbf{B} = \mathbf{A}^T$ yields $C(\mathbf{A}^T) \oplus N(\mathbf{A}) = \mathbb{R}^m$, etc.

DATA PARTITIONING

1. Mappings as data

1.1. Vector spaces of mappings and matrix representations

A vector space \mathcal{L} can be formed from all linear mappings from the vector space $\mathcal{U} = (U, S, +, \cdot)$ to another vector space $\mathcal{V} = (V, S, +, \cdot)$

$$\mathcal{L} = \{L, S, +, \cdot\}, L = \{f | f: U \rightarrow V, f(au + bv) = af(u) + bf(v)\},$$

with addition and scaling of linear mappings defined by $(f + g)(u) = f(u) + g(u)$ and $(af)(u) = af(u)$. Let $B = \{u_1, u_2, \dots\}$ denote a basis for the domain U of linear mappings within \mathcal{L} , such that the linear mapping $f \in \mathcal{L}$ is represented by the matrix

$$\mathbf{A} = [f(u_1) \ f(u_2) \ \dots].$$

When the domain and codomain are the real vector spaces $U = \mathbb{R}^n$, $V = \mathbb{R}^m$, the above is a standard matrix of real numbers, $\mathbf{A} \in \mathbb{R}^{m \times n}$. For linear mappings between infinite dimensional vector spaces the matrix is understood in a generalized sense to contain an infinite number of columns that are elements of the codomain V . For example, the indefinite integral is a linear mapping between the vector space of functions that allow differentiation to any order,

$$\int: \mathcal{C}^\infty \rightarrow \mathcal{C}^\infty \quad v(x) = \int u(x) dx$$

and for the monomial basis $B = \{1, x, x^2, \dots\}$, is represented by the generalized matrix

$$\mathbf{A} = \left[x \quad \frac{1}{2}x^2 \quad \frac{1}{3}x^3 \quad \dots \right].$$

Truncation of the basis expansion $u(x) = \sum_{j=1}^{\infty} u_j x^j$ where $u_j \in \mathbb{R}$ to n terms, and sampling of $u \in \mathcal{C}^\infty$ at points x_1, \dots, x_m , forms a standard matrix of real numbers

$$\mathbf{A} = \left[x \quad \frac{1}{2}x^2 \quad \frac{1}{3}x^3 \quad \dots \right] \in \mathbb{R}^{m \times n}, \quad \mathbf{x}^j = \begin{bmatrix} x_1^j \\ \vdots \\ x_m^j \end{bmatrix}.$$

As to be expected, matrices can also be organized as vector space \mathcal{M} , which is essentially the representation of the associated vector space of linear mappings,

$$\mathcal{M} = (M, S, +, \cdot) \quad M = \{A | A = [f(u_1) \ f(u_2) \ \dots]\}.$$

The addition $C = A + B$ and scaling $S = aR$ of matrices is given in terms of the matrix components by

$$c_{ij} = a_{ij} + b_{ij}, s_{ij} = ar_{ij}.$$

1.2. Measurement of mappings

From the above it is apparent that linear mappings and matrices can also be considered as data, and a first step in analysis of such data is definition of functionals that would attach a single scalar label to each linear mapping of matrix. Of particular interest is the definition of a norm functional that characterizes in an appropriate sense the size of a linear mapping.

Consider first the case of finite matrices with real components $A \in \mathbb{R}^{m \times n}$ that represent linear mappings between real vector spaces $f: \mathbb{R}^m \rightarrow \mathbb{R}^n$. The columns $\mathbf{a}_1, \dots, \mathbf{a}_n$ of $A \in \mathbb{R}^{m \times n}$ could be placed into a single column vector \mathbf{c} with mn components

$$\mathbf{c} = \begin{bmatrix} \mathbf{a}_1 \\ \vdots \\ \mathbf{a}_n \end{bmatrix}.$$

Subsequently the norm of the matrix A could be defined as the norm of the vector \mathbf{c} . An example of this approach is the Frobenius norm

$$\|A\|_F = \|\mathbf{c}\|_2 = \left(\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2 \right)^{1/2}.$$

A drawback of the above approach is that the structure of the matrix and its close relationship to a linear mapping is lost. A more useful characterization of the size of a mapping is to consider the amplification behavior of linear mapping. The motivation is readily understood starting from linear mappings between the reals $f: \mathbb{R} \rightarrow \mathbb{R}$, that are of the form $f(x) = ax$. When given an argument of unit magnitude $|x| = 1$, the mapping returns a real number with magnitude $|a|$. For mappings $f: \mathbb{R}^2 \rightarrow \mathbb{R}^2$ within the plane, arguments that satisfy $\|x\|_2 = 1$ are on the unit circle with components $x = [\cos \theta \ \sin \theta]$ have images through f given analytically by

$$f(x) = Ax = [\mathbf{a}_1 \ \mathbf{a}_2] \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix} = \cos \theta \mathbf{a}_1 + \sin \theta \mathbf{a}_2,$$

and correspond to ellipses.

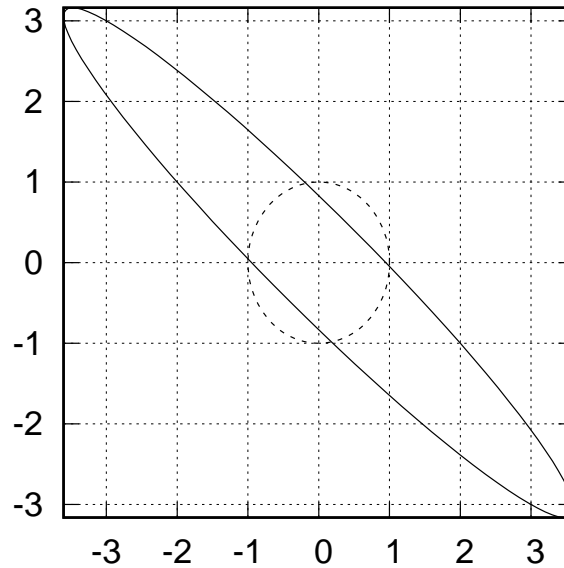


Figure 1. Mapping of unit circle by $f(x) = Ax$, $A = \begin{bmatrix} 2 & 3 \\ -1 & -3 \end{bmatrix}$.

From the above the mapping associated A amplifies some directions more than others. This suggests a definition of the size of a matrix or a mapping by the maximal amplification unit norm vectors within the domain.

DEFINITION. For vector spaces U, V with norms $\|\cdot\|_U: U \rightarrow \mathbb{R}_+$, $\|\cdot\|_V: V \rightarrow \mathbb{R}_+$, the *induced norm* of $f: U \rightarrow V$ is

$$\|f\| = \sup_{\|x\|_U=1} \|f(x)\|_V.$$

DEFINITION. For vector spaces $\mathbb{R}^n, \mathbb{R}^m$ with norms $\|\cdot\|^{(n)}: U \rightarrow \mathbb{R}_+$, $\|\cdot\|^{(m)}: V \rightarrow \mathbb{R}_+$, the *induced norm* of matrix $A \in \mathbb{R}^{m \times n}$ is

$$\|A\| = \sup_{\|x\|^{(n)}=1} \|Ax\|^{(m)}.$$

In the above, any vector norm can be used within the domain and codomain.

2. The Singular Value Decomposition (SVD)

The fundamental theorem of linear algebra partitions the domain and codomain of a linear mapping $f: U \rightarrow V$. For real vectors spaces $U = \mathbb{R}^n$, $V = \mathbb{R}^m$ the partition properties are stated in terms of spaces of the associated matrix A as

$$C(A) \oplus N(A^T) = \mathbb{R}^m \quad C(A) \perp N(A^T) \quad C(A^T) \oplus N(A) = \mathbb{R}^n \quad C(A^T) \perp N(A).$$

The dimension of the column and row spaces $r = \dim C(A) = \dim C(A^T)$ is the rank of the matrix, $n - r$ is the nullity of A , and $m - r$ is the nullity of A^T . A infinite number of bases could be defined for the domain and codomain. It is of great theoretical and practical interest bases with properties that facilitate insight or computation.

2.1. Orthogonal matrices

The above partitions of the domain and codomain are orthogonal, and suggest searching for orthogonal bases within these subspaces. Introduce a matrix representation for the bases

$$U = [u_1 \ u_2 \ \dots \ u_m] \in \mathbb{R}^{m \times m}, V = [v_1 \ v_2 \ \dots \ v_n] \in \mathbb{R}^{n \times n},$$

with $C(U) = \mathbb{R}^m$ and $C(V) = \mathbb{R}^n$. Orthogonality between columns u_i, u_j for $i \neq j$ is expressed as $u_i^T u_j = 0$. For $i = j$, the inner product is positive $u_i^T u_i > 0$, and since scaling of the columns of U preserves the spanning property $C(U) = \mathbb{R}^m$, it is convenient to impose $u_i^T u_i = 1$. Such behavior is concisely expressed as a matrix product

$$U^T U = I_m,$$

with I_m the identity matrix in \mathbb{R}^m . Expanded in terms of the column vectors of U the first equality is

$$[u_1 \ u_2 \ \dots \ u_m]^T [u_1 \ u_2 \ \dots \ u_m] = \begin{bmatrix} u_1^T \\ u_2^T \\ \vdots \\ u_m^T \end{bmatrix} [u_1 \ u_2 \ \dots \ u_m] = \begin{bmatrix} u_1^T u_1 & u_1^T u_2 & \dots & u_1^T u_m \\ u_2^T u_1 & u_2^T u_2 & \dots & u_2^T u_m \\ \vdots & \vdots & \ddots & \vdots \\ u_m^T u_1 & u_m^T u_2 & \dots & u_m^T u_m \end{bmatrix} = I_m.$$

It is useful to determine if a matrix X exists such that $UX = I_m$, or

$$UX = U [x_1 \ x_2 \ \dots \ x_m] = [e_1 \ e_2 \ \dots \ e_m].$$

The columns of X are the coordinates of the column vectors of I_m in the basis U , and can readily be determined

$$Ux_j = e_j \Rightarrow U^T Ux_j = U^T e_j \Rightarrow I_m x_j = \begin{bmatrix} u_1^T \\ u_2^T \\ \vdots \\ u_m^T \end{bmatrix} e_j \Rightarrow x_j = (U^T)_j,$$

where $(U^T)_j$ is the j^{th} column of U^T , hence $X = U^T$, leading to

$$U^T U = I = U U^T.$$

Note that the second equality

$$[u_1 \ u_2 \ \dots \ u_m] [u_1 \ u_2 \ \dots \ u_m]^T = [u_1 \ u_2 \ \dots \ u_m] \begin{bmatrix} u_1^T \\ u_2^T \\ \vdots \\ u_m^T \end{bmatrix} = u_1 u_1^T + u_2 u_2^T + \dots + u_m u_m^T = I$$

acts as normalization condition on the matrices $U_j = u_j u_j^T$.

DEFINITION. A square matrix U is said to be orthogonal if $U^T U = U U^T = I$.

2.2. Intrinsic basis of a linear mapping

Given a linear mapping $f: U \rightarrow V$, expressed as $y = f(x) = Ax$, the simplest description of the action of A would be a simple scaling, as exemplified by $g(x) = ax$ that has as its associated matrix aI . Recall that specification of a vector is typically done in terms of the identity matrix $b = Ib$, but may be more insightfully given in some other basis $Ax = Ib$. This suggests that especially useful bases for the domain and codomain would reduce the action of a linear mapping to scaling along orthogonal directions, and evaluate $y = Ax$ by first re-expressing y in another basis U , $Us = Iy$ and re-expressing x in another basis V , $Vr = Ix$. The condition that the linear operator reduces to simple scaling in these new bases is expressed as $s_i = \sigma_i r_i$ for $i = 1, \dots, \min(m, n)$, with σ_i the scaling coefficients along each direction which can be expressed as a matrix vector product $s = \Sigma r$, where $\Sigma \in \mathbb{R}^{m \times n}$ is of the same dimensions as A and given by

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

Imposing the condition that U, V are orthogonal leads to

$$Us = y \Rightarrow s = U^T y, Vr = x \Rightarrow r = V^T x,$$

which can be replaced into $s = \Sigma r$ to obtain

$$U^T y = \Sigma V^T x \Rightarrow y = U \Sigma V^T x.$$

From the above the orthogonal bases U, V and scaling coefficients Σ that are sought must satisfy $A = U \Sigma V^T$.

THEOREM. Every matrix $A \in \mathbb{R}^{m \times n}$ has a *singular value decomposition (SVD)*

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

with properties:

1. $U \in \mathbb{R}^{m \times m}$ is an orthogonal matrix, $U^T U = I_m$;
2. $V \in \mathbb{R}^{n \times n}$ is an orthogonal matrix, $V^T V = I_n$;
3. $\Sigma \in \mathbb{R}^{m \times n}$ is diagonal, $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_p)$, $p = \min(m, n)$, and $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$.

Proof. The proof of the SVD makes use of properties of the norm, concepts from analysis and complete induction. Adopting the 2-norm set $\sigma_1 = \|A\|_2$,

$$\sigma_1 = \sup_{\|x\|_2=1} \|Ax\|_2.$$

The domain $\|x\|_2 = 1$ is compact (closed and bounded), and the extreme value theorem implies that $f(x) = Ax$ attains its maxima and minima, hence there must exist some vectors u_1, v_1 of unit norm such that $\sigma_1 u_1 = Av_1 \Rightarrow \sigma_1 = u_1^T Av_1$. Introduce orthogonal bases U_1, V_1 for $\mathbb{R}^m, \mathbb{R}^n$ whose first column vectors are u_1, v_1 , and compute

$$U_1^T A V_1 = \begin{bmatrix} u_1^T \\ \vdots \\ u_m^T \end{bmatrix} [Av_1 \dots Av_n] = \begin{bmatrix} \sigma_1 & w^T \\ \mathbf{0} & B \end{bmatrix} = C.$$

In the above w^T is a row vector with $n-1$ components $u_1^T Av_j$, $j=2, \dots, n$, and $u_i^T Av_1$ must be zero for u_1 to be the direction along which the maximum norm $\|Av_1\|$ is obtained. Introduce vectors

$$y = \begin{bmatrix} \sigma_1 \\ w \end{bmatrix}, z = Cy = \begin{bmatrix} \sigma_1^2 + w^T w \\ Bw \end{bmatrix},$$

and note that $\|z\|_2 \geq \|y\|_2^2 = \sigma_1^2 + w^T w$. From $\|U_1^T A V_1\| = \|A\| = \sigma_1 = \|C\| \geq \sigma_1^2 + w^T w$ it results that $w = \mathbf{0}$. By induction, assume that B has a singular value decomposition, $B = U_2 \Sigma_2 V_2^T$, such that

$$U_1^T A V_1 = \begin{bmatrix} \sigma_1 & \mathbf{0}^T \\ \mathbf{0} & U_2 \Sigma_2 V_2^T \end{bmatrix} = \begin{bmatrix} 1 & \mathbf{0}^T \\ \mathbf{0} & U_2 \end{bmatrix} \begin{bmatrix} \sigma_1 & \mathbf{0}^T \\ \mathbf{0} & \Sigma_2 \end{bmatrix} \begin{bmatrix} 1 & \mathbf{0}^T \\ \mathbf{0} & V_2^T \end{bmatrix},$$

and the orthogonal matrices arising in the singular value decomposition of A are

$$U = U_1 \begin{bmatrix} 1 & \mathbf{0}^T \\ \mathbf{0} & U_2 \end{bmatrix}, V^T = \begin{bmatrix} 1 & \mathbf{0}^T \\ \mathbf{0} & V_2^T \end{bmatrix} V_1^T.$$

□

The scaling coefficients σ_j are called the *singular values* of A . The columns of U are called the *left singular vectors*, and those of V are called the *right singular vectors*.

The fact that the scaling coefficients are norms of A and submatrices of A , $\sigma_1 = \|A\|$, is crucial importance in applications. Carrying out computation of the matrix products

$$A = [\mathbf{u}_1 \ \mathbf{u}_2 \ \dots \ \mathbf{u}_r \ \mathbf{u}_{r+1} \ \dots \ \mathbf{u}_m] \begin{bmatrix} \sigma_1 & 0 & \dots & 0 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & 0 & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_r & 0 & \dots & 0 \\ 0 & 0 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1^T \\ \mathbf{v}_2^T \\ \vdots \\ \mathbf{v}_r^T \\ \vdots \\ \mathbf{v}_n^T \end{bmatrix} = [\mathbf{u}_1 \ \mathbf{u}_2 \ \dots \ \mathbf{u}_r \ \mathbf{u}_{r+1} \ \dots \ \mathbf{u}_m] \begin{bmatrix} \sigma_1 \mathbf{v}_1^T \\ \sigma_2 \mathbf{v}_2^T \\ \vdots \\ \sigma_r \mathbf{v}_r^T \\ \vdots \\ 0 \end{bmatrix}$$

leads to a representation of A as a sum

$$A = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T, r \leq \min(m, n).$$

Each product $\mathbf{u}_i \mathbf{v}_i^T$ is a matrix of rank one, and is called a rank-one update. Truncation of the above sum to p terms leads to an approximation of A

$$A \approx A_p = \sum_{i=1}^p \sigma_i \mathbf{u}_i \mathbf{v}_i^T.$$

In very many cases the singular values exhibit rapid, exponential decay, $\sigma_1 \gg \sigma_2 \gg \dots$, such that the approximation above is an accurate representation of the matrix A .



Figure 2. Successive SVD approximations of Frida Kahlo's (1907-1954) painting, *Portrait of a Lady in White* (1929), with $k = 10, 20, 40$ rank-one updates.

2.3. SVD solution of linear algebra problems

The SVD can be used to solve common problems within linear algebra.

Change of coordinates. To change from vector coordinates \mathbf{b} in the canonical basis $\mathbf{I} \in \mathbb{R}^{m \times m}$ to coordinates \mathbf{x} in some other basis $\mathbf{A} \in \mathbb{R}^{m \times m}$, a solution to the equation $\mathbf{I}\mathbf{b} = \mathbf{A}\mathbf{x}$ can be found by the following steps.

1. Compute the SVD, $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \mathbf{A}$;

2. Find the coordinates of \mathbf{b} in the orthogonal basis \mathbf{U} , $\mathbf{c} = \mathbf{U}^T \mathbf{b}$;
3. Scale the coordinates of \mathbf{c} by the inverse of the singular values $y_i = c_i / \sigma_i$, $i = 1, \dots, m$, such that $\mathbf{\Sigma} \mathbf{y} = \mathbf{c}$ is satisfied;
4. Find the coordinates of \mathbf{y} in basis \mathbf{V}^T , $\mathbf{x} = \mathbf{V} \mathbf{y}$.

Best 2-norm approximation. In the above \mathbf{A} was assumed to be a basis, hence $r = \text{rank}(\mathbf{A}) = m$. If columns of \mathbf{A} do not form a basis, $r < m$, then $\mathbf{b} \in \mathbb{R}^m$ might not be reachable by linear combinations within $C(\mathbf{A})$. The closest vector to \mathbf{b} in the norm is however found by the same steps, with the simple modification that in Step 3, the scaling is carried out only for non-zero singular values, $y_i = c_i / \sigma_i$, $i = 1, \dots, r$.

The pseudo-inverse. From the above, finding either the solution of $\mathbf{A} \mathbf{x} = \mathbf{I} \mathbf{b}$ or the best approximation possible if \mathbf{A} is not of full rank, can be written as a sequence of matrix multiplications using the SVD

$$(\mathbf{U} \mathbf{\Sigma} \mathbf{V}^T) \mathbf{x} = \mathbf{b} \Rightarrow \mathbf{U} (\mathbf{\Sigma} \mathbf{V}^T \mathbf{x}) = \mathbf{b} \Rightarrow (\mathbf{\Sigma} \mathbf{V}^T \mathbf{x}) = \mathbf{U}^T \mathbf{b} \Rightarrow \mathbf{V}^T \mathbf{x} = \mathbf{\Sigma}^+ \mathbf{U}^T \mathbf{b} \Rightarrow \mathbf{x} = \mathbf{V} \mathbf{\Sigma}^+ \mathbf{U}^T \mathbf{b},$$

where the matrix $\mathbf{\Sigma}^+ \in \mathbb{R}^{n \times m}$ (notice the inversion of dimensions) is defined as a matrix with elements σ_i^{-1} on the diagonal, and is called the pseudo-inverse of $\mathbf{\Sigma}$. Similarly the matrix

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

that allows stating the solution of $\mathbf{A} \mathbf{x} = \mathbf{b}$ simply as $\mathbf{x} = \mathbf{A}^+ \mathbf{b}$ is called the *pseudo-inverse* of \mathbf{A} . Note that in practice \mathbf{A}^+ is not explicitly formed. Rather the notation \mathbf{A}^+ is simply a concise reference to carrying out steps 1-4 above.