# Important definitions and results

## 1. Algebra and geometry of vectors

**Definition 1.1.** *A <u>linear combination</u> of vectors $\mathbf{v}_1, \ldots, \mathbf{v}_k \in \mathbf{R}^n$ is a vector of the form*

$$c_1 \mathbf{v}_1 + \cdots + c_k \mathbf{v}_k$$

*where $c_1, \ldots, c_k \in \mathbf{R}$ are scalars. The <u>span</u> of $\mathbf{v}_1, \ldots, \mathbf{v}_k$ is the set $\mathrm{span}(\mathbf{v}_1, \ldots, \mathbf{v}_k)$ of all possible linear combinations of $\mathbf{v}_1, \ldots, \mathbf{v}_k$.*

**Definition 1.2.** *Let $p \in \mathbf{R}^n$ be a point and $\mathbf{v} \in \mathbf{R}^n$ be a non-zero vector. The <u>line</u> through $p$ in direction $\mathbf{v}$ is the set*

$$L = \{p + t\mathbf{v} \in \mathbf{R}^n : t \in \mathbf{R}\}.$$

**Definition 1.3.** *The <u>dot product</u> of two vectors $\mathbf{v}, \mathbf{w} \in \mathbf{R}^n$ is the quantity*

$$\mathbf{v} \cdot \mathbf{w} := \sum_{j=1}^{n} v_j w_j.$$

*The <u>length</u> of a vector $\mathbf{v} \in \mathbf{R}^n$ is the quantity $\|\mathbf{v}\| := \sqrt{\mathbf{v} \cdot \mathbf{v}}$.*

**Definition 1.4.** *Two vectors $\mathbf{v}, \mathbf{w} \in \mathbf{R}^n$ are <u>parallel</u> if one is a scalar multiple of the other. They are <u>orthogonal</u> (or <u>perpendicular</u>) if $\mathbf{v} \cdot \mathbf{w} = 0$.*

**Theorem 1.5** (Orthogonal decomposition)**.** *Let $\mathbf{v}, \mathbf{w} \in \mathbf{R}^n$ be vectors such that $\mathbf{w} \neq \mathbf{0}$. Then there are unique vectors $\mathbf{v}_{\parallel}, \mathbf{v}_{\perp} \in \mathbf{R}^n$ such that*

- $\mathbf{v} = \mathbf{v}_{\parallel} + \mathbf{v}_{\perp}$;
- *$\mathbf{v}_{\parallel}$ is parallel to $\mathbf{w}$ and $\mathbf{v}_{\perp}$ is orthogonal to $\mathbf{w}$.*

The proof of the theorem gives a formula for $\mathbf{v}_{\parallel}$, and I use this to define orthogonal projection of $\mathbf{v}$ onto $\mathbf{w}$.

**Definition 1.6.** *Let $\mathbf{v}, \mathbf{w} \in \mathbf{R}^n$ be vectors such that $\mathbf{w} \neq \mathbf{0}$ The <u>orthogonal projection</u> of $\mathbf{v}$ onto $\mathbf{w}$ is the vector*

$$\mathrm{proj}_{\mathbf{w}}(\mathbf{v}) := \frac{\mathbf{v} \cdot \mathbf{w}}{\|\mathbf{w}\|^2} \mathbf{w}.$$

**Theorem 1.7** (Cauchy-Schwarz Inequality)**.** *For any vectors $\mathbf{v}, \mathbf{w} \in \mathbf{R}^n$ we have*

$$|\mathbf{v} \cdot \mathbf{w}| \leq \|\mathbf{v}\| \, \|\mathbf{w}\|$$

*with equality if and only if the vectors are parallel.*

**Theorem 1.8** (Triangle Inequality)**.** *For any vectors $\mathbf{v}, \mathbf{w} \in \mathbf{R}^n$ we have*

$$\|\mathbf{v} + \mathbf{w}\| \leq \|\mathbf{v}\| + \|\mathbf{w}\| \, .$$

## 2. LINEAR SYSTEMS

An $\underline{m \times n \text{ matrix}}$ is an array $A = (a_{ij})$ where $a_{ij} \in \mathbf{R}$ for each $1 \le i \le m$ and $1 \le j \le n$.

$$A = (a_{ij}) := \begin{bmatrix} a_{11} & a_{12} & \ldots & a_{1n} \\ a_{21} & a_{22} & \ldots & a_{2n} \\ & & \vdots & \\ a_{m1} & a_{m2} & \ldots & a_{mn} \end{bmatrix}$$

The number $a_{ij}$ is called the $\underline{ij\text{-entry}}$ of $A$. The vector

$$\mathbf{a}_j := \begin{bmatrix} a_{1j} \\ \ldots \\ a_{mj} \end{bmatrix}.$$

is called the $j$th column of $A$. Similarly, we call the vector

$$[a_{i1}, \ldots, a_{in}]$$

the $i$th row of $A$. We will often write a matrix $A = [\mathbf{a}_1, \ldots, \mathbf{a}_n]$ in terms of its columns. Later we will sometimes write a matrix in terms of its rows.

We will let $\mathcal{M}_{m \times n}$ denote the set of all $m \times n$ matrices. A matrix is called $\underline{square}$ if $m = n$, i.e. if it has the same number of rows as columns.

**Definition 2.1.** *The product of a matrix $A = [\mathbf{a}_1, \ldots, \mathbf{a}_n] \in \mathcal{M}_{m \times n}$ with a vector $\mathbf{x} \in \mathbf{R}^n$ is the vector*

$$A\mathbf{x} := x_1\mathbf{a}_1 + \cdots + x_n\mathbf{a}_n.$$

Note that in this definition it is important that the matrix and vector have compatible sizes, i.e. that the number of columns of $A$ equals the number of entries of $\mathbf{x}$. Also, it is important to write $A\mathbf{x}$ instead of $\mathbf{x}A$. The reason for this convention will become apparent later.

There are two matrices that are quite special from the point of view of matrix/vector multiplication. Specifically

- Let $0 \in \mathcal{M}_{m \times n}$ denote the matrix with all zero entries. Then $0\mathbf{x} = \mathbf{0}$ for every $\mathbf{x} \in \mathbf{R}^n$.
- Let $I \in \mathcal{M}_{n \times n}$ denote the square matrix with all 'diagonal' entries $a_{ii} = 1$ and all other entries equal to 0. Then $I\mathbf{x} = \mathbf{x}$ for all $\mathbf{x} \in \mathbf{R}^n$.

The matrix $I$ is called the $\underline{identity}$ matrix.

**Definition 2.2.** *An $\underline{m \times n \text{ linear system}}$ is a matrix/vector equation $A\mathbf{x} = \mathbf{b}$, where $A$ is an $m \times n$ matrix (the $\underline{\text{coefficient matrix}}$), $\mathbf{x} \in \mathbf{R}^n$ is a vector with unknown entries and $\mathbf{b} \in \mathbf{R}^m$ is a vector with given entries. Moreover,*

- *The $\underline{\text{augmented matrix}}$ for the system is the $m \times (n+1)$ matrix $\begin{bmatrix} A & \mathbf{b} \end{bmatrix}$.*
- *The $\overline{\text{system is}}$ $\underline{\text{homogeneous}}$ if $\mathbf{b} = \mathbf{0}$.*
- *The system is $\underline{\text{consistent}}$ if there exists at least one solution $\mathbf{x} \in \mathbf{R}^n$.*

Note that $\mathbf{x} = \mathbf{0}$ is always a solution (called the $\underline{\textit{trivial solution}}$ of a homogeneous system $A\mathbf{x} = \mathbf{0}$. In class I discussed the notion of an $\underline{\textit{elementary row operation}}$, which turns one $m \times n$ matrix $A$ into another $m \times n$ matrix $B$.

**Definition 2.3.** *Two matrices $A, B \in \mathcal{M}_{m \times n}$ are* <u>row equivalent</u> *if one can be transformed into the other by a finite sequence of elementary* <u>row operations</u>. *In this case, we write $A \sim B$.*

**Definition 2.4.** *Let $A$ be an $m \times n$ matrix. The leftmost non-vanishing entry (if one exists) in each row of $A$ is called the* <u>pivot</u> *entry for that row. We say that $A$ is in* <u>echelon form</u> *if the pivot in each row lies to the right of the pivot in the previous row. In particular any row with a non-zero entry is above any row with all zero entries. We further say that $A$ is in* <u>reduced echelon form</u> *if additionally*

- *Each pivot entry is equal to 1; and*
- *Each pivot entry is the only non-zero entry in its column.*

Using the method of Gaussian elimination, one shows the following

**Theorem 2.5.** *Every matrix is row equivalent to a matrix in reduced echelon form.*

In fact, it can be shown (as Shifrin does but I won't) that any given matrix $A$ is row equivalent to *exactly one* matrix in reduced echelon form. This is perhaps surprising since the sequence of row operations needed to get to reduced echelon form is far from unique.

**Definition 2.6.** *Let $\begin{bmatrix} A & \mathbf{b} \end{bmatrix}$ be the augmented matrix of a linear system and $\begin{bmatrix} \tilde{A} & \tilde{\mathbf{b}} \end{bmatrix}$ be a row equivalent augmented matrix in reduced echelon form. The* <u>pivot variables</u> *of $A\mathbf{x} = \mathbf{b}$ are the entries $x_j$ of $\mathbf{x}$ corresponding to pivot columns of $\begin{bmatrix} \tilde{A} & \tilde{\mathbf{b}} \end{bmatrix}$. The other entries of $\mathbf{x}$ are called* <u>free variables</u>.

Note that (with the notation in this definition)

- $A\mathbf{x} = \mathbf{b}$ is consistent if and only if the last column of $\begin{bmatrix} \tilde{A}\tilde{\mathbf{b}} \end{bmatrix}$ has no pivots.
- A consistent system $A\mathbf{x} = \mathbf{b}$ has exactly one solution if there are no free variables.
- A consistent system $A\mathbf{x} = \mathbf{b}$ has infinitely many solutions if there are free variables.

Since a matrix in (reduced) echelon form has at most one pivot in each row/column, linear systems $A\mathbf{x} = \mathbf{b}$ behave best only when $A$ is a square matrix. More precisely,

- If $A$ is an $m \times n$ matrix with $n > m$, then solutions of $A\mathbf{x} = \mathbf{b}$ are never unique.
- If $A$ is an $m \times n$ matrix with $m > n$, there always exist vectors $\mathbf{b} \in \mathbf{R}^m$ such that $A\mathbf{x} = \mathbf{b}$ is inconsistent.

For square matrices, life is *usually* better.

**Definition 2.7.** *A square matrix $A$ is* <u>non-singular</u> *if it is row equivalent to the identity matrix.*

**Proposition 2.8.** *Let $A$ be an $n \times n$ matrix. Then the following are equivalent.*

- *$A$ is non-singular.*
- *$A\mathbf{x} = \mathbf{0}$ has only the trivial solution.*
- *$A\mathbf{x} = \mathbf{b}$ is consistent for any $\mathbf{b} \in \mathbf{R}^n$.*
- *$A\mathbf{x} = \mathbf{b}$ has a unique solution for some $\mathbf{b} \in \mathbf{R}^n$.*
- *$A\mathbf{x} = \mathbf{b}$ has a unique solution for any $\mathbf{b} \in \mathbf{R}^n$.*

In particular, when $A$ is a square matrix, whether or not $A\mathbf{x} = \mathbf{b}$ has a unique solution depends on only $A$ and not at all on $\mathbf{b}$.

## 3. Linear Transformations

Recall that a function $f : X \to Y$ from a set $X$ to a set $Y$ is a 'rule' that associates each element $x \in X$ to *exactly one element* $f(x) \in Y$. The sets $X$ and $Y$ are known as the *source* (i.e. domain) and *target* (i.e. codomain) of the function. The set

$$f(X) := \{y \in Y : y = f(x) \text{ for some } x \in X\}$$

is called the *image* (i.e. range) of the function.

**Definition 3.1.** *A* linear transformation *is a function* $T : \mathbf{R}^n \to \mathbf{R}^m$ *satisfying*

- $T(c\mathbf{x}) = cT(\mathbf{x})$
- $T(\mathbf{x} + \mathbf{y}) = T(\mathbf{x}) + T(\mathbf{y})$

*for all vectors* $\mathbf{x}, \mathbf{y} \in \mathbf{R}^n$ *and scalars* $c \in \mathbf{R}$.

This definition can be restated by saying that $T$ is linear if and only if $T$ distributes over linear combinations, i.e.

$$T(c_1\mathbf{v}_1 + \ldots c_k\mathbf{v}_k) = c_1 T(\mathbf{v}_1) + \ldots c_k T(\mathbf{v}_k).$$

The most important example of a linear transformation is given by

**Proposition 3.2.** *Let $A$ be an $m \times n$ matrix and $T : \mathbf{R}^n \to \mathbf{R}^m$ be the function given by $T(\mathbf{x}) = A\mathbf{x}$. Then $T$ is a linear transformation.*

Thus we can rewrite a linear system $A\mathbf{x} = \mathbf{b}$ as $T(\mathbf{x}) = \mathbf{b}$ where $T$ is a linear transformation. This leads us to some information about the structure of the set of solutions of a linear system.

**Corollary 3.3.** *Given a matrix $A \in \mathcal{M}_{m \times n}$, suppose that the vectors $\mathbf{x}_1, \ldots, \mathbf{x}_k \in \mathbf{R}^n$ each solve the homogeneous system $A\mathbf{x} = \mathbf{0}$. Then any linear combination $\mathbf{x} = c_1\mathbf{x}_1 + \cdots + c_k\mathbf{x}_k$ also solves $A\mathbf{x} = \mathbf{0}$.*

**Corollary 3.4.** *Given a matrix $A \in \mathcal{M}_{m \times n}$, a vector $\mathbf{b} \in \mathbf{R}^m$ and a solution $\mathbf{x}_0 \in \mathbf{R}^n$ of $A\mathbf{x} = \mathbf{b}$, we have that*

- *for any other vector $\mathbf{x}_1 \in \mathbf{R}^n$, $\mathbf{x} = \mathbf{x}_1$ solves $A\mathbf{x} = \mathbf{b}$ if and only if $\mathbf{x} = \mathbf{x}_1 - \mathbf{x}_0$ solves $A\mathbf{x} = \mathbf{0}$; and*
- *in particular $\mathbf{x}_1$ is the only solution of $A\mathbf{x} = \mathbf{b}$ if and only if $A\mathbf{x} = \mathbf{0}$ has only the trivial solution.*

It should be pointed out that different matrices $A$ and $B$ define different linear transformations $T(\mathbf{x}) = A\mathbf{x}$ and $S(\mathbf{x}) = B\mathbf{x}$. If the sizes of $A$ and $B$ are different then $S$ and $T$ don't even have the same target and source. If $A$ and $B$ are different matrices of the same size, then one of the columns $\mathbf{a}_j$ of $A$ must be different from the corresponding column $\mathbf{b}_j$ of $B$. Now if $\mathbf{e}_j$ is the *jth standard basis vector*, i.e. the vector whose $j$th coordinate is 1 and all of whose other coordinates vanish, then we have $T(\mathbf{e}_j) = \mathbf{a}_j \neq \mathbf{b}_j = S(\mathbf{e}_j)$. So again $T \neq S$.

The idea of picking off columns of $A$ by applying $T$ to the standard basis vectors $\mathbf{e}_j$ allows us to show that *all* linear transformations (as I have defined them here) are represented by matrices.

**Proposition 3.5.** *Let $T : \mathbf{R}^n \to \mathbf{R}^m$ be a linear transformation and $A$ be the matrix whose $j$th column is $T(\mathbf{e}_j)$. Then $T(\mathbf{x}) = A\mathbf{x}$ for all $\mathbf{x} \in \mathbf{R}^n$.*

We will call the matrix $A$ in this proposition the *standard matrix for $T$*. Note in particular, that the *identity transformation* id : $\mathbf{R}^n \to \mathbf{R}^n$ given by id($\mathbf{x}$) = $\mathbf{x}$ has standard matrix $I$ and the *zero transformation* 0 : $\mathbf{R}^n \to \mathbf{R}^m$ given by 0($\mathbf{x}$) = $\mathbf{0}$ has standard matrix 0.

### 3.1. Algebra of linear transformations.

**Proposition 3.6.** *Let $S, T : \mathbf{R}^n \to \mathbf{R}^m$ be linear transformations and $c \in \mathbf{R}$ be a scalar. Then $T + S$ and $cT$ are both linear transformations. If, moreover, $A, B \in \mathcal{M}_{m \times n}$ are the standard matrices for $T$ and $S$, respectively, then*

- *the standard matrix for $T + S$ is the matrix $A + B$ obtained from $A$ and $B$ by adding corresponding entries of these two matrices*
- *the standard matrix for $cT$ is the matrix $cA$ obtained by multiplying each entry of $A$ by $c$.*

Note that addition and scalar multiplication of matrices works just like it does for vectors. Hence all the rules that work for addition and scalar multiplication of vectors also work for matrices. But for matrices we have a different sort of operation.

**Proposition 3.7.** *Let $S : \mathbf{R}^p \to \mathbf{R}^n$ and $T : \mathbf{R}^n \to \mathbf{R}^m$ be linear transformations. Then the composition $T \circ S : \mathbf{R}^p \to \mathbf{R}^m$ is a linear transformation. If $A$ and $B$ are the standard matrices for $T$ and $S$, then the standard matrix for $T \circ S$ is the $m \times p$ matrix $AB$ given by*

$$AB := \begin{bmatrix} A\mathbf{b}_1 & \dots & A\mathbf{b}_p \end{bmatrix},$$

*where $\mathbf{b}_1, \dots, \mathbf{b}_p$ are the columns of $B$.*

Note that for a product $AB$ to make sense, we need the number of rows of $B$ to equal the number of columns of $A$. Matrix multiplication behaves like regular multiplication in several important ways. All of the following facts can be established by appealing to facts about composition of linear transformations.

**Proposition 3.8.** *Given matrices $A, A' \in \mathcal{M}_{m \times n}$, $B, B' \in \mathcal{M}_{n \times p}$, $C \in \mathcal{M}_{p \times q}$ and a scalar $c \in \mathbf{R}$, we have*

- *$I_{m \times m} A = A I_{n \times n} = A$;*
- *$(A + A')B = AB + A'B$ and $A(B + B') = AB + AB'$;*
- *$c(AB) = (cA)B = A(cB)$; and*
- *$A(BC) = (AB)C$.*

The statement that is missing from this list (because it is false) is $AB = BA$. First of all, the product on one side might make sense while the product on the other might not. Secondly, even when both products make sense, they might result in matrices of different sizes. And finally, even when $A$ and $B$ are square matrices of the same size, so that both $AB$ and $BA$ make sense and are the same size, it is usually the case that $AB \neq BA$. So watch out!

**Definition 3.9.** *A linear transformation $T : \mathbf{R}^n \to \mathbf{R}^m$ is* invertible *if there is a linear transformation $S : \mathbf{R}^m \to \mathbf{R}^n$ such that $S \circ T(\mathbf{x}) = \mathbf{x}$ and $T \circ S(\mathbf{y}) = \mathbf{y}$ for all $\mathbf{x} \in \mathbf{R}^n$ and $\mathbf{y} \in \mathbf{R}^m$. We call $S$ the* inverse *of $T$ and write $T^{-1} = S$.*

**Proposition 3.10.** *Let $T : \mathbf{R}^n \to \mathbf{R}^m$ be a linear transformation and $A \in \mathcal{M}_{m \times n}$ be its standard matrix. Then the following statements are equivalent.*

(1) *$T$ is invertible.*

   (2) *There exists a matrix $B \in \mathcal{M}_{n \times m}$ such that $AB = BA = I$.*

   (3) *For any $\mathbf{b} \in \mathbf{R}^m$ the linear system $A\mathbf{x} = \mathbf{b}$ has exactly one solution.*

   (4) *$A$ is a non-singular square (i.e. $m = n$) matrix.*

**Definition 3.11.** *A square matrix $A \in \mathcal{M}_{n \times n}$ is* invertible *if there is a matrix $B \in \mathcal{M}_{n \times n}$ such that $AB = BA = I$. We call $B$ the* inverse *of $A$ and write $A^{-1} = B$.*

It is both interesting and convenient that one doesn't need to check both $AB$ and $BA$ to verify that $A^{-1} = B$.

**Proposition 3.12.** *Let $A, B \in \mathcal{M}_{n \times n}$ be square matrices. Then $AB = I$ if and only if $BA = I$.*

Here are a couple more useful, albeit straightforward observations

**Proposition 3.13.** *If $A \in \mathcal{M}_{n \times n}$ is invertible, then so is $A^{-1}$, and the inverse is given by $(A^{-1})^{-1} = A$. If $A_1, \ldots, A_k \in \mathcal{M}_{n \times n}$ are invertible, then so is the product $A_1 \ldots A_k$, and the inverse is given by*

$$(A_1 \ldots A_k)^{-1} = A_k^{-1} \ldots A_1^{-1}.$$

**Definition 3.14.** *An* elementary matrix *$E \in \mathcal{M}_{n \times n}$ is one obtained by performing a single elementary row operation on the identity matrix $I_{n \times n}$.*

One can check case-by-case that if $A \in \mathcal{M}_{m \times n}$ is a matrix, then performing a given elementary row operation on $A$ results in the matrix $EA$, where $E \in \mathcal{M}_{m \times m}$ is the elementary matrix for that row operation. Since, $EA = \begin{bmatrix} E\mathbf{a}_1 & \ldots & E\mathbf{a}_n \end{bmatrix}$, it suffices to check that this works for any vector $\mathbf{a} \in \mathbf{R}^m$. Since any row operation can be undone by some other row operation, it follows that an elementary matrix $E$ is invertible and that the inverse $E^{-1}$ is also an elementary matrix. The fact that any matrix is row equivalent to a matrix in reduced echelon form can therefore be restated as follows.

**Proposition 3.15.** *Any matrix $A \in \mathcal{M}_{m \times n}$ can be decomposed $A = E_k \ldots E_1 \tilde{A}$ into a product of elementary matrices $E_j \in \mathcal{M}_{m \times m}$ and a matrix $\tilde{A}$ in reduced echelon form. In particular, any non-singular square matrix $A$ can be written $A = E_k \ldots E_1$ as a product of elementary matrices.*

## 4. Subspaces

The following definition formalizes and generalizes the notion of a line or plane through the origin.

**Definition 4.1.** *A set of vectors $V \subset \mathbf{R}^n$ is a* subspace *if*
    (1) $\mathbf{0} \in V$;
    (2) *for any $\mathbf{v} \in V$ and $c \in \mathbf{R}$, we have $c\mathbf{v} \in V$; and*
    (3) *for any $\mathbf{v}, \mathbf{w} \in V$, we have $\mathbf{v} + \mathbf{w} \in V$.*

Note that $\mathbf{R}^n$ is a subspace of $\mathbf{R}^n$. So is the *trivial subspace* $\{\mathbf{0}\} \subset \mathbf{R}^n$. Beyond this, there are mainly only two other ways of obtaining subspaces of $\mathbf{R}^n$.

**Proposition 4.2.** *Let $\mathbf{v}_1, \ldots, \mathbf{v}_k \in \mathbf{R}^n$ be a list of vectors. Then $\mathrm{span}(\mathbf{v}_1, \ldots, \mathbf{v}_k)$ is a subspace of $\mathbf{R}^n$.*

**Definition 4.3.** *Let $S \subset \mathbf{R}^n$ be any set of vectors. The* orthogonal complement *of $S$ is the set*
$$S^\perp := \{\mathbf{v} \in \mathbf{R}^n : \mathbf{v} \cdot \mathbf{w} = 0 \text{ for all } \mathbf{w} \in S\}.$$

**Proposition 4.4.** *The orthogonal complement $S^\perp$ of a set $S \subset \mathbf{R}^n$ is a subspace.*

The set $S$ in this proposition need not be a subspace itself. However, if it is, it can be easier to test whether a vector belongs to $S^\perp$.

**Proposition 4.5.** *Given vectors $\mathbf{v}_1, \ldots, \mathbf{v}_k \in \mathbf{R}^n$, let $V = \mathrm{span}\{\mathbf{v}_1, \ldots, \mathbf{v}_k\}$. Then $\mathbf{v} \in V^\perp$ if and only if $\mathbf{v} \cdot \mathbf{v}_1 = \cdots = \mathbf{v} \cdot \mathbf{v}_k = 0$.*

Each matrix gives rise to several different subspaces.

**Definition 4.6.** *Let $A = \begin{bmatrix} \mathbf{a}_1 & \ldots & \mathbf{a}_n \end{bmatrix} \in \mathcal{M}_{m \times n}$ be a matrix.*
- *The* column space *of $A$ is the span $\mathrm{col}\, A \subset \mathbf{R}^m$ of the columns $\mathbf{a}_1, \ldots, \mathbf{a}_n$ of $A$.*
- *The* row space *of $A$ is likewise the span $\mathrm{row}\, A \subset \mathbf{R}^n$ of the rows of $A$.*
- *The* null space *of $A$ is the set*
$$\mathrm{nul}\, A := \{\mathbf{x} \in \mathbf{R}^n : A\mathbf{x} = \mathbf{0}\}.$$

It is useful to observe that

**Proposition 4.7.** *For any $A \in \mathcal{M}_{m \times n}$, we have $\mathrm{nul}\, A = (\mathrm{row}\, A)^\perp$*

**Corollary 4.8.** *The column, row and null spaces of a matrix are all subspaces.*

Word of caution: remember that the column space of an $m \times n$ matrix is a subspace of $\mathbf{R}^m$, whereas row and null spaces are subspaces of $\mathbf{R}^n$. This is perhaps easier to keep straight if you think in terms of the linear transformation $T : \mathbf{R}^n \to \mathbf{R}^m$ associated to $A$. The row and null spaces are subspaces of the *source* $\mathbf{R}^n$ of $T$, whereas the column space of $A$ is contained in the *target* $\mathbf{R}^m$ of $T$. In fact, if you think it through, you'll find that the column space of $A$ is the same as the *image* (i.e. range) $T(\mathbf{R}^n)$ of $T$.

## 5. Bases and dimension

**Definition 5.1.** *A finite sequence* $\mathbf{v}_1, \ldots, \mathbf{v}_k \in \mathbf{R}^n$ *is* linearly independent *if the only scalars* $c_1, \ldots, c_k \in \mathbf{R}$ *that satisfy*

$$c_1 \mathbf{v}_1 + \cdots + c_k \mathbf{v}_k = \mathbf{0}$$

*are* $c_1 = \cdots = c_k = 0$. *Otherwise,* $\mathbf{v}_1, \ldots, \mathbf{v}_k$ *are said to be* linearly dependent.

The condition in this definition is sometimes stated sans equation as '*the only linear combination of* $\mathbf{v}_1, \ldots, \mathbf{v}_k$ *that vanishes is the trivial one.*'

**Definition 5.2.** *Let* $V \subset \mathbf{R}^n$ *be a subspace. A finite sequence* $\mathbf{v}_1, \ldots, \mathbf{v}_k \in V$ *is called a* basis *for* $V$ *if it is linearly independent and spans* $V$.

**Proposition 5.3.** *Let* $V \subset \mathbf{R}^n$ *be a subspace with basis* $\mathbf{v}_1, \ldots, \mathbf{v}_k \in V$. *Then for each* $\mathbf{w} \in V$, *there are unique scalars* $c_1, \ldots, c_k$ *such that*

$$\mathbf{w} = c_1 \mathbf{v}_1 + \cdots + c_k \mathbf{v}_k.$$

The scalars $c_1, \ldots, c_k$ in this proposition are called the *coordinates* of $\mathbf{w}$ with respect to $\mathbf{v}_1, \ldots, \mathbf{v}_k$.

Note that the trivial subspace has no basis. The next theorem and its corollaries show that non-trivial subspaces always have bases and that, while there might (in fact, certainly will) be more than one basis for a given non-trivial subspace, any two such bases must have the same size.

**Theorem 5.4.** *Let* $V \subset \mathbf{R}^n$ *be a subspace. If* $\mathbf{v}_1, \ldots, \mathbf{v}_k \in V$ *are linearly independent vectors and* $\mathbf{w}_1, \ldots, \mathbf{w}_\ell \in V$ *span* $V$, *then* $k \leq \ell$.

I include a proof for this one, because I like my way of saying it better than Shifrin's.

*Proof.* Let $A = \begin{bmatrix} \mathbf{v}_1 & \ldots & \mathbf{v}_k \end{bmatrix} \in \mathcal{M}_{n \times k}$ and $B = \begin{bmatrix} \mathbf{w}_1 & \ldots & \mathbf{w}_\ell \end{bmatrix} \in \mathcal{M}_{n \times \ell}$. By hypothesis, there is no non-trivial linear combination of $\mathbf{v}_1, \ldots, \mathbf{v}_k$ that vanishes. That is, there is no non-trivial solution $\mathbf{x} \in \mathbf{R}^k$ of the homogeneous linear system $A\mathbf{x} = \mathbf{0}$.

Similarly, the hypothesis $\mathbf{v}_j \in V = \operatorname{span}(\mathbf{w}_1, \ldots, \mathbf{w}_\ell)$ means that there exists a solution $\mathbf{y} = \mathbf{c}_j \in \mathbf{R}^\ell$ of $B\mathbf{y} = \mathbf{v}_j$. Let $C = \begin{bmatrix} \mathbf{c}_1 & \ldots & \mathbf{c}_k \end{bmatrix} \in \mathcal{M}_{\ell \times k}$. Then $BC = A$.

I claim that there is no non-trivial solution $\mathbf{x} \in \mathbf{R}^k$ of $C\mathbf{x} = \mathbf{0}$. Indeed, if there were, then we would have that $A\mathbf{x} = BC\mathbf{x} = \mathbf{0}$, so that $\mathbf{x}$ would also be a non-trivial solution of $A\mathbf{x} = \mathbf{0}$. The first paragraph rules this out, so my claim holds. That $C\mathbf{x} = \mathbf{0}$ has no non-trivial solution means that when I use row operations to put $C$ into reduced echelon form $\tilde{C}$, the resulting matrix will have no free variables. That is, it will have a pivot in each column. Since there is at most one pivot per row, $\tilde{C}$ must have at least as many rows as columns. In short $\ell \geq k$. $\qquad\qquad\square$

**Corollary 5.5.** *Any non-trivial subspace* $V \subset \mathbf{R}^n$ *has a basis.*

**Corollary 5.6.** *Any two bases for the same subspace* $V \subset \mathbf{R}^n$ *have the same number of vectors.*

**Definition 5.7.** *The* dimension *of a non-trivial subspace* $V \subset \mathbf{R}^n$ *is the number* $\dim V$ *of elements in a basis for* $V$. *The trivial subspace* $\{\mathbf{0}\} \subset \mathbf{R}^n$ *is said to have dimension* 0.

In particular, $\dim \mathbf{R}^n = n$.

**Corollary 5.8.** *If $V, W \subset \mathbf{R}^n$ are subspaces and $V \subset W$, then $\dim V \leq \dim W$. Equality holds if and only if $V = W$.*

**Theorem 5.9.** *Suppose that $A, \tilde{A} \in \mathcal{M}_{m \times n}$ are row equivalent matrices and $\tilde{A}$ is in reduced echelon form.*

- *The columns of $A$ corresponding to pivots of $\tilde{A}$ form a basis for $\operatorname{col} A$.*
- *The non-zero rows of $\tilde{A}$ form a basis for $\operatorname{row} A$.*
- *The solutions of $A\mathbf{x} = \mathbf{0}$ obtained by setting one free variable equal to $1$ and the others to $0$ form a basis for $\operatorname{nul} A$.*

*In particular, both $\dim \operatorname{row} A$ and $\dim \operatorname{col} A$ are equal to the number of pivots of $\tilde{A}$ and $\dim \operatorname{nul} A$ is equal to the number of free variables of $\tilde{A}$.*

In applying this theorem, be careful to note that you use columns of $A$ to get a basis for $\operatorname{col} A$ but rows of $\tilde{A}$ to get a basis for $\operatorname{row} A$. Indeed $\operatorname{row} A = \operatorname{row} \tilde{A}$ in general, but $\operatorname{col} A$ is usually not the same as $\operatorname{col} \tilde{A}$.

**Definition 5.10.** *The* rank *of a matrix is the dimension of its column space.*

**Corollary 5.11** (Rank Theorem)**.** *For any $m \times n$ matrix $A$, we have*

$$\dim \operatorname{col} A + \dim \operatorname{nul} A = n.$$

**Corollary 5.12.** *For any subspace $V \subset \mathbf{R}^n$ we have*

- $\dim V + \dim V^{\perp} = n$*;*
- $(V^{\perp})^{\perp} = V$*; and*
- *for any $\mathbf{x} \in \mathbf{R}^n$ there are unique vectors $\mathbf{x}_{\parallel} \in V$ and $\mathbf{x}_{\perp} \in V^{\perp}$ such that $\mathbf{x} = \mathbf{x}_{\parallel} + \mathbf{x}_{\perp}$.*

**Corollary 5.13.** *Given a matrix $A \in \mathcal{M}_{m \times n}$ and a vector $\mathbf{b} \in \mathbf{R}^m$, suppose that $A\mathbf{x} = \mathbf{b}$ is consistent. Then there is a unique solution $\mathbf{x} = \mathbf{x}_p \in \operatorname{row} A$, and the set of all solutions of $A\mathbf{x} = \mathbf{b}$ is given by*

$$\{\mathbf{x}_p + \mathbf{x}_h \in \mathbf{R}^n : \mathbf{x}_h \in \operatorname{nul} A\}.$$

When a linear system $A\mathbf{x} = \mathbf{b}$ fails to be consistent, there is a 'next best thing' to a solution.

**Definition 5.14.** *Given $A \in \mathcal{M}_{m \times n}$ and a vector $\mathbf{b} \in \mathbf{R}^m$, we call a vector $\mathbf{x} \in \mathbf{R}^n$ a* least squares solution *of $A\mathbf{x} = \mathbf{b}$ if the quantity $\|A\mathbf{x} - \mathbf{b}\|$ is minimal; i.e. if*

$$\|A\mathbf{x} - \mathbf{b}\| \leq \|A\tilde{\mathbf{x}} - \mathbf{b}\|$$

*for any (other) $\tilde{\mathbf{x}} \in \mathbf{R}^n$.*

Note that if $\mathbf{x}$ actually solves $A\mathbf{x} = \mathbf{b}$, then $\mathbf{x}$ is a least squares solution of $A\mathbf{x} = \mathbf{b}$, and that a least squares solution of $A\mathbf{x} = \mathbf{b}$ is an actual solution if and only if the minimum value of $\|A\mathbf{x} - \mathbf{b}\|$ is zero. Unlike ordinary solutions, least squares solutions *always* exist. Perhaps more surprisingly, there is a convenient way to find them by solving a different but related linear system.

**Corollary 5.15.** *For any $A \in \mathcal{M}_{m \times n}$ and $\mathbf{b} \in \mathbf{R}^m$, the linear system $A\mathbf{x} = \mathbf{b}$ has at least one least squares solution. Moreover, $\mathbf{x} \in \mathbf{R}^n$ is a least squares solution of $A\mathbf{x} = \mathbf{b}$ if and only if $\mathbf{x}$ satisfies*

$$A^T A \mathbf{x} = A^T \mathbf{b}.$$

The linear system $A^T A\mathbf{x} = A^T \mathbf{b}$ is called the 'normal equation' associated to $A\mathbf{x} = \mathbf{b}$. The corollary implies that the normal equation is always consistent. But its solution does not have to be unique.

**Proposition 5.16.** *The least squares solution of a linear system $A\mathbf{x} = \mathbf{b}$ is unique if and only if the nullspace of $A$ is trivial.*

## 6. Limits and Continuity

It is often the case that a non-linear function of $n$-variables $\mathbf{x} = (x_1, \ldots, x_n)$ is not really defined on all of $\mathbf{R}^n$. For instance $f(x_1, x_2) = \frac{x_1 x_2}{x_1^2 - x_2^2}$ is not defined when $x_1 = \pm x_2$. However, I will adopt a convention from the vector calculus notes of Jones and write $F : \mathbf{R}^n \to \mathbf{R}^m$ regardless, meaning only that the source of $F$ is some subset of $\mathbf{R}^n$. While a bit imprecise, this will not cause any big problems and will simplify many statements.

I will often distinguish between functions $f : \mathbf{R}^n \to \mathbf{R}$ that are *scalar-valued* and functions $F : \mathbf{R}^n \to \mathbf{R}^m$, $m \geq 2$ that are *vector-valued*, using lower-case letters to denote the former and upper case letters to denote the latter. Note that any vector-valued function $F : \mathbf{R}^n \to \mathbf{R}^m$ may be written $F = (F_1, \ldots, F_m)$ where $F_j : \mathbf{R}^n \to \mathbf{R}$ are scalar-valued functions called the *components* of $F$. For example, $F : \mathbf{R}^2 \to \mathbf{R}^2$ given by $F(x_1, x_2) = (x_1 x_2, x_1 + x_2)$ is a vector valued function with components $F_1(x_1, x_2) = x_1 x_2$ and $F_2(x_1, x_2) = x_1 + x_2$.

**Definition 6.1.** *Let $\mathbf{a} \in \mathbf{R}^n$ be a point and $r > 0$ be a positive real number. The* open ball *of radius $r$ about $\mathbf{a}$ is the set*

$$B(\mathbf{a}, r) := \{\mathbf{x} \in \mathbf{R}^n : \|\mathbf{x} - \mathbf{a}\| < r\}.$$

I will also use $B^*(\mathbf{a}, r)$ to denote the set of all $\mathbf{x} \in B(\mathbf{a}, r)$ except $\mathbf{x} = \mathbf{a}$. Extending my above convention, I will say that a function $F : \mathbf{R}^n \to \mathbf{R}^m$ is defined *near* a point $\mathbf{a} \in \mathbf{R}^n$ if there exists $r > 0$ such that $F(\mathbf{x})$ is defined for all points $\mathbf{x} \in B^*(\mathbf{a}, r)$, except possibly the center $\mathbf{x} = \mathbf{a}$. The following definition of 'limit' is one of the most important in all of mathematics. In differential calculus, it is the key to relating non-linear (i.e. hard) functions to linear (i.e. easier) functions.

**Definition 6.2.** *Suppose that $F : \mathbf{R}^n \to \mathbf{R}^m$ is a function defined near a point $\mathbf{a} \in \mathbf{R}^n$. We say that $F(\mathbf{x})$ has* limit $\mathbf{b} \in \mathbf{R}^m$ *as $\mathbf{x}$ approaches $\mathbf{a}$, i.e.*

$$\lim_{\mathbf{x} \to \mathbf{a}} F(\mathbf{x}) = \mathbf{b} \in \mathbf{R}^m,$$

*if for each $\epsilon > 0$ there exists $\delta > 0$ such that $0 < \|\mathbf{x} - \mathbf{a}\| < \delta$ implies $\|F(\mathbf{x}) - \mathbf{b}\| < \epsilon$.*

Notice that the final phrase in this definition can be written in terms of balls instead of magnitudes: *for any $\epsilon > 0$ there exists $\delta > 0$ such that $x \in B^*(\mathbf{a}, \delta)$ implies $F(\mathbf{x}) \in B(\mathbf{b}, \epsilon)$.*

A function might or might not have a limit as $\mathbf{x}$ approaches some given point $\mathbf{a}$, but it never has more than one.

**Proposition 6.3** (uniqueness of limits)**.** *If $F : \mathbf{R}^n \to \mathbf{R}^m$ is defined near $\mathbf{a} \in \mathbf{R}^n$, then there is at most one point $\mathbf{b} \in \mathbf{R}^m$ such that $\lim_{\mathbf{x} \to \mathbf{a}} F(\mathbf{x}) = \mathbf{b}$.*

**Definition 6.4.** *We say that a function $F : \mathbf{R}^n \to \mathbf{R}^m$ is* continuous *at $\mathbf{a} \in \mathbf{R}^n$ if $F$ is defined near and at $\mathbf{a}$ and*

$$\lim_{\mathbf{x} \to \mathbf{a}} F(\mathbf{x}) = F(\mathbf{a}).$$

*If $F$ is continuous at all points in its domain, we say simply that $F$ is continuous.*

**Proposition 6.5.** *The following are continuous functions.*

- *The constant function $F : \mathbf{R}^n \to \mathbf{R}^m$, given by $F(\mathbf{x}) = \mathbf{b}$ for some fixed $\mathbf{b} \in \mathbf{R}^m$ and all $\mathbf{x} \in \mathbf{R}^n$.*
- *The magnitude function $f : \mathbf{R}^n \to \mathbf{R}$ given by $f(\mathbf{x}) = \|\mathbf{x}\|$.*
- *The addition function $f : \mathbf{R}^2 \to \mathbf{R}$ given by $f(x_1, x_2) = x_1 + x_2$.*

- *The multiplication function $f : \mathbf{R}^2 \to \mathbf{R}$ given by $f(x_1, x_2) = x_1 x_2$.*
- *The reciprocal function $f : \mathbf{R} \to \mathbf{R}$ given by $f(x) = 1/x$.*

**Theorem 6.6.** *Linear transformations $T : \mathbf{R}^n \to \mathbf{R}^m$ are continuous.*

The remaining results in this section are aimed at allowing us to answer questions about limits without going all the way back to the definition of limit. The first result says that limits of vector-valued functions can always be reduced to questions about scalar-valued functions.

**Proposition 6.7.** *Suppose that $F : \mathbf{R}^n \to \mathbf{R}^m$ is a vector-valued function $F = (F_1, \ldots, F_m)$ defined near $\mathbf{a} \in \mathbf{R}^n$. Then the following are equivalent.*

(a) $\lim_{\mathbf{x} \to \mathbf{a}} F(\mathbf{x}) = \mathbf{b} \in \mathbf{R}^m$ .
(b) $\lim_{\mathbf{x} \to \mathbf{a}} \|F(\mathbf{x}) - \mathbf{b}\| = 0$.
(c) $\lim_{\mathbf{x} \to \mathbf{a}} F_j(\mathbf{x}) = b_j$ *for $1 \le j \le m$.*

The following theorem is sometimes paraphrased by saying that limits commute with continuous functions.

**Theorem 6.8** (limits commute with continuous functions)**.** *Let $F : \mathbf{R}^n \to \mathbf{R}^m$ and $G : \mathbf{R}^m \to \mathbf{R}^p$ be functions and $\mathbf{a} \in \mathbf{R}^n$, $\mathbf{b} \in \mathbf{R}^m$ be points such that $\lim_{\mathbf{x} \to \mathbf{a}} F(\mathbf{x}) = \mathbf{b}$ and $G$ is continuous at $\mathbf{b}$. Then*

$$\lim_{\mathbf{x} \to \mathbf{a}} G \circ F(\mathbf{x}) = G(\mathbf{b}).$$

**Corollary 6.9.** *Let $F : \mathbf{R}^n \to \mathbf{R}^m$ and $G : \mathbf{R}^m \to \mathbf{R}^p$ be continuous functions. Then $G \circ F$ is continuous.*

**Corollary 6.10.** *Let $f, g : \mathbf{R}^n \to \mathbf{R}$ be functions with limits $\lim_{\mathbf{x} \to \mathbf{a}} f(\mathbf{x}) = b$ and $\lim_{\mathbf{x} \to \mathbf{a}} g(\mathbf{x}) = c$ at some point $\mathbf{a} \in \mathbf{R}^n$. Then*

- $\lim_{\mathbf{x} \to \mathbf{a}} |f(\mathbf{x})| = |b|$.
- $\lim_{\mathbf{x} \to \mathbf{a}} f(\mathbf{x}) + g(\mathbf{x}) = b + c$;
- $\lim_{\mathbf{x} \to \mathbf{a}} f(\mathbf{x})g(\mathbf{x}) = bc$;
- $\lim_{\mathbf{x} \to \mathbf{a}} \frac{1}{f(\mathbf{x})} = \frac{1}{b}$, *provided $b \ne 0$.*

*Hence a sum or product of continuous functions is continuous, as is the reciprocal of a continuous function.*

Actually, the corollary extends to dot products, magnitudes and sums of vector-valued functions $F, G : \mathbf{R}^n \to \mathbf{R}^m$, too. I'll let you write down the statements of these facts.

When used with the fact that functions can't have more than one limit at a given point, Theorem 6.8 leads to a useful criterion for establishing that a limit *doesn't* exist.

**Definition 6.11.** *A parametrized curve is a continuous function $\gamma : \mathbf{R} \to \mathbf{R}^n$.*

**Corollary 6.12.** *Given a function $F : \mathbf{R}^n \to \mathbf{R}^m$ defined near a point $\mathbf{a} \in \mathbf{R}^n$, suppose that $\gamma_1, \gamma_2 : \mathbf{R} \to \mathbf{R}^n$ are parametrized curves such that $\gamma_1(t) = \gamma_2(t) = \mathbf{a}$ if and only if $t = 0$. If the limits $\lim_{t \to 0} F \circ \gamma_1(t)$ and $\lim_{t \to 0} F \circ \gamma_2(t)$ are not equal, then $\lim_{\mathbf{x} \to \mathbf{a}} F(\mathbf{x})$ does not exist.*

There is one more fact about limits that will prove useful for us.

**Theorem 6.13** (The Squeeze Theorem)**.** *Suppose that $F : \mathbf{R}^n \to \mathbf{R}^m$ and $g : \mathbf{R}^n \to \mathbf{R}$ are functions defined near $\mathbf{a} \in \mathbf{R}^n$. Suppose there exists $r > 0$ such that*

- $\|F(\mathbf{x})\| \leq |g(\mathbf{x})|$ *for all* $\mathbf{x} \in B(\mathbf{a}, r)$, *except possibly* $\mathbf{x} = \mathbf{a}$;
- $\lim_{\mathbf{x} \to \mathbf{a}} g(\mathbf{x}) = 0$.

*Then* $\lim_{\mathbf{x} \to \mathbf{a}} F(\mathbf{x}) = \mathbf{0}$.

## 7. Differentiability

Recall the definition of derivative from one variable calculus

**Definition 7.1.** *We say that $f : \mathbf{R} \to \mathbf{R}$ is* differentiable *at a point $a \in \mathbf{R}$ if the quantity*

$$f'(a) := \lim_{h \to 0} \frac{f(a + h) - f(a)}{h}$$

*exists. We then call $f'(a)$ the* derivative *of $f$ at $a$.*

One way to transfer this definition to higher dimensions is via 'directional' derivatives.

**Definition 7.2.** *The* directional derivative *of a function $F : \mathbf{R}^n \to \mathbf{R}^m$ at a point $\mathbf{a} \in \mathbf{R}^n$ in the direction $\mathbf{v} \in \mathbf{R}^n$ is the quantity (if it exists)*

$$D_{\mathbf{v}} F(\mathbf{a}) := \lim_{t \to 0} \frac{F(\mathbf{a} + t\mathbf{v}) - F(\mathbf{a})}{t}$$

*When $\mathbf{v} = \mathbf{e}_j$ is a standard basis vector, we write $\frac{\partial F}{\partial x_j}(\mathbf{a}) := D_{\mathbf{e}_j} F(a)$ and call this quantity the* partial derivative *of $F$ with respect to $x_j$.*

Another way of stating this definition is that $D_{\mathbf{v}} F(\mathbf{a}) = h'(0)$ where $h : \mathbf{R} \to \mathbf{R}^m$ is the composite function

$$h(t) := F(\mathbf{a} + t\mathbf{v})$$

obtained by restricting $F$ to the line through $\mathbf{a}$ in the direction $\mathbf{v}$. This way of formulating directional derivatives is quite useful when you actually have to compute one!

A shortcoming of directional derivatives is that they don't always do a very good job of controlling the behavior of $F$ near a given point $\mathbf{a}$ (see Section 3.1 e.g. 2 in Shifrin for a good illustration of this). One needs a little bit more restrictive notion of derivative in order to guarantee this sort of control.

**Definition 7.3.** *We say that a function $F : \mathbf{R}^n \to \mathbf{R}^m$ is* differentiable *at a point $\mathbf{a} \in \mathbf{R}^n$ if there exists a linear transformation $T : \mathbf{R}^n \to \mathbf{R}^m$ such that*

$$(1) \qquad \lim_{\mathbf{h} \to \mathbf{0}} \frac{F(\mathbf{a} + \mathbf{h}) - F(\mathbf{a}) - T\,\mathbf{h}}{\|\mathbf{h}\|} = \mathbf{0}.$$

*If such a $T$ exists, then we call it the* derivative *of $F$ at $\mathbf{a}$ write $DF(\mathbf{a}) := T$.*

So under this definition, the derivative $DF(\mathbf{a})$ of $F$ at $\mathbf{a}$ is not a number but rather a linear transformation. This is not so strange if you remember any linear transformation $T : \mathbf{R}^n \to \mathbf{R}^m$ has a standard matrix $A \in \mathcal{M}_{m \times n}$, so you can think of the derivative of $F$ at $\mathbf{a}$ more concretely as a matrix, i.e. as a collection of $mn$ numbers that describe the way all the different components of $F = (F_1, \ldots, F_m)$ are changing in all the different directions one can approach $\mathbf{a}$. I'm sort of doing that already when I suppress parentheses in $T(\mathbf{h})$ and write $T\mathbf{h}$ instead.

In particular, if $f : \mathbf{R} \to \mathbf{R}$ is just a scalar function of a single variable, then the number $f'(a)$ above is just the lone entry in the $1 \times 1$ matrix for the linear transformation $T : \mathbf{R} \to \mathbf{R}$ given by $T(h) = f'(a)h$.

Note that Equation (1) can be written in several slightly different but equivalent ways. For instance, one could take the magnitude of the numerator and write instead (I'll use $DF(\mathbf{a})$ in place of $T$ now).

$$\lim_{\mathbf{h} \to \mathbf{0}} \frac{\|F(\mathbf{a} + \mathbf{h}) - F(\mathbf{a}) - DF(\mathbf{a})\mathbf{h}\|}{\|\mathbf{h}\|} = 0.$$

Or one could set $\mathbf{x} := \mathbf{a} + \mathbf{h}$ and rewrite the limit as

$$\lim_{\mathbf{x} \to \mathbf{a}} \frac{F(\mathbf{x}) - F(\mathbf{a}) - DF(\mathbf{a})(\mathbf{x} - \mathbf{a})}{\|\mathbf{x} - \mathbf{a}\|}.$$

Another very useful way to restate (1) is to say that

$$F(\mathbf{a} + \mathbf{h}) = F(\mathbf{a}) + DF(\mathbf{a})\mathbf{h} + E(\mathbf{h}),$$

where the 'error term' $E(\mathbf{h})$ satisfies $\lim_{\mathbf{h} \to \mathbf{0}} \frac{\|E(\mathbf{h})\|}{\|\mathbf{h}\|} = 0$.

The first result of this section indicates that differentiability of $F$ at $\mathbf{a}$ gives us some control of $F$ at nearby values of $\mathbf{a}$.

**Theorem 7.4.** *If $F : \mathbf{R}^n \to \mathbf{R}^m$ is differentiable at $\mathbf{a}$, then $F$ is continuous at $\mathbf{a}$.*

The second fact about our new notion of derivative $DF(\mathbf{a})$ is that it's not that far from partial and directional derivatives.

**Proposition 7.5.** *Suppose that $F : \mathbf{R}^n \to \mathbf{R}^m$ is differentiable at a point $\mathbf{a} \in \mathbf{R}^n$. Then the directional derivative of $F$ at $\mathbf{a}$ in direction $\mathbf{v} \in \mathbf{R}^n$ exists and is given by*

$$D_{\mathbf{v}}F(\mathbf{a}) = DF(\mathbf{a})\mathbf{v}.$$

*In particular, the matrix for the linear transformation $DF(\mathbf{a}) : \mathbf{R}^n \to \mathbf{R}^m$ is given column-wise by*

$$\begin{bmatrix} \frac{\partial F}{\partial x_1}(\mathbf{a}) & \cdots & \frac{\partial F}{\partial x_n}(\mathbf{a}) \end{bmatrix}$$

Among other things, this proposition tells us that there is only one candidate for $DF(\mathbf{a})$ and gives us a practical means for finding out what it is (by taking partial derivatives). It does not, however, tell us how to determine whether our candidate is a winner, i.e. whether $F$ is actually differentiable at $\mathbf{a}$. For most purposes, the following condition suffices for that purpose.

**Definition 7.6.** *A function $F : \mathbf{R}^n \to \mathbf{R}^m$ is said to be* continuously differentiable *at $\mathbf{a} \in \mathbf{R}^n$ if all partial derivatives $\frac{\partial F}{\partial x_j}$ of $F$ (exist and) are continuous at $\mathbf{a}$. If $F$ is continuously differentiable at all points in its domain, they we say simply that '$F$ is continuously differentiable' (without reference to any point).*

Continuously differentiable functions are often (alternatively) called '$C^1$ functions'.

**Theorem 7.7.** *If $F : \mathbf{R}^n \to \mathbf{R}^m$ is continuously differentiable at $\mathbf{a}$, then $F$ is differentiable at $\mathbf{a}$.*

The proof of this theorem depends on two further results, both of some interest in their own right. The first reduces the problem of differentiability for vector-valued functions to the simpler case of scalar-valued functions.

**Proposition 7.8.** *Let $F : \mathbf{R}^n \to \mathbf{R}^m$ be a vector-valued function with (scalar-valued) components $F_i : \mathbf{R}^n \to \mathbf{R}$, $1 \le i \le m$. Then $F$ is differentiable at $\mathbf{a} \in \mathbf{R}^n$ if and only if all the components $F_i$ are differentiable at $\mathbf{a}$. Moreover, the $i$th row of the standard matrix for $DF(\mathbf{a})$ is equal to the standard matrix for $DF_i(\mathbf{a})$.*

The second is a (the?) central result from one variable calculus.

**Theorem 7.9** (Mean Value Theorem). *Suppose that $f : (a, b) \to \mathbf{R}$ is a differentiable function on an open interval $(a, b)$ and $x, y \in (a, b)$ are two different points. Then there is a number $c$ between $x$ and $y$ such that*

$$f'(c) = \frac{f(y) - f(x)}{y - x}.$$

We will use the mean value theorem again later (e.g. to show equality of mixed partial derivatives).

## 8. More about derivatives

In one variable calculus, there are many facts about differentiation that allow one to compute derivatives algebraically, without resorting to the actual limit definition. There are similar facts in multi-variable calculus.

**Theorem 8.1.** *Suppose that $F, G : \mathbf{R}^n \to \mathbf{R}^m$ are vector-valued functions and $f, g : \mathbf{R}^n \to \mathbf{R}$ are scalar-valued functions, all differentiable at some point $\mathbf{a} \in \mathbf{R}^n$. Then*

(a) *$F + G$ is differentiable at $\mathbf{a}$ and $D(F + G)(\mathbf{a}) = DF(\mathbf{a}) + DG(\mathbf{a})$;*
(b) *$fg$ is differentiable at $\mathbf{a}$ and $D(fg)(\mathbf{a}) = f(\mathbf{a})Dg(\mathbf{a}) + g(\mathbf{a})Df(\mathbf{a})$;*
(c) *$F \cdot G$ is differentiable at $\mathbf{a}$ and $D(F \cdot G)(\mathbf{a})\mathbf{h} = G(\mathbf{a}) \cdot (DF(\mathbf{a})\mathbf{h}) + (DG(\mathbf{a})\mathbf{h}) \cdot F(\mathbf{a})$.*

Note that this is almost the same as Proposition 3.3.1 in Shifrin, but the second item here is a bit different (simpler and less general) than in Shifrin. These facts get used less often in multivariable calculus than they do in one variable calculus, but they are occasionally useful. The *really* important fact is the next one. Stating it in a straightforward and general way depends very much on using total rather than partial derivatives.

**Theorem 8.2** (Chain Rule). *Suppose that $G : \mathbf{R}^n \to \mathbf{R}^m$ and $F : \mathbf{R}^m \to \mathbf{R}^\ell$ are functions such that $G$ is differentiable at $\mathbf{a} \in \mathbf{R}^n$ and $F$ is differentiable at $G(\mathbf{a}) \in \mathbf{R}^m$. Then $F \circ G$ is differentiable at $\mathbf{a} \in \mathbf{R}^n$ and*

$$D(F \circ G)\mathbf{a} = DF(G(\mathbf{a})) \circ DG(\mathbf{a}).$$

In other words, the derivative of a composition is the composition (or product, if you think in matrix terms) of the derivatives. To see why this should be true on an intuitive level, one should think in terms of linear approximations. That is, for $\mathbf{x} \in \mathbf{R}^n$ near $\mathbf{a}$ we have

$$\mathbf{y} := G(\mathbf{x}) \approx G(\mathbf{a}) + DG(\mathbf{a})(\mathbf{x} - \mathbf{a}),$$

and for $\mathbf{y} \in \mathbf{R}^m$ near $G(\mathbf{a})$ we have

$$\mathbf{z} := F(\mathbf{y}) \approx F(G(\mathbf{a})) + DF(G(\mathbf{a}))(\mathbf{y} - G(\mathbf{a})).$$

So taking $\mathbf{y} = G(\mathbf{x})$, we know that $\mathbf{y}$ is close to $G(\mathbf{a})$ when $\mathbf{x}$ is close to $\mathbf{a}$ (why?). Therefore

$$
\begin{aligned}
F \circ G(\mathbf{x}) &\approx F(G(\mathbf{a})) + DF(G(\mathbf{a}))(G(\mathbf{x}) - G(\mathbf{a})) \\
&\approx F(G(\mathbf{a})) + DF(G(\mathbf{a}))DG(\mathbf{a})(\mathbf{x} - \mathbf{a}).
\end{aligned}
$$

Note for the final approximation, I have simply replaced $G(\mathbf{x})$ with its linear approximation and the $G(\mathbf{a})$ terms cancelled. Anyhow, it's easy to believe that this last expression should be the linear approximation of $F \circ G$ near $\mathbf{a}$. The first term is clearly the constant term, and the second term therefore corresponds to the derivative of $F \circ G$ at $\mathbf{a}$.

**Definition 8.3.** *A function $f : \mathbf{R}^n \to \mathbf{R}$ has a* local maximum *at $\mathbf{a} \in \mathbf{R}^n$ if there exists $\delta > 0$ such that for all $\mathbf{x} \in B(\mathbf{a}, \delta)$, $f$ is defined at $\mathbf{x}$ and $f(\mathbf{x}) \leq f(\mathbf{a})$.*

**Definition 8.4.** *A point $\mathbf{a} \in \mathbf{R}^n$ is* critical *for a scalar-valued funciton $f : \mathbf{R}^n \to \mathbf{R}$ if $f$ is differentiable at $\mathbf{a}$ and $\nabla f(\mathbf{a}) = \mathbf{0}$.*

**Proposition 8.5** (First derivative test). *If a scalar-valued function $f : \mathbf{R}^n \to \mathbf{R}$ has a local extremum at $\mathbf{a} \in \mathbf{R}^n$ and $f$ is differentiable at $\mathbf{a}$, then $\mathbf{a}$ is a critical point for $f$.*

Sometimes one is interested in finding maxima and minima of a function not on all of $\mathbf{R}^n$ but just along the level set of another function. For instance, instead of wondering where the hottest place in the universe is, one might ask where the hottest place is on the surface of the earth. So mathematically the question is to figure out where the temperature function is maximal on a big sphere.

**Definition 8.6.** *Let* $f, g : \mathbf{R}^n \to \mathbf{R}$ *be scalar-valued functions and* $c \in \mathbf{R}$ *be a number. We say that* $f$ *has a local maximum at* $\mathbf{a} \in \mathbf{R}^n$ *subject to the constraint* $g = c$ *if*

- $g(\mathbf{a}) = c$; *and*
- *there exists* $\delta > 0$ *such that* $f(\mathbf{x}) \leq f(\mathbf{a})$ *for all* $\mathbf{x} \in B(\mathbf{a}, \delta)$ *such that* $g(\mathbf{x}) = c$.

The 'method of Lagrange multipliers' is a proceedure for finding constrained local maxima and minima of differentiable functions.

**Theorem 8.7** (First derivative test, with one constraint)**.** *Suppose that* $f, g : \mathbf{R}^n \to \mathbf{R}$ *are scalar-valued functions and that* $f$ *has a local maximum at* $\mathbf{a} \in \mathbf{R}^n$ *subject to the constraint* $g = c$*. If* $f$ *and* $g$ *are* $C^1$ *at* $\mathbf{a}$*, then the gradients* $\nabla f(\mathbf{a})$ *and* $\nabla g(\mathbf{a})$ *are parallel.*

Shifrin presents a more general version of this theorem, involving multiple constraints (or rather, a vector-valued constraint function $g$).

## 9. Second order derivatives

Here I limit the discussion to scalar-valued functions.

**Definition 9.1.** *A function $f : \mathbf{R}^n \to \mathbf{R}$ is $C^2$ at $\mathbf{a} \in \mathbf{R}^n$ if all second order partial derivatives exist and are continuous at $\mathbf{a} \in \mathbf{R}^n$. If $f$ is $C^2$ at all points in its domain, then we say simply that '$f$ is $C^2$.'*

It is implicit in this definition that $f$ is twice differentiable at points near $\mathbf{a}$; otherwise it wouldn't make sense to say that second order partial derivatives are continuous at $\mathbf{a}$. It is also implicit that $f$ and it's first order partial derivatives exist and are continuous near $\mathbf{a}$; otherwise we couldn't take second order partial derivatives.

The first main fact about second order partial derivatives is that the order of differentiation is irrelevant.

**Theorem 9.2.** *Suppose that $f : \mathbf{R}^n \to \mathbf{R}$ is $C^2$ at $\mathbf{a}$. Then for any indices $1 \le i, j \le n$, we have*

$$\frac{\partial^2 f}{\partial x_i \, \partial x_j}(\mathbf{a}) = \frac{\partial^2 f}{\partial x_j \, \partial x_i}(\mathbf{a})$$

**Definition 9.3.** *Suppose that $f : \mathbf{R}^n \to \mathbf{R}$ is $C^2$ at $\mathbf{a} \in \mathbf{R}^n$. The* Hessian *of $f$ at $\mathbf{a}$ is the $n \times n$ matrix $Hf(\mathbf{a})$ with $ij$-entry equal to $\frac{\partial^2 f}{\partial x_i \, \partial x_j}(\mathbf{a})$.*

The previous theorem says that $Hf(\mathbf{a})$ is a *symmetric* matrix; i.e. it is equal to its own transpose. The second main fact about second order partial derivatives is that they afford extra control over the local behavior of a function.

**Theorem 9.4.** *Suppose that $f : \mathbf{R}^n \to \mathbf{R}$ is $C^2$ at $\mathbf{a}$. Then*

$$f(\mathbf{a} + \mathbf{h}) = f(\mathbf{a}) + \nabla f(\mathbf{a}) \cdot \mathbf{h} + \frac{1}{2}\mathbf{h}^T Hf(\mathbf{a})\mathbf{h} + E_2(\mathbf{h}),$$

*where $\lim_{\mathbf{h} \to \mathbf{0}} \frac{E_2(\mathbf{h})}{\|\mathbf{h}\|^2} = \mathbf{0}$.*

The *quadratic form* associated to a symmetric $n \times n$ matrix $A$ is the function $Q : \mathbf{R}^n \to \mathbf{R}$ given by

$$Q(\mathbf{x}) = \mathbf{x}^T A \mathbf{x}.$$

Note that $Q(c\mathbf{x}) = c^2 Q(\mathbf{x})$ for any scalar $c \in \mathbf{R}$.

**Definition 9.5.** *Let $A \in \mathcal{M}_{n \times n}$ be a symmetric square matrix and $Q : \mathbf{R}^n \to \mathbf{R}$ be the associated quadratic form. We say that*

- *$A$ is positive definite if $Q(\mathbf{x}) > 0$ for all non-zero $\mathbf{x} \in \mathbf{R}^n$;*
- *$A$ is negative definite if $Q(\mathbf{x}) < 0$ for all non-zero $\mathbf{x} \in \mathbf{R}^n$;*
- *$A$ is indefinite if there exist $\mathbf{x}, \mathbf{y} \in \mathbf{R}^n$ such that $Q(\mathbf{x}) < 0 < Q(\mathbf{y})$.*

A symmetric matrix can satisfy at most one of these three conditions, but it's hard to tell just by looking which, if any, holds for a given matrix. For $2 \times 2$ matrices, there is a fairly convenient condition one can apply.

**Theorem 9.6.** *A $2 \times 2$ symmetric matrix $A = \begin{bmatrix} a & b \\ b & c \end{bmatrix}$ is*

- *positive definite if and only if $a > 0$ and $ac > b^2$;*
- *negative definite if and only if $a < 0$ and $ac > b^2$;*

- *indefinite if and only if $ac < b^2$.*

For larger square matrices, there are ways to check definiteness by computing eigenvalues or by computing determinants of 'diagonal minors'. This requires more linear algebra than I am assuming at present, so I do not discuss these things here.

**Theorem 9.7** (Second derivative test). *Suppose that $f : \mathbf{R}^n \to \mathbf{R}$ is $C^2$ at some critical point $\mathbf{a} \in \mathbf{R}^n$ for $f$. Then $f$ has*

- *a local minimum at $\mathbf{a}$ if $Hf(\mathbf{a})$ is positive definite;*
- *a local maximum at $\mathbf{a}$ if $Hf(\mathbf{a})$ is negative definite;*
- *neither a local maximum nor a local minimum if $Hf(\mathbf{a})$ is indefinite.*

## 10. Topology of $\mathbf{R}^n$

In one variable calculus, one can (mostly) limit attention to functions whose domains are open or closed intervals in $\mathbf{R}$. To do multivariable calculus one needs a more general notion of open and closed. Roughly speaking, a closed set is one that contains its edge and an open set is one that omits its edge.

**Definition 10.1.** *Let $X \subset \mathbf{R}^n$ be a set. We say that $\mathbf{a} \in \mathbf{R}^n$ is*

- *an* interior point *of $X$ if there exists $\delta > 0$ such that $B(\mathbf{a}, \delta) \subset X$;*
- *an* exterior point *of $X$ if there exists $\delta > 0$ such that $B(\mathbf{a}, \delta) \cap X = \emptyset$; and*
- *a* boundary point *of $X$ if for any $\delta > 0$, the ball $B(\mathbf{a}, \delta)$ intersects both $X$ and $\mathbf{R}^n \backslash X$.*

Note that $\mathbf{a}$ is an exterior point of $X$ if and only if $\mathbf{a}$ is an interior point of the complement $\mathbf{R}^n \setminus X$; and $\mathbf{a}$ is a boundary point of $X$ if and only if $\mathbf{a}$ is interior to neither $X$ nor $\mathbf{R}^n \setminus X$. In particular, boundary points of $X$ coincide with boundary points of $\mathbf{R}^n \setminus X$, and any point $\mathbf{a} \in \mathbf{R}^n$ is exactly one of the three types (interior, exterior or boundary) relative to $X$.

**Definition 10.2.** *A set $X \subset \mathbf{R}^n$ is*

- *open if every $\mathbf{a} \in X$ is an interior point of $X$;*
- *closed if every boundary point of $X$ is contained in $X$.*

Note that $\mathbf{R}^n$ is both open and closed. So is the empty set. A set $\{\mathbf{a}\}$ containing a single point $\mathbf{a} \in \mathbf{R}^n$ is closed but not open. A (non-empty) ball $B(\mathbf{a}, r) \subset \mathbf{R}^n$ is open but not closed.

It is more or less immediate from definitions that

**Proposition 10.3.** *$X \subset \mathbf{R}^n$ is open if and only if the complement $\mathbf{R}^n \setminus X$ is closed, and vice versa.*

Open-ness and closed-ness interact well with other set operations.

**Proposition 10.4.** *Both the union and intersection of finitely many open subsets of $\mathbf{R}^n$ are open. Similarly, unions and intersections of finitely many closed subsets are closed.*

Many other open and closed sets are furnished by the following useful result.

**Proposition 10.5.** *Suppose $f : \mathbf{R}^n \to \mathbf{R}$ is continuous. Then for any $c \in \mathbf{R}$, the 'sub-level set'*

$$\{f < c\} := \{\mathbf{x} \in \mathbf{R}^n : f(\mathbf{x}) < c\}$$

*is open. So are the sets $\{f > c\}$ and $\{f \neq c\}$. The sets $\{f \leq c\}$, $\{f \geq c\}$ and $\{f = c\}$ are all closed.*

**Definition 10.6.** *A set $X \subset \mathbf{R}^n$ is* bounded *if there exists $R > 0$ such that $X \subset B(\mathbf{0}, R)$.*

**Definition 10.7.** *A subset of $\mathbf{R}^n$ is* compact *if it is closed and bounded.*

**Theorem 10.8** (Extreme Value Theorem)**.** *Suppose that $f : \mathbf{R}^n \to \mathbf{R}$ is continuous and that $K \subset \mathbf{R}^n$ is a compact subset of the domain of $f$. Then there exist $\mathbf{a}, \mathbf{b} \in K$ such that*

$$f(\mathbf{a}) \leq f(\mathbf{x}) \leq f(\mathbf{b}).$$

*for all $\mathbf{x} \in K$.*