



in planar transformations. In Section 6 we propose an approach to derivatives which makes more natural the generalization to derivatives as linear functions represented by Jacobian matrices. The proposed definition has the advantage of reducing the process of evaluating limits to evaluating zero limits and avoids the restriction of  $h$  to nonzero values. An added advantage is that the equation of the tangent line (or hyperplanes in higher dimensions) appears naturally as the linear approximation. A further benefit is that the proof of the chain rule is unchanged in higher dimensions.

## 2 Linear functions

A function  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$  is said to be *linear* if

1.  $f(\mathbf{x} + \mathbf{y}) = f(\mathbf{x}) + f(\mathbf{y})$  for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ ; and
2.  $f(c\mathbf{x}) = cf(\mathbf{x})$  for all  $\mathbf{x} \in \mathbb{R}^n$  and  $c \in \mathbb{R}$ .

Equivalently,  $f$  is linear if it satisfies  $f(c_1\mathbf{x} + c_2\mathbf{y}) = c_1f(\mathbf{x}) + c_2f(\mathbf{y})$  for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$  and  $c_1, c_2 \in \mathbb{R}$ . We shall see how matrices are used to represent such linear functions. Linear functions are also called *linear transformations* and we shall use these terms interchangeably.

Before proceeding to the main part of our discussion we remark that the term ‘linear’ is used in many contexts and can be confusing to students. Here are some examples:

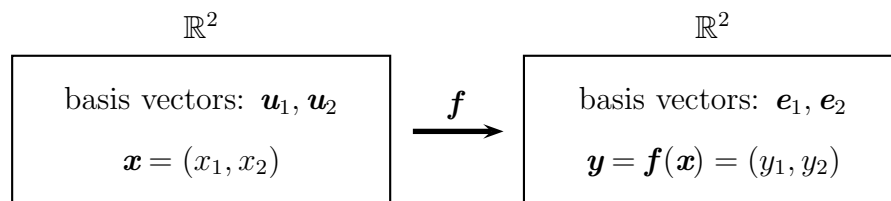
1. The polynomial  $p(x) = ax + b$ , although referred to as ‘linear’ is not a linear function in the above sense if  $b \neq 0$ .
2. A real-valued function of two variables,  $f(x, y) = ax + by + c$ , is not a linear function if  $c \neq 0$ . We do say ‘ $f$  is linear’ in  $x$  and  $y$ . The terms ‘homogeneous’ and ‘inhomogeneous’ are also used.

A property of linear functions,  $f(\mathbf{0}) = \mathbf{0}$ , becomes important in identifying when a function  $f$  is linear. Hence a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  is linear if its straight-line graph goes through the origin; a function  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$  is linear if its graph, a plane in  $\mathbb{R}^3$ , goes through the origin.

## 3 Planar transformations

In the first year it is sufficient to illustrate the ideas involved with transformations on the plane. Let the basis vectors in the planes of the domain and co-domain be denoted by  $\mathbf{u}_1, \mathbf{u}_2$  and  $\mathbf{e}_1, \mathbf{e}_2$  respectively. If  $\mathbf{x} = (x_1, x_2)$  is a point in the domain-plane  $\mathbb{R}^2$ , then we can write

$$\mathbf{x} = x_1\mathbf{u}_1 + x_2\mathbf{u}_2. \tag{3.1}$$



If  $f: \mathbb{R}^2 \rightarrow \mathbb{R}^2$  is linear then, by definition,

$$f(\mathbf{x}) = x_1f(\mathbf{u}_1) + x_2f(\mathbf{u}_2). \tag{3.2}$$

Hence knowledge of  $f(\mathbf{u}_1)$  and  $f(\mathbf{u}_2)$  completely determines how  $\mathbf{x}$  transforms under  $f$ . Since  $f(\mathbf{u}_j)$  is in the co-domain  $\mathbb{R}^2$ , we have the linear combinations

$$f(\mathbf{u}_1) = a_{11}\mathbf{e}_1 + a_{21}\mathbf{e}_2, \quad f(\mathbf{u}_2) = a_{12}\mathbf{e}_1 + a_{22}\mathbf{e}_2, \tag{3.3}$$

where the  $a_{ij}$  are real numbers. The  $a_{ij}$ 's completely determine the transformation  $f$  since, by (3.1), (3.2) and (3.3), we have

$$\begin{aligned} \mathbf{f}(\mathbf{x}) &= x_1(a_{11}\mathbf{e}_1 + a_{21}\mathbf{e}_2) + x_2(a_{12}\mathbf{e}_1 + a_{22}\mathbf{e}_2) = (a_{11}x_1 + a_{12}x_2)\mathbf{e}_1 + (a_{21}x_1 + a_{22}x_2)\mathbf{e}_2 \\ &= y_1\mathbf{e}_1 + y_2\mathbf{e}_2, \quad \text{where } y_i = a_{i1}x_1 + a_{i2}x_2. \end{aligned} \quad (3.4)$$

How do we represent this transformation?

If the standard basis in each plane are represented by column vectors,

$$\mathbf{u}_1 = \mathbf{e}_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad \mathbf{u}_2 = \mathbf{e}_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix},$$

then points in  $\mathbb{R}^2$  are also represented by column vectors; (3.1) becomes

$$\mathbf{x} = x_1 \begin{bmatrix} 1 \\ 0 \end{bmatrix} + x_2 \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} x_1 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ x_2 \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix},$$

and, similarly, (3.3) becomes

$$\mathbf{f}(\mathbf{u}_1) = \begin{bmatrix} a_{11} \\ a_{21} \end{bmatrix}, \quad \mathbf{f}(\mathbf{u}_2) = \begin{bmatrix} a_{12} \\ a_{22} \end{bmatrix}.$$

Hence, by (3.4),

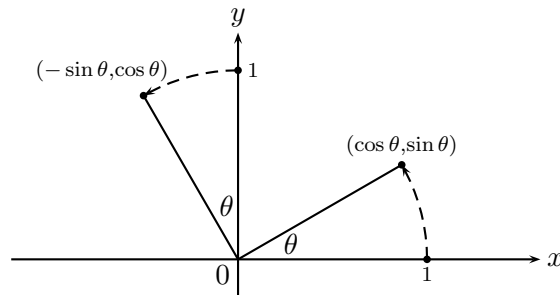
$$\mathbf{f}(\mathbf{x}) = \begin{bmatrix} a_{11}x_1 + a_{12}x_2 \\ a_{21}x_1 + a_{22}x_2 \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}.$$

Thus  $\mathbf{f}$  is represented by the matrix

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix},$$

and the value of  $\mathbf{f}$  at  $\mathbf{x}$  is represented by a matrix-vector product  $A\mathbf{x}$ . Note that the first column of  $A$  is  $\mathbf{f}(\mathbf{u}_1)$  while the second column is  $\mathbf{f}(\mathbf{u}_2)$ .

**Example 1** Determine the matrix representing a rotation through angle  $\theta$  anticlockwise centred at the origin. Now from the diagram (showing both sets of vectors on the same axes)



we see that

$$\begin{bmatrix} 1 \\ 0 \end{bmatrix} \mapsto \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix} \quad \text{and} \quad \begin{bmatrix} 0 \\ 1 \end{bmatrix} \mapsto \begin{bmatrix} -\sin \theta \\ \cos \theta \end{bmatrix}.$$

Hence the rotation matrix is

$$\begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}.$$

□

## 4 Composition as matrix multiplication

The algebraic rules for matrices and vectors follow from the rules for the corresponding linear functions and linear vector space. We now establish an important link between the composition of linear functions and matrix multiplication of their corresponding matrix representations. The key result is expressed in the following theorem.

**Theorem 1** *If  $\mathbf{f}: \mathbb{R}^2 \rightarrow \mathbb{R}^2$  and  $\mathbf{g}: \mathbb{R}^2 \rightarrow \mathbb{R}^2$  are linear transformations represented by the matrices  $A$  and  $B$  respectively, then the composition  $\mathbf{g} \circ \mathbf{f}$  is linear and is represented by the matrix product  $BA$ .*

**Proof:**  $\mathbf{g} \circ \mathbf{f}$  is linear since for all  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^2$  and  $c_1, c_2 \in \mathbb{R}$ , we have

$$\begin{aligned} \mathbf{g} \circ \mathbf{f}(c_1 \mathbf{x} + c_2 \mathbf{y}) &= \mathbf{g}(\mathbf{f}(c_1 \mathbf{x} + c_2 \mathbf{y})) && \text{(by definition of composition)} \\ &= \mathbf{g}(c_1 \mathbf{f}(\mathbf{x}) + c_2 \mathbf{f}(\mathbf{y})) && \text{(by linearity of } \mathbf{f} \text{)} \\ &= c_1 \mathbf{g}(\mathbf{f}(\mathbf{x})) + c_2 \mathbf{g}(\mathbf{f}(\mathbf{y})) && \text{(by linearity of } \mathbf{g} \text{)} \\ &= c_1 \mathbf{g} \circ \mathbf{f}(\mathbf{x}) + c_2 \mathbf{g} \circ \mathbf{f}(\mathbf{y}). \end{aligned}$$

Now let  $\mathbf{x} = [x_1, x_2]^T \in \mathbb{R}^2$  and  $\mathbf{f}(\mathbf{x}) = \mathbf{y} = [y_1, y_2]^T \in \mathbb{R}^2$ . From Section 3 we have

$$\begin{aligned} \mathbf{y} &= \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} a_{11}x_1 + a_{12}x_2 \\ a_{21}x_1 + a_{22}x_2 \end{bmatrix}, \\ \mathbf{g} \circ \mathbf{f}(\mathbf{x}) &= \mathbf{g}(\mathbf{f}(\mathbf{x})) = \mathbf{g}(\mathbf{y}) = B\mathbf{y} \\ &= \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} b_{11}y_1 + b_{12}y_2 \\ b_{21}y_1 + b_{22}y_2 \end{bmatrix} = \begin{bmatrix} b_{11}(a_{11}x_1 + a_{12}x_2) + b_{12}(a_{21}x_1 + a_{22}x_2) \\ b_{21}(a_{11}x_1 + a_{12}x_2) + b_{22}(a_{21}x_1 + a_{22}x_2) \end{bmatrix} \\ &= \begin{bmatrix} (b_{11}a_{11} + b_{12}a_{21})x_1 + (b_{11}a_{12} + b_{12}a_{22})x_2 \\ (b_{21}a_{11} + b_{22}a_{21})x_1 + (b_{21}a_{12} + b_{22}a_{22})x_2 \end{bmatrix} = \begin{bmatrix} b_{11}a_{11} + b_{12}a_{21} & b_{11}a_{12} + b_{12}a_{22} \\ b_{21}a_{11} + b_{22}a_{21} & b_{21}a_{12} + b_{22}a_{22} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}. \end{aligned}$$

Hence the matrix representing  $\mathbf{g} \circ \mathbf{f}$  is

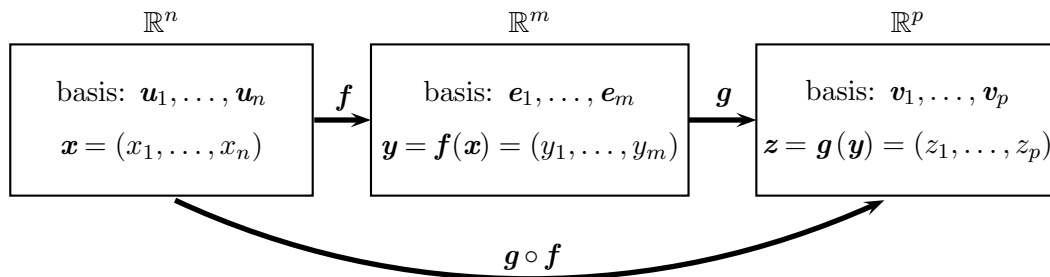
$$\begin{bmatrix} b_{11}a_{11} + b_{12}a_{21} & b_{11}a_{12} + b_{12}a_{22} \\ b_{21}a_{11} + b_{22}a_{21} & b_{21}a_{12} + b_{22}a_{22} \end{bmatrix} = \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} = BA. \quad \square$$

**Remarks:** In the general case where  $\mathbf{f}: \mathbb{R}^n \rightarrow \mathbb{R}^m$  is represented by an  $m \times n$  matrix  $A$ , and  $\mathbf{g}: \mathbb{R}^m \rightarrow \mathbb{R}^p$  by a  $p \times m$  matrix  $B$ , then the composition  $\mathbf{g} \circ \mathbf{f}: \mathbb{R}^n \rightarrow \mathbb{R}^p$  is represented by a  $p \times n$  matrix  $C$ , and  $C = BA$ .

1. The proof of the linearity of  $\mathbf{g} \circ \mathbf{f}$  given above is unchanged.
2. The proof that  $C = BA$  is easy to generalize, but can it be shortened? Since

$$\begin{array}{ccc} \mathbf{x} & \xrightarrow{\mathbf{f}} & \mathbf{y} & \xrightarrow{\mathbf{g}} & \mathbf{z} \\ & \searrow & \swarrow & & \\ & & \mathbf{g} \circ \mathbf{f} & & \end{array} \iff \begin{array}{ccc} \mathbf{x} & \xrightarrow{A} & \mathbf{y} & \xrightarrow{B} & \mathbf{z} \\ & \searrow & \swarrow & & \\ & & C & & \end{array}$$

we have  $\mathbf{g} \circ \mathbf{f}(\mathbf{x}) = \mathbf{g}(\mathbf{f}(\mathbf{x})) = \mathbf{g}(\mathbf{y}) = B\mathbf{y} = B(\mathbf{f}(\mathbf{x})) = B(A\mathbf{x})$ , and since  $\mathbf{g} \circ \mathbf{f}(\mathbf{x}) = C\mathbf{x}$  we have  $B(A\mathbf{x}) = C\mathbf{x}$  for all  $\mathbf{x} \in \mathbb{R}^n$ . Can we now conclude  $C = BA$ ? Why not? To show  $(BA)\mathbf{x} = B(A\mathbf{x})$ , we need to go through the steps exemplified above. But this generalization is now a simple extension of the summation notation.



Since  $\mathbf{f}$  is linear and  $\mathbf{f}(\mathbf{u}_j) \in \mathbb{R}^m$ , we have

$$\mathbf{f}(\mathbf{x}) = \sum_{j=1}^n x_j \mathbf{f}(\mathbf{u}_j) = \sum_{j=1}^n x_j \sum_{i=1}^m a_{ij} \mathbf{e}_i = \sum_{i=1}^m \left( \sum_{j=1}^n a_{ij} x_j \right) \mathbf{e}_i = \sum_{i=1}^m y_i \mathbf{e}_i = \mathbf{y},$$

where  $y_i = \sum_{j=1}^n a_{ij} x_j$ ,  $i = 1, \dots, m$ . Similarly, since  $\mathbf{g}$  is linear and  $\mathbf{g}(\mathbf{e}_i) \in \mathbb{R}^p$ , we have

$$\mathbf{g}(\mathbf{y}) = \sum_{i=1}^m y_i \mathbf{g}(\mathbf{e}_i) = \sum_{i=1}^m y_i \sum_{k=1}^p b_{ki} \mathbf{v}_k = \sum_{k=1}^p \left( \sum_{i=1}^m b_{ki} y_i \right) \mathbf{v}_k = \sum_{k=1}^p z_k \mathbf{v}_k,$$

where

$$z_k = \sum_{i=1}^m b_{ki} y_i = \sum_{i=1}^m b_{ki} \sum_{j=1}^n a_{ij} x_j = \sum_{j=1}^n \left( \sum_{i=1}^m b_{ki} a_{ij} \right) x_j.$$

Also, since  $\mathbf{g} \circ \mathbf{f}$  is linear and  $\mathbf{g} \circ \mathbf{f}(\mathbf{u}_j) \in \mathbb{R}^p$ ,

$$\mathbf{g} \circ \mathbf{f}(\mathbf{x}) = \sum_{j=1}^n x_j (\mathbf{g} \circ \mathbf{f})(\mathbf{u}_j) = \sum_{j=1}^n x_j \sum_{k=1}^p c_{kj} \mathbf{v}_k = \sum_{k=1}^p \left( \sum_{j=1}^n c_{kj} x_j \right) \mathbf{v}_k = \sum_{k=1}^p z_k \mathbf{v}_k,$$

where  $z_k = \sum_{j=1}^n c_{kj} x_j$ . Hence  $C = BA$ , where

$$c_{kj} = \sum_{i=1}^m b_{ki} a_{ij}, \quad k = 1, \dots, p; \quad j = 1, \dots, n.$$

3. The proof gives the rule for matrix multiplication.

Since  $\mathbf{f} \circ \mathbf{g} \neq \mathbf{g} \circ \mathbf{f}$  in general, it follows as a corollary that  $AB \neq BA$ . That is, the non-commutativity of matrix multiplication is a direct consequence of the non-commutativity of function composition.

The properties of matrix multiplication follow directly from the properties of function composition, for example, the associative property,

$$\mathbf{f} \circ (\mathbf{g} \circ \mathbf{h}) = (\mathbf{f} \circ \mathbf{g}) \circ \mathbf{h} \quad \Rightarrow \quad A(BC) = (AB)C.$$

## 5 Determinant as a scale factor of area

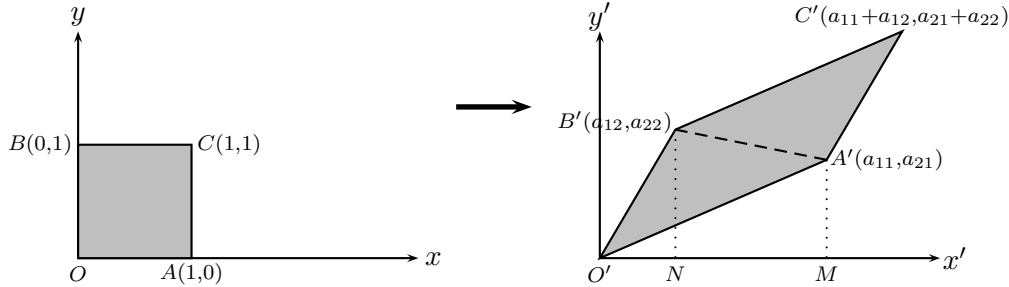
What does the determinant of a square matrix measure? The scale factor of area in linear transformations of the plane provides motivation for the study of determinants which are often introduced devoid of any motivation. Consider the image of the unit square with vertices at  $(0, 0)$ ,  $(1, 0)$ ,  $(0, 1)$  and  $(1, 1)$ . We have

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 0 & a_{11} & a_{12} & a_{11} + a_{12} \\ 0 & a_{21} & a_{22} & a_{21} + a_{22} \end{bmatrix}$$

and so

$$\begin{bmatrix} 0 \\ 0 \end{bmatrix} \mapsto \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad \begin{bmatrix} 1 \\ 0 \end{bmatrix} \mapsto \begin{bmatrix} a_{11} \\ a_{21} \end{bmatrix}, \quad \begin{bmatrix} 0 \\ 1 \end{bmatrix} \mapsto \begin{bmatrix} a_{12} \\ a_{22} \end{bmatrix}, \quad \begin{bmatrix} 1 \\ 1 \end{bmatrix} \mapsto \begin{bmatrix} a_{11} + a_{12} \\ a_{21} + a_{22} \end{bmatrix}.$$

The image of the unit square is a parallelogram as shown below.



By symmetry, the area of the parallelogram  $OA'C'B'$  is twice the area of the triangle  $OA'B'$ . Now

$$\begin{aligned} \text{area of } \triangle OA'B' &= |\text{area of } \triangle OB'N + \text{area of trapezium } B'NMA' - \text{area of } \triangle OA'M| \\ &= \left| \frac{1}{2}a_{12}a_{22} + \frac{1}{2}(a_{21} + a_{22})(a_{11} - a_{12}) - \frac{1}{2}a_{11}a_{21} \right| \\ &= \left| \frac{1}{2}(a_{11}a_{22} - a_{12}a_{21}) \right|. \end{aligned}$$

Hence the area of the parallelogram  $OA'C'B'$  is  $|a_{11}a_{22} - a_{12}a_{21}|$  which is the absolute value of the determinant of matrix  $A$ .

If  $A$  is a  $3 \times 3$  matrix it can be shown that the absolute value of its determinant,  $|\det(A)|$ , is a measure of the volume of a parallelepiped whose edges are the columns of  $A$ . This can be shown equal to  $|\mathbf{a} \cdot (\mathbf{b} \times \mathbf{c})|$ , where the columns (or rows) of  $A$  are  $\mathbf{a}$ ,  $\mathbf{b}$  and  $\mathbf{c}$ .

## 6 An approach to derivatives

A standard definition given in multivariable calculus states that a function  $\mathbf{f}: D \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$  is *differentiable* at  $\mathbf{x} \in D$  if there exists a linear transformation  $\mathbf{T}_{\mathbf{x}}: \mathbb{R}^n \rightarrow \mathbb{R}^m$  and a function  $\mathbf{E}_{\mathbf{x}}$  defined in a neighbourhood of  $\mathbf{0}$  and continuous at  $\mathbf{0}$ , such that

$$\mathbf{f}(\mathbf{x} + \mathbf{h}) = \mathbf{f}(\mathbf{x}) + \mathbf{T}_{\mathbf{x}}(\mathbf{h}) + \|\mathbf{h}\|\mathbf{E}_{\mathbf{x}}(\mathbf{h}). \quad (6.1)$$

If  $\mathbf{f}$  is differentiable at  $\mathbf{x}$  its *derivative*, denoted by  $\mathbf{f}'(\mathbf{x})$ , is the linear transformation  $\mathbf{T}_{\mathbf{x}}$ . As we have seen in Sections 3 and 4, this linear transformation is represented by a  $m \times n$  matrix  $J$ . We shall see that  $J$  is a matrix of partial derivatives called the Jacobian matrix of  $\mathbf{f}$  at  $\mathbf{x}$ . The value  $\mathbf{T}_{\mathbf{x}}(\mathbf{h})$  is represented by the matrix-vector product  $J\mathbf{h}$ , where  $\mathbf{h}$  is represented by a  $n \times 1$  column vector. Some applied mathematicians adopt a pragmatic view and define the derivative of  $\mathbf{f}$  as the Jacobian matrix  $J$ . This is sensible if the question of differentiability does not arise.

On the other hand, nearly all textbooks give the definition: A function  $f: (a, b) \rightarrow \mathbb{R}$  is differentiable at  $x \in (a, b)$  if

$$\lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \quad (6.2)$$

exists. If the limit exists it is called the derivative of  $f$  at  $x$  and is denoted by  $f'(x)$ . This definition is motivated by rates of change and the slope of a tangent line as a limiting value of the slope of secants. However, it bears little resemblance to the definition in higher dimensions and is a source of difficulty for students. What is needed is a definition in one variable that generalizes more naturally, especially when the concept of a derivative has already been well motivated.

**Theorem 2** Let  $K \in \mathbb{R}$  be a constant depending possibly on  $x$  but independent of  $h$ , and let  $I_0$  denote an open interval centred at 0. The following statements are equivalent.

- (a)  $\lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} = K$ .
- (b) There exists a function  $E : I_0 \rightarrow \mathbb{R}$ , continuous at 0, such that  $f(x+h) = f(x) + Kh + hE(h)$ .
- (c) There exists a function  $F : I_0 \rightarrow \mathbb{R}$ , continuous at 0, such that  $f(x+h) = f(x) + Kh + |h|F(h)$ .

**Proof:** (a)  $\implies$  (b) Define  $E$  on an interval  $I_0$  by

$$E(h) = \begin{cases} \frac{f(x+h) - f(x)}{h} - K, & \text{if } h \neq 0, \\ 0, & \text{if } h = 0. \end{cases}$$

It then follows that

$$\lim_{h \rightarrow 0} E(h) = \lim_{h \rightarrow 0} \left( \frac{f(x+h) - f(x)}{h} - K \right) = 0 = E(0),$$

and  $E$  is continuous at 0. If  $h \neq 0$  the formula for  $E(h)$  can be rearranged to give

$$f(x+h) = f(x) + Kh + hE(h).$$

This expression which avoids division is also true when  $h = 0$ .

(b)  $\implies$  (c) Define  $F : I_0 \rightarrow \mathbb{R}$  by

$$F(h) = \begin{cases} \frac{h}{|h|}E(h), & \text{if } h \neq 0, \\ 0, & \text{if } h = 0. \end{cases}$$

Now as  $h \rightarrow 0^+$ ,  $F(h) = E(h) \rightarrow 0$ , and as  $h \rightarrow 0^-$ ,  $F(h) = -E(h) \rightarrow 0$ . Hence  $F(h) \rightarrow 0 = F(0)$  as  $h \rightarrow 0$ , and  $F$  is continuous at 0. Therefore,  $hE(h) = |h|F(h)$  and

$$f(x+h) = f(x) + Kh + hE(h) \implies f(x+h) = f(x) + Kh + |h|F(h).$$

(c)  $\implies$  (a) Rearranging the expression yields

$$\lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} = \lim_{h \rightarrow 0} \left( K + \frac{|h|}{h}F(h) \right) = K,$$

because, since  $F$  is continuous at 0, as  $h \rightarrow 0^+$ ,  $(|h|/h)F(h) = F(h) \rightarrow 0$ , and as  $h \rightarrow 0^-$ ,  $(|h|/h)F(h) = -F(h) \rightarrow 0$  so that  $(|h|/h)F(h) \rightarrow 0$  as  $h \rightarrow 0$ .  $\square$

**Remark:** Both the statements (b) and (c) do not involve division by  $h$  explicitly and either could be used as an alternative to statement (a) in the definition of differentiability at a point. For example, the statement (b) is used for ease of application, whereas the equivalent statement (c) is a special case of (6.1). We now present some simple examples based on statement (b).

**Example 2** Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be defined by  $f(x) = x^3$ . By the binomial theorem,

$$f(x+h) = (x+h)^3 = x^3 + 3x^2h + 3xh^2 + h^3 = f(x) + Kh + hE(h),$$

where  $K = 3x^2$  and  $E(h) = h(3x+h)$ . Since  $E(h) \rightarrow 0 = E(0)$  as  $h \rightarrow 0$  it follows by definition that  $f$  is differentiable with derivative  $f'(x) = K = 3x^2$ .  $\square$

**Example 3** Let  $f : \mathbb{R} \setminus \{0\} \rightarrow \mathbb{R}$  be defined by  $f(x) = 1/x$ . Here we have

$$f(x+h) - f(x) = \frac{1}{x+h} - \frac{1}{x} = \frac{-h}{x(x+h)} = \frac{-h}{x^2} + \left( \frac{h}{x^2} - \frac{h}{x(x+h)} \right) = Kh + hE(h),$$

where  $K = -1/x^2$  and  $E(h) = 1/x^2 - 1/x(x+h) \rightarrow 0 = E(0)$  as  $h \rightarrow 0$ . Hence  $f$  is differentiable at  $x$  with derivative  $f'(x) = -1/x^2$ .  $\square$

**Example 4** If  $f : [0, \infty) \rightarrow \mathbb{R}$  is defined by  $f(x) = \sqrt{x}$  we have

$$\begin{aligned} f(x+h) - f(x) &= \sqrt{x+h} - \sqrt{x} = \frac{(x+h) - x}{\sqrt{x+h} + \sqrt{x}} = \frac{h}{2\sqrt{x}} + \left( \frac{h}{\sqrt{x+h} + \sqrt{x}} - \frac{h}{2\sqrt{x}} \right) \\ &= Kh + hE(h), \end{aligned}$$

where  $K = \frac{1}{2\sqrt{x}}$  and  $E(h) = \frac{1}{\sqrt{x+h} + \sqrt{x}} - \frac{1}{2\sqrt{x}} \rightarrow 0 = E(0)$  as  $h \rightarrow 0$ . Hence  $f$  is differentiable with derivative  $f'(x) = \frac{1}{2\sqrt{x}}$ . Note that  $x \neq 0$  even though  $0 \in \text{dom}(f)$  and also that  $x+h \geq 0$ . This means the interval on which  $E$  is defined is restricted in this case.  $\square$

**Example 5** Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be defined by  $f(x) = |x|$ . Since  $\lim_{x \rightarrow 0} f(x) = 0 = f(0)$ ,  $f$  is indeed continuous at 0 but its graph shows it is not differentiable there. Let us apply the definition. Now

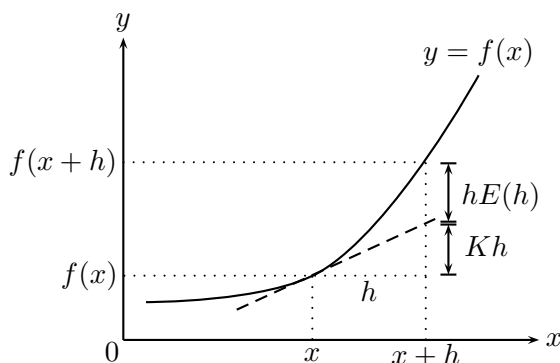
$$f(h) = |h| = \begin{cases} h, & \text{if } h \geq 0, \\ -h, & \text{if } h < 0. \end{cases}$$

If  $h > 0$ , we have  $f(h) = h = f(0) + Kh + hE(h)$ , where  $K = 1$  and  $E(h) = 0 \rightarrow 0$  as  $h \rightarrow 0^+$ . If  $h < 0$ ,  $f(h) = -h = f(0) + Kh + hE(h)$  with  $K = -1$  and  $E(h) = 0 \rightarrow 0$  as  $h \rightarrow 0^-$ . Since  $K$  is not unique, it follows that  $f'(0)$  does not exist and  $f$  is not differentiable at 0.

Suppose  $x > 0$  and  $x+h > 0$ . Then  $f(x+h) = |x+h| = x+h = f(x) + Kh + E(h)$ , where  $K = 1$  and  $E(h) = 0 \rightarrow 0$  as  $h \rightarrow 0$ . Hence  $f$  is differentiable at  $x$  with derivative  $f'(x) = 1$ . Similarly,  $f'(x) = -1$  if  $x < 0$ .  $\square$

As can be seen from the examples, there is little essential difference between using the statement (a) or (b) in the calculation of derivatives. What is different is that the evaluation of limits in (b) is now reduced to finding zero limits (and showing continuity at 0). We remark that a definition based on (b) avoids division by  $h$  and therefore the 'sticky and tricky' process of taking limits as  $h \rightarrow 0$  without allowing  $h = 0$  never arises. The statement (b) does not restrict  $h$  to nonzero values.

A discussion of the merits of an approach to derivatives based on statement (b) includes its geometrical interpretation and its connection with differentials for a function differentiable at  $x$ .



1. The slope or gradient of the tangent line at  $x$  (shown dashed) is  $K = f'(x)$ .



2. The equation of the tangent line is given by the ‘linear’ part of the expression in (b), that is,  $y = f(x) + Kh$  and is obtained for free. Putting  $x = a$  and then  $h = x - a$  gives the equation of the tangent line at  $a$ ,  $y = f(a) + K(x - a)$ .
3. The function  $E$  can be called an ‘error function’ because it is a measure of the difference of the function from its linear approximation. To show differentiability of  $f$  at  $x$  is to show that  $E$  is continuous at 0.
4. In some books the increment  $h$  is denoted by  $\delta x$  or  $\Delta x$  with the corresponding change in the function denoted by  $\delta y$  or  $\Delta y$ . Since  $x$  is the ‘independent’ variable we are free to choose the ‘differential of  $x$ ’, denoted by  $dx$ , to equal  $h$ . In this case the differential of  $y$ , denoted by  $dy$ , is represented by the change to the tangent line, that is,

$$dy = Kh = f'(x) dx.$$

The differentials  $dx$  and  $dy$  are finite quantities, and exist as algebraic elements independently of the notion of ‘infinitesimals’ used in physics. So it is legitimate to regard the derivative  $f'(x)$  as a quotient of two differentials as in

$$\frac{dy}{dx} = f'(x).$$

From the diagram, we then have

$$\Delta y = f(x + h) - f(x) = dy + dx E(dx) \Rightarrow \frac{\Delta y}{dx} = \frac{dy}{dx} + E(dx) = f'(x) + E(dx),$$

and hence, on recalling that  $h = dx = \Delta x$ , we have

$$E(h) = \frac{\Delta y}{\Delta x} - \frac{dy}{dx},$$

the difference between the slope of the secant line and the slope of the tangent line. So differentiability is related to

$$\lim_{\Delta x \rightarrow 0} \frac{\Delta y}{\Delta x} = \frac{dy}{dx}.$$

We now return to consider the Jacobian matrix for a differentiable function  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ . For simplicity, we shall consider functions whose domains are the whole of the space. We shall do this in several steps. We have already discussed derivatives of real-valued functions of a single variable. We next consider real and vector-valued functions of several variables. We remark that the terms ‘totally differentiable’ and ‘total derivative’ are used widely for ‘differentiable’ and ‘derivative’ respectively as a means of distinguishing between various types of partial and directional derivatives.

## 6.1 Real-valued functions of several variables

First we consider a simple example.

**Example 6** Let  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  be defined by  $f(x, y) = x^2 y$ . Show  $f$  is differentiable at  $(x, y)$  and find its derivative  $f'(x, y)$ . Now

$$\begin{aligned} f(x + h, y + k) &= (x + h)^2 (y + k) = (x^2 + 2xh + h^2)(y + k) \\ &= x^2 y + (2xyh + x^2 k) + yh^2 + 2xhk + h^2 k \\ &= f(x, y) + T_{(x,y)}(h, k) + \|(h, k)\| E(h, k), \end{aligned}$$

where  $T_{(x,y)}(h, k) = 2xyh + x^2k$  and  $E(h, k)$  is defined by

$$\|(h, k)\|E(h, k) = yh^2 + 2xhk + h^2k.$$

Now, by the triangle inequality, we have

$$\begin{aligned} \|(h, k)\|E(h, k) &= |yh^2 + 2xhk + h^2k| \leq |y|h^2 + 2|x||h||k| + h^2|k| \\ &\leq |y|(h^2 + k^2) + 2|x|(h^2 + k^2) + (h^2 + k^2)^{3/2} \\ &= \|(h, k)\|^2 \left( |y| + 2|x| + \|(h, k)\| \right), \\ \implies |E(h, k)| &\leq \|(h, k)\| \left( |y| + 2|x| + \|(h, k)\| \right), \end{aligned}$$

since  $|h| \leq \sqrt{h^2 + k^2}$  and  $|k| \leq \sqrt{h^2 + k^2}$ . It follows that as  $(h, k) \rightarrow (0, 0)$ ,  $|E(h, k)| \rightarrow 0$  and hence  $E(h, k) \rightarrow 0$ . It is continuous at  $(0, 0)$  if we define  $E(0, 0) = 0$ . Therefore, by definition,  $f$  is differentiable at  $(x, y)$  with derivative given by  $f'(x, y) = (2xy, x^2)$ .  $\square$

Consider a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  with  $\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$  and  $\mathbf{h} = (h_1, \dots, h_n) \in \mathbb{R}^n$ . By definition (6.1), if  $f$  is differentiable at  $\mathbf{x}$  then

$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + T_{\mathbf{x}}(\mathbf{h}) + \|\mathbf{h}\|E_{\mathbf{x}}(\mathbf{h}), \quad (6.3)$$

where  $T_{\mathbf{x}}$  is a linear function and  $E_{\mathbf{x}}(\mathbf{h}) \rightarrow 0 = E_{\mathbf{x}}(\mathbf{0})$  as  $\mathbf{h} \rightarrow \mathbf{0}$ . Let  $\mathbf{u}_1, \dots, \mathbf{u}_n$  be the unit coordinate vectors in  $\mathbb{R}^n$ . Then, by linearity, we have

$$\mathbf{h} = \sum_{j=1}^n h_j \mathbf{u}_j \implies T_{\mathbf{x}}(\mathbf{h}) = \sum_{j=1}^n h_j T_{\mathbf{x}}(\mathbf{u}_j). \quad (6.4)$$

Consider  $\mathbf{h} \rightarrow \mathbf{0}$  along the  $j$ -th axis, that is,  $\mathbf{h} = h_j \mathbf{u}_j$  with  $h_j \rightarrow 0$ . Then  $T_{\mathbf{x}}(\mathbf{h}) = h_j T_{\mathbf{x}}(\mathbf{u}_j)$  and

$$f(\mathbf{x} + h_j \mathbf{u}_j) = f(\mathbf{x}) + h_j T_{\mathbf{x}}(\mathbf{u}_j) + |h_j| E_{\mathbf{x}}(h_j \mathbf{u}_j)$$

with  $E_{\mathbf{x}}(h_j \mathbf{u}_j) \rightarrow 0 = E_{\mathbf{x}}(\mathbf{0})$  as  $h_j \rightarrow 0$ . Hence, by Theorem 2 (c), the *partial derivative* of  $f$  along the  $j$ -th direction exists and equals  $T_{\mathbf{x}}(\mathbf{u}_j)$  for each  $j = 1, \dots, n$ ,

$$T_{\mathbf{x}}(\mathbf{u}_j) = \frac{\partial f(\mathbf{x})}{\partial x_j}.$$

Substituting into (6.4) we therefore obtain

$$T_{\mathbf{x}}(\mathbf{h}) = \sum_{j=1}^n h_j \frac{\partial f(\mathbf{x})}{\partial x_j} = \nabla f(\mathbf{x}) \cdot \mathbf{h},$$

where  $\nabla f(\mathbf{x})$  is the gradient of  $f(\mathbf{x})$  represented by the  $1 \times n$  row-vector,

$$\nabla f(\mathbf{x}) = \left[ \frac{\partial f(\mathbf{x})}{\partial x_1}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_j}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_n} \right].$$

The dot product is represented by matrix multiplication if  $\mathbf{h}$  is represented by an  $n \times 1$  column vector.

Recall that in  $\mathbb{R}^2$  the equation of a tangent line to the curve  $y = f(x)$  at  $x_0$  is given by

$$y = f(x_0) + f'(x_0)(x - x_0).$$

In  $\mathbb{R}^3$  this generalizes to the equation of a tangent plane to the surface  $z = f(x, y)$  at  $(x_0, y_0)$  given by

$$z = f(x_0, y_0) + \nabla f(x_0, y_0) \cdot (x - x_0, y - y_0) = f(x_0, y_0) + (x - x_0) \frac{\partial f}{\partial x} \Big|_{(x_0, y_0)} + (y - y_0) \frac{\partial f}{\partial y} \Big|_{(x_0, y_0)}.$$

## 6.2 Vector-valued functions of several variables

Finally, we consider  $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ . Let  $\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$ ,  $\mathbf{h} = (h_1, \dots, h_n) \in \mathbb{R}^n$ , and let  $\mathbf{u}_1, \dots, \mathbf{u}_n$  and  $\mathbf{e}_1, \dots, \mathbf{e}_m$  be the unit coordinate vectors in  $\mathbb{R}^n$  and  $\mathbb{R}^m$  respectively. Then  $\mathbf{T}_{\mathbf{x}} : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is linear and  $\mathbf{T}_{\mathbf{x}}(\mathbf{h}) \in \mathbb{R}^m$  and so

$$\mathbf{T}_{\mathbf{x}}(\mathbf{h}) = \sum_{i=1}^m T_{\mathbf{x}}^i(\mathbf{h}) \mathbf{e}_i.$$

Since  $\mathbf{f}(\mathbf{x}) \in \mathbb{R}^m$ , we can write  $\mathbf{f}(\mathbf{x}) = \sum_{i=1}^m f_i(\mathbf{x}) \mathbf{e}_i$ , where  $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$ . Similarly,  $\mathbf{E}_{\mathbf{x}}(\mathbf{h}) = \sum_{i=1}^m E_{\mathbf{x}}^i(\mathbf{h}) \mathbf{e}_i$ . If  $\mathbf{f}$  is differentiable at  $\mathbf{x}$ , then

$$\sum_{i=1}^m f_i(\mathbf{x} + \mathbf{h}) \mathbf{e}_i = \sum_{i=1}^m f_i(\mathbf{x}) \mathbf{e}_i + \sum_{i=1}^m T_{\mathbf{x}}^i(\mathbf{h}) \mathbf{e}_i + \|\mathbf{h}\| \sum_{i=1}^m E_{\mathbf{x}}^i(\mathbf{h}) \mathbf{e}_i,$$

with  $E_{\mathbf{x}}^i(\mathbf{h}) \rightarrow 0 = E_{\mathbf{x}}^i(\mathbf{0})$  as  $\mathbf{h} \rightarrow \mathbf{0}$  for each  $i = 1, \dots, m$ . Hence

$$f_i(\mathbf{x} + \mathbf{h}) = f_i(\mathbf{x}) + T_{\mathbf{x}}^i(\mathbf{h}) + \|\mathbf{h}\| E_{\mathbf{x}}^i(\mathbf{h})$$

which has the same form as (6.3). Noting that  $T_{\mathbf{x}}^i : \mathbb{R}^n \rightarrow \mathbb{R}$  is linear for each  $i = 1, \dots, m$  we follow the same procedure as in Section 6.1 and conclude

$$T_{\mathbf{x}}^i(\mathbf{h}) = \nabla f_i(\mathbf{x}) \cdot \mathbf{h},$$

and, in terms of matrix representation, we have

$$\mathbf{T}_{\mathbf{x}}(\mathbf{h}) = \begin{bmatrix} \nabla f_1(\mathbf{x}) \cdot \mathbf{h} \\ \vdots \\ \nabla f_i(\mathbf{x}) \cdot \mathbf{h} \\ \vdots \\ \nabla f_m(\mathbf{x}) \cdot \mathbf{h} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \frac{\partial f_1(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial f_1(\mathbf{x})}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_i(\mathbf{x})}{\partial x_1} & \frac{\partial f_i(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial f_i(\mathbf{x})}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m(\mathbf{x})}{\partial x_1} & \frac{\partial f_m(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial f_m(\mathbf{x})}{\partial x_n} \end{bmatrix} \begin{bmatrix} h_1 \\ \vdots \\ h_j \\ \vdots \\ h_n \end{bmatrix}.$$

Hence the Jacobian matrix is the  $m \times n$  matrix of partial derivatives,

$$J(\mathbf{f}(\mathbf{x})) = \frac{\partial(f_1, \dots, f_m)}{\partial(x_1, \dots, x_n)} = \begin{bmatrix} \nabla f_1(\mathbf{x}) \\ \vdots \\ \nabla f_i(\mathbf{x}) \\ \vdots \\ \nabla f_m(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \frac{\partial f_1(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial f_1(\mathbf{x})}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_i(\mathbf{x})}{\partial x_1} & \frac{\partial f_i(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial f_i(\mathbf{x})}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m(\mathbf{x})}{\partial x_1} & \frac{\partial f_m(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial f_m(\mathbf{x})}{\partial x_n} \end{bmatrix}.$$

## 6.3 The chain rule

**Theorem 3** *If  $f$  is differentiable at  $x$  with derivative  $f'(x)$  and  $g$  is differentiable at  $f(x)$  with derivative  $g'(f(x))$ , then the composition  $g \circ f$  is differentiable at  $x$  with derivative  $(g \circ f)'(x) = g'(f(x))f'(x)$ .*

**Proof:** Suppose  $f$  is differentiable at  $x \in \text{dom}(f)$  and let  $x + h \in \text{dom}(f)$ . We then have

$$f(x + h) = f(x) + f'(x)h + hF_x(h),$$

with  $F_x$  continuous at 0. If  $f$  maps its domain into the domain of  $g$ , then let  $y = f(x)$  and  $y + k = f(x + h)$  so that, by the continuity of  $f$  at  $x$ ,

$$k = f(x + h) - f(x) \rightarrow 0 = k(0) \text{ as } h \rightarrow 0.$$

If  $g$  is also differentiable at  $y$ , then

$$g(y + k) = g(y) + g'(y)k + kG_y(k),$$

with  $G_y$  continuous at 0. Hence

$$\begin{aligned} g(f(x + h)) &= g(f(x)) + g'(f(x))\left(f'(x)h + hF_x(h)\right) + kG_y(k) \\ &= g(f(x)) + g'(f(x))f'(x)h + hE_x(h), \end{aligned}$$

where

$$hE_x(h) = hg'(f(x))F_x(h) + kG_y(k).$$

Now, by the triangular inequality,

$$\begin{aligned} |k| &= |f(x + h) - f(x)| \leq |h|(|f'(x)| + |F_x(h)|) \\ \Rightarrow |h||E_x(h)| &\leq |h|\left(|g'(f(x))||F_x(h)| + |G_y(k)|(|f'(x)| + |F_x(h)|)\right) \\ \Rightarrow |E_x(h)| &\leq |g'(f(x))||F_x(h)| + |G_y(k)|(|f'(x)| + |F_x(h)|). \end{aligned}$$

As  $h \rightarrow 0$  we have  $k \rightarrow 0 = k(0)$ ,  $F_x(h) \rightarrow 0 = F_x(0)$  and  $G_y(k) \rightarrow 0 = G_y(0)$  and therefore  $E_x(h) \rightarrow 0 = E_x(0)$ . Hence  $g \circ f$  is differentiable at  $x$  with derivative  $(g \circ f)'(x) = g'(f(x))f'(x)$ .  $\square$

The proof given generalizes in an obvious way to the general case as we illustrate. Let  $\mathbf{f}: \text{dom}(\mathbf{f}) \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$  and  $\mathbf{g}: \text{dom}(\mathbf{g}) \subset \mathbb{R}^m \rightarrow \mathbb{R}^p$  be functions with the range of  $\mathbf{f}$ ,  $\mathbf{f}(\text{dom}(\mathbf{f}))$ , a subset of the domain of  $\mathbf{g}$ . The chain rule then reads:

**Theorem 4** *If  $\mathbf{f}$  is differentiable at  $\mathbf{x}$  with derivative  $\mathbf{f}'(\mathbf{x})$  and  $\mathbf{g}$  is differentiable at  $\mathbf{f}(\mathbf{x})$  with derivative  $\mathbf{g}'(\mathbf{f}(\mathbf{x}))$ , then the composition  $\mathbf{g} \circ \mathbf{f}$  is differentiable at  $\mathbf{x}$  with derivative  $(\mathbf{g} \circ \mathbf{f})'(\mathbf{x}) = \mathbf{g}'(\mathbf{f}(\mathbf{x}))\mathbf{f}'(\mathbf{x})$ .*

**Proof:** Suppose  $\mathbf{f}$  is differentiable at  $\mathbf{x}$  and let  $\mathbf{x} + \mathbf{h} \in \text{dom}(\mathbf{f})$ . We then have

$$\mathbf{f}(\mathbf{x} + \mathbf{h}) = \mathbf{f}(\mathbf{x}) + \mathbf{f}'(\mathbf{x})(\mathbf{h}) + \|\mathbf{h}\|\mathbf{F}_x(\mathbf{h}),$$

with  $\mathbf{F}_x$  continuous at  $\mathbf{0}$ . If  $\mathbf{f}$  maps its domain into the domain of  $\mathbf{g}$ , then let  $\mathbf{y} = \mathbf{f}(\mathbf{x})$  and  $\mathbf{y} + \mathbf{k} = \mathbf{f}(\mathbf{x} + \mathbf{h})$  so that, by the continuity of  $\mathbf{f}$  at  $\mathbf{x}$ ,

$$\mathbf{k} = \mathbf{f}(\mathbf{x} + \mathbf{h}) - \mathbf{f}(\mathbf{x}) \rightarrow \mathbf{0} = \mathbf{k}(\mathbf{0}) \text{ as } \mathbf{h} \rightarrow \mathbf{0}.$$

If  $\mathbf{g}$  is also differentiable at  $\mathbf{y}$ , then

$$\mathbf{g}(\mathbf{y} + \mathbf{k}) = \mathbf{g}(\mathbf{y}) + \mathbf{g}'(\mathbf{y})(\mathbf{k}) + \|\mathbf{k}\|\mathbf{G}_y(\mathbf{k}),$$

with  $\mathbf{G}_y$  continuous at  $\mathbf{0}$ . Hence

$$\begin{aligned} \mathbf{g}(\mathbf{f}(\mathbf{x} + \mathbf{h})) &= \mathbf{g}(\mathbf{f}(\mathbf{x})) + \mathbf{g}'(\mathbf{f}(\mathbf{x}))\left(\mathbf{f}'(\mathbf{x})(\mathbf{h}) + \|\mathbf{h}\|\mathbf{F}_x(\mathbf{h})\right) + \|\mathbf{k}\|\mathbf{G}_y(\mathbf{k}) \\ &= \mathbf{g}(\mathbf{f}(\mathbf{x})) + \mathbf{g}'(\mathbf{f}(\mathbf{x}))\mathbf{f}'(\mathbf{x})(\mathbf{h}) + \|\mathbf{h}\|\mathbf{E}_x(\mathbf{h}), \end{aligned}$$

where

$$\|\mathbf{h}\|\mathbf{E}_x(\mathbf{h}) = \|\mathbf{h}\|\mathbf{g}'(\mathbf{f}(\mathbf{x}))\mathbf{F}_x(\mathbf{h}) + \|\mathbf{k}\|\mathbf{G}_y(\mathbf{k}).$$

Now, by the triangular and the Cauchy-Schwarz inequality,

$$\begin{aligned} \|\mathbf{k}\| &= \|\mathbf{f}(\mathbf{x} + \mathbf{h}) - \mathbf{f}(\mathbf{x})\| \leq \|\mathbf{h}\|(\|\mathbf{f}'(\mathbf{x})\| + \|\mathbf{F}_x(\mathbf{h})\|) \\ \Rightarrow \|\mathbf{h}\|\|\mathbf{E}_x(\mathbf{h})\| &\leq \|\mathbf{h}\|\left(\|\mathbf{g}'(\mathbf{f}(\mathbf{x}))\|\|\mathbf{F}_x(\mathbf{h})\| + \|\mathbf{G}_y(\mathbf{k})\|(\|\mathbf{f}'(\mathbf{x})\| + \|\mathbf{F}_x(\mathbf{h})\|)\right) \\ \Rightarrow \|\mathbf{E}_x(\mathbf{h})\| &\leq \|\mathbf{g}'(\mathbf{f}(\mathbf{x}))\|\|\mathbf{F}_x(\mathbf{h})\| + \|\mathbf{G}_y(\mathbf{k})\|(\|\mathbf{f}'(\mathbf{x})\| + \|\mathbf{F}_x(\mathbf{h})\|). \end{aligned}$$

As  $\mathbf{h} \rightarrow \mathbf{0}$  we have  $\mathbf{k} \rightarrow \mathbf{0} = \mathbf{k}(\mathbf{0})$ ,  $\mathbf{F}_x(\mathbf{h}) \rightarrow \mathbf{0} = \mathbf{F}_x(\mathbf{0})$  and  $\mathbf{G}_y(\mathbf{k}) \rightarrow \mathbf{0} = \mathbf{G}_y(\mathbf{0})$  and therefore  $\mathbf{E}_x(\mathbf{h}) \rightarrow \mathbf{0} = \mathbf{E}_x(\mathbf{0})$ . Hence  $\mathbf{g} \circ \mathbf{f}$  is differentiable at  $\mathbf{x}$  with derivative  $(\mathbf{g} \circ \mathbf{f})'(\mathbf{x}) = \mathbf{g}'(\mathbf{f}(\mathbf{x}))\mathbf{f}'(\mathbf{x})$ .  $\square$

**Remark:** As we can see this proof is identical with that given in the one-dimensional case, with the exception that the Cauchy-Schwarz inequality is needed here.

## 6.4 Matrix form of the chain rule

Let  $\mathbf{y} = \mathbf{f}(\mathbf{x})$ ,  $\mathbf{z} = \mathbf{g}(\mathbf{y})$  and denote  $\mathbf{g} \circ \mathbf{f} = \mathbf{h}$ . When the derivatives are represented by Jacobian matrices the chain rule takes the form

$$\begin{aligned} \mathbf{h}'(\mathbf{x}) &= \mathbf{g}'(\mathbf{f}(\mathbf{x}))\mathbf{f}'(\mathbf{x}), \\ J(\mathbf{h}(\mathbf{x})) &= J(\mathbf{g}(\mathbf{y}))J(\mathbf{f}(\mathbf{x})), \\ \frac{\partial(h_1, \dots, h_p)}{\partial(x_1, \dots, x_n)} &= \frac{\partial(g_1, \dots, g_p)}{\partial(y_1, \dots, y_m)} \frac{\partial(f_1, \dots, f_m)}{\partial(x_1, \dots, x_n)}, \\ \begin{bmatrix} \frac{\partial z_1}{\partial x_1} & \dots & \frac{\partial z_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial z_p}{\partial x_1} & \dots & \frac{\partial z_p}{\partial x_n} \end{bmatrix} &= \begin{bmatrix} \frac{\partial z_1}{\partial y_1} & \dots & \frac{\partial z_1}{\partial y_m} \\ \vdots & & \vdots \\ \frac{\partial z_p}{\partial y_1} & \dots & \frac{\partial z_p}{\partial y_m} \end{bmatrix} \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial y_m}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}. \end{aligned}$$

## 6.5 Some miscellaneous results

We consider some immediate results following from the differentiability of a function at a point. We do this for a function  $\mathbf{f}: \mathbb{R}^n \rightarrow \mathbb{R}^m$ . The first generalizes the one dimensional result that differentiability implies continuity.

**Theorem 5** *If  $\mathbf{f}$  is differentiable at  $\mathbf{x}$ , it is continuous at  $\mathbf{x}$ .*

**Proof:** If  $\mathbf{f}$  is differentiable at  $\mathbf{x}$ , then

$$\mathbf{f}(\mathbf{x} + \mathbf{h}) = \mathbf{f}(\mathbf{x}) + \mathbf{f}'(\mathbf{x})(\mathbf{h}) + \|\mathbf{h}\|\mathbf{E}_x(\mathbf{h})$$

with  $\mathbf{E}_x$  continuous  $\mathbf{0}$ . It follows that  $\mathbf{f}(\mathbf{x} + \mathbf{h}) \rightarrow \mathbf{f}(\mathbf{x})$  as  $\mathbf{h} \rightarrow \mathbf{0}$ . Hence  $\mathbf{f}$  is continuous at  $\mathbf{x}$ .  $\square$

The next result on continuity of a composition follows as a corollary of this result and by the chain rule of Theorem 4.

**Theorem 6** *If  $\mathbf{f}$  is continuous at  $\mathbf{x}$  and  $\mathbf{g}$  is continuous at  $\mathbf{f}(\mathbf{x})$ , then the composition  $\mathbf{g} \circ \mathbf{f}$  is continuous at  $\mathbf{x}$ .*

**Definition 1** (*Directional derivative*) *The directional derivative of  $\mathbf{f}$  at  $\mathbf{x}$  in the direction of  $\mathbf{s}$  is defined to be the limit*

$$\lim_{t \rightarrow 0} \frac{\mathbf{f}(\mathbf{x} + t\mathbf{s}) - \mathbf{f}(\mathbf{x})}{t}$$

*if it exists in which case it is denoted by  $\mathbf{f}'(\mathbf{x}; \mathbf{s})$ .*

**Theorem 7** *The directional derivative of  $\mathbf{f}$  at  $\mathbf{x}$  in the direction of  $\mathbf{s}$  exists if and only if there exist  $\mathbf{K}(\mathbf{x}; \mathbf{s}) \in \mathbb{R}^m$  and a function  $\mathbf{E} : I_0 \rightarrow \mathbb{R}^m$ , where  $I_0$  is some interval centred at 0, with  $\lim_{t \rightarrow 0} \mathbf{E}(t) = \mathbf{0} = \mathbf{E}(0)$  such that*

$$\mathbf{f}(\mathbf{x} + t\mathbf{s}) = \mathbf{f}(\mathbf{x}) + t\mathbf{K}(\mathbf{x}; \mathbf{s}) + t\mathbf{E}(t) \text{ for all } t \in I_0.$$

**Proof:** ( $\implies$ ) If  $\mathbf{f}'(\mathbf{x}; \mathbf{s})$  exists, then choose  $\mathbf{K}(\mathbf{x}; \mathbf{s}) = \mathbf{f}'(\mathbf{x}; \mathbf{s})$  and define  $\mathbf{E}$  by

$$\mathbf{E}(t) = \begin{cases} \frac{\mathbf{f}(\mathbf{x} + t\mathbf{s}) - \mathbf{f}(\mathbf{x})}{t} - \mathbf{f}'(\mathbf{x}; \mathbf{s}), & \text{if } t \neq 0, \\ \mathbf{0}, & \text{if } t = 0. \end{cases}$$

Hence we have

$$\lim_{t \rightarrow 0} \mathbf{E}(t) = \lim_{t \rightarrow 0} \left( \frac{\mathbf{f}(\mathbf{x} + t\mathbf{s}) - \mathbf{f}(\mathbf{x})}{t} - \mathbf{f}'(\mathbf{x}; \mathbf{s}) \right) = \mathbf{0} = \mathbf{E}(0),$$

and the required expression then follows from a simple rearrangement of the defining formula for  $\mathbf{E}$ .  
 (  $\Leftarrow$  ) Conversely, given  $\mathbf{f}(\mathbf{x} + t\mathbf{s}) = \mathbf{f}(\mathbf{x}) + t\mathbf{K}(\mathbf{x}; \mathbf{s}) + t\mathbf{E}(t)$  with  $\mathbf{E}(t) \rightarrow \mathbf{0} = \mathbf{E}(0)$  as  $t \rightarrow 0$ , we have

$$\lim_{t \rightarrow 0} \frac{\mathbf{f}(\mathbf{x} + t\mathbf{s}) - \mathbf{f}(\mathbf{x})}{t} = \lim_{t \rightarrow 0} \left( \mathbf{K}(\mathbf{x}; \mathbf{s}) + \mathbf{E}(t) \right) = \mathbf{K}(\mathbf{x}; \mathbf{s}).$$

Hence  $\mathbf{f}$  has a directional derivative at  $\mathbf{x}$  in the direction  $\mathbf{s}$ ,  $\mathbf{f}'(\mathbf{x}; \mathbf{s}) = \mathbf{K}(\mathbf{x}; \mathbf{s})$ . □

**Remarks:**

1. Theorem 7 may be used instead of Definition 1 to define directional derivatives.
2. The direction vector  $\mathbf{s}$  is usually taken to be a unit vector.
3. Let  $\mathbf{s}$  be any vector in  $\mathbb{R}^n$  and let  $\mathbf{g}$  denote the composition  $t \mapsto \mathbf{x} + t\mathbf{s} \mapsto \mathbf{f}(\mathbf{x} + t\mathbf{s})$ , so that  $\mathbf{g}(t) = \mathbf{f}(\mathbf{x} + t\mathbf{s})$ . Then  $\mathbf{f}$  is differentiable at  $\mathbf{x}$  if and only if  $\mathbf{g}$  is differentiable at 0, that is,

$$\mathbf{g}(t) = \mathbf{g}(0) + \mathbf{g}'(0)t + t\mathbf{G}(t)$$

with  $\mathbf{G}$  continuous at 0. Now  $\mathbf{g}'(0) = \mathbf{f}'(\mathbf{x})\mathbf{s}$  by the chain rule and, by Theorem 7, it follows that  $\mathbf{f}'(\mathbf{x}; \mathbf{s}) = \mathbf{f}'(\mathbf{x})\mathbf{s} = J(\mathbf{f}(\mathbf{x})) \cdot \mathbf{s}$ . In particular, if  $\mathbf{f}$  is real-valued we have  $f'(\mathbf{x}; \mathbf{s}) = \nabla f(\mathbf{x}) \cdot \mathbf{s}$  and if  $\mathbf{s} = \mathbf{u}_j$ , the unit coordinate vector along the  $j$ -th axis then the directional derivative is the partial derivative,  $f'(\mathbf{x}; \mathbf{u}_j) = \partial f(\mathbf{x})/\partial x_j$ . Hence directional derivatives at  $\mathbf{x}$  exist in all directions for a function differentiable at  $\mathbf{x}$ .