# 2

# Vectors and Vector Spaces

In this chapter we discuss a wide range of basic topics related to vectors of real numbers. Some of the properties carry over to vectors over other fields, such as complex numbers, but the reader should not assume this. Occasionally, for emphasis, we will refer to "real" vectors or "real" vector spaces, but unless it is stated otherwise, we are assuming the vectors and vector spaces are real. The topics and the properties of vectors and vector spaces that we emphasize are motivated by applications in the data sciences.

## 2.1 Operations on Vectors

The elements of the vectors we will use in the following are real numbers, that is, elements of $\mathbb{R}$. We call elements of $\mathbb{R}$ *scalars*. Vector operations are defined in terms of operations on real numbers.

Two vectors can be added if they have the same number of elements. The sum of two vectors is the vector whose elements are the sums of the corresponding elements of the vectors being added. Vectors with the same number of elements are said to be *conformable* for addition. A vector all of whose elements are 0 is the *additive identity* for all conformable vectors.

We overload the usual symbols for the operations on the reals to signify the corresponding operations on vectors or matrices when the operations are defined. Hence, "+" can mean addition of scalars, addition of conformable vectors, or addition of a scalar to a vector. This last meaning of "+" may not be used in many mathematical treatments of vectors, but it is consistent with the semantics of modern computer languages such as Fortran, R, and Matlab. By the *addition of a scalar and a vector*, we mean the addition of the scalar to each element of the vector, resulting in a vector of the same number of elements.

A *scalar multiple of a vector* (that is, the product of a real number and a vector) is the vector whose elements are the multiples of the corresponding elements of the original vector. Juxtaposition of a symbol for a scalar and a symbol for a vector indicates the multiplication of the scalar with each element of the vector, resulting in a vector of the same number of elements.

The basic operation in working with vectors is the addition of a scalar multiple of one vector to another vector,

$$z = ax + y, \tag{2.1}$$

where $a$ is a scalar and $x$ and $y$ are vectors conformable for addition. Viewed as a single operation with three operands, this is called an *axpy operation* for obvious reasons. (Because the Fortran versions of BLAS to perform this operation were called `saxpy` and `daxpy`, the operation is also sometimes called "saxpy" or "daxpy". See Sect. 12.2.1 on page 555, for a description of the BLAS.)

The axpy operation is a *linear combination*. Such linear combinations of vectors are the basic operations in most areas of linear algebra. The composition of axpy operations is also an axpy; that is, one linear combination followed by another linear combination is a linear combination. Furthermore, any linear combination can be decomposed into a sequence of axpy operations.

A special linear combination is called a *convex combination*. For vectors $x$ and $y$, it is the combination

$$ax + by, \tag{2.2}$$

where $a, b \geq 0$ and $a + b = 1$. A set of vectors that is closed with respect to convex combinations is said to be *convex*.

### 2.1.1 Linear Combinations and Linear Independence

If a given vector can be formed by a linear combination of one or more vectors, the set of vectors (including the given one) is said to be linearly dependent; conversely, if in a set of vectors no one vector can be represented as a linear combination of any of the others, the set of vectors is said to be *linearly independent*. In equation (2.1), for example, the vectors $x$, $y$, and $z$ are not linearly independent. It is possible, however, that any two of these vectors are linearly independent.

Linear independence is one of the most important concepts in linear algebra.

We can see that the definition of a linearly independent set of vectors $\{v_1, \ldots, v_k\}$ is equivalent to stating that if

$$a_1 v_1 + \cdots a_k v_k = 0, \tag{2.3}$$

then $a_1 = \cdots = a_k = 0$. If the set of vectors $\{v_1, \ldots, v_k\}$ is not linearly independent, then it is possible to select a *maximal linearly independent subset*;

that is, a subset of $\{v_1, \ldots, v_k\}$ that is linearly independent and has maximum cardinality. We do this by selecting an arbitrary vector, $v_{i_1}$, and then seeking a vector that is independent of $v_{i_1}$. If there are none in the set that is linearly independent of $v_{i_1}$, then a maximum linearly independent subset is just the singleton, because all of the vectors must be a linear combination of just one vector (that is, a scalar multiple of that one vector). If there is a vector that is linearly independent of $v_{i_1}$, say $v_{i_2}$, we next seek a vector in the remaining set that is independent of $v_{i_1}$ and $v_{i_2}$. If one does not exist, then $\{v_{i_1}, v_{i_2}\}$ is a maximal subset because any other vector can be represented in terms of these two and hence, within any subset of three vectors, one can be represented in terms of the two others. Thus, we see how to form a maximal linearly independent subset, and we see that the maximum cardinality of any subset of linearly independent vectors is unique however they are formed.

It is easy to see that the maximum number of $n$-vectors that can form a set that is linearly independent is $n$. (We can see this by assuming $n$ linearly independent vectors and then, for any $(n + 1)^{\text{th}}$ vector, showing that it is a linear combination of the others by building it up one by one from linear combinations of two of the given linearly independent vectors. In Exercise 2.1, you are asked to write out these steps.)

Properties of a set of vectors are usually invariant to a permutation of the elements of the vectors if the same permutation is applied to all vectors in the set. In particular, if a set of vectors is linearly independent, the set remains linearly independent if the elements of each vector are permuted in the same way.

If the elements of each vector in a set of vectors are separated into subvectors, linear independence of any set of corresponding subvectors implies linear independence of the full vectors. To state this more precisely for a set of three $n$-vectors, let $x = (x_1, \ldots, x_n)$, $y = (y_1, \ldots, y_n)$, and $z = (z_1, \ldots, z_n)$. Now let $\{i_1, \ldots, i_k\} \subseteq \{1, \ldots, n\}$, and form the $k$-vectors $\tilde{x} = (x_{i_1}, \ldots, x_{i_k})$, $\tilde{y} = (y_{i_1}, \ldots, y_{i_k})$, and $\tilde{z} = (z_{i_1}, \ldots, z_{i_k})$. Then linear independence of $\tilde{x}$, $\tilde{y}$, and $\tilde{z}$ implies linear independence of $x$, $y$, and $z$. (This can be shown directly from the definition of linear independence. It is related to equation (2.19) on page 20, which you are asked to prove in Exercise 2.5.)

### 2.1.2 Vector Spaces and Spaces of Vectors

Let $V$ be a set of $n$-vectors such that any linear combination of the vectors in $V$ is also in $V$. Such a set together with the usual vector algebra is called a *vector space*. A vector space is a *linear space*, and it necessarily includes the additive identity (the zero vector). (To see this, in the axpy operation, let $a = -1$ and $y = x$.) A vector space is necessarily convex.

The set consisting only of the additive identity, along with the axpy operation, is a vector space. It is called the "null vector space". Some people define "vector space" in a way that excludes it, because its properties do not conform to many general statements we can make about other vector spaces.

The "usual algebra" is a *linear algebra* consisting of two operations: vector addition and scalar times vector multiplication, which are the two operations comprising an axpy. It has closure of the space under the combination of those operations, commutativity and associativity of addition, an additive identity and inverses, a multiplicative identity, distribution of multiplication over both vector addition and scalar addition, and associativity of scalar multiplication and scalar times vector multiplication.

A vector space can also be composed of other objects, such as matrices, along with their appropriate operations. The key characteristic of a vector space is a linear algebra.

We generally use a calligraphic font to denote a vector space; $\mathcal{V}$ or $\mathcal{W}$, for example. Often, however, we think of the vector space merely in terms of the set of vectors on which it is built and denote it by an ordinary capital letter; $V$ or $W$, for example. A vector space is an *algebraic structure* consisting of a set together with the axpy operation, with the restriction that the set is closed under the operation. To indicate that it is a structure, rather than just a set, we may write

$$\mathcal{V} = (V, \circ),$$

where $V$ is just the set and $\circ$ denotes the axpy operation, or a similar linear operation under which the set is closed.

### 2.1.2.1 Generating Sets

Given a set $G$ of vectors of the same order, a vector space can be formed from the set $G$ together with all vectors that result from the axpy operation being applied to all combinations of vectors in $G$ and all values of the real number $a$; that is, for all $v_i, v_j \in G$ and all real $a$,

$$\{av_i + v_j\}.$$

This set together with the axpy operation itself is a vector space. It is called the *space generated by* $G$. We denote this space as

$$\mathrm{span}(G).$$

We will discuss generating and spanning sets further in Sect. 2.1.3.

### 2.1.2.2 The Order and the Dimension of a Vector Space

The vector space consisting of all $n$-vectors with real elements is denoted $\mathbb{R}^n$. (As mentioned earlier, the notation $\mathbb{R}^n$ can also refer to just the *set* of $n$-vectors with real elements; that is, to the set over which the vector space is defined.)

The *dimension of a vector space* is the maximum number of linearly independent vectors in the vector space. We denote the dimension by

$$\dim(\cdot),$$

which is a mapping $\mathrm{I\!R}^n \to \mathbb{Z}_+$ (where $\mathbb{Z}_+$ denotes the positive integers).

The *order of a vector space* is the order of the vectors in the space. Because the maximum number of $n$-vectors that can form a linearly independent set is $n$, as we showed above, the order of a vector space is greater than or equal to the dimension of the vector space.

Both the order and the dimension of $\mathrm{I\!R}^n$ are $n$. A set of $m$ linearly independent $n$-vectors with real elements can generate a vector space within $\mathrm{I\!R}^n$ of order $n$ and dimension $m$.

We also may use the phrase *dimension of a vector* to mean the dimension of the vector space of which the vector is an element. This term is ambiguous, but its meaning is clear in specific contexts, such as *dimension reduction*, that we will discuss later.

### 2.1.2.3 Vector Spaces with an Infinite Number of Dimensions

It is possible that no finite set of vectors span a given vector space. In that case, the vector space is said to be of infinite dimension.

Many of the properties of vector spaces that we discuss hold for those with an infinite number of dimensions; but not all do, such as the equivalence of norms (see page 29).

Throughout this book, however, unless we state otherwise, we assume the vector spaces have a finite number of dimensions.

### 2.1.2.4 Essentially Disjoint Vector Spaces

If the only element in common between two vector spaces $\mathcal{V}$ and $\mathcal{W}$ is the additive identity, the spaces are said to be *essentially disjoint*. Essentially disjoint vector spaces necessarily have the same order.

If the vector spaces $\mathcal{V}$ and $\mathcal{W}$ are essentially disjoint, it is clear that any element in $\mathcal{V}$ (except the additive identity) is linearly independent of any set of elements in $\mathcal{W}$.

### 2.1.2.5 Some Special Vectors: Notation

We denote the additive identity in a vector space of order $n$ by $0_n$ or sometimes by $0$. This is the vector consisting of all zeros:

$$0_n = (0, \ldots, 0). \tag{2.4}$$

We call this the *zero vector*, or the *null vector*. (A vector $x \neq 0$ is called a "nonnull vector".) This vector by itself is sometimes called the *null vector space*. It is not a vector space in the usual sense; it would have dimension 0. (All linear combinations are the same.)

Likewise, we denote the vector consisting of all ones by $1_n$ or sometimes by 1:

$$1_n = (1, \ldots, 1). \tag{2.5}$$

We call this the *one vector* and also the "summing vector" (see page 34). This vector and all scalar multiples of it are vector spaces with dimension 1. (This is true of any single nonzero vector; all linear combinations are just scalar multiples.) Whether 0 and 1 without a subscript represent vectors or scalars is usually clear from the context.

The zero vector and the one vector are both instances of *constant vectors*; that is, vectors all of whose elements are the same. In some cases we may abuse the notation slightly, as we have done with "0" and "1" above, and use a single symbol to denote both a scalar and a vector all of whose elements are that constant; for example, if "$c$" denotes a scalar constant, we may refer to the vector all of whose elements are $c$ as "$c$" also. These notational conveniences rarely result in any ambiguity. They also allow another interpretation of the definition of addition of a scalar to a vector that we mentioned at the beginning of the chapter.

The $i^{\text{th}}$ *unit vector*, denoted by $e_i$, has a 1 in the $i^{\text{th}}$ position and 0s in all other positions:

$$e_i = (0, \ldots, 0, 1, 0, \ldots, 0). \tag{2.6}$$

Another useful vector is the *sign vector*, which is formed from signs of the elements of a given vector. It is denoted by "sign($\cdot$)" and for $x = (x_1, \ldots, x_n)$ is defined by

$$\begin{aligned}
\text{sign}(x)_i &= 1 && \text{if } x_i > 0, \\
&= 0 && \text{if } x_i = 0, \\
&= -1 && \text{if } x_i < 0.
\end{aligned} \tag{2.7}$$

### 2.1.2.6 Ordinal Relations Among Vectors

There are several possible ways to form a rank ordering of vectors of the same order, but no complete ordering is entirely satisfactory. (Note the unfortunate overloading of the words "order" and "ordering" here.) If $x$ and $y$ are vectors of the same order and for corresponding elements $x_i > y_i$, we say $x$ is *greater than $y$* and write

$$x > y. \tag{2.8}$$

In particular, if all of the elements of $x$ are positive, we write $x > 0$.

If $x$ and $y$ are vectors of the same order and for corresponding elements $x_i \geq y_i$, we say $x$ is *greater than or equal to $y$* and write

$$x \geq y. \tag{2.9}$$

This relationship is a *partial ordering* (see Exercise 8.2a on page 396 for the definition of partial ordering).

The expression $x \geq 0$ means that all of the elements of $x$ are nonnegative.

### 2.1.2.7 Set Operations on Vector Spaces

The ordinary operations of subsetting, intersection, union, direct sum, and direct product for sets have analogs for vector spaces, and we use some of the same notation to refer to vector spaces that we use to refer to sets. The set operations themselves are performed on the individual sets to yield a set of vectors, and the resulting vector space is the space generated by that set of vectors.

Unfortunately, there are many inconsistencies in terminology used in the literature regarding operations on vector spaces. When I use a term and/or symbol, such as "union" or "$\cup$", for a structure such as a vector space, I use it in reference to the *structure*. For example, if $\mathcal{V} = (V, \circ)$ and $\mathcal{W} = (W, \circ)$ are vector spaces, then $V \cup U$ is the ordinary union of the sets; however, $\mathcal{V} \cup \mathcal{W}$ is the union of the vector spaces, and is not necessarily the same as $(U \cup W, \circ)$, which may not even be a vector space. Occasionally in the following discussion, I will try to point out common variants in usage.

The convention that I follow allows the wellknown relationships among common set operations to hold for the corresponding operations on vector spaces; for example, if $\mathcal{V}$ and $\mathcal{W}$ are vector spaces, $\mathcal{V} \subseteq \mathcal{V} \cup \mathcal{W}$, just as for sets $V$ and $W$.

The properties of vector spaces are proven the same way that properties of sets are proven, after first requiring that the axpy operation have the same meaning in the different vector spaces. For example, to prove that one vector space is a subspace of another, we show that any given vector in the first vector space is necessarily in the second. To prove that two vector spaces are equal, we show that each is a subspace of the other. Some properties of vector spaces and subspaces can be shown more easily using "basis sets" for the spaces, which we discuss in Sect. 2.1.3, beginning on page 21.

Note that if $(V, \circ)$ and $(W, \circ)$ are vector spaces of the same order and $U$ is some set formed by an operation on $V$ and $W$, then $(U, \circ)$ may not be a vector space because it is not closed under the axpy operation, $\circ$. We sometimes refer to a set of vectors of the same order together with the axpy operator (whether or not the set is closed with respect to the operator) as a "space of vectors" (instead of a "vector space").

### 2.1.2.8 Subpaces

Given a vector space $\mathcal{V} = (V, \circ)$, if $W$ is any subset of $V$, then the vector space $\mathcal{W}$ generated by $W$, that is, $\text{span}(W)$, is said to be a *subspace* of $\mathcal{V}$, and we denote this relationship by $\mathcal{W} \subseteq \mathcal{V}$.

If $\mathcal{W} \subseteq \mathcal{V}$ and $\mathcal{W} \neq \mathcal{V}$, then $\mathcal{W}$ is said to be a *proper subspace* of $\mathcal{V}$. If $\mathcal{W} = \mathcal{V}$, then $\mathcal{W} \subseteq \mathcal{V}$ and $\mathcal{V} \subseteq \mathcal{W}$, and the converse is also true.

The maximum number of linearly independent vectors in the subspace cannot be greater than the maximum number of linearly independent vectors in the original space; that is, if $\mathcal{W} \subseteq \mathcal{V}$, then

$$\dim(\mathcal{W}) \leq \dim(\mathcal{V}) \tag{2.10}$$

(Exercise 2.2). If $\mathcal{W}$ is a proper subspace of $\mathcal{V}$, then $\dim(\mathcal{W}) < \dim(\mathcal{V})$.

### 2.1.2.9 Intersections of Vector Spaces

For two vector spaces $\mathcal{V}$ and $\mathcal{W}$ of the same order with vectors formed from the same field, we define their *intersection*, denoted by $\mathcal{V} \cap \mathcal{W}$, to be the set of vectors consisting of the intersection of the sets in the individual vector spaces together with the axpy operation.

The intersection of two vector spaces of the same order that are not essentially disjoint is a vector space, as we can see by letting $x$ and $y$ be any vectors in the intersection $\mathcal{U} = \mathcal{V} \cap \mathcal{W}$, and showing, for any real number $a$, that $ax + y \in \mathcal{U}$. This is easy because both $x$ and $y$ must be in both $\mathcal{V}$ and $\mathcal{W}$.

Note that if $\mathcal{V}$ and $\mathcal{W}$ are essentially disjoint, then $\mathcal{V} \cap \mathcal{W} = (0, \circ)$, which, as we have said, is not a vector space in the usual sense.

Also note that

$$\dim(\mathcal{V} \cap \mathcal{W}) \leq \min(\dim(\mathcal{V}), \dim(\mathcal{W})) \tag{2.11}$$

(Exercise 2.2).

### 2.1.2.10 Unions and Direct Sums of Vector Spaces

Given two vector spaces $\mathcal{V}$ and $\mathcal{W}$ of the same order, we define their *union*, denoted by $\mathcal{V} \cup \mathcal{W}$, to be the vector space generated by the union of the sets in the individual vector spaces together with the axpy operation. If $\mathcal{V} = (V, \circ)$ and $\mathcal{W} = (W, \circ)$, this is the vector space generated by the set of vectors $V \cup W$; that is,

$$\mathcal{V} \cup \mathcal{W} = \mathrm{span}(V \cup W). \tag{2.12}$$

The union of the sets of vectors in two vector spaces may not be closed under the axpy operation (Exercise 2.3b), but the union of vector spaces *is* a vector space by definition.

The vector space generated by the union of the sets in the individual vector spaces is easy to form. Since $(V, \circ)$ and $(W, \circ)$ are vector spaces (so for any vector $x$ in either $V$ or $W$, $ax$ is in that set), all we need do is just include all simple sums of the vectors from the individual sets, that is,

$$\mathcal{V} \cup \mathcal{W} = \{v + w, \text{ s.t. } v \in \mathcal{V}, \ w \in \mathcal{W}\}. \tag{2.13}$$

It is easy to see that this is a vector space by showing that it is closed with respect to axpy. (As above, we show that for any $x$ and $y$ in $\mathcal{V} \cup \mathcal{W}$ and for any real number $a$, $ax + y$ is in $\mathcal{V} \cup \mathcal{W}$.)

(Because of the way the union of vector spaces can be formed from simple addition of the individual elements, some authors call the vector space in

equation (2.13) the "sum" of $\mathcal{V}$ and $\mathcal{W}$, and write it as $\mathcal{V} + \mathcal{W}$. Other authors, including myself, call this the *direct sum*, and denote it by $\mathcal{V} \oplus \mathcal{W}$. Some authors define "direct sum" only in the cases of vector spaces that are essentially disjoint. Still other authors define "direct sum" to be what I will call a "direct product" below.)

Despite the possible confusion with other uses of the notation, I often use the notation $\mathcal{V} \oplus \mathcal{W}$ because it points directly to the nice construction of equation (2.13). *To be clear:* to the extent that I use "direct sum" and "$\oplus$" for vector spaces $\mathcal{V}$ and $\mathcal{W}$, I will mean the direct sum

$$\mathcal{V} \oplus \mathcal{W} \equiv \mathcal{V} \cup \mathcal{W}, \tag{2.14}$$

as defined above.

Note that

$$\dim(\mathcal{V} \oplus \mathcal{W}) = \dim(\mathcal{V}) + \dim(\mathcal{W}) - \dim(\mathcal{V} \cap \mathcal{W}) \tag{2.15}$$

(Exercise 2.4). Therefore

$$\dim(\mathcal{V} \oplus \mathcal{W}) \geq \max(\dim(\mathcal{V}), \dim(\mathcal{W}))$$

and

$$\dim(\mathcal{V} \oplus \mathcal{W}) \leq \dim(\mathcal{V}) + \dim(\mathcal{W}).$$

### 2.1.2.11 Direct Sum Decomposition of a Vector Space

In some applications, given a vector space $\mathcal{V}$, it is of interest to find essentially disjoint vector spaces $\mathcal{V}_1, \ldots, \mathcal{V}_n$ such that

$$\mathcal{V} = \mathcal{V}_1 \oplus \cdots \oplus \mathcal{V}_n.$$

This is called a *direct sum decomposition* of $\mathcal{V}$. (As I mentioned above, some authors who do not use "direct sum" as I do would use the term in this context because the individual matrices are essentially disjoint.)

It is clear that if $\mathcal{V}_1, \ldots, \mathcal{V}_n$ is a direct sum decomposition of $\mathcal{V}$, then

$$\dim(\mathcal{V}) = \sum_{i=1}^{n} \dim(\mathcal{V}_i) \tag{2.16}$$

(Exercise 2.4).

A collection of essentially disjoint vector spaces $\mathcal{V}_1, \ldots, \mathcal{V}_n$ such that $\mathcal{V} = \mathcal{V}_1 \oplus \cdots \oplus \mathcal{V}_n$ is said to be *complementary with respect to $\mathcal{V}$*.

An important property of a direct sum decomposition is that it allows a unique representation of a vector in the decomposed space in terms of a sum of vectors from the individual essentially disjoint spaces; that is, if $\mathcal{V} = \mathcal{V}_1 \oplus \cdots \oplus \mathcal{V}_n$ is a direct sum decomposition of $\mathcal{V}$ and $v \in \mathcal{V}$, then there exist unique vectors $v_i \in \mathcal{V}_i$ such that

$$v = v_1 + \cdots + v_n. \tag{2.17}$$

We will prove this for the case $n = 2$. This is without loss, because additional spaces in the decomposition add nothing different.

Given the direct sum decomposition $\mathcal{V} = \mathcal{V}_1 \oplus \mathcal{V}_2$, let $v$ be any vector in $\mathcal{V}$. Because $\mathcal{V}_1 \oplus \mathcal{V}_2$ can be formed as in equation (2.13), there exist vectors $v_1 \in \mathcal{V}_1$ and $v_2 \in \mathcal{V}_2$ such that $v = v_1 + v_2$. Now all we need to do is to show that they are unique.

Let $u_1 \in \mathcal{V}_1$ and $u_2 \in \mathcal{V}_2$ be such that $v = u_1 + u_2$. Now we have $(v - u_1) \in \mathcal{V}_2$ and $(v - v_1) \in \mathcal{V}_2$; hence $(v_1 - u_1) \in \mathcal{V}_2$. However, since $v_1, u_1 \in \mathcal{V}_1$, $(v_1 - u_1) \in \mathcal{V}_1$. Since $\mathcal{V}_1$ and $\mathcal{V}_2$ are essentially disjoint, and $(v_1 - u_1)$ is in both, it must be the case that $(v_1 - u_1) = 0$, or $u_1 = v_1$. In like manner, we show that $u_2 = v_2$; hence, the representation $v = v_1 + v_2$ is unique.

An important fact is that for any vector space $\mathcal{V}$ with dimension 2 or greater, a direct sum decomposition exists; that is, there exist essentially disjoint vector spaces $\mathcal{V}_1$ and $\mathcal{V}_2$ such that $\mathcal{V} = \mathcal{V}_1 \oplus \mathcal{V}_2$.

This is easily shown by first choosing a proper subspace $\mathcal{V}_1$ of $\mathcal{V}$ and then constructing an essentially disjoint subspace $\mathcal{V}_2$ such that $\mathcal{V} = \mathcal{V}_1 \oplus \mathcal{V}_2$. The details of these steps are made simpler by use of basis sets which we will discuss in Sect. 2.1.3, in particular the facts listed on page 22.

### 2.1.2.12 Direct Products of Vector Spaces and Dimension Reduction

The set operations on vector spaces that we have mentioned so far require that the vector spaces be of a fixed order. Sometimes in applications, it is useful to deal with vector spaces of different orders.

The *direct product* of the vector space $\mathcal{V}$ of order $n$ and the vector space $\mathcal{W}$ of order $m$ is the vector space of order $n + m$ on the set of vectors

$$\{(v_1, \ldots, v_n, w_1, \ldots, w_m), \text{ s.t. } (v_1, \ldots, v_n) \in \mathcal{V}, (w_1, \ldots, w_m) \in \mathcal{W}\}, \tag{2.18}$$

together with the axpy operator defined as the same operator in $\mathcal{V}$ and $\mathcal{W}$ applied separately to the first $n$ and the last $m$ elements. The direct product of $\mathcal{V}$ and $\mathcal{W}$ is denoted by $\mathcal{V} \otimes \mathcal{W}$.

Notice that while the direct sum operation is commutative, the direct product is not commutative in general.

The vectors in $\mathcal{V}$ and $\mathcal{W}$ are sometimes called "subvectors" of the vectors in $\mathcal{V} \otimes \mathcal{W}$. These subvectors are related to projections, which we will discuss in more detail in Sect. 2.2.2 (page 36) and Sect. 8.5.2 (page 358).

We can see that the direct product is a vector space using the same method as above by showing that it is closed under the axpy operation.

Note that

$$\dim(\mathcal{V} \otimes \mathcal{W}) = \dim(\mathcal{V}) + \dim(\mathcal{W}) \tag{2.19}$$

(Exercise 2.5).

Note that for integers $0 < p < n$,

$$\mathbb{R}^n = \mathbb{R}^p \otimes \mathbb{R}^{n-p}, \tag{2.20}$$

where the operations in the space $\mathbb{R}^n$ are the same as in the component vector spaces with the meaning adjusted to conform to the larger order of the vectors in $\mathbb{R}^n$. (Recall that $\mathbb{R}^n$ represents the algebraic structure consisting of the set of $n$-tuples of real numbers plus the special axpy operator.)

In statistical applications, we often want to do "dimension reduction". This means to find a smaller number of coordinates that cover the relevant regions of a larger-dimensional space. In other words, we are interested in finding a lower-dimensional vector space in which a given set of vectors in a higher-dimensional vector space can be approximated by vectors in the lower-dimensional space. For a given set of vectors of the form $x = (x_1, \ldots, x_n)$ we seek a set of vectors of the form $z = (z_1, \ldots, z_p)$ that almost "cover the same space". (The transformation from $x$ to $z$ is called a projection.)

### 2.1.3 Basis Sets for Vector Spaces

If each vector in the vector space $\mathcal{V}$ can be expressed as a linear combination of the vectors in some set $G$, then $G$ is said to be a *generating set* or *spanning set* of $\mathcal{V}$. The number of vectors in a generating set is at least as great as the dimension of the vector space.

If all linear combinations of the elements of $G$ are in $\mathcal{V}$, the vector space is the *space generated by $G$* and is denoted by $\mathcal{V}(G)$ or by $\mathrm{span}(G)$, as we mentioned on page 14. We will use either notation interchangeably:

$$\mathcal{V}(G) \equiv \mathrm{span}(G). \tag{2.21}$$

Note that $G$ is also a generating or spanning set for $\mathcal{W}$ where $\mathcal{W} \subseteq \mathrm{span}(G)$.

A *basis* for a vector space is a set of linearly independent vectors that generate or span the space. For any vector space, a generating set consisting of the minimum number of vectors of any generating set for that space is a basis set for the space. A basis set is obviously not unique.

Note that the linear independence implies that a basis set cannot contain the 0 vector.

An important fact is

- The representation of a given vector in terms of a given basis set is unique.

To see this, let $\{v_1, \ldots, v_k\}$ be a basis for a vector space that includes the vector $x$, and let

$$x = c_1 v_1 + \cdots c_k v_k.$$

Now suppose

$$x = b_1 v_1 + \cdots b_k v_k,$$

so that we have
$$0 = (c_1 - b_1)v_1 + \cdots + (c_k - b_k)v_k.$$

Since $\{v_1, \ldots, v_k\}$ are independent, the only way this is possible is if $c_i = b_i$ for each $i$.

A related fact is that if $\{v_1, \ldots, v_k\}$ is a basis for a vector space of order $n$ that includes the vector $x$ and $x = c_1v_1 + \cdots c_kv_k$, then $x = 0_n$ if and only if $c_i = 0$ for each $i$.

For any vector space, the order of the vectors in a basis set is the same as the order of the vector space.

Because the vectors in a basis set are independent, the number of vectors in a basis set is the same as the dimension of the vector space; that is, if $B$ is a basis set of the vector space $\mathcal{V}$, then

$$\dim(\mathcal{V}) = \#(B). \tag{2.22}$$

A simple basis set for the vector space $\mathbb{R}^n$ is the set of unit vectors $\{e_1, \ldots, e_n\}$, defined on page 16.

### 2.1.3.1 Properties of Basis Sets of Vector Subspaces

There are several interesting facts about basis sets for vector spaces and various combinations of the vector spaces. Verifications of these facts all follow similar arguments, and most are left as exercises.

- If $B_1$ is a basis set for $\mathcal{V}_1$, $B_2$ is a basis set for $\mathcal{V}_2$, and $\mathcal{V}_1$ and $\mathcal{V}_2$ are essentially disjoint, then $B_1 \cap B_2 = \emptyset$.
  This fact is easily seen by assuming the contrary; that is, assume that $b \in B_1 \cap B_2$. (Note that $b$ cannot be the 0 vector.) This implies, however, that $b$ is in both $\mathcal{V}_1$ and $\mathcal{V}_2$, contradicting the hypothesis that they are essentially disjoint.
- If $B$ is a basis set for $\mathcal{V}$ and $\mathcal{V}_1 \subseteq \mathcal{V}$, then there exists $B_1$, with $B_1 \subseteq B$, such that $B_1$ is a basis set for $\mathcal{V}_1$.
- If $B_1$ is a basis set for $\mathcal{V}_1$ and $B_2$ is a basis set for $\mathcal{V}_2$, then $B_1 \cup B_2$ is a generating set for $\mathcal{V}_1 \oplus \mathcal{V}_2$.
  (We see this easily from the definition of $\oplus$ because any vector in $\mathcal{V}_1 \oplus \mathcal{V}_2$ can be represented as a linear combination of vectors in $B_1$ plus a linear combination of vectors in $B_2$.)
- If $\mathcal{V}_1$ and $\mathcal{V}_2$ are essentially disjoint, $B_1$ is a basis set for $\mathcal{V}_1$, and $B_2$ is a basis set for $\mathcal{V}_2$, then $B_1 \cup B_2$ is a basis set for $\mathcal{V} = \mathcal{V}_1 \oplus \mathcal{V}_2$.
  This is the case that $\mathcal{V}_1 \oplus \mathcal{V}_2$ is a direct sum decomposition of $\mathcal{V}$.
- Suppose $\mathcal{V}_1$ is a real vector space of order $n_1$ (that is, it is a subspace of $\mathbb{R}^{n_1}$) and $B_1$ is a basis set for $\mathcal{V}_1$. Now let $\mathcal{V}_2$ be a real vector space of order $n_2$ and $B_2$ be a basis set for $\mathcal{V}_2$. For each vector $b_1$ in $B_1$ form the vector
  $$\tilde{b}_1 = (b_1|0, \ldots, 0) \quad \text{where there are } n_2 \text{ 0s,}$$

and let $\widetilde{B}_1$ be the set of all such vectors. (The order of each $\tilde{b}_1 \in \widetilde{B}_1$ is $n_1 + n_2$.) Likewise, for each vector $b_2$ in $B_2$ form the vector

$$\tilde{b}_2 = (0, \ldots, 0 | b_2) \quad \text{where there are } n_1 \text{ 0s},$$

and let $\widetilde{B}_2$ be the set of all such vectors. Then $\widetilde{B}_1 \cup \widetilde{B}_2$ is a basis for $\mathcal{V}_1 \otimes \mathcal{V}_2$.

### 2.1.4 Inner Products

A useful operation on two vectors $x$ and $y$ of the same order is the *inner product*, which we denote by $\langle x, y \rangle$ and define as

$$\langle x, y \rangle = \sum_i x_i \bar{y}_i, \tag{2.23}$$

where $\bar{z}$ represents the complex conjugate of $z$; that is, if $z = a + bi$, then $\bar{z} = a - bi$. In general, throughout this book unless stated otherwise, I assume that we are working with real numbers, and hence, $\bar{z} = z$. Most statements will hold whether the numbers are real or complex. When the statements only hold for reals, I will generally include the exception in the statement. The main differences have to do with inner products and an important property defined in terms of an inner product, called *orthogonality*.

In the case of vectors with real elements, we have

$$\langle x, y \rangle = \sum_i x_i y_i. \tag{2.24}$$

In that case (which is what we generally assume throughout this book), the inner product is a mapping

$$\mathbb{R}^n \times \mathbb{R}^n \to \mathbb{R}.$$

The inner product is also called the dot product or the scalar product. The dot product is actually a special type of inner product, and there is some ambiguity in the terminology. The dot product is the most commonly used inner product in the applications we consider, and so we will use the terms synonymously.

The inner product is also sometimes written as $x \cdot y$, hence the name dot product. Yet another notation for the inner product for real vectors is $x^\mathrm{T} y$, and we will see later that this notation is natural in the context of matrix multiplication. So for real vectors, we have the equivalent notations

$$\langle x, y \rangle \equiv x \cdot y \equiv x^\mathrm{T} y. \tag{2.25}$$

(I will mention one more notation that is equivalent for real vectors. This is the "bra·ket" notation originated by Paul Dirac, and is still used in certain

areas of application. Dirac referred to $x^{\mathrm{T}}$ as the "bra $x$", and denoted it as $\langle x|$. He referred to an ordinary vector $y$ as the "ket $y$", and denoted it as $|y\rangle$. He then denoted the inner product of the vectors as $\langle x||y\rangle$, or, omitting one vertical bar, as $\langle x|y\rangle$.)

In general, the inner product is a mapping from a real vector space $\mathcal{V}$ to $\mathbb{R}$ that has the following properties:

1.  Nonnegativity and mapping of the additive identity:
    if $x \neq 0$, then $\langle x, x \rangle > 0$ and $\langle 0, x \rangle = \langle x, 0 \rangle = \langle 0, 0 \rangle = 0$.
2.  Commutativity:
    $\langle x, y \rangle = \langle y, x \rangle$.
3.  Factoring of scalar multiplication in dot products:
    $\langle ax, y \rangle = a\langle x, y \rangle$ for real $a$.
4.  Relation of vector addition to addition of dot products:
    $\langle x + y, z \rangle = \langle x, z \rangle + \langle y, z \rangle$.

These properties in fact define an *inner product* for mathematical objects for which an addition, an additive identity, and a multiplication by a scalar are defined. Notice that the operation defined in equation (2.24) is not an inner product for vectors over the complex field because, if $x$ is complex, we can have $\langle x, x \rangle = 0$ when $x \neq 0$.

A vector space together with an inner product is called an *inner product space*.

Inner products are also defined for matrices, as we will discuss on page 97. We should note in passing that there are two different kinds of multiplication used in property 3. The first multiplication is scalar multiplication, that is, an operation from $\mathbb{R} \times \mathbb{R}^n$ to $\mathbb{R}^n$, which we have defined above, and the second multiplication is ordinary multiplication in $\mathbb{R}$, that is, an operation from $\mathbb{R} \times \mathbb{R}$ to $\mathbb{R}$. There are also two different kinds of addition used in property 4. The first addition is vector addition, defined above, and the second addition is ordinary addition in $\mathbb{R}$. The dot product can reveal fundamental relationships between the two vectors, as we will see later.

A useful property of inner products is the *Cauchy-Schwarz inequality*:

$$\langle x, y \rangle \leq \langle x, x \rangle^{\frac{1}{2}} \langle y, y \rangle^{\frac{1}{2}}. \tag{2.26}$$

This relationship is also sometimes called the Cauchy-Bunyakovskii-Schwarz inequality. (Augustin-Louis Cauchy gave the inequality for the kind of discrete inner products we are considering here, and Viktor Bunyakovskii and Hermann Schwarz independently extended it to more general inner products, defined on functions, for example.) The inequality is easy to see, by first observing that for every real number $t$,

$$\begin{aligned}
0 &\leq \langle (tx + y), (tx + y) \rangle \\
&= \langle x, x \rangle t^2 + 2\langle x, y \rangle t + \langle y, y \rangle \\
&= at^2 + bt + c,
\end{aligned}$$

where the constants $a$, $b$, and $c$ correspond to the dot products in the preceding equation. This quadratic in $t$ cannot have two distinct real roots. Hence the discriminant, $b^2 - 4ac$, must be less than or equal to zero; that is,

$$\left(\frac{1}{2}b\right)^2 \leq ac.$$

By substituting and taking square roots, we get the Cauchy-Schwarz inequality. It is also clear from this proof that equality holds only if $x = 0$ or if $y = rx$, for some scalar $r$.

Two vectors $x$ and $y$ such that $\langle x, y \rangle = 0$ are said to be *orthogonal*. This term has such an intuitive meaning that we may use it prior to a careful definition and study, so I only introduce it here. We will discuss orthogonality more thoroughly in Sect. 2.1.8 beginning on page 33.

### 2.1.5 Norms

We consider a set of objects $S$ that has an addition-type operator, $+$, a corresponding additive identity, 0, and a scalar multiplication; that is, a multiplication of the objects by a real (or complex) number. On such a set, a *norm* is a function, $\|\cdot\|$, from $S$ to $\mathbb{R}$ that satisfies the following three conditions:

1.  Nonnegativity and mapping of the additive identity:
    if $x \neq 0$, then $\|x\| > 0$, and $\|0\| = 0$.
2.  Relation of scalar multiplication to real multiplication:
    $\|ax\| = |a|\,\|x\|$ for real $a$.
3.  Triangle inequality:
    $\|x + y\| \leq \|x\| + \|y\|$.

(If property 1 is relaxed to require only $\|x\| \geq 0$ for $x \neq 0$, the function is called a *seminorm*.) Because a norm is a function whose argument is a vector, we also often use a functional notation such as $\rho(x)$ to represent a norm of the vector $x$.

Sets of various types of objects (functions, for example) can have norms, but our interest in the present context is in norms for vectors and (later) for matrices. (The three properties above in fact define a more general norm for other kinds of mathematical objects for which an addition, an additive identity, and multiplication by a scalar are defined. Norms are defined for matrices, as we will discuss later. Note that there are two different kinds of multiplication used in property 2 and two different kinds of addition used in property 3.)

A vector space together with a norm is called a *normed space.*

For some types of objects, a norm of an object may be called its "length" or its "size". (Recall the ambiguity of "length" of a vector that we mentioned at the beginning of this chapter.)

### 2.1.5.1 Convexity

A function $f(\cdot)$ over a convex domain $S$ into a range $R$, where both $S$ and $R$ have an addition-type operator, $+$, corresponding additive identities, and scalar multiplication, is said to be *convex*, if, for any $x$ and $y$ in $S$, and $a$ such that $0 \leq a \leq 1$,

$$f(ax + (1-a)y) \leq af(x) + (1-a)f(y). \tag{2.27}$$

If, for $x \neq y$ and $a$ such that $0 < a < 1$, the inequality in (2.27) is sharp, then the function is said to be *strictly convex*.

It is clear from the triangle inequality that a norm is convex.

### 2.1.5.2 Norms Induced by Inner Products

There is a close relationship between a norm and an inner product. For any inner product space with inner product $\langle \cdot, \cdot \rangle$, a norm of an element of the space can be defined in terms of the square root of the inner product of the element with itself:

$$\|x\| = \sqrt{\langle x, x \rangle}. \tag{2.28}$$

Any function $\| \cdot \|$ defined in this way satisfies the properties of a norm. It is easy to see that $\|x\|$ satisfies the first two properties of a norm, nonnegativity and scalar equivariance. Now, consider the square of the right-hand side of the triangle inequality, $\|x\| + \|y\|$:

$$\begin{aligned}
(\|x\| + \|y\|)^2 &= \langle x, x \rangle + 2\sqrt{\langle x, x \rangle \langle y, y \rangle} + \langle y, y \rangle \\
&\geq \langle x, x \rangle + 2\langle x, y \rangle + \langle y, y \rangle \\
&= \langle x+y, \, x+y \rangle \\
&= \|x+y\|^2;
\end{aligned} \tag{2.29}$$

hence, the triangle inequality holds. Therefore, given an inner product, $\langle x, y \rangle$, then $\sqrt{\langle x, x \rangle}$ is a norm.

Equation (2.28) defines a norm given any inner product. It is called the *norm induced by the inner product.*

Norms induced by inner products have some interesting properties. First of all, they have the Cauchy-Schwarz relationship (inequality (2.26)) with their associated inner product:

$$|\langle x, y \rangle| \leq \|x\|\|y\|. \tag{2.30}$$

In the sequence of equations above for an induced norm of the sum of two vectors, one equation (expressed differently) stands out as particularly useful in later applications:

$$\|x + y\|^2 = \|x\|^2 + \|y\|^2 + 2\langle x, y \rangle. \tag{2.31}$$

If $\langle x, y \rangle = 0$ (that is, the vectors are orthogonal), equation (2.31) becomes the *Pythagorean theorem*:

$$\|x + y\|^2 = \|x\|^2 + \|y\|^2.$$

Another useful property of a norm induced by an inner product is the *parallelogram equality*:

$$2\|x\|^2 + 2\|y\|^2 = \|x + y\|^2 + \|x - y\|^2. \tag{2.32}$$

This is trivial to show, and you are asked to do so in Exercise 2.7. (It is also the case that if the parallelogram equality holds for every pair of vectors in the space, then the norm is necessarily induced by an inner product. This fact is both harder to show and less useful than its converse; I state it only because it is somewhat surprising.)

A vector space whose norm is induced by an inner product has an interesting structure; for example, the geometric properties such as projections, orthogonality, and angles between vectors that we discuss in Sect. 2.2 are defined in terms of inner products and the associated norm.

### 2.1.5.3 $L_p$ Norms

There are many norms that could be defined for vectors. One type of norm is called an $L_p$ norm, often denoted as $\| \cdot \|_p$. For $p \geq 1$, it is defined as

$$\|x\|_p = \left( \sum_i |x_i|^p \right)^{\frac{1}{p}}. \tag{2.33}$$

This is also sometimes called the *Minkowski norm* and also the *Hölder norm*. An $L_p$ norm is also called a *p*-norm, or 1-norm, 2-norm, or $\infty$-norm (defined by a limit) in those special cases.

It is easy to see that the $L_p$ norm satisfies the first two conditions above. For general $p \geq 1$ it is somewhat more difficult to prove the triangular inequality (which for the $L_p$ norms is also called the Minkowski inequality), but for some special cases it is straightforward, as we will see below.

The most common $L_p$ norms, and in fact the most commonly used vector norms, are:

- $\|x\|_1 = \sum_i |x_i|$, also called the *Manhattan norm* because it corresponds to sums of distances along coordinate axes, as one would travel along the rectangular street plan of Manhattan (except for Broadway and a few other streets and avenues).
- $\|x\|_2 = \sqrt{\sum_i x_i^2}$, also called the *Euclidean norm*, the *Euclidean length*, or just the *length* of the vector. The $L_2$ norm is induced by an inner product; it is the square root of the inner product of the vector with itself: $\|x\|_2 = \sqrt{\langle x, x \rangle}$. It is the only $L_p$ norm induced by an inner product. (See Exercise 2.9.)

- $\|x\|_\infty = \max_i |x_i|$, also called the *max norm* or the *Chebyshev norm*. The $L_\infty$ norm is defined by taking the limit in an $L_p$ norm, and we see that it is indeed $\max_i |x_i|$ by expressing it as

$$\|x\|_\infty = \lim_{p\to\infty} \|x\|_p = \lim_{p\to\infty} \left( \sum_i |x_i|^p \right)^{\frac{1}{p}} = m \lim_{p\to\infty} \left( \sum_i \left| \frac{x_i}{m} \right|^p \right)^{\frac{1}{p}}$$

with $m = \max_i |x_i|$. Because the quantity of which we are taking the $p^{\text{th}}$ root is bounded above by the number of elements in $x$ and below by 1, that factor goes to 1 as $p$ goes to $\infty$.

It is easy to see that, for any $n$-vector $x$, the $L_p$ norms have the relationships

$$\|x\|_\infty \ \leq\ \|x\|_2 \ \leq\ \|x\|_1. \tag{2.34}$$

More generally, for given $x$ and for $p \geq 1$, we see that $\|x\|_p$ is a nonincreasing function of $p$.

We also have bounds that involve the number of elements in the vector:

$$\|x\|_\infty \ \leq\ \|x\|_2 \ \leq\ \sqrt{n}\|x\|_\infty, \tag{2.35}$$

and

$$\|x\|_2 \ \leq\ \|x\|_1 \ \leq\ \sqrt{n}\|x\|_2. \tag{2.36}$$

The triangle inequality obviously holds for the $L_1$ and $L_\infty$ norms. For the $L_2$ norm it can be seen by expanding $\sum(x_i + y_i)^2$ and then using the Cauchy-Schwarz inequality (2.26) on page 24. Rather than approaching it that way, however, we will show below that the $L_2$ norm can be defined in terms of an inner product, and then we will establish the triangle inequality for any norm defined similarly by an inner product; see inequality (2.29). Showing that the triangle inequality holds for other $L_p$ norms is more difficult; see Exercise 2.11.

A generalization of the $L_p$ vector norm is the *weighted $L_p$ vector norm* defined by

$$\|x\|_{wp} = \left( \sum_i w_i |x_i|^p \right)^{\frac{1}{p}}, \tag{2.37}$$

where $w_i \geq 0$ and $\sum_i w_i = 1$.

In the following, if we use the unqualified symbol $\|\cdot\|$ for a vector norm and do not state otherwise, we will mean the $L_2$ norm; that is, the Euclidean norm, the induced norm.

### 2.1.5.4 Basis Norms

If $\{v_1, \ldots, v_k\}$ is a basis for a vector space that includes a vector $x$ with $x = c_1 v_1 + \cdots + c_k v_k$, then

$$\rho(x) = \left( \sum_i c_i^2 \right)^{\frac{1}{2}} \tag{2.38}$$

is a norm. It is straightforward to see that $\rho(x)$ is a norm by checking the following three conditions:

- $\rho(x) \geq 0$ and $\rho(x) = 0$ if and only if $x = 0$ because $x = 0$ if and only if $c_i = 0$ for all $i$.
- $\rho(ax) = \left( \sum_i a^2 c_i^2 \right)^{\frac{1}{2}} = |a| \left( \sum_i c_i^2 \right)^{\frac{1}{2}} = |a|\rho(x)$.
- If also $y = b_1 v_1 + \cdots + b_k v_k$, then

$$\rho(x+y) = \left( \sum_i (c_i + b_i)^2 \right)^{\frac{1}{2}} \leq \left( \sum_i c_i^2 \right)^{\frac{1}{2}} + \left( \sum_i b_i^2 \right)^{\frac{1}{2}} = \rho(x) + \rho(y).$$

The last inequality is just the triangle inequality for the $L_2$ norm for the vectors $(c_1, \cdots, c_k)$ and $(b_1, \cdots, b_k)$.

In Sect. 2.2.5, we will consider special forms of basis sets in which the norm in equation (2.38) is identically the $L_2$ norm. (This is called Parseval's identity, equation (2.60) on page 41.)

### 2.1.5.5 Equivalence of Norms

There is an equivalence among any two norms over a normed finite-dimensional linear space in the sense that if $\| \cdot \|_a$ and $\| \cdot \|_b$ are norms, then there are positive numbers $r$ and $s$ such that for any $x$ in the space,

$$r\|x\|_b \leq \|x\|_a \leq s\|x\|_b. \tag{2.39}$$

Expressions (2.35) and (2.36) are examples of this general equivalence for three $L_p$ norms.

We can prove inequality (2.39) by using the norm defined in equation (2.38). We need only consider the case $x \neq 0$, because the inequality is obviously true if $x = 0$. Let $\| \cdot \|_a$ be any norm over a given normed linear space and let $\{v_1, \ldots, v_k\}$ be a basis for the space. (Here's where the assumption of a vector space with finite dimensions comes in.) Any $x$ in the space has a representation in terms of the basis, $x = c_1 v_1 + \cdots + c_k v_k$. Then

$$\|x\|_a = \left\| \sum_{i=1}^k c_i v_i \right\|_a \leq \sum_{i=1}^k |c_i| \, \|v_i\|_a.$$

Applying the Cauchy-Schwarz inequality to the two vectors $(c_1, \cdots, c_k)$ and $(\|v_1\|_a, \cdots, \|v_k\|_a)$, we have

$$\sum_{i=1}^k |c_i| \, \|v_i\|_a \leq \left( \sum_{i=1}^k c_i^2 \right)^{\frac{1}{2}} \left( \sum_{i=1}^k \|v_i\|_a^2 \right)^{\frac{1}{2}}.$$

Hence, with $\tilde{s} = (\sum_i \|v_i\|_a^2)^{\frac{1}{2}}$, which must be positive, we have

$$\|x\|_a \le \tilde{s}\rho(x).$$

Now, to establish a lower bound for $\|x\|_a$, let us define a subset $C$ of the linear space consisting of all vectors $(u_1, \ldots, u_k)$ such that $\sum |u_i|^2 = 1$. This set is obviously closed. Next, we define a function $f(\cdot)$ over this closed subset by

$$f(u) = \left\| \sum_{i=1}^k u_i v_i \right\|_a .$$

Because $f$ is continuous, it attains a minimum in this closed subset, say for the vector $u_*$; that is, $f(u_*) \le f(u)$ for any $u$ such that $\sum |u_i|^2 = 1$. Let

$$\tilde{r} = f(u_*),$$

which must be positive, and again consider any $x$ in the normed linear space and express it in terms of the basis, $x = c_1 v_1 + \cdots c_k v_k$. If $x \ne 0$, we have

$$\|x\|_a = \left\| \sum_{i=1}^k c_i v_i \right\|_a$$

$$= \left( \sum_{i=1}^k c_i^2 \right)^{\frac{1}{2}} \left\| \sum_{i=1}^k \left( \frac{c_i}{\left( \sum_{i=1}^k c_i^2 \right)^{\frac{1}{2}}} \right) v_i \right\|_a$$

$$= \rho(x) f(\tilde{c}),$$

where $\tilde{c} = (c_1, \cdots, c_k)/(\sum_{i=1}^k c_i^2)^{1/2}$. Because $\tilde{c}$ is in the set $C$, $f(\tilde{c}) \ge r$; hence, combining this with the inequality above, we have

$$\tilde{r}\rho(x) \le \|x\|_a \le \tilde{s}\rho(x).$$

This expression holds for any norm $\|\cdot\|_a$ and so, after obtaining similar bounds for any other norm $\|\cdot\|_b$ and then combining the inequalities for $\|\cdot\|_a$ and $\|\cdot\|_b$, we have the bounds in the equivalence relation (2.39). (This is an equivalence relation because it is reflexive, symmetric, and transitive. Its transitivity is seen by the same argument that allowed us to go from the inequalities involving $\rho(\cdot)$ to ones involving $\|\cdot\|_b$.)

As we have mentioned, there are some differences in the properties of vector spaces that have an infinite number of dimensions and those with finite dimensions. The equivalence of norms is one of those differences. The argument above fails in the properties of the continuous function $f$. (Recall, however, as we have mentioned, unless we state otherwise, we assume that the vector spaces we discuss have finite dimensions.)

### 2.1.6 Normalized Vectors

The Euclidean norm of a vector corresponds to the length of the vector $x$ in a natural way; that is, it agrees with our intuition regarding "length". Although, as we have seen, this is just one of many vector norms, in most applications it is the most useful one. (I must warn you, however, that occasionally I will carelessly but naturally use "length" to refer to the order of a vector; that is, the number of elements. This usage is common in computer software packages such as R and SAS IML, and software necessarily shapes our vocabulary.)

Dividing a given vector by its length *normalizes* the vector, and the resulting vector with length 1 is said to be *normalized*; thus

$$\tilde{x} = \frac{1}{\|x\|}x \tag{2.40}$$

is a normalized vector. Normalized vectors are sometimes referred to as "unit vectors", although we will generally reserve this term for a special kind of normalized vector (see page 16). A normalized vector is also sometimes referred to as a "normal vector". I use "normalized vector" for a vector such as $\tilde{x}$ in equation (2.40) and use "normal vector" to denote a vector that is orthogonal to a subspace (as on page 34).

### 2.1.6.1 "Inverse" of a Vector

Because the mapping of an inner product takes the elements of one space into a different space (the inner product of vectors takes elements of $\mathbb{R}^n$ into $\mathbb{R}$), the concept of an inverse for the inner product does not make sense in the usual way. First of all, there is no identity with respect to the inner product.

Often in applications, however, inner products are combined with the usual scalar-vector multiplication in the form $\langle x, y\rangle z$; therefore, for given $x$, it may be of interest to determine $y$ such that $\langle x, y\rangle = 1$, the multiplicative identity in $\mathbb{R}$. For $x$ in $\mathbb{R}^n$ such that $x \neq 0$, the additive identity in $\mathbb{R}^n$,

$$y = \frac{1}{\|x\|^2}x = \frac{1}{\|x\|}\tilde{x} \tag{2.41}$$

uniquely satisfies $\langle x, y\rangle = 1$. Such a $y$ is called the *Samelson inverse* of the vector $x$ and is sometimes denoted as $x^{-1}$ or as $[x]^{-1}$. It is also sometimes called the Moore-Penrose vector inverse because it satisfies the four properties of the definition the Moore-Penrose inverse. (See page 128, where, for example, the first property is interpreted both as $\langle x, [x]^{-1}\rangle x$ and as $x\langle [x]^{-1}, x\rangle$.)

The norm in equation (2.41) is obviously the Euclidean norm (because of the way we defined $\tilde{x}$), but the idea of the inverse could also be extended to other norms associated with other inner products.

The Samelson inverse has a nice geometric interpretation: it is the inverse point of $x$ with respect to the unit sphere in $\mathbb{R}^n$. This inverse arises in the vector $\epsilon$-algorithm used in accelerating convergence of vector sequences in numerical computations (see Wynn 1962).

### 2.1.7 Metrics and Distances

It is often useful to consider how far apart two objects are; that is, the "distance" between them. A reasonable distance measure would have to satisfy certain requirements, such as being a nonnegative real number.

A function $\Delta$ that maps any two objects in a set $S$ to $\mathbb{R}$ is called a *metric* on $S$ if, for all $x$, $y$, and $z$ in $S$, it satisfies the following three conditions:

1.  $\Delta(x, y) > 0$ if $x \neq y$ and $\Delta(x, y) = 0$ if $x = y$;
2.  $\Delta(x, y) = \Delta(y, x)$;
3.  $\Delta(x, y) \leq \Delta(x, z) + \Delta(z, y)$.

These conditions correspond in an intuitive manner to the properties we expect of a distance between objects.

A vector space together with a metric defined on it is called a *metric space.* A normed vector space is a metric space because the norm can induce a metric. In the following, we may speak almost interchangeably of an inner product space, a normed space, or a metric space, but we must recognize that none is a special case of another. (Recall that a normed space whose norm is the $L_1$ norm is not equivalent to an inner product space, for example.)

### 2.1.7.1 Metrics Induced by Norms

If subtraction and a norm are defined for the elements of $S$, the most common way of forming a metric is by using the norm. If $\|\cdot\|$ is a norm, we can verify that

$$\Delta(x, y) = \|x - y\| \tag{2.42}$$

is a metric by using the properties of a norm to establish the three properties of a metric above (Exercise 2.12).

The norm in equation (2.42) may, of course, be induced by an inner product.

The general inner products, norms, and metrics defined above are relevant in a wide range of applications. The sets on which they are defined can consist of various types of objects. In the context of real vectors, the most common inner product is the dot product; the most common norm is the Euclidean norm that arises from the dot product; and the most common metric is the one defined by the Euclidean norm, called the Euclidean distance.

### 2.1.7.2 Convergence of Sequences of Vectors

A sequence of real numbers $a_1, a_2, \ldots$ is said to converge to a finite number $a$ if for any given $\epsilon > 0$ there is an integer $M$ such that, for $k > M$, $|a_k - a| < \epsilon$, and we write $\lim_{k \to \infty} a_k = a$, or we write $a_k \to a$ as $k \to \infty$.

We define convergence of a sequence of vectors in a normed vector space in terms of the convergence of a sequence of their norms, which is a sequence of

real numbers. We say that a sequence of vectors $x_1, x_2, \ldots$ (of the same order) converges to the vector $x$ with respect to the norm $\|\cdot\|$ if the sequence of real numbers $\|x_1 - x\|, \|x_2 - x\|, \ldots$ converges to 0. Because of the bounds (2.39), the choice of the norm is irrelevant (for finite dimensional vector spaces), and so convergence of a sequence of vectors is well-defined without reference to a specific norm. (This is one reason that equivalence of norms is an important property.)

A sequence of vectors $x_1, x_2, \ldots$ in the metric space $\mathcal{V}$ that come arbitrarily close to one another (as measured by the given metric) is called a *Cauchy sequence*. (In a Cauchy sequence $x_1, x_2, \ldots$, for any $\epsilon > 0$ there is a number $N$ such that for $i, j > N$, $\Delta(x_i, x_j) < \epsilon$.) Intuitively, such a sequence should converge to some fixed vector in $\mathcal{V}$, but this is not necessarily the case. A metric space in which every Cauchy sequence converges to an element in the space is said to be a *complete metric space* or just a *complete space*. The space $\mathbb{R}^n$ (with any norm) is complete.

A complete normed space is called a *Banach space*, and a complete inner product space is called a *Hilbert space*. It is clear that a Hilbert space is a Banach space (because its inner product induces a norm). As we have indicated, a space with a norm induced by an inner product, such as a Hilbert space, has an interesting structure. Most of the vector spaces encountered in statistical applications are Hilbert spaces. The space $\mathbb{R}^n$ with the $L_2$ norm is a Hilbert space.

### 2.1.8 Orthogonal Vectors and Orthogonal Vector Spaces

Two vectors $v_1$ and $v_2$ such that

$$\langle v_1, v_2 \rangle = 0 \tag{2.43}$$

are said to be *orthogonal*, and this condition is denoted by $v_1 \perp v_2$. (Sometimes we exclude the zero vector from this definition, but it is not important to do so.) Normalized vectors that are all orthogonal to each other are called *orthonormal* vectors.

**An Aside: Complex Vectors**

If the elements of the vectors are from the field of complex numbers, orthogonality and normality are also defined as above; however, the inner product in the definition (2.43) must be as defined in equation (2.23), and the expression $x^T y$ in equation (2.25) is not equivalent to the inner product. We will use a different notation in this case: $x^H y$. The relationship between the two notations is

$$x^H y = \bar{x}^T y.$$

With this interpretation of the inner product, all of the statements below about orthogonality hold for complex numbers as well as for real numbers.

A set of nonzero vectors that are mutually orthogonal are necessarily linearly independent. To see this, we show it for any two orthogonal vectors and then indicate the pattern that extends to three or more vectors. First, suppose $v_1$ and $v_2$ are nonzero and are orthogonal; that is, $\langle v_1, v_2 \rangle = 0$. We see immediately that if there is a scalar $a$ such that $v_1 = av_2$, then $a$ must be nonzero and we have a contradiction because $\langle v_1, v_2 \rangle = a\langle v_2, v_2 \rangle \neq 0$. Hence, we conclude $v_1$ and $v_2$ are independent (there is no $a$ such that $v_1 = av_2$). For three mutually orthogonal vectors, $v_1$, $v_2$, and $v_3$, we consider $v_1 = av_2 + bv_3$ for $a$ or $b$ nonzero, and arrive at the same contradiction.

Two vector spaces $\mathcal{V}_1$ and $\mathcal{V}_2$ are said to be *orthogonal*, written $\mathcal{V}_1 \perp \mathcal{V}_2$, if each vector in one is orthogonal to every vector in the other. If $\mathcal{V}_1 \perp \mathcal{V}_2$ and $\mathcal{V}_1 \oplus \mathcal{V}_2 = \mathbb{R}^n$, then $\mathcal{V}_2$ is called the *orthogonal complement* of $\mathcal{V}_1$, and this is written as $\mathcal{V}_2 = \mathcal{V}_1^\perp$. More generally, if $\mathcal{V}_1 \perp \mathcal{V}_2$ and $\mathcal{V}_1 \oplus \mathcal{V}_2 = \mathcal{V}$, then $\mathcal{V}_2$ is called the orthogonal complement of $\mathcal{V}_1$ with respect to $\mathcal{V}$. This is obviously a symmetric relationship; if $\mathcal{V}_2$ is the orthogonal complement of $\mathcal{V}_1$, then $\mathcal{V}_1$ is the orthogonal complement of $\mathcal{V}_2$.

A vector that is orthogonal to all vectors in a given vector space is said to be orthogonal to that space or *normal* to that space. Such a vector is called a *normal vector* to that space.

If $B_1$ is a basis set for $\mathcal{V}_1$, $B_2$ is a basis set for $\mathcal{V}_2$, and $\mathcal{V}_2$ is the orthogonal complement of $\mathcal{V}_1$ with respect to $\mathcal{V}$, then $B_1 \cup B_2$ is a basis set for $\mathcal{V}$. It is a basis set because since $\mathcal{V}_1$ and $\mathcal{V}_2$ are orthogonal, it must be the case that $B_1 \cap B_2 = \emptyset$. (See the properties listed on page 22.)

If $\mathcal{V}_1 \subset \mathcal{V}$, $\mathcal{V}_2 \subset \mathcal{V}$, $\mathcal{V}_1 \perp \mathcal{V}_2$, and $\dim(\mathcal{V}_1) + \dim(\mathcal{V}_2) = \dim(\mathcal{V})$, then

$$\mathcal{V}_1 \oplus \mathcal{V}_2 = \mathcal{V}; \tag{2.44}$$

that is, $\mathcal{V}_2$ is the orthogonal complement of $\mathcal{V}_1$. We see this by first letting $B_1$ and $B_2$ be bases for $\mathcal{V}_1$ and $\mathcal{V}_2$. Now $\mathcal{V}_1 \perp \mathcal{V}_2$ implies that $B_1 \cap B_2 = \emptyset$ and $\dim(\mathcal{V}_1) + \dim(\mathcal{V}_2) = \dim(\mathcal{V})$ implies $\#(B_1) + \#(B_2) = \#(B)$, for any basis set $B$ for $\mathcal{V}$; hence, $B_1 \cup B_2$ is a basis set for $\mathcal{V}$.

The intersection of two orthogonal vector spaces consists only of the zero vector (Exercise 2.14).

A set of linearly independent vectors can be mapped to a set of mutually orthogonal (and orthonormal) vectors by means of the Gram-Schmidt transformations (see equation (2.56) below).

### 2.1.9 The "One Vector"

The vector with all elements equal to 1 that we mentioned previously is useful in various vector operations. We call this the "one vector" and denote it by 1 or by $1_n$. The one vector can be used in the representation of the sum of the elements in a vector:

$$1^{\mathrm{T}}x = \sum x_i. \tag{2.45}$$

The one vector is also called the "summing vector".

### 2.1.9.1 The Mean and the Mean Vector

Because the elements of $x$ are real, they can be summed; however, in applications it may or may not make sense to add the elements in a vector, depending on what is represented by those elements. If the elements have some kind of essential commonality, it may make sense to compute their sum as well as their *arithmetic mean*, which for the $n$-vector $x$ is denoted by $\bar{x}$ and defined by

$$\bar{x} = 1_n^\mathrm{T} x/n. \tag{2.46}$$

We also refer to the arithmetic mean as just the "mean" because it is the most commonly used mean.

It is often useful to think of the mean as an $n$-vector all of whose elements are $\bar{x}$. The symbol $\bar{x}$ is also used to denote this vector; hence, we have

$$\bar{x} = \bar{x} 1_n, \tag{2.47}$$

in which $\bar{x}$ on the left-hand side is a vector and $\bar{x}$ on the right-hand side is a scalar. We also have, for the two different objects,

$$\|\bar{x}\|^2 = n\bar{x}^2. \tag{2.48}$$

The meaning, whether a scalar or a vector, is usually clear from the context. In any event, an expression such as $x - \bar{x}$ is unambiguous; the addition (subtraction) has the same meaning whether $\bar{x}$ is interpreted as a vector or a scalar. (In some mathematical treatments of vectors, addition of a scalar to a vector is not defined, but here we are following the conventions of modern computer languages.)

## 2.2 Cartesian Coordinates and Geometrical Properties of Vectors

Points in a Cartesian geometry can be identified with vectors, and several definitions and properties of vectors can be motivated by this geometric interpretation. In this interpretation, vectors are directed line segments with a common origin. The geometrical properties can be seen most easily in terms of a Cartesian coordinate system, but the properties of vectors defined in terms of a Cartesian geometry have analogues in Euclidean geometry without a coordinate system. In such a system, only length and direction are defined, and two vectors are considered to be the same vector if they have the same length and direction. Generally, we will not assume that there is a "location" or "position" associated with a vector.

### 2.2.1 Cartesian Geometry

A Cartesian coordinate system in $d$ dimensions is defined by $d$ unit vectors, $e_i$ in equation (2.6), each with $d$ elements. A unit vector is also called a *principal axis* of the coordinate system. The set of unit vectors is orthonormal. (There is an implied number of elements of a unit vector that is inferred from the context. Also parenthetically, we remark that the phrase "unit vector" is sometimes used to refer to a vector the sum of whose squared elements is 1, that is, whose length, in the Euclidean distance sense, is 1. As we mentioned above, we refer to this latter type of vector as a "normalized vector".)

The sum of all of the unit vectors is the one vector:

$$\sum_{i=1}^{d} e_i = 1_d. \tag{2.49}$$

A point $x$ with Cartesian coordinates $(x_1, \ldots, x_d)$ is associated with a vector from the origin to the point, that is, the vector $(x_1, \ldots, x_d)$. The vector can be written as the linear combination

$$x = x_1 e_1 + \ldots + x_d e_d \tag{2.50}$$

or, equivalently, as

$$x = \langle x, e_1 \rangle e_1 + \ldots + \langle x, e_d \rangle e_d.$$

(This is a Fourier expansion, equation (2.58) below.)

### 2.2.2 Projections

The *projection* of the vector $y$ onto the nonnull vector $x$ is the vector

$$\hat{y} = \frac{\langle x, y \rangle}{\|x\|^2} x. \tag{2.51}$$

This definition is consistent with a geometrical interpretation of vectors as directed line segments with a common origin. The projection of $y$ onto $x$ is the inner product of the normalized $x$ and $y$ times the normalized $x$; that is, $\langle \tilde{x}, y \rangle \tilde{x}$, where $\tilde{x} = x/\|x\|$. Notice that the order of $y$ and $x$ is the same.

An important property of a projection is that when it is subtracted from the vector that was projected, the resulting vector, called the "residual", is orthogonal to the projection; that is, if

$$\begin{aligned} r &= y - \frac{\langle x, y \rangle}{\|x\|^2} x \\ &= y - \hat{y} \end{aligned} \tag{2.52}$$

then $r$ and $\hat{y}$ are orthogonal, as we can easily see by taking their inner product (see Fig. 2.2.2). Notice also that the Pythagorean relationship holds:
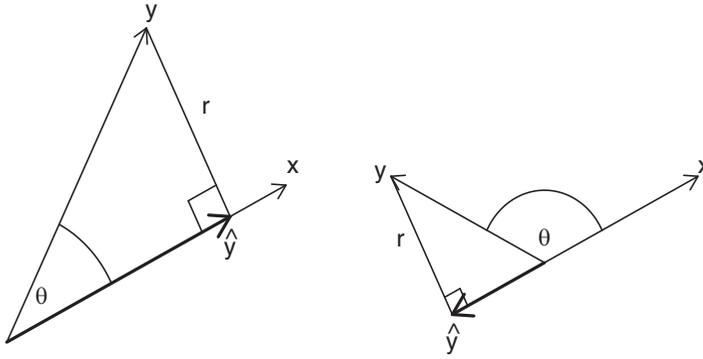
**Figure 2.1.** Projections and angles

$$\|y\|^2 = \|\hat{y}\|^2 + \|r\|^2. \tag{2.53}$$

As we mentioned on page 35, the mean $\bar{y}$ can be interpreted either as a scalar or as a vector all of whose elements are $\bar{y}$. As a vector, it is the projection of $y$ onto the one vector $1_n$,

$$\frac{\langle 1_n, y \rangle}{\|1_n\|^2} 1_n = \frac{1_n^{\mathrm{T}} y}{n} 1_n$$
$$= \bar{y}\, 1_n,$$

from equations (2.46) and (2.51).

We will consider more general projections (that is, projections onto planes or other subspaces) on page 352, and on page 409 we will view linear regression fitting as a projection onto the space spanned by the independent variables.

### 2.2.3 Angles Between Vectors

The *angle* between the nonnull vectors $x$ and $y$ is determined by its cosine, which we can compute from the length of the projection of one vector onto the other. Hence, denoting the angle between the nonnull vectors $x$ and $y$ as angle$(x, y)$, we define

$$\text{angle}(x, y) = \cos^{-1}\left(\frac{\langle x, y \rangle}{\|x\|\|y\|}\right), \tag{2.54}$$

with $\cos^{-1}(\cdot)$ being taken in the interval $[0, \pi]$. The cosine is $\pm\|\hat{y}\|/\|y\|$, with the sign chosen appropriately; see Fig. 2.2.2. Because of this choice of $\cos^{-1}(\cdot)$, we have that angle$(y, x) = $ angle$(x, y)$—but see Exercise 2.19e on page 54.

The word "orthogonal" is appropriately defined by equation (2.43) on page 33 because orthogonality in that sense is equivalent to the corresponding geometric property. (The cosine is 0.)

Notice that the angle between two vectors is invariant to scaling of the vectors; that is, for any positive scalar $a$, $\text{angle}(ax, y) = \text{angle}(x, y)$.

A given vector can be defined in terms of its length and the angles $\theta_i$ that it makes with the unit vectors. The cosines of these angles are just the scaled coordinates of the vector:

$$\cos(\theta_i) = \frac{\langle x, e_i \rangle}{\|x\| \|e_i\|}$$
$$= \frac{1}{\|x\|} \, x_i. \tag{2.55}$$

These quantities are called the *direction cosines* of the vector.

Although geometrical intuition often helps us in understanding properties of vectors, sometimes it may lead us astray in high dimensions. Consider the direction cosines of an arbitrary vector in a vector space with large dimensions. If the elements of the arbitrary vector are nearly equal (that is, if the vector is a diagonal through an orthant of the coordinate system), the direction cosine goes to 0 as the dimension increases. In high dimensions, any two vectors are "almost orthogonal" to each other; see Exercise 2.16.

The geometric property of the angle between vectors has important implications for certain operations both because it may indicate that rounding in computations will have deleterious effects and because it may indicate a deficiency in the understanding of the application.

We will consider more general projections and angles between vectors and other subspaces on page 359. In Sect. 5.3.3, we will consider rotations of vectors onto other vectors or subspaces. Rotations are similar to projections, except that the length of the vector being rotated is preserved.

### 2.2.4 Orthogonalization Transformations: Gram-Schmidt

Given $m$ nonnull, linearly independent vectors, $x_1, \ldots, x_m$, it is easy to form $m$ orthonormal vectors, $\tilde{x}_1, \ldots, \tilde{x}_m$, that span the same space. A simple way to do this is sequentially. First normalize $x_1$ and call this $\tilde{x}_1$. Next, project $x_2$ onto $\tilde{x}_1$ and subtract this projection from $x_2$. The result is orthogonal to $\tilde{x}_1$; hence, normalize this and call it $\tilde{x}_2$. These first two steps, which are illustrated in Fig. 2.2.4, are

$$\tilde{x}_1 = \frac{1}{\|x_1\|} \, x_1,$$

$$\tilde{x}_2 = \frac{1}{\|x_2 - \langle \tilde{x}_1, x_2 \rangle \tilde{x}_1\|} \, (x_2 - \langle \tilde{x}_1, x_2 \rangle \tilde{x}_1). \tag{2.56}$$

These are called *Gram-Schmidt transformations*.

The Gram-Schmidt transformations have a close relationship to least squares fitting of overdetermined systems of linear equations and to least

squares fitting of linear regression models. For example, if equation (6.33) on page 290 is merely the system of equations in one unknown $x_1 b = x_2 - r$, and it is approximated by least squares, then the "residual vector" $r$ is $\tilde{x}_2$ above, and of course it has the orthogonality property shown in equation (6.37) for that problem.

The Gram-Schmidt transformations can be continued with all of the vectors in the linearly independent set. There are two straightforward ways equations (2.56) can be extended. One method generalizes the second equation in an obvious way:

for $k = 2, 3 \ldots$,

$$\tilde{x}_k = \left( x_k - \sum_{i=1}^{k-1} \langle \tilde{x}_i, x_k \rangle \tilde{x}_i \right) \Big/ \left\| x_k - \sum_{i=1}^{k-1} \langle \tilde{x}_i, x_k \rangle \tilde{x}_i \right\|. \tag{2.57}$$

In this method, at the $k^{\text{th}}$ step, we orthogonalize the $k^{\text{th}}$ vector by computing its residual with respect to the plane formed by all the previous $k - 1$ orthonormal vectors.

Another way of extending the transformation of equations (2.56) is, at the $k^{\text{th}}$ step, to compute the residuals of all remaining vectors with respect just to the $k^{\text{th}}$ normalized vector. If the initial set of vectors are linearly independent, the residuals at any stage will be nonzero. (This is fairly obvious, but you are asked to show it in Exercise 2.17.) We describe this method explicitly in Algorithm 2.1.

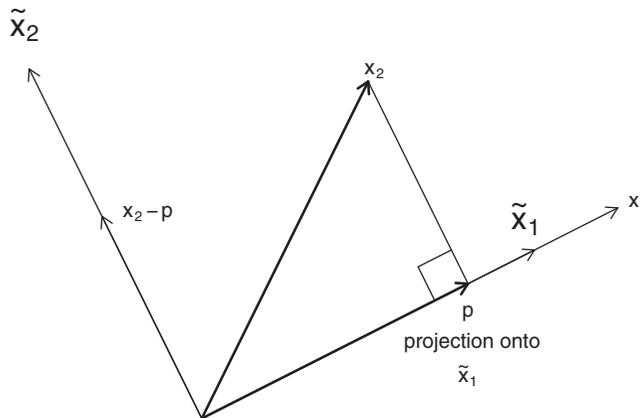**Algorithm 2.1 Gram-Schmidt orthonormalization of a set of linearly independent vectors, $x_1, \ldots, x_m$**



**Figure 2.2.** Orthogonalization of $x_1$ and $x_2$

0. For $k = 1, \ldots, m$,
   {
   set $\tilde{x}_k = x_k$.
   }
1. Ensure that $\tilde{x}_1 \neq 0$;
   set $\tilde{x}_1 = \tilde{x}_1 / \|\tilde{x}_1\|$.
2. If $m > 1$, for $k = 2, \ldots, m$,
   {
      for $j = k, \ldots, m$,
      {
         set $\tilde{x}_j = \tilde{x}_j - \langle \tilde{x}_{k-1}, \tilde{x}_j \rangle \tilde{x}_{k-1}$.
      }
      ensure that $\tilde{x}_k \neq 0$;
      set $\tilde{x}_k = \tilde{x}_k / \|\tilde{x}_k\|$.
   }
∎

Although the method indicated in equation (2.57) is mathematically equivalent to this method, the use of Algorithm 2.1 is to be preferred for computations because it is less subject to rounding errors. (This may not be immediately obvious, although a simple numerical example can illustrate the fact—see Exercise 11.1c on page 537. We will not digress here to consider this further, but the difference in the two methods has to do with the relative magnitudes of the quantities in the subtraction. The method of Algorithm 2.1 is sometimes called the "modified Gram-Schmidt method", although I call it the "Gram-Schmidt method". I will discuss this method again in Sect. 11.2.1.3.) This is an instance of a principle that we will encounter repeatedly: *the form of a mathematical expression and the way the expression should be evaluated in actual practice may be quite different.*

These orthogonalizing transformations result in a set of orthogonal vectors that span the same space as the original set. They are not unique; if the order in which the vectors are processed is changed, a different set of orthogonal vectors will result.

Orthogonal vectors are useful for many reasons: perhaps to improve the stability of computations; or in data analysis to capture the variability most efficiently; or for dimension reduction as in principal components analysis (see Sect. 9.4 beginning on page 424); or in order to form more meaningful quantities as in a vegetative index in remote sensing. We will discuss various specific orthogonalizing transformations later.

### 2.2.5 Orthonormal Basis Sets

A basis for a vector space is often chosen to be an orthonormal set because it is easy to work with the vectors in such a set.

If $u_1, \ldots, u_n$ is an orthonormal basis set for a space, then a vector $x$ in that space can be expressed as

$$x = c_1 u_1 + \cdots + c_n u_n, \tag{2.58}$$

and because of orthonormality, we have

$$c_i = \langle x, \, u_i \rangle. \tag{2.59}$$

(We see this by taking the inner product of both sides with $u_i$.) A representation of a vector as a linear combination of orthonormal basis vectors, as in equation (2.58), is called a *Fourier expansion*, and the $c_i$ are called *Fourier coefficients*.

By taking the inner product of each side of equation (2.58) with itself, we have *Parseval's identity*:

$$\|x\|^2 = \sum c_i^2. \tag{2.60}$$

This shows that the $L_2$ norm is the same as the norm in equation (2.38) (on page 29) for the case of an orthogonal basis.

Although the Fourier expansion is not unique because a different orthogonal basis set could be chosen, Parseval's identity removes some of the arbitrariness in the choice; no matter what basis is used, the sum of the squares of the Fourier coefficients is equal to the square of the norm that arises from the inner product. ("The" inner product means the inner product used in defining the orthogonality.)

Another useful expression of Parseval's identity in the Fourier expansion is

$$\left\| x - \sum_{i=1}^{k} c_i u_i \right\|^2 = \langle x, \, x \rangle - \sum_{i=1}^{k} c_i^2 \tag{2.61}$$

(because the term on the left-hand side is 0).

The expansion (2.58) is a special case of a very useful expansion in an orthogonal basis set. In the finite-dimensional vector spaces we consider here, the series is finite. In function spaces, the series is generally infinite, and so issues of convergence are important. For different types of functions, different orthogonal basis sets may be appropriate. Polynomials are often used, and there are some standard sets of orthogonal polynomials, such as Jacobi, Hermite, and so on. For periodic functions especially, orthogonal trigonometric functions are useful.

### 2.2.6 Approximation of Vectors

In high-dimensional vector spaces, it is often useful to approximate a given vector in terms of vectors from a lower dimensional space. Suppose, for example, that $\mathcal{V} \subseteq \mathbb{R}^n$ is a vector space of dimension $k$ (necessarily, $k \leq n$) and $x$ is a given $n$-vector, not necessarily in $\mathcal{V}$. We wish to determine a vector $\tilde{x}$ in $\mathcal{V}$ that approximates $x$. Of course if $\mathcal{V} = \mathbb{R}^n$, then $x \in \mathcal{V}$, and so the problem is not very interesting. The interesting case is when $\mathcal{V} \subset \mathbb{R}^n$.

### 2.2.6.1 Optimality of the Fourier Coefficients

The first question, of course, is what constitutes a "good" approximation. One obvious criterion would be based on a norm of the difference of the given vector and the approximating vector. So now, choosing the norm as the Euclidean norm, we may pose the problem as one of finding $\tilde{x} \in \mathcal{V}$ such that

$$\|x - \tilde{x}\| \le \|x - v\| \quad \forall\, v \in \mathcal{V}. \tag{2.62}$$

This difference is a *truncation error*.

Let $u_1, \ldots, u_k$ be an orthonormal basis set for $\mathcal{V}$, and let

$$\tilde{x} = c_1 u_1 + \cdots + c_k u_k, \tag{2.63}$$

where the $c_i$ are the Fourier coefficients of $x$, $\langle x,\, u_i \rangle$.

Now let $v = a_1 u_1 + \cdots + a_k u_k$ be any other vector in $\mathcal{V}$, and consider

$$
\begin{aligned}
\|x - v\|^2 &= \left\| x - \sum_{i=1}^{k} a_i u_i \right\|^2 \\
&= \left\langle x - \sum_{i=1}^{k} a_i u_i,\ x - \sum_{i=1}^{k} a_i u_i \right\rangle \\
&= \langle x,\, x \rangle - 2 \sum_{i=1}^{k} a_i \langle x,\, u_i \rangle + \sum_{i=1}^{k} a_i^2 \\
&= \langle x,\, x \rangle - 2 \sum_{i=1}^{k} a_i c_i + \sum_{i=1}^{k} a_i^2 + \sum_{i=1}^{k} c_i^2 - \sum_{i=1}^{k} c_i^2 \\
&= \langle x,\, x \rangle + \sum_{i=1}^{k} (a_i - c_i)^2 - \sum_{i=1}^{k} c_i^2 \\
&= \left\| x - \sum_{i=1}^{k} c_i u_i \right\|^2 + \sum_{i=1}^{k} (a_i - c_i)^2 \\
&\ge \left\| x - \sum_{i=1}^{k} c_i u_i \right\|^2. \tag{2.64}
\end{aligned}
$$

Therefore we have $\|x - \tilde{x}\| \le \|x - v\|$, and so $\tilde{x}$ formed by the Fourier coefficients is the best approximation of $x$ with respect to the Euclidean norm in the $k$-dimensional vector space $\mathcal{V}$. (For some other norm, this may not be the case.)

### 2.2.6.2 Choice of the Best Basis Subset

Now, posing the problem another way, we may seek the best $k$-dimensional subspace of $\mathrm{I\!R}^n$ from which to choose an approximating vector. This question

is not well-posed (because the one-dimensional vector space determined by $x$ is the solution), but we can pose a related interesting question: suppose we have a Fourier expansion of $x$ in terms of a set of $n$ orthogonal basis vectors, $u_1, \ldots, u_n$, and we want to choose the "best" $k$ basis vectors from this set and use them to form an approximation of $x$. (This restriction of the problem is equivalent to choosing a coordinate system.) We see the solution immediately from inequality (2.64): we choose the $k$ $u_i$s corresponding to the $k$ largest $c_i$s in absolute value, and we take

$$\tilde{x} = c_{i_1} u_{i_1} + \cdots + c_{i_k} u_{i_k}, \qquad (2.65)$$

where $\min(\{|c_{i_j}| \; : \; j = 1, \ldots, k\}) \geq \max(\{|c_{i_j}| \; : \; j = k+1, \ldots, n\})$.

### 2.2.7 Flats, Affine Spaces, and Hyperplanes

Given an $n$-dimensional vector space of order $n$, $\mathbb{R}^n$ for example, consider a system of $m$ linear equations in the $n$-vector variable $x$,

$$c_1^\mathrm{T} x = b_1$$
$$\vdots \quad \vdots$$
$$c_m^\mathrm{T} x = b_m,$$

where $c_1, \ldots, c_m$ are linearly independent $n$-vectors (and hence $m \leq n$). The set of points defined by these linear equations is called a *flat*. Although it is not necessarily a vector space, a flat is also called an *affine space*. An intersection of two flats is a flat.

If the equations are *homogeneous* (that is, if $b_1 = \cdots = b_m = 0$), then the point $(0, \ldots, 0)$ is included, and the flat is an $(n - m)$-dimensional subspace (also a vector space, of course). Stating this another way, a flat through the origin is a vector space, but other flats are not vector spaces.

If $m = 1$, the flat is called a *hyperplane*. A hyperplane through the origin is an $(n - 1)$-dimensional vector space.

If $m = n-1$, the flat is a line. A line through the origin is a one-dimensional vector space.

### 2.2.8 Cones

A set of vectors that contains all nonnegative scalar multiples of any vector in the set is called a *cone*. A cone always contains the zero vector. (Some authors define a cone as a set that contains all positive scalar multiples, and in that case, the zero vector may not be included.) If a set of vectors contains all scalar multiples of any vector in the set, it is called a *double cone*.

Geometrically, a cone is just a set, possibly a finite set, of lines or half-lines. (A double cone is a set of lines.) In general, a cone may not be very interesting, but certain special cones are of considerable interest.

Given two (double) cones over the same vector space, both their union and their intersection are (double) cones. A (double) cone is in general not a vector space.

### 2.2.8.1 Convex Cones

A set of vectors $C$ in a vector space $\mathcal{V}$ is a *convex cone* if, for all $v_1, v_2 \in C$ and all nonnegative real numbers $a, b \geq 0$, $av_1 + bv_2 \in C$. Such a cone is called a homogeneous convex cone by some authors. (An equivalent definition requires that the set $C$ be a cone, and then, more in keeping with the definition of convexity, includes the requirement $a + b = 1$ along with $a, b \geq 0$ in the definition of a convex cone.)

If $C$ is a convex cone and if $v \in C$ implies $-v \in C$, then $C$ is called a *double convex cone*.

A (double) convex cone is in general not a vector space because, for example, $v_1 + v_2$ may not be in $C$.

It is clear that a (double) convex cone is a (double) cone; in fact, a convex cone is the most important type of cone. A convex cone corresponds to a solid geometric object with a single finite vertex.

An important convex cone in an $n$-dimensional vector space with a Cartesian coordinate system is the positive orthant together with the zero vector. This convex cone is not closed, in the sense that it does not contain some limits. The closure of the positive orthant (that is, the nonnegative orthant) is also a convex cone.

A generating set or spanning set of a cone $C$ is a set of vectors $S = \{v_i\}$ such that for any vector $v$ in $C$ there exists scalars $a_i \geq 0$ so that $v = \sum a_i v_i$. If, in addition, for any scalars $b_i \geq 0$ with $\sum b_i v_i = 0$, it is necessary that $b_i = 0$ for all $i$, then $S$ is a basis set for the cone. The concept of a generating set is of course more interesting in the case of a convex cone.

If a generating set of a convex cone has a finite number of elements, the cone is a *polyhedron*. For the common geometric object in three dimensions with elliptical contours and which is the basis for "conic sections", any generating set has an uncountable number of elements. Cones of this type are sometimes called "Lorentz cones".

It is easy to see from the definition that if $C_1$ and $C_2$ are convex cones over the same vector space, then $C_1 \cap C_2$ is a convex cone. On the other hand, $C_1 \cup C_2$ is not necessarily a convex cone. Of course the union of two cones, as we have seen, is a cone.

### 2.2.8.2 Dual Cones

Given a set of vectors $S$ in a given vector space (in cases of interest, $S$ is usually a cone, but not necessarily), the *dual cone* of $S$, denoted $C^*(S)$, is defined as

$$C^*(S) = \{v^* \text{ s.t. } \langle v^*, v \rangle \geq 0 \text{ for all } v \in S\}.$$

See Fig. 2.2.8.2 in which $S = \{v_1, v_2, v_3\}$. Clearly, the dual cone is a cone, and also $S \subseteq C^*(S)$.

If, as in the most common cases, the underlying set of vectors is a cone, say $C$, we generally drop the reference to an underlying set of vectors, and just denote the dual cone of $C$ as $C^*$.

Geometrically, the dual cone $C^*$ of $S$ consists of all vectors that form nonobtuse angles with the vectors in $S$.

Notice that for a given set of vectors $S$, if $-S$ represents the set of vectors $v$ such that $-v \in S$, then $C^*(-S) = -(C^*(S))$, or just $-C^*(S)$, which represents the set of vectors $v^*$ such that $-v^* \in C^*(S)$.

Further, from the definition, we note that if $S_1$ and $S_2$ are sets of vectors in the same vector space such that $S_1 \subseteq S_2$ then $C^*(S_1) \subseteq C^*(S_2)$, or $C_1^* \subseteq C_2^*$.

A dual cone $C^*(S)$ is a closed convex cone. We see this by considering any $v_1^*, v_2^* \in C^*$ and real numbers $a, b \geq 0$. For any $v \in S$, it must be the case that $\langle v_1^*, v \rangle \geq 0$ and $\langle v_2^*, v \rangle \geq 0$; hence, $\langle (av_1^* + bv_2^*), v \rangle \geq 0$, that is, $av_1^* + bv_2^* \in C^*$, so $C^*$ is a convex cone. The closure property comes from the $\geq$ condition in the definition.
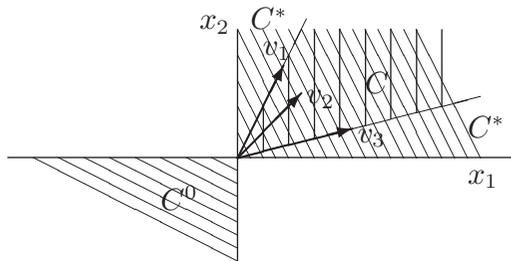


**Figure 2.3.** A set of vectors $\{v_1, v_2, v_3\}$, and the corresponding convex cone $C$, the dual cone $C^*$, and the polar cone $C^0$

### 2.2.8.3 Polar Cones

Given a set of vectors $S$ in a given vector space (in cases of interest, $S$ is usually a cone, but not necessarily), the *polar cone* of $S$, denoted $C^0(S)$, is defined as

$$C^0(S) = \{v^0 \text{ s.t. } \langle v^0, v \rangle \leq 0 \text{ for all } v \in S\}.$$

See Fig. 2.2.8.2.

We generally drop the reference to an underlying set of vectors, and just denote the dual cone of the set $C$ as $C^0$.

From the definition, we note that if $S_1$ and $S_2$ are sets of vectors in the same vector space such that $S_1 \subseteq S_2$ then $C^0(S_1) \subseteq C^0(S_2)$, or $C_1^0 \subseteq C_2^0$.

The polar cone and the dual cone of a double cone are clearly the same.

From the definitions, it is clear in any case that the polar cone $C^0$ can be formed by multiplying all of the vectors in the corresponding dual cone $C^*$ by $-1$, and so $C^0 = -C^*$.

The relationships of the polar cone to the dual cone and the properties we have established for a dual cone immediately imply that a polar cone is also a convex cone.

Another interesting property of polar cones is that for any set of vectors $S$ in a given vector space, $S \subseteq (C^0)^0$. We generally write $(C^0)^0$ as just $C^{00}$. (The precise notation of course is $C^0(C^0(S))$.) We see this by first taking any $v \in S$. Therefore, if $v^0 \in C^0$ then $\langle v, v^0 \rangle \leq 0$, which implies $v \in (C^0)^0$, because
$$C^{00} = \{v \text{ s.t. } \langle v, v^0 \rangle \leq 0 \text{ for all } v^0 \in C^0\}.$$

### 2.2.8.4 Additional Properties

As noted above, a cone is a very loose and general structure. In my definition, the vectors in the set do not even need to be in the same vector space. A convex cone, on the other hand is a useful structure, and the vectors in a convex cone must be in the same vector space.

Most cones of interest, in particular, dual cones and polar cones are not necessarily vector spaces.

Although the definitions of dual cones and polar cones can apply to any set of vectors, they are of the most interest in the case in which the underlying set of vectors is a cone in the nonnegative orthant of a Cartesian coordinate system on $\mathbb{R}^n$ (the set of $n$-vectors all of whose elements are nonnegative). In that case, the dual cone is just the full nonnegative orthant, and the polar cone is just the nonpositive orthant (the set of all vectors all of whose elements are nonpositive).

The whole nonnegative orthant itself is a convex cone, and as we have seen for any convex cone within that orthant, the dual cone is the full nonnegative orthant.

Because the nonnegative orthant is its own dual, and hence is said to be "self-dual". (There is an extension of the property of self-duality that we will not discuss here.)

Convex cones occur in many optimization problems. The feasible region in a linear programming problem is generally a convex polyhedral cone, for example.

### 2.2.9 Cross Products in $\mathbb{R}^3$

The vector space $\mathbb{R}^3$ is especially interesting because it serves as a useful model of the real world, and many physical processes can be represented as vectors in it.

For the special case of the vector space $\mathbb{R}^3$, another useful vector product is the cross product, which is a mapping from $\mathbb{R}^3 \times \mathbb{R}^3$ to $\mathbb{R}^3$. Before proceeding,

we note an overloading of the term "cross product" and of the symbol "×" used to denote it. If $A$ and $B$ are sets, the *set cross product* or the *set Cartesian product* of $A$ and $B$ is the set consisting of all doubletons $(a, b)$ where $a$ ranges over all elements of $A$, and $b$ ranges independently over all elements of $B$. Thus, $\mathbb{R}^3 \times \mathbb{R}^3$ is the set of all pairs of all real 3-vectors.

The *vector cross product* of the 3-vectors

$$x = (x_1, x_2, x_3)$$

and

$$y = (y_1, y_2, y_3),$$

written $x \times y$, is defined as

$$x \times y = (x_2 y_3 - x_3 y_2, \ x_3 y_1 - x_1 y_3, \ x_1 y_2 - x_2 y_1). \tag{2.66}$$

(We also use the term "cross products" in a different way to refer to another type of product formed by several inner products; see page 359.) The cross product has the following properties, which are immediately obvious from the definition:

1. Self-nilpotency:
   $x \times x = 0$, for all $x$.
2. Anti-commutativity:
   $x \times y = -y \times x$.
3. Factoring of scalar multiplication;
   $ax \times y = a(x \times y)$ for real $a$.
4. Relation of vector addition to addition of cross products:
   $(x + y) \times z = (x \times z) + (y \times z)$.

The cross product has the important property (sometimes taken as the definition),

$$x \times y = \|x\|\|y\| \sin(\text{angle}(y, x))e, \tag{2.67}$$

where $e$ is a vector such that $\|e\| = 1$ and $\langle e, x \rangle = \langle e, y \rangle = 0$, and $\text{angle}(y, x)$ is interpreted as the "smallest angle through which $y$ would be rotated to become a nonnegative multiple of $x$". (See Exercise 2.19e on page 54.)

In the definition of angles between vectors given on page 37, $\text{angle}(y, x) = \text{angle}(x, y)$. As we pointed out there, sometimes it is important to distinguish the direction of the angle, and this is the case in equation (2.67), as in many applications in $\mathbb{R}^3$. The direction of angles in $\mathbb{R}^3$ often is used to determine the orientation of the principal axes in a coordinate system. The coordinate system is often defined to be "right-handed" (see Exercise 2.19f).

The cross product is useful in modeling phenomena in nature, which naturally are often represented as vectors in $\mathbb{R}^3$. The cross product is also useful in "three-dimensional" computer graphics for determining whether a given surface is visible from a given perspective and for simulating the effect of lighting on a surface.

## 2.3 Centered Vectors and Variances and Covariances of Vectors

In this section, we define some scalar-valued functions of vectors that are analogous to functions of random variables averaged over their probabilities or probability density. The functions of vectors discussed here are the same as the ones that define sample statistics. This short section illustrates the properties of norms, inner products, and angles in terms that should be familiar to the reader.

These functions, and transformations using them, are useful for applications in the data sciences. It is important to know the effects of various transformations of data on data analysis.

### 2.3.1 The Mean and Centered Vectors

When the elements of a vector have some kind of common interpretation, the sum of the elements or the mean (equation (2.46)) of the vector may have meaning. In this case, it may make sense to *center* the vector; that is, to subtract the mean from each element. For a given vector $x$, we denote its centered counterpart as $x_c$:

$$x_c = x - \bar{x}. \tag{2.68}$$

We refer to any vector whose sum of elements is 0 as a centered vector; note, therefore, for any centered vector $x_c$,

$$1^T x_c = 0;$$

or, indeed, for any constant vector $a$, $a^T x_c = 0$.

From the definitions, it is easy to see that

$$(x + y)_c = x_c + y_c \tag{2.69}$$

(see Exercise 2.20). Interpreting $\bar{x}$ as a vector, and recalling that it is the projection of $x$ onto the one vector, we see that $x_c$ is the residual in the sense of equation (2.52). Hence, we see that $x_c$ and $\bar{x}$ are orthogonal, and the Pythagorean relationship holds:

$$\|x\|^2 = \|\bar{x}\|^2 + \|x_c\|^2. \tag{2.70}$$

From this we see that the length of a centered vector is less than or equal to the length of the original vector. (Notice that equation (2.70) is just the formula familiar to data analysts, which with some rearrangement is $\sum (x_i - \bar{x})^2 = \sum x_i^2 - n\bar{x}^2$.)

For any scalar $a$ and $n$-vector $x$, expanding the terms, we see that

$$\|x - a\|^2 = \|x_c\|^2 + n(a - \bar{x})^2, \tag{2.71}$$

where we interpret $\bar{x}$ as a scalar here. An implication of this equation is that for all values of $a$, $\|x - a\|$ is minimized if $a = \bar{x}$.

Notice that a nonzero vector when centered may be the zero vector. This leads us to suspect that some properties that depend on a dot product are not invariant to centering. This is indeed the case. The angle between two vectors, for example, is not invariant to centering; that is, in general,

$$\text{angle}(x_c, y_c) \neq \text{angle}(x, y) \tag{2.72}$$

(see Exercise 2.21).

### 2.3.2 The Standard Deviation, the Variance, and Scaled Vectors

We also sometimes find it useful to scale a vector by both its length (normalize the vector) and by a function of its number of elements. We denote this *scaled* vector as $x_s$ and define it as

$$x_s = \sqrt{n - 1}\ \frac{x}{\|x_c\|}. \tag{2.73}$$

For comparing vectors, it is usually better to center the vectors prior to any scaling. We denote this *centered and scaled* vector as $x_{cs}$ and define it as

$$x_{cs} = \sqrt{n - 1}\ \frac{x_c}{\|x_c\|}. \tag{2.74}$$

Centering and scaling is also called *standardizing*. Note that the vector is centered before being scaled. The angle between two vectors is not changed by scaling (but, of course, it may be changed by centering).

The multiplicative inverse of the scaling factor,

$$s_x = \|x_c\| / \sqrt{n - 1}, \tag{2.75}$$

is called the *standard deviation* of the vector $x$. The standard deviation of $x_c$ is the same as that of $x$; in fact, the standard deviation is invariant to the addition of any constant. The standard deviation is a measure of how much the elements of the vector vary. If all of the elements of the vector are the same, the standard deviation is 0 because in that case $x_c = 0$.

The square of the standard deviation is called the *variance*, denoted by V:

$$\begin{aligned} \text{V}(x) &= s_x^2 \\ &= \frac{\|x_c\|^2}{n - 1}. \end{aligned} \tag{2.76}$$

(In perhaps more familiar notation, equation (2.76) is just $\text{V}(x) = \sum(x_i - \bar{x})^2 / (n - 1)$.) From equation (2.70), we see that

$$\text{V}(x) = \frac{1}{n - 1}\left(\|x\|^2 - \|\bar{x}\|^2\right).$$

(The terms "mean", "standard deviation", "variance", and other terms we will mention below are also used in an analogous, but slightly different, manner to refer to properties of *random variables*. In that context, the terms to refer to the quantities we are discussing here would be preceded by the word "sample", and often for clarity I will use the phrases "sample standard deviation" and "sample variance" to refer to what is defined above, especially if the elements of $x$ are interpreted as independent realizations of a random variable. Also, recall the two possible meanings of "mean", or $\bar{x}$; one is a vector, and one is a scalar, as in equation (2.47).)

If $a$ and $b$ are scalars (or $b$ is a vector with all elements the same), the definition, together with equation (2.71), immediately gives

$$\mathrm{V}(ax + b) = a^2 \mathrm{V}(x).$$

This implies that for the scaled vector $x_\mathrm{s}$,

$$\mathrm{V}(x_\mathrm{s}) = 1.$$

If $a$ is a scalar and $x$ and $y$ are vectors with the same number of elements, from the equation above, and using equation (2.31) on page 26, we see that the variance following an axpy operation is given by

$$\mathrm{V}(ax + y) = a^2 \mathrm{V}(x) + \mathrm{V}(y) + 2a \frac{\langle x_\mathrm{c},\ y_\mathrm{c} \rangle}{n - 1}. \tag{2.77}$$

While equation (2.76) appears to be relatively simple, evaluating the expression for a given $x$ may not be straightforward. We discuss computational issues for this expression on page 502. This is an instance of a principle that we will encounter repeatedly: *the form of a mathematical expression and the way the expression should be evaluated in actual practice may be quite different.*

### 2.3.3 Covariances and Correlations Between Vectors

If $x$ and $y$ are $n$-vectors, the *covariance* between $x$ and $y$ is

$$\mathrm{Cov}(x, y) = \frac{\langle x - \bar{x},\ y - \bar{y} \rangle}{n - 1}. \tag{2.78}$$

By representing $x - \bar{x}$ as $x - \bar{x}1$ and $y - \bar{y}$ similarly, and expanding, we see that $\mathrm{Cov}(x, y) = (\langle x,\ y \rangle - n\bar{x}\bar{y})/(n - 1)$. Also, we see from the definition of covariance that $\mathrm{Cov}(x, x)$ is the variance of the vector $x$, as defined above.

From the definition and the properties of an inner product given on page 24, if $x$, $y$, and $z$ are conformable vectors, we see immediately that

- $\mathrm{Cov}(x, y) = 0$
  if $\mathrm{V}(x) = 0$ or $\mathrm{V}(y) = 0$;
- $\mathrm{Cov}(ax, y) = a\mathrm{Cov}(x, y)$
  for any scalar $a$;

- $\text{Cov}(y, x) = \text{Cov}(x, y)$;
- $\text{Cov}(y, y) = \text{V}(y)$; and
- $\text{Cov}(x + z, y) = \text{Cov}(x, y) + \text{Cov}(z, y)$,
  in particular,
  - $\text{Cov}(x + y, y) = \text{Cov}(x, y) + \text{V}(y)$, and
  - $\text{Cov}(x + a, y) = \text{Cov}(x, y)$
    for any scalar $a$.

Using the definition of the covariance, we can rewrite equation (2.77) as

$$\text{V}(ax + y) = a^2\text{V}(x) + \text{V}(y) + 2a\text{Cov}(x, y). \tag{2.79}$$

The covariance is a measure of the extent to which the vectors point in the same direction. A more meaningful measure of this is obtained by the covariance of the centered and scaled vectors. This is the *correlation* between the vectors, which if $\|x_\text{c}\| \neq 0$ and $\|y_\text{c}\| \neq 0$,

$$\begin{aligned}
\text{Cor}(x, y) &= \text{Cov}(x_\text{cs}, y_\text{cs}) \\
&= \left\langle \frac{x_\text{c}}{\|x_\text{c}\|}, \ \frac{y_\text{c}}{\|y_\text{c}\|} \right\rangle \\
&= \frac{\langle x_\text{c}, \ y_\text{c} \rangle}{\|x_\text{c}\|\|y_\text{c}\|}.
\end{aligned} \tag{2.80}$$

If $\|x_\text{c}\| = 0$ or $\|y_\text{c}\| = 0$, we define $\text{Cor}(x, y)$ to be 0. We see immediately from equation (2.54) that the correlation is the cosine of the angle between $x_\text{c}$ and $y_\text{c}$:

$$\text{Cor}(x, y) = \cos(\text{angle}(x_\text{c}, y_\text{c})). \tag{2.81}$$

(Recall that this is not the same as the angle between $x$ and $y$.)

An equivalent expression for the correlation, so long as $\text{V}(x) \neq 0$ and $\text{V}(y) \neq 0$, is

$$\text{Cor}(x, y) = \frac{\text{Cov}(x, y)}{\sqrt{\text{V}(x)\text{V}(y)}}. \tag{2.82}$$

It is clear that the correlation is in the interval $[-1, 1]$ (from the Cauchy-Schwarz inequality). A correlation of $-1$ indicates that the vectors point in opposite directions, a correlation of 1 indicates that the vectors point in the same direction, and a correlation of 0 indicates that the vectors are orthogonal.

While the covariance is equivariant to scalar multiplication, the absolute value of the correlation is invariant to it; that is, the correlation changes only as the sign of the scalar multiplier,

$$\text{Cor}(ax, y) = \text{sign}(a)\text{Cor}(x, y), \tag{2.83}$$

for any scalar $a$.

## Exercises

2.1. Write out the step-by-step proof that the maximum number of $n$-vectors that can form a set that is linearly independent is $n$.

2.2. Prove inequalities (2.10) and (2.11).

2.3. a) Give an example of a vector space and a subset of the set of vectors in it such that that subset together with the axpy operation is *not* a vector space.

b) Give an example of two vector spaces such that the union of the sets of vectors in them together with the axpy operation is *not* a vector space.

2.4. Prove the equalities (2.15) and (2.16).
*Hint*: Use of basis sets makes the details easier.

2.5. Prove (2.19).

2.6. Let $\{v_i\}_{i=1}^n$ be an orthonormal basis for the $n$-dimensional vector space $\mathcal{V}$. Let $x \in \mathcal{V}$ have the representation

$$x = \sum b_i v_i.$$

Show that the Fourier coefficients $b_i$ can be computed as

$$b_i = \langle x, v_i \rangle.$$

2.7. Show that if the norm is induced by an inner product that the parallelogram equality, equation (2.32), holds.

2.8. Let $p = \frac{1}{2}$ in equation (2.33); that is, let $\rho(x)$ be defined for the $n$-vector $x$ as

$$\rho(x) = \left( \sum_{i=1}^n |x_i|^{1/2} \right)^2.$$

Show that $\rho(\cdot)$ is not a norm.

2.9. Show that the $L_1$ norm is not induced by an inner product.
*Hint*: Find a counterexample that does not satisfy the parallelogram equality (equation (2.32)).

2.10. Prove equation (2.34) and show that the bounds are sharp by exhibiting instances of equality. (Use the fact that $\|x\|_\infty = \max_i |x_i|$.)

2.11. Prove the following inequalities.

a) Prove Hölder's inequality: for any $p$ and $q$ such that $p \geq 1$ and $p + q = pq$, and for vectors $x$ and $y$ of the same order,

$$\langle x, y \rangle \leq \|x\|_p \|y\|_q.$$

b) Prove the triangle inequality for any $L_p$ norm. (This is sometimes called Minkowski's inequality.)
*Hint*: Use Hölder's inequality.

2.12. Show that the expression defined in equation (2.42) on page 32 is a metric.

2.13. Show that equation (2.53) on page 37 is correct.

2.14. Show that the intersection of two orthogonal vector spaces consists only of the zero vector.

2.15. From the definition of direction cosines in equation (2.55), it is easy to see that the sum of the squares of the direction cosines is 1. For the special case of $\mathbb{R}^3$, draw a sketch and use properties of right triangles to show this geometrically.

2.16. In $\mathbb{R}^2$ with a Cartesian coordinate system, the diagonal directed line segment through the positive quadrant (orthant) makes a 45° angle with each of the positive axes. In 3 dimensions, what is the angle between the diagonal and each of the positive axes? In 10 dimensions? In 100 dimensions? In 1000 dimensions? We see that in higher dimensions any two lines are almost orthogonal. (That is, the angle between them approaches 90°.) What are some of the implications of this for data analysis?

2.17. Show that if the initial set of vectors are linearly independent, all residuals in Algorithm 2.1 are nonzero. (For given $k \geq 2$, all that is required is to show that

$$\tilde{x}_k - \langle \tilde{x}_{k-1}, \tilde{x}_k \rangle \tilde{x}_{k-1} \neq 0$$

if $\tilde{x}_k$ and $\tilde{x}_{k-1}$ are linearly independent. Why?)

2.18. Convex cones.

a) I defined a convex cone as a set of vectors (not necessarily a cone) such that for any two vectors $v_1, v_2$ in the set and for any nonnegative real numbers $a, b \geq 0$, $av_1 + bv_2$ is in the set. Then I stated that an equivalent definition requires first that the set be a cone, and then includes the requirement $a + b = 1$ along with $a, b \geq 0$ in the definition of a convex cone. Show that the two definitions are equivalent.

b) The restriction that $a + b = 1$ in the definition of a convex cone is the kind of restriction that we usually encounter in definitions of convex objects. Without this restriction, it may seem that the linear combinations may get "outside" of the object. Show that this is not the case for convex cones.

In particular in the two-dimensional case, show that if $x = (x_1, x_2)$, $y = (y_1, y_2)$, with $x_1/x_2 < y_1/y_2$ and $a, b \geq 0$, then

$$x_1/x_2 \leq (ax_1 + by_1)/(ax_2 + by_2) \leq y_1/y_2.$$

This should also help to give a geometrical perspective on convex cones.

c) Show that if $C_1$ and $C_2$ are convex cones over the same vector space, then $C_1 \cap C_2$ is a convex cone. Give a counterexample to show that $C_1 \cup C_2$ is not necessarily a convex cone.

2.19. $\mathbb{R}^3$ and the cross product.

   a) Is the cross product associative? Prove or disprove.

   b) For $x, y \in \mathbb{R}^3$, show that the area of the triangle with vertices $(0, 0, 0)$, $x$, and $y$ is $\|x \times y\|/2$.

   c) For $x, y, z \in \mathbb{R}^3$, show that

$$\langle x, \ y \times z \rangle = \langle x \times y, \ z \rangle.$$

   This is called the "triple scalar product".

   d) For $x, y, z \in \mathbb{R}^3$, show that

$$x \times (y \times z) = \langle x, \ z \rangle y - \langle x, \ y \rangle z.$$

   This is called the "triple vector product". It is in the plane determined by $y$ and $z$.

   e) The magnitude of the angle between two vectors is determined by the cosine, formed from the inner product. Show that in the special case of $\mathbb{R}^3$, the angle is also determined by the sine and the cross product, and show that this method can determine both the magnitude and the *direction* of the angle; that is, the way a particular vector is rotated into the other.

   f) In a Cartesian coordinate system in $\mathbb{R}^3$, the principal axes correspond to the unit vectors $e_1 = (1, 0, 0)$, $e_2 = (0, 1, 0)$, and $e_3 = (0, 0, 1)$. This system has an indeterminate correspondence to a physical three-dimensional system; if the plane determined by $e_1$ and $e_2$ is taken as horizontal, then $e_3$ could "point upward" or "point downward". A simple way that this indeterminacy can be resolved is to require that the principal axes have the orientation of the thumb, index finger, and middle finger of the right hand when those digits are spread in orthogonal directions, where $e_1$ corresponds to the index finger, $e_2$ corresponds to the middle finger, and $e_3$ corresponds to the thumb. This is called a "right-hand" coordinate system. Show that in a right-hand coordinate system, if we interpret the angle between $e_i$ and $e_j$ to be measured in the direction from $e_i$ to $e_j$, then $e_3 = e_1 \times e_2$ and $e_3 = -e_2 \times e_1$.

2.20. Using equations (2.46) and (2.68), establish equation (2.69).

2.21. Show that the angle between the centered vectors $x_c$ and $y_c$ is not the same in general as the angle between the uncentered vectors $x$ and $y$ of the same order.

2.22. Formally prove equation (2.77) (and hence equation (2.79)).

2.23. Let $x$ and $y$ be any vectors of the same order over the same field.

   a) Prove

$$(\mathrm{Cov}(x, y))^2 \leq \mathrm{V}(x)\mathrm{V}(y).$$

   b) Hence, prove

$$-1 \leq \mathrm{Cor}(x, y)) \leq 1.$$