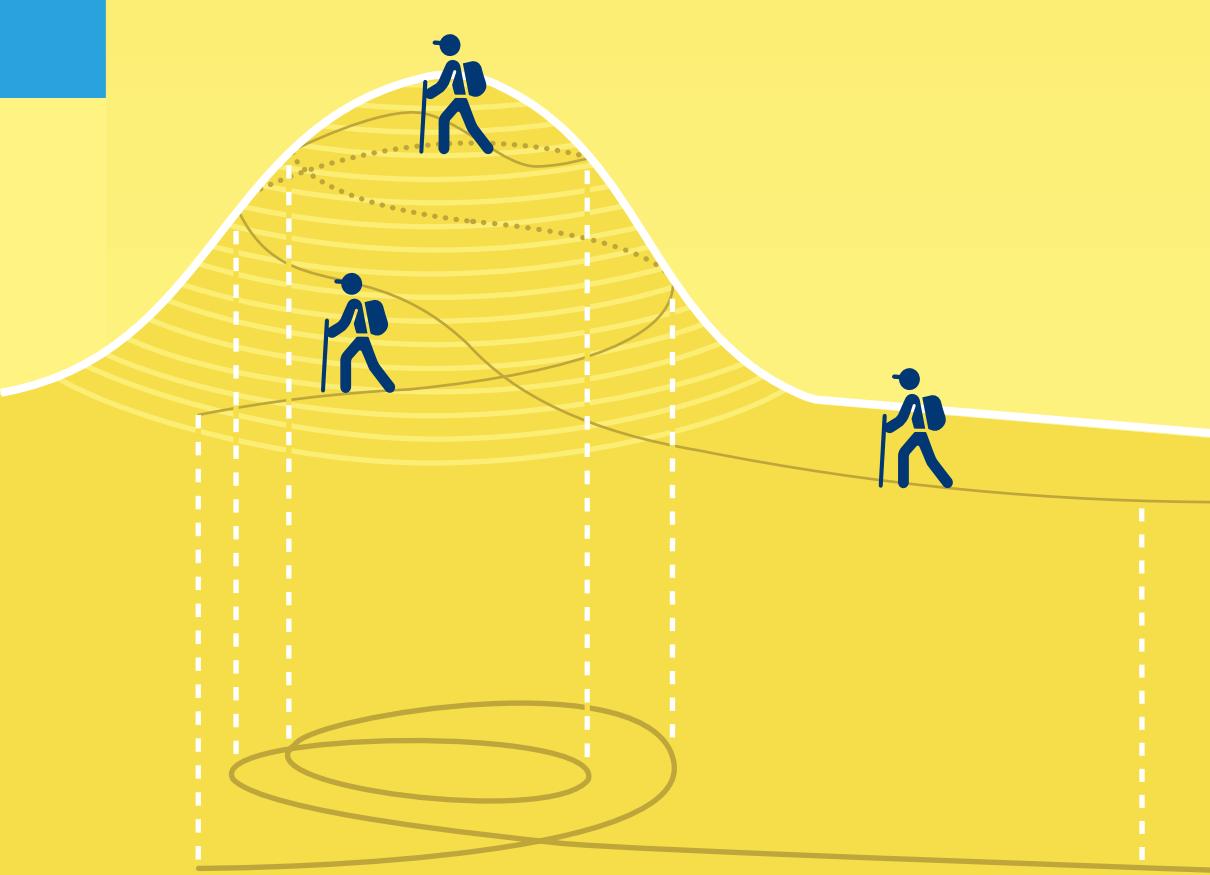


An Illustrative Guide to Multivariable and Vector Calculus

Stanley J. Miklavcic



An Illustrative Guide to Multivariable and Vector Calculus

Stanley J. Miklavcic

An Illustrative Guide to Multivariable and Vector Calculus

In collaboration with Ross A. Frick



Springer

Stanley J. Miklavcic
University of South Australia
(Mawson Lakes Campus)
Adelaide, SA, Australia

ISBN 978-3-030-33458-1 ISBN 978-3-030-33459-8 (eBook)
<https://doi.org/10.1007/978-3-030-33459-8>

Mathematics Subject Classification (2010): 26B05, 26B10, 26B12, 26B15, 26B20

© Springer Nature Switzerland AG 2020

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, expressed or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

This Springer imprint is published by the registered company Springer Nature Switzerland AG
The registered company address is: Gewerbestrasse 11, 6330 Cham, Switzerland

Dedication

To my children Arya, Nadia, Jacob, and David.

Preface

This book originated as a set of lectures prepared for courses given by me at the University of Linköping in Sweden and at the University of South Australia in Australia. At Linköping University the material (apart from Section 3.E) was delivered in a second year, single semester course (14 weeks, 2 two-hour lectures per week) to engineering students, with the first half focused on the differential calculus of real-valued multivariable functions, while the second half was divided between integral calculus and vector calculus. At the University of South Australia the subject was delivered in two separate semester courses (12 weeks, 2 two-hour lectures per week), the first of which was offered to second year engineering, science and mathematics students and featured differential and integral calculus, including an introduction to partial differential equations. The second course, taken mostly by third year mathematics and science students, dealt with vector calculus, although only the first five weeks of that course was covered by the material in this book.

The lectures generally were so well-received by students that it was thought the material might appeal to a wider audience. Having taken the decision to convert my notes into a book, I aimed for a document of manageable size rather than generate yet another bulky tome on calculus. The result is a book that students can carry easily to and from class, can take out and leaf through on the library lawn, or in a booth of a pub, or while lying on the banks of a river waiting for the fish to bite.

Very many ideas in mathematics are more readily conveyed and more easily appreciated when presented visually. This is certainly true of multivariable and vector calculus, and as my lecture material took advantage of many visual devices, I sought to capture the spirit if not the body of these same devices in printed form. Consequently, the majority of concepts are introduced and explained with the support of figures and graphics as well as the generous use of colour. Indeed, colour is used to highlight specific pieces of information, to emphasize relationships between variables in different

equations, and to distinguish different roles or actions. The inevitable issue of colour blindness was raised in the course of the book's development. To minimize difficulties, colour typesetting has been configured to allow for some degree of differentiation even by those readers with impaired colour vision. In addition, colour has been implemented so as only to benefit one's understanding, and *not* as an essential condition for understanding.

The book is self-contained and complete as an introduction to the theory of the differential and integral calculus of both real-valued and vector-valued multivariable functions. The entire material is suitable as a textbook in its own right for one to two, semester-long courses in either the second year or third year of University studies, and for students who have already completed courses in single variable calculus and linear algebra. Some selection of content may be necessary depending on student need and time available. For instance, as the topic of partial differential equations (PDEs) is normally offered as a separate course to mathematics students, Section 3.E can be skipped in a multivariable calculus course. On the other hand, a course in PDEs is not always included in engineering and science curricula, so Section 3.E is a pragmatic, albeit brief, introduction to the subject, particularly as its focus is on solving PDEs in simple cases. Alternatively, because of its illustrative emphasis, the book can also perform the role of a reference text to complement one of the more standard textbooks in advanced calculus, such as [1], thus providing the student with a different visual perspective.

Consequent to the ambition of producing a portable book, the reader should not be surprised that some areas of the calculus are not covered in detail. One other notable sacrifice is mathematical rigour. There are very few proofs included and those that have been are deliberately sketchy, included only to give students a rational justification for, or to illustrate the origin of, an idea. Consequently, students of pure mathematics may want to complement this book with one that offers a deeper analysis, such as [2].

Within each chapter is a sequence of Mastery Checks, exercises on the topic under discussion that are usually preceded by solved examples. Students are encouraged to attempt these Mastery Checks and keep a record of their solutions for future reference. To reinforce the ideas, additional exercises appear at the end of each chapter to supplement the Mastery Checks. Solutions to both sets of exercises are available to instructors upon request. I have limited the number of problems in order to restrict the size of the book, assuming that students would have access to auxiliary exercises in more standard treatises. All the same, the book contains over 90 Mastery Checks and over 120 Supplementary Exercises, many with multiple parts.

The reader should be aware that I have made use of mathematical symbols (such as \Rightarrow and \exists) and abbreviations (w.r.t., 3D) in place of text, a common

practice in mathematics texts and research literature. A glossary of definitions can be found at the end of the book. Wherever they appear in the book they should be read as the pieces of text they replace. Finally, for easy reference a list of Important Formulae, covering various topics in multi-variable and vector calculus, is given on page xiii.

Acknowledgements

In drafting this book I had great pleasure in working closely with my colleague Ross Frick who was instrumental in turning my original lecture material and supplementary notes into book form. His skill with L^AT_EX and MATLAB[®] was critical in this endeavour. I would also like to thank Dr. Loretta Bartolini, Mathematics Editor at Springer, for her strong support and encouragement of this venture and for her efficient handling of the publication of this book.

I will forever be indebted to Julie for her patience and enduring support over the many, many months of editing and re-editing to which this book was subjected. It is no exaggeration to say that without her understanding the task of completing this book would have been a far greater challenge than it has been.

Lastly, I would like to thank the students who have taken my course over the years, particularly those (now graduate) students who gave feedback on the notes prior to their publication. Their general enthusiasm has been an absolutely essential factor in getting the book to this point. I hope that future students of this important area of mathematics will also enjoy and be inspired by what this little volume has to offer.

Adelaide, Australia
December 2019

Stanley J. Miklavcic

Contents

1	Vectors and functions	1
1.A	Some vector algebra essentials	2
1.B	Introduction to sets	9
1.C	Real-valued functions	17
1.D	Coordinate systems	25
1.E	Drawing or visualizing surfaces in \mathbb{R}^3	27
1.F	Level sets	38
1.G	Supplementary problems	43
2	Differentiation of multivariable functions	49
2.A	The derivative	49
2.B	Limits and continuity	53
2.C	Partial derivatives	62
2.D	Differentiability of $f : \mathbb{R}^n \rightarrow \mathbb{R}$	67
2.E	Directional derivatives and the gradient	74
2.F	Higher-order derivatives	80
2.G	Composite functions and the chain rule	84
2.H	Implicit functions	101
2.I	Taylor's formula and Taylor series	113
2.J	Supplementary problems	119
3	Applications of the differential calculus	125
3.A	Extreme values of $f : \mathbb{R}^n \rightarrow \mathbb{R}$	125

3.B	Extreme points: The complete story	133
3.C	Differentials and error analysis	145
3.D	Method of least squares	146
3.E	Partial derivatives in equations: Partial differential equations	152
3.F	Supplementary problems	171
4	Integration of multivariable functions	177
4.A	Multiple integrals	177
4.B	Iterated integration in \mathbb{R}^2	184
4.C	Integration over complex domains	187
4.D	Generalized (improper) integrals in \mathbb{R}^2	193
4.E	Change of variables in \mathbb{R}^2	198
4.F	Triple integrals	204
4.G	Iterated integration in \mathbb{R}^3	207
4.H	Change of variables in \mathbb{R}^3	211
4.I	n -tuple integrals	213
4.J	Epilogue: Some practical tips for evaluating integrals	215
4.K	Supplementary problems	217
5	Vector calculus	223
5.A	Vector-valued functions	223
5.B	Vector fields	238
5.C	Line integrals	246
5.D	Surface integrals	260
5.E	Gauss's theorem	273
5.F	Green's and Stokes's theorems	281
5.G	Supplementary problems	293
Glossary of symbols		301
Bibliography		305
Index		307

Important Formulae

Multivariable calculus

- Suppose $f, g, h \in C^1$ and $w = f(u, v)$ where $u = g(x, y)$ and $v = h(x, y)$, then the partial derivative of w with respect to x is given by $\frac{\partial w}{\partial x} = \frac{\partial f}{\partial u} \frac{\partial g}{\partial x} + \frac{\partial f}{\partial v} \frac{\partial h}{\partial x}$, and similarly for the partial derivative of w with respect to y .
- Suppose $f \in C^3$ at the point (a, b) , then for points (x, y) in a neighbourhood of (a, b) the function f has the following Taylor approximation:

$$f(x, y) = f(a, b) + f_x(a, b)\Delta x + f_y(a, b)\Delta y + \frac{1}{2}Q(\Delta x, \Delta y) + R(\Delta x, \Delta y)$$

where

$$Q(\Delta x, \Delta y) = f_{xx}(a, b)(\Delta x)^2 + 2f_{xy}(a, b)\Delta x \Delta y + f_{yy}(a, b)(\Delta y)^2$$

and R is a remainder term of order $((\Delta x^2 + \Delta y^2)^{3/2})$.

- $f = f(x, y, z) \in C^2$ is a solution of Laplace's equation in domain $D \subset \mathbb{R}^3$ if f satisfies

$$\nabla^2 f = \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2} + \frac{\partial^2 f}{\partial z^2} = 0.$$

- The Jacobian for the transformation $(x, y) \mapsto (u, v)$:

$$J = \frac{\partial(x, y)}{\partial(u, v)} = \begin{vmatrix} \frac{\partial x}{\partial u} & \frac{\partial y}{\partial u} \\ \frac{\partial x}{\partial v} & \frac{\partial y}{\partial v} \end{vmatrix}$$

- Orthogonal curvilinear coordinate systems and their corresponding change of variables:

1. Polar coordinates: $r \geq 0$ and $0 \leq \theta \leq 2\pi$; Jacobian, $J = r$.

$$\{x = r \cos \theta, y = r \sin \theta\}$$

2. Cylindrical polar coordinates: $r \geq 0$, $0 \leq \theta \leq 2\pi$ and $z \in \mathbb{R}$; Jacobian, $J = r$.

$$\{x = r \cos \theta, y = r \sin \theta, z = z\}$$

3. Spherical polar coordinates: $\rho \geq 0$, $0 \leq \phi \leq \pi$, $0 \leq \theta \leq 2\pi$; Jacobian, $J = \rho^2 \sin \phi$.

$$\{x = \rho \sin \phi \cos \theta, y = \rho \sin \phi \sin \theta, z = \rho \cos \phi\}$$

- A level set to a function $f(x, y)$ is the set $\{(x, y) : \text{s.t. } f(x, y) = C\}$ for some real constant C . The gradient of f is always normal to a level set of f .
- An iterated integral of a function of two variables over a y -simple domain $D = \{(x, y) : a \leq x \leq b; g_1(x) \leq y \leq g_2(x)\}$

$$\iint_D f(x, y) dA = \int_a^b dx \int_{g_1(x)}^{g_2(x)} f(x, y) dy.$$

- An iterated integral of a function of two variables over a x -simple domain $D = \{(x, y) : c \leq y \leq d; h_1(y) \leq x \leq h_2(y)\}$

$$\iint_D f(x, y) dA = \int_c^d dy \int_{h_1(y)}^{h_2(y)} f(x, y) dx.$$

- For a bijective transformation $D \ni (x, y) \mapsto (u, v) \in E$ with Jacobian determinant, $J \neq 0$, the double integral of $f(x, y) = F(u, v)$ is

$$\iint_D f(x, y) dx dy = \iint_E F(u, v) \left| \frac{\partial(x, y)}{\partial(u, v)} \right| du dv.$$

Vector calculus

- A C^1 vector field $\mathbf{F} = \mathbf{F}(\mathbf{x})$ defined in some domain $D \subset \mathbb{R}^3$ is said to be
 1. solenoidal in D , if $\nabla \cdot \mathbf{F} = 0$ in D ;
 2. irrotational in D , if $\nabla \times \mathbf{F} = \mathbf{0}$ in D ;
 3. conservative in D , if $\mathbf{F} = \nabla\phi$ in D for some C^2 real-valued function $\phi = \phi(\mathbf{x})$.
- In terms of a 3D curvilinear coordinate system $\{\xi_1, \xi_2, \xi_3\}$, with unit vectors $\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3$, and scale factors h_1, h_2, h_3 , the gradient, divergence and curl operations on scalar ($\phi \in C^1$) and vector ($\mathbf{F} \in C^1$) fields, respectively, take the form

$$\nabla\phi = \frac{1}{h_1} \frac{\partial\phi}{\partial\xi_1} \mathbf{a}_1 + \frac{1}{h_2} \frac{\partial\phi}{\partial\xi_2} \mathbf{a}_2 + \frac{1}{h_3} \frac{\partial\phi}{\partial\xi_3} \mathbf{a}_3$$

$$\nabla \cdot \mathbf{F} = \frac{1}{h_1 h_2 h_3} \left[\frac{\partial}{\partial\xi_1} (h_2 h_3 F_1) + \frac{\partial}{\partial\xi_2} (h_1 h_3 F_2) + \frac{\partial}{\partial\xi_3} (h_1 h_2 F_3) \right],$$

$$\nabla \times \mathbf{F} = \frac{1}{h_1 h_2 h_3} \begin{vmatrix} h_1 \mathbf{a}_1 & h_2 \mathbf{a}_2 & h_3 \mathbf{a}_3 \\ \frac{\partial}{\partial\xi_1} & \frac{\partial}{\partial\xi_2} & \frac{\partial}{\partial\xi_3} \\ h_1 F_1 & h_2 F_2 & h_3 F_3 \end{vmatrix}$$

For Cartesian coordinates $\{\xi_1, \xi_2, \xi_3\} = \{x_1, x_2, x_3\} = \{x, y, z\}$, $h_1 = h_2 = h_3 = 1$, and ($\mathbf{a}_1 = \mathbf{e}_1 = \mathbf{i}$; $\mathbf{a}_2 = \mathbf{e}_2 = \mathbf{j}$; $\mathbf{a}_3 = \mathbf{e}_3 = \mathbf{k}$).

- Some useful vector identities. Suppose $\psi, \phi : R^3 \rightarrow R$ and $h : R \rightarrow R$ are C^1 scalar-valued functions, $\mathbf{f}, \mathbf{g} : R^3 \rightarrow R^3$ are C^1 vector-valued functions, $\mathbf{x} = (x, y, z)$ is a position vector of length $r = |\mathbf{x}| = \sqrt{x^2 + y^2 + z^2}$ and \mathbf{c} is a constant vector.

- (1) $\nabla(\phi\psi) = \psi\nabla\phi + \phi\nabla\psi$
- (2) $\nabla \cdot (\phi\mathbf{f}) = \phi\nabla \cdot \mathbf{f} + \mathbf{f} \cdot \nabla\phi$
- (3) $\nabla \times (\phi\mathbf{f}) = \phi\nabla \times \mathbf{f} + \nabla\phi \times \mathbf{f}$
- (4) $\nabla(\mathbf{f} \cdot \mathbf{g}) = (\mathbf{f} \cdot \nabla)\mathbf{g} + (\mathbf{g} \cdot \nabla)\mathbf{f} + \mathbf{f} \times (\nabla \times \mathbf{g}) + \mathbf{g} \times (\nabla \times \mathbf{f})$
- (5) $\nabla \cdot (\mathbf{f} \times \mathbf{g}) = \mathbf{g} \cdot (\nabla \times \mathbf{f}) - \mathbf{f} \cdot (\nabla \times \mathbf{g})$
- (6) $\nabla \times (\mathbf{f} \times \mathbf{g}) = \mathbf{f}(\nabla \cdot \mathbf{g}) - \mathbf{g}(\nabla \cdot \mathbf{f}) + (\mathbf{g} \cdot \nabla)\mathbf{f} - (\mathbf{f} \cdot \nabla)\mathbf{g}$
- (7) $\nabla \times (\nabla\phi) = \mathbf{0}$
- (8) $\nabla \cdot (\nabla \times \mathbf{f}) = 0$
- (9) $\nabla \times (\nabla \times \mathbf{f}) = \nabla(\nabla \cdot \mathbf{f}) - \nabla^2\mathbf{f}$
- (10) $\nabla \cdot \mathbf{x} = 3$
- (11) $\nabla h(r) = \frac{dh}{dr} \frac{\mathbf{x}}{r}$

$$(12) \quad \nabla \cdot (h(r)\mathbf{x}) = 3h(r) + r \frac{dh}{dr}$$

$$(13) \quad \nabla \times (h(r)\mathbf{x}) = \mathbf{0}$$

$$(14) \quad \nabla(\mathbf{c} \cdot \mathbf{x}) = \mathbf{c}$$

$$(15) \quad \nabla \cdot (\mathbf{c} \times \mathbf{x}) = 0$$

$$(16) \quad \nabla \times (\mathbf{c} \times \mathbf{x}) = 2\mathbf{c}$$

- Vector integration

1. Line integral of $\mathbf{f} = \mathbf{f}(\mathbf{r})$ over $\Gamma = \{\mathbf{r} = \mathbf{r}(t) : a \leq t \leq b\}$:

$$\int_{\Gamma} \mathbf{f} \cdot d\mathbf{r} = \int_a^b \mathbf{f}(\mathbf{r}(t)) \cdot \frac{d\mathbf{r}}{dt} dt$$

2. Surface integral of $\mathbf{f} = \mathbf{f}(\mathbf{r})$ over $S = \{\mathbf{r} = \mathbf{r}(u, v) : (u, v) \in D \subset \mathbb{R}^2\}$ with unit surface normal \mathbf{N} :

$$\iint_S \mathbf{f} \cdot d\mathbf{S} = \iint_S \mathbf{f} \cdot \mathbf{N} dS = \iint_D \mathbf{f}(\mathbf{r}(u, v)) \cdot \left(\frac{\partial \mathbf{r}}{\partial u} \times \frac{\partial \mathbf{r}}{\partial v} \right) du dv$$

3. Green's theorem for $\mathbf{f} = (f_1, f_2) \in C^1$ over a finite 2D region D bounded by a positively oriented closed curve Γ :

$$\oint_{\Gamma} \mathbf{f} \cdot d\mathbf{r} = \oint_{\Gamma} (f_1 dx + f_2 dy) = \iint_D \left(\frac{\partial f_2}{\partial x} - \frac{\partial f_1}{\partial y} \right) dA$$

4. Gauss's theorem (divergence theorem) for $\mathbf{f} \in C^1$ over a finite 3D region V bounded by smooth closed surface S with outward pointing, unit surface normal, \mathbf{N} :

$$\iiint_V (\nabla \cdot \mathbf{f}) dV = \iint_S \mathbf{f} \cdot \mathbf{N} dS$$

5. Stokes's theorem for $\mathbf{f} \in C^1$ defined on a smooth surface S with unit surface normal \mathbf{N} and bounded by a positively oriented, closed curve Γ :

$$\iint_S (\nabla \times \mathbf{f}) \cdot d\mathbf{S} = \iint_S (\nabla \times \mathbf{f}) \cdot \mathbf{N} dS = \oint_{\Gamma} \mathbf{f} \cdot d\mathbf{r}$$



Chapter 1

Vectors and functions

Many mathematical properties possessed by functions of several variables are couched in geometric terms and with reference to elementary set theory. In this introductory chapter I will revisit some of the concepts that will be needed in later chapters. For example, vector calculus springs naturally from vector algebra so it is appropriate to begin the review with the latter topic. This is followed by a short review of elementary set theory, which will be referred to throughout the book and will indeed help establish many foundation concepts in both the differential and integral calculus. Coordinate systems and the notion of level sets are also discussed. Once again, both topics find application in differential and integral multivariable calculus, as well as in vector calculus.

It goes without saying that a review of single-variable functions is helpful. This begins in this chapter (Section 1.C), but continues in Chapters 2, 3 and 4 as needed.

To help appreciate the behaviour of multivariable functions defined on two or higher dimensional domains, It is useful to at least visualize their domains of definition. Sometimes, though, it is possible, as well as necessary, to visualize the entire graph of a function, or some approximation to it. Some people are more hard-wired to visual cues and visual information, while others are more comfortable with abstract ideas. Whatever your preference, being able to draw figures is always useful. Consequently, in this chapter we also review some basic 3D structures and show how to draw them using MATLAB®. Of course, other software will serve equally well. In the event of the reader being unable to access software solutions, there is included a subsection which may hopefully illustrate, by example, how one can obtain a picture of a region

or of a function graph directly from a mathematical formula or equation. Although it is not possible to offer a general procedure that works in all cases, some of the steps may be applicable in other instances.

1.A Some vector algebra essentials

Unit vectors in 3-space.

Let $a > 0$ be a scalar, and let

$$\begin{aligned}\mathbf{v} &= (\alpha, \beta, \gamma) \\ &= \alpha\mathbf{i} + \beta\mathbf{j} + \gamma\mathbf{k} \\ &= \alpha\mathbf{e}_1 + \beta\mathbf{e}_2 + \gamma\mathbf{e}_3\end{aligned}$$

be a vector in \mathbb{R}^3 (see Section 1.B) with x -, y -, and z -components α , β , and γ .

This vector has been written in the three most common forms appearing in current texts. The sets $\{\mathbf{i}, \mathbf{j}, \mathbf{k}\}$ and $\{\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3\}$ represent the same set of unit vectors in mutually orthogonal directions in \mathbb{R}^3 . The first form simply shows the components along the three orthogonal directions without reference to the unit vectors themselves, although the unit vectors and the coordinate system are implicit in this notation. The reader should be aware that we shall have occasion to refer to vectors using any of the three formats. The choice will depend on what is most convenient at that time without compromising understanding.

Multiplying a vector \mathbf{v} with a scalar will return a new vector with either the same direction if the scalar is positive or the opposite direction if the scalar is negative. In either case the resulting vector has different magnitude (Figure 1.1). This re-scaling will be a feature in Chapter 5 where we will need vectors of unit magnitude. For $a\mathbf{v}$, with $a \in \mathbb{R}$, to be a unit vector we must have

$$|a\mathbf{v}| = |a||\mathbf{v}| = a\sqrt{\alpha^2 + \beta^2 + \gamma^2} = 1, \text{ i.e., } a = \frac{1}{\sqrt{\alpha^2 + \beta^2 + \gamma^2}}.$$

Therefore, to construct a unit vector in the direction of a specific vector \mathbf{v} we simply divide \mathbf{v} by its length:

$$\mathbf{N} = \frac{\mathbf{v}}{|\mathbf{v}|}.$$

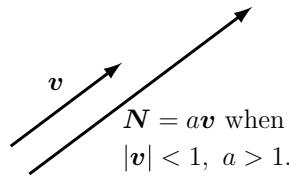


Figure 1.1 The unit vector.

The product of two vectors in 3-space.

Let \mathbf{u} and \mathbf{v} be two non-parallel vectors in \mathbb{R}^3 :

$$\mathbf{u} = (a_1, a_2, a_3) \quad \mathbf{v} = (b_1, b_2, b_3).$$

There are two particular product operations that we will utilize on many occasions. These are the vector and scalar products. From them very useful information can be extracted.

- (a) A vector *perpendicular to both* \mathbf{u} and \mathbf{v} is

$$\begin{aligned} \mathbf{w} = \mathbf{u} \times \mathbf{v} &= \begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ a_1 & a_2 & a_3 \\ b_1 & b_2 & b_3 \end{vmatrix} \\ &= (a_2 b_3 - a_3 b_2, a_3 b_1 - a_1 b_3, a_1 b_2 - a_2 b_1) \\ &= -\mathbf{v} \times \mathbf{u}. \end{aligned}$$

This is called the “vector” or “cross” product. Note that $\mathbf{u} \times \mathbf{v}$ is antiparallel to $\mathbf{v} \times \mathbf{u}$. The relationship between the three vectors is shown in Figure 1.5.

- (b) The magnitude of the vector (cross) product of two vectors

$$|\mathbf{u} \times \mathbf{v}| = |\mathbf{w}| = \left\| \begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ a_1 & a_2 & a_3 \\ b_1 & b_2 & b_3 \end{vmatrix} \right\| = \sqrt{(a_2 b_3 - a_3 b_2)^2 + \dots}$$

gives the *area* of a plane parallelogram whose side lengths are $|\mathbf{u}|$ and $|\mathbf{v}|$ (Figure 1.2).

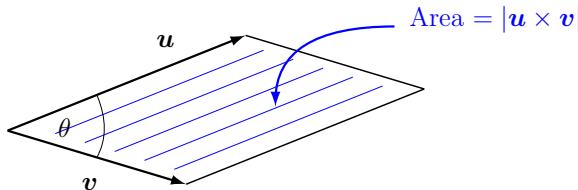


Figure 1.2 The $\mathbf{u} \times \mathbf{v}$ parallelogram.

The magnitude of the vector product is then given by

$$|\mathbf{w}| = |\mathbf{u}||\mathbf{v}| \sin \theta,$$

where θ is the angle between \mathbf{u} and \mathbf{v} lying in the plane defined by \mathbf{u} and \mathbf{v} .

- (c) The angle θ between the vectors \mathbf{u} and \mathbf{v} (Figure 1.3) can also be (and usually is) obtained from the “scalar” or “dot” product, defined as

$$\begin{aligned}\mathbf{u} \cdot \mathbf{v} &= a_1 b_1 + a_2 b_2 + a_3 b_3 \\ &= |\mathbf{u}||\mathbf{v}| \cos \theta.\end{aligned}$$

So we have

$$\cos \theta = \frac{\mathbf{u} \cdot \mathbf{v}}{|\mathbf{u}||\mathbf{v}|}.$$

If $\theta = 0$, then the vectors are parallel and $\mathbf{u} \cdot \mathbf{v} = |\mathbf{u}||\mathbf{v}|$. If $\theta = \frac{\pi}{2}$, then the vectors are orthogonal and $\mathbf{u} \cdot \mathbf{v} = 0$. For example, in 2(a) above, $\mathbf{w} \cdot \mathbf{u} = \mathbf{w} \cdot \mathbf{v} = 0$ as \mathbf{w} is orthogonal to both \mathbf{u} and \mathbf{v} .

We will make extensive use of these products in Chapter 2 (Sections 2.E and 2.G) and throughout Chapter 5.

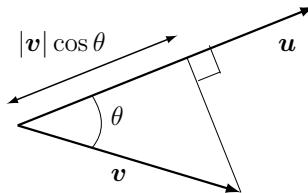


Figure 1.3 The projection of \mathbf{v} onto \mathbf{u} .

A plane in 3-space.

The equation of a plane in \mathbb{R}^3 , expressed mathematically as

$$P = \{(x, y, z) : ax + by + cz = d; a, b, c \text{ not all zero.}\}$$

can be determined knowing either

- (a) three non-collinear points on the plane; or
- (b) one point and two non-parallel vectors lying parallel to the plane.

Consider alternative 3(a). Let $\mathbf{x}_i = (x_i, y_i, z_i)$, $i = 1, 2, 3$, be the three points. Construct two vectors, \mathbf{u} and \mathbf{v} in the plane (Figure 1.4):

$$\begin{aligned}\mathbf{u} &= \mathbf{x}_2 - \mathbf{x}_1 = (x_2 - x_1, y_2 - y_1, z_2 - z_1) \\ \mathbf{v} &= \mathbf{x}_3 - \mathbf{x}_1 = (x_3 - x_1, y_3 - y_1, z_3 - z_1)\end{aligned}$$

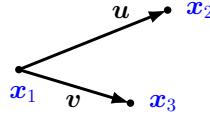


Figure 1.4 Construction of vectors \mathbf{u} and \mathbf{v} .

As long as \mathbf{x}_1 , \mathbf{x}_2 and \mathbf{x}_3 are not collinear, then \mathbf{u} and \mathbf{v} will not be superimposed or even parallel and

$$\mathbf{w} = \mathbf{u} \times \mathbf{v} = (\alpha, \beta, \gamma)$$

will be a vector normal (perpendicular) to \mathbf{u} and \mathbf{v} and thus normal to the plane in which the \mathbf{x}_i lie.

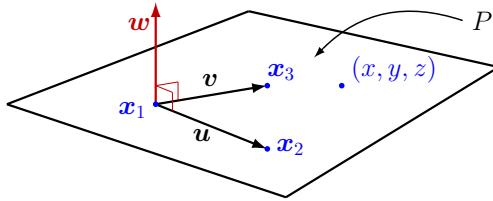


Figure 1.5 Construction of the plane P .

By convention, the direction of $\mathbf{w} = \mathbf{u} \times \mathbf{v}$ is given by the **right-hand rule**: *Using your right hand, point the index finger in the direction of \mathbf{u} and the middle finger in the direction of \mathbf{v} , then the thumb shows the direction of \mathbf{w} .* The vector product and its various geometric properties will play central roles in Sections 2.C and 5.D.

Now consider alternative 3(b):

Let (x, y, z) be any point in the plane P except for the given points (x_i, y_i, z_i) , $i = 1, 2, 3$. We construct the plane vector $(x - x_1, y - y_1, z - z_1)$ that joins this new point with the point \mathbf{x}_1 . Using concept 2(c) we have

$$\begin{aligned}\mathbf{w} \cdot (x - x_1, y - y_1, z - z_1) &= 0 \\ \implies \alpha(x - x_1) + \beta(y - y_1) + \gamma(z - z_1) &= 0 \\ \implies \alpha x + \beta y + \gamma z &= \mathcal{K}.\end{aligned}$$

The scalar product is thus instrumental in giving us the equation for the plane P with $a = \alpha$, $b = \beta$, $c = \gamma$, and $d = \mathcal{K}$.

Alternative 3(b) is actually a version of alternative 3(a), except we are here given \mathbf{u} and \mathbf{v} with which to create the orthogonal vector \mathbf{w} . This method of deriving the equation of a plane will be utilized in Sections 2.C and 5.D.

A line in 3-space.

The general equation of a line,

$$L = \left\{ (x, y, z) : \frac{x - x_0}{a} = \frac{y - y_0}{b} = \frac{z - z_0}{c} \right\}, \quad (1.1)$$

can be derived using analogous reasoning. We need to know either

- (a) two points on the line L , or
- (b) one point and one vector parallel to L .

Consider alternative 4(a).

Let (x_i, y_i, z_i) , $i = 0, 1$, be the two given points. Construct the vector \mathbf{u} directed from one point to the other:

$$\mathbf{u} = (x_1 - x_0, y_1 - y_0, z_1 - z_0) = (\alpha, \beta, \gamma).$$

As the two points lie in the straight line L so too must the vector \mathbf{u} as shown in Figure 1.6. Note that as in Point 3, this construction leads directly to alternative 4(b) where \mathbf{u} is given.

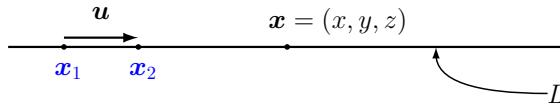
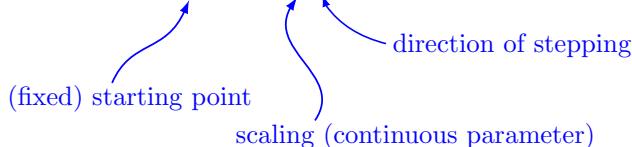


Figure 1.6 Vector \mathbf{u} parallel to line L .

Using either alternative set of information, any point (x, y, z) on L can be determined by simple vector addition,

$$\mathbf{x} = \mathbf{x}_0 + t \mathbf{u}, \quad t \in \mathbb{R}. \quad (1.2)$$



As indicated, this equation says that to determine any new point on the line we begin at a convenient starting point \mathbf{x}_0 and continue in the direction \mathbf{u} a distance determined by the scalar t .

This simple vector equation is equivalent to the general equation given in Equation (1.1) above. We get the latter by splitting Equation (1.2) into its components and solving each component equation for the common scalar variable t .

A particularly important feature of this equation, which is linear in the parameter t , emerges from the single-variable derivative of each component

$$\frac{dx}{dt} = \alpha, \quad \frac{dy}{dt} = \beta, \quad \frac{dz}{dt} = \gamma.$$

Combining these into a vector equation we have that

$$\frac{d\mathbf{x}}{dt} = \mathbf{u} \quad \text{— the tangent vector.}$$

This last result is elementary, but has important applications in Sections 5.A and 5.C, where the straight line concepts are generalized to the case of nonlinear curves.

The scalar triple product $\mathbf{u} \cdot (\mathbf{v} \times \mathbf{w})$.

Let $\mathbf{u}, \mathbf{v}, \mathbf{w}$ be three non-parallel vectors. These define the edges of a parallelepiped as shown in Figure 1.7.

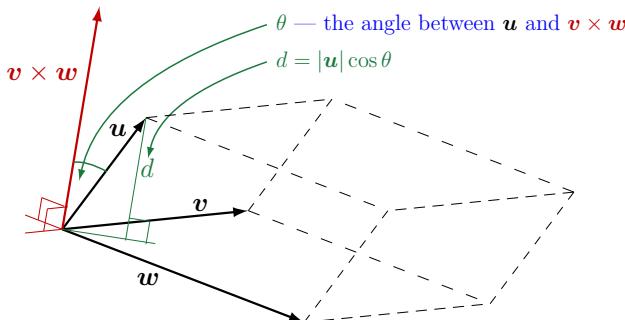


Figure 1.7 The $\mathbf{u}, \mathbf{v}, \mathbf{w}$ parallelepiped.

To form the scalar triple product, we first form the vector product of \mathbf{v} and \mathbf{w} , $\mathbf{v} \times \mathbf{w}$, and then form the scalar product of that result and \mathbf{u} . The magnitude of the scalar triple product, which is found using Point 2(a), is given by

$$\begin{aligned}
 |\mathbf{u} \cdot (\mathbf{v} \times \mathbf{w})| &= |\mathbf{u}| |\mathbf{v} \times \mathbf{w}| |\cos \theta| \\
 &= |\mathbf{v} \times \mathbf{w}| |\mathbf{u}| |\cos \theta| \\
 &= \underbrace{|\mathbf{v} \times \mathbf{w}|}_{\text{area of the } (\mathbf{v}, \mathbf{w}) \text{ parallelogram}} \underbrace{|\mathbf{u}| \cos \theta}_{\text{perpendicular height, } d}
 \end{aligned}$$

This gives the *volume of the parallelepiped* formed by the vectors \mathbf{u} , \mathbf{v} , and \mathbf{w} . Let

$$\begin{aligned}
 \mathbf{u} &= (a_1, a_2, a_3) = a_1\mathbf{i} + a_2\mathbf{j} + a_3\mathbf{k}, \\
 \mathbf{v} &= (b_1, b_2, b_3) = b_1\mathbf{i} + b_2\mathbf{j} + b_3\mathbf{k}, \text{ and} \\
 \mathbf{w} &= (c_1, c_2, c_3) = c_1\mathbf{i} + c_2\mathbf{j} + c_3\mathbf{k}.
 \end{aligned}$$

Using the determinant expression in 2(a) we have

$$\mathbf{v} \times \mathbf{w} = \begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ b_1 & b_2 & b_3 \\ c_1 & c_2 & c_3 \end{vmatrix} = (b_2c_3 - b_3c_2)\mathbf{i} + (b_3c_1 - b_1c_3)\mathbf{j} + (b_1c_2 - b_2c_1)\mathbf{k},$$

and therefore

$$\begin{aligned}
 \mathbf{u} \cdot (\mathbf{v} \times \mathbf{w}) &= (a_1\mathbf{i} + a_2\mathbf{j} + a_3\mathbf{k}) \cdot ((b_2c_3 - b_3c_2)\mathbf{i} + (b_3c_1 - b_1c_3)\mathbf{j} + (b_1c_2 - b_2c_1)\mathbf{k}) \\
 &= a_1(b_2c_3 - b_3c_2) + a_2(b_3c_1 - b_1c_3) + a_3(b_1c_2 - b_2c_1) \\
 &= \begin{vmatrix} a_1 & a_2 & a_3 \\ b_1 & b_2 & b_3 \\ c_1 & c_2 & c_3 \end{vmatrix}.
 \end{aligned}$$

The scalar triple product can be written succinctly in determinant form. Make the important note that the determinant notation *does not* mean that we take absolute values! So this result could be negative or positive. Remember, we are dealing here with vectors and angles.

The scalar triple product and the interpretation of its magnitude as the volume of a parallelepiped is a central feature of multiple integrals in Section 4.H.

The vector triple product $\mathbf{u} \times (\mathbf{v} \times \mathbf{w})$.

As we noted on Page 3, the vector $(\mathbf{v} \times \mathbf{w})$ is perpendicular to both \mathbf{v} and \mathbf{w} . Now suppose vector \mathbf{u} is not coplanar with \mathbf{v} and \mathbf{w} . What happens if we form the vector cross product of this vector with $(\mathbf{v} \times \mathbf{w})$?

Following the same line of reasoning, the result is a new vector which is perpendicular to both \mathbf{u} and $(\mathbf{v} \times \mathbf{w})$. Given that we have only three dimensions to play with (in \mathbb{R}^3), $\mathbf{u} \times (\mathbf{v} \times \mathbf{w})$ *must lie in the plane defined by the original vectors \mathbf{v} and \mathbf{w}* .

Consequently, $\mathbf{u} \times (\mathbf{v} \times \mathbf{w})$ must be a linear combination of \mathbf{v} and \mathbf{w} .

In fact, by twice applying the determinant formula for the cross product it can easily be verified that

$$\mathbf{u} \times (\mathbf{v} \times \mathbf{w}) = (\mathbf{u} \cdot \mathbf{w})\mathbf{v} - (\mathbf{u} \cdot \mathbf{v})\mathbf{w}.$$

Try it!

The subject of this section may seem elementary within the context of straight lines, but the concepts will prove to be quite important when generalized to multivariable scalar and vector function settings in which we deal with *tangent vectors* to more general differential curves.

1.B Introduction to sets

We begin this section with some useful definitions. The reader may refer back to their notes from linear algebra. Alternatively, a good reference is [16].

Definition 1.1

Given that $\mathbb{R}^n = \mathbb{R} \times \mathbb{R} \times \cdots \times \mathbb{R}$ is the set of all points \mathbf{x} having n independent real components, then $\mathbf{x} = (x_1, x_2, \dots, x_n) \in \mathbb{R}^n$ (where $x_i \in \mathbb{R}$) defines a point in n -dimensional Cartesian space.

Other notations in common use for a point in \mathbb{R}^n are: x , \vec{x} , \bar{x} , and \mathbf{x} .

■ Example 1.1:

The set of points in \mathbb{R}^2 is given as $\mathbb{R}^2 = \{(\mathbf{x}, \mathbf{y}) : \mathbf{x} \in \mathbb{R}, \mathbf{y} \in \mathbb{R}\}$; a single point with its defining pair of coordinates is shown in Figure 1.8.

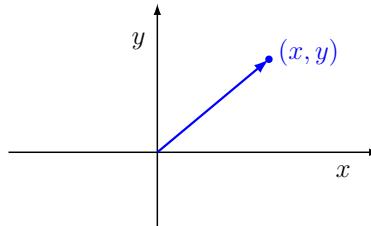


Figure 1.8 A point in 2D Cartesian space. ■

■ Example 1.2:

The set of points in \mathbb{R}^3 is given as $\mathbb{R}^3 = \{(\mathbf{x}, \mathbf{y}, \mathbf{z}) : \mathbf{x} \in \mathbb{R}, \mathbf{y} \in \mathbb{R}, \mathbf{z} \in \mathbb{R}\}$; a single point with its defining triad of coordinates is shown in Figure 1.9.

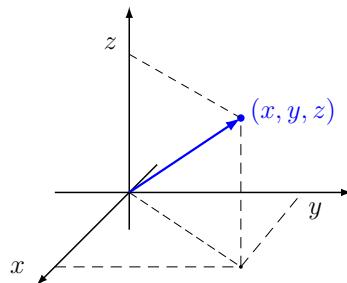


Figure 1.9 A point in 3D Cartesian space. ■

It is not possible to provide a picture of $\mathbf{x} \in \mathbb{R}^n$ for $n > 3$ (nor \mathbb{R}^n itself). However, there should be no cause for concern as points in \mathbb{R}^n behave the same as points in \mathbb{R}^2 and \mathbb{R}^3 . That is, they follow the same set of rules. So it is enough to be familiar with points and point operations in \mathbb{R}^2 and \mathbb{R}^3 , and then being able to generalize their properties. The most important point operations are listed below.

Vector algebra laws.

Let $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ and $\lambda \in \mathbb{R}$.

At its most basic description, \mathbb{R}^n is an example of a linear vector space which is characterized by the two properties of addition and scalar multiplication:

- (a) $\mathbf{x} + \mathbf{y} = (x_1 + y_1, x_2 + y_2, \dots, x_n + y_n) \in \mathbb{R}^n$.
- (b) $\lambda \mathbf{x} = (\lambda x_1, \lambda x_2, \dots, \lambda x_n) \in \mathbb{R}^n$.

As a direct generalization of the scalar product of 2(c), points in \mathbb{R}^n satisfy the so-called *inner product*,

- (c) $\mathbf{x} \cdot \mathbf{y} = x_1 y_1 + \dots + x_n y_n \in \mathbb{R}$ — \mathbb{R}^n is called an *inner product space*.

Finally, there are the following generalizations to \mathbb{R}^n of the two fundamental geometric measures:

- (d) $|\mathbf{x}| = \sqrt{\mathbf{x} \cdot \mathbf{x}} = \sqrt{x_1^2 + \dots + x_n^2}$ — the length of \mathbf{x} .
- (e) $|\mathbf{x} - \mathbf{y}| = \sqrt{(x_1 - y_1)^2 + \dots + (x_n - y_n)^2}$ — the distance between points.

With this distance property \mathbb{R}^n is also a so-called *metric space*, since the distance between points is one measure or metric that allows a geometric characterization of a space.

Using the above definitions one can prove (See Mastery Checks 1.1 and 1.2) some fundamental relations satisfied by position vectors in \mathbb{R}^n . These are useful in analysis to establish order relations between vector quantities.

- (f) Cauchy-Schwarz inequality:

$$|\mathbf{x} \cdot \mathbf{y}| \leq |\mathbf{x}| |\mathbf{y}|; \quad x_1 y_1 + \dots + x_n y_n \leq \sqrt{x_1^2 + \dots + x_n^2} \sqrt{y_1^2 + \dots + y_n^2}.$$
- (g) Triangle inequality: $|\mathbf{x} + \mathbf{y}| \leq |\mathbf{x}| + |\mathbf{y}|$.
- (h) The Cauchy-Schwarz inequality (f) means $-1 \leq \frac{\mathbf{x} \cdot \mathbf{y}}{|\mathbf{x}| |\mathbf{y}|} \leq 1$.

Just as in \mathbb{R}^2 and \mathbb{R}^3 , property (h) allows us to define an angle θ between vectors \mathbf{x} and \mathbf{y} in \mathbb{R}^n such that

$$\cos \theta = \frac{\mathbf{x} \cdot \mathbf{y}}{|\mathbf{x}| |\mathbf{y}|}.$$

☞ Mastery Check 1.1:

Prove the Cauchy-Schwarz inequality (f).

Hint: Consider the case $n = 2$ before tackling the general theorem. If you square each side of the inequality, you can show that the difference between the results for each side is ≥ 0 . This suggests that a proof might be possible by working backwards.



☞ Mastery Check 1.2:

Prove the triangle inequality (g).

Hint: This is achieved easily using (f). Begin by squaring $|\mathbf{x} + \mathbf{y}|$.



Points and sets

In general \mathbb{R}^n we can use property (e), involving the distance between two points, $|\mathbf{x} - \mathbf{y}|$, to help generalize the “interval” concept to \mathbb{R}^n .

In \mathbb{R} , the inequality $|x - a| < \epsilon$ (which is equivalent to saying $a - \epsilon < x < a + \epsilon$) describes the set of all $x \in \mathbb{R}$ which lie within ϵ of a . In \mathbb{R}^n we have the analogous case:

Definition 1.2

Given $\mathbf{a} \in \mathbb{R}^n$, an **open sphere** $S_r(\mathbf{a}) \subset \mathbb{R}^n$ centred at \mathbf{a} and of radius r is the set of all points $\mathbf{x} \in \mathbb{R}^n$ that satisfy $|\mathbf{x} - \mathbf{a}| < r$:

$$S_r(\mathbf{a}) = \{\mathbf{x} : |\mathbf{x} - \mathbf{a}| < r\}$$

for some $r \in \mathbb{R}$. That is, the set of all points $\mathbf{x} \in \mathbb{R}^n$ which are no further than r from the given point \mathbf{a} .

Remarks

- * The open sphere is non-empty, since it contains \mathbf{a} at least.
- * In this context \mathbf{a} is called the *centre*, and r the *radius* of the set.

The open sphere $S_r(\cdot)$ may now be used to define other point and set properties.

Definition 1.3

A point x is called

- an **interior point** of a set $M \subset \mathbb{R}^n$ if there is an open sphere $S_r(x) \subset M$ for some $r > 0$;
- an **exterior point** of a set $M \subset \mathbb{R}^n$ if there is an open sphere $S_r(x) \not\subset M$ for some $r > 0$;
- a **boundary point** of a set $M \subset \mathbb{R}^n$ if for any $r > 0$ (no matter how small), $S_r(x)$ contains points in M and points not in M .

These point definitions are illustrated in Figure 1.10.

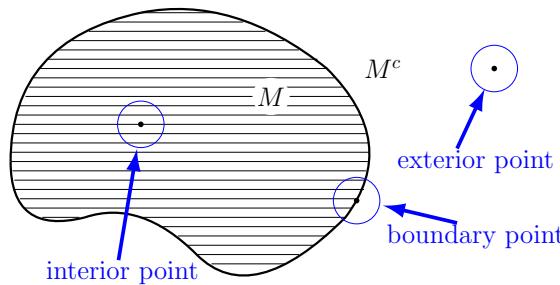


Figure 1.10 Interior, exterior, and boundary points to a set.

The reader should work through the following Mastery Checks to consolidate their understanding of these point definitions before going on to set-level concepts.

Mastery Check 1.3:

Let $M = \{(x, y) : 0 < x < 1, 0 < y < 1\}$

Draw a diagram of M on the Cartesian plane, showing the points

$P_1(\frac{5}{6}, \frac{5}{6})$, $P_2(1, \frac{1}{3})$, and $P_3(-1, 1)$.

Classify the points as interior, exterior, or boundary.



Mastery Check 1.4:

Let $M = \{(x, y) : \frac{(x - 1)^2}{4} + (y - 1)^2 \leq 1\}$.

Draw a diagram of M showing the points $P_1(2, 1)$, $P_2(3, 1)$, and $P_3(-1, 2)$.

Classify the points as interior, exterior, or boundary.



We now establish a framework within which to categorize points that possess common properties. We begin by grouping points according to Definition 1.3.

Definition 1.4

- The set of all interior points of a set M is called the **interior** of M , and denoted $\text{Int}(M)$:

$$\text{Int}(M) = \{\mathbf{x} : \mathbf{x} \in M \text{ and } S_r(\mathbf{x}) \subset M \text{ for some } r > 0\}.$$

- The set of all points not in M is called the **complement** of M , and denoted M^c :

$$M^c = \{\mathbf{x} : \mathbf{x} \notin M\}.$$

- The set of all boundary points of a set M is called the **boundary** of M , and denoted ∂M .

It follows from these definitions that $\text{Int}(M^c) \subset M^c$, and if \mathbf{x} is an exterior point to M , then $\mathbf{x} \in \text{Int}(M^c) \subset M^c$.

The concept introduced in the next definition will play an important role in our analysis of multivariable functions.

Definition 1.5

A set M is called **open** if it contains **only** interior points.

Accordingly, a set M is an open set if for every point $\mathbf{x} \in M$ a positive radius r can be found such that $S_r(\mathbf{x})$ contains only other points in M . Also, only under the specific condition of set M being open is $\text{Int}(M) = M$. Finally, an *open neighbourhood* of a point $\mathbf{a} \in M$ is an open set $W \in M$.

Although we can utilize the notion of an open set to define a closed set (see Supplementary problem 7), it proves useful to invoke an independent concept to define a closed set, that of so-called *limit points*. In this way we

can introduce a notion that is central to our forthcoming discussion on limits.

Definition 1.6

*A point \mathbf{a} of a set M is said to be a **limit point** of M if every open sphere $S_r(\mathbf{a})$ contains at least one point of M different from \mathbf{a} .*

This means that there are points in M that are arbitrarily close to \mathbf{a} . Hence, in approaching \mathbf{a} from within M we are always guaranteed to encounter other points in M .

Introducing limit points not only allows for a meaningful definition of a closed set, but it also allows one to readily prove a number of facts about closed sets, some of which are included as exercises at the end of this chapter. First the definition:

Definition 1.7

*A set M of \mathbb{R}^n is said to be **closed** if it contains all its limit points.*

And, intimately related to this definition is the concept of set *closure*. For our purposes we invoke the following definition.

Definition 1.8

*The **closure** of a set M , denoted \overline{M} , is the union of $\text{Int}(M)$ with its boundary:*

$$\overline{M} = \text{Int}(M) \cup \partial M = \{\mathbf{x} : \mathbf{x} \in \text{Int}(M) \text{ or } \mathbf{x} \in \partial M\}.$$

Alternatively, \overline{M} can be defined as the union of M and the set of all its limit points, L_M . It can be shown (see Supplementary problem 7) from this definition that (a) a closed set M is equal to its closure \overline{M} , and (b) that a set is closed if and only if it contains its boundary. Along this same line of thought, an alternative consequence of Definitions 1.5–1.8 is that the boundary of a set M contains those points that are in common with the closure of a set M and the closure of its complement, M^c . In other words

$$\partial M = \{\mathbf{x} : \mathbf{x} \in \overline{M} \cap \overline{M^c}\}.$$

The concept of set *boundedness* arises in both contexts of differential and integral calculus of multivariable functions. As the term suggests it essentially relates to a set being limited in geometric extent.

Definition 1.9

A set $M \subset \mathbb{R}^n$ is called **bounded** if there exists a $K \in \mathbb{R}$ such that $|\mathbf{x}| < K$ for all $\mathbf{x} \in M$.

Definition 1.5 is utilized in the definition of derivatives in Chapter 2, while Definitions 1.4–1.9 are invoked in Chapter 3 and 4, although they are also used implicitly elsewhere.

☞ Mastery Check 1.5:

- (1) Let $M = \{(x, y) : 0 < x < 1, 0 < y < 1\}$. What is the set ∂M ?
- (2) Let $M = \left\{ (x, y) : \frac{(x-1)^2}{4} + (y-1)^2 \leq 1 \right\}$. What is the set ∂M ?



The next definition is most useful when invoked together with function continuity to establish conditions that guarantee certain function behaviour. We shall see this employed in practice in Sections 3.B, 4.A and 4.D, but also 1.C.

Definition 1.10

A set $M \subset \mathbb{R}^n$ is called **compact** if it is both **closed** and **bounded**.

☞ Mastery Check 1.6:

For each of the sets M given below, answer the following questions:

Is M bounded? If it is, find a number K such that $|\mathbf{x}| < K$ for all $\mathbf{x} \in M$?

Is M compact? If it is not, write down the closure $\overline{M} = M \cup \partial M$.

Then draw a diagram showing M , ∂M , and K .

- (1) Let $M = \{(x, y) : 0 < x < 1, 0 < y < 1\}$.
- (2) Let $M = \left\{ (x, y) : \frac{(x-1)^2}{4} + (y-1)^2 \leq 1 \right\}$.



1.C Real-valued functions

Basic concepts and definitions.

In Chapters 2, 3, and 4, we focus attention almost exclusively on scalar-valued functions of many variables, while in Chapter 5 we extend the ideas to vector-valued functions. In both contexts the following introduction to fundamental properties of multi-valued functions is invaluable. To start, we introduce some more notation and a pictorial view of what functions do.

In single-variable calculus we have the following scenario:

Let $y = f(x)$. The “graph” of f is the set of ordered pairs $\{(x, f(x))\} \in \mathbb{R}^2$. This is shown graphically in Figure 1.11 where the independent variable x and dependent variable y are plotted on mutually orthogonal axes.

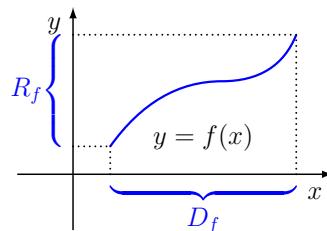


Figure 1.11 The Cartesian graph of $y = f(x)$.

This way of visualizing functions of one variable was introduced in the early 17th century by René Descartes [17], and is named the Cartesian representation in recognition. It is quite a useful means of illustrating function dependence and function properties, especially for functions of one or two variables.

It ceases to be as useful, however, for functions of more than two variables. For the latter cases one resorts to simply considering a set-mapping picture. For the case $y = f(x)$ this is a simple interval-to-interval map as shown in Figure 1.12.

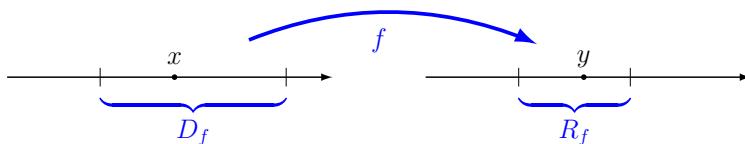


Figure 1.12 The set map of $D_f \subset \mathbb{R} \longrightarrow R_f \subset \mathbb{R}$.

For real-valued functions of many variables: $y = f(\mathbf{x}) = f(x_1, x_2, \dots, x_n)$, the corresponding illustrative representation is shown in Figure 1.13. The left-hand x -interval in the single-variable calculus case is replaced by a more general \mathbf{x} -region for the multivariable case.

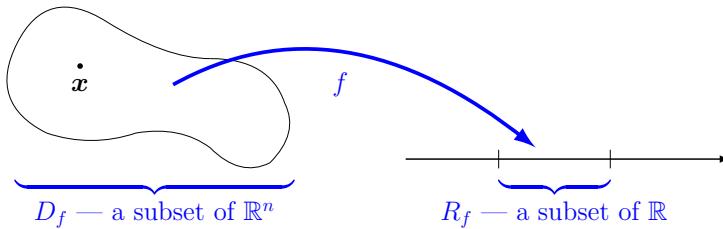


Figure 1.13 The set map of $D_f \subset \mathbb{R}^n \longrightarrow R_f \subset \mathbb{R}$.

Definition 1.11

Consider a real-valued function f of one or more variables, whose graph is the point set $\{(\mathbf{x}, f(\mathbf{x}))\}$.

The **domain** D_f of f is the set of all values of \mathbf{x} for which f is defined (that is, for which f makes sense).

The **range** R_f of f is the set of all possible values of $y = f(\mathbf{x})$ for all $\mathbf{x} \in D_f$.

In our multivariable setting $\mathbf{x} \in \mathbb{R}^n$ is the independent multivariable and $y \in \mathbb{R}$ is (still) the dependent variable. It is common to find the following terminology used in text books. The independent variable, here $\mathbf{x} \in \mathbb{R}^n$, is sometimes referred to as the *pre-image* of y , while the dependent variable, here $y \in \mathbb{R}$, is called the *image* under f . The function f is alternatively referred to as a *mapping* from one set to another, or an *operation* that takes a point, here \mathbf{x} , to y , or a *rule* that associates y with a point \mathbf{x} . As far as mathematical notation is concerned, the mapping under f is either described *pointwise*

$$f : \mathbf{x} \longmapsto y$$

or in a *set-wise* sense:

$$f : D_f \subseteq \mathbb{R}^n \longrightarrow R_f \subseteq \mathbb{R}.$$

Both references will be used in this book.

■ **Example 1.3:**

Consider $f(x, y) = \sqrt[3]{y - x^2}$.

Now, recall from single-variable calculus that $s = t^{1/3}$ is defined for all t , since

if $t > 0$, then $s > 0$;

if $t < 0$, then $s < 0$;

if $t = 0$ then $s = 0$.

In addition, we readily see that $y - x^2$ can take any real value. Combining these facts we deduce that f is defined everywhere. This implies that $D_f = \mathbb{R}^2$ and $R_f = \mathbb{R}$.



☞ **Mastery Check 1.7:**

Find the (implied) domain and range of the following functions:

1. $f(x) = \sqrt{16 - x^4}$;

4. $f(x, y) = \sqrt{9 - x^2 - y^2}$;

2. $f(x) = \frac{3 - x^2}{1 + x^2}$;

5. $f(x, y) = \ln(x - y)$;

3. $gd(x) = \sin^{-1}(\tanh(x))$;

6. $f(x, y, z) = \ln(|z - x^2 - y^2|)$.

(The function $gd(x)$ is known as the Gudermannian function.)



■ **Example 1.4:**

Suppose $f(x, y) = \sin^{-1}(x^2 + y^2)$.

Before considering this multivariable function, recall from single-variable calculus that within the intervals $-\frac{\pi}{2} \leq z \leq \frac{\pi}{2}$ and $-1 \leq w \leq 1$,

$$z = \sin^{-1} w \iff w = \sin z.$$

The graphs of these *inverse* functions are shown in Figure 1.14.

Note that in our case $w = x^2 + y^2 \geq 0$, and therefore so is $z \geq 0$. So,

$$|w| \leq 1 \implies |x^2 + y^2| \leq 1 \implies 0 \leq x^2 + y^2 \leq 1.$$

This defines the unit disc in the xy -plane. That is, D_f is the unit disc (the unit circle and its interior).

Similarly, $|z| \leq \frac{\pi}{2} \implies 0 \leq z \leq \frac{\pi}{2}$ since $z \geq 0$. So, $R_f = \left[0, \frac{\pi}{2}\right]$.

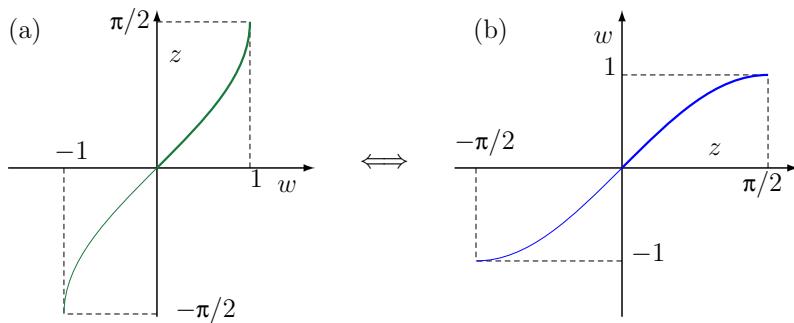


Figure 1.14 Graphs of the functions (a) $z = \sin^{-1} w$ and (b) $w = \sin z$.

■

☞ Mastery Check 1.8:

Consider the function $f(x, y) = \ln(2 - e^{x^2+y^2})$. Find the implied domain and range.



Although we will pay considerable attention to real-valued functions of several variables, we shall see in Chapter 5 that there is another important class of functions of several variables. These are vector-valued multivariable functions. Examples include:

(a) $\mathbf{f} : t \longmapsto \mathbf{y} \in \mathbb{R}^p$

— vector-valued functions of one real variable, t .

$$\mathbf{f}(t) = (f_1(t), f_2(t), \dots, f_p(t))$$

(b) $\mathbf{f} : \mathbf{x} \in \mathbb{R}^n \longmapsto \mathbf{y} \in \mathbb{R}^m$

— vector-valued functions of a vector variable
(several real variables).

$$\mathbf{f}(\mathbf{x}) = (f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_m(\mathbf{x}))$$

Limits and continuity.

In the next chapter we introduce and explore the concept of partial differentiation. In the lead up to that discussion it will be necessary to explain a number of concepts we shall then take for granted. Most importantly there is the notion of *function continuity*. For multivariable functions this will be

discussed in detail in Section 2.B, but we can set the stage here with a short review of the subject as it relates to functions of one variable.

Function continuity is defined in terms of limiting processes. Mention has already been made of limit points of closed sets. We said that a point a is a limit point if any open sphere centred on a , no matter how small in radius, contains points other than a .

Similarly, segments of the real line possess the property that any open *interval* I , no matter how small, centred on a point a , contain points x in I different from a . The real line and any of its finite segments are therefore said to be *complete*: containing no gaps. This conjures up the notion of a set *continuum*, moving smoothly from one real value to another, never meeting any holes.

This notion gives critical meaning to the formalism $x \rightarrow a$ as the process of approaching a real value a along the real line. To be even more precise, we specify $x \rightarrow a^-$ and $x \rightarrow a^+$ as meaning the respective approaches to a along the real line from “below” a ($x < a$) and from “above” a ($x > a$).

Now with thought given to single-variable functions defined on a domain $D_f \subset \mathbb{R}$, the different approaches $x \rightarrow a^-$ and $x \rightarrow a^+$ for $a, x \in D_f$ can have all manner of implications for the function. Assuming $a, x \in D_f$ we define the *process* of taking a limit of a function, which we denote either by

$$\lim_{x \rightarrow a^-} f(x), \lim_{x \rightarrow a^+} f(x), \text{ or } \lim_{x \rightarrow a} f(x)$$

as considering the sequence of values f progressively takes as $x \rightarrow a^-$, $x \rightarrow a^+$, or in their combination. These considerations are of course separate to the question of what value f actually takes at a . To summarize all of these ideas we have the following definition.

Definition 1.12

A function is said to be **continuous** at an interior point a of its domain $D_f \subset \mathbb{R}$ if

$$\lim_{x \rightarrow a} f(x) = f(a).$$

If either the equality is not satisfied, or the limit fails to exist, then f is said to be *discontinuous* at a .

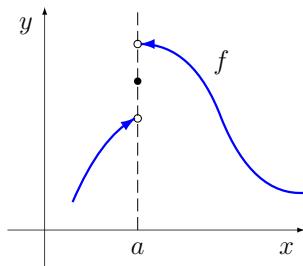
To reiterate, in the context of functions of a single variable the above limit is equivalent to the requirement that the limits approaching a from below ($x < a$) and from above ($x > a$) exist and are equal. That is, the single expression $\lim_{x \rightarrow a} f(x)$ means that

$$\underbrace{\lim_{x \rightarrow a^-} f(x)}_{\substack{\text{left-hand-side} \\ \text{limit of } f(x)}} = \underbrace{\lim_{x \rightarrow a^+} f(x)}_{\substack{\text{right-hand-side} \\ \text{limit of } f(x)}} = A.$$

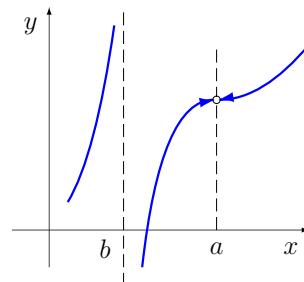
If the above equality is not satisfied we say that the limit does not exist. Definition 1.12 then also stipulates that for continuity the common limiting value, the aforementioned A , must also equal the value of the function f at $x = a$, $f(a)$.

The following example demonstrates graphically some different circumstances under which a limit of a function exists or does not exist, and how these relate to the left and right limits. Note the sole case of the function value actually being specified (solid dot) in the left-most graph in Figure 1.15. Is either function continuous at a ?

■ Example 1.5:



The limit does not exist
at $x = a$.



The limit *does* exist
at $x = a$, not at b .

Figure 1.15 When does a limit exist?

From the definition and subsequent discussion we are led to an important theorem of pointwise continuity.

Theorem 1.1

A function f is continuous at an interior point a if and only if

$$\lim_{x \rightarrow a^-} f(x) = \lim_{x \rightarrow a^+} f(x) = f(a).$$

In slightly more practical mathematical language the statement of Definition 1.12 and Theorem 1.1 can be expressed by the following:

$$0 < |x - a| < \delta, \quad x, a \in D_f \quad \implies \quad |f(x) - A| < \epsilon \text{ for some } \delta = \delta(\epsilon).$$

Graphically, this limit definition can be represented as in Figure 1.16 below.

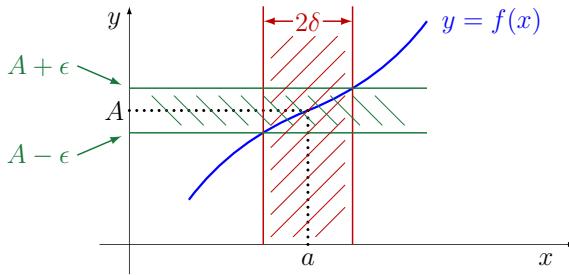


Figure 1.16 The ϵ - δ condition.

The concepts developed above will be employed in Section 2.B. For now, this pointwise concept can be extended to the entire function domain.

Definition 1.13

*A function f that is continuous at every point of its domain $D_f \subset \mathbb{R}$ is said to be **continuous** over that domain.*

We denote by $C(D_f)$ the set of all continuous functions defined on the domain D_f .

Still on the single-variable case, we will have need in Sections 3.C, 4.A and 4.B of the following important theorem, which combines the concepts of function continuity and domain compactness to give an important result.

Theorem 1.2

Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be continuous on a closed and bounded interval $D_f \subset \mathbb{R}$. Then f attains an absolute maximum and an absolute minimum value in D_f . That is, there exist values $x_L \in D_f$ and $x_U \in D_f$ such that

$$f(x_L) \leq f(x) \leq f(x_U) \text{ for all } x \in D_f.$$

A moment's thought and possibly some simple sketches will make this theorem intuitively obvious. The self-evidence of the theorem, however, does not diminish its significance.

We end this section with a short catalogue of well-established results that can assist us in evaluating limit processes for both single-variable and multi-variable functions, where the latter cases comprise single-variable functions. Three of the squeeze relations listed below are featured in Figure 1.17.

Some useful standard limits:

- * $\lim_{x \rightarrow 0} \frac{\sin x}{x} = 1.$
- * $\lim_{x \rightarrow \infty} \frac{\ln x}{x^\alpha} = 0 \quad \text{for constant } \alpha > 0.$
- * $\lim_{x \rightarrow 0} \frac{\tan x}{x} = 1.$
- * $\lim_{x \rightarrow 0} x^\alpha \ln x = 0 \quad \text{for constant } \alpha > 0.$
- * $\lim_{t \rightarrow 0} \frac{\sin(xt)}{t} = x.$
- * $\lim_{t \rightarrow 0} \frac{\cos t - 1}{t} = 0.$

Some useful squeeze relations:

- (a) $\sin x < x < \tan x, \quad 0 < x < \frac{\pi}{2}.$
- (b) $x < e^x - 1 < \frac{x}{1-x}, \quad 0 < x < 1.$
- (c) $\frac{x}{1+x} < \ln(1+x) < x, \quad \text{for } x > -1, \neq 0.$
- (d) $e^x > 1 + x \quad \forall x \neq 0.$

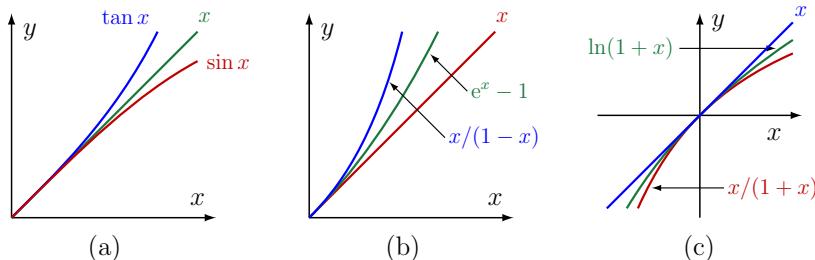


Figure 1.17 Graphs of various squeeze relations relative to $y = x$.

1.D Coordinate systems

Up until now we have represented points in \mathbb{R}^2 and \mathbb{R}^3 in terms of *Cartesian* coordinates, (x, y) as in Example 1.1 and (x, y, z) as in Example 1.2. However, problems arise that are better described in other coordinate systems. Such problems arise in both the differential and integral calculus (Sections 3.E, 4.E, and 4.H) and are usually associated with the geometry of the region under consideration. The most common coordinate systems that we will encounter are the *polar coordinate* system in \mathbb{R}^2 , and the *cylindrical* and *spherical coordinate* systems in \mathbb{R}^3 . Note that there are other standard systems that can be useful in specific cases (see [15]) and even non-standard systems may be needed to solve some problems (see Section 4.E).

There are three general features to note. First, the 2D Cartesian and polar coordinate systems have the same origin. Similarly, the 3D Cartesian and cylindrical or spherical coordinate systems have a common origin. Second, the non-Cartesian coordinates are designed to uniquely identify and represent every point in \mathbb{R}^2 or \mathbb{R}^3 , as do their Cartesian counterparts. That is, these coordinate systems span the whole of \mathbb{R}^2 and \mathbb{R}^3 , respectively. Finally, the individual coordinate variables within a given non-Cartesian system are independent of each other, just as the individual Cartesian coordinates are independent variables in the Cartesian system.

Polar coordinates

Consider an arbitrary point P in the plane with Cartesian coordinates (x, y) . P 's distance from the origin is

$$r = \sqrt{x^2 + y^2},$$

while the angle between P 's position vector $\mathbf{r} = (x, y)$ and the x -axis is given by

$$\tan \theta = \frac{y}{x}.$$

The unique inverse relation is given by the pair of equations

$$\left. \begin{array}{l} x = r \cos \theta \\ y = r \sin \theta \end{array} \right\} \quad \text{for } 0 \leq \theta \leq 2\pi.$$

Thus, every point in \mathbb{R}^2 can be uniquely represented by the pair of so-called *polar coordinates* (r, θ) defined on the domain $[0, \infty) \times [0, 2\pi]$. The relationship between the two coordinate representations is shown in Figure 1.18(a).

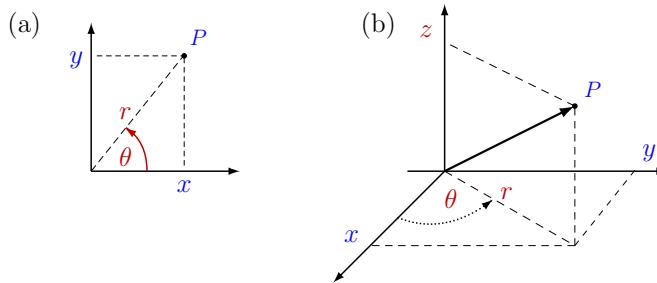


Figure 1.18 (a) 2D polar and (b) 3D cylindrical coordinates.

The distance D between two arbitrary points $P(x_1, y_1)$ and $Q(x_2, y_2)$ is then

$$\begin{aligned} D &= [(x_1 - x_2)^2 + (y_1 - y_2)^2]^{1/2} \\ &= [r_1^2 + r_2^2 - 2r_1 r_2 \cos(\theta_1 - \theta_2)]^{1/2}, \end{aligned}$$

where $x_i = r_i \cos \theta_i$ and $y_i = r_i \sin \theta_i$, $i = 1, 2$.

Cylindrical polar coordinates

An arbitrary point P in 3D is defined by Cartesian coordinates (x, y, z) . The preceding case of plane polar coordinates is thus easily generalized to *cylindrical polar coordinates* in 3D by the inclusion of the Cartesian coordinate z to account for the third dimension.

We therefore have the relations

$$x = r \cos \theta, \quad y = r \sin \theta, \quad z = z.$$

Figure 1.18(b) shows the point P represented by the two alternative coordinate systems (x, y, z) and (r, θ, z) . The distance between any two points P and Q generalizes to

$$\begin{aligned} D &= [(x_1 - x_2)^2 + (y_1 - y_2)^2 + (z_1 - z_2)^2]^{1/2} \\ &= [r_1^2 + r_2^2 - 2r_1 r_2 \cos(\theta_1 - \theta_2) + (z_1 - z_2)^2]^{1/2}. \end{aligned}$$

Spherical polar coordinates

The second generalization to 3D of polar coordinates is the *spherical polar* coordinate system. This is based on the notion of defining a point on a sphere in terms of latitude and longitude angles. To be precise, an arbitrary point P in 3D with Cartesian coordinates (x, y, z) is identified by the triplet of

independent variables (ρ, ϕ, θ) defined by

$$\left. \begin{aligned} x &= \rho \sin \phi \cos \theta, \\ y &= \rho \sin \phi \sin \theta, \\ z &= \rho \cos \phi, \end{aligned} \right\} \quad \begin{aligned} 0 &\leq \rho < \infty, \\ 0 &\leq \phi \leq \pi, \\ 0 &\leq \theta \leq 2\pi, \end{aligned}$$

with the inverse relations

$$\begin{aligned} \rho^2 &= x^2 + y^2 + z^2 = r^2 + z^2, \quad \tan \theta = \frac{y}{x}, \\ \cos \phi &= \frac{z}{(r^2 + z^2)^{1/2}} \quad \text{or} \quad \sin \phi = \frac{r}{\rho}. \end{aligned}$$

Figure 1.19 illustrates how the variables are related geometrically. The origins of the angles ϕ and θ as z -axis and x -axis, respectively, are also indicated.

The distance between two arbitrary points P and Q in 3D is now expressed

$$D^2 = \rho_1^2 + \rho_2^2 - 2\rho_1\rho_2 \cos(\phi_1 - \phi_2) - 2\rho_1\rho_2 \sin \phi_1 \sin \phi_2 (\cos(\theta_1 - \theta_2) - 1).$$

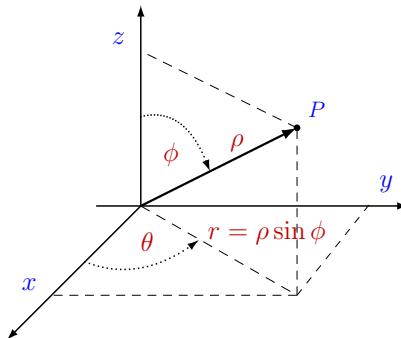


Figure 1.19 3D spherical coordinates.

1.E Drawing or visualizing surfaces in \mathbb{R}^3

Throughout the book and indeed throughout the subject generally we will need to recognize, but also to sketch or otherwise visualize, areas and volumes embedded in \mathbb{R}^2 and \mathbb{R}^3 , respectively. The ability to conceptualize regions in 2D and 3D makes the task of understanding multivariable function behaviour easier. Also, in the case of integration in Chapter 4, it simplifies the task of

establishing integration limits that define the boundaries of a region. Lastly, being able to visualize structures in 3D can be helpful when checking the reasonableness of possible solutions to mathematical exercises.

Most of the examples and exercises to follow utilize MATLAB[®] in the visualization of various surfaces (using the “surf” function). However, the reader with access to other graphing software should be able to translate the programming information shown below into relevant code for their own graphing tool.

As an alternative the first example that follows shows how to piece together a picture of a graph with little or no reliance on software. A more detailed discussion of this manual approach (restricted to functions of a single variable) can be found in Chapter 5 of [1].

■ Example 1.6:

We shall consider here the graph of the function

$$f(x, y) = \frac{4x}{1 + x^2 + y^2} \quad \text{for } (x, y) \in \mathbb{R}^2.$$

This function features in an exercise in a later chapter. For now we are just interested in determining the form taken by the function’s graph,

$$G = \{(x, y, z) : (x, y) \in \mathbb{R}^2, z = f(x, y)\}.$$

In the steps that follow we will in essence dissect the function, and with the pieces we obtain we will build up a picture of the graph.

Step 1: The first thing to note is the domain of definition. What you would be looking for are the limits on the independent variables as well as possible points where the function is not defined. In our case, the function is defined everywhere so the domain is the entire xy -plane.

Step 2: The second thing to do is to look for any zeros of the function. That is, we look for intercept points in the domain at which the function takes the value zero. Here, $f = 0$ when $x = 0$, that is, at all points along the y -axis.

Step 3: We now look for any symmetry. We note that the function is odd in x but even in y . The symmetry in x means that for any fixed y — which means taking a cross-section of the graph parallel to the x -axis — howsoever the graph appears for $x > 0$, it will be inverted in the xy -plane for $x < 0$.

The symmetry in y means that for any fixed x (that is, a cross-section parallel to the y -axis) the graph will look the same on the left of $y = 0$ as on the right. Note, however, that because of the oddness in x , the graph will sit *above* the xy -plane for $x > 0$, but *below* the plane for $x < 0$.

So far, we have the impressions shown in Figure 1.20.

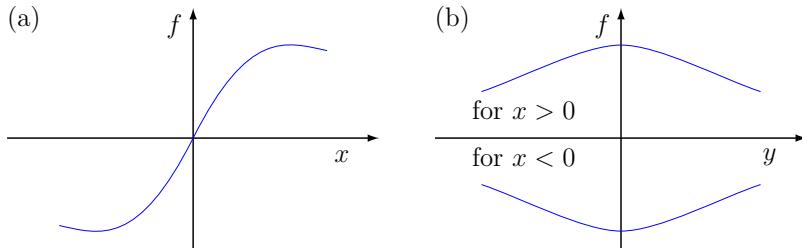


Figure 1.20 The function f is odd in x (a), but even in y (b).

Step 4: It is often instructive to look at small values of x and y .

Supposing $|x|$ to be very small compared to both 1 and y we see that f behaves as

$$f \approx \left(\frac{4}{1+y^2} \right) x \quad \text{as } |x| \rightarrow 0.$$

(The symbol \approx means “approximately equal to” and indicates very close correspondence.) We see that f behaves linearly with respect to x , with a coefficient that depends on y . This is consistent with the fact that f is odd in x (Figure 1.20(a)).

Next, supposing $|y|$ to be very small compared to 1 and x , the function will tend to

$$f \approx \frac{4x}{1+x^2} \quad \text{as } |y| \rightarrow 0.$$

That is, f behaves very much as a constant with respect to y , dependent only on the given value of x . Again this is consistent with our finding that f is even in y : we would expect the function to be approximately constant for small $|y|$, as in Figure 1.20(b).

Step 5: We now consider the behaviour of f for large values of x and y , the “asymptotic” behaviour of f .

Again fixing y and taking $|x|$ very large compared with either 1 or y , the function will tend to behave as

$$f \approx \frac{4}{x} \rightarrow 0 \quad \text{as } x \rightarrow \pm\infty.$$

The approach to zero will depend on the sign of x : approaching zero from *above* for $x > 0$ and from *below* for $x < 0$. (See Figure 1.21(a).)

On the other hand, fixing x instead and taking $|y|$ very large compared to either 1 or x , we find that

$$f \approx \frac{4x}{y^2} \longrightarrow 0 \quad \text{as} \quad y \longrightarrow \pm\infty.$$

So the function again approaches zero. Note again the sign difference for positive and negative x (Figure 1.21(b)).

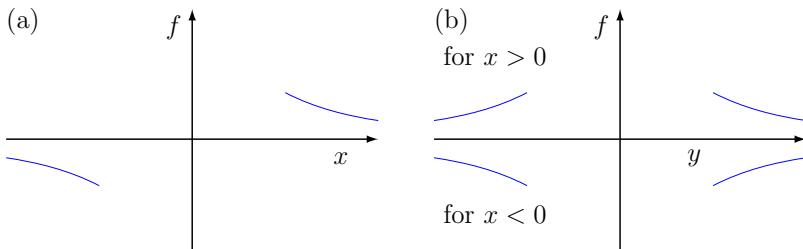


Figure 1.21 The behaviour of f for large $|x|$ (a), and large $|y|$ (b).

Now let's see if we can deduce something from this incomplete picture.

The function *is* zero along $x = 0$, and it *tends* to zero at large $|x|$ (and large $|y|$) and is nonzero in between.

Since the function does not have any singular behaviour anywhere in \mathbb{R}^2 , we can conclude that there must be at least one point along $x > 0$ where the function peaks at some positive value, and similarly there must be at least one point along $x < 0$ where the function bottoms at some negative value. We are thus led to ...

Step 6: Investigate f for maxima and minima. This step we will leave until we have at our disposal the differentiation tools developed in the next chapter and applied in Chapter 3 (see Mastery Check 3.7). But the above information is enough to put together the sketch shown in our final Figure 1.22.

If we look at the curves running parallel to the x -axis we see that the extrema (large $|x|$) match the predictions in Figure 1.21(a), while the middle sections agree with the curves in Figure 1.20(a). Regarding the curves parallel to the y -axis, the extremes (large $|y|$) concur with Figure 1.21(b), while the sections crossing the x -axis agree with the lines in Figure 1.20(b).

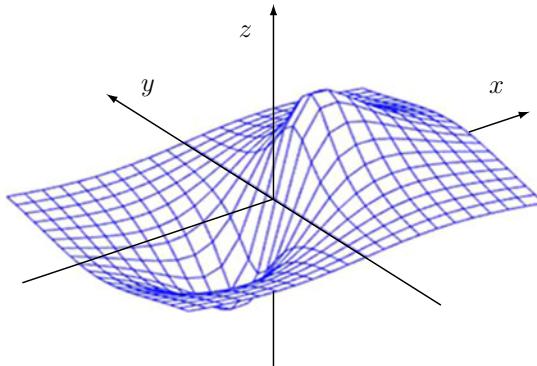


Figure 1.22 Putting it all together. ■

This next example similarly shows how one can visualize a surface without using graphing software.

Example 1.7:

Consider $S = \{(x, y, z) : x^2 + y^2 + z^2 = a^2, a > 0\}$. This is a surface in \mathbb{R}^3 , Figure 1.23; it is a surface because there exists a relation between the three variables (x, y, z) . They are no longer completely independent: one variable can be considered a function of the other two.

Now set $z = 0$. This simplifies to the subset satisfying $x^2 + y^2 = a^2$ which is a curve (circle) in the xy -plane. Note that these two equations for the three variables, which is equivalent to setting two conditions on the three variables, generate a curve in \mathbb{R}^2 .

A consistent interpretation is that of the intersection of two surfaces: The plane $z = 0$ and the sphere S giving rise to the subset of points the surfaces have in common — the circle of radius a in the xy -plane.

Suppose that $a > 2$, say, in S . Then setting

$$\begin{aligned} z = 0 &\implies x^2 + y^2 = a^2, \\ z = 1 &\implies x^2 + y^2 = a^2 - 1 < a^2, \\ z = 2 &\implies x^2 + y^2 = a^2 - 4 < a^2 - 1 < a^2, \\ z = a &\implies x^2 + y^2 = a^2 - a^2 = 0 \iff x = y = 0. \end{aligned}$$

These are examples of *level sets* defining circles in the xy -plane. We will come back to discuss these in detail in Section 1.F.

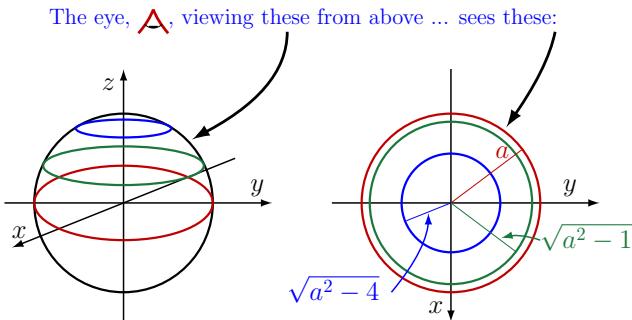


Figure 1.23 The sphere of radius a and a few of its level sets.



■ Example 1.8:

The same example as Example 1.7, but now using MATLAB[®]: Figure 1.24.

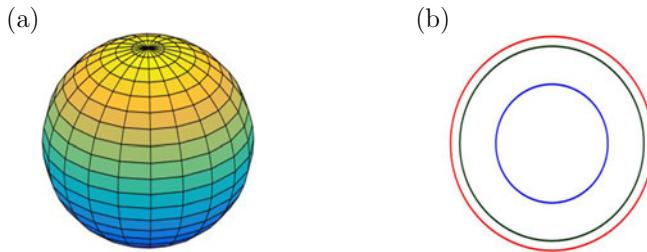


Figure 1.24 The sphere of radius 2.4 (a) and its level sets (b).

This version of the MATLAB[®] code produces figures without axes.

```
figure(1)
[X,Y,Z]=sphere; % generates three 21X21 matrices for a unit
sphere
X=2.4*X; Y=2.4*Y; Z=2.4*Z; % sphere now has radius 2.4
surf(X,Y,Z), axis tight, axis off
figure(2)
[X,Y,Z]=sphere(100); % better definition with 100 segments
X=2.4*X; Y=2.4*Y; Z=2.4*Z;
map=[1,0,0;0,0.2,0;0,0,1]; % colours are red, dark green, blue
contour(X,Y,Z,[0,1,2]), axis off
colormap(map);
```



☞ Mastery Check 1.9:

Set up your own matrices in MATLAB[®] for the `surf` plot, as follows (or otherwise), and draw the sphere again:

```
r=2.4; theta=linspace(0,2*pi,21); phi=linspace(0,pi,21);
X=r*sin(phi')*cos(theta);
Y=r*sin(phi')*sin(theta);
Z=r*cos(phi')*ones(1,21);
surf(X,Y,Z)
```



■ Example 1.9:

A circular cone: $S = \{(x, y, z) : z^2 = x^2 + y^2, -1 \leq x, y \leq 1\}$

Let $0 \leq \theta \leq 2\pi$ and $0 \leq r \leq 1$, $\begin{cases} x = r \cos \theta \\ y = r \sin \theta \end{cases} \implies z = \pm r$.

This example illustrates why care should be exercised in cases involving squares. It is easy to forget that there is some ambiguity when taking the square root: see Figure 1.25.

The MATLAB[®] default figure format has tick marks with labels on the axes which suit most purposes, and there are simple functions for producing labels for the axes themselves, as shown in the sample code that follows.

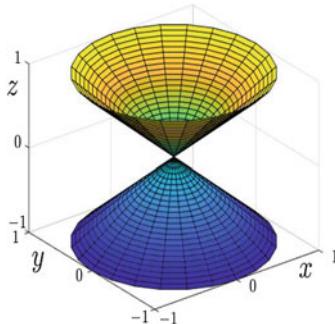


Figure 1.25 The graph of $z = \pm\sqrt{x^2 + y^2}$.

The MATLAB[®] code:

```
theta=linspace(0,2*pi,25);
r=linspace(0,1,25); % 25 intervals from 0 to 1
x=cos(theta')*r; y=sin(theta')*r; Z=sqrt(x.^2+y.^2);
```

```
surf(x,y,Z), hold on, surf(x,y,-Z)
xlabel('x'), ylabel('y'), zlabel('z')
```

However, the following may be used in place of the last line of code to produce clearer labels, given here for the x -axis, and easily adapted for the other two axes.

```
ax=gca; % Get the axis handle, call it 'ax'
xticks([-1,0,1]); % This sets the tick marks
% Place the new larger axis label at (0.6,-1.2,-1)
text(0.6,-1.2,-1,'$x$', 'interpreter', 'latex', 'fontsize', 24)
% Overwrite the tick labels with blanks
ax.XAxis.TickLabels=' ', ' ', ' ';
% Place the new x-tick label '-1' at (-1,-1.2,-1)
text(-1,-1.2,-1,'-1', 'fontsize', 16)
text(0,-1.2,-1,'0', 'fontsize', 16)
text(1,-1.2,-1,'1', 'fontsize', 16)
```



☞ Mastery Check 1.10:

Consider these conic sections for the case in Example 1.9:

1. Set $y = 0 \implies z^2 = x^2 \implies z = \pm x$ — a pair of straight lines.
2. Set $y = 0.5 \implies z^2 = x^2 + 0.25 \implies (z - x)(z + x) = 0.5^2$
— a hyperbola.
3. Set $z = 0.6 \implies x^2 + y^2 = 0.6^2$ — a circle, radius 0.6.
4. Set $z = y + 0.25 \implies x^2 + y^2 = y^2 + 0.5y + 0.0625$
 $\implies y = 2x^2 - 0.125$ — a parabola.

Each of these curves may be generated using MATLAB®.

Your task is to add in turn the groupings of lines of the following code to the end of the code for the cone, then use the “Rotate 3D” button on the MATLAB® figure to view the curves in space.

```
% the line pair
title('$z^2=x^2+y^2$', 'interpreter', 'latex')
x=linspace(-1,1,11); y=zeros(11,1); Z=x*ones(1,11);
surf(x,y,Z)
% the hyperbola
title('$z^2=x^2+y^2$', 'interpreter', 'latex')
```

```

x=linspace(-1,1,11)'; y=0.5*ones(11,1); Z=x*ones(1,11);
surf(x,y,Z)
% the circle
title('$z^2=x^2+y^2$', 'interpreter', 'latex')
x=linspace(-1,1,11); y=x; Z=0.6*ones(11);
surf(x,y,Z)
% the parabola
title('$z^2=x^2+y^2$', 'interpreter', 'latex')
x=linspace(-1,1,11)'; y=x; Z=(y+0.25)*ones(1,11);
surf(x,y,Z)

```

Figures 1.26(a) and (b) are for the last one of these, at two different aspects (see if you can get these plots):

`view(0,90)`, and `view(-90,4)`

Conic sections: The parabola.

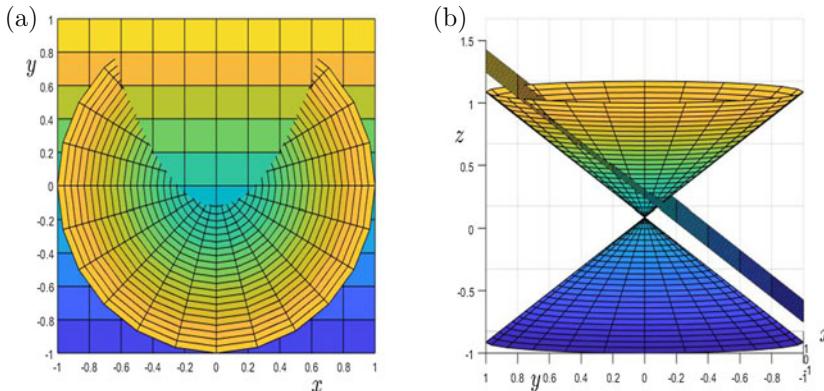


Figure 1.26 Two views of the intersection of the graphs of $z^2 = x^2 + y^2$ and $y = x$: (a) top view, (b) side view.

■ Example 1.10:

The hyperbolic paraboloid of Figure 1.27: $S = \{(x, y, z) : z = 1 + x^2 - y^2\}$.

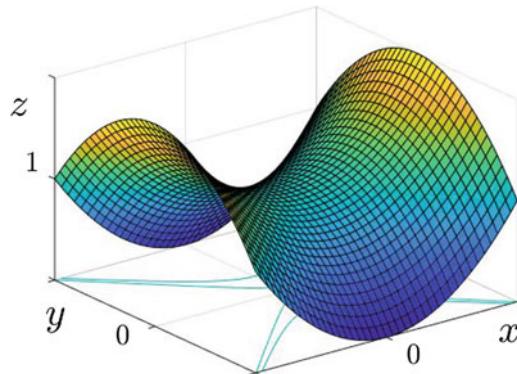


Figure 1.27 The graph of $z = 1 + x^2 - y^2$ with two level sets.

Setting $z = \text{constant}$ will give a hyperbola. For example, set $z = 0.9$, $y^2 - x^2 = 0.1$, or set $z = 0.99$, $y^2 - x^2 = 0.01$.

These curves shown in the figure are projections onto the xy -plane ($z = 0$). The MATLAB[®] code:

```
X=linspace(-1,1,41)*ones(1,41); Y=ones(41,1)*linspace(-1,1,41);
Z=1+X.^2-Y.^2; surf(X,Y,Z)
xticks([-1,0,1]), yticks([-1,0,1]), zticks([0,1,2])
text(0.7,-1.2,-0.1,'$x$', 'interpreter', 'latex', 'fontsize', 28)
text(-1.3,0.7,-0.1,'$y$', 'interpreter', 'latex', 'fontsize', 28)
text(-1.2,1.2,1.7,'$z$', 'interpreter', 'latex', 'fontsize', 28)
hold on
contour(X,Y,Z, [0.9,0.99])
ax=gca;
ax.XAxis.TickLabels=' ', ' ', ' ', ' ';
ax.YAxis.TickLabels=' ', ' ', ' ', ' ';
ax.ZAxis.TickLabels=' ', ' ', ' ', ' ';
text(0,-1.2,-0.1,'$0$', 'interpreter', 'latex', 'fontsize', 20)
text(-1.3,0,0,'$0$', 'interpreter', 'latex', 'fontsize', 20)
text(-1,1.3,1,'$1$', 'interpreter', 'latex', 'fontsize', 20)
```

☞ Mastery Check 1.11:

The hyperbolic paraboloid: $S = \{(x, y, z) : z = 1 + x^2 - y^2\}$.

Setting $y = 0$ will give the parabola $z = 1 + x^2$. Setting $x = 0$ will give the parabola $z = 1 - y^2$.

Produce 3D plots for each of these.



■ Example 1.11:

How to draw an ellipsoid in MATLAB[®].

The following code produces the graph of the ellipsoid $\frac{x^2}{4^2} + \frac{y^2}{3^2} + \frac{z^2}{2^2} = 1$ shown in Figure 1.28, whose semi-axes are $a = 4$, $b = 3$, and $c = 2$. The MATLAB[®] code uses an elliptical parametrization,

$$x = 4 \sin \phi \cos \theta, \quad y = 3 \sin \phi \sin \theta, \quad z = 2 \cos \phi,$$

which bears some similarities to the spherical coordinate transformation (Page 26) of the Cartesian coordinates.

This plot has been made partially transparent using the ‘‘FaceAlpha’’ property. The `line` commands are used to set x -, y -, z -axes.

```
theta=linspace(0,2*pi,41); phi=linspace(0,pi,41);
X=4*sin(phi')*cos(theta); Y=3*sin(phi')*sin(theta);
Z=2*cos(phi')*ones(1,41);
surf(X,Y,Z,'FaceAlpha',0.6), hold on
line([-5,0,0;5,0,0],[0,-5,0;0,5,0],[0,0,-5;0,0,5],...
'color','k','linewidth',2)
axis off
```

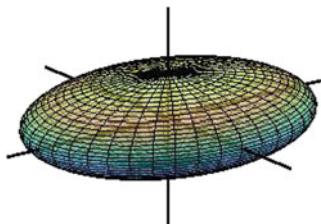


Figure 1.28 The graph of $\frac{x^2}{4^2} + \frac{y^2}{3^2} + \frac{z^2}{2^2} = 1$.



1.F Level sets

In many areas of mathematics, physics, and engineering, there arise equations of the form

$$f(\mathbf{x}) = f(x_1, \dots, x_n) = c$$

where c is a constant real scalar. Although this expression appears to place interest on the outcome of a function evaluation, it actually raises the question of what values (points) of the argument, \mathbf{x} , give rise to the specific value of f . This so-called *inverse* problem leads to the notion of a *level set*.

Definition 1.14

The set of all points $\mathbf{x} \in \mathbb{R}^n$ which give the constant value c for f is called a **level set**, L , or more precisely, the **c -level set** of f :

$$L = \{\mathbf{x} \in D_f : f(\mathbf{x}) = c\}$$

Students often confuse a level set with the graph of a function. But as we have said, this concerns the specific subset of points $\mathbf{x} = (x_1, x_2, \dots) \in D_f$, rather than what comes out of a function evaluation (except, of course, for the value c !).

■ Example 1.12:

In \mathbb{R}^2 , $\mathbf{x} = (x, y)$, and

$$L = \{(x, y) \in D_f : f(x, y) = c\}$$

is a *level curve*. For example, the level set $f(x, y) = x^2 + y^2 = 4$ is the set of points on the circle in the xy -plane (in \mathbb{R}^2) with centre $(0, 0)$ and radius 2. In contrast, the graph of $z = f(x, y)$ is a 3D object in \mathbb{R}^3 .

In \mathbb{R}^3 , $\mathbf{x} = (x, y, z)$, and

$$L = \{(x, y, z) \in D_f : f(x, y, z) = c\}$$

is a *level surface*. For example, the level set $f(x, y, z) = x^2 + y^2 + z^2 = 4$ is the set of points on the surface of the sphere in \mathbb{R}^3 with centre $(0, 0, 0)$ and radius 2.



By construction, determining the level set from the expression $f(\mathbf{x}) = c$ is an inverse problem. Sometimes when f is given explicitly, as in Example 1.12, we are able to “solve” for one variable in terms of the others. In the above 2D example $x^2 + y^2 = 4$, we obtain $y = \sqrt{4 - x^2}$, a semicircle curve passing through $(0, 2)$, and $y = -\sqrt{4 - x^2}$, a semicircle curve passing through $(0, -2)$.

The next example shows how a 3D surface can give rise to level sets in different 2D planes.

■ Example 1.13:

Consider the circular paraboloid of Figure 1.29: $S = \{(x, y, z) : z = x^2 + y^2; -1 \leq x, y \leq 1\}$.

Horizontal level sets occur at fixed values of z . The paraboloid is shown in Figure 1.29(a) together with the level sets for $z = r^2$, $r = 0.5, 0.6, 0.7, 0.8, 0.9$.

Vertical level sets occur at fixed values of x or y . Shown in Figure 1.29(b) is the level set for $y = 0$.

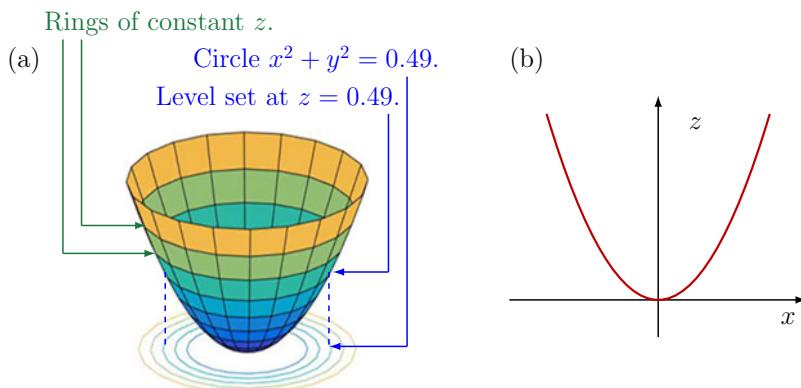


Figure 1.29 The paraboloid $z = x^2 + y^2$ (a), with level sets of two kinds in (a) and (b).

The basic figure in (a) was generated using this MATLAB[®] code:

```
theta=linspace(0,2*pi,21); % 20 intervals around the circle
r=linspace(0,1,11); % 10 intervals from 0 to 1
x=cos(theta')*r; y=sin(theta')*r; z=x.^2+y.^2;
% x, y and z are 21X11 matrices
surf(x,y,z), axis off
hold on % wait for the contours corresponding to
```

```
% the level sets at z=r^2; r=.5, .6, .7, .8, .9
contour(x,y,z,[0.25,0.36,0.49,0.64,0.81])
hold off.
```



Example 1.14 draws our attention to the fact that an expression involving only two variables may still describe a function in 3D, although any level sets may have a simple form.

■ Example 1.14:

Consider the parabolic cylinder of Figure 1.30: $S = \{(x, y, z) : z = 4x^2\}$.

Even though there is no specific y -dependence, this is a surface in 3D as opposed to the 2D parabola of the last example which we found by setting $y = 0$. The lack of a y -dependence means that the shape persists for all values of y . The curves of constant z (the level sets) are therefore lines parallel to the y -axis.

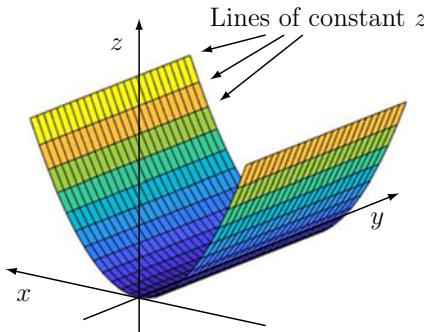


Figure 1.30 The graph of $z = 4x^2$.

The figure may be generated using MATLAB[®] code similar to this:

```
x=linspace(-2,2,25); y=linspace(-2,2,25);
[X,Y]=meshgrid(x,y); Z=4*X.^2;
surf(X,Y,Z)
text(1.0,-2.4,-0.2,'$x$', 'interpreter', 'latex', 'fontsize', 32)
text(2.4,1.0,-0.2,'$y$', 'interpreter', 'latex', 'fontsize', 32)
text(-2.4,-2.4,14,'$z$', 'interpreter', 'latex', 'fontsize', 32)
xticks([]), yticks([]), zticks([0,4,8,12,16])
view(36.5,22)
```



A level set of the form $g(x, y, z) = \text{a constant}$ is equivalent to declaring a function of two variables, which can in principle at least be plotted in a 3D diagram as we see in the next example. Note that the points in this 3D set lie in the domain of g , *not the graph of g* .

■ **Example 1.15:**

Lastly, consider the surfaces shown in Figure 1.31: $w = g(x, y, z) = z - f(x, y) = z - x^2 + y/5 = k$. Let k take the values 1, 2, 3, 4.

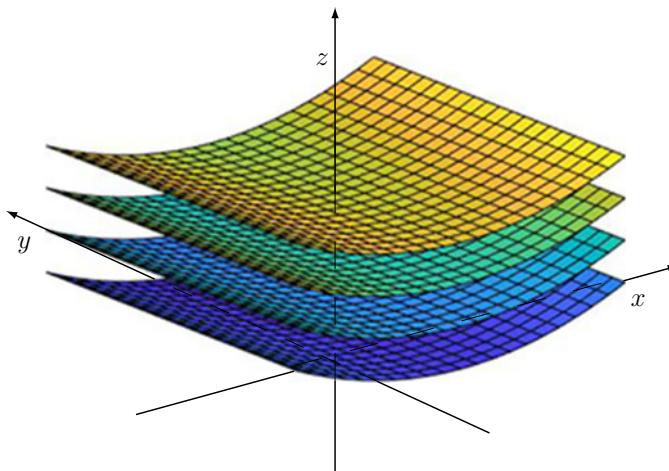


Figure 1.31 The level sets of $z - x^2 + y/5 = k$ for $k = 1, 2, 3, 4$.

The MATLAB[®] code:

```
X=ones(1,21)’*linspace(-1,1,21); Y=X’;
Z=zeros(21,21);
for k=1:4
Z=X.^2+k-Y/5;
surf(X,Y,Z), hold on
end
xticks([-1,0,1]); yticks([-1,0,1]); zticks([0,1,2,3,4,5,6]);
text(0.6,-1.2,0,’$x$’,’interpreter’,’latex’,’fontsize’,24)
text(-1.2,0.7,-0.1,’$y$’,’interpreter’,’latex’,’fontsize’,24)
text(-1.2,1.1,4.5,’$z$’,’interpreter’,’latex’,’fontsize’,24)
```

```
axis([-1 1 -1 1 0 6])
hold off
```



☞ Mastery Check 1.12:

Consider the function $f(x, y) = \frac{x^3 - y^2}{x^2 - x}$ in the domain $D_f = \{(x, y) : x \geq 0, y \geq 0\}$.

For what values of the constant k can the level set $f(x, y) = k$ be solved for y as a function of x throughout D_f ? Choose two such values for k , and use MATLAB[®] to plot the resulting curves for $0 \leq x \leq 2$ on the same figure.



☞ Mastery Check 1.13:

Consider the function $f(x, y, z) = x^2 - y^2 + z^2$ defined for $D_f = \{(x, y, z) : |x| \leq 2, |y| < \infty, z \geq 0\}$.

Show that we may solve the level set $f(x, y, z) = k$ for z in terms of x and y for all $k \geq 4$. For the cases $k = 4, 6, 8$, draw the graphs for $|x| \leq 2, |y| \leq 2$, on the same figure.



☞ Mastery Check 1.14:

Consider $F(w, x, y, z) = 36x^2 - 36y^2 - 4z^2 - 9w^2 = 0$. This is a level set in \mathbb{R}^4 .

Describe in words the precise subsets of this level set that arise from setting $x = 0, 1, 2, 3$.

Use MATLAB[®] to draw the graphs for cases $x = 1, 2$, on the same axes. You may wish to use the ‘‘FaceAlpha’’ property (see Page 37).



Although in this book we shall generally work with problems and examples which admit closed-form inversion of the level set equation, the reader should acknowledge that this will not always be possible in practical situations.

1.G Supplementary problems

Section 1.A

- Suppose three vectors \mathbf{u} , \mathbf{v} , \mathbf{w} are such that $\mathbf{u} + \mathbf{v} + \mathbf{w} = \mathbf{0}$. Show that $\mathbf{u} \times \mathbf{v} = \mathbf{v} \times \mathbf{w} = \mathbf{w} \times \mathbf{u}$. With the aid of a diagram describe what this result means.
- Let $\mathbf{x}_i = (x_i, y_i, z_i)$, $i = 0, 1, 2, 3$ be four non-coplanar points in \mathbb{R}^3 , and let vectors $\mathbf{u}_i = \mathbf{x}_i - \mathbf{x}_0$, $i = 1, 2, 3$, be edges of the tetrahedron formed by those points. Consider the four vectors $\mathbf{a}, \mathbf{b}, \mathbf{c}, \mathbf{d}$ with magnitudes in turn equal to twice the area of the four faces of the tetrahedron, and directions outwards and normal to those faces. Express these vectors in terms of the \mathbf{u}_i and hence show that $\mathbf{a} + \mathbf{b} + \mathbf{c} + \mathbf{d} = \mathbf{0}$.

(In Figure 1.32, regard the point (x_2, y_2, z_2) as being to the rear, without any loss of generality. Normal vectors \mathbf{a} and \mathbf{d} are shown.)

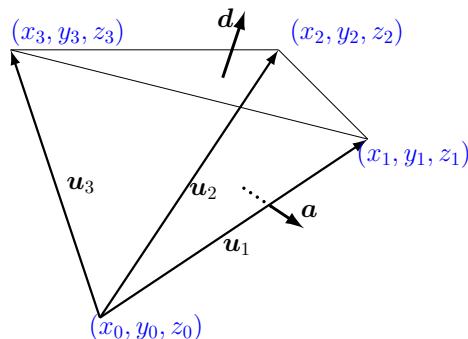


Figure 1.32 Four vectors in 3D space.

- Let $\mathbf{u} = (a_1, a_2, a_3)$, $\mathbf{v} = (b_1, b_2, b_3)$, and $\mathbf{w} = (c_1, c_2, c_3)$ be vectors in 3-D space.
 - Show that $\mathbf{u} \cdot (\mathbf{v} \times \mathbf{w}) = (\mathbf{u} \times \mathbf{v}) \cdot \mathbf{w}$. That is, show that in the scalar triple product the “dot” and the “cross” can change places.
 - Show that $\mathbf{u} \times (\mathbf{v} \times \mathbf{w}) = (\mathbf{u} \cdot \mathbf{w})\mathbf{v} - (\mathbf{u} \cdot \mathbf{v})\mathbf{w}$.

Section 1.B

4. Consider the three points $\mathbf{a} = (1, -1, 2, 2)$, $\mathbf{b} = (3, 1, -1, 1)$, $\mathbf{c} = (-2, 0, 2, -1)$ in \mathbb{R}^4 .
 - (a) Find the distances $|\mathbf{b} - \mathbf{a}|$, $|\mathbf{c} - \mathbf{b}|$, and $|\mathbf{a} - \mathbf{c}|$.
 - (b) Do either of the points \mathbf{b} and \mathbf{c} lie inside the open sphere $S_{3\sqrt{2}}(\mathbf{a})$?
 - (c) Find the angle θ between $\mathbf{b} - \mathbf{a}$ and $\mathbf{c} - \mathbf{a}$.
5. In \mathbb{R}^n , for what values of k is $\mathbf{b} = (k, k, k, \dots)$ inside $S_r(\mathbf{a})$ when $\mathbf{a} = (1, 1, 1, \dots)$?
6. Sketch the following regions and determine their boundaries. Also establish whether the regions are open or closed or neither.
 - (a) $\{(x, y) : |x| + |y| < 1\}$.
 - (b) $\{(x, y) : \max(|x|, |y|) < 1\}$.
 - (c) $\{(x, y) : x^2 \leq y \leq \sqrt{x}\}$.
 - (d) $\{(x, y) : 1 < (x - 1)^2 + (y + 1)^2 \leq 2\}$.
 - (e) $\{(x, y) : |x + 3y| \leq 3\}$.
7. Prove the following statements about sets:
 - (a) The boundary of a set M is a closed set.
 - (b) A set M is closed $\iff M = \overline{M}$.
 - (c) A set M is open $\iff M = \text{Int}(M)$.
 - (d) A set M is closed \iff its complement M^c is open.
 - (e) The union of any number of open sets is open, and any finite intersection of open sets is open.
 - (f) The intersection of any number of closed sets is closed, and any finite union of closed sets is closed.
8. If $\mathbf{x} = (x, y, z)$ and $|\mathbf{x}| = \sqrt{x^2 + y^2 + z^2}$, show that

$$\max(|x|, |y|, |z|) \leq |\mathbf{x}| \leq |x| + |y| + |z| \leq \sqrt{3}|\mathbf{x}| \leq 3 \max(|x|, |y|, |z|)$$

Section 1.C

9. Describe the implied domain D_f for each of the following functions:
Is it (i) closed?, (ii) finite?, (iii) compact?

- (a) $f(x, y, z) = \arcsin(x^2 + y^2 + z^2)$.
- (b) $f(x, y, z) = \arcsin(x^2 + y^2)$.
- (c) $f(x, y) = \arctan(x^2 + y^2)$.
- (d) $f(x, y, z) = \ln(1 - |x + y + z|)$.

Section 1.D

10. MATLAB[®] is able to plot functions expressed in 2D polar coordinates using a plotter called `ezpolar`. Use this function to plot the curves $r = 2 \sin n\theta$, $0 \leq \theta \leq 2\pi$, for $n = 2, 3, 4$, on separate graphs.
11. (a) Express the 2D polar function $r = 2a \sin \theta$, $0 \leq \theta \leq \pi$, a constant, in Cartesian coordinates, and describe the curve.
(b) What is the curve defined by $r = 2a \sin \theta$, $0 \leq \theta \leq 2\pi$?
12. A curve in \mathbb{R}^3 can be sufficiently prescribed in terms of one independent parameter.

Represent all points on the curve of intersection of the plane $ax + by + cz = d$ with the unit sphere centred at the point (x_0, y_0, z_0) in terms of spherical coordinates, if the plane passes through the sphere's centre.

Give conditions that must be satisfied by the constants a, b, c, d for the intersection to be possible and your representation valid.

Hints: Set the origin for the spherical coordinate system to be at the sphere's centre. Select the longitude angle θ as the independent variable, then the latitude angle ϕ becomes a function of θ .

13. The surfaces $z^2 = 2x^2 + 2y^2$ and $z = y - 1$ intersect to create a closed curve. Use cylindrical coordinates to represent points on this curve.
14. A surface in \mathbb{R}^3 can be sufficiently prescribed in terms of two independent parameters.

Represent all points on the plane $ax + by + cz = d$ within and coinciding with the sphere $x^2 + y^2 + z^2 = R^2$ in terms of spherical coordinates defined with respect to the origin at the sphere's centre.

Give conditions that must be satisfied by the constants a, b, c, d, R for the intersection to be possible and your representation valid.

15. The surfaces $z = 1 + x + y$, $z = 0$, and $x^2 + y^2 = 1$ bound a closed volume, V . Represent all points in V and on its boundary in spherical coordinates. Be mindful of the domains of the respective independent variables.

Section 1.E

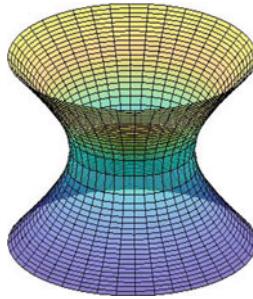


Figure 1.33 The graph of $\frac{x^2}{a^2} + \frac{y^2}{b^2} - \frac{z^2}{c^2} = 1$.

16. Figure 1.33 shows the elliptic hyperbola of one sheet,

$$\frac{x^2}{a^2} + \frac{y^2}{b^2} - \frac{z^2}{c^2} = 1.$$

Use MATLAB[®] to reproduce this plot.

17. Use MATLAB[®] to reproduce Figure 1.34 which shows the graphs of the cylinders $x^2 - 2x + y^2 = 0$ and $z^2 - 2x = 0$ for $0 \leq x \leq 2$, $-2 \leq y \leq 2$, $-2 \leq z \leq 2$, plotted on the same axes.

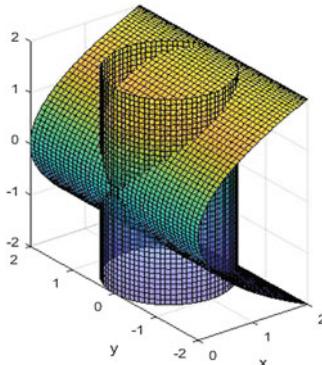


Figure 1.34 Two intersecting cylinders.

18. Use MATLAB[®] to draw a sketch graph of the cone $x^2 + y^2 = xz$ for $-4 \leq z \leq 4$.
19. Describe and sketch the graphs of the following functions using ideas analogous to those described in Example 1.6.
- $f(x, y) = \sqrt{1 - x^2 - y^2}$, for $x^2 + y^2 \leq 1$.
 - $f(x, y) = \sqrt{1 - y^2}$, for $|y| \leq 1, x \in \mathbb{R}$.
 - $f(x, y) = 1 - (x^2 + y^2)$, for $(x, y) \in \mathbb{R}^2$.
 - $f(x, y) = \frac{1 - \sqrt{x^4 + y^4}}{x^2 + y^2 + 1}$, for $(x, y) \in \mathbb{R}^2$.

Section 1.F

20. Consider the function $f(x, y, z) = \sin(xy) + \cos(yz)$, $-\frac{\pi}{2} \leq x, y, z \leq \frac{\pi}{2}$. Use MATLAB[®] to draw the graphs of the level sets $x = 1, x = 2; y = 1, y = 2$; and $z = 1, z = 2$.
21. Sketch the level curves of the following functions and determine the conditions for the allowed constant values of f .
- $f(x, y) = x^2 + y^2 - 4x + 2y$.
 - $f(x, y) = x^2y$.
 - $f(x, y) = x\sqrt{y^2 + 1}$.
 - $f(x, y) = \frac{x}{x^2 + y^2}$.
 - $f(x, y) = \frac{4x}{1 + x^2 + y^2}$.



Chapter 2

Differentiation of multivariable functions

What would calculus be without derivatives? In this chapter we cover the theory of the differential calculus, beginning with the limit concept as it pertains to functions of many variables and their derivatives. Considerable emphasis is placed on the geometric meaning of partial derivatives and of differentiability in general. The discussion also covers higher-order derivatives and introduces the new concept of the *gradient* of a function.

The focus of attention is then directed to composite functions, the chain rule of partial differentiation, and to implicit functions. The dedication of considerable space to these latter topics is motivated partly by their level of complication which is much greater than in the case of functions of one variable, and partly by the simple fact that they are commonly encountered in practice.

Having established these foundation concepts we shall put them to practical use in Chapter 3, where we discuss a fair assortment of applications.

2.A The derivative

Differentiation is all about *limiting processes* and *linear approximations*. To prepare us for that discussion we turn to the 1D case for inspiration. Many of the necessary basic features and results of limits of functions of one variable appearing below were covered in Chapter 1. The reader may wish to refer back to Section 1.C for details.

Definition 2.1

Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a continuous function over an open interval domain. That is, in more mathematical notation, let $f \in C(D_f)$, where the domain of f $D_f \subset \mathbb{R}$ is open.

Let x_0 and $x_0 + h \in D_f$. If $\lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h}$ exists, it is called the (first) derivative of f at x_0 , and we write

$$\frac{dy}{dx} \Big|_{x_0} = f'(x_0) = \lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h},$$

where $\frac{dy}{dx} \Big|_{x_0}$ is the **slope** of the tangent line to the **graph** of f at the point $(x_0, y_0 = f(x_0))$.

An equivalent definition:

$$\frac{dy}{dx} \Big|_{x_0} = \lim_{x \rightarrow x_0} \frac{f(x_0) - f(x)}{x_0 - x}; \quad x_0, x \in D_f.$$

Note that either explicitly or implicitly we have assumed the following properties which are essential criteria for the existence of the limit:

(i) $x_0, x_0 + h \in D_f$

— both points belong to the domain, D_f

(ii) $\lim_{h \rightarrow 0^-} f(x_0 + h) = \lim_{h \rightarrow 0^+} f(x_0 + h) = f(x_0)$

— the left limit equals the right limit which equals the function value at x_0

Thus, $\lim_{h \rightarrow 0} f(x_0 + h)$ exists and is equal to $f(x_0)$.

(iii) Similarly, $\lim_{h \rightarrow 0^-} \frac{f(x_0 + h) - f(x_0)}{h} = \lim_{h \rightarrow 0^+} \frac{f(x_0 + h) - f(x_0)}{h}$

— the left and right limits of these ratios exist and are equal

The reasons why these conditions are essential for the definition of a derivative are demonstrated in the following two classic examples of problem cases.

■ **Example 2.1:**

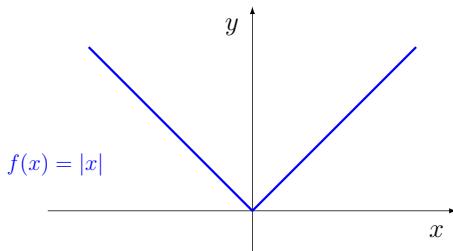


Figure 2.1 A function not everywhere differentiable.

For the function shown in Figure 2.1, (i) and (ii) are satisfied everywhere. At $x = 0$, however, although the left and right limits of (iii) exist, they are not equal, implying that no derivative exists there. Everywhere else (iii) is satisfied. ■

■ **Example 2.2:**

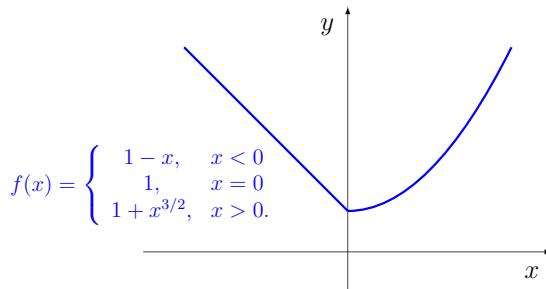


Figure 2.2 Another function not everywhere differentiable.

For the function shown in Figure 2.2, the only problem appears at $x = 0$. Conditions (i) and (ii) are satisfied, but in the case of condition (iii) we have that

$$\lim_{h \rightarrow 0^-} \frac{f(0+h) - f(0)}{h} = \lim_{h \rightarrow 0^-} \frac{(1-h) - 1}{h} = -1.$$

$$\lim_{h \rightarrow 0^+} \frac{f(0+h) - f(0)}{h} = \lim_{h \rightarrow 0^+} \frac{(1+h^{3/2}) - 1}{h} = \lim_{h \rightarrow 0^+} h^{1/2} = 0.$$

That is, the left limit is not equal to the right limit, implying that no derivative exists there. ■

Point (ii) is the definition of continuity at x_0 (Section 1.C). A function for which any of the equalities is not satisfied is said to be discontinuous at x_0 . What we are now saying is that continuity is a necessary but not a sufficient condition for differentiability. Functions for which (iii) is not satisfied at any point, x_0 , such as those of the foregoing examples, are said to be *singular* at that point.

Now let us apply what we have learnt for a function of one variable to the case of a function of two variables. The most obvious analogous expression of a limit generalized to some function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ of two variables is:

$$\lim_{P_1 \rightarrow P_0} \frac{f(x_0, y_0) - f(x_1, y_1)}{\sqrt{(x_1 - x_0)^2 + (y_1 - y_0)^2}}. \quad (2.1)$$

If this limit exists, should we call it “the” derivative of f ? Alongside this question we also need to ask what are the generalizations of criteria (i)–(iii) to \mathbb{R}^2 (or $\mathbb{R}^3, \dots, \mathbb{R}^n$)?

The graphical foundation for the limit expression (2.1) is shown in Figure 2.3. The things to note are, firstly, the graph of f is suspended in 3D; secondly, the domain D_f lies in the xy -plane; thirdly, the points P_0 and P_1 in D_f give rise to values z_0 and z_1 , respectively; and finally, the line in the domain joining P_0 and P_1 traces out the black curve in the graph of f .

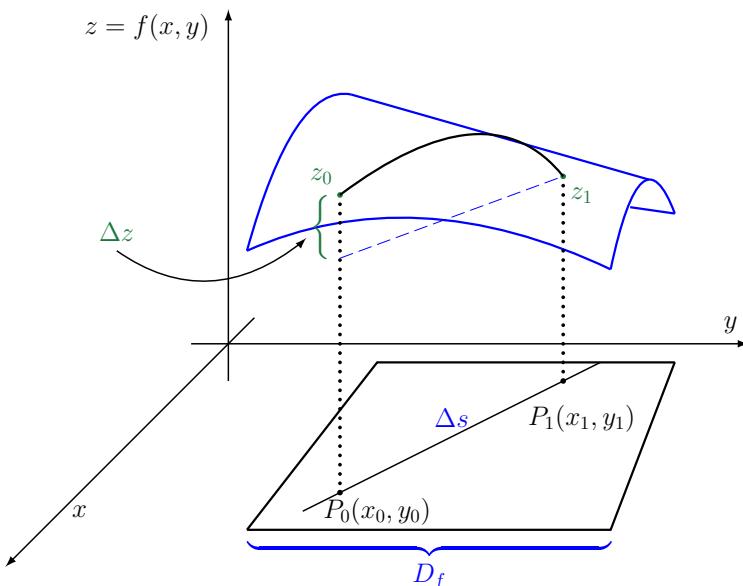


Figure 2.3 A 3D graph of a function of two variables.

In direct analogy with the 1D derivative we have $\Delta z = f(x_0, y_0) - f(x_1, y_1)$, and $\Delta s = \sqrt{\Delta x^2 + \Delta y^2}$.

Referring to Figure 2.3, for the limit process to make sense several related questions need to be addressed:

- * does $f(x_0, y_0)$ exist, or equivalently, does (x_0, y_0) belong to D_f ?
- * is f continuous at (x_0, y_0) ?
- * does the limit $\lim_{P_1 \rightarrow P_0}$ exist if $P_1 \in D_f$?

Over the next few pages we consider these questions with the aim of establishing a set of conditions for the existence of derivatives and possibly a set of guidelines that can be followed in applications.

2.B Limits and continuity

In Chapter 1 the pointwise limit of a function of one variable was explained in some detail, culminating in the ϵ - δ condition shown in Figure 1.16. The latter condition permits a direct generalization of the limit concept to functions of several variables. Keep in mind that the discussion below pertains to given points in \mathbb{R}^n , *i.e.* it too holds pointwise.

Definition 2.2

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ approaches a limit L as $\mathbf{x} \rightarrow \mathbf{a}$ for all points \mathbf{x} in a neighbourhood of \mathbf{a} belonging to D_f if:

Given any small positive number ϵ , another positive number δ , which may depend on ϵ , can be found such that if \mathbf{x} is within a radius δ of \mathbf{a} , then f will be within a radius ϵ of L .

In mathematical notation: Given any $\epsilon > 0$ there exists a $\delta > 0$ such that wherever $0 < |\mathbf{x} - \mathbf{a}| < \delta$ then $|f(\mathbf{x}) - L| < \epsilon$.

We write $\lim_{\mathbf{x} \rightarrow \mathbf{a}} f(\mathbf{x}) = L$.

If it can be established that two functions, f and g , satisfy the conditions of this fundamental definition at a given point, then a number of important results involving their combination follow. We refer to these as *laws* and apply them to all well-behaved functions.

Limit laws: If $\lim_{x \rightarrow a} f(x) = L$, $\lim_{x \rightarrow a} g(x) = M$, then the following sum, product, quotient, convergence and composition results can be proved.

$$(a) \lim_{x \rightarrow a} (f(x) + g(x)) = L + M$$

$$(b) \lim_{x \rightarrow a} (f(x) \cdot g(x)) = L \cdot M$$

$$(c) \lim_{x \rightarrow a} \frac{f(x)}{g(x)} = \frac{L}{M} \quad (M \neq 0)$$

(d) $\lim_{x \rightarrow a} f(x) = \lim_{x \rightarrow a} g(x)$ and $f(x) \leq h(x) \leq g(x)$ means that $\lim_{x \rightarrow a} h(x)$ exists and equals L which equals M (a “squeeze theorem”).

(e) If $F(t)$ is a continuous function at $t = L$ then

$$\lim_{x \rightarrow a} F(f(x)) = F(L) = F\left(\lim_{x \rightarrow a} f(x)\right).$$

That is, for continuous functions, we may interchange the limit and function composition operations.

■ Example 2.3:

Here is a proof of limit law (a) using the ϵ - δ concept in Definition 2.2.

We may assume that, given $\epsilon_1 > 0$ and $\epsilon_2 > 0$, we have found numbers $\delta_1 > 0$ and $\delta_2 > 0$ such that $|f(x) - L| < \epsilon_1$ whenever $|x - a| < \delta_1$, and $|g(x) - M| < \epsilon_2$ whenever $|x - a| < \delta_2$.

For given arbitrarily small $\epsilon > 0$, let $\epsilon_1 = \epsilon_2 = \frac{1}{2}\epsilon$. Then we have

$$\begin{aligned} |(f(x) + g(x)) - (L + M)| &= |(f(x) - L) + (g(x) - M)| \\ &\leq |f(x) - L| + |g(x) - M| \\ &\quad \text{by the triangle inequality,} \\ &< \epsilon_1 + \epsilon_2 = \epsilon \quad \text{provided both} \\ &\quad |x - a| < \delta_1 \text{ and } |x - a| < \delta_2. \end{aligned}$$

Now we may choose $\delta = \min(\delta_1, \delta_2)$, and we then have

$$|(f(x) + g(x)) - (L + M)| < \epsilon \text{ whenever } |x - a| < \delta.$$

Thus, we have proved that $\lim_{x \rightarrow a} (f(x) + g(x)) = L + M$. ■

☞ Mastery Check 2.1:

Prove the limit laws (b)–(e).

Hint: For law (b), with the assumptions in Example 2.3, assume $|\mathbf{x} - \mathbf{a}| < \delta = \min(\delta_1, \delta_2)$, and write $f(\mathbf{x}) = L + e_1(\mathbf{x})$, which implies $|e_1(\mathbf{x})| < \epsilon_1$, and similarly for g . Expand $f(\mathbf{x}) \cdot g(\mathbf{x})$ in terms of e_1 and e_2 . Let $\epsilon_1 = \frac{1}{3}\epsilon/|M|$ and $\epsilon_2 = \frac{1}{3}\epsilon/|L|$ and consider $|f(\mathbf{x}) \cdot g(\mathbf{x}) - L \cdot M| < \epsilon$.

For law (c), prove first that $\lim_{\mathbf{x} \rightarrow \mathbf{a}} \frac{1}{g(\mathbf{x})} = \frac{1}{M}$ and then invoke law (b).



■ Example 2.4:

Consider $L = \lim_{(x,y) \rightarrow (1,\pi)} \frac{\cos(xy)}{1 - x - \cos y} = \lim_{(x,y)} \frac{g(x,y)}{h(x,y)}$, noting in particular that $\lim_{(x,y) \rightarrow (1,\pi)} h(x,y) \neq 0$. Applying the standard rules we find that

$$L = \frac{\lim \cos(xy)}{\lim(1 - x - \cos y)} = \frac{-1}{+1} = -1.$$

Here, we have used the sum, product, quotient, and composition laws.



In evaluating limits of any well-behaved $f : \mathbb{R}^n \rightarrow \mathbb{R}$ for $n > 2$, we follow the exact same process as implied in the above example: besides using the limit laws, the reader can also make use of results from the study of limits of functions of one variable, some of which are listed on Page 24. However, the simple statement made in the limit definition hides considerable detail that we need to confront in more complicated cases. Definition 2.2 implicitly means that

- * $\lim_{\mathbf{x} \rightarrow \mathbf{a}} f(\mathbf{x})$ exists and is equal to L if $f \rightarrow L$ independently of how \mathbf{x} approaches \mathbf{a} !
- * The limit L , if it exists, is unique!
- * No limit of f exists if f has different limits when \mathbf{x} approaches \mathbf{a} along different curves!

The graphical depiction of Definition 2.2, in analogy with Figure 1.16, is shown in Figure 2.4 on the next page.

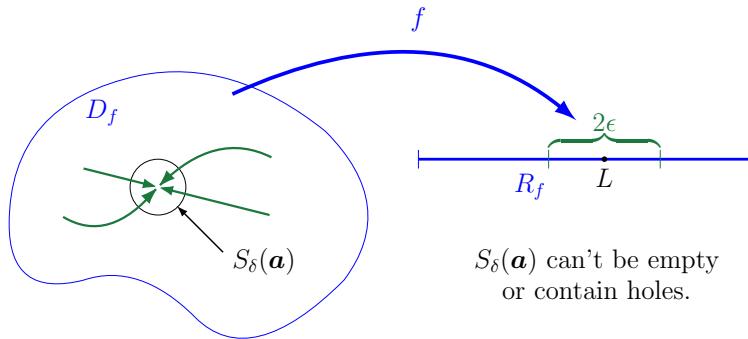


Figure 2.4 Schematic of the limit definition in 2D.

As in the 1D cases (Examples 2.1 and 2.2), the idea of considering limit values from multiple approaches is best illustrated by an example which fails to meet one or more criteria, such as the example below.

■ Example 2.5:

Consider $\lim_{(x,y) \rightarrow (1,1)} \frac{x-y}{x-1}$. (This is of the form $\frac{0}{0}$. The function is therefore undefined at $(1,1)$.) We attempt to evaluate the limit by approaching the point $(1,1)$ along four different paths as shown in Figure 2.5.

$$\begin{aligned}
 y = x : \quad & \lim_{(x,y) \rightarrow (1,1)} \frac{x-y}{x-1} = \lim_{(x,y) \rightarrow (1,1)} \frac{0}{x-1} = 0. \\
 y = 2 - x : \quad & \lim_{(x,y) \rightarrow (1,1)} \frac{x-y}{x-1} = \lim_{(x,y) \rightarrow (1,1)} \frac{2(x-1)}{x-1} = 2. \\
 y = x^2 : \quad & \lim_{(x,y) \rightarrow (1,1)} \frac{x-y}{x-1} = \lim_{(x,y) \rightarrow (1,1)} \frac{x-x^2}{x-1} = -1. \\
 y = 1 : \quad & \lim_{(x,y) \rightarrow (1,1)} \frac{x-y}{x-1} = \lim_{(x,y) \rightarrow (1,1)} \frac{x-1}{x-1} = 1.
 \end{aligned}$$

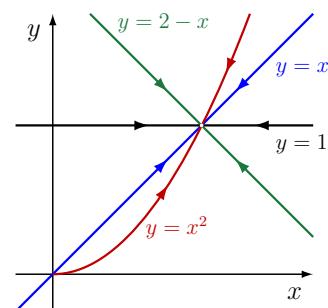


Figure 2.5 Different paths of approach to a limit point.

The resulting limiting values found by following these different paths are all different. From this we conclude that no limit exists. Note that it is enough for *one* of these cases to give a different result for us to conclude that no limit exists. The graph of the function is shown in Figure 2.6.

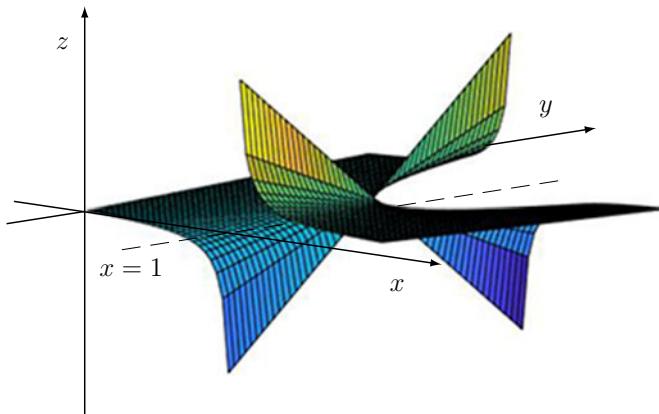


Figure 2.6 The graph of $z = \frac{x - y}{x - 1}$, $0 \leq x, y \leq 2$.

The point $(1, 1)$ is the cusp where the two sheets of the function meet (along the dashed line).

■

☞ Mastery Check 2.2:

Consider the function $f(x, y) = \frac{x^2 + 2y^2}{2x^2 + y^2}$. Does $\lim_{(x,y) \rightarrow (0,0)} f(x, y)$ exist? Hint:

Note that $f(x, y)$ is of the form $\frac{0}{0}$ at $(0, 0)$ making $f(x, y)$ undefined at $(0, 0)$. Take limits along the lines $x = 0$, $y = x$, and finally along the line $y = kx$, and see if your limits agree.

☞

From the perspective of effort expended, cases where a limit fails to exist are most often the least taxing, a few well-chosen approach paths will suffice. Now we need to ask, what about the cases where we get the same result for a few different trials? Do we need to try all the infinite number of approach directions to be convinced? One solution is proposed in the next example.

■ Example 2.6:

Consider $\lim_{(x,y) \rightarrow (0,0)} \frac{x^3 - x^2 y}{x^2 + y^2 + xy}$. Note once again that the function is undefined at the origin.

First we evaluate the limit along a few simple paths.

$$\text{Along } y = 0 : \lim_{(0,0)} \frac{x^3 - x^2 y}{x^2 + y^2 + xy} = \lim_{(0,0)} \frac{x^3}{x^2} = \lim_{x \rightarrow 0} x = 0.$$

$$\text{Along } x = 0 : \lim_{(0,0)} \frac{x^3 - x^2 y}{x^2 + y^2 + xy} = \lim_{y \rightarrow 0} \frac{0}{y^2} = 0.$$

We get the same result along any straight line $y = kx$. If the limit exists, it must be 0!

So, consider an arbitrary curve $r = f(\theta) > 0$, where $x = r \cos \theta$, $y = r \sin \theta$, and let $r \rightarrow 0$. (Shown in Figure 2.7 is one of the cases $r \rightarrow 0$ as θ increases, but any path in the plane will do.)

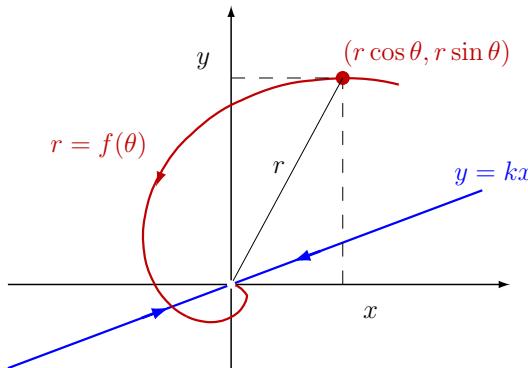


Figure 2.7 Example of conventional and unconventional paths to a limit point.

Substitute the polar functions for x and y in the definition of the function limit.

$$\begin{aligned} \left| \frac{x^3 - x^2 y}{x^2 + y^2 + xy} - 0 \right| &= \left| \frac{r^3 \cos^3 \theta - r^3 \cos^2 \theta \sin \theta}{r^2 \cos^2 \theta + r^2 \cos^2 \theta + r^2 \cos \theta \sin \theta} \right| \\ &= \left| \frac{r^3 \cos^2 \theta (\cos \theta - \sin \theta)}{r^2 (1 + \cos \theta \sin \theta)} \right| \\ &= r \cos^2 \theta \frac{|\cos \theta - \sin \theta|}{|1 + \sin \theta \cos \theta|} \end{aligned}$$

(It is actually sufficient to stop here: the denominator is not zero and the numerator is bounded and proportional to r which $\rightarrow 0$.)

$$\begin{aligned} \left| \frac{x^3 - x^2 y}{x^2 + y^2 + xy} - 0 \right| &= r \cos^2 \theta \frac{\sqrt{2} |\cos \theta \cos(\pi/4) - \sin \theta \sin(\pi/4)|}{|1 + \frac{1}{2} 2 \sin \theta \cos \theta|} \\ &\leq r \sqrt{2} \frac{|\cos(\theta + \pi/4)|}{|1 + \frac{1}{2} \sin(2\theta)|} \\ &\leq r \frac{\sqrt{2}}{1/2} = 2\sqrt{2}r \rightarrow 0 \text{ as } r \rightarrow 0. \end{aligned}$$

In these steps we have used only known properties of the trigonometric functions for arbitrary angles.

Thus, given $\epsilon > 0$, however, small, we can find a δ , a function of ϵ (choose $\delta = \frac{\epsilon}{2\sqrt{2}}$), such that

$$\left| \frac{x^3 - x^2 y}{x^2 + y^2 + xy} - 0 \right| < \epsilon \quad \text{whenever} \quad r < \delta.$$

Given that we have invoked an arbitrary curve whose sole requirement is to pass through the limit point (the origin) the result is general, the limit exists, and is indeed 0. The surface itself is reproduced in Figure 2.8.

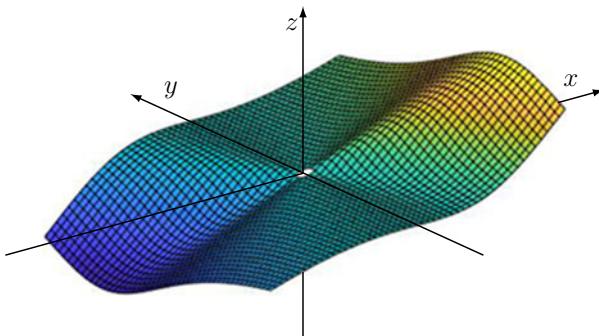


Figure 2.8 The graph of $z = \frac{x^3 - x^2 y}{x^2 + y^2 + xy}$.



 **Mastery Check 2.3:**

Consider the function $f(x, y) = \frac{x+y}{\ln(x^2+y^2)}$ (which is undefined at $(0, 0)$).

Does $\lim_{(x,y) \rightarrow (0,0)} f(x, y)$ exist?

Hint: Use the approach featured in Example 2.6, and standard limits.



Not all limits behave nicely!

Up to this point, the student may be inclined to think that taking limits along lines $y = kx$, $0 \leq k \leq \infty$, is sufficient to determine whether a limit at $(x, y) = (0, 0)$ exists or not. That is, that limit operations in \mathbb{R}^2 are a straightforward (no pun intended) extensions of the essential criteria listed on Page 50 for limits in \mathbb{R} .

But this is not so! The following case is a counterexample to Example 2.6:

 **Mastery Check 2.4:**

Consider the function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ defined by

$$f(x, y) = \frac{x^4 y^2}{(x^4 + y^2)^2}, \quad (x, y) \neq (0, 0).$$

Your task is to show that $\lim_{(x,y) \rightarrow (0,0)} f(x, y)$ does not exist (without first drawing the graph).

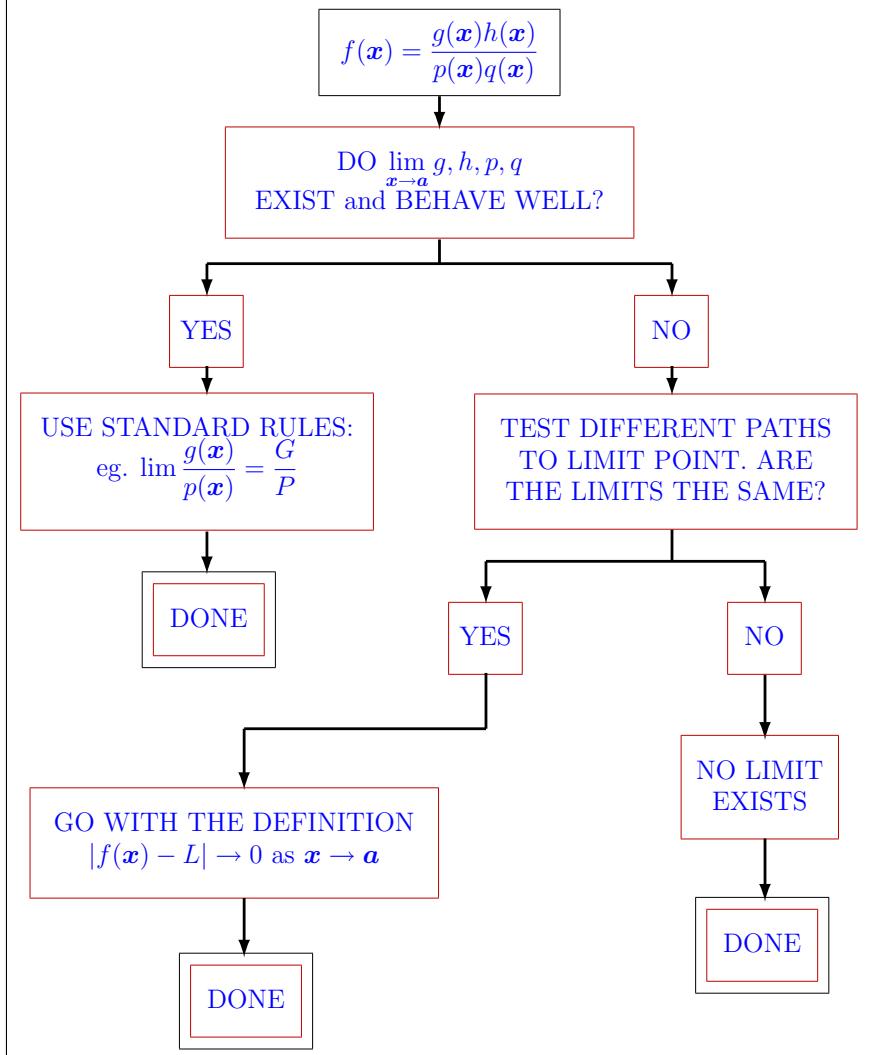
Hint: Consider another class of curves through the origin.



The student reader might get a better idea of all that is involved in evaluating limits with the following summary flowchart.

Flowchart 2.1: How to work through a limit problem

Let $\mathbf{x} = (x_1, x_2, \dots, x_n)$ be a variable in \mathbb{R}^n . Suppose we have



2.C Partial derivatives

Once the intricacies of limit processes for functions of many variables are understood, the application of these same principles to *ratios of differences* is relatively straightforward. In fact, the concept of partial derivatives becomes a simple extension of the 1D ordinary derivative.

Definition 2.3

Let \mathbf{x}_0 be an **interior point** or a **boundary point** of D_f of a continuous function $f : \mathbb{R}^n \rightarrow \mathbb{R}$.

- * If the limit $\lim_{h \rightarrow 0} \frac{f(\mathbf{x}_0 + h\mathbf{e}_j) - f(\mathbf{x}_0)}{h}$
 $= \lim_{h \rightarrow 0} \frac{f(x_{0,1}, \dots, x_{0,j} + h, \dots, x_{0,n}) - f(x_{0,1}, \dots, x_{0,n})}{h}$
exists, we call it the **first partial derivative** of f w.r.t. x_j at the point \mathbf{x}_0 and denote it $\frac{\partial f}{\partial x_j}(\mathbf{x}_0)$.
- * If all $\frac{\partial f}{\partial x_j}(\mathbf{x}_0)$, $j = 1, \dots, n$, exist then we say that f is **partially differentiable** at \mathbf{x}_0 .
- * We say f is **partially differentiable in D_f** if it is partially differentiable at every point $\mathbf{x}_0 \in D_f$.

Equivalent notations for partial derivatives are:

$$\frac{\partial f}{\partial x_j}(\mathbf{x}_0), \quad D_j f(\mathbf{x}_0), \quad f_{x_j}(\mathbf{x}_0), \quad f_j(\mathbf{x}_0).$$

Analogous to the 1D case, functions for which any of the n partial derivative limits fails to exist at a point are said to be singular at that point. That is, a multivariable function may be continuous everywhere in its domain of definition, but need not be differentiable at every point in its domain.

☞ Mastery Check 2.5:

Let $f(x, y, z) = xy + z \sin(yz)$. Using Definition 2.3, determine $\frac{\partial f}{\partial y}$ at an arbitrary point (x, y, z) .

Hint: You may need to use standard limits (see Page 24).

In solving this Mastery Check problem you will have noticed that you could

have and would have arrived at the same result had you used the rules of differentiation for functions of one variable, provided you treated x and z as if they were constants! In actual fact, Definition 2.3 effectively states that in taking the limit with respect to one variable, we do keep *all other* variables fixed. It should not come as a surprise that we find this equivalence. We demonstrate this very convenient operational equivalence with an example and leave it to Mastery Check 2.6 to reinforce the procedure.

■ Example 2.7:

Let $f(x, y, z) = \ln(1 + e^{xyz}) = g(h(x, y, z))$. We wish to calculate $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$.

In each case we assume two variables are constant and differentiate w.r.t. the third using the chain rule of single-variable calculus:

$$\begin{aligned}\frac{\partial f}{\partial x} &= \frac{dg}{dh} \frac{\partial h}{\partial x} = \frac{1}{1 + e^{xyz}} yze^{xyz}, \\ \frac{\partial f}{\partial y} &= \frac{1}{1 + e^{xyz}} xze^{xyz}, \\ \frac{\partial f}{\partial z} &= \frac{1}{1 + e^{xyz}} xye^{xyz}.\end{aligned}$$

■

☛ Mastery Check 2.6:

Find the (first-order) partial derivatives of the following functions with respect to each variable:

1. $f(x, y, z) = \frac{x^2 + y^2}{x^2 - z^2};$
2. $f(x, y, u, v) = x^2 \sin(2y) \ln(2u + 3v);$
3. $f(s, t, u) = \sqrt{s^2t + stu + tu^2};$
4. $f(x, y, z) = y \sin^{-1}(x^2 - z^2);$
5. $f(x, y, z, u) = \sin(3x) \cosh(2y) - \cos(3z) \sinh(2u);$
6. $f(u, v, w) = u^2 e^{uv^2 w}.$



Now that we can evaluate them, what are partial derivatives?

Let's look more closely at Figure 2.3 (Page 52). Given the foregoing discussion and particularly Definition 2.3 we consider two specific cases of that graph of the function of two variables.

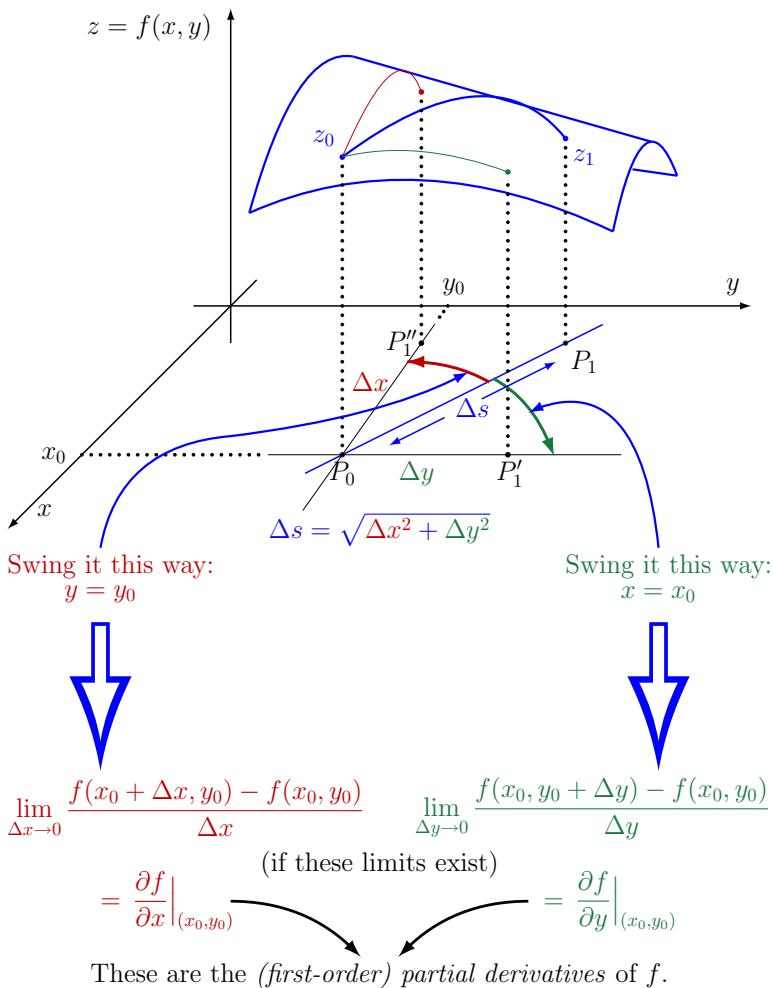


Figure 2.9 Partial derivatives of a function of two variables.

From Figure 2.9 we deduce geometric interpretations for $\frac{\partial f}{\partial x}|_0$, $\frac{\partial f}{\partial y}|_0$:

$\frac{\partial f}{\partial x}|_0$ — the *slope* of the *tangent line* L_1 to the curve $f(x, y_0)$ at (x_0, y_0) .

$\frac{\partial f}{\partial y}|_0$ — the *slope* of the *tangent line* L_2 to the curve $f(x_0, y)$ at (x_0, y_0) .

In fact, by taking two orthogonal cross sections through the point $(x_0, y_0, f(x_0, y_0))$ on the graph of f , one parallel to the x -axis and one parallel to the y -axis, the following facts can be obtained.

In the cross section parallel to the xz -plane, a vector parallel to line L_1 is \mathbf{v}_1 shown in Figure 2.10.

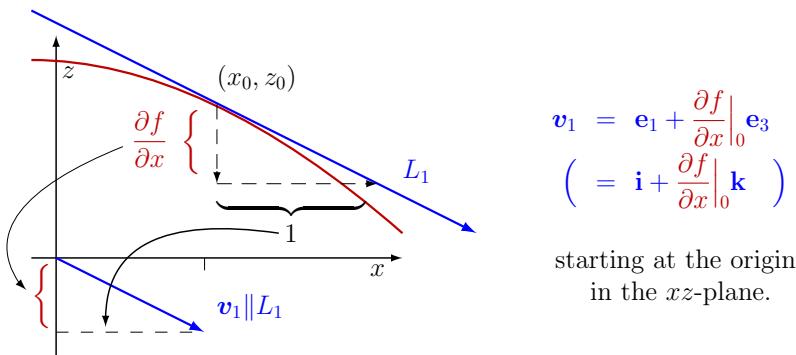


Figure 2.10 A tangent vector and line in the x -direction.

In the cross section parallel to the yz -plane, a vector parallel to line L_2 is \mathbf{v}_2 shown in Figure 2.11.

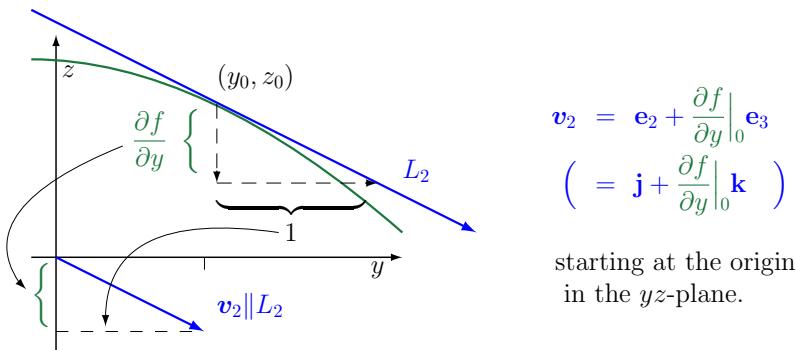


Figure 2.11 A tangent vector and line in the y -direction.

In the full 3D graph the tangent lines L_1 and L_2 , and corresponding vectors \mathbf{v}_1 and \mathbf{v}_2 , appear as in the figure below. (Compare the corresponding lines in Figure 2.12 below with those in Figures 2.9, 2.10, and 2.11 above.)

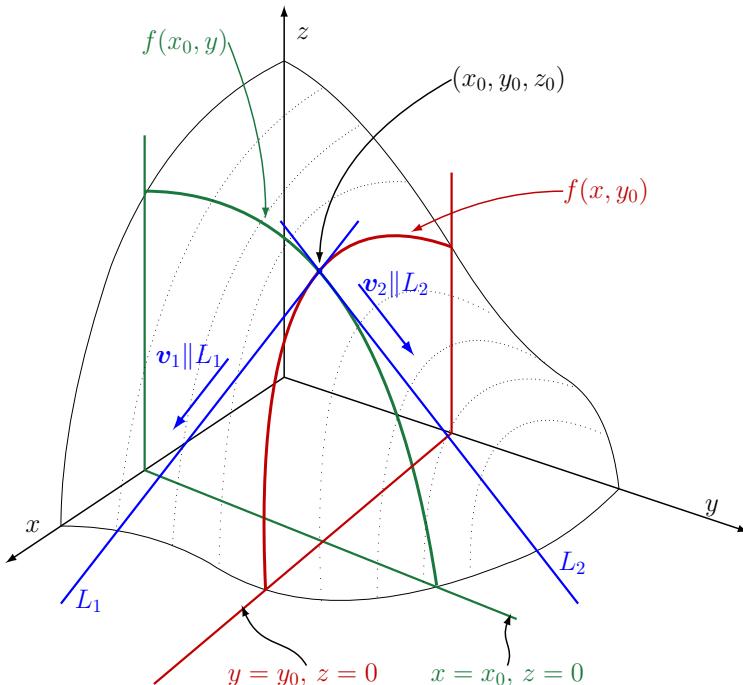


Figure 2.12 Tangent vectors and lines to a 3D function graph.

Notice that \mathbf{v}_1 and \mathbf{v}_2 are not parallel to each other! This is useful as the vectors \mathbf{v}_1 and \mathbf{v}_2 define a *tangent plane*, T , with normal vector \mathbf{n} :

$$\mathbf{n} = \mathbf{v}_1 \times \mathbf{v}_2 = \begin{vmatrix} \mathbf{e}_1 & \mathbf{e}_2 & \mathbf{e}_3 \\ 1 & 0 & \frac{\partial f}{\partial x} \Big|_0 \\ 0 & 1 & \frac{\partial f}{\partial y} \Big|_0 \end{vmatrix} = -\frac{\partial f}{\partial x} \Big|_0 \mathbf{e}_1 - \frac{\partial f}{\partial y} \Big|_0 \mathbf{e}_2 + \mathbf{e}_3.$$

The plane defined by L_1 and L_2 is *tangent* to the surface $z = f$ at (x_0, y_0, z_0) , is spanned by \mathbf{v}_1 and \mathbf{v}_2 , and, of course, has the same normal as the normal to the graph of $z = f(x, y, z)$ at (x_0, y_0, z_0) .

The equation of this tangent plane can be found from the scalar vector product (Page 4):

$$\mathbf{n} \cdot (\mathbf{x} - \mathbf{x}_0) = 0.$$

an arbitrary point on the plane. the specific point on the plane that is also on the surface, f .

$$\begin{aligned} \Rightarrow \quad & \frac{\partial f}{\partial x} \Big|_0 (x - x_0) + \frac{\partial f}{\partial y} \Big|_0 (y - y_0) - (z - z_0) = 0 \\ \Rightarrow \quad & z - z_0 = \frac{\partial f}{\partial x} \Big|_0 (x - x_0) + \frac{\partial f}{\partial y} \Big|_0 (y - y_0). \end{aligned}$$

✳️ Mastery Check 2.7:

Let $z = f(x, y) = \arcsin(xy)$.

Find the normal vector to the surface generated by $f(x, y)$, and the equation of the tangent plane, at $(1, \frac{1}{2}, \frac{\pi}{6})$.



2.D Differentiability of $f : \mathbb{R}^n \rightarrow \mathbb{R}$

We can now use the developments of the last section to establish a convenient definition of differentiability, extending the following geometric argument from single-variable calculus.

Consider the function $f : \mathbb{R} \rightarrow \mathbb{R}$. In saying that f is differentiable at a point $x_0 \in D_f$ we mean, *geometrically*, on the one hand, that there exists a tangent line to $f(x)$ at the point x_0 (Figure 2.13):

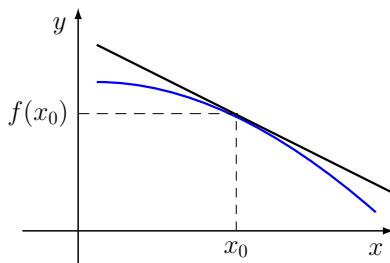


Figure 2.13 Tangent line to the graph of $f(x)$.

and *analytically* on the other hand:

$$\begin{aligned}
 \lim_{x \rightarrow x_0} \frac{f(x) - f(x_0)}{x - x_0} &= c \quad \left(= \frac{df}{dx} \Big|_0 \right) \\
 \iff \left| \frac{f(x) - f(x_0) - c(x - x_0)}{x - x_0} \right| &\rightarrow 0 \quad \text{as } x \rightarrow x_0 \\
 \iff \frac{f(x) - f(x_0) - c(x - x_0)}{x - x_0} &= \rho(x - x_0) \rightarrow 0 \quad \text{as } x \rightarrow x_0.
 \end{aligned}$$

The variable ρ is a function of $(x - x_0)$ which $\rightarrow 0$ as $x \rightarrow x_0$.

Rewriting this last result, we conclude that $f(x)$ is differentiable at x_0

$$\begin{aligned}
 \iff f(x) &= f(x_0) + c(x - x_0) + |\Delta x| \rho(\Delta x) \text{ and } \lim_{\Delta x \rightarrow 0} \rho(\Delta x) = 0 \\
 \iff f &\text{ can be approximated by a line.}
 \end{aligned}$$

The generalization of this argument to a function of two variables $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ is reasonably straightforward with the tangent line approximation being replaced by a tangent plane approximation.

We say $f(x, y)$ is differentiable at (x_0, y_0)

$$\begin{aligned}
 \iff f(x, y) &= f(x_0, y_0) + c_1 \Delta x + c_2 \Delta y + |\Delta \mathbf{x}| \rho(\Delta \mathbf{x}) \\
 \iff f &\text{ can be approximated by a plane at } (x_0, y_0).
 \end{aligned}$$

We formalize this reasoning in an even more general definition for a function of n variables.

Definition 2.4

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$. f is said to be **differentiable** at $\mathbf{x}_0 \in D_f$ if there exists a **linear approximation**, $\phi(\mathbf{x} - \mathbf{x}_0)$, such that

$$f(\mathbf{x}) = f(\mathbf{x}_0) + \phi(\mathbf{x} - \mathbf{x}_0) + |\Delta \mathbf{x}| \rho(\Delta \mathbf{x}; \mathbf{x}, \mathbf{x}_0) \quad \text{near } \mathbf{x}_0$$

with $\phi(\mathbf{x} - \mathbf{x}_0) = c_1(x_1 - x_{0,1}) + \dots + c_n(x_n - x_{0,n})$ and $\lim_{\Delta \mathbf{x} \rightarrow 0} \rho(\Delta \mathbf{x}; \mathbf{x}, \mathbf{x}_0) = 0$.
A function for which no linear approximation can be defined at a point is said to be **singular** at that point.

Although it is not critical to the discussion here (instead see Section 2.I and Chapter 3), a word or two about the function ρ (relevant to both the 1D

and nD cases) is warranted.

The condition that the function f is differentiable at \mathbf{x}_0 is equivalent to the condition of the existence of a tangent plane, ϕ , at that point. For $\mathbf{x} \neq \mathbf{x}_0$, a rearrangement of the equation in Definition 2.4 then defines the function $\rho(\Delta\mathbf{x}; \mathbf{x}, \mathbf{x}_0)$ as the ratio of the difference (effectively) of f and ϕ to $|\Delta\mathbf{x}|$,

$$\rho(\Delta\mathbf{x}; \mathbf{x}, \mathbf{x}_0) = \frac{f(\mathbf{x}) - f(\mathbf{x}_0) - \phi(\mathbf{x} - \mathbf{x}_0)}{|\Delta\mathbf{x}|},$$

which should only be a nonlinear contribution. Definition 2.4 then states that as a further condition for differentiability, this function must vanish in the limit $|\Delta\mathbf{x}| \rightarrow 0$. Essentially, for differentiability $f(\mathbf{x}) - f(\mathbf{x}_0)$ must behave as a linear function of the independent variables. We clarify this explanation with an example.

■ Example 2.8:

Consider the function $f(x, y) = xy^2$ and the point $\mathbf{x}_0 = (-2, 1)$. We introduce $\mathbf{h} = (h, k)$ so that $\mathbf{x} = \mathbf{x}_0 + \mathbf{h} = (-2 + h, 1 + k)$. Then

$$\begin{aligned} f(\mathbf{x}_0 + \mathbf{h}) - f(\mathbf{x}_0) &= (-2 + h)(1 + k)^2 - (-2)1^2 \\ &= (-2 + h)(1 + 2k + k^2) + 2 \\ &= (h - 4k) + (-2k^2 + 2hk + hk^2). \end{aligned}$$

In the last expression on the right-hand side, the two pairs of parentheses separate the linear approximation, $\phi(h, k) = h - 4k$, from the remaining nonlinear terms. From the latter terms we then form our ρ function

$$\rho(\mathbf{h}; \mathbf{x}_0 + \mathbf{h}, \mathbf{x}_0) = \frac{-2k^2 + 2hk + hk^2}{\sqrt{h^2 + k^2}}$$

which vanishes in the limit $|\mathbf{h}| = \sqrt{h^2 + k^2} \rightarrow 0$ since the numerator is at least quadratic in h and k while the denominator is of linear order. From the linear approximation, ϕ , we read off that $\left(\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}\right) \Big|_{(-2,1)} = (1, -4)$. ■

More generally, the linear function $\phi(\mathbf{x} - \mathbf{x}_0)$ in Definition 2.4 is

$$\phi(\mathbf{x} - \mathbf{x}_0) = \frac{\partial f}{\partial x_1} \Big|_0 (x_1 - x_{1,0}) + \cdots + \frac{\partial f}{\partial x_n} \Big|_0 (x_n - x_{n,0}).$$

That is, the coefficients c_1, \dots, c_n of the respective linear factors are simply the n partial derivatives of $f(\mathbf{x})$ evaluated at the point \mathbf{x}_0 , which means that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is differentiable at \mathbf{x}_0 if the limit vanishes:

$$\lim_{\Delta \mathbf{x} \rightarrow 0} \frac{f(x_{0,1} + \Delta x_1, \dots, x_{0,n} + \Delta x_n) - f(x_{0,1}, \dots, x_{0,n}) - \sum_{i=1}^n \frac{\partial f}{\partial x_i} \Big|_{\mathbf{x}_0} \Delta x_i}{|\Delta \mathbf{x}|} = 0.$$

Theorem 2.1

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ with continuous partial derivatives $\frac{\partial f}{\partial x_i}$ in the neighbourhood of a point \mathbf{x}_0 is differentiable at \mathbf{x}_0 .

This theorem on the continuity of partial derivatives as a condition for differentiability inspires the pictorial interpretation in Figure 2.14.

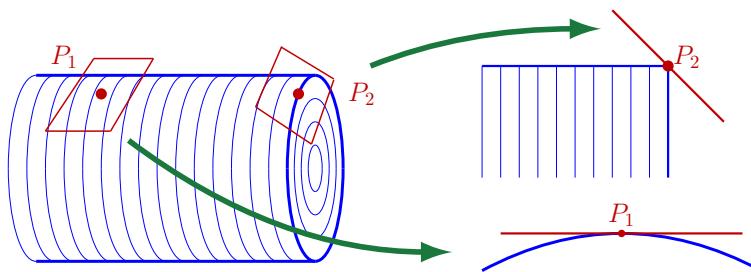


Figure 2.14 The relation between partial differentiability and differentiability.

At P_1 the surface is partially differentiable and differentiable, but at P_2 it has limited partial differentiability and so is not differentiable there.

Now for two important theorems (for proofs see [1] or similar texts):

Theorem 2.2

A differentiable function is continuous.

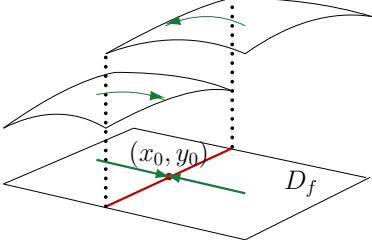
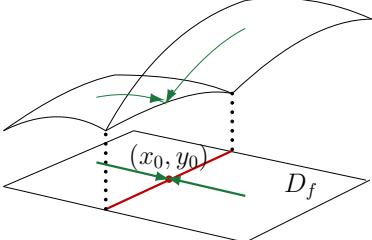
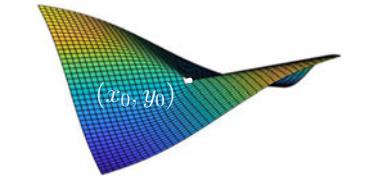
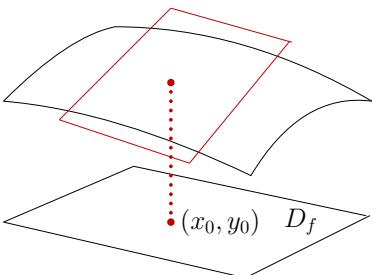
A continuous function is not necessarily a differentiable function.

Theorem 2.3

A differentiable function is partially differentiable.

A partially differentiable function is not necessarily a differentiable function.

Table 2.1: A pictorial table of differentiable functions (not exhaustive)

Function	Diff'ble at (x_0, y_0) ?	Why or why not?
	No	Limit does not exist at points on the red line. Function is not continuous. No tangent plane.
	No	Limit exists and function is continuous across the red line. But not <i>all</i> partial derivatives exist. No tangent plane.
	No	Function is continuous and all partial derivatives exist, but they are <i>not continuous</i> at one point. No tangent plane! (See Mastery Check 2.8.)
	Yes	Function is continuous and <i>all</i> partial derivatives exist and are continuous at (x_0, y_0) . There exists a unique tangent plane!

So, it appears that for $f : \mathbb{R} \rightarrow \mathbb{R}$, the function is differentiable at a point if the derivative exists, but for $f : \mathbb{R}^2 \rightarrow \mathbb{R}$, the partial derivatives have to exist *and be continuous* in an open circle about the point. The following Example and Mastery Check make this clear.

■ Example 2.9:

Consider the function $f : \mathbb{R}^2 \mapsto \mathbb{R}$, $f(x, y) = x \arctan(y/x)$, $f(0, y) = 0$. We wish to discuss the continuity of f , and the existence and continuity of f_x and f_y at points on the x -axis.

We have $\lim_{x \rightarrow 0} f(x, y) = \lim_{x \rightarrow 0} x \arctan(y/x) = 0 = f(0, y)$ for all y , since $|\arctan(y/x)| < \pi/2$.

The function is continuous for all points on the x -axis.

We have, for $y \neq 0$, $f_x = \arctan(y/x) - \frac{xy}{x^2 + y^2}$.

Then for $y > 0$, $\lim_{x \rightarrow 0^+} f_x = \frac{\pi}{2} - 0 = \frac{\pi}{2}$, $\lim_{x \rightarrow 0^-} f_x = -\frac{\pi}{2} + 0 = -\frac{\pi}{2}$;

and for $y < 0$, $\lim_{x \rightarrow 0^+} f_x = -\frac{\pi}{2} + 0 = -\frac{\pi}{2}$, $\lim_{x \rightarrow 0^-} f_x = \frac{\pi}{2} - 0 = \frac{\pi}{2}$.

Thus f_x is not continuous along $x = 0$ for $y \neq 0$. Also,

$$f_y = \frac{x^2}{x^2 + y^2}, \quad \lim_{x \rightarrow 0^+} f_y = \lim_{x \rightarrow 0^-} f_y = 0.$$

If we define $f_y(0, 0) = 0$, then f_y exists and is continuous on $x = 0$.



☛ Mastery Check 2.8:

Consider the function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ defined by

$$z = f(x, y) = \frac{xy}{\sqrt{x^2 + y^2}} \text{ if } x^2 + y^2 \neq 0, \quad f(0, 0) = 0,$$

whose graph appears at the end of this exercise in Figure 2.15. We wish to investigate the behaviour of f near $(0, 0)$.

- (1) Find the two partial derivatives for $(x, y) \neq (0, 0)$.
- (2) Show using Definition 2.3 that both partial derivatives are zero at $(0, 0)$.
Thus, the tangent plane at $(0, 0)$, if it exists, *must coincide with the plane* $z = 0$.

- (3) But the partial derivatives are not continuous. Show that these limits are not all the same:

$$\lim_{x \rightarrow 0} \left(\lim_{y \rightarrow 0} \frac{\partial f}{\partial x} \right), \quad \lim_{y \rightarrow 0^-} \left(\lim_{x \rightarrow 0} \frac{\partial f}{\partial x} \right), \quad \lim_{y \rightarrow 0^+} \left(\lim_{x \rightarrow 0} \frac{\partial f}{\partial x} \right).$$

(A similar result holds for the other derivative.)

So, we do not expect the tangent plane to exist.

- (4) Now recall the properties of a tangent plane as outlined in Definition 2.4. See if you can construct the expression

$\Delta z = f(\mathbf{x}) - \phi(\mathbf{x} - \mathbf{x}_0) = f(\mathbf{x}_0) + |\Delta \mathbf{x}| \rho(\Delta \mathbf{x})$ for the special case that \mathbf{x} lies on the line $y = x$, at distance

$$|\Delta \mathbf{x}| = \sqrt{\Delta x^2 + \Delta y^2} \text{ from } \mathbf{x}_0 = (0, 0), \text{ with } \Delta y = \Delta x.$$

That is, find $\rho(\Delta \mathbf{x})$.

Use this result to decide whether a tangent plane exists.

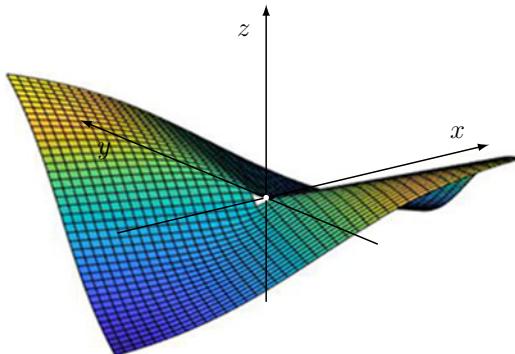


Figure 2.15 The graph of $z = \frac{xy}{\sqrt{x^2 + y^2}}$, $x, y \neq 0$.



2.E Directional derivatives and the gradient

The directional derivative

Thus far we have established that the partial derivatives of a function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ have the properties that:

$\frac{\partial f}{\partial x} \Big|_{\text{some point}} = \text{the rate of change (that is, the slope) of } f \text{ in the (positive) } x\text{-direction at "some point".}$

$\frac{\partial f}{\partial y} \Big|_{\text{some point}} = \text{the rate of change (that is, the slope) of } f \text{ in the (positive) } y\text{-direction at "some point".}$

These interpretations now beg the question: What if we wanted to find the rate of change of f in some other direction, such as \mathbf{u} depicted in Figure 2.16?

Suppose f is given and we know it is differentiable at a point (x_0, y_0) and we wanted the rate of change of f in that particular direction \mathbf{u} . We may now combine all the ingredients that go into the limit definition of the derivative in Equation (2.1) on Page 52 and suppose, in addition, that $\mathbf{u} = (u, v)$ is a given *unit* vector in the xy -plane.

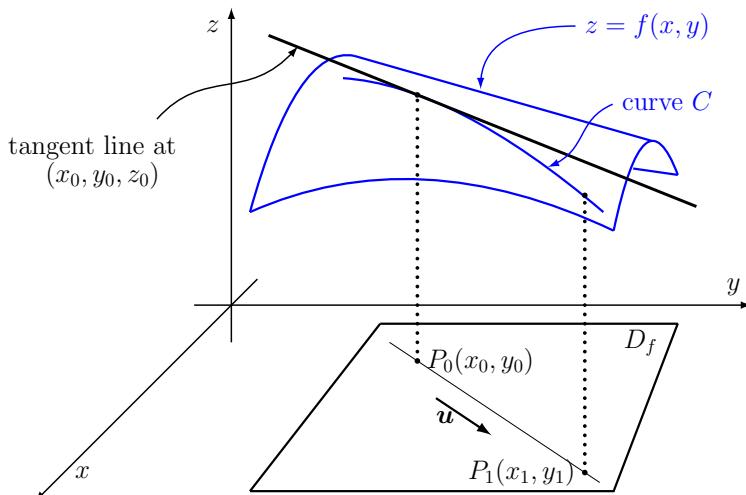


Figure 2.16 The tangent vector in an arbitrary direction.

We notice first that the two points P_0 and P_1 appearing in Equation (2.1) define a vector in the xy -plane, parallel to our given \mathbf{u} . Hence,

$$\Delta \mathbf{x} = \overrightarrow{P_0 P_1} = (x_1 - x_0, y_1 - y_0) = t \cdot \mathbf{u} = t \cdot (u, v).$$


Now we may re-consider the general expression for the derivative of f in Equation (2.1) which we rewrite for this special case.

Definition 2.5

The derivative limit:

$$\begin{aligned} \lim_{\Delta x \rightarrow 0} \frac{f(x_1, y_1) - f(x_0, y_0)}{|\Delta x|} &= \lim_{\Delta x \rightarrow 0} \frac{f(x_0 + tu, y_0 + tv) - f(x_0, y_0)}{|t\mathbf{u}|} \\ &= \lim_{t \rightarrow 0} \frac{f(x_0 + tu, y_0 + tv) - f(x_0, y_0)}{|t|}, \quad (|\mathbf{u}|=1), \end{aligned}$$

if it exists, is called the **directional derivative** of f in the direction of \mathbf{u} at (x_0, y_0) .

Commonly used notations for the directional derivative include

$D_u f(x_0)$ and $\frac{df}{dy}(x_0)$.

To calculate the directional derivative, there are two alternatives: Either we use the above definition (which may be necessary if the function is *not continuous or not differentiable*), or defer to the following theorem.

Theorem 2.4

If $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ is a differentiable function, then

$$D_{\mathbf{u}} f(\mathbf{x}_0) = \frac{df}{d\mathbf{u}}(\mathbf{x}_0) = \frac{\partial f}{\partial x}\Big|_0 u + \frac{\partial f}{\partial y}\Big|_0 v.$$

For the conditions stated, this theorem is easy to prove.

■ **Proof:**

In light of Definition 2.4, f differentiable means

$$\begin{aligned}\frac{f(\mathbf{x}_0 + t\mathbf{u}) - f(\mathbf{x}_0)}{t} &= \frac{1}{t} \left(f(\mathbf{x}_0) + \sum_{i=1}^n \frac{\partial f}{\partial x_i} \Big|_0 t u_i + t \rho(t) - f(\mathbf{x}_0) \right) \\ &= \sum_{i=1}^n \frac{\partial f}{\partial x_i} \Big|_0 u_i + \rho(t) \rightarrow \sum_{i=1}^n \frac{\partial f}{\partial x_i} \Big|_0 u_i \quad \text{as } t \rightarrow 0 \text{ since } \rho(t) \rightarrow 0.\end{aligned}$$

■

The above simple proof suggests that we can easily extend the definition of a directional derivative and its convenient form to $f : \mathbb{R}^n \rightarrow \mathbb{R}$. For these functions we have

$$D_{\mathbf{u}} f(\mathbf{x}_0) = \frac{\partial f}{\partial x_1} \Big|_0 u_1 + \frac{\partial f}{\partial x_2} \Big|_0 u_2 + \cdots + \frac{\partial f}{\partial x_n} \Big|_0 u_n \quad (2.2)$$

where $\mathbf{u} = (u_1, u_2, \dots, u_n)$ and $|\mathbf{u}| = 1$.

For a differentiable function what the directional derivative gives us is the slope of the tangent plane in the direction \mathbf{u} !

Gradient of a scalar function

Let's look a little more closely at what we use to calculate the directional derivative of a differentiable function. In the general case of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ we interpret Equation (2.2) as the scalar product of two vectors:

$$\begin{aligned}D_{\mathbf{u}} f(\mathbf{x}_0) &= \sum_{i=1}^n \frac{\partial f}{\partial x_i} \Big|_0 u_i \\ &= \underbrace{\left(\frac{\partial f}{\partial x_1} \Big|_0, \frac{\partial f}{\partial x_2} \Big|_0, \dots, \frac{\partial f}{\partial x_n} \Big|_0 \right)}_{\text{a new vector called ...}} \cdot \underbrace{(u_1, u_2, \dots, u_n)}_{\text{direction of interest}}\end{aligned}$$

Definition 2.6

The gradient vector function of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at a point $\mathbf{x} \in D_f$ is defined as

$$\begin{aligned}\text{grad } f(\mathbf{x}) &= \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right) \equiv \nabla f(\mathbf{x}) \\ &\text{— we say “grad } f\text{” or “del } f\text{”.}\end{aligned}$$

Therefore, for a differentiable $f : \mathbb{R}^n \rightarrow \mathbb{R}$, the directional derivative of $f(\mathbf{x}_0)$ in the direction \mathbf{u} is the scalar product of the gradient of f evaluated at \mathbf{x}_0 and the *unit* vector \mathbf{u} .

$$D_{\mathbf{u}} f(\mathbf{x}_0) = \text{grad } f(\mathbf{x}_0) \cdot \mathbf{u} = \nabla f(\mathbf{x}_0) \cdot \mathbf{u}.$$

■ Example 2.10:

Consider the function $f(x, y, z) = xy^3 + yz^2$. What is the directional derivative at the point $(1, 2, -1)$ in the direction $\mathbf{u} = 2\mathbf{e}_1 + \mathbf{e}_2 + 2\mathbf{e}_3$?

The gradient of f is

$$\nabla f = y^3\mathbf{e}_1 + (3xy^2 + z^2)\mathbf{e}_2 + 2yz\mathbf{e}_3.$$

At $(1, 2, -1)$ this is $\nabla f = 8\mathbf{e}_1 + 13\mathbf{e}_2 - 4\mathbf{e}_3$.

The direction of $\mathbf{u} = 2\mathbf{e}_1 + \mathbf{e}_2 + 2\mathbf{e}_3$ is

$$\mathbf{n} = \frac{\mathbf{u}}{|\mathbf{u}|} = \frac{2}{3}\mathbf{e}_1 + \frac{1}{3}\mathbf{e}_2 + \frac{2}{3}\mathbf{e}_3.$$

The directional derivative we require is the scalar product of these:

$$\nabla f \cdot \mathbf{n} = (8\mathbf{e}_1 + 13\mathbf{e}_2 - 4\mathbf{e}_3) \cdot \left(\frac{2}{3}\mathbf{e}_1 + \frac{1}{3}\mathbf{e}_2 + \frac{2}{3}\mathbf{e}_3\right) = 7.$$



☞ Mastery Check 2.9:

What is the unit normal to the surface $xy^2z^3 = 4$ at the point $(1, 2, -1)$?



☞ Mastery Check 2.10:

Calculate the directional derivative of $f(x, y, z) = xy + e^{yz} + z$ in the direction $\mathbf{u} = (\alpha, \beta, \gamma)$, where $\alpha^2 + \beta^2 + \gamma^2 = 1$, at the point $(1, 1, 0)$.

When you have found your directional derivative, answer this question:

How should α, β, γ be chosen so that this derivative should be maximal?



Remarks — Some all-important facts about the gradient, ∇f :

- * $\nabla f(\mathbf{x})$ is the generalization to functions of several variables of $\frac{dg}{dx}$ for a function g of one variable, x .
- * In $1D$, $\frac{dg}{dx} = 0$ for all $x \in D_g \implies g(x) = \text{const.}$

For $f : \mathbb{R}^n \rightarrow \mathbb{R}$, if

$$\nabla f = \mathbf{0} \text{ for all } \mathbf{x} \text{ in an open set in } D_f,$$

then $f(\mathbf{x})$ is constant in that set.

$$\nabla f(x, y, z) = \mathbf{0} \quad \forall \mathbf{x} \in D_f \implies \left\{ \begin{array}{l} \frac{\partial f}{\partial x} = 0 \Rightarrow f = f(y, z) \\ \frac{\partial f}{\partial y} = 0 \Rightarrow f = f(x, z) \\ \frac{\partial f}{\partial z} = 0 \Rightarrow f = f(x, y) \end{array} \right\} \implies f(x, y, z) = \text{const.}$$

But, quite often $\nabla f = \mathbf{0}$ at some *isolated* point \mathbf{x}_0 : f is then *not* constant. See Section 3.A for further discussion.

- * At a point $\mathbf{x} \in D_f$, the differentiable function

$$f : \mathbb{R}^n \rightarrow \mathbb{R} \quad \left\{ \begin{array}{l} \text{increases} \\ \text{decreases} \end{array} \right\} \text{most rapidly in the direction of } \pm \nabla f(\mathbf{x}).$$

The maximum rate of change of f is given by $|\nabla f(\mathbf{x})|$.

In Section 3.E we will encounter the gradient again, while later in Chapter 5 the “del” operator appears in other guises.

- * If $\nabla f(\mathbf{x}) \neq \mathbf{0}$ for a differentiable $f : \mathbb{R}^n \rightarrow \mathbb{R}$ then $\nabla f(\mathbf{x})$ is a vector which is normal to a level set, that is, a $\left\{ \begin{array}{l} \text{level curve in } \mathbb{R}^2 \\ \text{level surface in } \mathbb{R}^3 \end{array} \right\}$ of f .

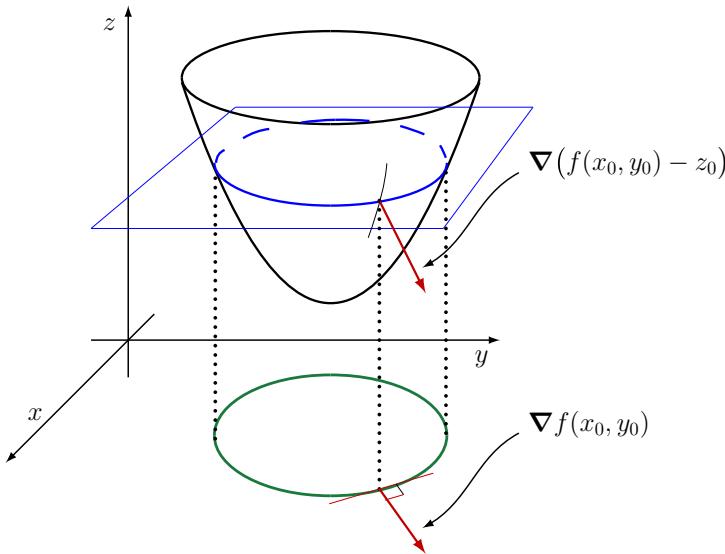


Figure 2.17 Comparison of the gradient applied in 2D and 3D circumstances.

In Figure 2.17, it can be seen that $\nabla(f(x_0, y_0) - z_0) = \frac{\partial f}{\partial x}\Big|_0 \mathbf{e}_1 + \frac{\partial f}{\partial y}\Big|_0 \mathbf{e}_2 - \mathbf{e}_3$ is normal to the 3D surface $f(x, y) - z = \text{const.}$ at (x_0, y_0, z_0) , while $\nabla f(x_0, y_0) = \frac{\partial f}{\partial x}\Big|_0 \mathbf{e}_1 + \frac{\partial f}{\partial y}\Big|_0 \mathbf{e}_2$ is normal to the 2D level curve $f(x, y) = \text{const.}$ at (x_0, y_0) .

☞ Mastery Check 2.11:

Find the equation of the tangent plane to the level surface of $w = f(x, y, z)$ when $w = 2$ at the point $(1, 1, \pi)$, where $f(x, y, z) = xy \cos(z) + 3x$.



☞ Mastery Check 2.12:

Find the equation of the tangent plane to the level surface of $w = f(x, y, z)$ at the point $(1, 1, \pi)$, where $f(x, y, z) = xy \cos(z) + 3x$.

Hint: We are now working in four dimensions. Consider the level set

$$g(x, y, z, w) = f(x, y, z) - w = 0.$$



2.F Higher-order derivatives

By now the reader will have correctly surmised that, just as in the single-variable case, higher-order derivatives are possible for functions of many variables.

Indeed, if

- (a) $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a continuous function of x_1, x_2, \dots, x_n , and
- (b) some given partial derivative $\frac{\partial f}{\partial x_i}$ exists and is itself a continuous (not necessarily differentiable) function of x_1, x_2, \dots, x_n ,

then $\frac{\partial f}{\partial x_i}$ can itself be considered a function of \mathbf{x} (independent of f).

For a convenient explanation, we shall refer to this particular derivative as $g(\mathbf{x})$ ($\equiv \frac{\partial f}{\partial x_i}(\mathbf{x})$). We can now consider the partial derivatives of g just as we had done with f :

Definition 2.7

If $\frac{\partial g}{\partial x_j} \equiv \lim_{h \rightarrow 0} \frac{g(x_1, \dots, x_j + h, \dots, x_n) - g(x_1, \dots, x_j, \dots, x_n)}{h}$ exists it is called a **second-order partial derivative of f** . More specifically, it is a **second-order mixed partial derivative of f w.r.t. x_i and x_j** .

In terms of the original f we have that

$$\frac{\partial g}{\partial x_j} = \frac{\partial}{\partial x_j} \left(\frac{\partial f}{\partial x_i} \right) \equiv \frac{\partial^2 f}{\partial x_j \partial x_i}$$


— this is the notation used
(mostly) in this book

differentiate first w.r.t. x_i then w.r.t. x_j

Note that other notations for the second derivative of $f(x, y)$ are in common use such as

$$\frac{\partial^2 f}{\partial y \partial x} = f_{xy} = f_{12},$$

each of which describes a second-order partial derivative. First, a partial derivative w.r.t. x , then a partial derivative w.r.t. y . The reader should exercise some care in interpreting the different notations.

We are now implored to explain what higher partial derivatives are. It suffices to consider a function of two variables, $f(x, y)$. If $\frac{\partial f}{\partial x}\bigg|_{(x_0, y_0)}$ is the slope of the tangent to f at (x_0, y_0) in the direction of x , then, just as in the single-variable case, $\frac{\partial^2 f}{\partial x^2}\bigg|_{(x_0, y_0)}$ is the rate of change of the slope in this same direction. It is therefore a measure of the curvature of f in this direction. On the other hand, $\frac{\partial^2 f}{\partial y \partial x}\bigg|_{(x_0, y_0)}$ is the rate of change of the x -directional slope in the y -direction.

A convenient and useful result for so-called *smooth* functions which, apart from their applications in applied contexts (Chapters 3 and 5), relieves some of the stress of interpreting notation, is the following.

Theorem 2.5

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuous and $\frac{\partial f}{\partial x_i}$, $i = 1, 2, \dots, n$ exist and are continuous in $S_r(\mathbf{x}) \subset D_f$ and that both $\frac{\partial^2 f}{\partial x_i \partial x_j}$ and $\frac{\partial^2 f}{\partial x_j \partial x_i}$ exist and are continuous at $\mathbf{x} \in D_f$. Then $\frac{\partial^2 f}{\partial x_i \partial x_j} = \frac{\partial^2 f}{\partial x_j \partial x_i}$ at $\mathbf{x} \in D_f$.

(For the standard proof, see a standard text book such as [1] or [2].)

Note the conditions of the above theorem highlighted in Figure 2.18.

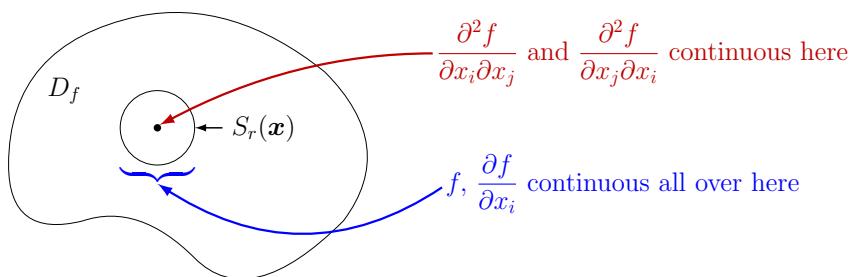


Figure 2.18 Conditions for the equivalence of mixed partial derivatives.

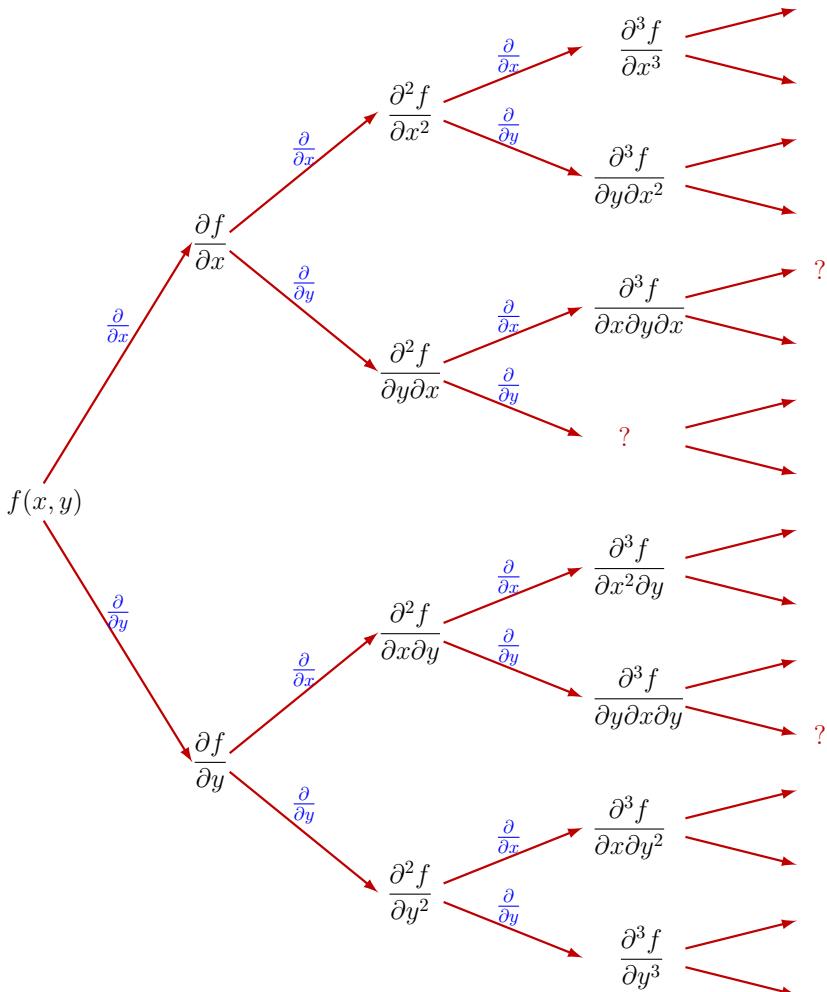


Figure 2.19 A chart of possible derivatives of $f: \mathbb{R}^2 \rightarrow \mathbb{R}$.

Figure 2.19 gives an indication of the scope of possibilities of higher-order partial derivatives for a function (*of two variables*) that is sufficiently differentiable. What constitutes “sufficient” will be defined shortly. In the meantime, referring to the arrow convention in Figure 2.19, what are the derivatives in the positions where question marks appear?

Definition 2.8

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ with **continuous derivatives up to and including order m** ($0, 1, 2, \dots, m$) in an open subset of D_f is said to be of class C^m .

A C^2 -function thus satisfies $\frac{\partial^2 f}{\partial x_i \partial x_j} = \frac{\partial^2 f}{\partial x_j \partial x_i}$.

☞ Mastery Check 2.13:

Verify that for $f(x, y) = 2 - x^2y - xy^2$, the mixed derivatives $\frac{\partial^2 f}{\partial x \partial y}$ and $\frac{\partial^2 f}{\partial y \partial x}$ are equal. Draw the graph of $f(x, y)$ using MATLAB®.

**☞ Mastery Check 2.14:**

Determine all partial derivatives of order 2 of $f(x, y) = \arctan\left(\frac{x}{y}\right)$.

Specify any points where f or its derivatives are not defined.

**☞ Mastery Check 2.15:**

Determine all C^2 functions $z = f(x, y)$ which satisfy the conditions

$$\frac{\partial z}{\partial x} = xye^x + 1, \quad \frac{\partial z}{\partial y} = (x - 1)e^x + 1.$$

Hint: First, check to see whether there exist such functions. Then, find possible antiderivatives to the two conditioning equations.

**☞ Mastery Check 2.16:**

Determine all C^2 functions $f(x, y)$ such that

a) $\frac{\partial f}{\partial x} = 2x \sin x^2, \quad \frac{\partial f}{\partial y} = \cos y.$

b) $\frac{\partial f}{\partial x} = 2x + y, \quad \frac{\partial f}{\partial y} = 2y + x.$

c) $\frac{\partial f}{\partial x} = x + 3yx^2, \quad \frac{\partial f}{\partial y} = x^3 + xy.$



2.G Composite functions and the chain rule

We now come to a topic which many find challenging. However, it is so important in multivariable calculus as well as in practice that we will devote some considerable space to it.

What are *composite functions*? These are functions *of functions*, of one or more independent variables. The relationships between the functions and their variable dependencies can be readily represented by *ball-and-stick diagrams* (see below). Although we will retain function names in our analyses, the ball-and-stick diagrams show the relationships between dependent and independent variables. However, the reader should bear in mind that the functions themselves actually provide the links between the variables. These links are illustrated with the help of *domain-and-range diagrams*, which seek to aid understanding not only of the dependencies but also of the conditions that must be satisfied for the composite functions to be defined.

What then is the *chain rule*? The simple truth that this is the process by which one differentiates composite functions is rather unhelpful at this point. It will be necessary to go through the various cases we will be considering, in order of complexity, for this statement to have meaning.

As for illustrating the chain rule — as distinct from visualizing composite functions — we take advantage of the notion that derivatives describe rates of change and so imagine a derivative to represent the regulated flow of water through a sluice gate or *floodgate* from one water reservoir to another. The chain rule, which involves sums of products of derivatives, we shall represent by *floodgate diagrams*: an arrangement of reservoirs fitted with floodgates to regulate water flow. The net flow of water out through the final gate will depend on which gates above it are open (the variables), by how much (the partial derivatives), how two or more flow rates reinforce (the products), and in what combinations (the sums).

In this context, probably more than any other, it is important to distinguish between the independent variable that is involved in a partial derivative and others that are held fixed. To this end we will use notation such as

$$\left(\frac{\partial f}{\partial x}\right)_y \quad \text{and} \quad \left(\frac{\partial F}{\partial u}\right)_v$$

to refer to partial derivatives (here w.r.t. x and u , respectively) and the independent variables that are kept constant (here y and v , respectively).

Case 1

This is the simplest example which the student would have encountered in their single-variable calculus course. It nevertheless exhibits all the features inherent in the more complicated multivariable cases to follow. Accordingly, the format we follow in this discussion is repeated in the latter cases. Within this format we itemize the variable dependence of the functions involved including their domains and ranges, the composite function and its domain and range, and finally the appropriate chain rule for derivatives of the composite function.

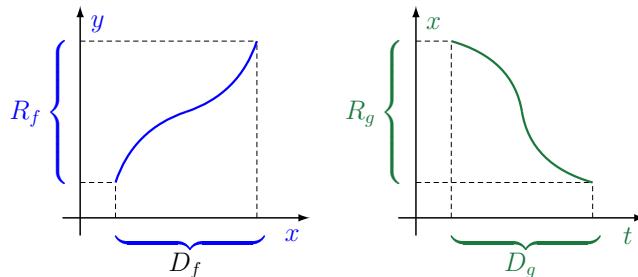


Figure 2.20 The graphs of $y = f(x)$ and $x = g(t)$.

Consider two functions $f, g \in C^1(\mathbb{R})$ of one variable.

$$\begin{aligned} f : \mathbb{R} &\longrightarrow \mathbb{R} \\ x \longmapsto y &= f(x), \end{aligned}$$

$$\begin{aligned} g : \mathbb{R} &\longrightarrow \mathbb{R} \\ t \longmapsto x &= g(t). \end{aligned}$$

The domains and ranges of these are shown in Figure 2.20 above. From these functions we form the composite function of the independent variable t :

$$y = F(t) = (f \circ g)(t) = f(g(t)).$$

The composite function may be represented schematically by the ball-and-stick diagram in Figure 2.21. The diagram (as with the more complex ones to follow) indicates that the variable y depends on x , which in turn depends on t . Thus, a variation in t leads to a variation in x , which leads to a variation in y .



Figure 2.21 Ball-and-stick model for $f(g(t))$.

The domain and range of F , which are based on the sets D_g , R_g , D_f , and R_f , must be such that F makes sense. Examine Figure 2.22, from left to right, noting the termini (start and end points) of the arrows.

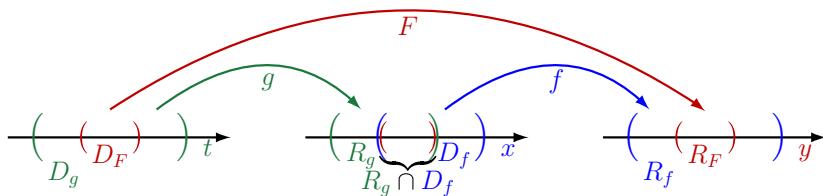


Figure 2.22 Conditional definition of D_F and R_F .

From the figure follow two set relations and one critical set constraint: the domain of F is a subset of the domain of g , the range of F is a subset of the range of f , and then there must be a non-empty intersection of the range of g and domain of f . In set notation these are summarized as follows.

$$\left. \begin{array}{l} \text{Domain of } F : D_F \subseteq D_g \\ \text{Range of } F : R_F \subseteq R_f \end{array} \right\} \text{An important condition: } R_g \cap D_f \neq \emptyset.$$

The derivative of F w.r.t. t is given by the chain rule and can be represented schematically by the *floodgate* diagram in Figure 2.23.

$$\begin{aligned}
 \frac{dF}{dt} &= \frac{d}{dt}(f \circ g)(t) \\
 &= \frac{df}{dx} \Big|_{x=g(t)} \cdot \frac{dg}{dt} \Big|_t \\
 &= \left\{ \begin{array}{l} \text{derivative} \\ \text{of } f \\ \text{evaluated} \\ \text{at } x = g(t) \end{array} \right\} \times \left\{ \begin{array}{l} \text{derivative} \\ \text{of } g \\ \text{evaluated} \\ \text{at } t \end{array} \right\}
 \end{aligned}$$

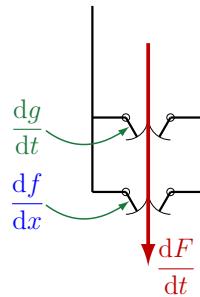


Figure 2.23 The floodgate diagram for dF/dt .

In constructing the relevant floodgate diagram for any given case we aim to compile and represent all the term- and factor-wise contributions to a chain rule derivative. The objective is the net outflow at the bottom. In this case the final flow rate out through the bottom gate of the reservoir ($\frac{dF}{dt}$) is dependent not only on the flow rate from the x -reservoir to the y -reservoir ($\frac{df}{dx}$) but also on the flow rate from the t -reservoir to the x -reservoir ($\frac{dg}{dt}$): one reinforces (multiplies) the other.

A final note concerns the notation used to represent and describe the chain rule derivative. The more commonly seen notation is

$$\frac{dy}{dt} = \frac{dy}{dx} \cdot \frac{dx}{dt}.$$

Although it is intuitive and appealing to express the chain rule in this way (admittedly it is convenient sometimes), this notation can be problematic in some cases.

Case 2

We now move on to more complicated functional arrangements. As was remarked earlier, the format for the discussion here remains the same as in Case 1.

Consider function f as before, but now suppose that g is a function of two variables.

$$\begin{aligned}
 f : \mathbb{R} &\longrightarrow \mathbb{R} \\
 x \longmapsto y &= f(x)
 \end{aligned}$$

$$\begin{aligned}
 g : \mathbb{R}^2 &\longrightarrow \mathbb{R} \\
 (s, t) \longmapsto x &= g(s, t)
 \end{aligned}$$

The respective domains and ranges are shown schematically in Figure 2.24.

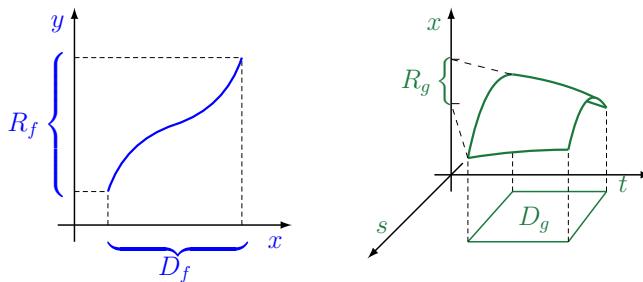


Figure 2.24 The graphs of $y = f(x)$ and $x = g(s, t)$.

The composite function of the two independent variables, s and t , is

$$y = F(s, t) = (f \circ g)(s, t) = f(g(s, t)).$$

The composite function can here too be represented schematically by a ball-and-stick diagram, but this time a *branch* diagram (Figure 2.25). The variable y depends on x , which depends on s and t . Consequently, a variation in s or t leads to a variation in x , which leads in turn to a variation in y .

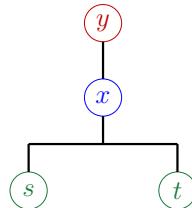


Figure 2.25 Ball-and-stick for $f(g(s, t))$.

The domain and range of F , which again are determined by D_g , R_g , D_f , and R_f , are such that F makes sense (Figure 2.26).

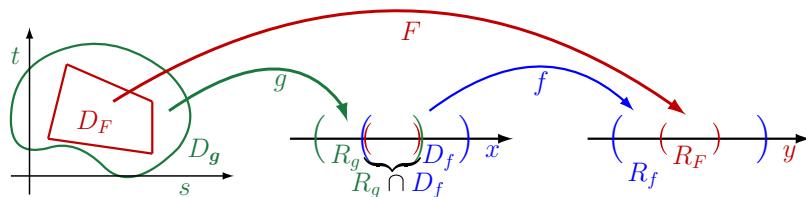


Figure 2.26 Conditional definition of D_F and R_F .

We once again find two set relations and the same set constraint.

$$\left. \begin{array}{l} \text{Domain of } F : D_F \subseteq D_g \\ \text{Range of } F : R_F \subseteq R_f \end{array} \right\} \text{An important condition: } R_g \cap D_f \neq \emptyset.$$

That is, the domain of F is not bigger than, and can be smaller than, the domain of g , and the range of F is not bigger and can be smaller than the range of f . All is dictated by the size of $R_g \cap D_f$.

The derivatives of F w.r.t. s and t are now *partial derivatives* given by the chain rule, and represented by their respective *floodgate diagrams*: Figures 2.27 and 2.28.

$$\left(\frac{\partial F}{\partial s} \right)_t = \frac{\partial}{\partial s} (f \circ g)(s, t)$$

t — held fixed

$$= \frac{df}{dx} \Big|_{x=g(s,t)} \cdot \left(\frac{\partial g}{\partial s} \right)_t$$

$$= \left\{ \begin{array}{l} \text{full} \\ \text{derivative} \\ \text{of } f \\ \text{evaluated} \\ \text{at } x = g(s, t) \end{array} \right\} \times \left\{ \begin{array}{l} \text{partial} \\ \text{derivative} \\ \text{of } g \\ \text{evaluated} \\ \text{at } (s, t) \end{array} \right\}$$

Figure 2.27 The floodgate diagram for $\partial F / \partial s$.

$$\left(\frac{\partial F}{\partial t} \right)_s = \frac{\partial}{\partial t} (f \circ g)(s, t)$$

s — held fixed

$$= \frac{df}{dx} \Big|_{x=g(s,t)} \cdot \left(\frac{\partial g}{\partial t} \right)_s$$

$$= \left\{ \begin{array}{l} \text{full} \\ \text{derivative} \\ \text{of } f \\ \text{evaluated} \\ \text{at } x = g(s, t) \end{array} \right\} \times \left\{ \begin{array}{l} \text{partial} \\ \text{derivative} \\ \text{of } g \\ \text{evaluated} \\ \text{at } (s, t) \end{array} \right\}$$

Figure 2.28 The floodgate diagram for $\partial F / \partial t$.

In this case (in contrast with Case 1) the partial derivatives of F mean that one variable is kept fixed, and thus its associated gate remains closed, giving

no contribution to the flow out through the bottom floodgate.

We demonstrate how this case works with an example and leave the reader with an exercise to consolidate their understanding.

■ Example 2.11:

Consider the functions $f : x \rightarrow y = f(x)$, and $g : (s, t) \rightarrow x = g(s, t)$.

We wish to find the domain D_F of the composite function

$F : (s, t) \rightarrow y = (f \circ g)(s, t)$, and the derivatives $\left(\frac{\partial F}{\partial s}\right)_t, \left(\frac{\partial F}{\partial t}\right)_s$, when $f(x) = \ln x$, $g(s, t) = s(1 - t^2)$.

The domain of f is $D_f = \{x : x > 0\}$, and the range of g is

$R_g = \{x : x \in \mathbb{R}\}$. The intersection is $\{x : x = s(1 - t^2) > 0\}$, that is, $D_F = \{(s, t) : (s > 0 \text{ and } |t| < 1) \cup (s < 0 \text{ and } |t| > 1)\}$.

$$\begin{aligned}\left(\frac{\partial F}{\partial s}\right)_t &= \frac{df}{dx} \frac{\partial g}{\partial s} = \frac{1}{x} \bigg|_{x=s(1-t^2)} (1-t^2) = \frac{1}{s}. \\ \left(\frac{\partial F}{\partial t}\right)_s &= \frac{df}{dx} \frac{\partial g}{\partial t} = \frac{1}{x} \bigg|_{x=s(1-t^2)} (-2st) = -\frac{2t}{1-t^2}.\end{aligned}$$



☛ Mastery Check 2.17:

Consider the function: $y = f(x) = \arcsin x$, where $x = g(s, t) = s^2 + \frac{1}{t}$. What are D_f , R_f , D_g , and R_g ?

Determine where $y = F(s, t)$ makes sense, and then find (if possible) the partial derivatives $\left(\frac{\partial F}{\partial s}\right)_t$ and $\left(\frac{\partial F}{\partial t}\right)_s$.

Note: the final result should be expressed in terms of s and t !



Case 3

Consider two functions $f, \mathbf{g} \in C^1(\mathbb{R}^2)$ of two variables.

$$\begin{aligned}f : \mathbb{R}^2 &\rightarrow \mathbb{R} \\ (x, y) &\mapsto z = f(x, y)\end{aligned}$$

$$\begin{aligned}\mathbf{g} : \mathbb{R}^2 &\rightarrow \mathbb{R}^2 \\ (s, t) &\mapsto (x = g_1(s, t), y = g_2(s, t))\end{aligned}$$

The composite function F of two variables s and t derived from f and \mathbf{g} is

$$z = F(s, t) = (f \circ \mathbf{g})(s, t) = f(g_1(s, t), g_2(s, t)).$$

This composite function is represented by the more elaborate branch model of dependent and independent variables shown in Figure 2.29.

This time z depends on x and y , and both x and y depend on s and t .

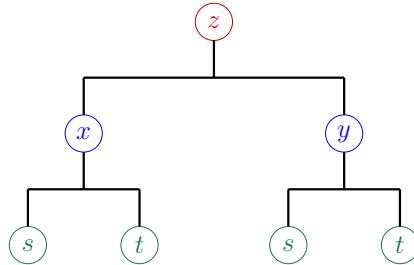


Figure 2.29 Ball-and-stick for $f(g_1(s, t), g_2(s, t))$.

The domain and range of F , which are dictated by D_f and R_f , and D_g and R_g , are such that F makes sense. In analogy with Cases 1 and 2, the same set conditions and set constraint can be established from the depiction in Figure 2.30:

$$\left. \begin{array}{l} \text{Domain of } F : D_F \subseteq D_g \\ \text{Range of } F : R_F \subseteq R_f \end{array} \right\} \text{An important condition: } R_g \cap D_f \neq \emptyset.$$

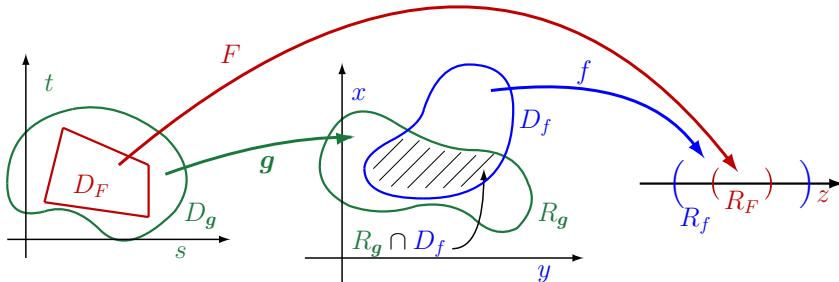


Figure 2.30 Conditional definition of D_F and R_F .

The partial derivatives of F w.r.t. s and t are given by the chain rule:

$$\left(\frac{\partial F}{\partial s}\right)_t = \left(\frac{\partial f}{\partial x}\right)_y \Bigg|_{\substack{x=g_1(s,t) \\ y=g_2(s,t)}} \cdot \left(\frac{\partial g_1}{\partial s}\right)_t + \left(\frac{\partial f}{\partial y}\right)_x \Bigg|_{\substack{x=g_1(s,t) \\ y=g_2(s,t)}} \cdot \left(\frac{\partial g_2}{\partial s}\right)_t$$

t — held constant

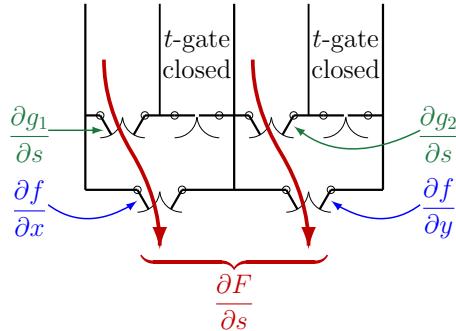


Figure 2.31 The floodgate diagram for $\partial F/\partial s$.

$$\left(\frac{\partial F}{\partial t}\right)_s = \left(\frac{\partial f}{\partial x}\right)_y \Bigg|_{\substack{x=g_1(s,t) \\ y=g_2(s,t)}} \cdot \left(\frac{\partial g_1}{\partial t}\right)_s + \left(\frac{\partial f}{\partial y}\right)_x \Bigg|_{\substack{x=g_1(s,t) \\ y=g_2(s,t)}} \cdot \left(\frac{\partial g_2}{\partial t}\right)_s$$

s — held constant

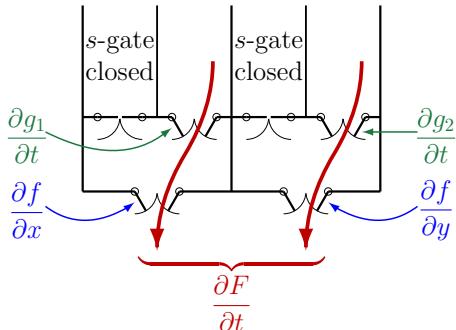


Figure 2.32 The floodgate diagram for $\partial F/\partial t$.

In each of the cases shown in Figures 2.31 and 2.32, both bottom floodgates ($\frac{\partial f}{\partial x}$ and $\frac{\partial f}{\partial y}$) are open and contribute to the total flow, but with strengths that are modulated by the floodgates above them, that is, by the partial derivatives g_1 and g_2 .

☞ Mastery Check 2.18:

Consider the following function:

$$z = f(x, y) = \sin(x^2 y), \text{ where } x = g_1(s, t) = st^2 \text{ and } y = g_2(s, t) = s^2 + \frac{1}{t}.$$

Let $F(s, t) = (f \circ \mathbf{g})(s, t)$. Find $\left(\frac{\partial F}{\partial s}\right)_t$ and $\left(\frac{\partial F}{\partial t}\right)_s$, if they make sense.



Case 4

As a final example, consider two functions $f \in C^1(\mathbb{R}^3)$, $\mathbf{g} \in C^1(\mathbb{R}^2)$. This is a mixed case where the “outer” function, f , depends on an independent variable both directly and indirectly.

$$f : \mathbb{R}^3 \longrightarrow \mathbb{R}$$

$$(x, y, t) \longmapsto z = f(x, y, t)$$

$$\mathbf{g} : \mathbb{R}^2 \longrightarrow \mathbb{R}^2$$

$$(s, t) \longmapsto (x = g_1(s, t), y = g_2(s, t))$$

The composite function F of two variables s and t formed from f and \mathbf{g} is

$$z = F(s, t) = f(g_1(s, t), g_2(s, t), t).$$

The ball-and-stick branch model appropriate for this example is shown below in Figure 2.33.

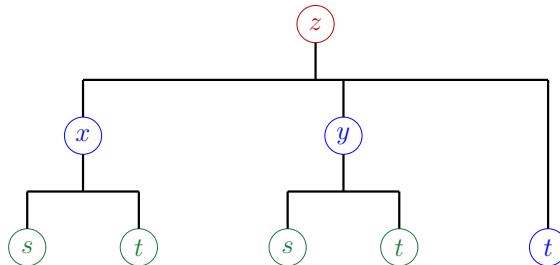


Figure 2.33 The complex ball-and-stick diagram for $f(g_1(s, t), g_2(s, t), t)$.

As in all previous examples, the sets D_f , R_f , D_g , and R_g establish the domain and range of F so that F makes sense. However, this time there is the added complication of the appearance of a common independent variable, t .

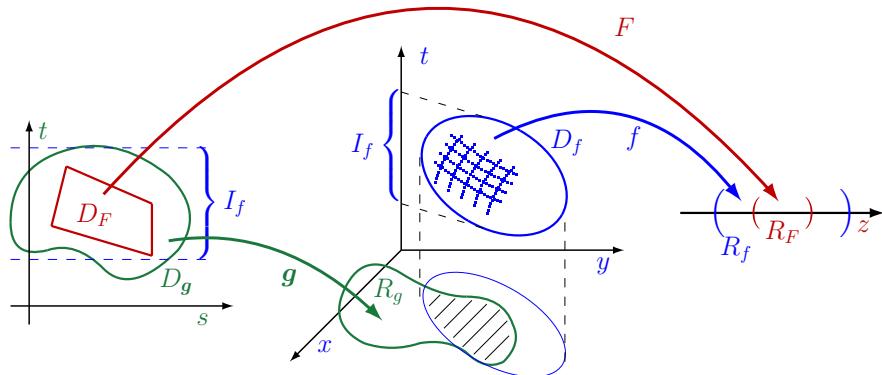


Figure 2.34 Conditional definition of D_F and R_F .

Note the complex intersections of domains and ranges in Figure 2.34. The particular complication here is the fact that t -values must lie in I_f , the t -interval making up *one* dimension of D_f , as well as in the *one* dimension within the domain of g .

To reconcile the different sets, let PD_f be the projection of D_f in the xy -plane and consider the infinite strip $I_f \times \mathbb{R} \subset \mathbb{R}^2$ in the st -plane. The rectangular strip $I_f \times \mathbb{R}$ is shown in the left-hand diagram in Figure 2.34.

We can now formally establish the range and domain of F :

$$\left. \begin{array}{l} \text{Domain of } F: D_F \subseteq D_g \cap (I_f \times \mathbb{R}) \\ \text{Range of } F: R_F \subseteq R_g \end{array} \right\}.$$

There are two important conditions: $\left\{ \begin{array}{l} R_g \cap PD_f \neq \emptyset \\ D_g \cap (I_f \times \mathbb{R}) \neq \emptyset \end{array} \right\}$.

The preceding two conditions are critically important for the validity of the composite function. The domain of F must be consistent with $R_g \cap PD_f$, but the allowed t -values in D_g must also be within I_f . Consequently, t in f is *not* independent of t in g ! Note also that the limits of the interval I_f *may* depend on x and y values in PD_f .

We will meet this idea again in another context in Section 4.G.

The partial derivative of F w.r.t. s is given by the chain rule, largely identical to the previous case:

$$\left(\frac{\partial F}{\partial s}\right)_t = \left(\frac{\partial f}{\partial x}\right)_{y,t} \Bigg|_{\substack{x=g_1(s,t) \\ y=g_2(s,t)}} \cdot \left(\frac{\partial g_1}{\partial s}\right)_t + \left(\frac{\partial f}{\partial y}\right)_{x,t} \Bigg|_{\substack{x=g_1(s,t) \\ y=g_2(s,t)}} \cdot \left(\frac{\partial g_2}{\partial s}\right)_t$$

t — held constant

The corresponding floodgate model for $\left(\frac{\partial F}{\partial t}\right)_s$ is also effectively as appears in the preceding case. On the other hand, the partial derivative of F w.r.t. t is given by a version of the chain rule that has three contributions (see Figure 2.35).

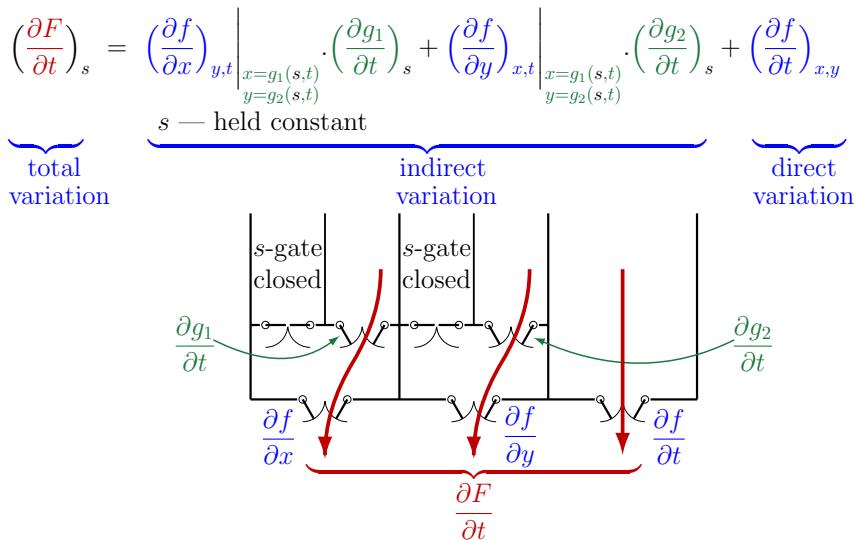


Figure 2.35 The floodgate diagram for $\partial F / \partial t$.

The indirect contribution to $\left(\frac{\partial F}{\partial t}\right)_s$ is as in the preceding case, but now there is an extra contribution from the direct dependence on t . This is reflected in Figure 2.35 by the feature of the reservoir above the floodgate $\left(\frac{\partial f}{\partial t}\right)_{x,y}$ having no other influences, while there are t -reservoirs that influence (multiply or reinforce) the other contributions.

The reader should compare the expression for the partial derivative $\left(\frac{\partial F}{\partial t}\right)_s$ with the expression one would write using the less precise notation that treats x, y, z , as both dependent and independent variables. With $z = F(s, t)$ and $z = f(x, y, t)$, the chain rule for the partial derivative with respect to t would then be written

$$\frac{\partial z}{\partial t} = \underbrace{\frac{\partial z}{\partial x} \cdot \frac{\partial x}{\partial t}}_{\substack{\text{total} \\ \text{variation}}} + \underbrace{\frac{\partial z}{\partial y} \cdot \frac{\partial y}{\partial t}}_{\substack{\text{indirect} \\ \text{variation}}} + \underbrace{\frac{\partial z}{\partial t}}_{\text{direct variation}}$$

We see that by being imprecise we arrive at an expression involving two terms with the same notation but which mean different things!

If the reader insists on using x, y, z as both functions and independent variables instead of just as independent variables, then they should at least write the derivatives as

$$\frac{\partial z}{\partial t} = \left(\frac{\partial z}{\partial x}\right)_{y,t} \cdot \frac{\partial x}{\partial t} + \left(\frac{\partial z}{\partial y}\right)_{x,t} \cdot \frac{\partial y}{\partial t} + \left(\frac{\partial z}{\partial t}\right)_{x,y}.$$

We have completed our coverage of composite functions and their first partial derivatives. Of course, these four cases do not exhaust all possibilities. However, the reader may well discover that more complicated cases of composite functions and their respective partial derivatives may be readily if not easily constructed by generalizing the exposition given here.

Mastery Check 2.19:

Consider the function

$$z = f(x, y, t) = xt \cos y,$$

where

$$x = g_1(s, t) = st + 1, \quad y = g_2(s, t) = s^2 + t^2.$$

Let $(D_f)_t \subseteq \mathbb{R}^2$ denote the domain of f when t is held constant.

Your task is first to establish that $(D_f)_t \cap R_g \neq \emptyset$, and then to determine the partial derivatives of $F(s, t) = (f \circ g)(s, t)$ w.r.t. s and t .



Second derivatives and the chain rule

Applications involving the chain rule are not limited to first partial derivatives. So, while it is essential to understand the chain rule in principle, a Master Class in the practical use of the chain rule cannot be complete without a discussion of higher-order partial derivatives. In the author's experience, this is what most students find challenging.

By way of demonstration consider [Case 3](#) again:

$f : \mathbb{R}^2 \rightarrow \mathbb{R}$, and $\mathbf{g} : \mathbb{R}^2 \rightarrow \mathbb{R}^2$, with

$$F(s, t) = (\mathbf{f} \circ \mathbf{g})(s, t) = \mathbf{f}(\mathbf{g}_1(s, t), \mathbf{g}_2(s, t)).$$

The first partial derivative of F w.r.t. s is (as on [Page 91](#)):

$$\frac{\partial F}{\partial s} = \frac{\partial f}{\partial x} \cdot \frac{\partial g_1}{\partial s} + \frac{\partial f}{\partial y} \cdot \frac{\partial g_2}{\partial s},$$

where for convenience we have suppressed parentheses and their subscripts.

If we now want a second derivative, say $\frac{\partial^2 F}{\partial t \partial s}$, then we must take note of two facts:

- (a) $\frac{\partial F}{\partial s}$ is the sum of *products* of functions!
- (b) $\frac{\partial f}{\partial x}$ and $\frac{\partial f}{\partial y}$ are *two new composite functions of s and t !* Let's denote these by K and H .

Hence, using

$$\frac{\partial f}{\partial x} = k(x(s, t), y(s, t)) = K(s, t), \quad \frac{\partial f}{\partial y} = h(x(s, t), y(s, t)) = H(s, t),$$

the first partial derivative of F w.r.t. s will become

$$\frac{\partial F}{\partial s} = K(s, t) \cdot \frac{\partial g_1}{\partial s} + H(s, t) \cdot \frac{\partial g_2}{\partial s}.$$

In this guise, the first partial derivative is more clearly seen to be a sum of products of functions of s and t . Consequently, in taking a second partial derivative — any second partial derivative for that matter — we must take the following steps in the order given:

Step 1: use the product rule of differentiation;

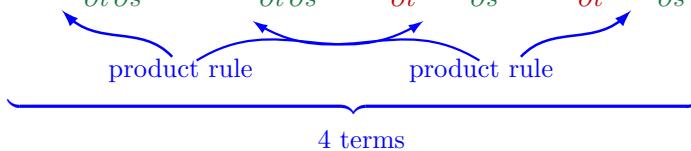
Step 2: use the chain rule *again*, this time on K and H ;

Step 3: express everything in terms of the independent variables of F .

As we said, these actions *must* be taken *in this order!*

Step 1: Differentiate the sum of products using the product rule

$$\begin{aligned} \frac{\partial^2 F}{\partial t \partial s} &= \frac{\partial}{\partial t} \left(\frac{\partial F}{\partial s} \right) = \frac{\partial}{\partial t} \left(K(s, t) \frac{\partial g_1}{\partial s} \right) + \frac{\partial}{\partial t} \left(H(s, t) \frac{\partial g_2}{\partial s} \right) \\ &= K(s, t) \frac{\partial^2 g_1}{\partial t \partial s} + H(s, t) \frac{\partial^2 g_2}{\partial t \partial s} + \frac{\partial K(s, t)}{\partial t} \cdot \frac{\partial g_1}{\partial s} + \frac{\partial H(s, t)}{\partial t} \cdot \frac{\partial g_2}{\partial s} \end{aligned}$$



The first two terms are finished.

Step 2: After the product rule, apply the *chain* rule (again) on the second two terms:

$$\begin{aligned} \frac{\partial K}{\partial t} &= \frac{\partial k}{\partial x} \cdot \frac{\partial g_1}{\partial t} + \frac{\partial k}{\partial y} \cdot \frac{\partial g_2}{\partial t} = \frac{\partial^2 f}{\partial x^2} \cdot \frac{\partial g_1}{\partial t} + \frac{\partial^2 f}{\partial y \partial x} \cdot \frac{\partial g_2}{\partial t} \\ \frac{\partial H}{\partial t} &= \frac{\partial h}{\partial x} \cdot \frac{\partial g_1}{\partial t} + \frac{\partial h}{\partial y} \cdot \frac{\partial g_2}{\partial t} = \frac{\partial^2 f}{\partial x \partial y} \cdot \frac{\partial g_1}{\partial t} + \frac{\partial^2 f}{\partial y^2} \cdot \frac{\partial g_2}{\partial t}. \end{aligned}$$

Step 3: Finally, replace all the K and H factors with the f , g_1 , and g_2 factors (with all derivatives of f evaluated at $x = g_1(s, t)$ and $y = g_2(s, t)$):

$$\begin{aligned} \frac{\partial^2 F}{\partial t \partial s} &= \frac{\partial f}{\partial x} \Bigg|_{\substack{x=g_1 \\ y=g_2}} \frac{\partial^2 g_1}{\partial t \partial s} + \frac{\partial f}{\partial y} \Bigg|_{\substack{x=g_1 \\ y=g_2}} \frac{\partial^2 g_2}{\partial t \partial s} \\ &\quad + \frac{\partial g_1}{\partial s} \left(\frac{\partial^2 f}{\partial x^2} \Bigg|_{\substack{x=g_1 \\ y=g_2}} \frac{\partial g_1}{\partial t} + \frac{\partial^2 f}{\partial y \partial x} \Bigg|_{\substack{x=g_1 \\ y=g_2}} \frac{\partial g_2}{\partial t} \right) \\ &\quad + \frac{\partial g_2}{\partial s} \left(\frac{\partial^2 f}{\partial x \partial y} \Bigg|_{\substack{x=g_1 \\ y=g_2}} \frac{\partial g_1}{\partial t} + \frac{\partial^2 f}{\partial y^2} \Bigg|_{\substack{x=g_1 \\ y=g_2}} \frac{\partial g_2}{\partial t} \right). \end{aligned}$$

For this example the second partial derivative has six terms in total!

 **Mastery Check 2.20:**

Consider the function $z = f(x, y) = x \cos y + y$,

where $x = u(s, t) = st$, $y = v(s, t) = s^2 + t^2$. Determine $\frac{\partial^2 z}{\partial s \partial t}$.



The Leibniz integral rule

A particularly important application of the chain rule is to differentiating an integral such as

$$\frac{d}{dt} \int_{a(t)}^{b(t)} h(x, t) dx,$$

with respect to a parameter. Suppose h , a and b are C^1 functions and the indicated integral, as well as the integral of $\partial h / \partial t$, exists. The integral itself produces a function $z = F(t) = (f \circ \mathbf{g})(t)$, which depends on t through three channels. Let's call these u , v and t , where u and v take the place of the upper and lower limits of the integral, and where

$$\begin{aligned} f : \mathbb{R}^3 &\longrightarrow \mathbb{R}, \\ (u, v, t) &\longmapsto z = \int_u^v h(x, t) dx, \\ \text{and } \mathbf{g} &= (u = a(t), v = b(t)). \end{aligned}$$

The branch model relevant to this is shown in Figure 2.36. It is a somewhat simplified version of Case 4 on Page 93.

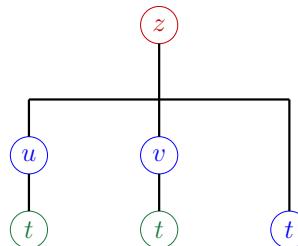


Figure 2.36 Ball-and-stick for $f(g_1(t), g_2(t), t)$.

Referring back to that Case 4, the derivative of F with respect to t is thus

$$\begin{aligned} \frac{dF}{dt} &= \left(\frac{\partial f}{\partial u} \right) \bigg|_{u=b(t), v=a(t), t} \cdot \left(\frac{db}{dt} \right) \\ &\quad + \left(\frac{\partial f}{\partial v} \right) \bigg|_{u=b(t), v=a(t), t} \cdot \left(\frac{da}{dt} \right) + \left(\frac{\partial f}{\partial t} \right) \bigg|_{u=b(t), v=a(t), t}. \end{aligned}$$

The partial derivatives of f with respect to u and v are straightforward using the fundamental theorem of integral calculus:

$$g(x) = \pm \frac{dG}{dx} \iff G(x) = \pm \int_c^x g(x') dx'$$

where c is some constant. Thus,

$$\left(\frac{\partial f}{\partial u} \right) \Big|_{(b(t), a(t), t)} = h(b(t), t) \quad \text{and} \quad \left(\frac{\partial f}{\partial v} \right) \Big|_{(b(t), a(t), t)} = -h(a(t), t).$$

For the partial derivative of f with respect to t we use the definition:

$$\frac{f(u, v, t + \Delta t) - f(u, v, t)}{\Delta t} = \int_u^v \left[\frac{h(x, t + \Delta t) - h(x, t)}{\Delta t} \right] dx.$$

where we are permitted to put everything under the one integral sign since f is C^1 . Now, taking the limit $\Delta t \rightarrow 0$ we get

$$\frac{\partial f}{\partial t} = \int_u^v \frac{\partial h}{\partial t}(x, t) dx.$$

All together, replacing u and v everywhere with $a(t)$ and $b(t)$, respectively, we have the very useful Leibniz rule

$$\begin{aligned} \frac{d}{dt} \int_{a(t)}^{b(t)} h(x, t) dx \\ = h(b(t), t) \cdot \left(\frac{db}{dt} \right) - h(a(t), t) \cdot \left(\frac{da}{dt} \right) + \int_{a(t)}^{b(t)} \frac{\partial h}{\partial t}(x, t) dx \end{aligned}$$

Example 2.12:

We apply this rule to the following integrals. Note the use of different independent variables.

$$(1) \text{ Suppose } F(x) = \int_0^{x^2} \sin(u^2) du. \text{ Then } F'(x) = 2x \sin(x^4).$$

$$(2) \text{ Suppose } F(u) = \int_{1-\ln u}^1 e^{t^2} dt. \text{ Then } F'(u) = \frac{1}{u} e^{(1-\ln u)^2}.$$

$$(3) \text{ Suppose } F(t) = \int_{\sin t}^{\cos t} e^{2xt} dx. \text{ Then}$$

$$F'(t) = -e^{2t \cos t} \sin t - e^{2t \sin t} \cos t + \int_{\sin t}^{\cos t} 2x e^{2xt} dx.$$



2.H Implicit functions

Suppose we are given the following task: In each of the cases below express the variable y as a function of the remaining variables:

- (a) $8y + 64x^2 = 0$;
- (b) $2y^2 + 8y + 16z \sin x = 0$;
- (c) $\ln |y| + y^3x + 20x^2 = w$.

I am as certain that you cannot complete task (c) as I am that you *can* complete tasks (a) and (b). Although task (c) is impossible, the equation suggests there is a functional relationship, in principle.

This introduces the notion of an implied or *implicit function*. In task (c) the equation implies that y can be a function f of the variables x and w . What we shall do in this section is establish conditions under which such a function is defined, *at least locally*. Along the way we will get, as reward, a linear approximation to this unknown function, in terms of the independent variables near a given point, and an explicit expression, and value, for the derivative (or derivatives) of this function at that point.

As before we explain by considering examples of increasing complexity. In each case we will also discuss an analogous linear problem. Since the argument we follow is based on linearization, we hope that the parallels will facilitate reader understanding. The purist reader may frown on the questionable rigour. However, the possibility of greater appreciation for the end result is worth sacrificing some degree of mathematical sophistication.

Suppose we are given the following three problems:

$$1) \ e^{x+y} + xy = 0 \implies F(x, y) = 0 \quad \text{— a level curve.}$$

$$2) \ e^{x+y+z} - (x + y + z)^2 = 1 \implies F(x, y, z) = 0 \quad \text{— a level surface.}$$

$$3) \ \begin{cases} e^{x+y+z} - (x + y + z)^2 - 1 = 0 \\ z \sin(xy) - x \cos(zy) = 0 \end{cases} \implies \begin{cases} F(x, y, z) = 0 \\ G(x, y, z) = 0 \end{cases} \quad \text{— a curve of intersection.}$$

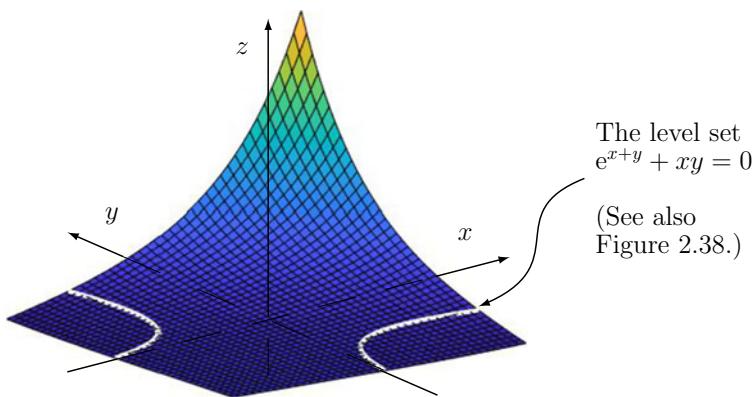


Figure 2.37 The graph of $z = e^{x+y} + xy$ and the level set $F(x, y) = 0$.

Consider Problem 1. $F(x, y) = e^{x+y} + xy = 0$.

This equation is nonlinear. On top of that, it cannot be manipulated to get y in terms of x . All the same, the level curve is shown in Figure 2.37.

But suppose we consider the linear approximation to $F(x, y) = e^{x+y} + xy$ for points (x, y) about a point (a, b) which lies on the level curve $F(x, y) = 0$. The linear approximation is shown in Figure 2.38 and developed on Page 103.

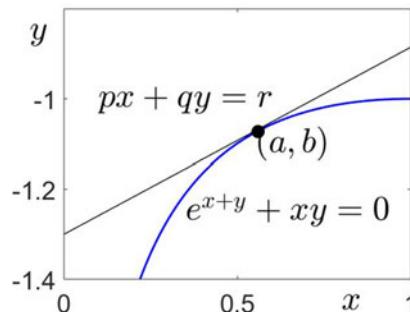


Figure 2.38 The level curve $e^{x+y} + xy = 0$ and its approximation.

As the analysis parallels the description of a straight line in 2D, it (the general analysis) is preceded by that simple geometric discussion. The student should compare the two mathematical arguments.

Consider the level set describing the general form for the equation of a 2D line (Figure 2.39):

$$ax + by = c, \quad a, b, c \neq 0$$

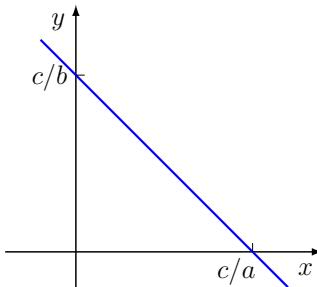


Figure 2.39 The line $ax + by = c$.

We divide this linear equation by b , the coefficient of y , and solve for y . Naturally, we get the equation of a line in standard form.

$$y = -\frac{a}{b}x + \frac{c}{b} = mx + k,$$

where $-\infty < m = -a/b < \infty$ as long as $b \neq 0$. This is an important point for what follows: *The coefficient of the variable we solved for cannot be zero.*

Now consider a more general and potentially nonlinear function, F , such as the one suggested here.

If $F \in C^1$, then for all (x, y) near (a, b) (see Page 68), we know that a linear approximation to F is obtained from

$$F(x, y) = F(a, b) + \frac{\partial F}{\partial x} \Big|_{(a,b)} \cdot (x - a) + \frac{\partial F}{\partial y} \Big|_{(a,b)} \cdot (y - b) + |\Delta x| \rho(\Delta x)$$

where, on the level curve itself, both $F(x, y)$ and $F(a, b)$ are zero:

$$0 = \frac{\partial F}{\partial x} \Big|_{(a,b)} \cdot (x - a) + \frac{\partial F}{\partial y} \Big|_{(a,b)} \cdot (y - b) + |\Delta x| \rho(\Delta x),$$

or, equivalently,

$$\frac{\partial F}{\partial x} \Big|_{(a,b)} \cdot x + \frac{\partial F}{\partial y} \Big|_{(a,b)} \cdot y = \frac{\partial F}{\partial x} \Big|_{(a,b)} \cdot a + \frac{\partial F}{\partial y} \Big|_{(a,b)} \cdot b - |\Delta x| \rho(\Delta x).$$

This is in the general form of the equation of a line $px + qy = r$. Therefore, just as we did in the simple linear case we obtain $y = -(p/q)x + r/q$, provided

$q \neq 0$. This last condition is a key point to remember: the solution for y for points (x, y) near (a, b) is valid *as long as* $q = \partial F / \partial y \neq 0$ at (a, b) . And, in that case we have the linear approximation to the curve $F(x, y)$:

$$y = -\frac{\frac{\partial F}{\partial x} \Big|_{(a,b)}}{\frac{\partial F}{\partial y} \Big|_{(a,b)}} x + \frac{r}{q}$$

What we are saying is that we have thus obtained an approximate representation for the level curve for points (x, y) near (a, b) that *defines* y locally as a function of x . This is given by the straight line in Figure 2.38. (The graph of the actual function is shown in blue.)

Also we have derived the precise value of the derivative of the unknown implicit function $y = f(x)$ at (a, b) even though we cannot write f out explicitly, and in the process we are provided with a condition for the *existence* of the implicit function ($q \neq 0$). The derivative of the implicit function is, in fact,

$$\frac{dy}{dx} \Big|_{x=a} = -\frac{\frac{\partial F}{\partial x} \Big|_{(a,b)}}{\frac{\partial F}{\partial y} \Big|_{(a,b)}} \iff F \in C^1 \text{ at } (a, b), \text{ and } q = \frac{\partial F}{\partial y} \Big|_{(a,b)} \neq 0.$$

■ Example 2.13:

Suppose the equation $x^3y + 2y^3x = 3$ defines y as a function f of x in the neighbourhood of the point $(1, 1)$. We wish to find the derivative of f at $x = 1$, and a linear approximation to f near the point. Let $F(x, y) = x^3y + 2y^3x - 3$.

Note that $F \in C^1 \forall (x, y) \in \mathbb{R}^2$. Then we have

$$\frac{\partial F}{\partial x} = 3x^2y + 2y^3, \quad \frac{\partial F}{\partial y} = x^3 + 6xy^2.$$

We note that $\frac{\partial F}{\partial y} \neq 0$ at $(1, 1)$. Thus, from our linear approximation we have $\frac{dy}{dx} \Big|_{(1,1)} = -\frac{3+2}{1+6} = -\frac{5}{7}$.

The linear approximation is $y = -\frac{5}{7}x + c$. To determine c , use the fact that the line passes through $(1, 1)$, giving $y = -\frac{5}{7}x + \frac{12}{7}$. ■

☛ Mastery Check 2.21:

Show that $x^y + \sin y = 1$ defines a function $y = f(x)$ in the neighbourhood

of $(1, 0)$, and find $\frac{dy}{dx}$. Find a linear approximation to f , valid near $(1, 0)$. 

Consider Problem 2. $F(x, y, z) = e^{x+y+z} - (x + y + z)^2 - 1$.

Suppose $\mathbf{a} = (a, b, c)$ is a point on the surface. That is, suppose $F(a, b, c) = 0$. We want to know if $F(x, y, z) = 0$ defines a function, f , so that the level surface has the form $z = f(x, y)$ for points $\mathbf{x} = (x, y, z)$ on that level surface near $\mathbf{a} = (a, b, c)$. In particular, does there exist a tangent plane approximation

$$z = c + \frac{\partial f}{\partial x} \Big|_{(a,b)} \cdot (x - a) + \frac{\partial f}{\partial y} \Big|_{(a,b)} \cdot (y - b) \quad (2.3)$$

to the surface at this point? The answer depends on the behaviour of $F(x, y, z)$ near the point, and on the existence of the linear approximation to F .

We again lead with an analogy from linear algebra. Consider the equation of the plane shown in Figure 2.40:

$$ax + by + cz = d, \quad a, b, c, d > 0.$$

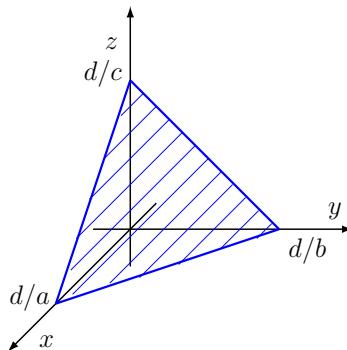


Figure 2.40 The plane $ax + by + cz = d$.

We divide this equation by c , the coefficient of z , and solve for z to get

$$z = \frac{d}{c} - \frac{a}{c}x - \frac{b}{c}y \quad \text{if } c \neq 0.$$

This plane has (partial) slope $-\frac{a}{c}$ in the x -direction and (partial) slope $-\frac{b}{c}$ in y -direction; both slopes will be finite as long as $c \neq 0$. This is an important point for our next consideration: *The coefficient of the variable we solve for*

cannot be zero. Note the outcome that, subject to this condition, we are able to express one variable, z , in terms of the other two variables, x and y .

Now let's consider the nonlinear function F of x, y, z , supposing that F is differentiable at \mathbf{a} . Under the latter condition we can obtain the linear approximation to F for points \mathbf{x} near \mathbf{a} by the methods learned earlier in this chapter:

$$F(x, y, z) = \cancel{F(\mathbf{a}, \overset{0}{b}, c)} + \frac{\partial F}{\partial x} \Big|_{(a, b, c)} \cdot (x - a) + \frac{\partial F}{\partial y} \Big|_{(a, b, c)} \cdot (y - b) + \frac{\partial F}{\partial z} \Big|_{(a, b, c)} \cdot (z - c) \\ + |\mathbf{x} - \mathbf{a}| \rho(\mathbf{x} - \mathbf{a}) \text{ (small if } \mathbf{x} \text{ near } \mathbf{a} \text{)} = 0.$$

This approximation can be rearranged to give

$$\frac{\partial F}{\partial z} \Big|_{(a, b, c)} \cdot (z - c) = - \frac{\partial F}{\partial x} \Big|_{(a, b, c)} \cdot (x - a) - \frac{\partial F}{\partial y} \Big|_{(a, b, c)} \cdot (y - b) \\ + \text{small terms for } (x, y, z) \text{ near } (a, b, c).$$

That is, for \mathbf{x} near \mathbf{a} we obtain

$$z = c - \frac{\frac{\partial F}{\partial x} \Big|_{(a, b, c)}}{\frac{\partial F}{\partial z} \Big|_{(a, b, c)}} \cdot (x - a) - \frac{\frac{\partial F}{\partial y} \Big|_{(a, b, c)}}{\frac{\partial F}{\partial z} \Big|_{(a, b, c)}} \cdot (y - b) \quad (2.4) \\ \iff \frac{\partial F}{\partial z} \Big|_{(a, b, c)} \neq 0.$$

This is the tangent plane approximation to $F(x, y, z) = 0$ at (a, b, c) .

If $\frac{\partial F}{\partial z} \Big|_{(a, b, c)} \neq 0$, then the tangent plane is well defined and is identical to the tangent plane to the surface $z = f(x, y)$.

That is, equating corresponding terms in Equations (2.3) and (2.4), we obtain explicit values of the partial derivatives of our implicit function f at (a, b) (with $z = c$).

$$\frac{\partial f}{\partial x} \Big|_{(a, b)} = - \frac{\frac{\partial F}{\partial x} \Big|_{(a, b, c)}}{\frac{\partial F}{\partial z} \Big|_{(a, b, c)}}, \quad \frac{\partial f}{\partial y} \Big|_{(a, b)} = - \frac{\frac{\partial F}{\partial y} \Big|_{(a, b, c)}}{\frac{\partial F}{\partial z} \Big|_{(a, b, c)}}.$$

So, if $\frac{\partial F}{\partial z} \Big|_{(a, b, c)} \neq 0$, then both $\frac{\partial f}{\partial x} \Big|_{(a, b)}$ and $\frac{\partial f}{\partial y} \Big|_{(a, b)}$ are well defined,

and $y = f(x, y)$ exists locally. We therefore have all that is needed to answer the question posed.

Consider Problem 3).
$$\begin{cases} F(x, y, z) = e^{x+y+z} - (x + y + z)^2 - 1 = 0 \\ G(x, y, z) = z \sin(xy) - x \cos(zy) = 0 \end{cases}.$$

The reader who might expect there to be a parallel with a linear algebraic system will not be disappointed to know that we preface the discussion with a review of two pertinent problems in linear algebra.

Suppose we have the situation posed in Figure 2.41. This represents two lines in a plane and a 2×2 system of equations for unknowns x and y .

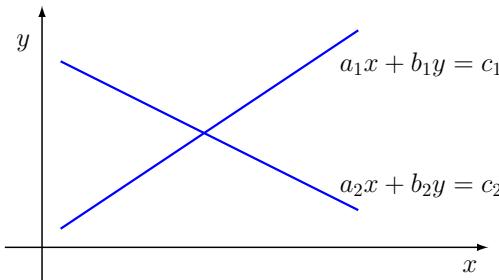


Figure 2.41 Two intersecting lines.

When conveniently expressed in matrix form, the system can be readily solved.

$$\begin{aligned} \begin{pmatrix} a_1 & b_1 \\ a_2 & b_2 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} &\Rightarrow \begin{pmatrix} x \\ y \end{pmatrix} = \frac{1}{a_1b_2 - a_2b_1} \begin{pmatrix} b_2 & -b_1 \\ -a_2 & a_1 \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} \\ \Rightarrow x = \frac{c_1b_2 - c_2b_1}{a_1b_2 - a_2b_1} &= \frac{\begin{vmatrix} c_1 & c_2 \\ b_1 & b_2 \end{vmatrix}}{\begin{vmatrix} a_1 & b_1 \\ a_2 & b_2 \end{vmatrix}}, \quad y = \frac{-c_1a_2 + c_2a_1}{a_1b_2 - a_2b_1} = \frac{\begin{vmatrix} a_1 & a_2 \\ c_1 & c_2 \end{vmatrix}}{\begin{vmatrix} a_1 & b_1 \\ a_2 & b_2 \end{vmatrix}} \end{aligned}$$

We know that this system of equations has a unique solution if and only if the determinant of the original matrix is not equal to zero. In that case the solution corresponds to a single point in \mathbb{R}^2 .

If, on the other hand, the determinant $a_1b_2 - a_2b_1 = 0$, then the two equations are linearly dependent and either:

- 1) the two lines are parallel and no solution exists (the numerator $\neq 0$),
or

- 2) the lines are parallel and superimposed, in which case there are an infinite number of solutions (**numerator = 0**).

As in the earlier linear problems, the condition of a nonzero determinant is the important point to note.

Now consider a second problem from linear algebra, that of two planes.

$$\left. \begin{array}{l} a_1x + b_1y + c_1z = d_1 \\ a_2x + b_2y + c_2z = d_2 \end{array} \right\} \quad \begin{array}{l} \text{two equations in} \\ \text{three unknowns } (x, y, z). \end{array}$$

Again, the objective is to determine whether the planes intersect or not, *i.e.* if there exists a solution to the pair of equations. If so, then the solution would correspond to a line of intersection since there are not enough equations to solve for all three variables (unknowns), but we can solve for two of them in terms of the third:

$$\begin{aligned} \left. \begin{array}{l} a_1x + b_1y = d_1 - c_1z \\ a_2x + b_2y = d_2 - c_2z \end{array} \right\} &\Rightarrow \begin{pmatrix} a_1 & b_1 \\ a_2 & b_2 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} d_1 - c_1z \\ d_2 - c_2z \end{pmatrix} \\ &\Rightarrow \begin{pmatrix} x \\ y \end{pmatrix} = \frac{1}{a_1b_2 - a_2b_1} \begin{pmatrix} b_2 & -b_1 \\ -a_2 & a_1 \end{pmatrix} \begin{pmatrix} d_1 - c_1z \\ d_2 - c_2z \end{pmatrix}. \end{aligned}$$

Again, this system has a unique solution

$$\iff \text{the determinant} = a_1b_2 - a_2b_1 \neq 0.$$

$$\begin{aligned} \text{For example, } x &= \frac{b_2(d_1 - c_1z) - b_1(d_2 - c_2z)}{a_1b_2 - a_2b_1}, \quad a_1b_2 - a_2b_1 \neq 0, \\ &\Rightarrow x = \frac{b_2d_1 - b_1d_2}{a_1b_2 - a_2b_1} - \frac{b_2c_1 - b_1c_2}{a_1b_2 - a_2b_1} \cdot z, \\ &\quad \left. \begin{array}{l} x = k_1 - m_1z \\ y = k_2 - m_2z \end{array} \right\} \quad \frac{x - k_1}{-m_1} = \frac{y - k_2}{-m_2} = z. \end{aligned}$$

As expected this is the equation of a line in 3D (as long as $0 < |m_1|, |m_2| < \infty$).

$$\text{The two equations for two planes: } \left. \begin{array}{l} a_1x + b_1y + c_1z = d_1 \\ a_2x + b_2y + c_2z = d_2 \end{array} \right\}$$

define a *line of intersection* if the determinant (of the coefficients of the variables we wish to solve for, x and y) is nonzero, *i.e.* $a_1b_2 - a_2b_1 \neq 0$. If the determinant is identically zero, then either:

- 1) the two planes are parallel and there is *no* solution, or

- 2) the two planes are parallel and superimposed, in which case there are an infinity of solutions.

We could go on to higher dimensional linear manifolds and consider systems of equations of many more variables, but the situations we would encounter would be the same:

- *The determinant of coefficients of the variables to be solved for cannot be zero if we want a unique solution.*
- *If the determinant is zero, then we have either no solution or an infinity of solutions.*

Now that we have reviewed these linear systems we are now ready to tackle the nonlinear problem.

Let $\mathbf{a} = (a, b, c)$ be a point on *both* surfaces. That is,

$$F(a, b, c) = G(a, b, c) = 0.$$

Just as in the linear problem on Page 108 these two equations define the set of points $\mathbf{x} = (x, y, z)$ which the two surfaces have in common. In other words,

this set is a curve of intersection.

Now we ask, when do these equations define a curve that can be expressed in the form $\begin{cases} x = f(z), \\ y = g(z), \\ z = z, \end{cases}$ with z as an independent variable?

The answer again depends on the existence of the linear approximations to $F(x, y, z)$ and $G(x, y, z)$ discussed in Section 2.D.

Suppose F and G are differentiable at $\mathbf{a} = (a, b, c)$; that is, $F, G \in C^1(\mathbb{R}^3)$. Then for points \mathbf{x} near \mathbf{a} on the curve of intersection, we find that

$$\begin{cases} F(x, y, z) = \cancel{F(\mathbf{a}, b, c)}^0 + \frac{\partial F}{\partial x} \Big|_{\mathbf{a}} \cdot (x - a) + \frac{\partial F}{\partial y} \Big|_{\mathbf{a}} \cdot (y - b) + \frac{\partial F}{\partial z} \Big|_{\mathbf{a}} \cdot (z - c) + \dots = 0 \\ G(x, y, z) = \cancel{G(\mathbf{a}, b, c)}^0 + \frac{\partial G}{\partial x} \Big|_{\mathbf{a}} \cdot (x - a) + \frac{\partial G}{\partial y} \Big|_{\mathbf{a}} \cdot (y - b) + \frac{\partial G}{\partial z} \Big|_{\mathbf{a}} \cdot (z - c) + \dots = 0. \end{cases}$$

Dropping the “+ . . .”, these equations can be approximated by

$$\begin{aligned}\left. \frac{\partial F}{\partial x} \right|_a \cdot (x - a) + \left. \frac{\partial F}{\partial y} \right|_a \cdot (y - b) &= -\left. \frac{\partial F}{\partial z} \right|_a \cdot (z - c) \\ \left. \frac{\partial G}{\partial x} \right|_a \cdot (x - a) + \left. \frac{\partial G}{\partial y} \right|_a \cdot (y - b) &= -\left. \frac{\partial G}{\partial z} \right|_a \cdot (z - c).\end{aligned}$$

These linear approximations form a matrix equation for $(x - a)$ and $(y - b)$:

$$\begin{pmatrix} \left. \frac{\partial F}{\partial x} \right|_a & \left. \frac{\partial F}{\partial y} \right|_a \\ \left. \frac{\partial G}{\partial x} \right|_a & \left. \frac{\partial G}{\partial y} \right|_a \end{pmatrix} \begin{pmatrix} (x - a) \\ (y - b) \end{pmatrix} = \begin{pmatrix} -\left. \frac{\partial F}{\partial z} \right|_a \cdot (z - c) \\ -\left. \frac{\partial G}{\partial z} \right|_a \cdot (z - c) \end{pmatrix}$$

This 2×2 system can be solved for $(x - a)$ and $(y - b)$ if and only if the determinant of the coefficient matrix, $\left. \frac{\partial F}{\partial x} \right|_a \cdot \left. \frac{\partial G}{\partial y} \right|_a - \left. \frac{\partial F}{\partial y} \right|_a \cdot \left. \frac{\partial G}{\partial x} \right|_a$, is not identically zero. This is analogous to our second linear algebraic system discussed earlier.

Incidentally, it bears noting that this determinant of derivatives appears in many related but also unrelated contexts (see Chapters 4 and 5 for more details). We take the opportunity here to assign to it a special notation. With it the results to follow are more concisely presented.

Definition 2.9

The determinant $\left. \frac{\partial F}{\partial x} \right|_a \cdot \left. \frac{\partial G}{\partial y} \right|_a - \left. \frac{\partial F}{\partial y} \right|_a \cdot \left. \frac{\partial G}{\partial x} \right|_a \equiv \left. \frac{\partial(F, G)}{\partial(x, y)} \right|_a$ is called a **Jacobian determinant**.

Inverting the coefficient matrix and using the Jacobian definition we have

$$(x - a) = -\frac{\left. \frac{\partial(F, G)}{\partial(z, y)} \right|_a}{\left. \frac{\partial(F, G)}{\partial(x, y)} \right|_a} \cdot (z - c), \quad (y - b) = -\frac{\left. \frac{\partial(F, G)}{\partial(x, z)} \right|_a}{\left. \frac{\partial(F, G)}{\partial(x, y)} \right|_a} \cdot (z - c),$$

for (x, y, z) very near (a, b, c) and provided $\left. \frac{\partial(F, G)}{\partial(x, y)} \right|_a \neq 0$.

Compare these expressions with their linear counterparts in our precursor problem on Page 108.

Therefore, provided $\frac{\partial(F, G)}{\partial(x, y)} \Big|_a \neq 0$, the set of equations

$$\begin{cases} x - a = \textcolor{teal}{m}_1(z - c) \\ y - b = \textcolor{teal}{m}_2(z - c) \\ z - c = z - c \end{cases}, \text{ with } \textcolor{teal}{m}_1 = -\frac{\frac{\partial(F, G)}{\partial(z, y)} \Big|_a}{\frac{\partial(F, G)}{\partial(x, y)} \Big|_a}, \quad \textcolor{red}{m}_2 = -\frac{\frac{\partial(F, G)}{\partial(x, z)} \Big|_a}{\frac{\partial(F, G)}{\partial(x, y)} \Big|_a}$$

is the *tangent line approximation* to the implied curve of intersection

$$x = f(z), \quad y = g(z), \quad z = z, \quad \text{near } (a, b, c). \quad (2.5)$$

We conclude that as long as m_1 and m_2 are finite the above curve of intersection is well defined. Moreover, since the tangent line to the curve given by Equation (2.5) at the point (a, b, c) ,

$$\begin{cases} x - a = \frac{df}{dz} \Big|_c \cdot (z - c) \\ y - b = \frac{dg}{dz} \Big|_c \cdot (z - c) \\ z - c = z - c, \end{cases}$$

is unique we can also deduce the following results:

$$\frac{df}{dz} \Big|_c = -\frac{\frac{\partial(F, G)}{\partial(z, y)} \Big|_a}{\frac{\partial(F, G)}{\partial(x, y)} \Big|_a}, \quad \frac{dg}{dz} \Big|_c = -\frac{\frac{\partial(F, G)}{\partial(x, z)} \Big|_a}{\frac{\partial(F, G)}{\partial(x, y)} \Big|_a}.$$

Once again, the conditions for these to be valid and for the implicit functions to exist are

- (i) F and G differentiable at (a, b, c) ;
- (ii) the Jacobian determinant $J = \frac{\partial(F, G)}{\partial(x, y)} \neq 0$ at (a, b, c) .

The important point to note in this example is that critical condition (ii) involves the matrix of coefficients of the dependent variables, x and y . This is consistent with the preceding examples and is a general rule of thumb with establishing the existence of any implicit functions!

 **Mastery Check 2.22:**

Suppose $x = h(u, v) = u^2 + v^2$ and $y = k(u, v) = uv$ are to be “solved” for u and v in terms of x and y . Find $\frac{\partial u}{\partial x}$, $\frac{\partial u}{\partial y}$, $\frac{\partial v}{\partial x}$, $\frac{\partial v}{\partial y}$, where possible. Show

$$\frac{\partial(u, v)}{\partial(x, y)} = \left(\frac{\partial(x, y)}{\partial(u, v)} \right)^{-1} \text{ provided the denominator } \neq 0. \text{ Hint: This problem}$$

is similar to Problem 3) above, but this time with four variables. It proves to be worthwhile here to be a little liberal with our notation convention, and refer to u and v as functions as well as variables.

- 1) Define suitable level sets $\begin{cases} F(x, y, u, v) = 0, \\ G(x, y, u, v) = 0, \end{cases}$ which we wish to solve for some functions $u = f(x, y)$, $v = g(x, y)$, the curves of intersection.
- 2) Set up the linear approximations to the level sets for (x, y, u, v) near (a, b, c, d) . Rewrite the approximations as a matrix equation in $\begin{pmatrix} u - c \\ v - d \end{pmatrix}$.
- 3) Solve the equation. (What appears as a denominator?)
- 4) Compare your solution to the true tangent lines and so obtain $\frac{\partial f}{\partial x}$, $\frac{\partial f}{\partial y}$, $\frac{\partial g}{\partial x}$, $\frac{\partial g}{\partial y}$, in terms of Jacobians.
- 5) Finally, compare $\frac{\partial(f, g)}{\partial(x, y)}$ with $\frac{\partial(h, k)}{\partial(u, v)}$.



In this last Mastery Check, we met a *fundamental property* of the Jacobian

$$(\text{where it exists}), \text{ namely, } \frac{\partial(u, v)}{\partial(x, y)} = \left(\frac{\partial(x, y)}{\partial(u, v)} \right)^{-1}.$$

In the single-variable case, it is true that $\frac{dy}{dx} = 1 / \frac{dx}{dy}$ if $\frac{dx}{dy} \neq 0$.

However in the context of multivariable functions and partial derivatives it is the *Jacobian* which takes the place of the ordinary derivative.

In the notation used in the Mastery Check,

$$\frac{\partial u}{\partial x} \left(= \frac{\partial f}{\partial x} \right) \neq 1 / \frac{\partial x}{\partial u} \left(= 1 / \frac{\partial h}{\partial u} \right).$$

 **Mastery Check 2.23:**

Let $\begin{cases} F(x, y, u, v) = xyuv - 1 = 0 \\ G(x, y, u, v) = x + y + u + v = 0, \end{cases}$
and consider points $P_0 = (1, 1, -1, -1)$, $P_1 = (1, -1, 1, -1)$.

Find $\left(\frac{\partial y}{\partial x}\right)_u$ at P_0 and P_1 .

Hint: Which are the independent variables?



 **Mastery Check 2.24:**

Show that the system of equations $\begin{cases} xy^2 + xzu + yv^2 = 3 \\ x^3yz + 2xv - u^2v^2 = 2 \end{cases}$ can be solved for u, v as functions of x, y, z near the point $P_0(1, 1, 1, 1, 1)$.

Find the value of $\frac{\partial v}{\partial y}$ for the solution at $(x, y, z) = (1, 1, 1)$.



2.I Taylor's formula and Taylor series

On the one hand, this next topic is but a natural extension of our earlier discussion on tangent plane approximations (Section 2.D). On the other hand, the subject of Taylor series and Taylor approximations is so incredibly useful in analysis and incredibly practical in computational applications that it is worth giving some consideration. The next chapter will highlight some examples of such applications. For the present purpose we consider this topic as a means of getting good approximations to functions, whether explicit or implicit. A convenient place to start is with the single-variable case.

Recall the properties of Taylor and Maclaurin polynomials for any function $F : \mathbb{R} \rightarrow \mathbb{R}$: Let F and $F^{(k)}$, $k = 1, 2, \dots, n$ be continuous on an open interval I including the point t_0 , and let $F^{(n+1)}(t)$ exist for all $t \in I$. The best polynomial approximation of order n to f near $t_0 \in I$ is the first contribution in the next equation:

$$\begin{aligned} F(t) &= P_n(t) + E_n(t; t_0) \\ &= F(t_0) + F'(t_0)(t - t_0) + \frac{F''(t_0)}{2}(t - t_0)^2 + \dots \\ &\quad + \frac{F^{(n)}(t_0)}{n!}(t - t_0)^n + E_n(t; t_0). \end{aligned}$$

$P_n(t)$ is referred to as the Taylor polynomial approximation to F of order n , while

$$\begin{aligned} E_n(t; t_0) &= F(t) - P_n(t) \\ &= \frac{F^{(n+1)}(a)}{(n+1)!}(t - t_0)^{n+1}, \quad a \in (\min(t_0, t), \max(t_0, t)) \end{aligned}$$

is the error term (the difference between the true value and its approximation).

Here are a few special cases for $n < \infty$:

- (i) $P_1(t) = F(t_0) + F'(t_0)(t - t_0)$ — linear approximation
- (ii) $P_2(t) = P_1(t) + \frac{F''(t_0)}{2}(t - t_0)^2$ — quadratic approximation
- (iii) $P_3(t) = P_2(t) + \frac{F'''(t_0)}{3!}(t - t_0)^3$ — cubic approximation
- (iv) $P_4(t) = P_3(t) + \frac{F^{(4)}(t_0)}{4!}(t - t_0)^4$ — quartic approximation

Specific cases to note

- * The existence of a linear approximation means there is a tangent line to $F(t)$ at $t = t_0$.
- * A quadratic approximation is useful for critical-point analysis when $F'(t_0) = 0$, meaning that

$$F(t) - F(t_0) \approx \frac{F''(t_0)}{2}(t - t_0)^2 \quad \begin{array}{l} > 0, \\ < 0, \end{array} \quad \begin{array}{l} \text{— a minimum point} \\ \text{— a maximum point.} \end{array}$$

- * A cubic approximation means that there is a cubic curve osculating $F(t)$ at $x = t_0$. (What is that?)

- * A quartic approximation may be useful in the uncommon cases when $F'(t_0) = F''(t_0) = F'''(t_0) = 0$, meaning that

$$F(t) - F(t_0) \approx \frac{F^{(4)}(t_0)}{4!}(t - t_0)^4 \quad \begin{array}{l} > 0, \\ < 0, \end{array} \quad \begin{array}{l} \text{— a minimum point} \\ \text{— a maximum point.} \end{array}$$

Some functions have derivatives of all orders. Therefore, we can consider extending n without limit, that is, $n \rightarrow \infty$. The above polynomial can then be developed to an infinite power series, provided it converges (absolutely) in some interval.

Functions that can be differentiated an indefinite number of times and whose interval of convergence is the whole real line, for example, $\sin t$, $\cos t$, e^t , are called *analytic*.

The Taylor series representation of a function $F(t)$ is defined as follows:

Definition 2.10

If there exists a $t_0 \in \mathbb{R}$ and an $R > 0$ such that $F(t) = \sum_{k=0}^{\infty} a_k(t - t_0)^k$ converges for $|t - t_0| < R$, then this is the Taylor series of F , and $a_k = \frac{F^{(k)}(t_0)}{k!} \forall k$.

What is especially important for us is the particular choice $t_0 = 0$. For this choice we have the well-known *Maclaurin polynomial*.

To be precise, if a single-variable function $F : \mathbb{R} \rightarrow \mathbb{R}$ has continuous derivatives of all orders less than or equal to $n + 1$, in an open interval I centred at $t = 0$, then F can be approximated by the Maclaurin polynomial,

$$F(t) = F(0) + F'(0)t + \frac{1}{2!}F''(0)t^2 + \cdots + \frac{1}{n!}F^{(n)}(0)t^n + R_n(\theta, t) \quad (2.6)$$

for all t in the interval I , and where

$$R_n(\theta, t) = \frac{1}{(n+1)!}F^{(n+1)}(\theta)t^{n+1}$$

is the measure of error in the approximation, with $0 < |\theta| < t$.

As with the Taylor series, if F has derivatives of all orders, that is, for $n \rightarrow \infty$, at $t = 0$, then we can define the Maclaurin series representation of F :

$$F(t) = \sum_{k=0}^{\infty} a_k t^k \text{ for } t \in I, \text{ with } 0 \in I, \text{ and } a_k = \frac{F^{(k)}(0)}{k!} \forall k.$$

For example, it is easy to verify by repeated differentiation and substitution

that the Maclaurin series for the sine function is

$$\sin t = \sum_{k=0}^{\infty} \frac{(-1)^{k+1}}{(2k+1)!} t^{2k+1}.$$

From this single-variable case we can derive corresponding versions of Taylor and Maclaurin polynomials for functions of several variables.

Let's consider a special function F and look at its value at $t = 1$:

Let $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ be defined and have continuous partial derivatives of orders $0, 1, \dots, n+1$, at the point (x_0, y_0) in the domain of f .

For fixed $(x, y) \in S_r(x_0, y_0)$ and therefore fixed $h = x - x_0$ and $k = y - y_0$, consider $f(x, y)$ at $x = x_0 + th, y = y_0 + tk$. We can therefore define a single-valued function of t

$$F(t) = f(x_0 + th, y_0 + tk)$$

whose value at $t = 0$ is $F(0) = f(x_0, y_0)$, and whose value at $t = 1$ is $F(1) = f(x_0 + h, y_0 + k)$.

We now develop the Maclaurin polynomial for F *via the chain rule*, which leads us to the Taylor polynomial for f .

The terms in the Maclaurin polynomial are

$$F(0) = f(x_0, y_0), \quad F'(0)t = \left(\frac{\partial f}{\partial x} \Big|_{(x_0, y_0)} h + \frac{\partial f}{\partial y} \Big|_{(x_0, y_0)} k \right) t, \\ \frac{1}{2} F''(0)t^2 = \frac{1}{2} \left(\frac{\partial^2 f}{\partial x^2} \Big|_{(x_0, y_0)} h^2 + 2 \frac{\partial^2 f}{\partial x \partial y} \Big|_{(x_0, y_0)} hk + \frac{\partial^2 f}{\partial y^2} \Big|_{(x_0, y_0)} k^2 \right) t^2 \dots$$

Inserting these in (2.6) and letting $t = 1$ give us the Taylor polynomial approximation to f :

$$f(x_0 + h, y_0 + k) = f(x_0, y_0) + \frac{\partial f}{\partial x} \Big|_0 h + \frac{\partial f}{\partial y} \Big|_0 k \\ + \frac{1}{2} \left(\frac{\partial^2 f}{\partial x^2} \Big|_0 h^2 + 2 \frac{\partial^2 f}{\partial x \partial y} \Big|_0 hk + \frac{\partial^2 f}{\partial y^2} \Big|_0 k^2 \right) + \dots \\ + \sum_{j=0}^n \frac{1}{j!(n-j)!} \frac{\partial^n f}{\partial x^j \partial y^{n-j}} \Big|_0 h^j k^{n-j} + R_n(\theta, h, k). \quad (2.7)$$

Alternative derivation:

Suppose we can approximate $f(x, y)$ in a neighbourhood $S_r(x_0, y_0)$ of (x_0, y_0)

by a general polynomial:

$$\begin{aligned} f(x, y) = & a_{00} + a_{01}(x - x_0) + a_{01}(y - y_0) \\ & + a_{20}(x - x_0)^2 + a_{11}(x - x_0)(y - y_0) + a_{02}(y - y_0)^2 + \dots \\ & + a_{n0}(x - x_0)^n + \dots + a_{0n}(y - y_0)^n + E_n(\mathbf{x}_0). \end{aligned}$$

(Approximation error will depend on \mathbf{x}_0 and n .)

If f is differentiable to order n , then taking partial derivatives of both sides and evaluating these results at (x_0, y_0) we identify

$$\begin{aligned} \frac{\partial f}{\partial x} \Big|_0 &= a_{10} & \frac{\partial f}{\partial y} \Big|_0 &= a_{01} \\ \frac{\partial^2 f}{\partial x^2} \Big|_0 &= 2a_{20} & \frac{\partial^2 f}{\partial y^2} \Big|_0 &= 2a_{02} & \frac{\partial^2 f}{\partial x \partial y} \Big|_0 &= a_{11} = \frac{\partial^2 f}{\partial y \partial x} \Big|_0 \\ \text{and generally } & \frac{\partial^{k+\ell} f}{\partial x^k \partial y^\ell} \Big|_0 &= k! \ell! a_{k\ell}. \end{aligned}$$

Substitution will give (2.7) again.

We end this section with two Mastery Check exercises involving Taylor polynomials, postponing until the next chapter a demonstration of the usefulness of Taylor approximations. However, it is appropriate first to make a few important comments.

- (1) With the trivial step of setting $x_0 = y_0 = 0$ we get the 2D version of a Maclaurin polynomial approximation of order n .
- (2) As in the single-variable case, if our multivariable function has partial derivatives of all orders, then the Taylor polynomial approximations of our $f(x, y)$ can be developed into a series representation — provided it converges, of course.
- (3) In Section 2.H we dealt with implicit functions and showed how one can calculate first derivatives of such functions even though the functions themselves could not be expressed explicitly. Now, with assistance from the Taylor and Maclaurin expansions, one can construct explicit polynomial or even full series representations of implicit functions by successively differentiating these first-order derivatives (with the help of the chain rule — Section 2.G). Our second Mastery Check exercise explores this possibility.

 **Mastery Check 2.25:**

Determine Taylor's polynomial of order 2 to the function

$$f(x, y) = \ln(2x^2 + y^2) - 2y$$

about the point $(0, 1)$, and evaluate an approximation to $f(0.1, 1.2)$.

 **Mastery Check 2.26:**

Determine the Taylor polynomial of order 2 about the point $(0, 0)$ to the implicit function $z = f(x, y)$ defined by the equation

$$e^{x+y+z} - (x + y + z)^2 - 1 = 0.$$

Hint: see Section 2.H Problem 2.



2.J Supplementary problems

Section 2.A

1. Consider the function $f : \mathbb{R} \rightarrow \mathbb{R}$, $x \mapsto f(x)$, where
 - $f(x) = x^{2/3}$;
 - $f(x) = x^{4/3}$;
 - $f(x) = x^{1/2}$.
 Decide in each case whether the derivative of f exists at $x = 0$ as follows:
 - (i) Check whether f exists at $x=0$ and at $x = 0+h$ (for h arbitrarily small).
 - (ii) Check right and left limits to see whether f is continuous at $x = 0$.
 - (iii) Check right and left limits to see whether f has a derivative at $x = 0$.
2. Consider the function $f : \mathbb{R} \rightarrow \mathbb{R}$, $x \mapsto (x - 1)^\alpha$, $\alpha \in \mathbb{R}$. Decide from first principles for what values of α the derivative of f exists at $x = 1$.

Section 2.B

3. Are the following functions continuous at $(0, 0)$?
 - (a) $f(x, y) = \frac{xy}{x^2 + y^2}$, $f(0, 0) = 0$.
 - (b) $f(x, y) = \frac{xy}{\sqrt{x^2 + y^2}}$, $f(0, 0) = 0$.
4. Establish whether the limits of the following functions exist at the given points, and if so determine their values.
 - (a) $f(x, y) = \frac{\tan(xy)}{x^2 + y^2}$, at $(x, y) = (0, 0)$.
 - (b) $f(x, y) = \frac{x^2 - x}{y^2 - y}$, at $(x, y) = (0, 0)$.
 - (c) $f(x, y) = \frac{x^4 + y \sin(x^3)}{x^4 + y^4 + x^2 y^2}$, at $(x, y) = (0, 0)$.
 - (d) $f(x, y) = \frac{x^2 \sin y}{x^2 + y^2}$, at $(x, y) = (0, 0)$.
5. Use the limit definition to show that

- (a) $\lim_{(x,y) \rightarrow (1,1)} \left(x + \frac{1}{y} \right) = 2.$
- (b) $\lim_{(x,y) \rightarrow (1,2)} (x^2 + 2y) = 5.$
- (c) $\lim_{(x,y) \rightarrow (0,0)} \left(\frac{\cos(2xy)}{1 + x^2 + y^2} \right) = 1.$

Repeat the limit calculation using the limit laws, where these are applicable.

Section 2.C

6. Find all first partial derivatives of the following functions:

- (a) $f(x, y) = \arctan(y/x).$
- (b) $f(x, y) = \arctan(x^2 + xy + y^2).$
- (c) $f(x, y) = (x^2 y^3 + 1)^3.$
- (d) $f(x, y) = \exp(x^2 \sin^2 y + 2xy \sin x \sin y + y^2).$
- (e) $f(x, y, z) = x^2 \sqrt{y^2 + xz}.$
- (f) $f(x_1, x_2, x_3) = \ln |x_1 x_2 + x_2 x_3 + x_1 x_3|.$

Section 2.D

7. Consider the function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$, $f(x, y) = x \arctan(y/x)$,
 $f(0, y) = 0$.
 Discuss (a) the continuity of f , and (b) the continuity and existence of f_x and f_y , at points on the y -axis, without first drawing the graph.
 Discuss (c) the existence of the second partial derivatives at these points, and determine them if they do.

Section 2.E

8. Consider the function $f(x, y, z) = xy^3 + yz^2$. What is the directional derivative at the point $(1, 2, -1)$ in the direction $\mathbf{u} = 2\mathbf{e}_1 + \mathbf{e}_2 + 2\mathbf{e}_3$?
9. Consider the surface $xy^3z^2 = 8$. What is the unit normal to this surface at the point $(1, 2, -1)$?

10. Determine the points on the flattened ellipsoid

$$x^2 + 13y^2 + 7z^2 - 6\sqrt{3}zy = 16$$

where the tangent plane is parallel to one of the coordinate planes.

11. Determine all points on the surface

$$x^2 + 2y^2 + 3z^2 + 2xy + 2yz = 1$$

where the tangent plane is parallel to the plane $x + y + z = 0$.

Section 2.F

12. Find the first and second partial derivatives of the following functions:

(a) $f(x, y) = (\sin xy)^2$.

(b) $f(x, y) = \ln(x^2 + 2y)$.

(c) $f(x, y) = \sqrt{1 + x^2 + y^3}$.

13. Suppose $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ is a C^3 harmonic function of variables (x, y) . That is, suppose f satisfies the 2D Laplace equation,

$$\Delta f = \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2} = 0$$

(see Section 3.E for more information on this topic). Show that the function $g = x \frac{\partial f}{\partial x} + y \frac{\partial f}{\partial y}$ is also a solution.

Section 2.G

14. Consider the functions $f : x \mapsto y = f(x)$, and $g : t \mapsto x = g(t)$.

In each example that follows, find the domain D_F of the composite function $F : t \mapsto y = (f \circ g)(t)$, and the derivative $\frac{dF}{dt}$.

(a) $f(x) = \frac{x^2}{1 + x^2}$, $g(t) = \sinh t$.

(Do this example in two ways: by finding the expression for $F(t)$ explicitly in terms of t , then differentiating; and by using the chain rule.)

(b) $f(x) = \arcsin x^2$, $g(t) = 3e^{-t}$.

15. Consider the functions $f : x \mapsto y = f(x)$, and $g : (s, t) \mapsto x = g(s, t)$. In each example that follows, find the domain D_F of the composite function $F : (s, t) \mapsto y = (f \circ g)(s, t)$, and the derivatives $\frac{\partial F}{\partial s}$, $\frac{\partial F}{\partial t}$.
- (a) $f(x) = \ln x$, $g(s, t) = s(1 - t^2)$.
- (b) $f(x) = \arccos(x)$, $g(s, t) = \sqrt{s^2 - t^2}$.
16. Consider the function $z = f(x, y) = e^{x^2 y} + xy$, and the composite function $z = F(t) = f(\cos t, \sin t)$. Decide whether $F(t)$ makes sense, and if so, find its domain and compute $\frac{dF}{dt}$.
17. Consider the function $z = f(x, y) = \arcsin xy$, where $x = s - 2t$ and $y = \frac{s}{t^2}$. Check that the composite function $z = F(s, t)$ makes sense, and if so, find its domain and compute $\frac{\partial F}{\partial s}$ and $\frac{\partial F}{\partial t}$.
18. By introducing new variables $u = x^2 - y$ and $v = x + y^2$ transform the differential equation
- $$(1 - 2y) \frac{\partial f}{\partial x} + (1 + 2x) \frac{\partial f}{\partial y} = 0.$$
19. Suppose $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ is a C^2 function of variables (x, y) . By introducing the change of variables $x = 2s + 3t$ and $y = 4s - 4t$ we define the C^2 function $F(s, t)$ from $f(x, y)$. Show that
- $$\frac{\partial^2 F}{\partial t^2} = 9 \frac{\partial^2 f}{\partial x^2} - 24 \frac{\partial^2 f}{\partial y \partial x} + 16 \frac{\partial^2 f}{\partial y^2}.$$
20. Suppose $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ is a C^2 harmonic function of variables (x, y) . By introducing 2D polar coordinates $x = r \cos \theta$ and $y = r \sin \theta$ show that Laplace's equation becomes
- $$\Delta F = \frac{\partial^2 F}{\partial r^2} + \frac{1}{r} \frac{\partial F}{\partial r} + \frac{1}{r^2} \frac{\partial^2 F}{\partial \theta^2} = 0$$
- where $F(r, \theta) = f(r \cos \theta, r \sin \theta)$.
21. Suppose $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ is a C^2 function of variables (x, y) , by introducing the new variables $s = x^2 + y$ and $t = 2x$ transform the expression
- $$\frac{\partial^2 f}{\partial x^2} + 2 \frac{\partial^2 f}{\partial y \partial x} + \frac{\partial^2 f}{\partial y^2}$$

into a form involving a function $F(s, t)$.

Section 2.H

22. Suppose the equation $x^3y + 2y^3x = 3$ defines y as a function f of x in the neighbourhood of the point $(1, 1)$. Find the derivative of f at $x = 1$.
23. Suppose the equation $z^3 + z(y^2 + 1) + x^3 - 3x + y^2 - 8 = 0$ defines z as a function f of x, y in the neighbourhood of the point $(2, -1, -1)$. Find the derivatives of $\frac{\partial f}{\partial x}$ and $\frac{\partial f}{\partial y}$ at this point.

Section 2.I

24. Find the Taylor polynomial approximations of order 1, 2, ..., to the function $F(t) = \cos(t^{3/2})$ about the point $t = t_0 = 0$. Use a suitable low-order approximation to determine whether $t = 0$ is a maximum or minimum point.
25. If a function may be approximated by a convergent Taylor series whose terms are alternating in sign, then the absolute value of the error term is bounded by the absolute value of the first omitted term in a finite polynomial approximation. Use this idea to decide how many terms are needed to find $\sin 3^\circ$ to ten decimal places.
26. Write down the Taylor series for $\cos t$ about $t_0 = \pi/3$, and use it to compute $\cos 55^\circ$ to seven decimal places. (First establish how many terms are needed.)
27. Write down all the terms up to order $n = 3$ in the Taylor series for $f(x, y) = \sin(x + y) + x^2$ about $(x_0, y_0) = (0, \pi/6)$.
28. The function $z = f(x, y)$ satisfies $f(1, 1) = 0$ and $e^{xz} + x^2y^3 + z = 2$. Establish that $z = f(x, y)$ indeed exists in a region near $(1, 1)$ and determine its Taylor series approximation up to order 2 valid near that point.
29. Determine the Taylor polynomial approximation up and including second-order terms of $F(x, y) = \int_0^y \ln(1 + x + t)dt$ about the point $(0, 0)$. Note that we assume $x + y > -1$.



Chapter 3

Applications of the differential calculus

In this chapter, the theory of multivariable functions and their partial derivatives as covered in the preceding chapter is applied to problems arising in four contexts: finding function maxima and minima, error analysis, least-square approximations, and partial differential equations. Although applications arise in a wide variety of forms, these are among the more common examples.

3.A Extreme values of $f : \mathbb{R}^n \rightarrow \mathbb{R}$

Occupying a central position in the vastness of the space of applications of differential calculus is the subject of *optimization*. At its most basic, the term refers to the task of finding those points in a function's domain that gives rise to maxima or minima of that (scalar) function, and of determining the corresponding values of that function.

Of special interest in the study are the so-called *extreme points* of $f(\mathbf{x})$, a subset of which are the so-called *critical points*. These are points where the function can exhibit either a local maximum or minimum, and even a global maximum or minimum.

To set the stage we require some basic infrastructure. We start with a few essential definitions.

Definition 3.1

Consider a continuous $f : \mathbb{R}^n \rightarrow \mathbb{R}$.

A point $\mathbf{a} \in D_f$ is called a local

- (i) **minimum point** if $f(\mathbf{x}) \geq f(\mathbf{a})$, $\forall \mathbf{x} \in S_r(\mathbf{a})$,
- (ii) **maximum point** if $f(\mathbf{x}) \leq f(\mathbf{a})$, $\forall \mathbf{x} \in S_r(\mathbf{a})$.

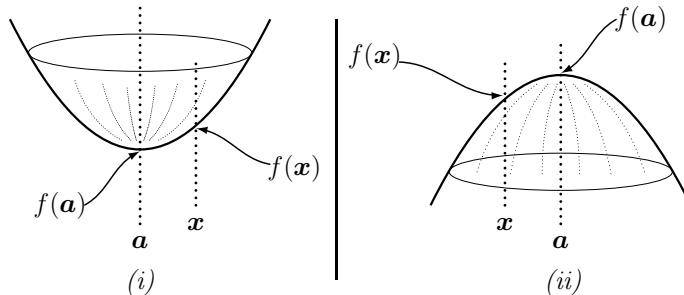


Figure 3.1 A function minimum and maximum.

We have here invoked the open sphere $S_r(\mathbf{a})$ to represent the set of points \mathbf{x} different from but near \mathbf{a} (the radius $r > 0$ is presumed small). We could equally well have referred to points \mathbf{x} in a larger “neighbourhood” of \mathbf{a} . However, that proves to be unnecessary and less convenient, it is enough to consider a small open sphere as we are defining *local* properties.

Points of local minimum (Figure 3.1(i)) and local maximum (Figure 3.1(ii)) are examples of *critical points*.

Definition 3.2

A **critical point** is an interior point $\mathbf{a} \in D_f$ at which $\nabla f|_{\mathbf{a}} = \mathbf{0}$ for $f \in C^1$.

While a local maximum point or a local minimum point must mean that $\nabla f = \mathbf{0}$ at that point the converse is not necessarily true, a critical point where $\nabla f = \mathbf{0}$ need not be either a point of maximum or minimum; there is a third alternative.

Definition 3.3

A critical point which is neither a maximum nor a minimum is called a **saddle point**.

Referring to Definition 3.1, in the case of a local minimum (left-hand figure) or a local maximum (right-hand figure) at $\mathbf{x} = \mathbf{a}$, the tangent plane is horizontal, which is a consequence of $\nabla f = \mathbf{0}$ at \mathbf{a} .

For a saddle point the tangent plane is still horizontal but neither of the figures in Definition 3.1 applies. Instead, around a saddle point a part of the function's graph is below the tangent plane and a part is above.

The following two simple examples convey the general idea of the above definitions.

Example 3.1:

Consider the function $z = f(x, y) = x^2 + y^2 - 2x$. We have

$$\nabla f = \begin{pmatrix} 2x - 2 \\ 2y \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \text{ at } (x, y) = (1, 0).$$

Now we examine f in the neighbourhood of this critical point $(1, 0)$. (Note that there is just one critical point in this example.) Let's consider the neighbouring point $(1 + h, 0 + k)$ in the domain of f . We have

$$f(1 + h, 0 + k) = (1 + h)^2 + k^2 - 2(1 + h) = h^2 + k^2 - 1,$$

while $f(1, 0) = 1 + 0 - 2 = -1$. We see that $f(1 + h, 0 + k) > f(1, 0)$ for all $h, k \neq 0$, since

$$f(1 + h, 0 + k) - f(1, 0) = h^2 + k^2 > 0.$$

As this is true for all (h, k) , that is, all (x, y) in the neighbourhood of $(1, 0)$, the point $(1, 0)$ is a minimum point. ■

Example 3.2:

Consider the function $z = f(x, y) = 1 + x^2 - y^2$. (See Example 1.10.) We have

$$\nabla f = \begin{pmatrix} 2x \\ -2y \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \text{ at } (x, y) = (0, 0).$$

We examine f in the neighbourhood of this critical point. At the neighbouring point $(0 + h, 0 + k)$ we have

$$f(0 + h, 0 + k) = 1 + h^2 - k^2,$$

which is > 0 along the line $y = k = 0$, but is < 0 along the line $x = h = 0$. The critical point $(0, 0)$ is neither a local maximum nor a local minimum. It is a *saddle point*. ■

The reader should now try their hand at a similar style problem.

Mastery Check 3.1:

Consider the function $z = f(x, y) = x^2 - y^2 - 2x$. Find the point (a, b) at which $\nabla f = \mathbf{0}$, and then find an expression for $f(a + h, b + k)$ for small $h, k \neq 0$. Use Definition 3.1 to decide whether the point (a, b) is a maximum or a minimum (or neither). 

Extreme values in two dimensions—general procedure

The functions in Example 3.1 and the last Mastery Check were nice ones to work with. We could use simple algebra to determine if a point of interest was a maximum point or a minimum point or neither. The question naturally arises, what do we do with functions which are more complicated?

The answer relies on the fact that since we are interested only in *local* extreme points we can make use of *local approximations* to functions.

In fact, all we usually ever need is Taylor's polynomial of second order which, as the next theorem states, is enough to represent a function locally.

Theorem 3.1

Let $f : D_f \subset \mathbb{R}^2 \rightarrow \mathbb{R}$ be a function with continuous derivatives of order 0, 1, 2, and 3 (that is, f is a class C^3 function) in some neighbourhood of a point $\mathbf{a} \in D_f$. Then,

$$\begin{aligned} f(a + h, b + k) &= f(a, b) + \frac{\partial f}{\partial x}(a, b).h + \frac{\partial f}{\partial y}(a, b).k \\ &+ \frac{1}{2} \left(\frac{\partial^2 f}{\partial x^2}(a, b).h^2 + 2 \frac{\partial^2 f}{\partial x \partial y}(a, b).hk + \frac{\partial^2 f}{\partial y^2}(a, b).k^2 \right) + (h^2 + k^2)^{3/2} B(h, k), \end{aligned}$$

where B is some bounded function in the neighbourhood of $(0, 0)$.

Using Taylor polynomials of order 2 (see Equation 2.7) results in a considerable simplification. The difficulties of critical point problems involving more complex functions are reduced to the level featured in Examples 3.1 and 3.2,

since the function approximations are algebraic.

The Taylor polynomial approximation of order 2 can be written more succinctly in a vector-matrix product form

$$f(\mathbf{a} + \mathbf{h}) = f(\mathbf{a}) + \text{grad } f(\mathbf{a}) \cdot \mathbf{h} + \frac{1}{2} \mathbf{h}^T \cdot H_{\mathbf{a}} f \cdot \mathbf{h} + \text{small terms,} \quad (3.1)$$

where

$$H_{\mathbf{a}} f = \begin{pmatrix} \frac{\partial^2 f}{\partial x^2}(a, b) & \frac{\partial^2 f}{\partial x \partial y}(a, b) \\ \frac{\partial^2 f}{\partial y \partial x}(a, b) & \frac{\partial^2 f}{\partial y^2}(a, b) \end{pmatrix}$$

is a symmetric matrix for C^2 functions called *the Hessian matrix*, and

$$\mathbf{h} = \begin{pmatrix} h \\ k \end{pmatrix}, \quad \mathbf{h}^T = (h, k), \quad \text{grad } f(\mathbf{a}) = \left(\frac{\partial f}{\partial x}(a, b), \frac{\partial f}{\partial y}(a, b) \right).$$

The vector-matrix expression, Equation 3.1, is a very convenient form to use as it is straightforward to generalize to functions of n variables. For $n \geq 3$ only the sizes of the vectors and the Hessian matrix increases, while the form stays the same. Try it?

We shall now see how Taylor's second-order polynomial can help us to examine the behaviour of a function in the neighbourhood of a critical point. Bear in mind that at such points the function's gradient vanishes.

Considering points in a small region around the critical point \mathbf{a} of $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ we have from Equation (3.1) (ignoring the small terms),

$$\begin{aligned} f(a + h, b + k) &\approx f(a, b) + \frac{1}{2} \left(\frac{\partial^2 f}{\partial x^2}(a, b) \cdot h^2 + 2 \frac{\partial^2 f}{\partial x \partial y}(a, b) \cdot hk + \frac{\partial^2 f}{\partial y^2}(a, b) \cdot k^2 \right) \\ &= f(a, b) + \frac{1}{2} Q(h, k). \end{aligned} \quad \text{— since the gradient term = 0}$$

Here $Q(h, k)$ is called a *quadratic form*. For a general function of n variables which has continuous derivatives of order 2, we can write

$$Q(\mathbf{h}) = \mathbf{h}^T \cdot H_{\mathbf{a}} f \cdot \mathbf{h}.$$

For $|\mathbf{h}| \ll 1$ the sign of Q determines whether \mathbf{a} is a maximum, a minimum, or a saddle point.

Suppose $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ has continuous derivatives of order no greater than 3 (less than or equal to 3) and $\mathbf{a} \in D_f$ is a critical point of f .

- 1) If $Q(h, k)$ is *positive definite*, then f has a local *minimum* value at \mathbf{a} :

$$Q(h, k) > 0 \text{ for all } 0 \neq |\mathbf{h}| \ll 1 \implies f(a + h, b + k) > f(a, b).$$
- 2) If $Q(h, k)$ is *negative definite*, then f has a local *maximum* value at \mathbf{a} :

$$Q(h, k) < 0 \text{ for all } 0 \neq |\mathbf{h}| \ll 1 \implies f(a + h, b + k) < f(a, b).$$
- 3) If $Q(h, k)$ is *indefinite*, then f has neither a maximum nor a minimum value at \mathbf{a} , and \mathbf{a} is a *saddle point*: For all $0 \neq |\mathbf{h}| \ll 1$,

$$Q(h, k) < 0 \text{ for some } \mathbf{h} \implies f(a + h, b + k) > f(a, b),$$

$$Q(h, k) > 0 \text{ for other } \mathbf{h} \implies f(a + h, b + k) < f(a, b).$$
- 4) If $Q(h, k)$ is *positive or negative semi-definite*:

$$Q > 0 \text{ or } Q < 0 \text{ and } Q = 0 \text{ for some } |\mathbf{h}| \neq 0,$$
then we cannot say anything.

A summary of this section on critical points and critical point classification appears below, with an aside reviewing the corresponding facts in the case of functions of one variable. The comparison between the 1D and the n D cases is quite instructive. The reader should note the similarities at corresponding points of the arguments. Readers will have an opportunity to test their understanding by solving Mastery Checks 3.2–3.5.

In the **one-dimensional case** an investigation into critical points of a function of one variable is summarized as follows.

Consider $y = f(x)$. Critical points are determined from solutions of the zero-derivative equation,

$$\frac{df}{dx} \Big|_a = 0 \implies x = a$$

As a consequence, we find that

$$\begin{aligned} \frac{\partial^2 f}{\partial x^2} \Big|_a > 0 &\quad \Rightarrow \text{ minimum} \\ \frac{\partial^2 f}{\partial x^2} \Big|_a < 0 &\quad \Rightarrow \text{ maximum} \\ \frac{\partial^2 f}{\partial x^2} \Big|_a = 0 &\quad \Rightarrow \text{ a stationary point} \end{aligned}$$

In some interval

$$I_r = \{x : 0 < |x - a| < r\}$$

about the critical point, we have the approximation

$$f(x) \approx f(a) + \frac{df}{dx}\Big|_a (x - a) + \frac{1}{2} \frac{d^2 f}{dx^2}\Big|_a (x - a)^2$$

Hence, for points $x = a + h$ near $x = a$ we have

$$\begin{aligned} f(a + h) - f(a) &\approx \frac{1}{2} \frac{d^2 f}{dx^2}\Big|_a h^2 \\ &< 0, \text{ maximum} \\ &> 0, \text{ minimum} \\ &= 0, \text{ saddle point.} \end{aligned}$$

In the **n -dimensional case**, the study of critical points of a function of n variables is very similar.

Consider $z = f(x_1, x_2, \dots, x_n)$. Critical points are determined by solving

$$\nabla f\Big|_a = 0 \Rightarrow \left\{ \begin{array}{l} \frac{\partial f}{\partial x_1} = 0 \\ \vdots \\ \frac{\partial f}{\partial x_n} = 0 \end{array} \right\} \Rightarrow \mathbf{x} = \mathbf{a}$$

Having identified a critical point, \mathbf{a} , we find that

$$\begin{aligned} f(\mathbf{x}) &= f(\mathbf{a} + \mathbf{h}) > f(\mathbf{a}) & \forall \mathbf{x} \in S_r(\mathbf{a}) & \Rightarrow \text{minimum} \\ f(\mathbf{x}) &= f(\mathbf{a} + \mathbf{h}) < f(\mathbf{a}) & \forall \mathbf{x} \in S_r(\mathbf{a}) & \Rightarrow \text{maximum} \\ f(\mathbf{x}) &\gtrless f(\mathbf{a}) & & \Rightarrow \text{a saddle point} \end{aligned}$$

In some neighbourhood,

$$S_r(\mathbf{a}) = \{\mathbf{x} : 0 < |\mathbf{x} - \mathbf{a}| < r\},$$

of the critical point, we have the approximation

$$f(\mathbf{x}) \approx f(\mathbf{a}) + \nabla f\Big|_a \cdot \mathbf{h} + \frac{1}{2} \mathbf{h}^T \cdot H(\mathbf{a}) \cdot \mathbf{h}$$

Hence, for points $\mathbf{x} = \mathbf{a} + \mathbf{h}$ near the critical point $\mathbf{x} = \mathbf{a}$ we have

$$\begin{aligned} f(\mathbf{a} + \mathbf{h}) - f(\mathbf{a}) &\approx \frac{1}{2} \mathbf{h}^T \cdot H(\mathbf{a}) \cdot \mathbf{h} \quad \left(= \frac{1}{2} Q(\mathbf{a}) \right) \\ &< 0, \text{ maximum} \\ &> 0, \text{ minimum} \\ &= 0, \text{ saddle point.} \end{aligned}$$

☞ Mastery Check 3.2:

Use the Taylor approximation of order 2 (Section 2.I, Equation (2.7), Page 116) to determine the nature of the critical point for each of the functions $f_1(x, y) = x^2 + y^2 - 2x$ and $f_2(x, y) = x^2 - y^2 - 2x$.

(These functions were the subjects of Example 3.1 and Mastery Check 3.1.)



☞ Mastery Check 3.3:

For $z = f(x, y) = \ln(2x^2 + y^2) - 2y$, verify that $\nabla f|_{(0,1)} = \mathbf{0}$ at the point $(0, 1)$, but that the function is neither a maximum nor a minimum at this point.

Hint: See Mastery Check 2.25. The graph is in Figure 3.2.

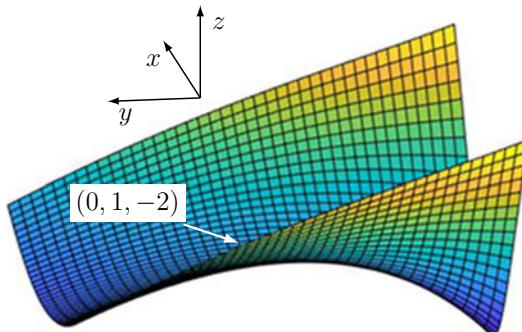


Figure 3.2 The graph of $z = \ln(2x^2 + y^2) - 2y$.



☞ Mastery Check 3.4:

Determine and classify all critical points of the function

$$f(x, y) = x^2 + y^2 + 2e^{xy+1}, \quad D_f = \mathbb{R}^2.$$



☞ Mastery Check 3.5:

Determine and classify all critical points of the function

$$f(x, y, z) = x^2y + y^2z + z^2 - 2x, \quad D_f = \mathbb{R}^3.$$



3.B Extreme points: The complete story

According to our discussion in the previous section, a critical point $\mathbf{a} \in D_f$ is a point of local maximum if $f(\mathbf{a}) \geq f(\mathbf{a} + \mathbf{h})$ for $|\mathbf{h}| \ll 1$ or a point of local minimum if $f(\mathbf{a}) \leq f(\mathbf{a} + \mathbf{h})$ for $|\mathbf{h}| \ll 1$.

Such a critical point is also called a point of *relative* maximum or *relative* minimum, respectively.

We contrast these references to *local* and *relative* quantities with the following definitions of *global* quantities.

Definition 3.4

A point $\mathbf{a} \in D_f$ is called a point of **absolute** $\begin{cases} \text{maximum} & \text{if} \\ \text{minimum} & \end{cases}$

$$f(\mathbf{a}) \begin{cases} \geq f(\mathbf{x}) & \text{for all } \mathbf{x} \in D_f \\ \leq f(\mathbf{x}) & \text{for all } \mathbf{x} \in D_f. \end{cases}$$

Remarks

- * The difference between Definition 3.1 and Definition 3.4 lies in the set of points considered. In Definition 3.1 only points in the immediate neighbourhood, $S_r(\mathbf{a})$, of \mathbf{a} are considered, while in Definition 3.4 all points in the domain, D_f , of the function are involved.
- * Definition 3.4 implies that a critical point, even if a point of local maximum or local minimum, need *not* be a point of *absolute* maximum or minimum.

Earlier we said that critical points are examples of extreme points. However, there are other types of extreme points which are *not* found using the gradient. These are

(i) *singular points of f* , points where ∇f does not exist (Figure 3.3):

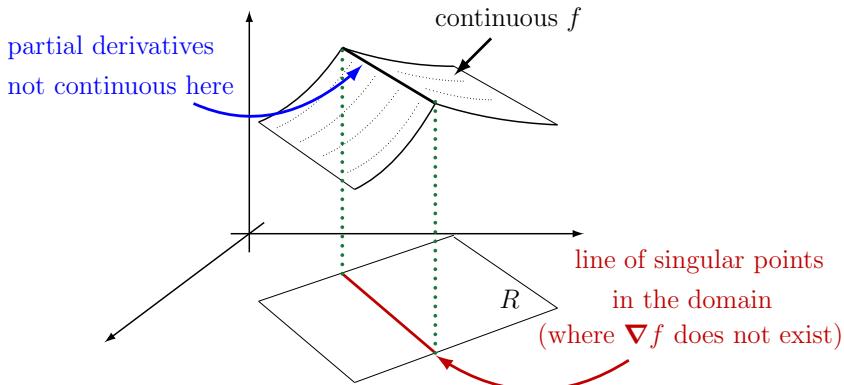


Figure 3.3 A function with singular points.

(ii) *boundary points of a restricted region $R \subset D_f$* (Figure 3.4):

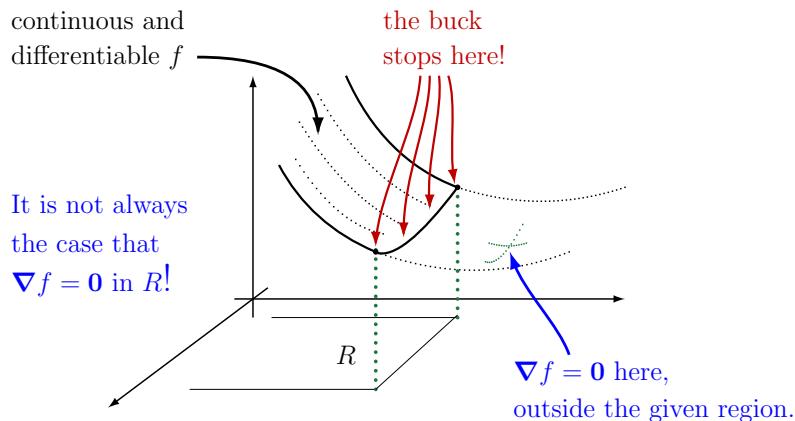


Figure 3.4 A function defined on a restricted region $R \subset D_f$.

Finding the absolute maxima and minima of functions $\mathbb{R}^n \rightarrow \mathbb{R}$ is part of the science of *optimization*. In general there are at least three categories of optimization problem:

- (1) optimizing over a compact domain;
- (2) optimizing completely free of restrictions;
- (3) optimizing functions under one or more constraints.

We shall study these in turn. We begin with optimizing over compact domains, assuming throughout that the functions involved are differentiable.

Optimization over compact domains

Recall our definition of a *compact set* (Definition 1.8): a set $\Omega \subseteq \mathbb{R}^n$ is said to be compact if it is *closed and bounded*.

For any function defined on a region $\Omega \subseteq D_f$ that is compact, we have the following very useful result.

Theorem 3.2

A continuous real-valued function defined on a compact region, Ω , obtains an absolute maximum and an absolute minimum value.

A few comments on this theorem are warranted.

Firstly, it is not necessary that the region being considered is the function's entire domain of definition, D_f , but it might be. The problem statement will usually specify this. If no region is given then the reader should assume the whole of D_f is implied.

Secondly, by Theorem 1.2, a continuous function defined on a closed and bounded region is necessarily bounded. This means that $|f(\mathbf{x})| < K$ for some $K \in \mathbb{R}$ and for all \mathbf{x} in that region. This simple result implies that we should expect f to exhibit an absolute minimum and an absolute maximum. In fact, this is the only time we are *guaranteed* that *absolute* maximum and minimum points exist.

The reader should always bear in mind that a *continuous* function is *not* necessarily differentiable everywhere. A consequence of this is that singular points can exist. These should then be inspected separately to any critical points. Naturally, the appealing notion of a closed and finite domain means

that the domain boundary (boundary points) need also to be considered separately.

We illustrate Theorem 3.2 in action with the following examples.

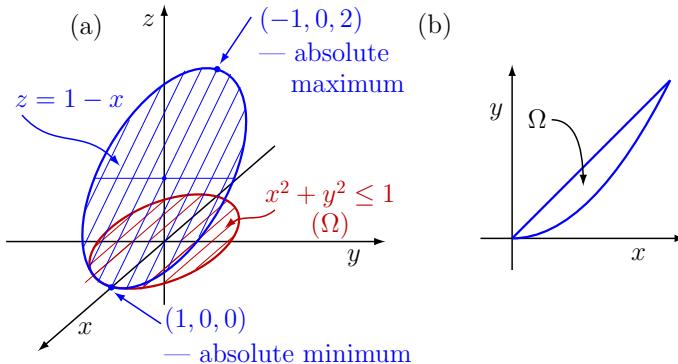


Figure 3.5 (a) The graph of $x + z = 1$ in Example 3.3;
 (b) The domain Ω in Example 3.4.

■ Example 3.3:

Consider the function

$$f(x, y, z) = x + z = 1,$$

and the region

$$\Omega = \{(x, y) : x^2 + y^2 \leq 1\}.$$

both of which are shown in Figure 3.5(a). In this case $\nabla f \neq 0$, but the function still has attained an absolute maximum and an absolute minimum.

■

■ Example 3.4:

Consider $f(x, y) = (y - x)e^{x^2 - y}$ for $x^2 \leq y \leq x$. Where does f achieve its maximum and minimum values?

We note that the domain Ω set by the above inequalities and shown in Figure 3.5(b) is non-empty provided $0 \leq x \leq 1$ and is compact, so we are sure to find

the extrema. The function is non-singular and we find that

$$\begin{aligned}\nabla f &= \begin{pmatrix} -1 + 2x(y-x) \\ 1 - (y-x) \end{pmatrix} e^{x^2-y} \\ &= \mathbf{0} \text{ when } \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} 1/2 \\ 3/2 \end{pmatrix},\end{aligned}$$

which is outside Ω . So, we need to inspect the boundaries.

On the boundary $y = x$, $f(x, y) = 0$.

On the boundary $y = x^2$, $f(x, y) = x^2 - x$, which has maximum 0 at $x = 0$ and $x = 1$, and minimum $-\frac{1}{4}$ at $x = \frac{1}{2}$, $y = \frac{1}{4}$.

Therefore, the absolute minimum is $f(1/2, 1/4) = -1/4$, and absolute maximum $f(x, x) = 0$.



Optimization free of constraints

In relaxing the condition of compactness, either by allowing the region $R \subseteq D_f$ to be unbounded or bounded but open, there is no longer any guarantee that points of finite maximum or minimum exist. For instance a function might become infinite at one or more points on the boundary of a bounded open set. Consider for example the case

$$f(x, y) = \frac{xy}{\sqrt{1 - x^2 - y^2}}, \quad R = D_f = \{(x, y) : x^2 + y^2 < 1\}.$$

The magnitude of this otherwise continuous function increases without bound as the independent variables approach the boundary of the unit disc: the function therefore attains neither an absolute maximum nor absolute minimum in D_f . In contrast, a continuous function on an unbounded domain may still attain a finite absolute maximum or minimum, as in the case of Mastery Check 3.7:

$$f(x, y) = \frac{4x}{1 + x^2 + y^2}, \quad R = D_f = \mathbb{R}^2.$$

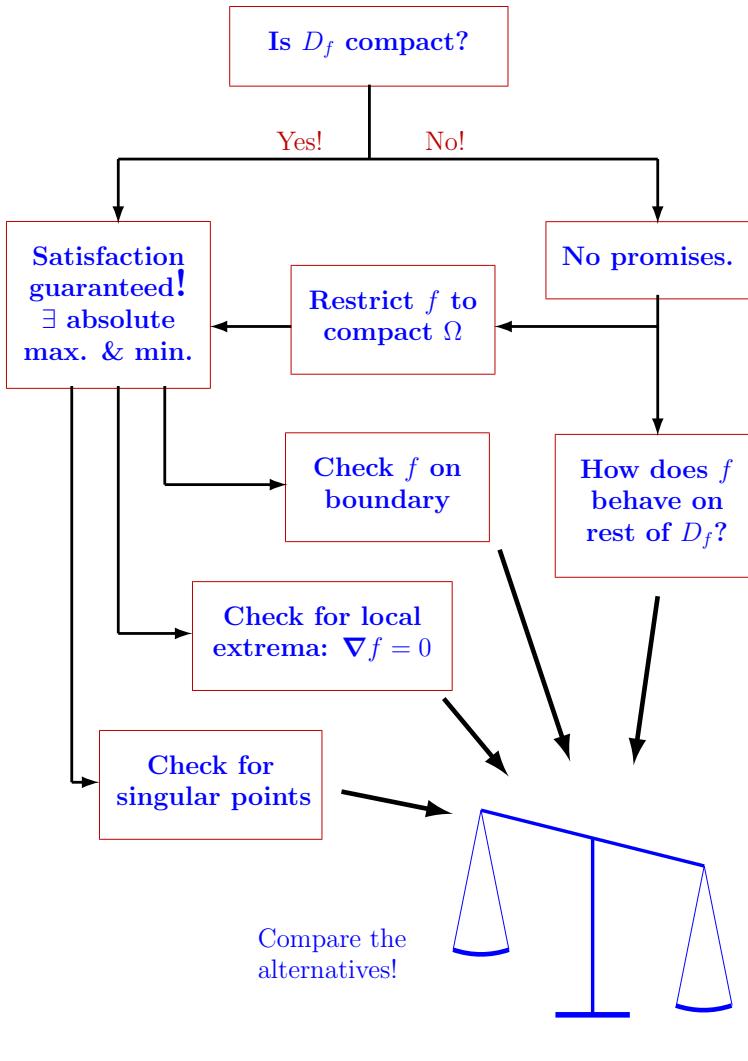
The function attains both an absolute maximum and an absolute minimum despite an unbounded domain of definition.

So, how does one proceed? We need only to modify the protocol for continuous functions on compact regions. This is the right-hand side of Flowchart 3.1.

We first ascertain if and where there are points on the edge of D_f at which the function f is discontinuous or not defined (that is, it “blows up”), and moreover whether f diverges to positive or negative infinity. This step will immediately answer the question of whether absolute extrema exist at all.

Flowchart 3.1: Optimization — how to play the game

Suppose f is continuous over D_f .



We next work with a convenient closed and bounded subregion of our own choosing over which f is continuous, and whose extent is easily characterized in terms of one or a few parameters; call this subset $\Omega \subset D_f$. Then, confining ourselves to Ω , we proceed as before and identify in Ω any critical points, points where the function's derivatives don't exist (are singular), and f 's behaviour on the boundary of Ω .

Finally, if D_f is unbounded or open we consider the function's behaviour outside of Ω over the rest of D_f .

The results of these three steps are then compared to determine which if any points in D_f are points of absolute maximum and absolute minimum.

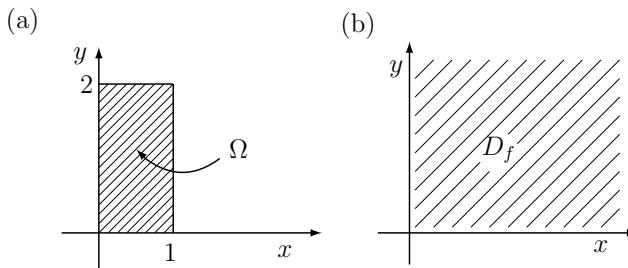


Figure 3.6 (a) Bounded domain, Ω , in Mastery Check 3.6;
 (b) Open domain, D_f , in Mastery Check 3.8.

✍ Mastery Check 3.6:

Determine the greatest and least values of

$$f(x, y) = y e^x - x y^2$$

in the region

$$\Omega = \{(x, y) : 0 \leq x \leq 1, 0 \leq y \leq 2\}$$

shown in Figure 3.6(a).

Hint: check for singular points and interior extreme points, then check boundary points — all of them! Then, and only then, should you plot the graph using MATLAB®.



Mastery Check 3.7:

Determine the greatest and least values of

$$f(x, y) = \frac{4x}{1 + x^2 + y^2} \text{ for } (x, y) \in D_f = \mathbb{R}^2.$$



Mastery Check 3.8:

Determine the greatest and least values of the function

$$f(x, y) = x + 8y + \frac{1}{xy}$$

on the open domain

$$D_f = \{(x, y) : x > 0, y > 0\}$$

shown in Figure 3.6(b).

Hint: check the behaviour of f for fixed y as $x \rightarrow 0, \infty$, and for fixed x as $y \rightarrow 0, \infty$, before proceeding to look for extrema.



Optimization subject to constraints

Function optimization under one or more constraints is a fairly common type of problem. It involves the task of maximizing or minimizing, generally expressed as optimizing, with respect to one or more variables, some quantity under special and restrictive conditions.

Such problems are generally expressed in a common way. For example, we:

minimize the physical dimensions *subject to* limited operating
 of electronic circuitry temperature;

$$\text{optimize} \quad f(\mathbf{x}) \quad \text{subject to} \quad g(\mathbf{x}) = 0$$

In some applications f is called the *objective* function while g is the *constraint*. In the case where there are more than one constraint, the additional constraints would similarly be expressed as equations such as $h(\mathbf{x}) = 0$.

The conceptual picture I like to impart to students is this: Suppose f were a mountain, while g gives rise to a walking track on the mountain side (see Figure 3.7). The constraint g is *not* itself the walking track, but a set of points in the plane of the mountain's base (the domain) that gives rise to the walking track. The absolute unconstrained maximum of f would be the mountain's peak, but the maximum of f subject to constraint $g = 0$ would give only the highest point on the track.

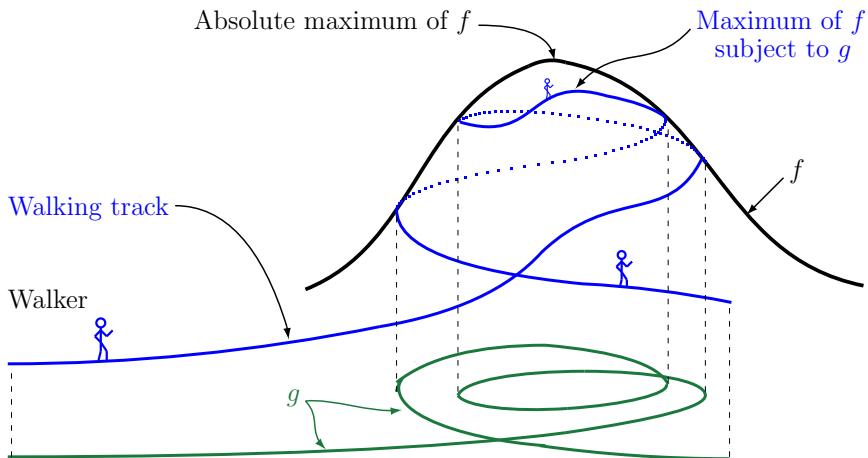


Figure 3.7 A constrained optimization problem.

In the following we will assume that both f and g have continuous partial derivatives of order 1 (at least). That is, f and g are of class C^1 .

We must also assume that *the level set*

$$L = \{\mathbf{x} : g(\mathbf{x}) = 0\}$$

is a non-empty subset of D_f , that is that it lies inside D_f . Moreover, we assume that L in the domain of g ($L \subset D_g$) has a dense set of points in common with D_f .

Although the concepts are applicable to any \mathbb{R}^n ($n < \infty$), it is convenient to restrict the following analysis and discussion to \mathbb{R}^2 .

Suppose

$$f : \mathbb{R}^2 \rightarrow \mathbb{R} \quad \text{and} \quad g : \mathbb{R}^2 \rightarrow \mathbb{R}$$

$$(x, y) \mapsto f(x, y) \quad (x, y) \mapsto g(x, y)$$

We presume $L \cap D_f \neq \emptyset$ and so we have the following picture, Figure 3.8. Since we consider a function of two variables, L is a curve in the \mathbb{R}^2 plane.

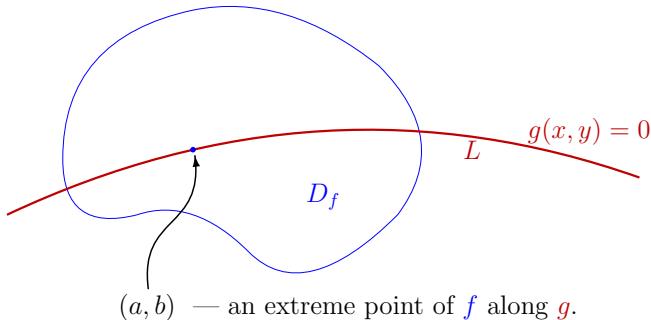


Figure 3.8 Domain D_f and the zero level set of g , L .

Suppose (a, b) is an interior local maximum or local minimum point of f when f is restricted to points on $g(x, y) = 0$ in \mathbb{R}^2 .

As g is a C^1 function there exists a continuous and differentiable parametrization of the level set of g in terms of a parameter t : $x = \phi(t)$, $y = \psi(t)$, and there exists a value t_0 such that $a = \phi(t_0)$, $b = \psi(t_0)$. Moreover, the level curve has a local tangent vector $(\phi'(t), \psi'(t))$ at any point $(\phi(t), \psi(t))$ along the curve. Finally, referring back to Section 2.E and the properties of the gradient, we deduce that $(\phi'(t), \psi'(t))$ is orthogonal to ∇g at (a, b) .

The single-variable function of t , $F(t) = f(\phi(t), \psi(t))$, is critical at t_0 . That is,

$$\begin{aligned} \frac{dF}{dt} \Big|_{t_0} &= \frac{\partial f}{\partial x} \frac{d\phi}{dt} \Big|_{t_0} + \frac{\partial f}{\partial y} \frac{d\psi}{dt} \Big|_{t_0} && \text{— a chain rule application} \\ &= \nabla f(a, b) \cdot (\phi'(t_0), \psi'(t_0)) \\ &= 0. \end{aligned}$$

The last equality implies that the plane vector $\nabla f(a, b)$ and the tangent vector to $g = 0$ at (a, b) , $(\phi'(t_0), \psi'(t_0))$ are orthogonal.

Since a tangent vector to the level curve $g = 0$ is orthogonal to the gradient of the function f , (Sections 1.F and 2.E), we have substantiated the following

theorem.

Theorem 3.3

If (a, b) is a point which lies in both D_f and D_g , and which is an extreme point of $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ under the constraint $g(x, y) = 0$, then the vectors $\nabla f(a, b)$ and $\nabla g(a, b)$ are parallel.

The reader should note the difference in conditions satisfied by the point (a, b) . It is a critical point of “ f subject to g ” (see the gradient equation below), but *not* a critical point of f alone.

We can now argue that if ∇g is not identically zero, there exists a $\lambda_0 \in \mathbb{R}$ such that

$$\nabla f(a, b) = -\lambda_0 \nabla g(a, b) \iff \nabla(f + \lambda g) \Big|_{a, b, \lambda_0} = 0,$$

which means that the 3D point $(a, b; \lambda_0)$ is a critical point of the new multi-variable function

$$F(x, y; \lambda) = f(x, y) + \lambda g(x, y).$$

This is called the *Lagrangian function* and λ is called the *Lagrange multiplier*.

The pictorial situation showing the relationship between the gradients of f and g , and the curves of constant f and of constant g , is shown in Figure 3.9.

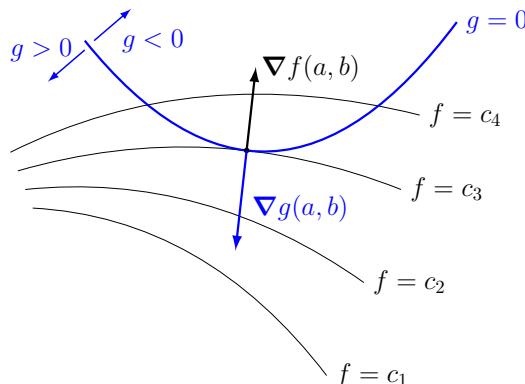


Figure 3.9 Level curves and gradient of f relative to the level curve $g = 0$ and ∇g .

Remarks

- * The critical points of $F(x, y; \lambda)$ are found from the set of equations

$$\left. \begin{array}{l} \frac{\partial F}{\partial x} = \frac{\partial f}{\partial x} + \lambda \frac{\partial g}{\partial x} = 0 \\ \frac{\partial F}{\partial y} = \frac{\partial f}{\partial y} + \lambda \frac{\partial g}{\partial y} = 0 \\ \frac{\partial F}{\partial \lambda} = g(x, y) = 0 \end{array} \right\} \quad \begin{array}{l} \text{— the condition } \nabla f \parallel \nabla g \text{ at } (a, b) \\ \text{— the constraint equation} \end{array}$$

- * What has been done, in effect, is that we have transformed a restricted 2D optimization problem into an *unrestricted* 3D optimization problem. We now need to find a, b , and λ_0 , to solve the full problem. Note that the actual value of λ_0 is not often needed, although it can be utilized in the numerical analysis of optimization problems.
- * The theory works beautifully for a general function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and constraint $g(x_1, \dots, x_n) = 0$:

$$\nabla f(\mathbf{a}) = -\lambda \nabla g(\mathbf{a}) \quad \text{if } \nabla g \neq 0.$$

This implies that the higher dimensional point $(\mathbf{a}; \lambda)$ is an extreme point of the $(n+1)$ -dimensional space on which the function $F : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$ is defined.

The theory of the Lagrangian function and Lagrange multiplier generalizes quite naturally to solve problems of optimizing under two or more constraints.

For example, suppose we wish to optimize $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ subject to the two constraints $g_1(x, y, z) = 0$ and $g_2(x, y, z) = 0$. Then the generalization of the 2D optimization problem would have $\nabla f(\mathbf{a})$, $\nabla g_1(\mathbf{a})$, and $\nabla g_2(\mathbf{a})$, be linearly dependent vectors. This means that there exist points \mathbf{a} and scalars $\lambda_0, \mu_0 \in \mathbb{R}$ such that

$$\nabla f(\mathbf{a}) + \lambda_0 \nabla g_1(\mathbf{a}) + \mu_0 \nabla g_2(\mathbf{a}) = 0$$

provided also that $\nabla g_1(\mathbf{a})$ and $\nabla g_2(\mathbf{a})$ are *not* co-parallel, which can be checked easily by showing that $\nabla g_1(\mathbf{a}) \times \nabla g_2(\mathbf{a}) \neq 0$.

In turn the linear dependency condition means that the higher dimensional point $(\mathbf{a}; \lambda_0, \mu_0)$ is a critical point of

$$F(\mathbf{x}; \lambda, \mu) = f(\mathbf{x}) + \lambda g_1(\mathbf{x}) + \mu g_2(\mathbf{x}).$$

Both λ and μ are Lagrange multipliers.

■ **Example 3.5:**

Suppose we wish to find the maximum value of $f(x, y) = e^{-(x^2+y^2)/2}$ subject to the constraint $g(x, y) = y - x^2 + 3 = 0$.

We set up the Lagrangian $F(x, y; \lambda) = e^{-(x^2+y^2)/2} + \lambda(y - x^2 + 3)$. Now

$$\nabla F = \begin{pmatrix} -xe^{-(x^2+y^2)/2} - 2\lambda x \\ -ye^{-(x^2+y^2)/2} + \lambda \\ y - x^2 + 3 \end{pmatrix} = \mathbf{0} \quad \text{when} \quad \begin{pmatrix} x \\ y \\ \lambda \end{pmatrix} = \begin{pmatrix} 0 \\ -3 \\ 3e^{-9/2} \end{pmatrix}.$$

The maximum value of f is $f(0, -3) = e^{-9/2}$.



☞ **Mastery Check 3.9:**

Find the maximum and minimum of $f(x, y, z) = x + 2y - 3z$ over the ellipsoid $x^2 + 4y^2 + 9z^2 \leq 108$. Try these two methods of solution:

- 1) Construct an argument that the maximum and minimum must both be on the surface of the ellipsoid. Then parameterize the ellipsoid in terms of the usual spherical coordinate angles (see Page 37ff) and look for critical points of f in terms of these angles.
- 2) Set up the Lagrangian $F = f + \lambda g$ for suitable g and look for critical points of F .

Use MATLAB[®] to draw the graph of the ellipsoid and the maximal level set.



3.C Differentials and error analysis

Definition 2.4 of a differentiable function allows the following interpretation:

$$f(\mathbf{x} + \Delta\mathbf{x}) - f(\mathbf{x}) = \text{grad } f(\mathbf{x}) \cdot \Delta\mathbf{x} + |\Delta\mathbf{x}|\rho(\Delta\mathbf{x}),$$

or, more simply, that $\Delta f(\mathbf{x}; \Delta\mathbf{x}) \approx \text{grad } f(\mathbf{x}) \cdot \Delta\mathbf{x}$ for $|\Delta\mathbf{x}| \ll 1$.

This leads to the idea of constructing a new function of both

$$\mathbf{x} = (x_1, x_2, \dots, x_n) \text{ and } \mathbf{dx} = (dx_1, dx_2, \dots, dx_n).$$

That is, it is a function of $2n$ variables.

Definition 3.5

The differential of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is defined as

$$\begin{aligned} df(\mathbf{x}, d\mathbf{x}) &= \nabla f(\mathbf{x}) \cdot d\mathbf{x} \\ &= \frac{\partial f}{\partial x_1}(\mathbf{x})dx_1 + \cdots + \frac{\partial f}{\partial x_n}(\mathbf{x})dx_n. \end{aligned}$$

The function df has the following three features:

- (1) It is an approximation to the change in f , Δf , coming from a change $\mathbf{x} \rightarrow \mathbf{x} + d\mathbf{x}$;
- (2) it is linear in $d\mathbf{x}$; and
- (3) it is a natural tool to use if considering overall error estimates when individual errors ($d\mathbf{x}$) are known.

The last feature identifies the differential's most useful application.

Suppose we have a quantity f whose value depends on many parameters, say x_1, \dots, x_n . Any errors incurred in measuring the x_i result in an error in the quantity f .

An estimate of the *maximum error* in f is thus given by

$$|df(\Delta\mathbf{x})| \leq \left| \frac{\partial f}{\partial x_1} \right|_{\mathbf{x}} |\Delta x_1| + \left| \frac{\partial f}{\partial x_2} \right|_{\mathbf{x}} |\Delta x_2| + \cdots + \left| \frac{\partial f}{\partial x_n} \right|_{\mathbf{x}} |\Delta x_n| \quad (3.2)$$

by the triangle inequality (Section 1.B).

The right-hand side of Equation (3.2) gives the *maximum* possible error only if one knows the *maximum* uncertainty of the individual x_i . If one knows the *exact* values of the dx_i (or Δx_i) including their signs then we use the differential $df(\mathbf{x}, d\mathbf{x})$ directly to give an approximate value to the change $\Delta f = f(\mathbf{x} + \Delta\mathbf{x}) - f(\mathbf{x})$.

3.D Method of least squares

In the year 1801 the world of astronomy was excited by the discovery of a new minor planet, Ceres, whose (rough) position in the night sky had been

noted a few times before it vanished from view. The young Carl Friedrich Gauss [19, 20] used a method — one he had worked out while still a student — to plot the orbit of the planet, and he was able to tell astronomers where in the sky to search for it. That method is the subject of this section.

In the first example that follows, we imagine the planet's orbit in the night sky to be a straight line which has to be fitted in some optimal fashion to a set of discrete pairs of observations which are subject to errors. In the second example, the method is applied to discrete observations on a supposed planetary orbit. In the third example, the method is extended to continuous domains.

The field of statistics deals with “observations” (measurements) on variables that are known to be subject to random errors. When we observe two or more variables at once, it is often appropriate to ask what is the relationship between these two variables. This question is the basis of the study known as “regression analysis”, which is outside the scope of this book. But the core of regression analysis is an application of the differential calculus called the method of least squares, invented independently by Gauss.

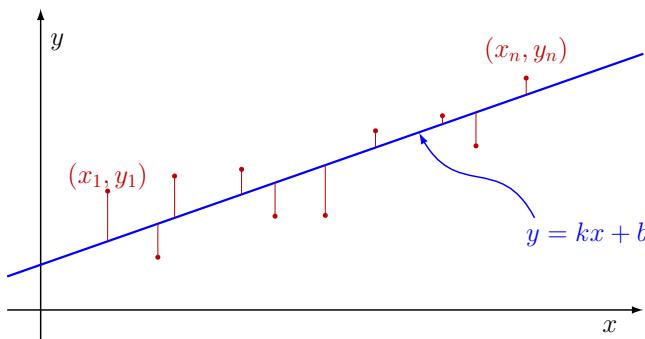


Figure 3.10 The line of best fit.

Fitting a straight line to observations

Suppose we believe that a variable y is in some way dependent on a variable x by the relation $y = kx + b$. In situations such as this it is implicitly assumed that the independent variable is *deterministic*, that is, given, and not subject to error, while the dependent variable is subject to observation or measurement errors.

To find the dependency relationship, we select a sequence of values $\{x_1, x_2, \dots, x_n\}$ of the independent variable, and measure the corresponding values $\{y_1, y_2, \dots, y_n\}$ of the dependent variable.

Because observed measurements always have some error associated with

them, the observations won't necessarily lie exactly on the straight line, but may fall above or below the line, as in Figure 3.10 above.

Problem: How to determine the “line of best fit” $y = k^*x + b^*$ through the “noisy” *discrete* experimental observations. (The values k^* and b^* become estimates of the true parameters k and b in the relation $y = kx + b$ which we believe connects x with y .)

Solution: We choose the line parameters k and b so that *the sum of the squares of the differences* between the observations and the fitted values is a minimum.

That is, we construct the function $S(k, b) = \sum_{i=1}^n (y_i - kx_i - b)^2$ of k and b from the known data and the desired model and seek its minimum to give us the optimal k and b values. We apply the techniques of the preceding chapter to get the critical points of S by solving the two equations

$$\frac{\partial S}{\partial k} = 0 \quad \text{and} \quad \frac{\partial S}{\partial b} = 0.$$

Because of its form S has no upper bound but does have a lower bound. Moreover, there will be only one critical point which will be the (k, b) point for which S has its minimum value. In fact, with the above explicit expression for S we get the equations

$$\sum_{i=1}^n 2(y_i - k^*x_i - b^*)(-x_i) = 0 \quad \text{and} \quad \sum_{i=1}^n 2(y_i - k^*x_i - b^*)(-1) = 0.$$

That is, $\sum_{i=1}^n x_i y_i - k^* \sum_{i=1}^n x_i^2 - b^* \sum_{i=1}^n x_i = 0$ and $\sum_{i=1}^n y_i - k^* \sum_{i=1}^n x_i - nb^* = 0$.

This pair of equations leads to the 2×2 matrix equation

$$\begin{pmatrix} \sum x_i^2 & \sum x_i \\ \sum x_i & n \end{pmatrix} \begin{pmatrix} k^* \\ b^* \end{pmatrix} = \begin{pmatrix} \sum x_i y_i \\ \sum y_i \end{pmatrix},$$

which can be solved to give the optimal k^* and b^* .

☞ Mastery Check 3.10:

Invert this matrix equation to obtain explicit estimates of k^* and b^* .

Hint: First define symbols $\bar{x} = \sum x_i/n$, etc.

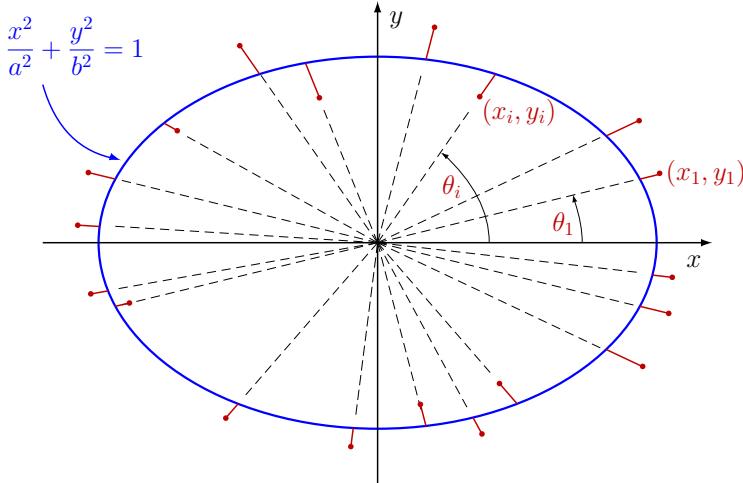


Figure 3.11 The ellipse of best fit.

Fitting a conic to observations

There is nothing unique about fitting straight lines: observations can follow other functional forms depending on the problem being considered.

Suppose, for example, the observations $\{(x_i, y_i)\}$ are scattered in the shape of an ellipse with the origin as centre. (Remember the case of the minor planet.) We want to determine that ellipse,

$$\frac{x^2}{a^2} + \frac{y^2}{b^2} = 1,$$

with optimal a and b that best represent the observations.

It is convenient at this point to convert to polar coordinates (Page 25)

$$(x_i, y_i) \longrightarrow (r_i, \theta_i) : \begin{cases} x_i = r_i \cos \theta_i, & i = 1, 2, \dots, n. \\ y_i = r_i \sin \theta_i, & \end{cases}$$

Choose as error function E the sum of squares of the distances of the observation pairs (x_i, y_i) from the corresponding points $(x(\theta_i), y(\theta_i))$ on the ellipse along the rays at the given θ_i , as in Figure 3.11.

What we are doing from here on is assuming that the angles θ_i have no associated errors, but that the radial distances r_i are subject to error. Analogous to the straight line example we form the sum of squares of differences

$$\begin{aligned} E &= \sum_{i=1}^n \left[(x(\theta_i) - x_i)^2 + (y(\theta_i) - y_i)^2 \right] \\ &= \sum_{i=1}^n \left[(a - r_i)^2 \cos^2 \theta_i + (b - r_i)^2 \sin^2 \theta_i \right] \end{aligned}$$

As before, look for the critical points of $E(a, b)$ with respect to a and b .

$$\frac{\partial E}{\partial a} = 0 \quad \text{and} \quad \frac{\partial E}{\partial b} = 0.$$

The 2×2 system of equations can be solved to give the optimal ellipse parameters

$$a^* = \frac{\sum_{i=1}^n r_i \cos^2 \theta_i}{\sum_{i=1}^n \cos^2 \theta_i} \quad \text{and} \quad b^* = \frac{\sum_{i=1}^n r_i \sin^2 \theta_i}{\sum_{i=1}^n \sin^2 \theta_i}.$$

Try to confirm these expressions for your own peace of mind!

Least-squares method and function approximations

The least-squares method admits a continuous version in which the discrete sum is replaced with an integral, and the discrete observations are replaced with a continuously varying function. In truth, the function need not be continuous, only integrable. However, for this introduction we'll stay with the more restrictive but less complicated case.

Suppose $f : \mathbb{R} \rightarrow \mathbb{R}$ is a continuous function on some interval $a \leq x \leq b$. The problem that is now posed is how do we “best” approximate f with an “approximating” function $g(x; \lambda_1, \dots, \lambda_n)$, where $\{\lambda_i\}_{i=1}^n$ is a set of n parameters?

In solving this problem we first note that “best” will depend on how this is measured. That is, this will depend on the choice of *distance function*. (See [9] and [10] for a more complete discussion.) Second, we note that the choice of the approximating function g is critical.

In the theme of least squares, the distance function most often considered is

$$I = \int_a^b (f(x) - g(x; \lambda_1, \dots, \lambda_n))^2 dx$$

which mathematicians refer to as the L_2 distance function.

The choice of the approximating function g is usually (but not always) a linear function of the λ_i . The general structure is

$$g(x; \lambda_1, \dots, \lambda_n) = \sum_{i=1}^n \lambda_i \phi_i(x),$$

where the functions $\phi_i(x)$ are chosen to satisfy some criteria specific to the problem being considered.

The best approximation to f is then the choice of the λ_i that minimize the distance function. Thus, we look for critical points of I :

$$\left. \begin{array}{l} \frac{\partial I}{\partial \lambda_1} = 0 \\ \vdots \\ \frac{\partial I}{\partial \lambda_n} = 0 \end{array} \right\} \quad \text{---} n \text{ equations in } n \text{ unknowns} \quad \implies \lambda_i^*, i = 1, \dots, n.$$

Some common and useful, but not exclusive, choices of g (that is, the ϕ_i) are:

$$g(x; \lambda_0, \dots, \lambda_n) = \lambda_0 + \lambda_1 x + \dots + \lambda_n x^n \quad \text{---} \text{polynomial approximation}$$

$$g(x; \lambda_1, \dots, \lambda_n) = \sum_{k=1}^n \lambda_n \sin kx \quad \text{---} \text{trigonometric (sine) approximation}$$

Notice that in these cases the λ_i appear as linear coefficients.

Mastery Check 3.11:

Find constants $a, b, c \in \mathbb{R}$ which minimize the integral

$$\begin{aligned} J &= \int_{-1}^1 (x - a \sin(\pi x) - b \sin(2\pi x) - c \sin(3\pi x))^2 dx \\ &= \int_{-1}^1 (x - g(x; a, b, c))^2 dx. \end{aligned}$$

Plot $f(x) = x$ and $g(x)$ over the interval $[-1, 1]$.



3.E Partial derivatives in equations: Partial differential equations

Although we treat the topic of partial differential equations (PDEs) as an application of differential calculus, the body of knowledge is immense and can (and sometimes does) easily cover several thick volumes of dedicated texts. Our coverage keeps with the idea of giving only some basic practical information that students of engineering and science will find beneficial in their undergraduate studies, and possibly later.

Although PDEs—which are simply scalar-valued equations that describe fixed relationships between the various partial derivatives of a scalar function—come in all shapes and sizes, the type that has been most thoroughly studied and on which we will focus is the *second-order, linear PDE*. The following definition features a scalar function u of two independent variables. However, second-order linear PDEs may depend on more than two variables.

Definition 3.6

A second-order, linear PDE with constant coefficients satisfied by a function $u \in C^2(\mathbb{R}^2)$ on a domain $D_u \subset \mathbb{R}^2$ is an equation of the form

$$Lu = Au_{xx} + Bu_{xy} + Cu_{yy} + Du_x + Eu_y + Fu + G = 0,$$

where the constant real numbers A, B, C, D, E, F satisfy the condition $A^2 + B^2 + C^2 \neq 0$, and $G(x, y)$ is a real-valued function defined on D_u .

Even in this restrictive case there is a quite a lot known. Unfortunately, the following remarks and other comments made in this section do not do justice to the discipline or to the effort that has gone into accumulating all that is known. Nevertheless, it is hoped that the reader will be somewhat better oriented with the subject with the salient terminology and information provided here.

Remarks

- * In this section subscripts x , y , and t on the dependent variable u will denote partial derivatives w.r.t. those variables.
- * The “ L ” as used here and in many texts will denote a “linear partial differential operator” with the property that $L(u + v) = Lu + Lv$.

- * “Second order” refers to the highest *order* of partial derivative appearing in an equation.
— a *third order* PDE would have one or more of $u_{xxx}, u_{xxy}, \dots, u_{yyy}$.
- * “Linear” refers to PDEs containing terms which are *at most* proportional to u or one of its derivatives. — $G(x, y)$ does not involve u .
- * The term “PDE” means an equation relating a function to its partial derivatives, or relating partial derivatives to one another.
- * The function u is called a *solution* of the PDE.
- * The PDE relationship applies locally in an open domain D_u but not on the boundary ∂D_u .
- * A *time-dependent* PDE generalizes Definition 3.6 to include an explicit dependence on the time variable t through additional terms involving partial derivatives of u w.r.t. t .

Definition 3.7

The quantity $B^2 - 4AC$ is called the *discriminant* of the PDE.

The reader may speculate that this discriminant is somehow connected with the discriminant appearing in the solution of a quadratic equation. The connection *is* there and *is* important but would take us too far afield to explain what it is about. The curious reader might refer to the bibliography [11–13, 15] for more information. However, for our purposes it is sufficient to just state without proof some important facts.

Firstly, the value of the discriminant (more significantly its sign) determines the type or nature of the PDE. A PDE is said to be of

- (a) hyperbolic type $\iff B^2 - 4AC > 0$;
- (b) elliptic type $\iff B^2 - 4AC < 0$;
- (c) parabolic type $\iff B^2 - 4AC = 0$.

These types of PDE exhibit distinctively different behaviour with regard to their properties and conditions. So too do their solutions.

Summary of some relevant facts

* Linearity

- Adding any number of possible solutions leads to a new solution. Solutions $u_1(x, y)$, $u_2(x, y)$, and constants α and β , can be combined to give a new solution: $\alpha u_1(x, y) + \beta u_2(x, y)$
- The solutions are not unique, since any scalar multiple of a solution is also a solution.
- To get uniqueness we impose additional conditions, which are used to specify the scalar multipliers needed. These conditions are classed as either *initial* conditions or *boundary* conditions.

* Initial conditions

Initial conditions are needed for the solution of time-dependent PDEs and inform how the solution begins. They are of the form:

- “displacement”: $u(\cdot, t = 0) = u_0(\cdot)$, with $u_0(\cdot)$ specified;
- “velocity”: $u_t(\cdot, t = 0) = v_0(\cdot)$, with $v_0(\cdot)$ specified;

where the “.” indicates that other independent variables apply.

* Boundary conditions

Boundary conditions on the domain boundary ∂D_u relate to spatial-dependent PDEs.

These are of three forms:

- Dirichlet conditions
 $u(\mathbf{r} = \mathbf{r}_b, t) = U_b(\mathbf{r}_b, t)$, with U_b specified at $\mathbf{r}_b \in \partial D_u$.
- Neumann conditions (Figure 3.12)
 Let $\mathbf{N}(\mathbf{r}_b)$ be the outward normal to the boundary ∂D_u . Then the Neumann condition fixes the normal component of the solution gradient $\nabla u(\mathbf{r} = \mathbf{r}_b, t) \cdot \mathbf{N}(\mathbf{r}_b) = V_b$, with V_b specified at $\mathbf{r}_b \in \partial D_u$.

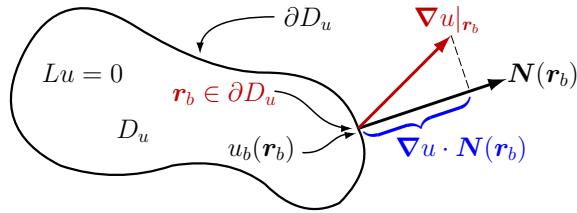


Figure 3.12 Neumann condition on the boundary of domain D_u .

- Cauchy conditions

These specify both the solution and the normal component of the solution gradient at the same point.

$$\begin{cases} u(\mathbf{r} = \mathbf{r}_b, t) = U_b(\mathbf{r}_b, t), \\ \nabla u(\mathbf{r} = \mathbf{r}_b, t) \cdot \mathbf{N}(\mathbf{r}_b) = V_b(\mathbf{r}_b, t), \end{cases}$$

with U_b and V_b specified at $\mathbf{r}_b \in \partial D_u$.

- * Mixed conditions

These are needed for solutions of PDEs that depend on both time and space in bounded spatial domains.

For example,

$$\begin{cases} u(\mathbf{r}, t = 0) = u_0(\mathbf{r}), \\ u(\mathbf{r} = \mathbf{r}_b, t) = U_b(\mathbf{r}_b), \\ \nabla u(\mathbf{r} = \mathbf{r}_b, t) \cdot \mathbf{N}(\mathbf{r}_b) = V_b(\mathbf{r}_b). \end{cases}$$

(U_b, V_b are not time-dependent. Contrast this set with the boundary condition set on Page 162, Equation (3.5).)

- * Rule of thumb for applying these conditions

- We need as many initial conditions as the order of the highest order time derivative;
- We need as many boundary conditions as the order of the highest order spatial derivative in *each independent variable*.

- * Types of PDE

- The discriminant (Definition 3.7) indicates the type of PDE under consideration. In turn, the type dictates what additional conditions are needed and allowed in order to determine a unique solution to the problem.

- It is easily shown [12] that any second-order, linear PDE can be reduced to one of three forms in the respective classes of *hyperbolic*, *elliptic*, and *parabolic* PDE. We illustrate these forms by means of the most common example from each class.

Classic equations of mathematical physics

Heeding the last bullet point, we focus attention on three classic equations of mathematical physics: the Laplace equation (*elliptic*), the diffusion equation (*parabolic*) and the wave equation (*hyperbolic*). The applications in which they arise are mentioned in their respective sections below. Suffice to say that their solutions play central roles in describing a wide range of physical phenomena.

The Laplace equation

The Laplace equation is perhaps the PDE with the longest history, dating back to the late 18th century when it was introduced by Pierre-Simon Laplace [17] in the context of astronomy. Apart from being the equation satisfied by a free-space gravitational potential, the PDE also arises in the context of the steady dynamics of an incompressible fluid, and is satisfied by the electrostatic potential in electrostatic theory and the velocity potential in linear acoustic theory [13, 15]. Despite their different physical origins, the mathematical properties of the potential functions are the same.

Definition 3.8

Suppose $u \in C^2(\mathbb{R}^3)$. The equation

$$u_{xx} + u_{yy} + u_{zz} = 0 \quad (3.3)$$

satisfied by u in an open domain $D_u \subset \mathbb{R}^3$ is called the (three-dimensional) **Laplace equation (potential equation)**.

Remarks

- * The Laplace equation is an example of an *elliptic* PDE.
- * The appropriate boundary conditions to use to establish a unique solution of Laplace's equation are of Dirichlet type (Page 154).

Definition 3.9***The 2D Dirichlet Problem***

Let D_u be a simply-connected bounded region in \mathbb{R}^2 and let ∂D_u be its boundary. Let $g(x, y)$ be a continuous function on ∂D_u . Then the Dirichlet problem is to find the (unique) function $u = f(x, y)$ such that

- (a) u is defined and continuous on \bar{D}_u ,
- (b) u satisfies the 2D Laplace equation, $u_{xx} + u_{yy} = 0 \forall (x, y) \in D_u$,
- (c) $u = g(x, y) \forall (x, y) \in \partial D_u$.

- * There exists a branch of pure mathematics called *harmonic analysis* that specializes on properties and behaviour of solutions of the 2D Laplace equation.

Definition 3.10

A function $u = f(x, y) \in C^2(\mathbb{R}^2)$ that is a solution of the 2D Laplace equation (3.3) is called a **harmonic function**.

- * Equation (3.3) is sometimes abbreviated $\Delta u \equiv \nabla^2 u \equiv \nabla \cdot \nabla u = 0$, where $\nabla = \left(\frac{\partial}{\partial x}, \frac{\partial}{\partial y}, \frac{\partial}{\partial z} \right)$ is the gradient operator, and $\Delta \equiv \nabla^2 \equiv \nabla \cdot \nabla$ — called the *Laplace operator* or *Laplacian operator*.
- * All terms appearing in the Laplace equation involve the unknown function u . It is therefore said to be a *homogeneous* PDE. If a term not involving u were present it would then be an example of an *inhomogeneous* equation.

Definition 3.11

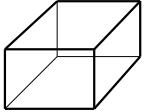
Suppose $u \in C^2(\mathbb{R}^3)$ and $g \in C(\mathbb{R}^3)$. The equation

$$\Delta u = u_{xx} + u_{yy} + u_{zz} = g(x, y, z)$$

satisfied by u in an open domain $D_u \cap D_g \subset \mathbb{R}^3$ is called a **Poisson equation** (inhomogeneous version of the Laplace equation).

- * The Laplacian operator (appearing in 3D diffusion, potential, and wave problems) can be expressed in different forms depending on the problem

geometry. The most common and the most studied forms are those shown in Figure 3.13 (see Section 1.D).

Rectangular problems	Cylindrical problems	Spherical problems
(x, y, z)	(r, θ, z)	(ρ, ϕ, θ)
		

$$\Delta u = u_{xx} + u_{yy} + u_{zz}.$$

$$\Delta u = \frac{1}{r} (r u_r)_r + \frac{1}{r^2} u_{\theta\theta} + u_{zz}.$$

$$\Delta u = \frac{1}{\rho^2} (\rho^2 u_\rho)_\rho + \frac{1}{\rho^2 \sin \phi} (\sin \phi u_\phi)_\phi + \frac{1}{\rho^2 \sin^2 \phi} u_{\theta\theta}.$$

Figure 3.13 The Laplacian expressed in different curvilinear coordinates.

In Figure 3.13 we use the notation introduced incidentally in the comment immediately following Definition 2.3, where subscripts denote partial differentiation with respect to the variable featured in the subscript. For example,

$$(r u_r)_r \equiv \frac{\partial}{\partial r} \left(r \frac{\partial u}{\partial r} \right).$$

The diffusion equation

The diffusion equation is also commonly referred to as the equation of conduction of heat. The equation is satisfied by a function describing the temporal and spatial development of temperature in uniform medium. In the early 1800s Jean-Baptiste Joseph Fourier [17] provided the first in-depth study of this equation in the context of heat transfer, and of its solution by an innovative (for that time) solution method [13, 15]. The equation, with the same form and therefore with the same mathematical properties of its solution, is also satisfied by a function describing the concentration of material diffusing by random processes. The diffusion equation also arises in problems as diverse as radiative transfer, insect migration and the spread of infection.

Definition 3.12

Suppose $u \in C^2(\mathbb{R}^3 \times [0, \infty))$. The equation

$$\frac{\partial u}{\partial t} - k \left(\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} + \frac{\partial^2 u}{\partial z^2} \right) = 0$$

satisfied by u in the domain $D_u \subset \mathbb{R}^3 \times [0, \infty)$ is called the (three-dimensional) **diffusion equation (heat equation)**.

Remarks

- * The diffusion equation is an example of a *parabolic* PDE.
- * The appropriate supplementary conditions for the diffusion equation are of mixed type, with one initial condition and Dirichlet or Neumann boundary conditions.
- * When u is the *temperature* at point r at time t the diffusion equation is called the *heat equation*.
- * The constant k is called the thermal *diffusivity*: $k = \frac{K}{s\rho}$, where K is thermal *conductivity*, ρ is *mass density*, and s is *specific heat*.

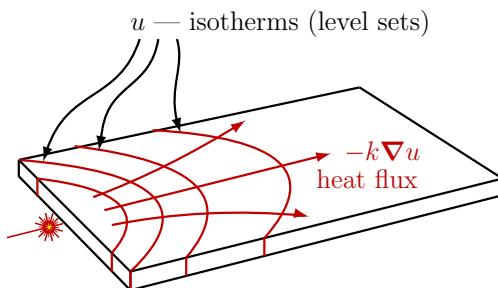


Figure 3.14 Illustration of heat diffusion in a slab.

The observant reader will have noticed our depiction here of the local heat flux being orthogonal to the isotherms. The gradient operator has that property, and moreover points in the direction of increasing scalar (isotherm) values. The negative sign inherent in the heat flux definition, shown in Figure 3.14, reverses that direction.

- * The 3D diffusion (heat) equation can be abbreviated

$$\frac{\partial u}{\partial t} = k \nabla \cdot \nabla u = k \Delta u$$

(↑
(Δ — the Laplacian operator. —))

- * In the limit $t \rightarrow \infty$, $\frac{\partial u}{\partial t} \rightarrow 0$, and $u(\mathbf{r}, t \rightarrow \infty) \rightarrow u_\infty(\mathbf{r})$,

the steady-state solution of the Laplace equation: $\Delta u_\infty = 0$. Thus, the Laplace equation is the temporal limit of the diffusion equation.

The boundary-value problem for the diffusion equation

For some $b > 0$ consider \bar{R} to be the closed set $\bar{R} = \{(x, t) : 0 \leq x \leq \pi, 0 \leq t \leq b\}$ with boundary $\partial R = \Gamma_1 \cup \Gamma_2 \cup \Gamma_3 \cup \Gamma_4$ where

$$\begin{aligned}\Gamma_1 &= \{(0, t) : 0 \leq t \leq b\}, \quad \Gamma_2 = \{(x, 0) : 0 \leq x \leq \pi\}, \\ \Gamma_3 &= \{(\pi, t) : 0 \leq t \leq b\}, \quad \Gamma_4 = \{(x, b) : 0 \leq x \leq \pi\}.\end{aligned}$$

Definition 3.13

The boundary-value problem for the 1D heat-conduction problem is to find a function $u = f(x, t)$ such that

- (a) f is continuous on \bar{R} ,
- (b) f satisfies the heat (diffusion) equation in R ,
- (c) f satisfies $f(0, t) = g_1$ on Γ_1 ; $f(\pi, t) = g_3$ on Γ_3 ; $f(x, 0) = g_2$ on Γ_2 .

Remarks

- * Note that there is no explicit condition to be satisfied on Γ_4 !
- * The upper bound b could be taken as large as desired. The only restriction is that $b > 0$. That is, we *solve forward in time*. In fact, we often take $b = \infty$ so the time variable $t \in [0, \infty)$.
- * The conditions to be applied on Γ_1 and Γ_3 are *boundary conditions*, while the condition on Γ_2 is an *initial condition*.

And, in terms of the time variable t and space variable x ,

- * The boundary conditions $u(0, t) = g_1(t)$ and $u(\pi, t) = g_3(t)$ describe the case of prescribed end temperatures, which may be time-dependent.

- * The boundary conditions on Γ_1 and Γ_3 may instead be $u_x(0, t) = h_1(t)$ and $u_x(\pi, t) = h_3(t)$ which describe the case of *prescribed end heat fluxes*: rates at which heat is conducted across the ends. If $h_1 = h_3 = 0$, the ends are *insulating*. That is, there is no heat conduction.

The wave equation

The wave equation is the third classic equation of mathematical physics we discuss. As the name suggests it is an equation governing wave-like phenomena; not simply propagation, but oscillatory motion. Its discovery is associated with the names of the 18th-century mathematicians Jean-Baptiste d'Alembert (1D version) and Leonard Euler (3D version) [17]. The equation governs the linear dynamics of water waves, sound waves in linear acoustics, linear elastic vibrations in solid mechanics, and light waves in electromagnetism [13, 15].

Definition 3.14

Suppose $u \in C^2(\mathbb{R}^3 \times \mathbb{R})$. The equation

$$u_{tt} - c^2(u_{xx} + u_{yy} + u_{zz}) = G(x, y, z, t) \quad (3.4)$$

satisfied by u in an open domain $D_u \subset \mathbb{R}^3 \times \mathbb{R}$ is called the (three-dimensional) **wave equation**. Here, c is called the **wave speed**, and G is defined on D_u .

Remarks

- * The wave equation is an example of a *hyperbolic* PDE.
- * The appropriate supplementary conditions for a unique solution of the wave equation are either of Cauchy type for unbounded domains or of mixed type for bounded domains.
- * Unlike the diffusion equation, the wave equation can be solved forwards or backwards in time.
- * If $G(x, y, z, t) \equiv 0$ then Equation (3.4) is the *homogeneous* wave equation.
- * The wave speed c is also called the phase speed of propagation of information.

* The general form of Equation (3.4) can be abbreviated

$$\frac{\partial^2 u}{\partial t^2} - c^2 \nabla \cdot \nabla u = \frac{\partial^2 u}{\partial t^2} - c^2 \Delta u = G(x, y, z, t)$$

(Δ — the Laplacian operator. 

A general-purpose method of solution: Separation of variables

The method we shall now describe is based on the fact that the linear PDEs just described are *separable* in a number of coordinate systems. This means that their solutions can be expressed in terms of factors involving only one variable. A necessary assumption is that one of the coordinates of the system of choice is constant over a surface on which boundary conditions are prescribed. The method is actually applicable to a wider variety of linear PDEs than the ones that are here highlighted, defined on bounded or semi-bounded domains. Consequently, all students should be familiar with this method of solution [12, 13, 15]. We illustrate the approach by applying it to a simple problem involving 1D wave propagation.

Consider the following *mixed boundary-value problem* on the space-time, semi-infinite strip $D_u = [a, b] \times [0, \infty)$ shown in Figure 3.15. Note that this problem involves one space dimension x and one time dimension t .

Definition 3.15

The 1D wave problem. Find a function $u=f(x, t)$ on $D_u \subset \mathbb{R} \times [0, \infty)$ such that for functions g_1 and g_2 defined on $[a, b]$ and continuous functions h_1 and h_2 on $[0, \infty)$, u satisfies the homogeneous wave equation

$$\left. \begin{aligned} u_{tt} - c^2 u_{xx} &= 0, \\ \text{with initial conditions } &\left\{ \begin{aligned} f(x, 0) &= g_1(x) \\ f_t(x, 0) &= g_2(x) \end{aligned} \right\} a \leq x \leq b \\ \text{and boundary conditions } &\left\{ \begin{aligned} f(a, t) &= h_1(t) \\ f(b, t) &= h_2(t) \end{aligned} \right\} t \geq 0. \end{aligned} \right\} \quad (3.5)$$

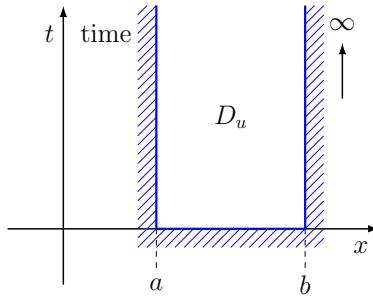


Figure 3.15 Domain D_u of the 1D wave problem.

Remarks

- * Equation (3.5) is a quite general form of a mixed b.v.p. for the wave equation in one space variable, but not the most general (see Supplementary problem 21).
- * The smoothness (that is, degree of differentiability) of the boundary data, h_1 and h_2 , and the initial data, g_1 and g_2 , determine the smoothness of the solution $u = f(x, t)$.

To illustrate the separation of variables method we consider the special case:

$$(h_1, h_2, g_1, g_2) = (0, 0, g, 0).$$

Similar analyses apply in the other special cases: $(0, 0, 0, g)$, $(0, h, 0, 0)$, . . .

For the most general case of h_1 , h_2 , g_1 , and g_2 all nonzero, we appeal to the *principle of superposition* (defined shortly). For the moment, we consider the simpler mixed b.v.p.

$$\left. \begin{aligned} u_{tt} - c^2 u_{xx} &= 0, \\ f(x, 0) &= g(x), \quad 0 \leq x \leq L \\ f_t(x, 0) &= 0, \quad 0 \leq x \leq L \\ f(0, t) &= f(L, t) = 0, \quad t \geq 0 \end{aligned} \right\} \quad (3.6)$$

Remarks

- * Equation (3.6) is the mathematical model of a vibrating string of length L and density ρ under tension τ , fixed at its ends and subject to small displacements elsewhere along its length (Figure 3.16).

- * $c = \sqrt{\frac{\tau}{\rho}}$; the *wave speed* is determined by the string properties.
- * If $g \equiv 0$ then the unique solution would be $f \equiv 0$.
 - a typical way of proving uniqueness

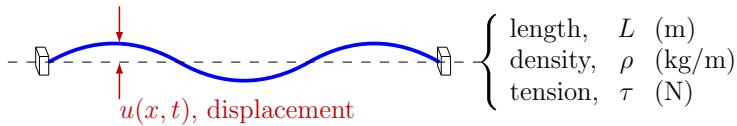


Figure 3.16 Schematic of a vibrating string.

For the vibrating string problem, u ($|u| \ll L$) is the local dynamic displacement of the stretched string at position x at time t . One sets the string vibrating either by “plucking” or “striking”, or any linear combination of these. In our example we have chosen the first means. Suppose the string is pulled aside at $t = 0$ a distance h from equilibrium at $x = b$ (Figure 3.17) and released. This condition defines the function $g(x)$.

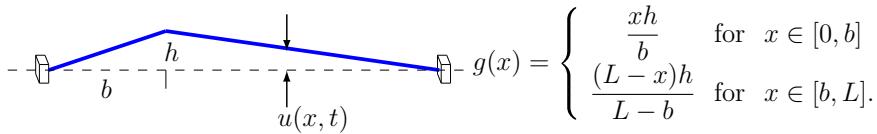


Figure 3.17 Profile $g(x)$ of a string stretched a distance h at $x = b$.

We proceed in steps through the separation of variables method of solution of this problem.

Step 1:

The separation of variables method always begins by assuming a nontrivial solution in the form of a product of functions of the independent variables. In this case we assume the form

$$f(\mathbf{x}, \mathbf{t}) = \mathbf{X}(\mathbf{x})\mathbf{T}(\mathbf{t}).$$

That is, the solution is a product of a function of \mathbf{x} only (\mathbf{X}) and a function of \mathbf{t} only (\mathbf{T}).

Substitution into Equation (3.6a) gives

$$\mathbf{X}(\mathbf{x}) \frac{d^2\mathbf{T}}{dt^2} - c^2 \mathbf{T}(t) \frac{d^2\mathbf{X}}{dx^2} = 0 \quad (3.7)$$

and assuming $X(x) \neq 0$ and $T(t) \neq 0$, Equation (3.7) implies

$$\frac{1}{X(x)} \frac{d^2X}{dx^2} = \frac{1}{c^2} \frac{1}{T(t)} \frac{d^2T}{dt^2} = \mu$$

$$\left\{ \begin{array}{l} \text{function} \\ \text{of } x\text{-only} \end{array} \right\} = \left\{ \begin{array}{l} \text{function} \\ \text{of } t\text{-only} \end{array} \right\} = \{\text{constant}\} \quad (3.8)$$

The most critical element of the separation of variables method is the fact that the two expressions on the left of Equation (3.8) can be equal *only* if both equal the same constant! This allows the separation to occur. Equation (3.8) implies two equations

$$\left. \begin{array}{l} \frac{d^2X}{dx^2} = \mu X \\ \frac{d^2T}{dt^2} = \mu c^2 T \end{array} \right\} \left(\begin{array}{l} \text{Variables } x \text{ and } t \\ \text{have been separated} \end{array} \right) \quad (3.9)$$

and we have the critical simplifying step:

One PDE in two variables becomes two ordinary differential equations (ODEs) in single variables!

This separation is not always possible with other types of problems. Moreover, with some boundaries and some boundary conditions such a product form may not be possible either, even if the PDE allows it.

From this point on we need only to consider the solutions of the ordinary differential equations.

Step 2:

Consider possible values of μ . The boundary conditions (Equation (3.6d)) imply that

$$X(0) = X(L) = 0. \quad (3.10)$$

These are necessary for a unique solution.

Case 1: $\mu = 0$

Equation (3.9a) implies that $X'' = 0$, which is true if and only if $X(x) = c_1x + c_2$. But the boundary conditions in Equation (3.10) imply $c_1 = c_2 = 0$. That is, $X(x) \equiv 0$

— the trivial solution.

Case 2: $\mu > 0$, that is, $\mu = k^2 > 0$.

Equation (3.9a) implies that $X'' - k^2 X = 0 \iff X(x) = c_1 e^{kx} + c_2 e^{-kx}$.

But the boundary conditions in Equation (3.10) imply

$$\left. \begin{aligned} x = 0 : c_1 + c_2 &= 0 \\ x = L : c_1 e^{kL} + c_2 e^{-kL} &= 0 \end{aligned} \right\} \implies \begin{pmatrix} 1 & 1 \\ e^{kL} & e^{-kL} \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} = 0$$

$$\implies c_1 = c_2 = 0 \implies X(x) \equiv 0 \quad \text{— the trivial solution again.}$$

Case 3: $\mu < 0$, that is $\mu = -k^2 < 0$.

Equation (3.9a) implies that $X'' + k^2 X = 0$, which is true if and only if $X(x) = c_1 \cos kx + c_2 \sin kx$. And Equation (3.10) gives us

$$\begin{aligned} X(0) = 0 &\implies c_1 \cdot 1 = 0 \implies c_1 = 0 \\ X(L) = 0 &\implies c_2 \sin kL = 0 \implies c_2 = 0 \text{ or } \sin kL = 0 \end{aligned}$$

A nontrivial solution is possible only if $kL = \pi, 2\pi, 3\pi, \dots$, an integer multiple of π , and $c_2 \neq 0$.

The only nontrivial solutions are multiples of

$$X_n(x) = \sin\left(\frac{n\pi x}{L}\right), \quad n = 1, 2, 3, \dots$$

Definition 3.16

$\mu = -k^2 = -\left(\frac{\pi}{L}\right)^2, -\left(\frac{2\pi}{L}\right)^2, -\left(\frac{3\pi}{L}\right)^2, \dots$ are called **eigenvalues** or **characteristic values**.

$X_n(x) = \sin\left(\frac{n\pi x}{L}\right), \quad n = 1, 2, 3, \dots$ are called **eigenfunctions** or **characteristic functions**.

Step 3:

This choice of μ has implications for the solution of Equation (3.9b):

$$T''(t) + k^2 c^2 T(t) = 0, \quad k = \frac{\pi}{L}, \frac{2\pi}{L}, \frac{3\pi}{L}, \dots$$

means that for each possible k we get an independent solution. The coefficients of independent solutions will also depend on the values of k :

$$T_n(t) = A_n \cos\left(\frac{n\pi c t}{L}\right) + B_n \sin\left(\frac{n\pi c t}{L}\right), \quad n = 1, 2, 3, \dots$$

meaning that for each possible k we get a solution of (3.6a) satisfying (3.6d).

$$u_n(x, t) = \left(A_n \cos\left(\frac{n\pi c t}{L}\right) + B_n \sin\left(\frac{n\pi c t}{L}\right) \right) \sin\left(\frac{n\pi x}{L}\right), \quad n = 1, 2, 3, \dots$$

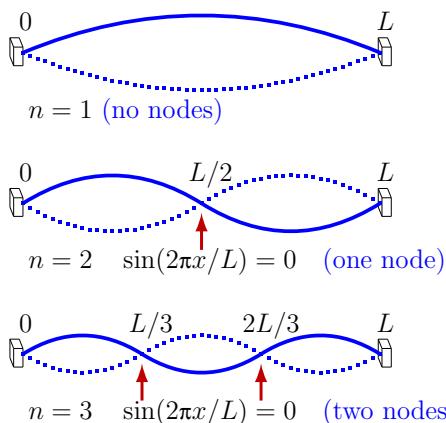
Since all the infinite n cases are possible solutions, called *harmonics*, the *principle of superposition* says that the most general solution is their linear combination. We have therefore an infinite series solution,

$$u = f(x, t) = \sum_{n=1}^{\infty} \sin\left(\frac{n\pi x}{L}\right) (A_n \cos \omega_n t + B_n \sin \omega_n t),$$

which describes the most general vibration of the ideal string fixed at its ends.

Remarks

- * The boundary conditions determine the *allowed* k 's and ω 's.
- * $k_n = \frac{n\pi}{L} = \frac{2\pi}{\lambda_n}$, where the λ_n are the allowed *wavelengths*.
- * $\omega_n = \frac{n\pi c}{L}$ are the allowed angular *frequencies* — *eigenfrequencies*.
- * $\omega_n = \frac{n\pi c}{L} = n\omega$. That is, ω_n is an integer multiple of the *fundamental frequency* $\omega = \frac{\pi c}{L}$ — the n^{th} *harmonic overtone*.
- * $u_n(x, t) = (A_n \cos \omega_n t + B_n \sin \omega_n t) \sin\left(\frac{n\pi x}{L}\right)$ describe the n^{th} *normal mode* of vibration, the first few examples of which are shown in Figure 3.18.



The *fundamental mode*
— greatest amplitude
— always set in motion?

The *first harmonic*
— lower amplitude
— twice fund. frequency
— not present when ...?

The *second harmonic*
— lower amplitude still
— thrice fund. frequency
— not present when ...?

Figure 3.18 The first few modes of vibration.

Remarks These are some points of a physical nature to note at this stage.

- * If the string is pulled aside (plucked) anywhere along its length then the fundamental mode, shown at the top of Figure 3.18, will be present in the solution, since the fundamental mode features displacement at all points.
- * On the other hand, a given harmonic is not present in the solution when the plucked point coincides with any of that harmonic's nodes. For example, if the string were to be plucked at its centre, the implied displacement would be inconsistent with that point being a node—a point of no motion. Consequently, the first harmonic mode (and any others that feature a node at the centre) cannot be included in the series expansion.

$$* \frac{\omega_n}{k_n} = \frac{n\pi c}{L} \cdot \frac{L}{n\pi} = c = \sqrt{\frac{\tau}{\rho}} \quad \text{— the same speed for all modes.}$$

Step 4:

To determine the unknown constants, $\{A_n\}_{n=1,2,\dots,\infty}$ and $\{B_n\}_{n=1,2,\dots,\infty}$, we apply the initial conditions of Equations (3.6b) and (3.6c). These imply that

$$f(x, 0) = \sum_{n=1}^{\infty} A_n \sin\left(\frac{n\pi x}{L}\right) \equiv g(x), \quad (3.11)$$

and

$$f_t(x, 0) = \sum_{n=1}^{\infty} B_n \cdot \omega_n \sin\left(\frac{n\pi x}{L}\right) = 0. \quad (3.12)$$

Equation (3.12) tells us that

$$B_n = 0 \quad \forall n,$$

and Equation (3.11) tells us that there exists a Fourier sine series for $g(x)$, with coefficients

$$A_n = \frac{2}{L} \int_0^L g(x) \sin\left(\frac{n\pi x}{L}\right) dx.$$

Given an initial form such as that represented in Figure 3.17, the A_n s can be determined and inserted in the series solution.

More general initial conditions of either plucked or struck strings are described by functions such as

$$\begin{cases} u(x, 0) = u_0(x) \text{ — string shape at } t = 0, \\ u_t(x, 0) = v_0(x) \text{ — string speed at } t = 0. \end{cases}$$

These conditions, or any linear combination of these conditions, are invoked to give $\{A_n, B_n\}$ in the more general case:

$$\begin{cases} A_n = \frac{2}{L} \int_0^L u_0(x) \sin\left(\frac{n\pi x}{L}\right) dx \\ B_n = \frac{2}{\omega_n L} \int_0^L v_0(x) \sin\left(\frac{n\pi x}{L}\right) dx \end{cases}$$

With all unknowns determined, the problem is solved.

Remark

- * The string system just considered is a good model for a string attached to an electric guitar. The electric guitar body is (usually) solid. It is therefore appropriate to assume the string ends are fixed. The body itself undergoes very little vibration of its own (if any) while the string is vibrating. Consequently, the notes registered by an electric guitar are almost as pure as those determined mathematically.

This differs fundamentally from the case of an acoustic guitar. The body of an acoustic guitar is hollow with the strings attached at one end to its flexible top plate. The vibrations of the strings are therefore transferred in part to the vibrations of the top plate (the ensuing air vibrations in the body are in fact responsible for a large proportion of the sound produced).

To describe the vibrations of a string on an acoustic guitar it is therefore more reasonable to adopt the model of one end of the string attached to moveable mass which is subject to an elastic restoring force. However, even in this case it can be shown that discrete vibration modes arise. Although related by a mathematical formula, the frequencies of the higher-order modes are not simple multiples of a fundamental frequency. (See Supplementary problem 21.)

Mastery Check 3.12:

For the problem of a string of length L and density ρ under tension τ , fixed at its ends and pulled aside a distance h at $x = b$, derive a closed-form expression

for the string's energy of vibration

$$E = \underbrace{\frac{1}{2}\rho \int_0^L \left(\frac{\partial u}{\partial t}\right)^2 dx}_{\text{kinetic}} + \underbrace{\frac{1}{2}\tau \int_0^L \left(\frac{\partial u}{\partial x}\right)^2 dx}_{\text{potential}},$$

using our separation of variables solution.

For the case $b = L/2$ show that the energy is conserved by comparing your result obtained above with the initial work done in pulling the string from its unstretched state.



☞ Mastery Check 3.13:

Using the separation of variables method, derive the unique solution to the mixed b.v.p.

$$\begin{cases} \text{(a)} \ u_t(x, t) = ku_{xx}(x, t) & 0 < x < a, \ t > 0 \\ \text{(b)} \ u(0, t) = 0 & t \geq 0 \\ \text{(c)} \ u(a, t) = 0 & t \geq 0 \\ \text{(d)} \ u(x, 0) = g(x) & 0 \leq x \leq a, \text{ where } g(0) = g(a) = 0. \end{cases}$$



3.F Supplementary problems

Section 3.A

1. Use a Taylor series approximation to find the nature of any critical points for the function

$$f(x, y) = (x + 2y)e^{-x^2-y^2},$$

for $\{(x, y) : x^2 + y^2 \leq 1\}$.

2. Locate and classify the critical points of the following functions:

(a) $f(x, y) = y^3 + 3x^2y - 3x^2 - 3y^2 + 2,$

(b) $f(x, y) = x^3 - y^3 - 2xy + 6,$

(c) $f(x, y) = e^{-x^2} (2xy + y^2).$

Section 3.B

3. Determine the maximum and minimum values of

$$f(x, y) = x^2 + y^2 - 3x,$$

in the region $|x| + |y| \leq 1$.

4. Show that in the region $|x| \leq 1, |y| \leq 1$ the function (Figure 3.19),

$$f(x, y) = e^{x^4-y^4},$$

has a stationary point which is a saddle point, then determine its maximum and minimum values in the region.

5. Determine the extreme points of the surface

$$z = \sin x \sin y \sin(x + y),$$

over the rectangular domain

$$D = \{(x, y) : 0 \leq x, y \leq \pi\}.$$

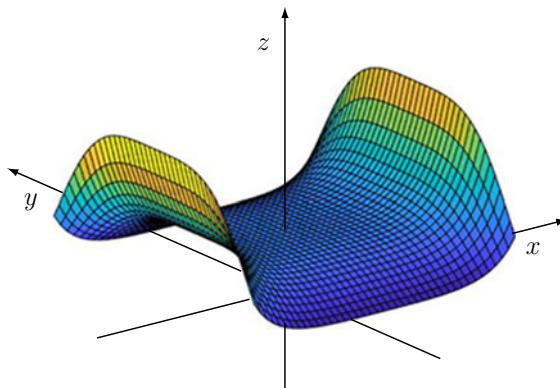


Figure 3.19 The graph of $z = e^{x^4 - y^4}$.

6. Determine the maximum and minimum values of

$$f(x, y) = xy \ln(1 + x^2 + y^2),$$

in the region $x^2 + y^2 \leq 1$.

7. Determine the maximum and minimum values of

$$f(x, y) = (x + y)e^{-x^2 - y^2},$$

(Figure 3.20) in the region $x^2 + y^2 \leq 1$.

8. Determine the maximum and minimum values of

$$f(x, y) = (x + y)e^{-x^2 - y^2},$$

(Figure 3.20) in the triangular region:

$$R = \{(x, y) : x \geq 0, y \geq 0, x + y \leq 2\}.$$

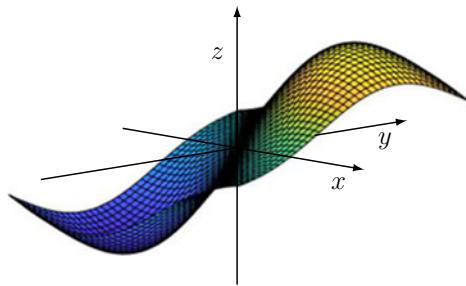


Figure 3.20 The graph of $z = (x + y)e^{-(x^2+y^2)}$.

9. Determine the maximum of the function

$$f(x, y, z) = \log x + \log y + 3 \log z,$$

over that portion of the sphere $x^2 + y^2 + z^2 = 5r^2$, in the first octant. Deduce that

$$abc^3 \leq 27 \left(\frac{a+b+c}{5} \right)^5.$$

10. Suppose C is the conic described by the equation

$$Ax^2 + 2Bxy + Cy^2 = 1,$$

where $A > 0$ and $B^2 < AC$. If we denote by p and P the distance from the origin to the closest and furthest point on the conic, respectively, show that

$$P^2 = \frac{A+C+\sqrt{(A-C)^2+4B^2}}{2(AC-B^2)},$$

with an analogous expression for p^2 .

Section 3.C

11. Surveyors have measured two sides and the angle between them of a triangular plot of land, for the purpose of finding the area of the plot. The area is given by $f = \frac{1}{2}ab \sin C$, where a and b are the lengths of the sides and C is the included angle. The measurements are all subject

to error. The measured values for the sides were $a = 152.60 \pm 0.005$ m and $b = 163.81 \pm 0.005$ m, and for the included angle $\theta = 43^\circ 26' \pm 0.2'$.

What is the largest possible error in the computation of the area? What is the largest possible percentage error?

Section 3.D

12. Suppose we wish to represent the function $f(x) = x^2$, $-1 \leq x \leq 1$, by a combination of cosine functions of the form

$$g(x; a, b, c, d) = a + b \cos \pi x + c \cos 2\pi x + d \cos 3\pi x.$$

Find the least-squares values of a, b, c, d . Plot the two functions on the same graph to check the fit.

Section 3.E

13. Show that $u(x, t) = \exp(-x^2/4t)/\sqrt{t}$ is a solution to the parabolic, 1D diffusion equation

$$\frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2}.$$

14. Show that $u(x, y, t) = t^{-1} \exp(-(x^2 + y^2)/4t)$ is a solution to the parabolic, 2D diffusion equation

$$\frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2}.$$

15. Suppose $u : \mathbb{R}^3 \rightarrow \mathbb{R}$ is a C^2 function of variables (x, y, z) . By introducing 3D cylindrical polar coordinates (Section 1.D) confirm the expression for the Laplacian shown in the centre of Figure 3.13, applied to $U(r, \theta, z) = u(r \cos \theta, r \sin \theta, z)$.

16. Suppose $u : \mathbb{R}^3 \rightarrow \mathbb{R}$ is a C^2 function of variables (x, y, z) . By introducing 3D spherical polar coordinates (Section 1.D) confirm the expression for the Laplacian shown in the right-hand side of Figure 3.13, applied to $U(\rho, \phi, \theta) = u(\rho \sin \phi \cos \theta, \rho \sin \phi \sin \theta, \rho \cos \phi)$.

17. By introducing new variables $\xi = x + y$ and $\eta = \ln(x - y)$ (such that $x - y > 0$) and using the chain rule, determine all C^2 functions that solve the differential equation

$$(x - y) \left(\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} \right) = 4(x + y)^2.$$

18. Suppose $u : \mathbb{R}^2 \rightarrow \mathbb{R}$ is a C^2 function of variables (x, t) . With the help of the change of variables $\xi = x + ct$ and $\eta = x - ct$, transform the hyperbolic, 1D wave equation

$$\frac{\partial^2 u}{\partial x^2} - \frac{1}{c^2} \frac{\partial^2 u}{\partial t^2} = 0$$

into a simpler form involving a function $U(\xi, \eta)$ and thereby determine all possible C^2 functions that solve the wave equation.

19. By introducing the change of variables $\xi = x + e^y$ and $\eta = x - e^y$, determine the most general differentiable solution to the partial differential equation

$$e^{2y} \frac{\partial^2 u}{\partial x^2} - \frac{\partial^2 u}{\partial y^2} + \frac{\partial u}{\partial y} = 0.$$

- 20 Use the method of separation of variables to find the general solution to the equation

$$x^2 u_{xx} + x u_x - u_y = 0.$$

Find the particular solution that satisfies the boundary conditions $u(1, 0) = 0$, $u(e, 0) = 2$.

Hint: At some stage you will need to substitute $x = e^v$.

- 21 Redo the 1D wave (vibrating string) problem for the initial and boundary condition sets

(a) $(h_1, h_2, g_1, g_2) = (0, 0, 0, g(x))$, $0 \leq x \leq 2$;

(b) initial conditions $(g_1, g_2) = (g(x), 0)$, and boundary conditions

$$\left(f(0, t) = 0, M \frac{\partial^2 f}{\partial t^2}(L, t) = -\tau \frac{\partial f}{\partial x}(L, t) - \kappa f(L, t) \right) \text{ for } M, \kappa > 0.$$

Case (b) models a string of length L under tension τ attached at $x = L$ to a mass M subject to an elastic restoring force of strength κ .



Chapter 4

Integration of multivariable functions

We have seen that differentiation of a function is the result of a limit process. This led to the definition of a tangent line to a curve in 2D and a tangent plane to a surface in 3D. In this chapter we shall again consider limit processes but with the aim of establishing the reverse operation, that of integration of a scalar function of many variables.

It should be emphasized that we focus attention here on integrals of functions over subsets of \mathbb{R}^n . That is, the integrals we consider are taken over zero curvature regions in \mathbb{R}^n : a straight line segment in \mathbb{R} , a planar region in \mathbb{R}^2 , a volume subset in \mathbb{R}^3 , etc. In this important context we shall be able to rely on some familiar concepts derived from geometry such as areas (under curves) and volumes (under surfaces) to assist our understanding. In Chapter 5, we will revisit 1D and 2D integration but in the sense of integrals over nonzero curvature geometries (curves and curved surfaces in \mathbb{R}^3). In that context, geometric interpretations will be replaced with more physical conceptions.

4.A Multiple integrals

As we've done before, we shall first revisit the 1D case. To be precise, we shall brush up on the Riemann theory of integration for a function of one variable. Readers interested in the more general measure theory of Riemann-Stieltjes integration are referred to [2]. A comparison between the foundations of single-variable integration and multivariable integration is particularly fruitful, since the latter case is essentially a direct generalization of the former case.

Almost every idea and theoretical result we discuss in this chapter is as valid for functions of an arbitrary number of variables as for functions of two variables. Therefore to simplify matters, we will present the theory (mostly) with reference to functions of *two* variables and discuss the more general case at the end of the chapter. Any important practical differences will be highlighted in that discussion.

As already mentioned, the definition of the integral of a function f of $\mathbf{x} \in \mathbb{R}^n$ is based on a limit process. We illustrate the first steps of the process in Figure 4.1.

Integration of $f : I \subset \mathbb{R} \rightarrow \mathbb{R}$

Suppose f is a continuous function of x and assume that the interval I is closed and bounded and lying in the function domain, D_f . That is,

$$I = \{x : a \leq x \leq b\} \subset D_f.$$

The graph of f is a curve in $\mathbb{R} \times \mathbb{R} = \mathbb{R}^2$ as shown in Figure 4.1 below.

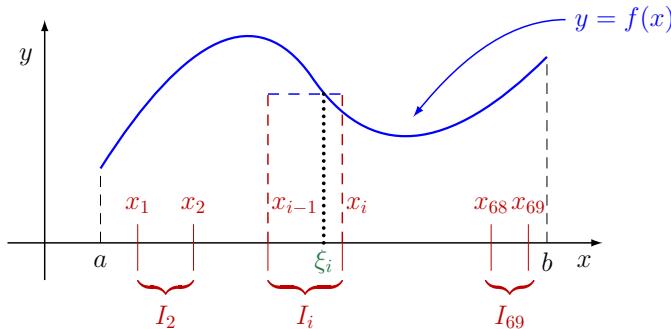


Figure 4.1 The graph of f and some subintervals of I .

First the interval I is cut I into small bits — this is called a *partition* of I :

$$a = x_0 < x_1 < x_2 < \dots < x_{n-1} < x_n = b$$

with n subintervals $I_i = \{x : x_{i-1} \leq x \leq x_i\}$, of width $\Delta x_i = x_i - x_{i-1}$. A few of these subintervals are shown in Figure 4.1.

Then, choosing some real number $\xi_i \in I_i$ from each subinterval, we form the sum

$$\sigma_n = \sum_{i=1}^n f(\xi_i) \Delta x_i.$$

This is called the *Riemann sum of f over I* . From its construction we see that it *must* depend on the partition of n subintervals.

The geometric interpretation of σ_n for $f : I \rightarrow \mathbb{R}$.

If $f \geq 0$, then $f(\xi_i)\Delta x_i$ is the *area* of the rectangle of height $f(\xi_i)$ and width Δx_i as shown here in Figure 4.2.

Hence, the sum σ_n is an approximation to the area “under” the curve $y = f(x)$ and over I .

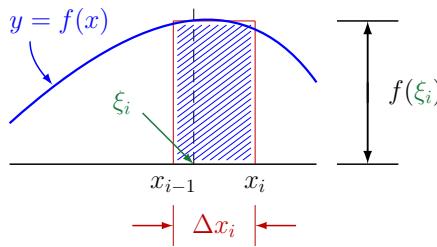


Figure 4.2 Rectangular area approximation.

To improve on this approximation we find numbers ℓ_i and m_i in each interval I_i such that $f(\ell_i) \leq f(x) \leq f(m_i)$ for all x in I_i .

For a given partition we form the upper and lower sums

$$R_{\min} = \sum_{i=1}^n f(\ell_i)\Delta x_i \leq \sum_{i=1}^n f(\xi_i)\Delta x_i \leq \sum_{i=1}^n f(m_i)\Delta x_i = R_{\max}.$$

In this process we have constructed upper and lower bounds on σ_n . That is,

$$R_{\min} \leq \sigma_n \leq R_{\max}.$$

We now take the simultaneous limit of the number of intervals $n \rightarrow \infty$ and the representative size of the intervals $\max(\Delta x_i) \rightarrow 0$. We find that, as $n \rightarrow \infty$, R_{\min} increases and R_{\max} decreases. If the dual limits exist and $\lim R_{\min} = \lim R_{\max}$, then an application of a squeeze theorem gives:

Definition 4.1

The integral of $f : \mathbb{R} \rightarrow \mathbb{R}$ over I is defined as the limit (if it exists)

$$\lim_{\substack{n \rightarrow \infty \\ \max(\Delta x_i) \rightarrow 0}} \sigma_n = \int_a^b f(x) \, dx.$$

Under the conditions of continuity of f and compactness of I this limit does exist. This is a unique number, called the *Riemann integral* ([1], [2]), which is independent of how we set the partition originally.

Now let's see how things work for a function $f(x, y)$ of two variables.

Integration of $f : R \subset \mathbb{R}^2 \rightarrow \mathbb{R}$

Suppose f is a continuous function of x and y and, for starters, we are given the closed and bounded rectangular region

$$R = \{(x, y) : a \leq x \leq b, c \leq y \leq d\}$$

which lies inside D_f . Note that $R \subset D_f$ and R is compact.

The *graph* of such a function f is a *surface* in $\mathbb{R}^2 \times \mathbb{R} = \mathbb{R}^3$ as shown in Figure 4.3 below.

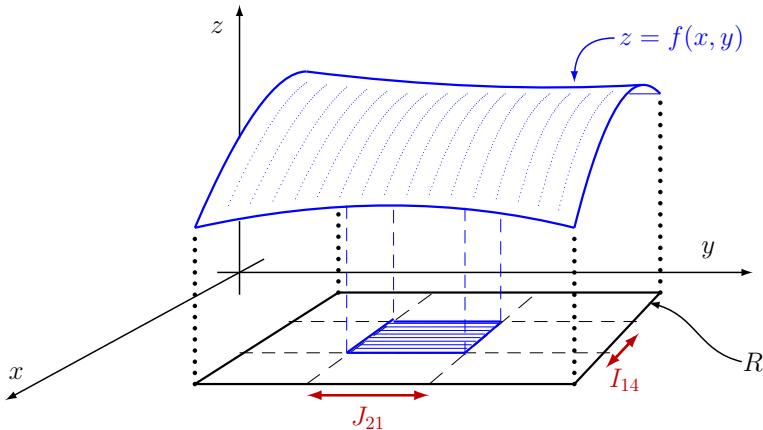


Figure 4.3 The graph of f over rectangle R , and subrectangle $I_{14} \times J_{21}$.

The construction of the integral of f over R is accomplished in perfect harmony with the 1D case except that rectangles replace intervals.

First, the rectangle R is partitioned into $n \times m$ rectangles $R_{ij} = I_i \times J_j$ of area ΔA_{ij} , $1 \leq i \leq n, 1 \leq j \leq m$. We then choose some point (ξ_i, η_j) from each rectangle R_{ij} and form the sum

$$\sigma_{nm} = \sum_{i=1}^n \sum_{j=1}^m f(\xi_i, \eta_j) \Delta A_{ij}.$$

This is called the *Riemann sum of f over R* . From its construction this Riemann sum *must* therefore depend on the partition.

The geometric interpretation of σ_{nm} for $f : R \subset \mathbb{R}^2 \rightarrow \mathbb{R}$.

If $f \geq 0$, then $f(\xi_i, \eta_j)\Delta A_{ij}$ is the *volume* of the rectangular block of height $f(\xi_i, \eta_j)$ and base area $\Delta A_{ij} = \Delta x_i \cdot \Delta y_j$ as shown in Figure 4.4 below.

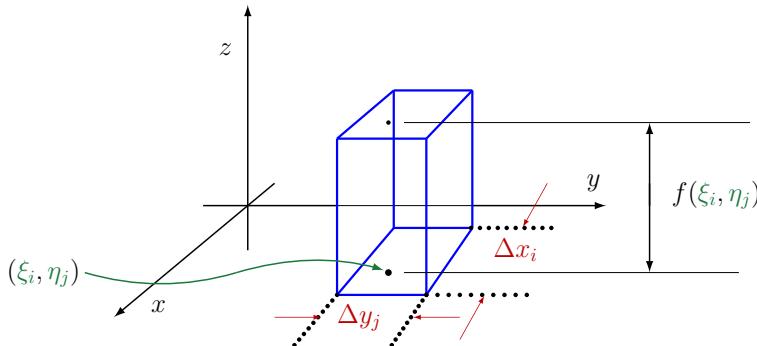


Figure 4.4 Rectangular block volume approximation.

Thus, the sum σ_{nm} is an approximation to the volume “under” the surface $y = f(x, y)$ and over R .

A completely analogous line of reasoning to the steps leading up to the definition of a 1D integral can now be applied as follows.

For a given partition we form upper and lower sums of the rectangular blocks. Between these two sums is the sum σ_{nm} . With each refinement of the partition, new upper and lower sums are determined with the new lower sum increased compared with its predecessor and the new upper sum decreased compared with its predecessor.

If in the dual limit process of $n, m \rightarrow \infty$ and $\max \left(\sqrt{(\Delta x_i)^2 + (\Delta y_j)^2} \right) \rightarrow 0$ the limits of upper and lower sums exist and are equal, an application of the squeeze theorem leads to our next definition.

Definition 4.2

The double integral of $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ over R is defined as the limit (if it exists)

$$\lim_{\substack{m,n \rightarrow \infty \\ \max \left(\sqrt{(\Delta x_i)^2 + (\Delta y_j)^2} \right) \rightarrow 0}} \sigma_{nm} = \iint_R f(x) \, dA.$$

An important theorem and some corollaries — 1D version.

Theorem 4.1 *All continuous functions are integrable on compact subsets of their domains of definition.*

Corollary 4.1.1 *If $f \geq 0$, then $\int_I f(x) dx$ is the area under the curve $y = f(x)$.*

Corollary 4.1.2 *If $f \geq g \geq 0$, then $\int_I (f(x) - g(x)) dx$ is the area between the curves (Figure 4.5).*

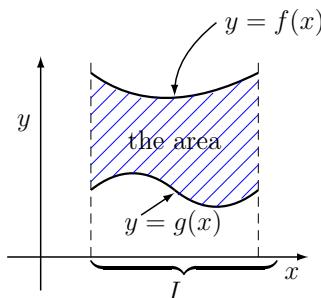


Figure 4.5 The area between two curves in 2D.

Corollary 4.1.3 $\int_I 1 dx$ is the length of interval I .

Corollary 4.1.4 *The average value of function f over I is*

$$\frac{1}{\text{length } I} \cdot \int_I f(x) dx = \frac{\int_I f(x) dx}{\int_I 1 dx}.$$

Corollary 4.1.5 *Linearity: If $a, b \in \mathbb{R}$ then*

$$\int_I (af(x) + bg(x)) dx = a \int_I f(x) dx + b \int_I g(x) dx.$$

Corollary 4.1.6 *Additivity: (very important)*

If $I \cap J = \{\}$ = \emptyset , then $\int_{I \cup J} f(x) dx = \int_I f(x) dx + \int_J f(x) dx$.

An important theorem and some corollaries — 2D version.

Theorem 4.1 All continuous functions are integrable on compact subsets of their domains of definition.

Corollary 4.1.1 If $f \geq 0$, then $\iint_R f(x, y) dA$ is the volume under the surface $z = f(x, y)$.

Corollary 4.1.2 If $f \geq g \geq 0$, then $\iint_R (f(x, y) - g(x, y)) dA$ is the volume of the body between the surfaces (Figure 4.6).

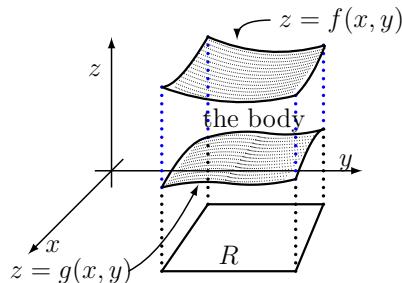


Figure 4.6 The volume between two surfaces in 3D.

Corollary 4.1.3 $\iint_R 1 dA$ is the area of R . — true not just for rectangles!!

Corollary 4.1.4

$$\frac{1}{\text{area } R} \cdot \iint_R f(x, y) dA = \frac{\iint_R f(x, y) dA}{\iint_R 1 dA} \text{ is the average value of } f \text{ over } R.$$

Corollary 4.1.5 Linearity: If $a, b \in \mathbb{R}$ then

$$\iint_R (af(x, y) + bg(x, y)) dA = a \iint_R f(x, y) dA + b \iint_R g(x, y) dA.$$

Corollary 4.1.6 Additivity: (very important)

If $R_1 \cap R_2 = \{\} \equiv \emptyset$, then

$$\iint_{R_1 \cup R_2} f(x, y) dA = \iint_{R_1} f(x, y) dA + \iint_{R_2} f(x, y) dA.$$

Some comments on Theorem 4.1 are warranted.

From Theorem 1.2 all continuous functions on compact domains are bounded. These functions do not exhibit singular behaviour and so in the integral definition all partitions have finite volume (area). The unique limits in Definitions 4.1 and 4.2 therefore exist.

Corollary 4.1.3 may seem trivial but is undeniably useful in some contexts. (See Mastery Check 5.24, Section 5.D, and Example 5.9.)

Corollary 4.1.5 allows us to split complex functions into sums of simpler functions and to integrate them individually.

Corollary 4.1.6 is useful when an integration domain is or can be described piecewise, especially if the pieces warrant different techniques of integration.

4.B Iterated integration in \mathbb{R}^2

Iterated integration is the workhorse of multiple integrals.

The definition of the multiple integral as the limit of a sum is not practical. Fortunately, there is an alternative. The suggestion is made that we calculate our “volumes” by *slicing* rather than by *dicing*.

Consider the thin slice of the “body” under f shown in Figure 4.7. The area of the left-hand side face, that is, the area under the curve of constant y , $y = y_0$, is $A(y_0) = \int_a^b f(x, y_0) dx$. Similarly, $A(y_0 + \Delta y) = \int_a^b f(x, y_0 + \Delta y) dx$ is the area of the right-hand side face.

If $|\Delta y|$ is a small increment then $A(y_0) \approx A(y_0 + \Delta y)$, which is easy to see by expanding $f(x, y_0 + \Delta y)$ in a Taylor series about (x, y_0) . Then, using the simple two-point trapezoidal rule approximation, the volume of the “slice” is approximately

$$\begin{aligned} V(y_0) &= \frac{1}{2} (A(y_0) + A(y_0 + \Delta y)) \Delta y \\ &= A(y_0) \Delta y + O(\Delta y^2). \end{aligned}$$

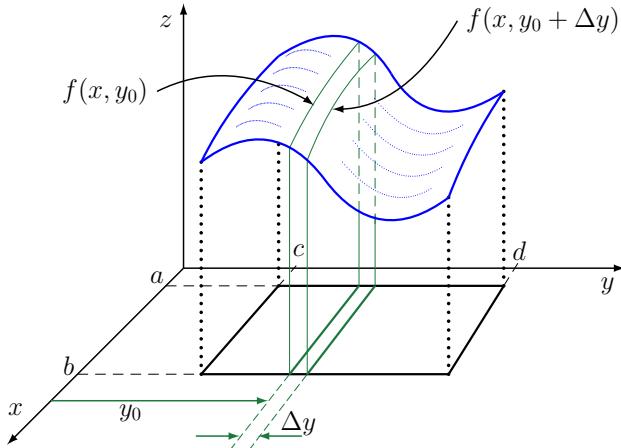


Figure 4.7 Adding slices to determine volumes.

The total volume of the “body” under $f(x, y)$ is then the limiting sum of these 1D volumes of slices (Definition 4.1) as $\Delta y \rightarrow 0$. That is, the volume under $f(x, y)$ over R is the Riemann integral of $A(y)$ over the interval $c \leq y \leq d$:

$$V = \int_c^d A(y) dy = \int_c^d \left(\int_a^b f(x, y) dx \right) dy.$$

Alternatively, slicing parallel to the y -axis instead of the above would give

$$V = \int_a^b A(x) dx = \int_a^b \left(\int_c^d f(x, y) dy \right) dx,$$

which must give the exact same value for the volume.

Hence, for integration over the rectangle $[a, b] \times [c, d]$ we have the important result

$$\underbrace{\iint_R f(x, y) dA}_{\text{double integral of } f \text{ over } R} = \underbrace{\int_a^b \left(\int_c^d f(x, y) dy \right) dx}_{\text{iterated integrals of } f \text{ over } R} = \int_c^d \left(\int_a^b f(x, y) dx \right) dy$$

The left-hand side is the definition of a double integral, while the two right-hand sides are the actual ways one can evaluate the double integral. These

are called the *iterated* integrals of f . In each case of iterated integral, the *inner* integral within the parentheses is to be evaluated first, only then is the *outer* integral evaluated. Notice how in the above equation the integration limits follow the variable being integrated.

At a very practical level, when evaluating the inner integral we treat the *outer integral variable* as if it were a constant parameter. In such a situation all single-variable techniques apply to each individual iterated integral. Example 4.1 illustrates this process.

Regarding notation, the above clearly specifies the order of operation. However, to skip the parentheses (and avoid the tedium of writing them) we have two alternative notations in common use:

$$\int_a^b \int_c^d f(x, y) \, dy \, dx \quad \text{— this borders on the ambiguous; the user must not confuse the order.}$$

$$\int_a^b dx \int_c^d f(x, y) \, dy \quad \text{— this is somewhat better; it is easier to interpret and better for complex regions (see next section).}$$

■ Example 4.1:

Determine the volume of the body lying under the surface $z = x^2 + y^2$ (Figure 4.8) and over the rectangle,

$$R = \{(x, y) : 0 \leq x \leq 1, 0 \leq y \leq 1\}, \quad (= [0, 1] \times [0, 1]).$$

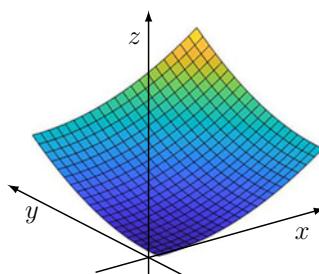


Figure 4.8 The graph of $z = x^2 + y^2$.

Solution:

$$\begin{aligned}
 V &= \iint_R (x^2 + y^2) \, dA \quad \text{--- the double integral} \\
 &= \int_0^1 dy \int_0^1 (x^2 + y^2) \, dx \quad \text{--- an iterated integral version} \\
 &= \int_0^1 dy \left[\frac{x^3}{3} + y^2 x \right]_{x=0}^1 \quad \text{--- } y \text{ is held "fixed" during the } x\text{-integration} \\
 &= \int_0^1 \left(\frac{1}{3} + y^2 \right) dy = \left[\frac{y}{3} + \frac{y^3}{3} \right]_0^1 = \frac{2}{3} \text{ volume units.}
 \end{aligned}$$

In this iterated integral y is given the role of the outer integration variable, while x plays the role of inner variable. We could just as easily have reversed the roles to arrive at the same result. Notice the very important fact, which will be contrasted with later, that in this example the bounds on the x -integral are constants; they are *not* functions of y !



Mastery Check 4.1:

Evaluate $\iint_R ye^{xy} \, dA$, where $R = \{(x, y) : 0 \leq x \leq 1, 0 \leq y \leq 2\}$.



Mastery Check 4.2:

Compute $\iint_R \frac{xe^{x\sqrt{y}}}{\sqrt{y}} \, dA$, where $R = \{(x, y) : 0 < x < 1, 0 < y < 1\}$.



4.C Integration over complex domains

In general practice a region of integration is more often *not* a rectangle. We therefore need a way to work with compact regions that are more complex. Luckily, we can treat this problem using our previous results, but after first rethinking the function, the domain, and the iteration definitions.

First, let's look at the function and its domain. For a function f defined on a non-rectangular, compact domain D , we introduce a new function and new domain through the following definition.

Definition 4.3

Suppose $f(x, y)$ is continuous over a compact domain $D \subseteq D_f$. Let \hat{f} be the extension of f beyond D :

$$\hat{f}(x, y) = \begin{cases} f(x, y) & \text{for } (x, y) \in D \\ 0 & \text{for } (x, y) \notin D. \end{cases}$$

Since D is bounded there exists a rectangle R such that $D \subset R$. Thus, if \hat{f} is integrable over R , then

$$\underbrace{\iint_D f(x, y) dA}_{\text{definition}} = \underbrace{\iint_R \hat{f}(x, y) dA}_{\text{calculable value}}$$

This last equation is true since all we have done is added zero to the original double integral. The picture we imagine is something like Figure 4.9 below.

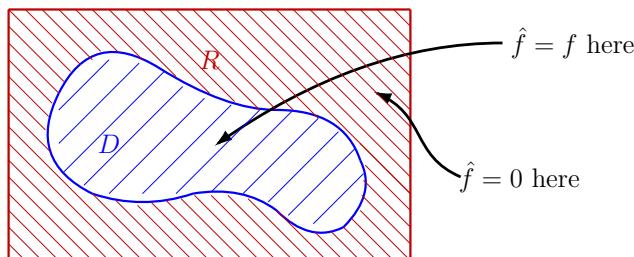


Figure 4.9 The extended function \hat{f} and its domain R .

Second, we examine the domain, D , a little further. What can these more complicated domains look like? In Figures 4.10–4.12 we define three main classes of regions into one class of which the domain D , or a piece of D , may be placed.

Type I:

Suppose the domain D is of the kind shown in Figure 4.10 and defined as

$$D = \{(x, y) : a \leq x \leq b, g_1(x) \leq y \leq g_2(x)\}.$$

This is called a *y-simple domain*, with the variable y bounded by two functions of x .

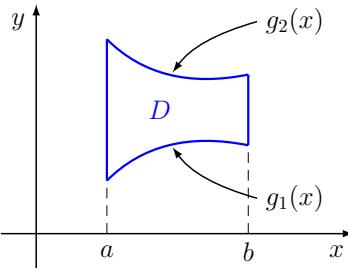


Figure 4.10 A *y*-simple domain.

Type II:

Suppose the domain D is of the kind shown in Figure 4.11 and defined as

$$D = \{(x, y) : c \leq y \leq d, h_1(y) \leq x \leq h_2(y)\}.$$

This is called a *x-simple domain*, with the variable x bounded by two functions of y .

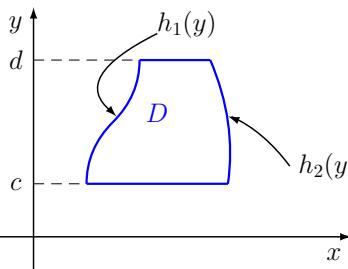


Figure 4.11 An *x*-simple domain.

Type III:

This is a domain D with an appearance like that shown in Figure 4.12, and with enough flexibility to have two interpretations. It could either be treated as a special case of **I** with $g_1(a) = g_2(a)$ and $g_1(b) = g_2(b)$, or as a special case of **II** with $h_1(c) = h_2(c)$ and $h_1(d) = h_2(d)$. That is, this sort of domain could be either *x*-simple or *y*-simple depending on how it is described.

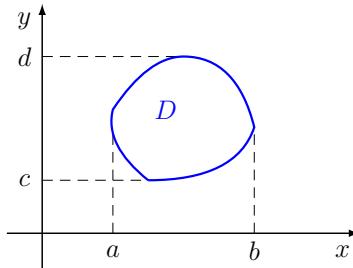


Figure 4.12 A domain that is both x -simple and y -simple.

Third, we bring these ideas together to arrive at a strategy for evaluating integrals over non-rectangular domains. We demonstrate this with a *y-simple domain* (Fig. 4.13).

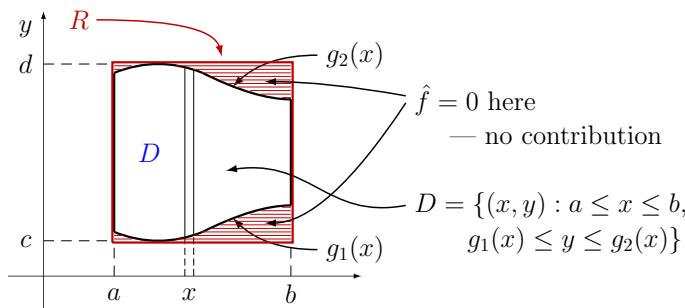


Figure 4.13 A y -simple domain in a rectangle R .

By construction we have the iterated integral of \hat{f} over $R = [a, b] \times [c, d] \supset D$,

$$\iint_D f(x, y) dA = \iint_R \hat{f}(x, y) dA = \int_a^b dx \int_c^d \hat{f}(x, y) dy.$$

However, for every value of the outer integral variable x , $\hat{f} = 0$ outside the interval $g_1(x) \leq y \leq g_2(x)$, and $\hat{f} = f$ in the interior of that interval. Hence,

$$\iint_D f(x, y) dA = \int_a^b dx \underbrace{\int_{g_1(x)}^{g_2(x)} f(x, y) dy}_{\text{iterated integral of } f \text{ over } D}.$$

We can now invoke the contrast alluded to earlier regarding the variable dependence of the limits of the inner integral. For all cases of non-rectangular domains, the limits of the inner integral will be functions of the outer integral variable. In the above example, the limits on the inner integral depend on x

and do *not* depend on y . Only in the case of rectangular domains will the limits of both the inner and outer integrals be constants!

For an x -simple domain, with $D = \{(x, y) : c \leq y \leq d, h_1(y) \leq x \leq h_2(y)\}$, we get the analogous result

$$\iint_D f(x, y) \, dA = \int_c^d dy \int_{h_1(y)}^{h_2(y)} f(x, y) \, dx.$$

Once again, the limits of the inner integral depend on the outer integral variable; the limits of the inner integral here depend on y and do *not* depend on x .

Through a very simple development we have arrived at very natural generalizations of iterated integrals over rectangles. Moreover, in the process we have done away with the extensions we used in this development.

The reader should now bear two things in mind. First, the order in which the iterated integrals are to be performed must be strictly adhered to. Second, interchanging the order will *always* involve a change in the limits. This is illustrated in Example 4.2 wherein a double integral is evaluated in two ways. The reader should note the limits on the two inner integrals and how they come about (see the vertical and horizontal bars in Figure 4.14).

■ Example 4.2:

Suppose D is that region bounded by the lines $y = x$, $y = 2x$, and $x = 1$.

Calculate the area of D as $\iint_D 1 \, dA$.

It cannot be stressed enough that one should always start solving a multiple integral problem by drawing a picture of the *region* in question. Here it is the one shown in Figure 4.14. Sketching the region of integration, if it is possible, not only allows us to get some perspective on the task involved, it also allows us to determine the limits of the integration variables. More importantly, it can potentially identify difficulties with setting limits, which formulae themselves may not do. An example of such complication arises when we treat D as x -simple.

The domain is both y -simple and piecewise x -simple.

We first treat D as y -simple.

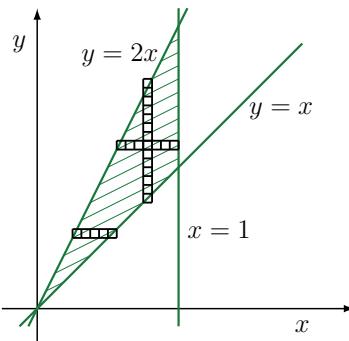


Figure 4.14 The y -simple domain D .

$$\begin{aligned}\text{Area of } D &= \iint_D 1 \, dA = \int_0^1 dx \int_x^{2x} 1 \, dy \\ &= \int_0^1 dx \left[y \right]_x^{2x} = \int_0^1 x \, dx = \left[\frac{x^2}{2} \right]_0^1 = \frac{1}{2} \text{ area units.}\end{aligned}$$

Now we treat D as x -simple. This is slightly more complex as Figure 4.14 suggests. We must break D up into two non-overlapping domains and invoke Corollary 4.1.6.

$$\begin{aligned}\text{Area of } D &= \iint_D 1 \, dA = \int_0^1 dy \int_{y/2}^y 1 \, dx + \int_1^2 dy \int_{y/2}^1 1 \, dx \\ &= \int_0^1 dy \left[x \right]_{y/2}^y + \int_1^2 dy \left[x \right]_{y/2}^1 \\ &= \int_0^1 \frac{y}{2} \, dy + \int_1^2 \left(1 - \frac{y}{2} \right) \, dy = \left[\frac{y^2}{4} \right]_0^1 + \left[y - \frac{y^2}{4} \right]_1^2 \\ &= \frac{1}{4} + (2 - 1) - \left(1 - \frac{1}{4} \right) = \frac{1}{2} \text{ area units.}\end{aligned}$$

■

☞ Mastery Check 4.3:

Determine $\iint_D (xy + y^2) \, dA$, over each of the following domains D .

(a) $D = \{(x, y) : 0 \leq x \leq 1, 0 \leq y \leq x^2\}$.

- (b) $D = \{(x, y) : 0 \leq x \leq \sqrt{y}, 0 \leq y \leq 1\}.$



Mastery Check 4.4:

Evaluate the iterated integral $I = \int_0^1 dx \int_{\sqrt{x}}^1 e^{y^3} dy.$

Hint: If you have trouble with this iterated integral, try thinking about it first as a double integral over some domain $D.$



4.D Generalized (improper) integrals in \mathbb{R}^2

The iterated integral approach only works if the corresponding multiple integral exists (that is, the limit of the sum σ_{mn} exists). It is very important to remember that thus far we have relied on the convenient assumptions of the function being continuous and the domain being bounded. However, in other cases the Riemann theory of integration on which the double integral was founded can break down. This leads us to consider so-called *improper integrals*, of which there are two types:

- (a) One type involves *unbounded domains*, e.g. $D = \{(x, y) : x > 0, y > 0\}.$
- (b) One type involves functions f , *not defined on part of the boundary*.
For example, $f(x, y) = 1/(x^2 + y^2)$ is not defined at $(0, 0).$

Under these more general circumstances, to answer the question

“Does $\iint_D f(x, y) dA$ exist?”,

we can rely somewhat on the combined action of the following theorems.

Theorem 4.2

If ...

$$\left. \begin{aligned} & \int dx \int |f(x, y)| dy \\ & \text{and } \int dy \int |f(x, y)| dx \end{aligned} \right\} \begin{aligned} & \text{exist} \\ & \text{and are} \\ & \text{equal} \end{aligned} \quad \text{then ...} \quad \iint_D |f(x, y)| dA \text{ exists.}$$

Theorem 4.3

$$\text{If } \iint_D |f(x, y)| \, dA \text{ exists, then } \iint_D f(x, y) \, dA \text{ exists.}$$

These two theorems may seem universally useful, and they are when any two iterated integrals *do not* give the same result, or if an iterated integral fails to converge. We would know then that the multiple integral does not exist. However, if only one ordered iterated integral can be evaluated, and it converges, and we cannot evaluate the other iterated integral, then there may still remain some doubt about the existence of the multiple integral.

Improper integrals in analysis

Nevertheless, guided by the above theorems we jump straight to the task of evaluating the iterated integral forms of a generalized multiple integral. In tackling improper iterated integrals we take advantage of the wisdom of single-variable calculus.

In single-variable calculus, the improper definite integral $\int_I f(x) \, dx$, for $f(x) > 0$ over the domain I , exists if

$\lim_{\epsilon \rightarrow 0^+} J(\epsilon) = \lim_{\epsilon \rightarrow 0^+} \int_{a+\epsilon}^b f(x) \, dx$, $a < b$ exists, when f is a function which diverges at the integration limit $x = a \in \bar{I}$ (Figure 4.15(a)); or

$\lim_{R \rightarrow \infty} J(R) = \lim_{R \rightarrow \infty} \int_a^R f(x) \, dx$ exists, when I is the unbounded domain, $I = [a, \infty)$ (Figure 4.15(b)).

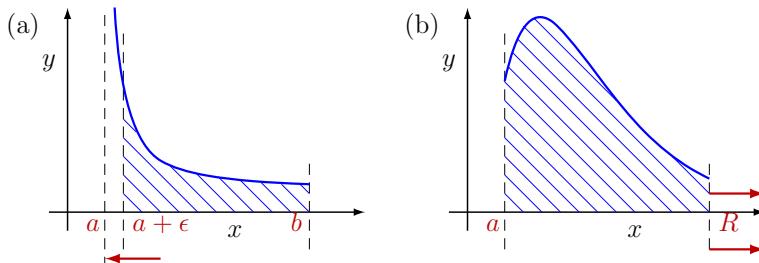


Figure 4.15 Two types of improper integral in single-variable calculus.

In Figure 4.15(a) the function f is singular at the lower limit $x = a$, while in Figure 4.15(b) the domain I is unbounded.

We notice that both cases rely on the Riemann theory of integration of a *continuous* function defined on a *bounded* sub-domain to give a finite number $J(\epsilon)$ and $J(R)$, respectively. And both cases subsequently test the convergence of limits, the first $\lim_{\epsilon \rightarrow 0} J(\epsilon)$ and the second $\lim_{R \rightarrow \infty} J(R)$, respectively, to define and provide the integrals wanted.

In multivariable calculus the improper multiple integral

$$\iint_D f(x, y) dA \quad (\text{for } f(x, y) > 0)$$

is similarly to be identified as a suitable limit. In practice we work with the iterated integral version of this double integral. However, the multiple integral version of the principle is more easily described using the double integral itself.

The student reader will no doubt note the similarity between the two scenarios. The common denominator is the integration domain D and in each case the sequence of smaller domains that do not present any difficulty. Note that the arguments below are valid and can substantiated *only* for the case of functions $f(x, y)$ that do not change sign over the integration region.

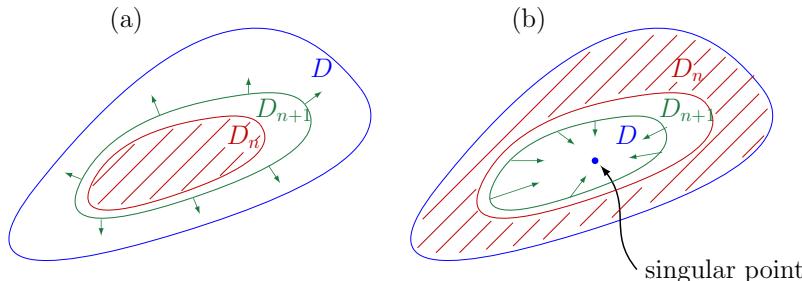


Figure 4.16 (a) Finite domains converging to unbounded D ;
 (b) Regular domains converging to singular D .

Suppose D is an *unbounded* domain, but contains *bounded* subsets, D_n , represented in Figure 4.16(a) such that

$$D_n \subset D_{n+1} \subset D \quad \forall n \text{ and } \bigcup_{n=1}^{\infty} D_n = D.$$

Every point in D belongs to at least one of the D_n .

On the other hand, suppose D is a domain containing a *singular point* of f , but which contains *bounded* subsets, D_n , as illustrated in Figure 4.16(b),

that *exclude* that point.

$$D_n \subset D_{n+1} \subset D \ \forall n \text{ and } \bigcup_{n=1}^{\infty} D_n = D.$$

Every point in D belongs to at least one of the D_n .

In either case we have the following useful result.

Theorem 4.4

Let f be a non-negative or non-positive function defined on a domain, D , which may be unbounded or contain points on which f is undefined. For the sequence of bounded subsets, D_n , satisfying the conditions $D_n \subset D_{n+1} \subset D$ and $\bigcup_{n=1}^{\infty} D_n = D$, if

$$\lim_{n \rightarrow \infty} J(n) = \lim_{n \rightarrow \infty} \iint_{D_n} f(x, y) \, dA$$

exists, then the improper integral $\iint_D f(x, y) \, dA$ exists and is equal to this limit.

An immediate corollary of this is the following.

Corollary 4.4.1 For the function and domain conditions of Theorem 4.4, if the improper integral $\iint_D f(x, y) \, dA$ exists, then the iterated integrals

$$\int dx \int f(x, y) \, dy \quad \text{and} \quad \int dy \int f(x, y) \, dx \quad \text{exist and are equal.}$$

Unfortunately, the only assertion that can be made for functions that change sign over the integral domain is Theorem 4.3. (See Mastery Check 4.8.) As we said, in practice we work with iterated integrals to determine our $J(n)$ integrals, then we take the limits of these, as illustrated by the next example.

■ **Example 4.3:**

Check whether the following integrals converge. If they do, compute them.

$$(a) \iint_D \frac{dx \, dy}{1 + (x + y)^2}, \quad D = \{(x, y) : x > 0, y > 0\}.$$

$$(b) \iint_D \frac{dx \, dy}{\sqrt{xy}}, \quad D = \{(x, y) : 0 < x < 1, 0 < y < 1\}.$$

Solution:

- (a) Take the x -simple integral over a finite domain

$$D_{AB} = \{(x, y) : 0 < x < A, 0 < y < B\}.$$

$$\begin{aligned} \iint_{D_{AB}} \frac{1}{1 + (x+y)^2} dx dy &= \int_0^B dy \int_0^A \frac{1}{1 + (x+y)^2} dx \\ &= \int_0^B dy \left[\arctan(x+y) \right]_{x=0}^A = \int_0^B (\arctan(y+A) - \arctan y) dy. \end{aligned}$$

$$\text{Now } \lim_{B \rightarrow \infty} \int_0^B \arctan y dy = \lim_{B \rightarrow \infty} \left[y \arctan y - \frac{1}{2} \ln(1+y^2) \right]_0^B,$$

which evidently does not exist. The integral over D does not converge.

- (b) The integrand is undefined on the boundary of D , at $x = 0$ and at $y=0$.

So integrate over the domain $D_\epsilon = \{(x, y) : \epsilon < x < 1, \epsilon < y < 1\}$.

$$\begin{aligned} \iint_{D_\epsilon} \frac{1}{\sqrt{xy}} dx dy &= \int_\epsilon^1 dy \int_\epsilon^1 \frac{1}{\sqrt{xy}} dx = \int_\epsilon^1 dy \left[2\sqrt{\frac{x}{y}} \right]_\epsilon^1 \\ &= \int_\epsilon^1 2 \frac{1 - \sqrt{\epsilon}}{\sqrt{y}} dy = 4(1 - \sqrt{\epsilon})^2 \longrightarrow 4 \text{ as } \epsilon \rightarrow 0. \end{aligned}$$

■

↳ Mastery Check 4.5:

Does the double integral $I = \iint_D f(x, y) dA$ converge or diverge when

$$f(x, y) = \frac{1}{x^2 + y^2} \quad \text{and} \quad D = \{x, y) : x \geq 1, 0 \leq y \leq x\}?$$

Hint: Draw a picture. Consider the sub-domain $D_R = \{x, y) : 1 \leq x \leq R, 0 \leq y \leq x\}$, and let $R \rightarrow \infty$. Write out both iterated integral versions, and choose the simpler of the two for analysis.

↳

↳ Mastery Check 4.6:

Dangers in ∞ !

Consider the two iterated integrals

$$I_1 = \int_0^1 dy \int_1^\infty (e^{-xy} - 2e^{-2xy}) dx, \quad \text{and} \quad I_2 = \int_1^\infty dx \int_0^1 (e^{-xy} - 2e^{-2xy}) dy.$$

Show that one of these must be < 0 , while the other must be > 0 , and ponder the implications in the context of Theorem 4.3.

↳

4.E Change of variables in \mathbb{R}^2

Sometimes, the multiple integrals are not so easy to evaluate when expressed in Cartesian coordinates, x and y . The problem might stem from the function being difficult, or from the domain being convoluted, or *both*. This is different from the 1D case where we only need to worry about the function.

In the 2D situation the problem could be avoided by changing variables from x and y to new variables, u and v , say. Here, we discuss a process for doing this, and on the way we point out when it is possible to do so, and when it is not. The 1D case provides an interesting comparison.

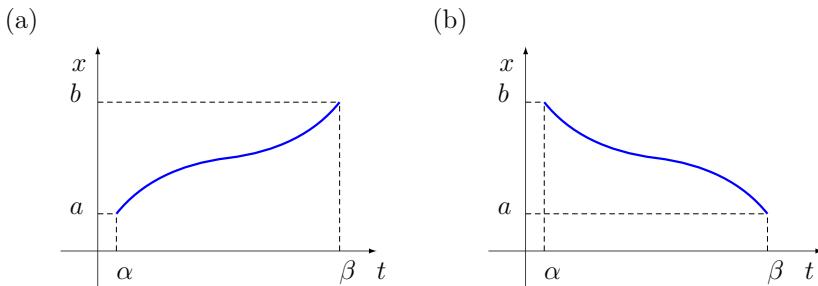


Figure 4.17 (a) $x(t)$ increasing; (b) $x(t)$ decreasing.

Change of variables in single integrals

Aim: We want to evaluate $\int_a^b f(x) dx$ by invoking a transform $x = x(t)$.

Suppose $x(t)$ is strictly *increasing* as in Figure 4.17(a), then $x'(t) > 0$, and $\frac{dF}{dx} = f$. Then

$$\begin{aligned} \int_{\alpha}^{\beta} f(x(t))x'(t) dt &= \int_{\alpha}^{\beta} F'(x(t))x'(t) dt \\ &= \int_{\alpha}^{\beta} \frac{d}{dt}F(x(t)) dt = F(x(\beta)) - F(x(\alpha)) \\ &= F(b) - F(a) = \int_a^b f(x) dx. \end{aligned}$$

On the other hand suppose $x(t)$ is strictly *decreasing* as in the case shown in

Figure 4.17(b), then $x'(t) < 0$, and $\frac{dF}{dx} = f$. Then

$$\begin{aligned} \int_{\alpha}^{\beta} f(x(t))(-x'(t)) dt &= - \int_{\alpha}^{\beta} F'(x(t))x'(t) dt \\ &= \int_{\beta}^{\alpha} \frac{d}{dt} F(x(t)) dt = F(x(\alpha)) - F(x(\beta)) \\ &= F(b) - F(a) = \int_a^b f(x) dx. \end{aligned}$$

Consequently,

$$\int_a^b f(x) dx = \int_{\alpha}^{\beta} f(x(t)) \left| \frac{dx}{dt} \right| dt.$$

So, in this case we see the integration interval has changed, but more significantly the change of variable $x \rightarrow t$ has introduced a positive factor $|x'|$ in the integral. This is called the *scale factor* since it scales up or down the interval size. We should expect a similar factor to appear in multiple integrals.

Change of variables in double integrals

For convenience we shall consider only bijective transformations:

$$\tau: \mathbf{u} \mapsto \mathbf{x}(\mathbf{u}) = \begin{cases} x = x(u, v) \\ y = y(u, v) \end{cases}$$

such that $\frac{\partial(x, y)}{\partial(u, v)} \neq 0$. The Jacobian determinant (Definition 2.9) is involved in the transformation of double integrals:

$$\iint_D f(x, y) dA \quad \text{becomes expressed as} \quad \iint_E g(u, v) dA'.$$

Geometrically the transformation affects areas both globally and locally.

To see how the transformation does this consider in Figure 4.18 the “parallelogram” in the xy -plane created by constant u and v contours. Suppose opposite sides of the parallelogram are separated by differences du and dv .

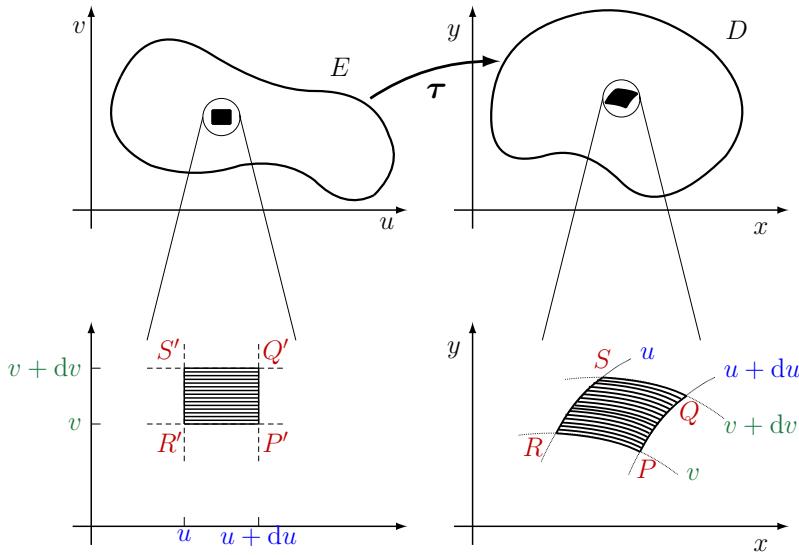


Figure 4.18 A geometrical view of a change of variables.

For small du and dv , the area of the element in D is given by the vector product

$$\begin{aligned} dA &= \left| \overrightarrow{RP} \times \overrightarrow{RS} \right| \quad \text{— geometric interpretation (see Page 3)} \\ &= \left| (dx_u \mathbf{e}_1 + dy_u \mathbf{e}_2) \times (dx_v \mathbf{e}_1 + dy_v \mathbf{e}_2) \right|. \\ &\quad \underbrace{\phantom{\left| (dx_u \mathbf{e}_1 + dy_u \mathbf{e}_2) \times (dx_v \mathbf{e}_1 + dy_v \mathbf{e}_2) \right|}}_{\text{along a line of}} \quad \underbrace{\phantom{\left| (dx_u \mathbf{e}_1 + dy_u \mathbf{e}_2) \times (dx_v \mathbf{e}_1 + dy_v \mathbf{e}_2) \right|}}_{\text{along a line of}} \\ &\quad \text{constant } v \quad \text{constant } u \end{aligned}$$

Thus,

$$\begin{aligned} dA &= \left| \left(\frac{\partial x}{\partial u} du \mathbf{e}_1 + \frac{\partial y}{\partial u} du \mathbf{e}_2 \right) \times \left(\frac{\partial x}{\partial v} dv \mathbf{e}_1 + \frac{\partial y}{\partial v} dv \mathbf{e}_2 \right) \right| \\ &\quad \text{— by the chain rule (Section 2.G)} \\ &= \left| \frac{\partial(x, y)}{\partial(u, v)} \right| du dv \quad \text{— cross product gives the Jacobian determinant} \end{aligned}$$

The reader should make particular note that it is the *absolute value* of the Jacobian determinant that appears here. This is reasonable since what we have done is transformed one area element to another, preserving the sign.

And so we have the following important theorem.

Theorem 4.5

Let $x(u, v)$ and $y(u, v)$ be a bijective [one-to-one and onto] and C^1 transformation of E in the uv -plane onto D in the xy -plane.

If $f(x, y)$ is integrable in D then $f(x(u, v), y(u, v)) = F(u, v)$ is integrable in E and

$$\iint_P f(x, y) \, dA = \iint_E F(u, v) \left| \frac{\partial(x, y)}{\partial(u, v)} \right| \, dA'.$$

Remarks

- * The *absolute value* of the Jacobian is the *scale factor* between the two area elements, dA in the xy -plane and dA' in the uv -plane; it takes the role played by $\left| \frac{dx}{dt} \right|$ in the single-variable case.

In other words $dA = \left| \frac{\partial(x, y)}{\partial(u, v)} \right| dA'$.

- * If $x(u, v)$ and $y(u, v)$ are bijective transformations, then $\frac{\partial(x, y)}{\partial(u, v)} \neq 0$.
 - * If $x(u, v)$ and $y(u, v)$ are C^1 , then $F(u, v) = f(x(u, v), y(u, v))$ is integrable over E whenever f is integrable over D .
 - * A change of variables in a *double* integral is NOT the same as a substitution in an *iterated* integral.

Before we demonstrate how one considers the change of variable to evaluate a double integral, we encourage the reader to verify the following Jacobian expressions.

Mastery Check 4.7:

Show that in transforming from (Hint: See Section 1.D.)

- (a) Cartesian to polar coordinates (r, θ) the Jacobian is r ;
 - (b) Cartesian to cylindrical coordinates (r, θ, z) the Jacobian is r ;

(c) Cartesian to spherical coordinates (ρ, ϕ, θ) the Jacobian is $\rho^2 \sin \phi$.



■ Example 4.4:

We wish to show that the volume of the right cone of Figure 4.19, of radius a and height h , is $V = \frac{1}{3}\pi a^2 h$.

We do this by integrating a variable height function $z = \frac{h}{a}(a - \sqrt{x^2 + y^2})$ over the region $R : x^2 + y^2 \leq a^2$.

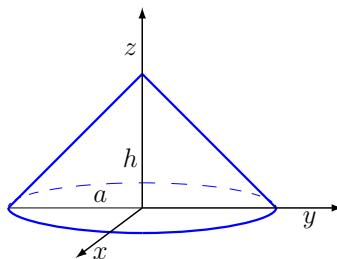


Figure 4.19 A right cone.

The volume is $V = \iint_R \frac{h}{a}(a - \sqrt{x^2 + y^2}) \, dx \, dy$.

Change to polar coordinates $(x, y) \rightarrow (r, \theta)$, $x = r \cos \theta$, $y = r \sin \theta$, $J = r$.

$$\begin{aligned} V &= \iint_R \frac{h}{a}(a - r)r \, dr \, d\theta = \int_0^a dr \int_0^{2\pi} \frac{h}{a}(ar - r^2) \, d\theta \\ &= \frac{2\pi h}{a} \left[a \frac{r^2}{2} - \frac{r^3}{3} \right]_0^a = \frac{\pi h a^2}{3}. \end{aligned}$$



✍ Mastery Check 4.8:

Does the integral $\iint_{\mathbb{R}^2} \frac{dx \, dy}{(1 + x^2 + y^2)^2}$ converge? If it does, find its value.

Hint: Find the integral over $D = \{(x, y) : x^2 + y^2 < A^2\}$, and let $A \rightarrow \infty$.



Philosophy of a change of variables

It is worthwhile pausing to reflect on the motivation behind a change of variables. This may help to guide the practitioner to choose the most suitable variables. Ultimately, we invoke a change of variables to simplify the evaluation of a multiple integral. One therefore chooses a transformation

$$(x(u, v), y(u, v)) : (x, y) \longmapsto (u, v)$$

to *either*

- (a) transform region D into a simpler region E (in the Example 4.4 E was a simple rectangle); *or*
- (b) transform the integrand $f(x, y)$ into some simpler $F(u, v)$.

Thus, in the case of (a), one is guided by the shape of the region D : What mathematical expressions determine the boundary, and can these be used to define the new boundary?

In the case of (b), is the form of $f(x, y)$ suggestive of a suitable transformation? For example, $f(x, y) = g(x^2 + y^2)$ suggests the polar coordinate transformation $x = r \cos \theta$, $y = r \sin \theta$, so that $x^2 + y^2 = r^2$.

In both cases always look for symmetry features. This having been said, there will always be consequences.

- (c) The region can be transformed into a more complicated one, even if the integrand becomes simpler;
- (d) The integrand may become more complex, even if the transformed region is simpler (recall that a Jacobian for the transformation needs to be considered);
- (e) Sometimes, however, we get lucky and *both* the function *and* the region become simpler.

✍ Mastery Check 4.9:

What is the image R_2 of the region R_1 bounded by the curves $x = y^2$, $y = \frac{1}{x}$, $y = \frac{2}{x}$, $x = \frac{y^2}{2}$, under the transformation $u = \frac{x}{y^2}$, $v = xy$?

Hint: Draw a picture for the region R_1 in the xy -plane and another for the

image R_2 in the uv -plane. Then stop to think: Is this OK?



💡 Mastery Check 4.10:

Calculate the volume of the solid defined by the intersection of the regions

$$x \geq 0, \quad z \geq 0, \quad z \leq x^2 + y^2, \quad 1 \leq x^2 + y^2 \leq 4.$$

Hint: Draw the graph of the solid (using MATLAB® is best). Sketch the region of integration in the xy -plane. Make an attempt at the integration in the x, y variables. Then think about a suitable transformation.



💡 Mastery Check 4.11:

This is an exercise with an ulterior motive — a bit like the last one, but more so!

Evaluate the double integral $I = \iint_T e^{(y-x)/(y+x)} dA$

where T is the triangle with corners at the points $(0, 0)$, $(a, 0)$, $(0, a)$.

Hint: An attempt to integrate in the xy -plane is likely to fail, but try anyway. Then make the simplest, most obvious, transformation: it works beautifully!

Finally, try the next most “obvious” transformation, polar coordinates. That works, too. You may need to recall a couple of trigonometric relations:

$$\tan(\theta - \pi/4) = \frac{\tan \theta - 1}{1 + \tan \theta}, \quad \cos(\theta - \pi/4) = \frac{\cos \theta + \sin \theta}{\sqrt{2}}.$$



4.F Triple integrals

To cement the ideas we've just introduced we now illustrate the case for functions of *three* variables integrated over regions of \mathbb{R}^3 .

Suppose $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ is a continuous function, defined (at least) over a rectangular box B : $B = \{(x, y, z) : a_1 \leq x \leq b_1, a_2 \leq y \leq b_2, a_3 \leq z \leq b_3\}$.
— we assume B to be closed and bounded, $a_i \leq b_i < \infty$.

B is shown in Figure 4.20.

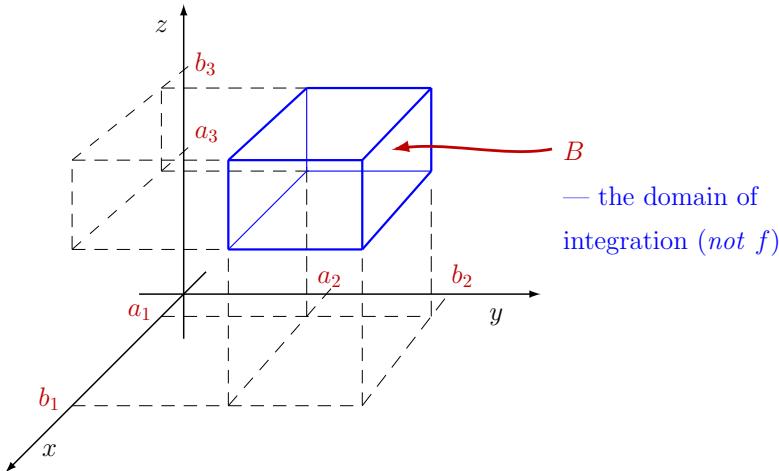


Figure 4.20 Box domain B in \mathbb{R}^3 .

The definition of a triple integral follows in analogy with those of double and single integrals. We just outline the relevant steps leading to the definition.

- Partition B into small rectangular blocks B_{ijk} of volumes $\Delta V_{ijk} = \Delta x_i \Delta y_j \Delta z_k$
- If $f \geq 0$ for all $\mathbf{x} \in B$, then we interpret the quantity f as a *density* so that in choosing $(\xi_i, \eta_j, \zeta_k) \in B_{ijk}$, $f(\xi_i, \eta_j, \zeta_k) \Delta V_{ijk}$ will be an approximation to the mass of block B_{ijk} .
- The Riemann sum of all such masses in the partition,

$$\sigma_{nml} = \sum_{i=1}^n \sum_{j=1}^m \sum_{k=1}^{\ell} f(\xi_i, \eta_j, \zeta_k) \Delta V_{ijk},$$

is an approximation to the total *mass* of the entire B .

- We then take the combined limits of the number of boxes to infinity with vanishing volumes. We therefore arrive at

Definition 4.4

The triple integral of f over B is defined as the limit (if it exists)

$$\lim_{\substack{n,m,\ell \rightarrow \infty \\ \max \sqrt{\Delta x_i^2 + \Delta y_j^2 + \Delta z_k^2} \rightarrow 0}} \sigma_{nml} = \iiint_B f(x, y, z) \, dV.$$

For regions of more general shape, starting with bounded regions, S , we can extend the definition of the integral of f over $S \subset \mathbb{R}^3$ in analogy with the 2D version:

- we enclose S in a closed and bounded box B : $S \subset B \subset \mathbb{R}^3$, and
- define $\hat{f}(x, y, z) = \begin{cases} f(x, y, z), & \mathbf{x} \in S \\ 0, & \mathbf{x} \notin S \end{cases}$

we then have

$$\iiint_S f(x, y, z) \, dV \equiv \iiint_B \hat{f}(x, y, z) \, dV.$$

This now sets the stage for Section 4.G where the practical evaluation of triple integrals over general regions is discussed. In the meantime we have, as in the 1D and 2D cases, the following useful theorem and its corollaries.

An important theorem and some corollaries — 3D version.

Theorem 4.1 *All continuous functions are integrable over compact subsets of their domains.*

Corollary 4.1.1 *If $f \geq 0$, then*

$\iiint_B f \, dV$ *is the “volume” of a 4-dimensional “solid” “under” f “over” B .*
— not a very helpful interpretation.

Corollary 4.1.2

If $f \geq 0$ is a $\frac{\text{mass}}{\text{charge}}$ density, then $\iiint_B f \, dV = \begin{cases} \text{total mass} \\ \text{total charge} \end{cases}$.
— a more helpful interpretation.

Corollary 4.1.3 *If $f = 1$, then $\iiint_B 1 \, dV = \text{volume of solid } B$.*

— even more useful, especially for more general regions.

Corollary 4.1.4 *Average of $f(x, y, z)$ over $B = \frac{\iiint_B f \, dV}{\iiint_B 1 \, dV} = \frac{\iiint_B f \, dV}{\text{vol. } B}$.*
— this is true even for more complex regions.

Corollary 4.1.5 *Linearity: If $a, b \in \mathbb{R}$ then*

$$\iiint_B (af + bg) \, dV = a \iiint_B f \, dV + b \iiint_B g \, dV.$$

Corollary 4.1.6 *Additivity: (very important)*

If $B_1 \cap B_2 = \{\} \equiv \emptyset$, then

$$\iiint_{B_1 \cup B_2} f \, dV = \iiint_{B_1} f \, dV + \iiint_{B_2} f \, dV.$$

As for the actual calculation of triple integrals the next section extends the ideas outlined in Section 4.C.

4.G Iterated integration in \mathbb{R}^3

In evaluating a triple integral, we again make use of the idea of “slicing”. The combination of a 3D integration domain and a function $f(x, y, z)$ results in a 4D graph, but we can visualize only the 3D domain. Keep in mind therefore that the first slice through the domain actually results in a *projection* of the graph of f onto 3D.

In Figure 4.21 we first take horizontal slices through the 3D integration domain (left panel) which results in 2D regions (right panel). We then take constant x - or constant y -slices through the horizontal z -slice.

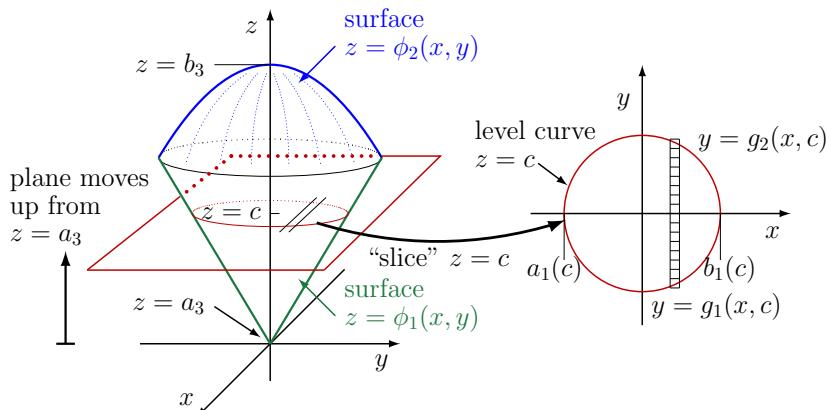


Figure 4.21 The process of domain slicing: first in z , then in x , last in y .

In choosing the order in which we take slices we also commit ourselves to the order in which we perform the iterated integrals (in reverse). For instance, in the case of Figure 4.21 the iterated integral is the following:

$$\iiint_S f \, dV = \int_{a_3}^{b_3} dz \int_{a_1(z)}^{b_1(z)} dx \int_{g_1(x,z)}^{g_2(x,z)} f(x, y, z) \, dy.$$

That is, one integrates with respect to y first, between limits that depend on x and z , then one integrates with respect to x , between limits that depend on z . Finally, one integrates with respect to z between two constants.

For the exact same problem we can consider vertical slices along the x -axis from a_1 to b_1 . For each x value we can take either y -slices or z -slices. Figure 4.22 shows the procedure corresponding to x -slices, then y -slices, and finally z -slices.

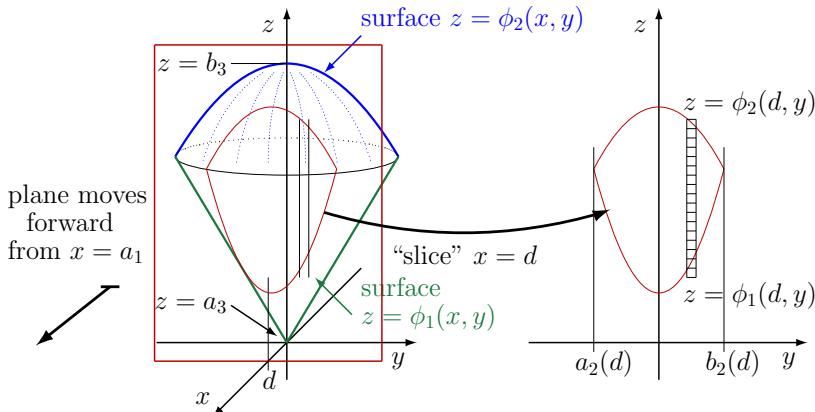


Figure 4.22 The process of domain slicing: first in x , then in y , last in z .

The iterated integral for the particular case of Figure 4.22 is this

$$\iiint_S f \, dV = \int_{a_1}^{b_1} dx \int_{a_2(x)}^{b_2(x)} dy \int_{\phi_1(x,y)}^{\phi_2(x,y)} f(x, y, z) \, dz.$$

That is, one integrates with respect to z first, between limits that depend on x and y , then one integrates with respect to y , between limits that depend on x , and finally one integrates with respect to x between two constants.

Generally, for any triple integral, the possible alternatives are these:

$$\begin{aligned}
 \iiint_S f \, dV : \quad & x\text{-slices} \quad \left\{ \begin{array}{l} \int_{a_1}^{b_1} dx \int_{a_2(x)}^{b_2(x)} dy \int_{a_3(x,y)}^{b_3(x,y)} f \, dz \\ \int_{a_1}^{b_1} dx \int_{a_3(x)}^{b_3(x)} dz \int_{a_2(x,z)}^{b_2(x,z)} f \, dy \end{array} \right. \\
 & y\text{-slices} \quad \left\{ \begin{array}{l} \int_{a_2}^{b_2} dy \int_{a_1(y)}^{b_1(y)} dx \int_{a_3(x,y)}^{b_3(x,y)} f \, dz \\ \int_{a_2}^{b_2} dy \int_{a_3(y)}^{b_3(y)} dz \int_{a_1(y,z)}^{b_1(y,z)} f \, dx \end{array} \right. \\
 & z\text{-slices} \quad \left\{ \begin{array}{l} \int_{a_3}^{b_3} dz \int_{a_1(z)}^{b_1(z)} dx \int_{a_2(x,z)}^{b_2(x,z)} f \, dy \\ \int_{a_3}^{b_3} dz \int_{a_2(z)}^{b_2(z)} dy \int_{a_1(y,z)}^{b_1(y,z)} f \, dx \end{array} \right.
 \end{aligned}$$

Thus, one triple integral gives rise to six possible iterated integrals. Note carefully the nature of the limits of each variable and how their dependencies vary from one integral to the next.

■ Example 4.5:

Evaluate the integral $\iiint_D xy \, dV$, where D is the interior of the sphere $x^2 + y^2 + z^2 = 1$ in the first octant, $0 \leq x, y, z \leq 1$.

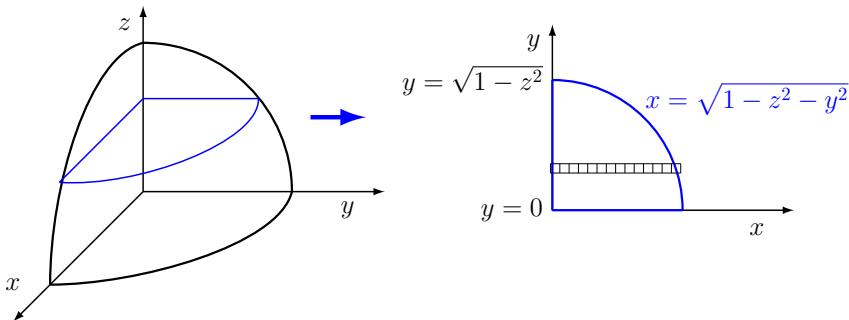


Figure 4.23 The domain D , and the projection of a z -slice onto the xy -plane.

Solution: As sketched in Figure 4.23, we will take horizontal slices, $z = a$

constant, for $0 \leq z \leq 1$, and (for fun) integrate w.r.t. x before we integrate w.r.t. y . The slices are $x^2 + y^2 \leq 1 - z^2$, $x \geq 0$, $y \geq 0$, for a given z .

We need the bounds for the x -integration as functions of y .

These are $x = 0, \sqrt{1 - z^2 - y^2}$.

The bounds for the y -integration are $0, \sqrt{1 - z^2}$, and for the z -integration they are $0, 1$. So we get

$$\begin{aligned}
 \iiint_D xy \, dV &= \int_0^1 dz \int_0^{\sqrt{1-z^2}} dy \int_0^{\sqrt{1-z^2-y^2}} xy \, dx \\
 &= \int_0^1 dz \int_0^{\sqrt{1-z^2}} y \left[\frac{1}{2} x^2 \right]_{x=0}^{\sqrt{1-z^2-y^2}} \, dy \\
 &= \int_0^1 dz \int_0^{\sqrt{1-z^2}} \frac{1}{2} y (1 - z^2 - y^2) \, dy \\
 &= - \int_0^1 \frac{1}{8} \left[(1 - z^2 - y^2)^2 \right]_{y=0}^{\sqrt{1-z^2}} \, dz \\
 &= - \int_0^1 \frac{1}{8} (0 - (1 - z^2)^2) \, dz = \frac{1}{15}.
 \end{aligned}$$



☞ Mastery Check 4.12:

Rewrite the iterated integral

$$I = \int_0^1 dz \int_z^1 dx \int_0^{x-z} f \, dy$$

as an iterated integral with the outermost integration w.r.t. x and innermost integration w.r.t. z .

Hint: Use the given limits to determine the 3D region of integration and then establish the limits of the new iterated integral.



 **Mastery Check 4.13:**

Let S be the volume of the body bounded by the planes $x = 0$, $y = 0$, and $z = 4$, and the surface $z = x^2 + y^2$. Calculate $I = \iiint_S x \, dV$.

 **Mastery Check 4.14:**

Compute the volume of that part of the cylinder $x^2 - 2x + y^2 = 0$ cut off by the cylinder $z^2 - 2x = 0$. The region defined by this intersection is shown in Figure 4.24. (See also Figure 1.34 on Page 46.)

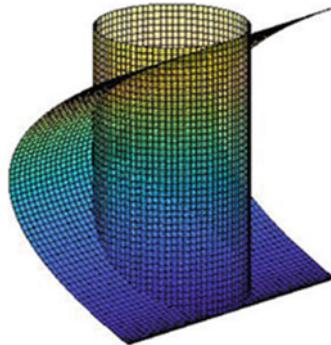


Figure 4.24 Two intersecting cylinders.

4.H Change of variables in \mathbb{R}^3

Consider the triple integral of a continuous function over a closed and bounded domain,

$$I = \iiint_S f(x, y, z) \, dV,$$

and the bijective C^1 -transformation:

$$\tau : \mathbb{R}^3 \longrightarrow \mathbb{R}^3; \quad \mathbf{u} \mapsto \mathbf{x}(\mathbf{u}).$$

This mapping transforms an element of volume in the xyz -domain to a volume element in the uvw -domain, as suggested graphically in Figure 4.25.

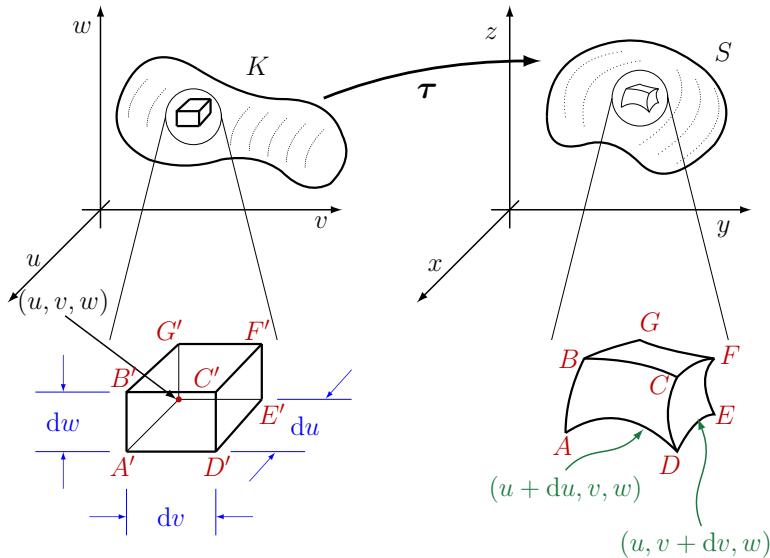


Figure 4.25 A geometrical view of a change of variables in a 3D domain.

If $du, dv, dw \ll 1$ then (according to Section 1.A) the volume element in S would be given by the absolute value of the scalar triple product

$$\begin{aligned} dV &= |(\overrightarrow{BA} \times \overrightarrow{BC}) \cdot \overrightarrow{BG}| \\ &= |((dx_u \mathbf{e}_1 + dy_u \mathbf{e}_2 + dz_u \mathbf{e}_3) \times (dx_v \mathbf{e}_1 + dy_v \mathbf{e}_2 + dz_v \mathbf{e}_3)) \\ &\quad \cdot (dx_w \mathbf{e}_1 + dy_w \mathbf{e}_2 + dz_w \mathbf{e}_3)| \end{aligned}$$

Invoking the chain rule, the scalar triple product can be written in determinant form (Section 1.A, Page 8).

$$dV = \left| \begin{array}{ccc} \frac{\partial x}{\partial u} du & \frac{\partial y}{\partial u} du & \frac{\partial z}{\partial u} du \\ \frac{\partial x}{\partial v} dv & \frac{\partial y}{\partial v} dv & \frac{\partial z}{\partial v} dv \\ \frac{\partial x}{\partial w} dw & \frac{\partial y}{\partial w} dw & \frac{\partial z}{\partial w} dw \end{array} \right| = \left| \frac{\partial(x, y, z)}{\partial(u, v, w)} \right| du dv dw.$$

As we found in the case of a change of variables in a double integral, a change of variables in a triple integral involves the Jacobian determinant for that transformation. The *absolute value* of the Jacobian is the scaling factor between the volume elements $dV = dx dy dz$ in xyz -space and $dV' = du dv dw$ in uvw -space.

Consequently, we have the end result that

$$\iiint_S f(x, y, z) \, dV = \iiint_K F(u, v, w) \left| \frac{\partial(x, y, z)}{\partial(u, v, w)} \right| \, dV',$$

where $F(u, v, w) = f(x(u, v, w), y(u, v, w), z(u, v, w))$ and S is the image of K under τ .

☞ Mastery Check 4.15:

Transform the iterated integral below into iterated integrals w.r.t. cylindrical coordinates, and w.r.t. spherical coordinates:

$$I = \int_0^1 dx \int_0^{\sqrt{1-x^2}} dy \int_0^{1+x+y} (x^2 - y^2) \, dz.$$

(Do not proceed to evaluate the integrals. If you skip to the very last note in this chapter, you will see why!) 

4.I n -tuple integrals

As we did in going from double integrals to triple integrals, all the preceding ideas, concepts and mathematical arguments can be generalized to n dimensions.

***n*-tuple integrals**

Suppose $S \subset \mathbb{R}^n$ is closed and bounded and we have $f : S \rightarrow \mathbb{R}$ (the graph of $f \subset \mathbb{R}^{n+1}$).

Enclose S in an n -dimensional box

$$[a_1, b_1] \times [a_2, b_2] \times \dots \times [a_n, b_n].$$

Partition the box into n -dimensional boxes of size

$$\Delta x_1 \times \Delta x_2 \times \Delta x_3 \times \dots \times \Delta x_n.$$

Choose $\xi_i \in [x_i, x_i + \Delta x_i] \quad i = 1, 2, \dots, n.$

Form the sum

$$\sum \sum \dots \sum f(\xi_1, \xi_2, \dots, \xi_n) \Delta x_1 \Delta x_2 \dots \Delta x_n.$$

If the limit of this sum as $n \rightarrow \infty$ and $|\Delta \mathbf{x}| \rightarrow 0$ exists we call it the n -dimensional integral of f over S

$$\iint \dots \int_S f(x_1, x_2, \dots, x_n) dV_n = I,$$

where dV_n is an n -dimensional volume element.

Iterated integrals

If S can be described as

$$S = \{ \mathbf{x} = (x_1, x_2, \dots, x_n) : (x_1, x_2, \dots, x_{n-1}) \in W \subset \mathbb{R}^{n-1} \text{ and } \phi_1(x_1, x_2, \dots, x_{n-1}) \leq x_n \leq \phi_2(x_1, x_2, \dots, x_{n-1}) \}$$

$$\text{then } I = \iint \dots \int_W dx_1 \dots dx_{n-1} \int_{\phi_1}^{\phi_2} f(\mathbf{x}) dx_n = \dots$$

and, in fact,

$$I = \int_{\alpha_1}^{\alpha_2} dx_1 \int_{\beta_1(x_1)}^{\beta_2(x_1)} dx_2 \int_{\gamma_1(x_1, x_2)}^{\gamma_2(x_1, x_2)} dx_3 \dots \int_{\phi_1(x_1, \dots, x_{n-1})}^{\phi_2(x_1, \dots, x_{n-1})} f(\mathbf{x}) dx_n,$$

which is just one of $n!$ alternative iterated integrals.

Change of variables

Consider a bijective C^1 transformation: $\tau : \mathbf{u} \mapsto \mathbf{x}(\mathbf{u})$, where the functions $x_i = x_i(u_1, u_2, \dots, u_n)$, $i = 1, 2, \dots, n$, are such that the Jacobian (Section 2.H)

$$J = \frac{\partial(x_1, \dots, x_n)}{\partial(u_1, \dots, u_n)} \neq 0.$$

The n -tuple integral I of $f(x_1, \dots, x_n)$ over S is equal to the n -tuple integral of the product of $|J|$ and

$$F(u_1, \dots, u_n) = f(x_1(u_1, u_2, \dots, u_n), \dots, x_n(u_1, u_2, \dots, u_n))$$

over the pre-image E of S under τ :

$$I = \iint \dots \int_S f(\mathbf{x}) dV_n = \iint \dots \int_E F(\mathbf{u}) |J| dV'_n.$$

See Section 5.A for an expansion of J .

4.J Epilogue: Some practical tips for evaluating integrals

Recall that ...

- * If $f(-x) = -f(x)$, then f is *odd*, and $\int_{-a}^a f(x) dx = 0$.
— for example: $\sin(x)$, x^3 , x^7 , $\arctan(x)$
- * If $f(-x) = f(x)$, then f is *even*, and $\int_{-a}^a f(x) dx = 2 \int_0^a f(x) dx$.
— for example: $\cos(x)$, x^2 , $\sin^2(x)$

Did you know that for functions of $n \geq 2$ variables there are other possible symmetry features?

- * If $f(-x, y) = -f(x, y)$, then f is *odd with respect to x* , which means that
$$\int_c^d dy \int_{-a}^a f(x, y) dx = 0. \quad \text{— for example: } xy^2, \sin(x).y, \arctan(xy)$$
- * If $f(-x, y) = f(x, y)$, then f is *even with respect to x* , which means that
$$\int_c^d dy \int_{-a}^a f(x, y) dx = 2 \int_c^d dy \int_0^a f(x, y) dx. \quad \text{— for example: } x^2y, \cos(x).y^3, \arctan(x^2y)$$
- * Similarly, we may have functions $f(x, y)$ which are *odd or even with respect to y* .
- * Now, for $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ we have a *new* possibility, indicated in Figure 4.26:

If $f(x, y) = f(y, x)$, then f is *symmetric* across $y = x$.

If $f(x, y) = -f(y, x)$, then f is *antisymmetric* across $y = x$.

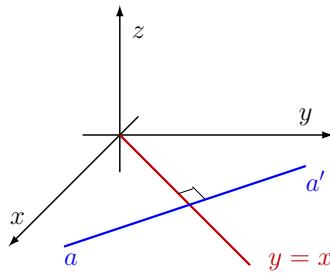


Figure 4.26 The oblique symmetry line $y = x$.

Some examples of functions *symmetric* from point a to point a' are

$$f(x, y) = x + y,$$

$$f(x, y) = x^2 + y^2,$$

$$f(x, y) = 1/\sqrt{x^2 + y^2},$$

$$f(x, y) = b + (y - x)^2 \text{ (Figure 4.27(a))}.$$

Some examples of functions *antisymmetric* from point a to point a' are

$$f(x, y) = x - y,$$

$$f(x, y) = x^2 - y^2,$$

$$f(x, y) = \sin(y - x) \text{ (Figure 4.27(b))}.$$

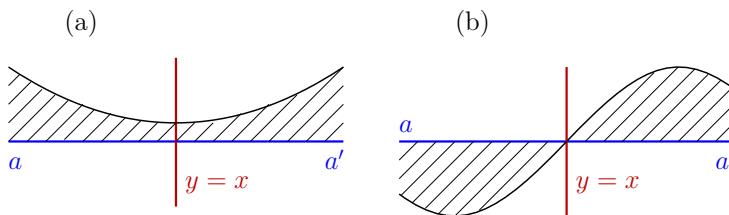


Figure 4.27 (a) Cross-section through $f(x, y) = b + (y - x)^2$; (b) Cross-section through $f(x, y) = \sin(y - x)$.

4.K Supplementary problems

Section 4.B

1. Determine $\iint_D x \cos xy \, dx \, dy$,
where $D = \{(x, y) \in \mathbb{R}^2 : 0 \leq x \leq \pi/2, 0 \leq y \leq 1\}$.
(Hint: Treat this as an iterated integral and integrate with respect to y first.)
2. Determine $\iint_D 2xy \sec^2(x^2y) \, dx \, dy$,
where $D = \{(x, y) \in \mathbb{R}^2 : 0 \leq x \leq 1, 0 \leq y \leq \pi/4\}$.

Section 4.C

3. Compute the integral $\iint_D \sqrt{xy} \, dx \, dy$,
where $D = \{(x, y) : 0 \leq x \leq 1, x^3 \leq y \leq x^2\}$.
4. Compute the integral $\iint_D xy \, dx \, dy$,
where $D = \{(x, y) : \cos y \leq x \leq \sin y, \frac{\pi}{4} \leq y \leq \frac{5\pi}{4}\}$.

Section 4.D

5. Check whether the following integrals converge. If they do, compute them.

- (a) $\iint_D \frac{dx \, dy}{1 + (x + y)^2}$, $D = \{(x, y) : x > 0, y > 0\}$.
- (b) $\iint_D \frac{dx \, dy}{\sqrt{xy}}$, $D = \{(x, y) : 0 < x < 1, 0 < y < 1\}$.
- (c) $\iint_D x e^{-(y+x^2/y)} \, dx \, dy$, where D is the first quadrant.

(Hint: Consider the rectangle $\{(x, y) : 0 < x < A, 0 < y < B\}$, and let $A \rightarrow \infty$, $B \rightarrow \infty$.)

Section 4.E

6. Determine $\iint_D ye^{-(x^2+y^2)} dx dy$, where $D = \{(x, y) \in \mathbb{R}^2 : x^2 + y^2 \leq 1, y > 0\}$.
7. A device commonly used to determine $I = \int_{-\infty}^{\infty} e^{-x^2/2} dx$ is as follows:

From symmetry (see Section 4.J) we may write

$$I^2 = 4 \left(\int_0^{\infty} e^{-x^2/2} dx \right) \left(\int_0^{\infty} e^{-y^2/2} dy \right) = \iint_R e^{-(x^2+y^2)/2} dx dy,$$

where R is the first quadrant on the Cartesian plane.

Use a polar coordinate transformation to evaluate this integral.

8. Let $D = \{(x, y) : x \geq 0, y \geq 0, x^3 + y^3 \leq 0\}$.

By introducing a change of variables $(x, y) \rightarrow (u = x^3, v = y^3)$ evaluate the double integral

$$\iint_D x^2 y^3 \sqrt{1 - x^3 - y^3} dA.$$

Express your answer in terms of $\Gamma(1/3)$ and $\Gamma(1/6)$.

Section 4.F

9. Recall from elementary physics that

- If a force $F = mg$ (where g is the acceleration due to gravity) acts on a mass m at a horizontal distance x from a (pivot) point O , then its *moment* about O is $F \cdot x$;
- If a mass m moves so that its distance from a point O is constrained to be x , then its *moment of inertia* about O is $m \cdot x^2$.

Now consider the integral $\iiint_V f(x, y, z) dx dy dz$.

Interpret this integral if

- $f(x, y, z) = 1$.
- $f(x, y, z)$ is the material density at point (x, y, z) .
- $f(x, y, z) = \sqrt{x^2 + y^2} \times$ the material density at point (x, y, z) .
- $f(x, y, z) = (x^2 + y^2) \times$ the material density at point (x, y, z) .

Section 4.G

10. For each of the following, sketch the region and evaluate the integral

$$(i) \int_0^1 dz \int_0^z dy \int_0^y dx$$

$$(ii) \int_0^2 dx \int_1^x dy \int_2^{x+y-1} y dz$$

$$(iii) \int_0^1 dx \int_x^{\sqrt{x}} dy \int_{1-x-y}^{1+x+y} xy dz$$

11. Let S be the solid in the first octant containing the origin and bounded by sections of the plane $y = 1 + 2x$ and the sphere $x^2 + y^2 + z^2 = 4$. For any integrable function f defined on S write the triple integral $\iiint_S f dV$ as an iterated integral with respect to x, y, z in six different ways.

That is, you are to determine the respective limits depicted in the iterated integral formulae on Page 209.

12. Determine the volume shown in Figure 4.28 enclosed by the two surfaces

$$z = 8 - x^2 - y^2 \quad \text{and} \quad z = x^2 + 3y^2.$$

Hint: Set up the integral as $V = \int_a^b dx \int_{a_1(x)}^{a_2(x)} dy \int_{a_3(x,y)}^{a_4(x,y)} 1 dz$, that is, take horizontal z -slices.

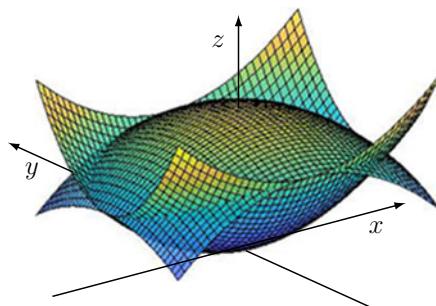


Figure 4.28 Two intersecting paraboloids.

Section 4.H

13. Evaluate the triple integral of the function

$$f(x, y, z) = x^2 y^2 z$$

over the region R bounded by the cone $x^2 + y^2 = xz$ and the planes $z = 0$ and $z = c$. (See Figure 4.29.)

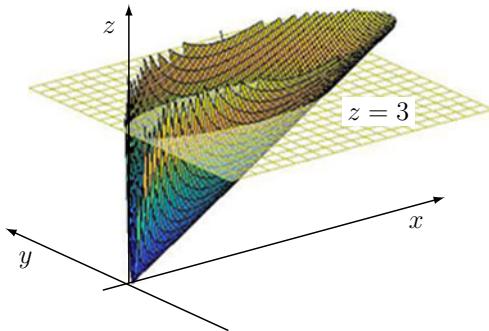


Figure 4.29 The graph of $x^2 + y^2 = xz$ intersected by $z = 3$.

14. Show that

$$\int_0^\infty dx \int_0^\infty xye^{-(x^2+y^2+2xy \cos \alpha)} dy = \frac{\sin \alpha - \alpha \cos \alpha}{4 \sin^3 \alpha}$$

where $0 < \alpha < \pi$.

15. Transform each of the following two iterated integrals into iterated integrals with respect to both cylindrical and spherical coordinates.

$$(i) \int_0^1 dx \int_0^{\sqrt{1-x^2}} dy \int_0^{1+x+y} dz$$

$$(ii) \int_{-1}^1 dx \int_{-\sqrt{1-x^2}}^{\sqrt{1-x^2}} dy \int_{\sqrt{x^2+y^2}}^1 dz$$

16. Determine the volume of the region in the first octant bounded by the surface $x^4 + y^4 + z^4 = c^4$ and the coordinate planes. Express your answer in terms of a Gamma function.

17. Let T be the tetrahedron with vertices $(a, 0, 0)$, $(0, a, 0)$, $(0, 0, a)$, and $(0, 0, 0)$. Show that

$$\iiint_T \frac{x^{1/2}y^{3/2}z^{5/2}}{(x+y+z)^{13/2}} dV = \frac{2\pi a}{3003}.$$

Hint: Find a transformation that maps T to the unit cube. You may need to utilize the integral and numerical properties of Gamma functions.

Section 4.I

18. Suppose $f(x, y, z, t)$ describes the rate of change of electric charge density at point (x, y, z) at time t throughout a volume V .

Give a meaning to the integral

$$I = \iiint_D f(x, y, z, t) dx dy dz dt$$

where $D = \{(x, y, z, t) : x^2 + y^2 + z^2 \leq a^2, 0 \leq t \leq T\}$ and indicate how the integration might be carried out.

Write down the result of the integration in the case that $f(x, y, z, t) = c$, a constant.

19. Show that

$$\begin{aligned} & \int_0^x dx_1 \int_0^{x_1} dx_2 \int_0^{x_2} dx_3 \cdots \int_0^{x_{n-1}} f(x_n) dx_n \\ &= \frac{1}{(n-1)!} \int_0^x (x-t)^{n-1} f(t) dt. \end{aligned}$$

20. Devise suitable n -dimensional spherical polar coordinates to satisfy

$$x_1^2 + x_2^2 + \cdots + x_n^2 = a^2$$

and, using integral properties of Gamma functions, derive the volume of the n -ball.

21. Using the results of the foregoing problem, determine the volume of the n -dimensional ellipsoid:

$$\frac{(x_1 - b_1)^2}{a_1^2} + \frac{(x_2 - b_2)^2}{a_2^2} + \cdots + \frac{(x_n - b_n)^2}{a_n^2} \leq 1.$$

Section 4.J

22. Give reasons why you may decide whether the following integrals are zero or not by inspection only, that is, without any computation.

(a) $\iint_D xe^{-x^2-y^2} dx dy$, where $D = \{(x, y) : |x| \leq 1, |y| \leq 1\}$.

(b) $\iint_D xe^{-x^2-y^2} dx dy$, where $D = \{(x, y) : |x| \leq 1, 0 \leq y \leq 1\}$.

(c) $\iint_D ye^{-x^2-y^2} dx dy$, where $D = \{(x, y) : |x| \leq 1, 0 \leq y \leq 1\}$.

(d) $\iint_D (x - y)e^{-x^2-y^2} dx dy$, where $D = \{(x, y) : |x| \leq 1, |y| \leq 1\}$.

(e) $\iint_D (x - y)e^{-x^2-y^2} dx dy$, where $D = \{(x, y) : |x| \leq 1, 0 \leq y \leq 1\}$.

(f) $\iint_D (x - y)^2 e^{-x^2-y^2} dx dy$, where $D = \{(x, y) : |x| \leq 1, |y| \leq 1\}$.

(g) $\iint_D (x - y)^2 \sin(x - y) dx dy$, where $D = \{(x, y) : |x| \leq 1, |y| \leq 1\}$.

(h) $\iint_D (x - y) \sin(x - y) dx dy$, where $D = \{(x, y) : |x| \leq 1, |y| \leq 1\}$.

(i) $\iint_D (x - y) \sin(x + y) dx dy$, where $D = \{(x, y) : |x| \leq 1, |y| \leq 1\}$.

23. Let S be the unit ball in \mathbb{R}^3 . Show that $\iiint_S f(x, y, z) dV = -4\pi$, where $f(x, y, z) = -3 + 2y + (x^4 + y^6 + z^8) \sin x^3$.



Chapter 5

Vector calculus

The majority of systems that arise in engineering and in the physical sciences fall into one of three camps: kinematic, dynamic, and static systems. Certainly in the first two, but even in the third camp, a system is only partially described by magnitudes of quantities. Systems in motion but also systems in a state of balance or equilibrium can only be completely characterized when directional dependencies are considered. The complete characterization of a system is therefore achieved by quantities that describe both direction and magnitude. These are vector-valued functions which vary with respect to specified independent variables. A force acting on an object is an example of a vector-valued function, as is the object's response in terms of its velocity of motion. Another example, this time of a distributed character, is the flow field of a fluid continuum.

This chapter brings together the concepts that were introduced in Chapters 2 and 4 for scalar-valued functions and extends their applications to functions that are themselves vectors. However, the result is not always a simple generalization from one to many dependent variables. We will discover a range of new results, new concepts, and new features, which hold specifically for vector-valued functions.

5.A Vector-valued functions

We have already named some examples of vector-valued functions, force and velocity, but it helps to have a more general definition that is not tied to a specific application.

Definition 5.1

A vector, \mathbf{f} , is called a vector-valued (m -dimensional) function of the vector variable \mathbf{x} if each of its m components, f_i , are real-valued functions of \mathbf{x} :

$$f_i = f_i(\mathbf{x}); \quad f_i : \mathbb{R}^n \longrightarrow \mathbb{R}, \quad i = 1, 2, \dots, m$$

and

$$\mathbf{f}(\mathbf{x}) = (f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_m(\mathbf{x})).$$

We say that $\mathbf{f} : \mathbb{R}^n \longrightarrow \mathbb{R}^m$ defines a transformation $\mathbf{x} \mapsto \mathbf{y} = \mathbf{f}(\mathbf{x})$:

$$\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \mapsto \begin{pmatrix} y_1 = f_1(\mathbf{x}) \\ y_2 = f_2(\mathbf{x}) \\ \vdots \\ y_m = f_m(\mathbf{x}) \end{pmatrix}$$

Figure 5.1 The vector mapping $\mathbf{x} \mapsto \mathbf{y} = \mathbf{f}(\mathbf{x})$.

The f_i , $i = 1, 2, \dots, m$, in Figure 5.1 are the components of \mathbf{f} in the corresponding orthogonal directions \mathbf{e}_i , $i = 1, 2, \dots, m$ in \mathbb{R}^m (that is, $\mathbf{i}, \mathbf{j}, \mathbf{k}$ in \mathbb{R}^3).

From this most general form we now consider three important specific classes of vector-valued functions.

I. Curves $\mathbf{f} : \mathbb{R} \rightarrow \mathbb{R}^m$

The dependence here is on a single variable. For applications in physics where we use (x, y, z) to denote position in \mathbb{R}^3 , we denote the independent variable by t . The vector function in general defines a transformation $t \mapsto \mathbf{f}(t)$ to a point in \mathbb{R}^m . As the independent variable varies over an interval domain I in \mathbb{R} , the point traces out a *curve* in \mathbb{R}^m , as illustrated in Figure 5.2 for the case $m = 2$.

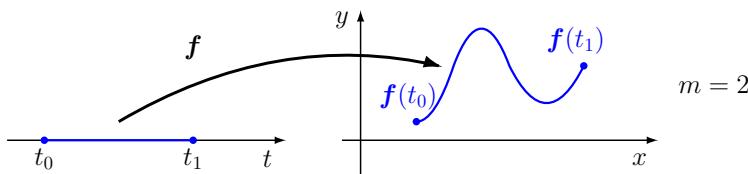


Figure 5.2 A curve in 2-space.

Physically, this mapping describes, for example, the *path* or *trajectory* of a

particle in motion. It is the foundation stone of the field of kinematics.

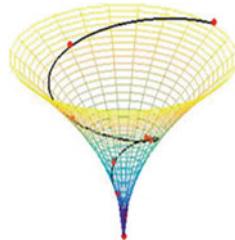


Figure 5.3 The trajectory in \mathbb{R}^3 .

■ **Example 5.1:**

A ball rolling down a funnel at uniform vertical speed, Figure 5.3. The position \mathbf{r} of the ball is time-dependent (t -dependent).

$$\begin{aligned} x &= \left(1 - \frac{t}{4\pi}\right)^2 \cos t \\ y &= \left(1 - \frac{t}{4\pi}\right)^2 \sin t \\ z &= 1 - \frac{t}{4\pi} \\ t &= [0, 4\pi] \end{aligned}$$

■

From the specific perspective of particle trajectories it is usual to replace f with \mathbf{r} to reflect the application of vector functions to position in space. We are thus motivated to write

$$\mathbf{r}(t) = (x_1(t), x_2(t), \dots, x_m(t))$$

to represent a curve in m -space. In \mathbb{R}^3 we use the notation

$$\mathbf{r}(t) = (x(t), y(t), z(t)).$$

The physical world places some restrictions on the types of vector functions that can be used to describe curves in space. The primary restrictions relate to continuity and differentiability.

Theorem 5.1

Suppose I is an open connected interval in \mathbb{R} . The vector $\mathbf{r}(t)$ is a continuous vector function of $t \in I$, with continuous first derivatives (that is, is a C^1 function) if all of its component functions are C^1 functions of $t \in I$.

As with all derivatives, the derivative vector is defined in terms of a converged limiting process, much the same as described in Chapter 1, but now applied to each component of $\mathbf{r}(t)$, and the results are recombined.

The preceding definition means that

$$\begin{aligned} & \lim_{\Delta t \rightarrow 0} \frac{\mathbf{r}(t + \Delta t) - \mathbf{r}(t)}{\Delta t} \\ &= \left(\lim_{\Delta t \rightarrow 0} \frac{x_1(t + \Delta t) - x_1(t)}{\Delta t}, \dots, \lim_{\Delta t \rightarrow 0} \frac{x_m(t + \Delta t) - x_m(t)}{\Delta t} \right) \\ &= \left(\frac{dx_1}{dt}(t), \dots, \frac{dx_m}{dt}(t) \right) \quad \text{— if all } x_i \text{ are } C^1 \\ &\equiv \underbrace{\frac{d\mathbf{r}}{dt}(t)}_{\text{definition}} = \underbrace{\mathbf{v}(t)}_{\text{connection to physics — the velocity of particle motion}} = \mathbf{r}'(t). \end{aligned}$$

Some finite steps in this limit process are shown in Figure 5.4.

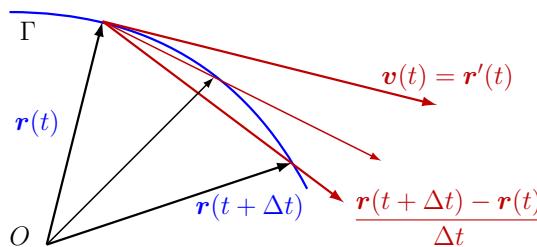


Figure 5.4 The limit process for the curve Γ in \mathbb{R}^n .

Rules for differentiating curve vectors

With the definition of the derivative of a vector-valued function of one variable being a generalization of the derivative of a scalar function of one variable, it is natural to expect that the rules for differentiating scalar functions also generalize. These can be proved by breaking the vector functions down into their components and applying concepts from single-variable calculus to each component, the results are combined thereafter.

For differentiable vector and scalar functions

$$\mathbf{u}, \mathbf{v} : \mathbb{R} \longrightarrow \mathbb{R}^n, \quad \phi : \mathbb{R} \longrightarrow \mathbb{R}, \quad f : \mathbb{R}^n \longrightarrow \mathbb{R},$$

it can be shown that

$$(a) \frac{d}{dt}(\mathbf{u}(t) + \mathbf{v}(t)) = \mathbf{u}'(t) + \mathbf{v}'(t)$$

$$(b) \frac{d}{dt}(\phi(t)\mathbf{u}(t)) = \frac{d}{dt}(\phi(t)u_1(t) + \cdots + \phi(t)u_m(t)) \\ = \phi'(t)\mathbf{u}(t) + \phi(t)\mathbf{u}'(t)$$

$$(c) \frac{d}{dt}(\mathbf{u}(t) \cdot \mathbf{v}(t)) = \mathbf{u}'(t) \cdot \mathbf{v}(t) + \mathbf{u}(t) \cdot \mathbf{v}'(t) \quad \text{— the scalar product rule}$$

$$(d) \frac{d}{dt}(f(\mathbf{u}(t))) = \nabla f(\mathbf{u}(t)) \cdot \mathbf{u}'(t) \quad \text{— the chain rule}$$

$$(e) \frac{d}{dt}(\mathbf{u}(\phi(t))) = \mathbf{u}'(\phi(t))\phi'(t) \quad \text{— another chain rule application}$$

and, for $n = 3$, we also have

$$(f) \frac{d}{dt}(\mathbf{u}(t) \times \mathbf{v}(t)) = \mathbf{u}'(t) \times \mathbf{v}(t) + \mathbf{u}(t) \times \mathbf{v}'(t)$$

— the vector product rule — the order of the vectors is important

☞ Mastery Check 5.1:

Prove the derivative laws (c)–(f).



Remark

* Within the context of particle mechanics one can make a number of other associations, this time for the derivatives of $\mathbf{r}(t)$:

$\mathbf{r}(t)$ — particle position at time t

$\mathbf{r}'(t) = \mathbf{v}(t)$ — particle velocity at time t

$\mathbf{r}''(t) = \mathbf{a}(t)$ — particle acceleration at t

$|\mathbf{v}(t)| = v(t) = \sqrt{v_1^2(t) + \cdots + v_m^2(t)}$ — particle speed at time t

Elementary differential and integral geometry of curves

Describing trajectories, which are fundamental to particle mechanics, is one application of a vector-valued function of one variable. In that case, the physical interest is divided between determining where the particle is at any point in time and computing the particle's position history, *i.e.* the entire path the particle took leading up to that point. Another area of interest focuses on the local properties of the path itself, as a turning and twisting curve in space. While this perspective can potentially increase our understanding of particle dynamics, its real value lies in its application to differential geometry, continuum mechanics of materials (fluid and solid) and general relativity. The fundamental property on which this perspective is based is that of the *tangent vector*.

In order for the tangent vector to be well defined at any point along the curve, the curve must be such that the limit process defining the derivative can be executed (continuity) and then for that limit to exist (differentiability). Beyond these conditions, which are applied in single-variable calculus, we have the further condition that the derivative of the vector is nonzero. This leads to the concept of curve *smoothness*.

Definition 5.2

Suppose I is an open connected interval in \mathbb{R} . A curve Γ , described by a vector function $\mathbf{r}(t)$ for $t \in I$, is called **smooth** if $\mathbf{r}(t)$ is a C^1 -function $\left(\frac{d\mathbf{r}}{dt} \text{ exists}\right)$ and $\frac{d\mathbf{r}}{dt}$ never vanishes for any $t \in I$. $\left[\frac{d\mathbf{r}}{dt} = \mathbf{v} \neq 0\right]$

Remarks

- * The condition $\mathbf{r}'(t) \neq 0$ means that not all components of the vector $\mathbf{r}'(t)$ can vanish simultaneously. It may still arise that one or more components vanish, but not all n at the same t value.
- * In many texts, points $\mathbf{r}(t)$ for $t \in I$ where $\mathbf{r}'(t)$ exists and is nonzero are also called *regular points*. A curve possessing only regular points is therefore also called *regular*.

Definition 5.3

Suppose I is an open connected interval in \mathbb{R} . For a smooth curve Γ defined by $\mathbf{r}(t)$ for $t \in I$, the vector function $\mathbf{v}(t) = \mathbf{r}'(t)$ is a **tangent vector** to Γ at the point $\mathbf{r}(t)$.

The tangent vector is the end result of the limit process illustrated schematically in Figure 5.4. While it already provides new information about the curve (indicating the direction the curve continues with increasing t) much more can be derived from it.

Definition 5.4

For a smooth curve Γ described by a vector function $\mathbf{r}(t)$ defined on an open connected interval I , let $\mathbf{T}(t) = \frac{\mathbf{r}'(t)}{|\mathbf{r}'(t)|}$ be the unit vector in the direction of the tangent vector $\mathbf{r}'(t)$.

Then at those points of Γ for which $\mathbf{T}'(t) \neq 0$ we define the **unit principal normal vector** to Γ as $\mathbf{N}(t) = \frac{\mathbf{T}'(t)}{|\mathbf{T}'(t)|}$, and we define the **binormal vector** to Γ as $\mathbf{B}(t) = \mathbf{T}(t) \times \mathbf{N}(t)$.

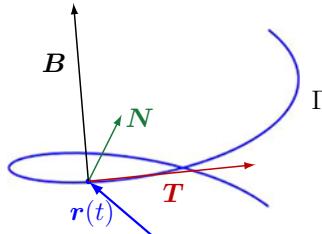


Figure 5.5 Local orthogonal vectors \mathbf{T} , \mathbf{N} , and \mathbf{B} .

Mastery Check 5.2:

Show that for any smooth curve in 3D, $\mathbf{T}(t)$ and $\mathbf{T}'(t)$ are orthogonal.

Hint: Differentiate $\mathbf{T} \cdot \mathbf{T}$ with respect to t (see Rule (c) on Page 227).

What conclusions can you then draw about $\mathbf{N}(t)$ and $\mathbf{N}'(t)$, and $\mathbf{B}(t)$ and $\mathbf{B}'(t)$?



Remarks

- * For each t , the orthogonal vectors $\mathbf{T}(t)$ and $\mathbf{N}(t)$ define the *osculating plane* whose orientation (in 3D) changes with t . The binormal vector $\mathbf{B}(t)$ is orthogonal to the osculating plane.
- * The vectors \mathbf{T} , \mathbf{N} , and \mathbf{B} , shown in Figure 5.5, define a local orthonormal set of vectors, much the same as the vectors $\mathbf{i}, \mathbf{j}, \mathbf{k}$. That is, they define a right-handed coordinate system. However, in contrast to the Cartesian unit vector set, the coordinate system moves and changes direction as t increases. The system is appropriately called a *moving trihedral*. The vectors $\{\mathbf{T}, \mathbf{N}, \mathbf{B}\}$ establish what is called the *Frenet* framework.
- * As a consequence of the foregoing remark, any 3D vector $\mathbf{u}(\mathbf{r}(t))$ relevant to the curve can be expressed in terms of the moving trihedral system,

$$\mathbf{u}(\mathbf{r}(t)) = \alpha(t)\mathbf{T}(t) + \beta(t)\mathbf{N}(t) + \gamma(t)\mathbf{B}(t).$$

In particular, this holds for $\mathbf{T}'(t)$, $\mathbf{N}'(t)$ and $\mathbf{B}'(t)$ as revealed by Definition 5.5 and Mastery Check 5.3.

Definition 5.5

Let Γ be a smooth curve described by $\mathbf{r}(t)$ defined on an open connected interval I . At any point $\mathbf{r}(t)$ on Γ the non-negative function $\kappa(t) = \frac{|\mathbf{T}'(t)|}{|\mathbf{r}'(t)|}$ is called the **curvature** of Γ at that point, and the real-valued function $\tau(t)$, defined such that $\mathbf{B}'(t) = -\tau(t)|\mathbf{r}'(t)|\mathbf{N}(t)$, is called the **torsion** of Γ at that point.

The function $\tau(t)$ gives a measure of the tendency of the curve to twist out of the osculating plane. It can be positive, zero, or negative. On the other hand, the curvature κ , which measures the extent to which the curve is bent at a point, is always non-negative ($\kappa(t) \geq 0$).

Mastery Check 5.3:

Prove that $\mathbf{N}'(t) = -\kappa(t)|\mathbf{r}'(t)|\mathbf{T}(t) + \tau(t)|\mathbf{r}'(t)|\mathbf{B}(t)$.

Hint: Differentiate $\mathbf{N}(t) = \mathbf{B}(t) \times \mathbf{T}(t)$ with respect to t . See Rule (f) on Page 227.



■ Example 5.2:

A straight line, $\mathbf{r}(t) = \mathbf{r}_0 + t\mathbf{u}$, has a constant tangent vector, $\mathbf{T}(t) = \mathbf{u}/|\mathbf{u}|$, and therefore has zero curvature, *i.e.* $\kappa(t) = 0$.

A circle in 2D, $\mathbf{r}(\theta) = \mathbf{r}_0 + a(\cos \theta, \sin \theta)$, or 3D,

$$\mathbf{r}(\theta) = \mathbf{r}_0 + a(\sin \phi \cos \theta, \sin \phi \sin \theta, \cos \phi) \quad \text{with} \quad \phi = f(\theta),$$

has constant curvature, $\kappa(\theta) = 1/a$, the reciprocal of the circle radius. The torsion also reduces to $\tau(\theta) = 0$ since the circle remains in the same plane. (For example, set $b = 0$ in Mastery Check 5.4. See also Supplementary problem 4.) ■

✉ Mastery Check 5.4:

Consider the helical curve in 3D described by the vector function

$$\mathbf{r}(t) = (a \cos t, a \sin t, bt), \quad 0 \leq t \leq 2\pi.$$

Compute $\mathbf{r}'(t)$, $\mathbf{r}''(t)$, $\mathbf{T}(t)$, $\mathbf{N}(t)$, and $\mathbf{B}(t)$.

Verify that $\kappa^2 = \frac{a^2}{(a^2 + b^2)^2}$, and that $\tau(t) = \frac{b}{a^2 + b^2}$ for all t .



From the differential properties of curves we move now to integral properties starting with curve length. Not surprisingly, the length of a curve is also the total path length travelled by a particle along that curve. For a formal justification see Section 5.C.

Definition 5.6

The **arc length** of a smooth curve Γ described by $\mathbf{r}(t)$ defined on an open connected interval I measured from $t = t_0 \in I$ to an arbitrary $t \in I$ is

$$s = \int_{t_0}^t v(\tau) d\tau. \quad (5.1)$$

[\tau is a “dummy” integration variable]

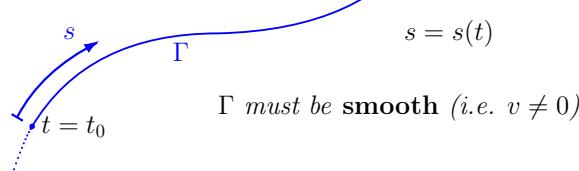
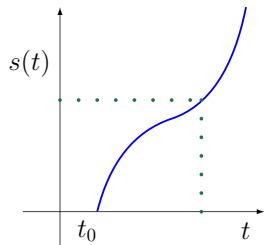


Figure 5.6 Arc length parameter, s .

By differentiating both sides of Equation (5.1) with respect to t , using Leibniz's rule (Page 99), we verify the fundamental theorem of calculus

$$\frac{ds}{dt} = |\mathbf{r}'(t)| = v(t) > 0.$$



s — a one-to-one and onto function of t , which means that t is a function of s , that is, $t = t(s)$

Figure 5.7 Arc length is a one-to-one function of t .

Given the one-to-one relationship between t and s (Figures 5.6 and 5.7), we can parameterize Γ with respect to s instead of t . We then find that

$$\frac{d\mathbf{r}}{ds} = \frac{d\mathbf{r}}{dt} \frac{dt}{ds} = \frac{\mathbf{r}'(t)}{|\mathbf{r}'(t)|}$$

— *the unit tangent vector to Γ (normalized velocity)*

which implies that

$$\left| \frac{d\mathbf{r}}{ds} \right| = 1.$$

— *the “speed” of travel as measured by arc length is constant*

☞ Mastery Check 5.5:

Calculate the length of the 3D helical curve described by the vector function $\mathbf{r}(t) = \cos t \mathbf{i} + \sin t \mathbf{j} + t/\pi \mathbf{k}$, $0 \leq t \leq 2\pi$.



II. Surfaces $\mathbf{r} : \mathbb{R}^2 \rightarrow \mathbb{R}^3$

Another important class of vector-valued functions comprises those that depend on *two* independent variables, which we denote generically by u and v . That is, we consider the class of vector-valued functions $\mathbf{f} : \mathbb{R}^2 \rightarrow \mathbb{R}^m$,

which map points (u, v) in the plane to points in m -space:

$$(u, v) \mapsto \mathbf{f}(u, v) = (f_1(u, v), f_2(u, v), \dots, f_m(u, v)).$$

Since functions from \mathbb{R}^2 to \mathbb{R}^3 have great physical significance and are most readily illustrated, we limit the following discussion to these. As in I, we present the function as a point in \mathbb{R}^3 , $\mathbf{r} = (x, y, z)$, all components of which depend on the two variables (u, v) :

$$\mathbf{r}(u, v) = (x(u, v), y(u, v), z(u, v)).$$

This vector function maps points in the uv -plane to points in 3-space and thereby traces out a *surface* shown here in Figure 5.8.

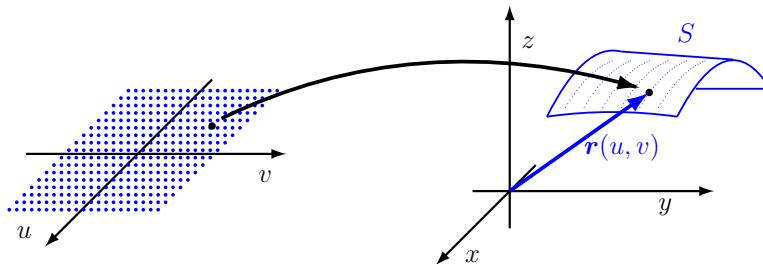


Figure 5.8 A surface in \mathbb{R}^3 .

It is worth reinforcing here that in I we saw that a curve $\mathbf{r}(t)$ in 3D depends on *one* variable (Example 5.2). Here we see that a surface $\mathbf{r}(u, v)$ in 3D depends on *two* variables. The distinction is worth remembering as we shall have occasion to invoke a dimensional reduction (see the Remarks immediately following). Incidentally, the expression above, as well as most of the analysis to follow, supposes that x, y, z are explicit functions of u & v . However, in practice this may not always be the case; a surface may be defined implicitly (see the discussion on implicit functions in Section 2.H).

■ Example 5.3:

The set of parametric equations

$$\left. \begin{array}{l} x = a \cos u \sin v \\ y = a \sin u \sin v \\ z = a \cos v \end{array} \right\} \text{ satisfies } x^2 + y^2 + z^2 = a^2.$$

This is a parametric representation of a sphere of radius a centred at the origin. ■

Remarks

- * Before continuing, the reader might find it useful to revisit the discussions on coordinate systems and visualization of surfaces in Sections 1.D and 1.E.
- * If we keep $u = u_0$ fixed we get $\mathbf{r} = \mathbf{r}(u_0, v)$, a vector function of one variable, v (Figure 5.9). That is, restricting the variable u results in a curve on S called the constant u curve. By the foregoing section this curve has a *tangent vector* given by

$$\mathbf{r}'_v(u_0, v) = \frac{\partial \mathbf{r}}{\partial v}(u_0, v) = \left(\frac{\partial x}{\partial v}(u_0, v), \frac{\partial y}{\partial v}(u_0, v), \frac{\partial z}{\partial v}(u_0, v) \right).$$

- * Similarly, if we keep $v = v_0$ fixed we get $\mathbf{r} = \mathbf{r}(u, v_0)$, a vector function of the single variable u . This too is a curve on S , called the constant v curve. Analogously, this curve has a tangent vector given by

$$\mathbf{r}'_u(u, v_0) = \frac{\partial \mathbf{r}}{\partial u}(u, v_0) = \left(\frac{\partial x}{\partial u}(u, v_0), \frac{\partial y}{\partial u}(u, v_0), \frac{\partial z}{\partial u}(u, v_0) \right).$$

- * If $\mathbf{r}'_u(u_0, v_0) \times \mathbf{r}'_v(u_0, v_0) \neq 0$, which is the case for independent variables, then $\mathbf{r}'_u \times \mathbf{r}'_v$ is a vector normal to S and normal to the *tangent plane* to S (at the point $\mathbf{r}(u_0, v_0)$) spanned by the vectors $\mathbf{r}'_u(u_0, v_0)$ and $\mathbf{r}'_v(u_0, v_0)$.

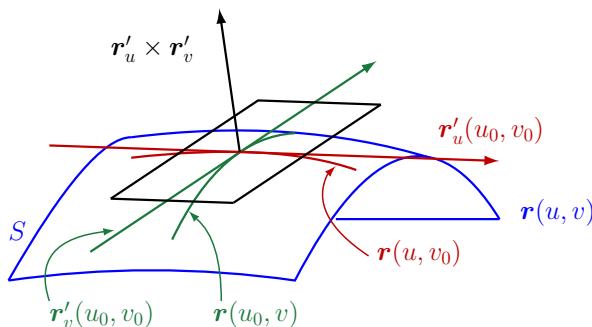


Figure 5.9 The tangent plane to S spanned by tangent vectors \mathbf{r}'_u and \mathbf{r}'_v .

III. The most general case $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$

Although applications arise in more general cases of $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, we derive no benefit by specializing any further. We can instead reflect on the parallels

that may be drawn between “projections” of a more general scenario and the special cases we have already discussed.

A differentiable vector function $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ with components f_i , $i = 1, 2, 3, \dots, m$, which are real-valued functions of \mathbf{x} , is a vector-valued function

$$\mathbf{f}(\mathbf{x}) = (f_1(\mathbf{x}), f_2(\mathbf{x}), f_3(\mathbf{x}), \dots, f_m(\mathbf{x})).$$

When we say that $\mathbf{x} \in \mathbb{R}^n$ is in the domain of \mathbf{f} , where $R_f \subset \mathbb{R}^m$, we mean that \mathbf{x} is in the domain of each component, the scalar functions f_i , $i = 1, 2, \dots, m$. Then if we can assume that each of the f_i is continuous and has continuous partial derivatives, we can compute the gradient for each component,

$$\nabla f_i(\mathbf{x}) = \left(\frac{\partial f_i}{\partial x_1}, \dots, \frac{\partial f_i}{\partial x_n} \right), \quad i = 1, 2, \dots, m.$$

We have met the gradient in Section 2.E, where it was used to determine the rate of change of a scalar function in a specified direction. So, for instance, for the case $n = 2$ and $m = 1$, we had $z = f(x, y)$ describing a surface in 3D space, and the rate of change of z at a point (x_0, y_0) in the direction of unit vector $\mathbf{u} = (u, v)$ was given by

$$D_{\mathbf{u}} f(x_0, y_0) = \nabla f|_0 \cdot \mathbf{u} = \left(\frac{\partial f}{\partial x}|_0 \mathbf{i} + \frac{\partial f}{\partial y}|_0 \mathbf{j} \right) \cdot (u\mathbf{i} + v\mathbf{j}),$$

or, in matrix notation,

$$D_{\mathbf{u}} f(x_0, y_0) = \begin{bmatrix} \frac{\partial f}{\partial x}|_0 & \frac{\partial f}{\partial y}|_0 \end{bmatrix} \begin{pmatrix} u \\ v \end{pmatrix}.$$

With vector functions $\mathbf{f}(\mathbf{x})$ we have the potential to simultaneously find the rate of change of more than one scalar function using matrix multiplication. So for the case $n = 2$ and $m = 3$ we would have

$$D_{\mathbf{u}} \mathbf{f}(x_0, y_0) = \begin{bmatrix} \frac{\partial f_1}{\partial x}|_0 & \frac{\partial f_1}{\partial y}|_0 \\ \frac{\partial f_2}{\partial x}|_0 & \frac{\partial f_2}{\partial y}|_0 \\ \frac{\partial f_3}{\partial x}|_0 & \frac{\partial f_3}{\partial y}|_0 \end{bmatrix} \begin{pmatrix} u \\ v \end{pmatrix}.$$

The 3×2 matrix on the right is an example of a *Jacobian* matrix. The Jacobian lies at the heart of every generalization of single-variable calculus to higher dimensions.

Definition 5.7

Let $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be a C^1 vector-valued function. The $m \times n$ matrix of first derivatives of \mathbf{f} is called the **Jacobian matrix**.

$$D\mathbf{f}(\mathbf{x}) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \frac{\partial f_1}{\partial x_3} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \frac{\partial f_2}{\partial x_3} & \cdots & \frac{\partial f_2}{\partial x_n} \\ \frac{\partial f_3}{\partial x_1} & \frac{\partial f_3}{\partial x_2} & \frac{\partial f_3}{\partial x_3} & \cdots & \frac{\partial f_3}{\partial x_n} \\ \vdots & \vdots & \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \frac{\partial f_m}{\partial x_3} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

*matrix of a vector field, $\mathbf{f} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$
— of special interest to
inhabitants of 3D space.*

gradient of $f_1(\mathbf{x}) : \nabla f_1(\mathbf{x})$.

*— derivative matrix of
vector function \mathbf{f} .*

*the tangent vector to the curve
 $\mathbf{f}(x_1, x_2^0, x_3^0, \dots, x_n^0)$.*

Figure 5.10 The Jacobian matrix of $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$.

Remarks

- * The j^{th} row of the Jacobian in Figure 5.10 is the *gradient* of a scalar function $f_j(\mathbf{x})$, and there are m of them.
- * The i^{th} column of the Jacobian in Figure 5.10 is a *tangent vector* to a curve in \mathbb{R}^m , and there are n of them.
- * The first three rows and first three columns correspond to the derivative matrix of a 3D vector field. (See Section 5.B.)

Example 5.4:

For a scalar function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{x} \mapsto f(\mathbf{x})$, the differential is

$$df = \frac{\partial f}{\partial x_1} dx_1 + \frac{\partial f}{\partial x_2} dx_2 + \cdots + \frac{\partial f}{\partial x_n} dx_n = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right) \begin{pmatrix} dx_1 \\ dx_2 \\ \vdots \\ dx_n \end{pmatrix}.$$

In the case of a vector function $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, $\mathbf{x} \mapsto \mathbf{f}(\mathbf{x})$, the differential generalizes to this:

$$d\mathbf{f} = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{pmatrix} \begin{pmatrix} dx_1 \\ dx_2 \\ \vdots \\ dx_n \end{pmatrix},$$

— a matrix-vector product giving a vector.



■ Example 5.5:

For a scalar function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{x} \mapsto f(\mathbf{x})$, the chain rule gives

$$\frac{\partial f}{\partial u_i} = \frac{\partial f}{\partial x_1} \frac{dx_1}{du_i} + \frac{\partial f}{\partial x_2} \frac{dx_2}{du_i} + \dots + \frac{\partial f}{\partial x_n} \frac{dx_n}{du_i} = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right) \begin{pmatrix} \frac{dx_1}{du_i} \\ \vdots \\ \frac{dx_n}{du_i} \end{pmatrix}.$$

In the case of a vector function $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^m$, $\mathbf{x} \mapsto \mathbf{f}(\mathbf{x})$, the chain rule generalizes to this:

$$\begin{aligned} D(\mathbf{f} \circ \mathbf{x})(\mathbf{u}) &= \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{pmatrix} \begin{pmatrix} \frac{\partial x_1}{\partial u_1} & \dots & \frac{\partial x_1}{\partial u_k} \\ \vdots & & \vdots \\ \frac{\partial x_n}{\partial u_1} & \dots & \frac{\partial x_n}{\partial u_k} \end{pmatrix} \\ &= D\mathbf{f}(\mathbf{x}) \cdot D\mathbf{x}(\mathbf{u}) \end{aligned}$$

— a matrix product giving a matrix.



5.B Vector fields

A 3D vector-valued function of a 3D vector variable,

$$\mathbf{f} : D_f \subseteq \mathbb{R}^3 \longrightarrow \mathbb{R}^3, \mathbf{x} \longmapsto \mathbf{y} = \mathbf{f}(\mathbf{x}),$$

has special significance in physics and engineering, and other applications in the real world. Hence, it is given a special name: a *vector field*. To be explicit, an arbitrary vector field has the form

$$\mathbf{f}(\mathbf{x}) = (f_1(x, y, z), f_2(x, y, z), f_3(x, y, z))$$

$$= f_1(x, y, z)\mathbf{i} + f_2(x, y, z)\mathbf{j} + f_3(x, y, z)\mathbf{k}$$

where the f_1 , f_2 , and f_3 are scalar functions of the three variables x, y, z .

Note that the subscripts “1”, “2”, and “3”, do *not* here refer to partial derivatives, they refer to the components of our vector field.

Unless otherwise stated, we shall assume that the vector fields we work with have *continuous partial derivatives of order $m \geq 2$* . We will often refer to these as *smooth* and presume the component functions are C^2 or better.

Some examples from physics

(i) Gravitational field

The gravitational force per unit mass (G is the gravitational constant) is a 3D vector field (Figure 5.11).

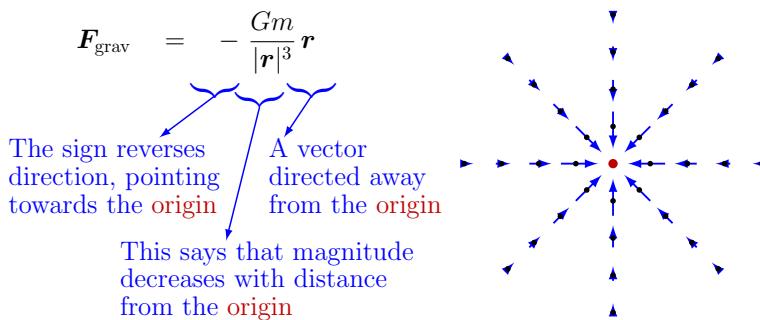


Figure 5.11 The gravitational field of a point mass.

(ii) **Electrostatic fields**

- (a) The 3D electrostatic field intensity (force per unit charge) is expressed in SI units by

$$\mathbf{E} = \frac{q}{4\pi\epsilon_0|\mathbf{r}|^3} \mathbf{r} \quad \text{— due to a point charge } q \text{ at the origin.}$$

The direction depends on the sign of q . The constant ϵ_0 is the “permittivity” of free space.

- (b) The corresponding 2D version has the form,

$$\mathbf{E} = \frac{\rho}{2\pi\epsilon_0|\mathbf{r}|^2} \mathbf{r}.$$

The latter case can be thought of as a field in 3D free space due to a uniformly charged wire of infinite length (Figure 5.12), with the quantity ρ being the charge per unit length of the wire.

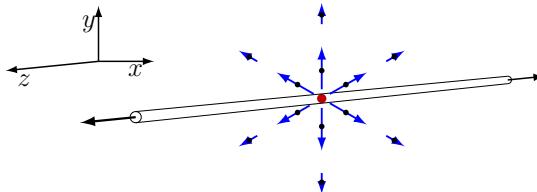


Figure 5.12 The electrostatic field near a charged wire.

(iii) **Gradient field**

A vector field can also be derived from any spatially dependent scalar function by taking its gradient. For example, if $T = T(\mathbf{r})$ is a spatially varying temperature field, then ∇T is the vector field, called the *temperature gradient*.

Gradient fields were discussed at some length in Section 2.E. They will recur often in this chapter, in both the differential and the integral contexts.



Figure 5.13 A field varying in direction and magnitude with position.

Remarks

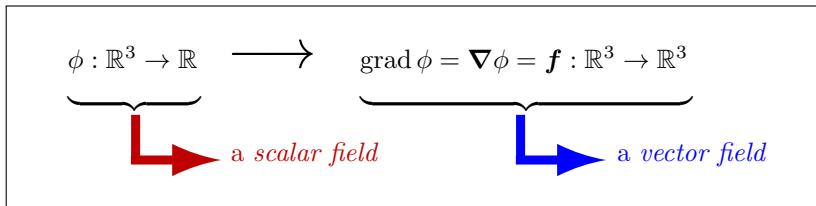
- * By now the reader will appreciate that both the magnitude and direction of a vector field will depend on position (Figure 5.13).
- * As mentioned, Examples (i) and (ii)(a) are 3D vector fields $\mathbb{R}^3 \rightarrow \mathbb{R}^3$.

- * In Example (ii)(b), there are just *two* components to \mathbf{E} : $E_1 = E_1(x, y)$ and $E_2 = E_2(x, y)$. This is an example of a *plane vector field*

$$\mathbf{E} = E_1(x, y)\mathbf{i} + E_2(x, y)\mathbf{j}.$$

These will feature from time to time and specifically in Section 5.F.

- * Example (iii) shows that one gets a vector field by taking the gradient of a real-valued scalar function, $\phi \in C^1(\mathbb{R}^3)$.



More information about vector fields

- **Divergence and curl**

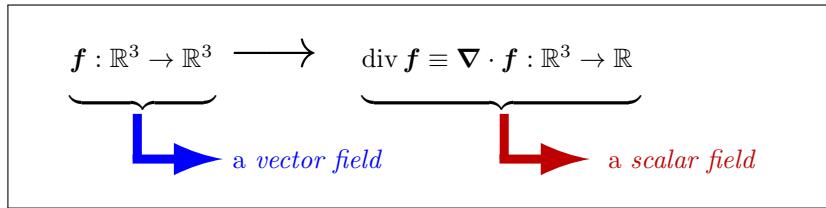
The physical significance of the gradient was explained on Pages 76–78. However, there are two siblings of the gradient that are worth defining now in an operational sense as they too have physical meaning and mathematical utility. These are the *divergence* and *curl* of a vector field.

Definition 5.8

Let $\mathbf{f} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ be a C^1 (at least) vector field. Denote and define the **divergence of \mathbf{f}** by

$$\text{div } \mathbf{f} = \frac{\partial f_1}{\partial x} + \frac{\partial f_2}{\partial y} + \frac{\partial f_3}{\partial z} \equiv \nabla \cdot \mathbf{f}.$$

The divergence thus operates on a vector field to give a *scalar* property of the field that applies locally. That is:



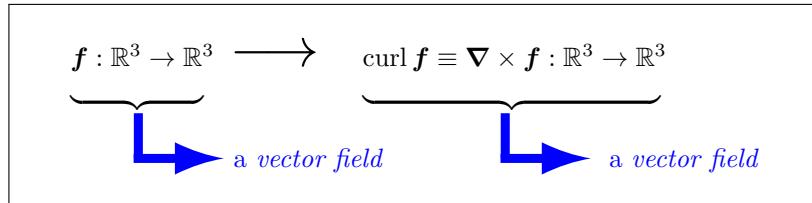
The notation employed in Definition 5.8 involving the gradient operator is standard and is used as a mnemonic to remind us of how the divergence is calculated, which is the expression shown in the central equality. The operator ∇ is treated as a vector in executing the scalar product with \mathbf{F} although no actual component-wise multiplication is performed; instead, the components of \mathbf{F} are partially differentiated with respect to the position variable corresponding to that component.

Definition 5.9

Let $\mathbf{f} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ be a C^1 (at least) vector field. Denote and define the **curl of \mathbf{f}** by

$$\begin{aligned} \operatorname{curl} \mathbf{f} &= \left(\frac{\partial f_3}{\partial y} - \frac{\partial f_2}{\partial z} \right) \mathbf{e}_1 + \left(\frac{\partial f_1}{\partial z} - \frac{\partial f_3}{\partial x} \right) \mathbf{e}_2 + \left(\frac{\partial f_2}{\partial x} - \frac{\partial f_1}{\partial y} \right) \mathbf{e}_3 \\ &= \begin{vmatrix} \mathbf{e}_1 & \mathbf{e}_2 & \mathbf{e}_3 \\ \frac{\partial}{\partial x} & \frac{\partial}{\partial y} & \frac{\partial}{\partial z} \\ f_1 & f_2 & f_3 \end{vmatrix} \equiv \nabla \times \mathbf{f}. \end{aligned}$$

In contrast to the divergence, the curl operates on a vector field to give a *vector* property of the field, also applied locally. That is:



As with the divergence operation, the notation employed in Definition 5.9 is standard form to help us remember how the curl is calculated, which is that shown on the right-hand side of the first equality. The operator ∇ is treated as a vector in the vector product with \mathbf{F} . Again, however, instead of pairwise multiplication, the components of \mathbf{F} are

partially differentiated with respect to the scalar variable corresponding to that location in the vector product.

Later we shall see that the *divergence* of a vector field is a measure of how much the field spreads out locally, while the *curl* of a vector field is a measure of how much the vector field turns or twists locally around an axis parallel to the direction of the vector $\nabla \times \mathbf{F}$.

Note that in the case of two real vectors, switching the order of the vectors in the scalar product doesn't change the result of that product, while switching the orders of the vectors in the vector product results only in a vector pointing in the opposite direction. In the case of the divergence and curl, switching the order of ∇ and \mathbf{F} , *i.e.* writing $\mathbf{F} \cdot \nabla$ and $\mathbf{F} \times \nabla$ makes no sense if these appear in isolation. However, either expression may be legitimate if used in the context of an operation on another scalar field or another vector field, the latter in the case of another scalar or vector product, say. For example, the following makes perfect sense:

$$(\mathbf{F} \cdot \nabla)\phi = \left(F_1 \frac{\partial}{\partial x} + F_2 \frac{\partial}{\partial y} + F_3 \frac{\partial}{\partial z} \right) \phi = F_1 \frac{\partial \phi}{\partial x} + F_2 \frac{\partial \phi}{\partial y} + F_3 \frac{\partial \phi}{\partial z}$$

Analogous meanings can be ascribed to expressions such as $(\mathbf{F} \times \nabla) \cdot \mathbf{G}$ and $(\mathbf{F} \times \nabla) \times \mathbf{G}$. Remember, though, that the operations must be carried out in the correct order.

A large number of general and specific results of applications of the gradient, the divergence, the curl, and their combinations have been established. Some of these are listed on the next page. The reader is invited in Mastery Check 5.6 to prove some of these by direct application of the definitions.

• Field lines

A useful concept particularly in fluid dynamics is that of field lines, also called *stream lines*. These are lines (actually *curves*) whose tangent vectors are parallel to the field vector at those points. That is, given that $\mathbf{r} : \mathbb{R} \rightarrow \mathbb{R}^3, t \mapsto \mathbf{r}(t)$ describes a curve in \mathbb{R}^3 (see page 224), then the field lines of a vector field \mathbf{f} are defined by the equality:

$$\underbrace{\frac{d\mathbf{r}}{dt}}_{\text{tangent vector}} = \lambda(t) \underbrace{\mathbf{f}(\mathbf{r}(t))}_{\text{vector field}}$$

proportionality constant that is a
function of position (not usually of interest)

This vector equation hides three equations that must be solved simultaneously. If we assume the conditions of a smooth curve ($|\mathbf{r}'(t)| \neq 0$) then we deduce that $\lambda \neq 0$. Hence, we can solve the three component equations for λ to get

$$\begin{aligned} \frac{1}{f_1(x, y, z)} \frac{dx}{dt} &= \frac{1}{f_2(x, y, z)} \frac{dy}{dt} = \frac{1}{f_3(x, y, z)} \frac{dz}{dt} \\ \Rightarrow \frac{dx}{f_1(x, y, z)} &= \frac{dy}{f_2(x, y, z)} = \frac{dz}{f_3(x, y, z)}, \quad f_i \neq 0 \end{aligned}$$

Solving these three simultaneous differential equations, if possible, gives the curve described by $\mathbf{r}(t)$.

• Some useful vector identities

Apart from the simple linear identities (e.g. $\nabla(\phi + \psi) = \nabla\phi + \nabla\psi$), the operations of gradient, divergence, and curl obey a number of standard relations when applied to differentiable fields.

Suppose vector fields $\mathbf{f}, \mathbf{g} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ are C^2 ; vector $\mathbf{r} = (x, y, z)$; scalar functions $\phi, \psi : \mathbb{R}^3 \rightarrow \mathbb{R}$ are C^2 ; scalar function $h : \mathbb{R} \rightarrow \mathbb{R}$ is C^1 ; and $\mathbf{c} \in \mathbb{R}^3$ is a constant vector. Then the following identities can be readily derived:

1. $\nabla(\phi\psi) = \psi\nabla\phi + \phi\nabla\psi$
2. $\nabla \cdot (\phi\mathbf{f}) = \phi\nabla \cdot \mathbf{f} + \mathbf{f} \cdot \nabla\phi$
3. $\nabla \times (\phi\mathbf{f}) = \phi\nabla \times \mathbf{f} + \nabla\phi \times \mathbf{f}$
4. $\nabla(\mathbf{f} \cdot \mathbf{g}) = (\mathbf{f} \cdot \nabla)\mathbf{g} + (\mathbf{g} \cdot \nabla)\mathbf{f} + \mathbf{f} \times (\nabla \times \mathbf{g}) + \mathbf{g} \times (\nabla \times \mathbf{f})$
5. $\nabla \cdot (\mathbf{f} \times \mathbf{g}) = \mathbf{g} \cdot (\nabla \times \mathbf{f}) - \mathbf{f} \cdot (\nabla \times \mathbf{g})$
6. $\nabla \times (\mathbf{f} \times \mathbf{g}) = \mathbf{f}(\nabla \cdot \mathbf{g}) - \mathbf{g}(\nabla \cdot \mathbf{f}) + (\mathbf{g} \cdot \nabla)\mathbf{f} - (\mathbf{f} \cdot \nabla)\mathbf{g}$
7. $\nabla \times (\nabla\phi) = \mathbf{0}$
8. $\nabla \cdot (\nabla \times \mathbf{f}) = 0$
9. $\nabla \times (\nabla \times \mathbf{f}) = \nabla(\nabla \cdot \mathbf{f}) - \nabla^2\mathbf{f}$
10. $\nabla \cdot \mathbf{r} = 3$
11. $\nabla h(r) = \frac{dh}{dr} \frac{\mathbf{r}}{r}$
12. $\nabla \cdot (h(r)\mathbf{r}) = 3h(r) + r \frac{dh}{dr}$
13. $\nabla \times (h(r)\mathbf{r}) = \mathbf{0}$

14. $\nabla(\mathbf{c} \cdot \mathbf{r}) = \mathbf{c}$
15. $\nabla \cdot (\mathbf{c} \times \mathbf{r}) = 0$
16. $\nabla \times (\mathbf{c} \times \mathbf{r}) = 2\mathbf{c}$

Mastery Check 5.6:

Confirm the vector identities 3, 5, 7, 8, 12 and 13.



• Conservative fields

From the perspectives of physical significance and mathematical simplicity, one of the most important classes of vector fields is the class of so-called *conservative fields*.

Definition 5.10

A vector field $\mathbf{f} : D \subseteq \mathbb{R}^3 \rightarrow \mathbb{R}^3$ is called a **conservative vector field** if there exists a C^1 scalar function $\phi : D \subseteq \mathbb{R}^3 \rightarrow \mathbb{R}$ such that

$$\mathbf{f}(x, y, z) = \nabla\phi(x, y, z).$$

The function $\phi(x, y, z)$ is called a **scalar potential** of \mathbf{f} .

Although scalar potentials and conservative fields arise in many areas of physics, it is far from true that all vector fields are conservative. That is, it is not generally true that all vector fields can be derived from scalar fields. In the next section we will discover an important and appealing mathematical property of conservative fields. For the moment we focus on the questions of establishing whether a vector field is conservative and if so what its scalar potential is.

To answer these questions we look at the properties of the scalar potential itself. Firstly, we see that for it to be a scalar potential, $\phi : D \rightarrow \mathbb{R}$ must be at least C^1 . Consequently, being a C^1 function, the differential of ϕ can be derived:

$$\begin{aligned} d\phi &= \frac{\partial \phi}{\partial x} dx + \frac{\partial \phi}{\partial y} dy + \frac{\partial \phi}{\partial z} dz \\ &= f_1 dx + f_2 dy + f_3 dz. \end{aligned}$$

The replacement of $\nabla\phi$ with \mathbf{f} in the last equation is valid since $\mathbf{f} = \nabla\phi$ by assumption. The right-hand side of this last equation is thus an

exact differential since it equals $d\phi$. That is, f is conservative if $f_1 dx + f_2 dy + f_3 dz$ is an exact differential. Moreover, if ϕ is a C^2 function, then (Definition 2.8, Page 83) we may conclude that

$$\frac{\partial^2 \phi}{\partial x \partial y} = \frac{\partial^2 \phi}{\partial y \partial x}, \quad \frac{\partial^2 \phi}{\partial x \partial z} = \frac{\partial^2 \phi}{\partial z \partial x}, \quad \frac{\partial^2 \phi}{\partial y \partial z} = \frac{\partial^2 \phi}{\partial z \partial y}.$$

Again, making the substitution $\nabla \phi = \mathbf{f}$, we see that if ϕ is a potential to \mathbf{f} , then the above equations are equivalent to:

$$\frac{\partial f_1}{\partial y} = \frac{\partial f_2}{\partial x}, \quad \frac{\partial f_3}{\partial x} = \frac{\partial f_1}{\partial z}, \quad \frac{\partial f_3}{\partial y} = \frac{\partial f_2}{\partial z}.$$

These are necessary conditions for $\mathbf{f} = \nabla \phi$ to be true and thus for \mathbf{f} to be a conservative field. That is, the components of a conservative field must satisfy these interrelations. (See also Pages 258–260 and 288.)

Mastery Check 5.7:

Determine whether the vector field $\mathbf{f} = \left(xy - \sin z, \frac{x^2}{2} - \frac{e^y}{z}, \frac{e^y}{z^2} - x \cos z \right)$ is conservative, and if so, determine a potential to \mathbf{f} .

Hints:

- (1) For what values of (x, y, z) is \mathbf{f} C^1 ?
- (2) For those points (x, y, z) confirm the necessary conditions for $\mathbf{f} = \nabla \phi$.
- (3) If possible, solve $\nabla \phi = \mathbf{f}$ for a possible scalar potential ϕ .
- (4) Does ϕ exist? Everywhere? Does ϕ have any arbitrary constants? Are they unique?



A scalar potential to a vector field \mathbf{f} is also a real-valued scalar function, and we have seen that (see Section 2.E) the *level surfaces* of $\phi : \phi(x, y, z) = c$ have normal vectors given by $\nabla \phi$. This means that for a conservative field, $\mathbf{f} = \nabla \phi$ is a vector normal to the surface $\phi = c$ at \mathbf{r} . The level surfaces of ϕ (defined in Section 1.F) are called *equipotential* surfaces of \mathbf{f} . See Figure 5.14.

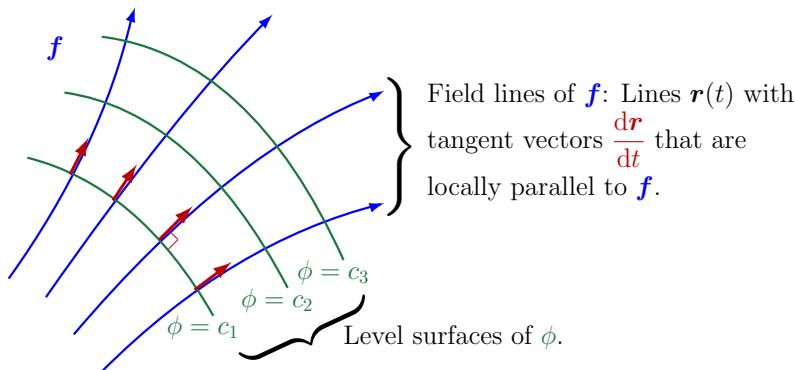


Figure 5.14 Conservative field lines and level surfaces of a scalar potential.

We leave differential vector calculus for the moment and consider some elements of the integral vector calculus, specifically so-called *curve* and *surface* integrals.

5.C Line integrals

Line integrals are somewhat more complicated versions of one-dimensional integrals. As we alluded to in Chapter 4, the most obvious generalization involves replacing the 1D *interval* over which an integral is evaluated, with a one-variable parameterized curve in 3D. As a consequence our intuitive view of a 1D integral as an area under a curve is no longer applicable. We need to rely on a visual idea. Given their respective physical applications, line integrals can be divided into two classes. Within each class it is possible to utilize a specific physical picture to help engender an appreciation for that class of integral.

I. Line integrals of real-valued functions

Physical motivation

In its simplest description these are concerned with one-dimensional integrals of scalar functions to evaluate the total measure of something that is distributed along a curve.

Recall from single-variable calculus that

$$\int_a^b f(x) dx$$

is motivated as giving the area “under” f and “over” straight-line interval $[a, b]$

— this is a geometric interpretation

Suppose instead that we interpret $f(x)$ as a variable linear density of some quantity (e.g. mass/unit length or charge/unit length) defined along the interval $I = [a, b]$. Then

$$\int_I f(x) dx$$

would give the total amount of that quantity, say total mass, attributed to I

— this is a physical interpretation

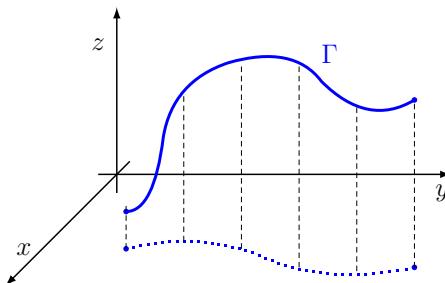


Figure 5.15 The curve Γ in 3D and its projection onto the 2D plane.

Now let’s extend the latter idea to higher dimensions. Suppose we were interested in evaluating the total mass of a nonlinear, one-dimensional object in 3D, one that is represented by the curve, Γ , in \mathbb{R}^3 (Figure 5.15).

What we want is an expression for the total mass of the curved object which corresponds to $\int_I f(x) dx$ for a straight-line interval.

Mathematical construction

Suppose Γ is a finite *smooth* curve in \mathbb{R}^3 . Then, from our discussion in Section 5.A.I, there exists a (non-unique) one-to-one parametrization of Γ ,

$$\mathbf{r}(t) = (x(t), y(t), z(t)), \quad t \in [a, b],$$

such that $\mathbf{r}(t)$ is continuously defined on the finite connected interval $[a, b]$,

and $|\mathbf{r}'(t)| \neq 0$.

Consider a partition of $[a, b]$, $a = t_0 < t_1 < t_2 < \dots < t_{n-1} < t_n = b$, which leads to a set of discrete points on Γ : $\{\mathbf{r}(t_0), \mathbf{r}(t_1), \dots, \mathbf{r}(t_n)\}$, as shown in Figure 5.16.

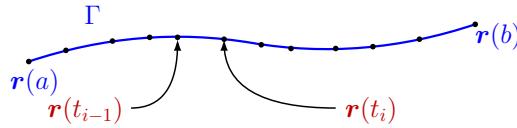


Figure 5.16 A parametrization of Γ .

The leading-order distance between nearest neighbour (in t) points on Γ is $|\Delta \mathbf{r}_i| = |\mathbf{r}(t_i) - \mathbf{r}(t_{i-1})|$ and adding all n such line segment contributions gives

$$\sigma_n = \sum_{i=1}^n |\Delta \mathbf{r}_i| = \sum_{i=1}^n |\mathbf{r}(t_i) - \mathbf{r}(t_{i-1})|$$

as an approximation to the total length of Γ . Since the curve is smooth w.r.t. t , we can apply the mean value theorem to each component function of \mathbf{r} to get

$$|\Delta \mathbf{r}_i| = |\mathbf{r}(t_i) - \mathbf{r}(t_{i-1})| = |(x'(\zeta_i), y'(\eta_i), z'(\xi_i))| \Delta t_i$$

where $\{\zeta_i, \eta_i, \xi_i\}$ are some values of t , not necessarily the same, in the interval (t_{i-1}, t_i) . We then have

$$\sigma_n = \sum_{i=1}^n |(x'(\zeta_i), y'(\eta_i), z'(\xi_i))| \Delta t_i.$$

Now taking the dual limit of an infinite number of partition intervals and of vanishing partition size we get the total length of Γ :

$$|\Gamma| = \lim_{\substack{n \rightarrow \infty \\ \max |\Delta t_i| \rightarrow 0}} \sum_{i=1}^n |\Delta \mathbf{r}_i| = \int_a^b |\mathbf{r}'(t)| dt$$

— provided the limit exists, which it should. (Why?)

We now extend the above argument to include a curve position-dependent function.

Let $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ be a continuous scalar function defined on some domain in \mathbb{R}^3 . When restricted to Γ in that domain it becomes a composite function of one variable:

$$f(x(t), y(t), z(t)) = f(\mathbf{r}(t)).$$

With this restriction f can be thought of as a local linear density of some property of Γ (again, mass or charge per unit length). Multiplying f with a segment of length $|\Delta \mathbf{r}|$ then gives the total amount of that property (mass or charge) possessed by that segment.

By an argument analogous to the above, but applied then to $f(\mathbf{r})|\Delta \mathbf{r}|$, we find that the integral measure of the property represented by f is given by ($t_i^* \in [t_{i-1}, t_i]$)

$$\lim_{\substack{n \rightarrow \infty \\ \max(|\Delta \mathbf{r}|) \rightarrow 0}} \sum_{i=1}^n f(x(t_i^*), y(t_i^*), z(t_i^*)) |\Delta \mathbf{r}_i| = \int_a^b f(x(t), y(t), z(t)) |\mathbf{r}'(t)| dt$$

— provided the limit exists, which it will if f is continuous and $[a, b]$ is bounded (Theorem 3.2).

Remarks

- * It should be obvious from our derivation of curve length and the extension to the total mass that integrals along a given curve of *scalar* functions have *no dependence* on the *direction* taken along the curve. We could have started at either end and arrived at the same result provided the parametrization is defined in a one-to-one manner, increasing from the start position to the end position. This is true of all curve integrals of scalar functions.
- * If we had considered a sub-interval $[a, \tau] \subset [a, b]$, the above argument would have led to Definition 5.6 (Page 231) for the *arc length* $s(\tau)$:

$$s(\tau) = \int_a^\tau |\mathbf{r}'(t)| dt.$$

with $s(a) = 0$ and $s(b) = |\Gamma|$. As shown on Page 231, an application of Leibniz's rule for the derivative of an integral gives

$$\frac{ds}{d\tau} = \left| \frac{d\mathbf{r}}{d\tau} \right|.$$

With a convenient renaming of independent variable $\tau \rightarrow t$ we are led to the differential arc length

$$ds = |\mathbf{d}\mathbf{r}| = \left| \frac{d\mathbf{r}}{dt} \right| dt$$

Consequently, in terms of this differential arc length (equivalent to the parametric integral derived above), the total length of Γ is

$$|\Gamma| = \int_{\Gamma} ds \quad \left(= s(b) = \int_a^b \left| \frac{d\mathbf{r}}{dt} \right| dt \right),$$

while the integral of $f(\mathbf{r})$ over Γ is

$$\int_{\Gamma} f(\mathbf{r}) ds \quad \left(= \int_a^b f(\mathbf{r}(t)) \left| \frac{d\mathbf{r}}{dt} \right| dt \right).$$

- * Neither integral can depend on how we define Γ , *i.e.* they cannot depend on what parametrization we choose.

■ Example 5.6:

Determine the integral $\int_{\Gamma} (xy + y) ds$ where Γ is the path along the 2D curve $y = \sqrt{x}$ from the point $(4, 2)$ to the point $(9, 3)$.

Let $\mathbf{r}(x, y)$ define a point on Γ . We parameterize $\mathbf{r} = (x, y)$ as (t, \sqrt{t}) , with $t : 4 \rightarrow 9$, and $\frac{d\mathbf{r}}{dt} = \left(1, \frac{1}{2\sqrt{t}}\right)$.

Then $\left| \frac{d\mathbf{r}}{dt} \right| = \sqrt{1 + \frac{1}{4t}}$ and

$$\begin{aligned} \int_{\Gamma} (xy + y) ds &= \int_4^9 (t^{3/2} + t^{1/2}) \sqrt{1 + \frac{1}{4t}} dt = \frac{1}{2} \int_4^9 (t + 1) \sqrt{4t + 1} dt \\ &= \left[(t + 1) \frac{1}{6} (4t + 1)^{3/2} \right]_4^9 - \int_4^9 \frac{1}{6} (4t + 1)^{3/2} dt \\ &= \left[(t + 1) \frac{1}{6} (4t + 1)^{3/2} - \frac{1}{60} (4t + 1)^{5/2} \right]_4^9 \\ &= \left[\frac{1}{20} (2t + 3) (4t + 1)^{3/2} \right]_4^9 = \frac{1}{20} (777\sqrt{37} - 187\sqrt{17}). \end{aligned}$$



☛ Mastery Check 5.8:

Evaluate $\int_C x^2 y^2 ds$ where C is the full circle $x^2 + y^2 = 2$.

Hint: A suitable parametrization is with polar coordinates. In this case $\mathbf{r}(t) = \sqrt{2} \cos t \mathbf{e}_1 + \sqrt{2} \sin t \mathbf{e}_2$, $t : 0 \rightarrow 2\pi$.



II. Line integrals of vector-valued functions

In this class of line integrals we not only deal with vector fields as opposed to scalar fields, but we take into consideration the relation between the direction of the integration path taken and the field direction.

Physical motivation

Imagine a block of mass m sitting on a table under gravity. Suppose a horizontal force of magnitude F is needed to slide the block against friction (Figure 5.17(a)).

To move the block a distance d under a constant force F and in the same direction as F , the work done is $W = \mathbf{F} \cdot \mathbf{d}$

If there is no friction, the only external force is gravity and no work is done in sliding the block horizontally.

However, in moving the block directly upwards against gravity (Figure 5.17(b)), the work done is proportional to the force required to overcome gravity. The work done is equal to the force required to overcome gravity (mg) times distance travelled (h): $W = mgh$.

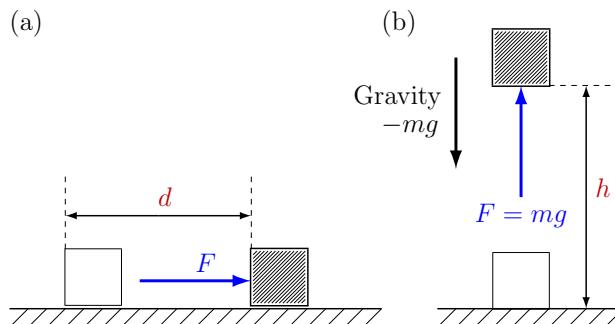


Figure 5.17 (a) Horizontal displacement opposing friction;
(b) Vertical displacement opposing gravity.

If we want to move the block a distance d in a straight line in an arbitrary direction when the only force is gravity, we would then have the scenario pictured in Figure 5.18. The work done depends *only* on the vertical component of the displacement: $W = mgh = mgd \sin \phi$.

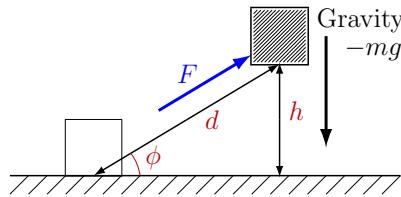


Figure 5.18 Displacement in an oblique direction.

In 3D space coordinates, where the block moves from A to B , and the only opposing force is gravity, the situation can be imagined as in Figure 5.19. The work done is now

$$\begin{aligned}
 W &= mg \cdot \Delta h = mg \cdot d \cdot \sin \phi \\
 &= mg \cdot d \cdot \sin \left(\frac{\pi}{2} - \theta \right) \\
 &= mg \cdot d \cdot \cos \theta \\
 &= \mathbf{F}_g \cdot (\mathbf{r}_B - \mathbf{r}_A).
 \end{aligned}$$

That is, it is simply but significantly the scalar product of the vector force needed to overcome gravity, and the vector displacement.

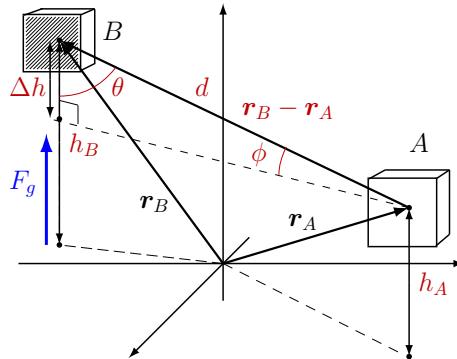


Figure 5.19 Displacement in 3D space.

Generally, the amount of work done depends on the direction of the displacement and the direction of the applied force.

If $F \perp$ displacement, then $W = 0$

If $F \parallel (+)$ displacement, then W is the maximum $= F \cdot d$

As we have seen, for a more general direction of displacement we have the intermediate case:

$$W = \mathbf{F} \cdot \mathbf{d} = Fd \cos \theta.$$

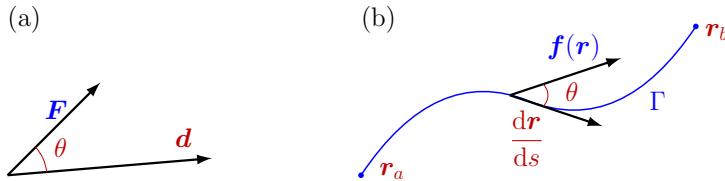


Figure 5.20 (a) Constant force and displacement;
(b) Incremental displacement along a curve.

So far, we have assumed the force and displacement are constant in terms of both direction and magnitude, Figure 5.20(a). But what happens when the force \mathbf{f} is not constant but a function of position, $\mathbf{f} = \mathbf{f}(\mathbf{r})$, and the displacement is along a more variable path Γ ?

Mathematical construction

Suppose $\mathbf{r} = \mathbf{r}(s)$ describes the path to be taken (Figure 5.20(b)). Consider a small segment of the curve, specifically the arc length differential ds at $\mathbf{r}(s)$.

Over the *infinitesimal* curve segment, ds , with unit *tangent vector*, $\frac{d\mathbf{r}}{ds}$, \mathbf{f} is approximately constant, and ds is approximately straight. Then the work increment to leading order in ds will be

$$\begin{aligned} dW &= |\mathbf{f}| \cos \theta ds \\ &= \mathbf{f}(\mathbf{r}) \cdot \frac{d\mathbf{r}}{ds} ds \\ &= \mathbf{f}(\mathbf{r}) \cdot d\mathbf{r}. \end{aligned}$$

This means that the total work done in moving from start to finish along Γ will be the integral

$$W = \int_{\Gamma} dW = \int_{\Gamma} \left(\mathbf{f}(\mathbf{r}) \cdot \frac{d\mathbf{r}}{ds} \right) ds = \int_{\Gamma} \mathbf{f}(\mathbf{r}) \cdot d\mathbf{r}.$$

Remarks

- * $\int_{\Gamma} \mathbf{f} \cdot d\mathbf{r}$ is called the line integral of the vector field \mathbf{f} (or its tangential component) along Γ .
- * Since $\left| \frac{d\mathbf{r}}{ds} \right| = 1$ (see Page 232), then $\mathbf{f}(\mathbf{r}) \cdot \frac{d\mathbf{r}}{ds} = |\mathbf{f}| \cos \theta$.
— the tangential component of \mathbf{f} along Γ

- * In Cartesian form

$$\int_{\Gamma} \mathbf{f} \cdot d\mathbf{r} = \int_{\Gamma} (f_1(x, y, z)dx + f_2(x, y, z)dy + f_3(x, y, z)dz).$$

— this is needed if the path is specified in x, y and z

- * For an arbitrary parametrization, $\mathbf{r}(t)$, with $t \in [a, b]$, $\mathbf{r}_a = \mathbf{r}(a)$, and $\mathbf{r}_b = \mathbf{r}(b)$, we have (for a smooth curve) the very practical expression

$$\int_{\Gamma} \mathbf{f} \cdot d\mathbf{r} = \int_a^b \mathbf{f}(\mathbf{r}(t)) \cdot \frac{d\mathbf{r}(t)}{dt} dt.$$

as a one-dimensional integral!

- * Work is positive or negative depending on the direction of motion; the sign dictating whether energy must be applied to achieve the displacement against the force (negative) or is gained in moving with the force (positive). Therefore, in contrast to line integrals of scalar functions, for line integrals of vector functions we *must always* specify the direction of displacement.

Since $\mathbf{f} \cdot (-d\mathbf{r}) = -\mathbf{f} \cdot d\mathbf{r}$, we find that

$$\int_{\Gamma} \mathbf{f} \cdot d\mathbf{r} = - \int_{-\Gamma} \mathbf{f} \cdot d\mathbf{r}$$

A curve with a specified direction is called *oriented* (Figure 5.21).

— this is important to remember in calculations.



Figure 5.21 Displacement directions along a curve.

- * $\int_{\Gamma} \mathbf{f} \cdot d\mathbf{r}$ is independent of the choice of parametrization — as long as we go in the *same* direction.
- * If $\Gamma = \Gamma_1 \cup \Gamma_2 \cup \Gamma_3 \cup \dots \cup \Gamma_n$, as in Figure 5.22, then

$$\int_{\Gamma} \mathbf{f} \cdot d\mathbf{r} = \int_{\Gamma_1} \mathbf{f} \cdot d\mathbf{r} + \int_{\Gamma_2} \mathbf{f} \cdot d\mathbf{r} + \int_{\Gamma_3} \mathbf{f} \cdot d\mathbf{r} + \dots + \int_{\Gamma_n} \mathbf{f} \cdot d\mathbf{r}.$$

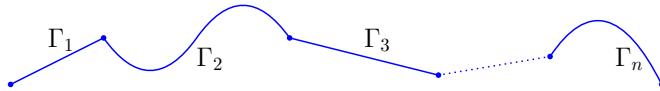


Figure 5.22 A connected, piecewise smooth, curve.

The total line integral over Γ is the sum of the line integrals over the *connected* pieces Γ_i that link the start (of Γ) to the end (of Γ).

- * If Γ is a closed curve (Figure 5.23), that is if $\mathbf{r}_a = \mathbf{r}_b$, then the line integral of \mathbf{f} along Γ is called the *circulation* of \mathbf{f} around Γ :

Circulation of \mathbf{f} around Γ is $\oint_{\Gamma} \mathbf{f} \cdot d\mathbf{r}$.

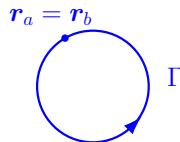


Figure 5.23 The circle: a simple closed curve.

In the next three examples we demonstrate the different ways of evaluating line integrals of vector fields.

■ Example 5.7:

Evaluating line integrals using Cartesian coordinates (x, y, z) .

We wish to evaluate $\int_{\Gamma} \mathbf{f} \cdot d\mathbf{r}$ where $\mathbf{f} = (x + y, y - x)$ and Γ is a curve with endpoints $(1, 1, 0)$ and $(4, 2, 0)$. We will choose two paths:

- $\Gamma = \{(x, y, z) : x = y^2, z = 0\}$ (Figure 5.24(a)), and
- $\Gamma = \Gamma_1 \cup \Gamma_2$, where Γ_1 and Γ_2 are paths of constant y and constant x as shown in Figure 5.24(b).

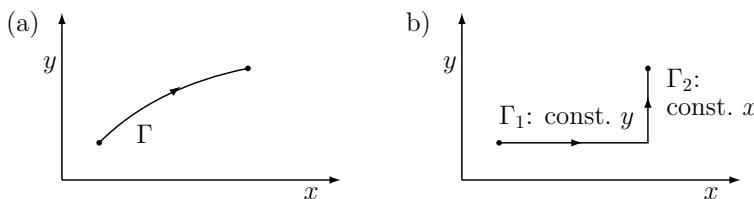


Figure 5.24 The paths (a) $x = y^2, z = 0$, and (b) $\Gamma = \Gamma_1 \cup \Gamma_2$.

- Let $y = t$, $t \in [1, 2]$. Then $x = t^2$, $\mathbf{f} = (f_1, f_2, f_3) = (t^2 + t, t - t^2, 0)$,

$\mathbf{r} = (t^2, t, 0)$, and $\frac{d\mathbf{r}}{dt} = (2t, 1, 0)$.

$$\text{Hence } \int_{\Gamma} \mathbf{f} \cdot d\mathbf{r} = \int_1^2 (t^2 + t, t - t^2, 0) \cdot (2t, 1, 0) dt = \frac{34}{3}.$$

- b) On Γ_1 let $x = t$, $t \in [1, 4]$. Then $y = 1$, $\mathbf{f} = (f_1, f_2, f_3) = (t + 1, 1 - t, 0)$, $\mathbf{r} = (t, 1, 0)$, and $\frac{d\mathbf{r}}{dt} = (1, 0, 0)$. On Γ_2 let $y = t$, $t \in [1, 2]$. Then $x = 4$, $\mathbf{f} = (f_1, f_2, f_3) = (4 + t, t - 4, 0)$, $\mathbf{r} = (4, t, 0)$, and $\frac{d\mathbf{r}}{dt} = (0, 1, 0)$.

$$\begin{aligned} \text{Thus } \int_{\Gamma} \mathbf{f} \cdot d\mathbf{r} &= \int_{\Gamma_1} \mathbf{f} \cdot d\mathbf{r} + \int_{\Gamma_2} \mathbf{f} \cdot d\mathbf{r} \\ &= \int_1^4 (t + 1, 1 - t, 0) \cdot (1, 0, 0) dt + \int_1^2 (4 + t, t - 4, 0) \cdot (0, 1, 0) dt \\ &= \int_1^4 (t + 1) dt + \int_1^2 (t - 4) dt = 8. \end{aligned}$$

■

Example 5.8:

Evaluating line integrals using curve the parametrization $(x(t), y(t), z(t))$.

Evaluate $\int_{\Gamma} \mathbf{f} \cdot d\mathbf{r} = \int_{\Gamma} (xz dx + y^2 dy + x^2 dz)$ from $(0, 0, 0)$ to $(1, 1, 1)$.

along the curves (a) $\mathbf{r}(t) = (t, t, t^2)$, and (b) $\mathbf{r}(t) = (t^2, t, t^2)$ (Figure 5.25).

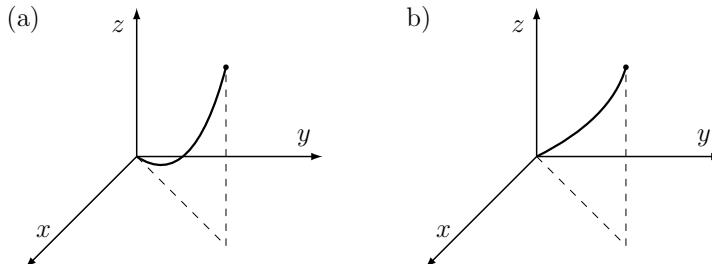


Figure 5.25 The paths (a) $\mathbf{r}(t) = (t, t, t^2)$, and (b) $\mathbf{r}(t) = (t^2, t, t^2)$.

In both cases $t \in [0, 1]$ will work.

- a) We have $x = t$, $y = t$, $z = t^2$, so $\mathbf{f} = (t^3, t^2, t^2)$, $\frac{d\mathbf{r}}{dt} = (1, 1, 2t)$, and

$$\int_C \mathbf{f} \cdot d\mathbf{r} = \int_0^1 (t^3, t^2, t^2) \cdot (1, 1, 2t) dt = \int_0^1 (t^3 + t^2 + 2t^3) dt = \frac{13}{12}.$$

b) We have $x = t^2$, $y = t$, $z = t^2$, so $\mathbf{f} = (t^4, t^2, t^4)$, $\frac{d\mathbf{r}}{dt} = (2t, 1, 2t)$, and

$$\int_C \mathbf{f} \cdot d\mathbf{r} = \int_0^1 (t^4, t^2, t^4) \cdot (2t, 1, 2t) dt = \int_0^1 (2t^5 + t^2 + 2t^5) dt = 1.$$

■

And here is the third example — for the reader to try.

☞ Mastery Check 5.9:

Evaluating curve integrals along curves defined by intersections of surfaces.

Evaluate $\oint_{\Gamma} \mathbf{f} \cdot d\mathbf{r}$ where $\mathbf{f} = y\mathbf{i} + \mathbf{k}$ and Γ is the curve of intersection of the cone

$$z^2 = 2x^2 + 2y^2$$

and the plane $y = z + 1$. See Figure 5.26.

Orientation is counterclockwise seen from $(0, 0, 117)$.

☞

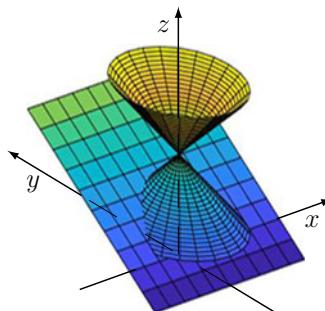


Figure 5.26 Intersection of a plane and a cone.

☞

☞ Mastery Check 5.10:

Calculate the work done by a force $\mathbf{f}(\mathbf{r}) = y\mathbf{e}_1 + 2x\mathbf{e}_2 - z\mathbf{e}_3$ along a helical curve Γ defined by $\mathbf{r}(t) = \cos t \mathbf{e}_1 + \sin t \mathbf{e}_2 + t \mathbf{e}_3$, $t : 0 \rightarrow 2\pi$.

☞

Remarks

- * Each of the line integrals in Examples 5.7 and 5.8 and Mastery Check 5.9 eventually involve some form of parametrization of the curve, with the result that the line integrals reduce to single integrals of functions of one variable, the parameter.
- * In Examples 5.7 and 5.8 the line integrals are of the same field along *different* curves joining the same endpoints. They give *different* results!

To help understand why this may happen it helps to defer to our physical motivation. The work done in going from \mathbf{r}_a to \mathbf{r}_b depends on the path taken; different values result from different routes. Here we can think of the action of *friction*.

The last remark raises an important point and an important question. What conditions do we need to impose on a vector field to get the same result for the line integral?

The satisfactory answer to this question is that the field must be *conservative*!

Theorem 5.2

Suppose $\mathbf{f}(\mathbf{x}) = (f_1(\mathbf{x}), f_2(\mathbf{x}), f_3(\mathbf{x}))$ is a conservative field with potential $\phi(\mathbf{x})$ defined on an open connected domain, D . Then

$$\int_{\Gamma} \mathbf{f} \cdot d\mathbf{r} = \phi(\mathbf{r}_b) - \phi(\mathbf{r}_a)$$

for **every** curve Γ lying entirely in D which joins the points $\mathbf{r}_a \rightarrow \mathbf{r}_b$.

Theorem 5.2 states that the line integral $\int_{\Gamma} \mathbf{f} \cdot d\mathbf{r}$ is independent of the choice of Γ if \mathbf{f} is *conservative* and is dependent only on the endpoints, \mathbf{r}_a and \mathbf{r}_b . To be precise, the integral is determined by, and dependent only on, the values of the potential ϕ at these points.

A proof of the theorem, and our consequential statements, follows readily from our earlier finding that for a conservative field $\mathbf{f} \cdot d\mathbf{r}$ is a perfect differential, $d\phi$. The integral of this perfect differential results in the difference in the values of the potential at the terminal points of the curve Γ .

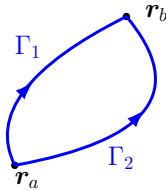


Figure 5.27 Two curves joining the same end points.

Corollary 5.2.1 *For the conditions of Theorem 5.2, let Γ_1 and Γ_2 be two curves in D that join \mathbf{r}_a and \mathbf{r}_b (Figure 5.27). Then*

$$\begin{aligned} \int_{\Gamma_1} \mathbf{f} \cdot d\mathbf{r} &= \int_{\Gamma_2} \mathbf{f} \cdot d\mathbf{r} \\ \implies \int_{\Gamma_1} \mathbf{f} \cdot d\mathbf{r} - \int_{\Gamma_2} \mathbf{f} \cdot d\mathbf{r} &= 0 \\ \implies \int_{\Gamma_1} \mathbf{f} \cdot d\mathbf{r} + \int_{-\Gamma_2} \mathbf{f} \cdot d\mathbf{r} &= 0 \\ \implies \int_{\Gamma_1 \cup (-\Gamma_2)} \mathbf{f} \cdot d\mathbf{r} &= \oint_{\Gamma} \mathbf{f} \cdot d\mathbf{r} = 0. \end{aligned}$$

Thus, the circulation of conservative field \mathbf{f} is zero: $\oint_{\Gamma} \mathbf{f} \cdot d\mathbf{r} = 0$.

In the statement of Theorem 5.2 we impose the condition that D must be a *connected domain*. This means that any pair of points \mathbf{r}_a and \mathbf{r}_b can be joined by a piecewise smooth curve which lies inside D .

Two examples of this, and a counterexample, are shown in Figure 5.28.

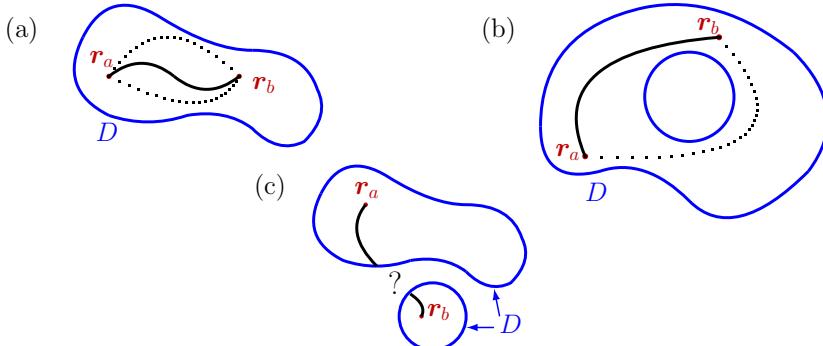


Figure 5.28 Connected domains (a) and (b).
(c) is a disjoint or disconnected domain.

Remark

* Theorem 5.2 and Corollary 5.2.1 are useful even for nonconservative fields — especially for *nearly* conservative fields. (See the next Mastery Check.)

 **Mastery Check 5.11:**

Let $\mathbf{f} = \left(Ax \ln z, By^2 z, \frac{x^2}{z} + y^3 \right)$.

(a) Determine A and B so that \mathbf{f} is conservative.

(b) Evaluate $\int_{\Gamma} (2x \ln z \, dx + 2y^2 z \, dy + y^3 \, dz)$

where Γ is the straight line $(1, 1, 1) \rightarrow (2, 1, 2)$.

 **Mastery Check 5.12:**

Evaluate the line integral of $\mathbf{f}(\mathbf{r}) = x\mathbf{e}_1 + e^y\mathbf{e}_2$ along a curve Γ where Γ is

- (a) a circular arc from $(0, 2)$ to $(2, 0)$, and
 (b) a straight line connecting $(0, 2)$ to $(2, 0)$.

Is \mathbf{f} conservative?



5.D Surface integrals

In complete analogy with line integrals, surface integrals are generalized versions of double integrals! As with line integrals, surface integrals fall into one of two classes dictated by the function being integrated.

Instead of a double integral with an integration over a planar domain, with its interpretation as the volume under the graph of a function, a surface integral is effectively a sum of some quantity that is distributed over the surface itself.

I. Surface integrals of real-valued functions

In this class the surface integrals operate on scalar functions and give the total amount of something that is distributed over a surface S in \mathbb{R}^3 .

Corollary 4.1.1 (2D version) for $f \geq 0$ gives the geometric interpretation of a double integral of f as

$$\iint_D f(x, y) \, dA = \text{volume of body under } f \text{ over } D \text{ (a planar subset of } \mathbb{R}^2\text{)}.$$

If, however, $f(x, y)$ were interpreted as, say, a charge/unit area (surface charge density) defined on \mathbb{R}^2 then the integral would give the total charge on the planar region D in the xy -plane: $\iint_D f(x, y) \, dA = Q$ (Figure 5.29).

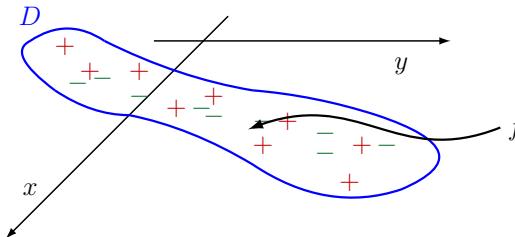


Figure 5.29 Physical interpretation of a double integral over a planar region D .

Now suppose that instead of a region in the plane we have a surface $S \subset \mathbb{R}^3$ defined as a set of points

$$S = \{\mathbf{r}(u, v) \in \mathbb{R}^3 : (u, v) \in D \subset \mathbb{R}^2\}.$$

For the present, we assume there is a one-to-one correspondence between points in D and points on S (this can be altered to consider piecewise mappings of parts of S to subdomains if a one-to-one mapping of the whole of S is not possible).

Moreover, suppose we are given a continuous scalar function $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ defined on a sub-domain in \mathbb{R}^3 . Restricting f to S within that domain, the function becomes a composite function of two variables:

$$f(x(u, v), y(u, v), z(u, v)) = f(\mathbf{r}(u, v)).$$

The function $f(\mathbf{r}(u, v))$ applied to a bounded surface S can be thought of as a summable surface area density (say, charge/unit area or mass/unit area), this time defined on S rather than D . Multiplying f with a segment of area ΔS gives the total amount of charge possessed by that segment.

Similarly to how we reasoned with line integrals, we partition the surface S into a number n of small segments, ΔS_i , $i = 1, 2, \dots, n$. (This can be done

most conveniently by considering a network of intersecting curves of constant u and constant v , as shown in Figure 5.30.) Within each segment a suitable point, (ξ_i, η_i, ζ_i) , is identified.

With this information, an approximation to the total amount carried by S of the property represented by f can be established:

$$Q \approx \sum_{i=1}^n f(\xi_i, \eta_i, \zeta_i) \Delta S_i$$

charge density at point (ξ_i, η_i, ζ_i) in ΔS_i

an element of area on S

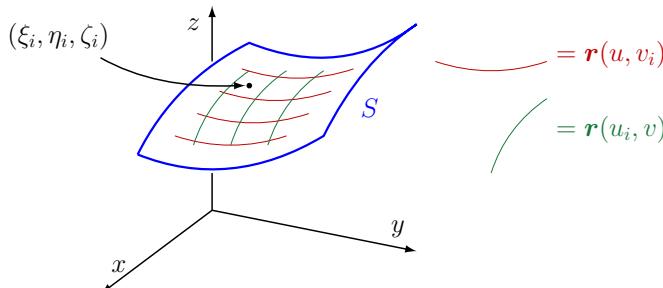


Figure 5.30 A partition of the surface S and an approximation to Q .

We can either assume that f is constant over ΔS_i or adopt a mean-value-theorem argument.

As we have done many times before (see Section 4.A), we refine the partition into smaller and smaller segments and take the limit as $\Delta S_i \rightarrow 0$ and $n \rightarrow \infty$. Then, provided the limit exists the result is what is designated to be the surface integral of f over S :

$$\iint_S f(x, y, z) dS$$

— the limit exists if f is continuous and S is bounded (Theorem 3.2).

Remarks

* It is easily established that (Corollary 4.1.6) if $S = S_1 \cup \dots \cup S_n$ then

$$\iint_S f dS = \sum_{i=1}^N \iint_{S_i} f dS.$$

This is a useful result if S needs partitioning into parts to ensure the existence of one-to-one mappings of those parts onto suitable regions of \mathbb{R}^2 , and the integral over S considered piecewise.

- * If $f = 1$ then, as with Corollary 4.1.3, $\iint_S f dS = \iint dS = \text{area of } S$.

Evaluating a surface integral

• General surface parametrization

To evaluate the surface integral, assuming it is tractable, requires rewriting the area element dS in terms of known quantities that define the surface. For example, recall from Section 5.A, Page 234, that for a surface parameterized by u and v , we can define tangent vectors $\frac{\partial \mathbf{r}}{\partial u}$ and $\frac{\partial \mathbf{r}}{\partial v}$ to constant v and constant u curves, respectively. Then, provided $\frac{\partial \mathbf{r}}{\partial u} \neq 0$, $\frac{\partial \mathbf{r}}{\partial v} \neq 0$, meaning that $\mathbf{r}(u, v_0)$ and $\mathbf{r}(u_0, v)$ are smooth curves, and S is a *smooth* surface, the

tangent vectors $\begin{cases} \frac{\partial \mathbf{r}}{\partial u} \\ \frac{\partial \mathbf{r}}{\partial v} \end{cases}$ lead to differential line elements: $\begin{cases} \frac{\partial \mathbf{r}}{\partial u} du \\ \frac{\partial \mathbf{r}}{\partial v} dv \end{cases}$ and

their cross product leads to the differential area element:

$$dS = \left| \frac{\partial \mathbf{r}}{\partial u} \times \frac{\partial \mathbf{r}}{\partial v} \right| du dv.$$

Consequently, we arrive at a pragmatic representation of the surface integral as a double integral over the planar domain, D , that defines the original S ($dA = du dv$).

$$Q = \iint_S f(\mathbf{r}) dS = \iint_D f(\mathbf{r}(u, v)) \left| \frac{\partial \mathbf{r}}{\partial u} \times \frac{\partial \mathbf{r}}{\partial v} \right| dA$$

$\overbrace{\phantom{\iint_S f(\mathbf{r}) dS}}$
 surface integral
 of f over S

 $\overbrace{\phantom{\iint_D f(\mathbf{r}(u, v)) \left| \frac{\partial \mathbf{r}}{\partial u} \times \frac{\partial \mathbf{r}}{\partial v} \right| dA}}$
 double integral of f . $\left| \dots \times \dots \right|$ over D
 leading to an iterated integral

— DEFINITION
— PRACTICE

 **Mastery Check 5.13:**

For a general C^1 parametrization of a surface $S : \mathbf{r} = \mathbf{r}(u, v)$, derive the expression for $\frac{\partial \mathbf{r}}{\partial u} \times \frac{\partial \mathbf{r}}{\partial v}$.

What then is the expression for dS ?

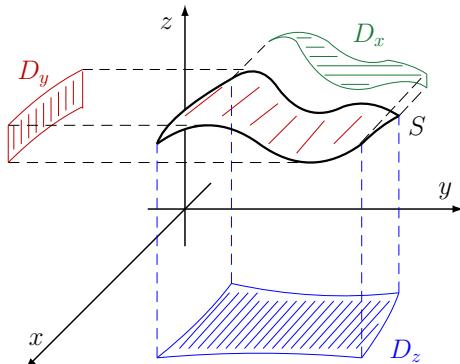


 **Mastery Check 5.14:**

Determine the areal moment of inertia about the z -axis of the parametric surface S given by

$$x = 2uv, \quad y = u^2 - v^2, \quad z = u^2 + v^2; \quad u^2 + v^2 \leq 1.$$

That is, evaluate the integral $\iint_S (x^2 + y^2) dS$.



- (a) $\mathbf{r} = (x, y, q(x, y))$
for the integral over D_z
- (b) $\mathbf{r} = (x, g(x, z), z)$
for the integral over D_y
- (c) $\mathbf{r} = (h(y, z), y, z)$
for the integral over D_x

Figure 5.31 Projections of S onto the three coordinate planes.

• **Cartesian coordinate representation**

The one and same surface *can* be parameterized, at least piecewise, in any of three different ways with respect to different pairwise combinations of the Cartesian coordinates, as shown in Figure 5.31.

From these different representations we have infinitesimal surface elements shown in Figure 5.32 below.

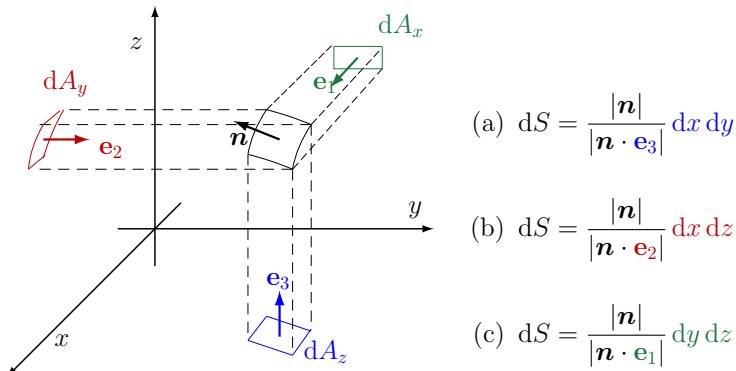


Figure 5.32 Projections of a surface area element dS onto the three coordinate planes.

☞ Mastery Check 5.15:

Evaluate the surface area element dS in the three cases, involving, respectively, the functions q , g , and h , as described in Figure 5.31.



II. Surface integrals of vector fields

The preceding discussion on surface integrals of scalar functions can be extended directly to surface integrals of vector fields to give new vector quantities. Suppose $\mathbf{f} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ is a vector field restricted to a smooth surface S . Then the surface integral of \mathbf{f} over S is simply

$$\iint_S \mathbf{f}(x, y, z) dS = \mathbf{i} \iint_S f_1(x, y, z) dS + \mathbf{j} \iint_S f_2(x, y, z) dS + \mathbf{k} \iint_S f_3(x, y, z) dS.$$

That is, by appealing to the linearity properties of vectors and of the integral operator one can determine the surface integral of a vector field as a vector of surface integrals of the components of that field. However, the following, more important variant of surface integrals of vector fields is the more usual.

One of the most useful qualities of vector fields is that of being able to describe collective movement, in terms of both direction and magnitude of motion. When things are moving, it is often of interest to know how much passes through a given region or across a given area. For example,

- vehicular traffic through an area of a city or an entire city,
- water through a semi-permeable membrane, such as a plant cell wall,

- magnetic flux through a steel cooking pot, or
- X-ray photons through a body.

This “how much” is called the *flux*. To evaluate this quantity we need two pieces of information:

- a vector description of the collective movement — that would be \mathbf{f} .
- a vector description of the region through which the collective movement is to pass — the surface (S) and the normal to the surface (\mathbf{N}).

Let’s start with a simple situation. Suppose \mathbf{f} is a constant vector field with $|\mathbf{f}|$ giving the number of photons per unit area, travelling in a constant direction, $\frac{\mathbf{f}}{|\mathbf{f}|}$. Now suppose we wanted to know the number of photons entering a body. A somewhat simple picture of the situation is shown in Figure 5.33.

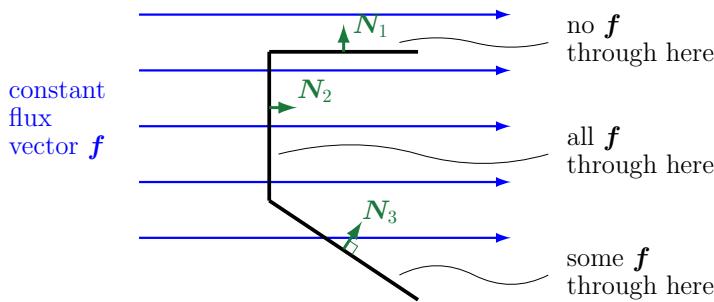


Figure 5.33 Constant flux \mathbf{f} and the boundary of a body.

We see from Figure 5.33 that to determine how much enters a body we need to know both the surface of the body and its relation to the uniform \mathbf{f} . For a flat surface of area A and *unit normal* \mathbf{N} the number of X-ray photons passing through will depend on the area A as well as the latter’s orientation with respect to the direction of \mathbf{f} .

If the vector \mathbf{f} makes an angle θ with the surface normal \mathbf{N} , then we have

$$\text{flux} = |\mathbf{f}| \cos \theta A = (\mathbf{f} \cdot \mathbf{N})A = \mathbf{f} \cdot \mathbf{A}.$$

Notice and remember that \mathbf{N} must be a unit normal! Why? Because the only feature that is needed is the cosine of the angle between \mathbf{f} and \mathbf{A} ; the

scalar product of two vectors involves the magnitudes of *both* vectors, but the flux only requires the magnitude of \mathbf{f} , which will be the case using the scalar product provided $|\mathbf{N}| = 1$. Unless otherwise stated in the remainder, \mathbf{N} will denote a unit normal vector.

Remarks

- * By combining A with \mathbf{N} we have made the surface area a vector, \mathbf{A} .
- * If $\mathbf{f} \cdot \mathbf{A} = -|\mathbf{f}||\mathbf{A}|$, we say the \mathbf{f} is *into* \mathbf{A} .
- * If $\mathbf{f} \cdot \mathbf{A} = |\mathbf{f}||\mathbf{A}|$, we say the \mathbf{f} is *out of* \mathbf{A} .

The above discussion assumes a constant field and planar surfaces — \mathbf{f} does not depend on position, and the surface normals are constant vectors. For non-uniform fields and smooth varying surfaces, Figure 5.34, the same reasoning can be applied, but locally. Suppose again that $\mathbf{f} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ is defined on a surface $S \subset \mathbb{R}^3$, and the surface is parameterized with respect to parameters (u, v) :

$$S = \{\mathbf{r}(u, v) : (u, v) \in D \subset \mathbb{R}^2\}.$$

Let dS be a differential element of area defined at the point $\mathbf{r} \in S$, with unit normal \mathbf{N} ; dS is sufficiently small (infinitesimal, even) that defined over it \mathbf{f} is uniform. Locally at $\mathbf{r} \in S$ the flux of \mathbf{f} through dS is

$$\mathbf{f}(\mathbf{r}) \cdot \mathbf{N}(\mathbf{r}) dS = \mathbf{f}(\mathbf{r}) \cdot d\mathbf{S}.$$

Accumulating such contributions over the entire surface we get

$$\text{total flux} = \iint_S (\mathbf{f} \cdot \mathbf{N}) dS = \iint_S \mathbf{f} \cdot d\mathbf{S}$$

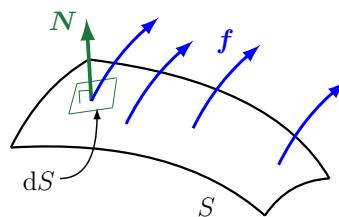


Figure 5.34 Non-constant \mathbf{f} and curved S .

The integral $\iint_S (\mathbf{f} \cdot \mathbf{N}) dS$ is in the form of a surface integral of a scalar field, $\mathbf{f} \cdot \mathbf{N}$. This means that we can use the different approaches to surface integrals discussed on Pages 263 – 264 to evaluate this integral. Two of these means are given in the Remarks below.

Remarks

- * If the surface S is defined by a level set equation (Section 1.F), that is, if

$$S = \{(x, y, z) : \phi(x, y, z) = C\},$$

then

$$\mathbf{N} = \frac{\nabla \phi}{|\nabla \phi|}.$$

- * If the surface S is defined parametrically as $S = \{\mathbf{r}(u, v) : (u, v) \in D\}$, then

$$\mathbf{N} = \frac{\frac{\partial \mathbf{r}}{\partial u} \times \frac{\partial \mathbf{r}}{\partial v}}{\left| \frac{\partial \mathbf{r}}{\partial u} \times \frac{\partial \mathbf{r}}{\partial v} \right|}.$$

— u and v can be Cartesian coordinates

In this case, recalling the formula for dS on Page 263 and inserting both expressions in the flux definition:

$$\iint_S \mathbf{f} \cdot \mathbf{N} dS = \iint_D \mathbf{f} \cdot \underbrace{\left(\frac{\frac{\partial \mathbf{r}}{\partial u} \times \frac{\partial \mathbf{r}}{\partial v}}{\left| \frac{\partial \mathbf{r}}{\partial u} \times \frac{\partial \mathbf{r}}{\partial v} \right|} \right)}_{\mathbf{N}} \underbrace{\left| \frac{\partial \mathbf{r}}{\partial u} \times \frac{\partial \mathbf{r}}{\partial v} \right|}_{dS} du dv,$$

we arrive at the conveniently simpler double integral form,

$$\iint_S \mathbf{f} \cdot \mathbf{N} dS = \iint_D \mathbf{f} \cdot \left(\frac{\partial \mathbf{r}}{\partial u} \times \frac{\partial \mathbf{r}}{\partial v} \right) du dv.$$

- * Clearly, in considering a surface integral of a vector field \mathbf{f} , we presume \mathbf{f} to be defined over all of S .
- * We also assume S is a smooth surface, at least in pieces, so that the existence of a tangent plane at every point of S , at least piecewise, implies in turn that at every point there exists a unit normal vector.

- * An important feature of this entire discussion concerns surface *orientation*: we assume the surface S is piecewise oriented. For a surface to be *orientable* it must possess a *continuously varying normal*, at least piecewise (Figure 5.35).

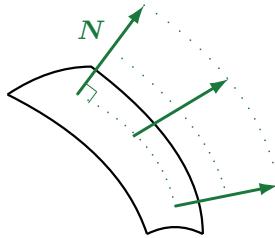


Figure 5.35 Continuously turning normal $\mathbf{N}(\mathbf{r})$.

- * From a very practical perspective the choice of parametrization of the surface determines its orientation; a different choice of parametrization can result in the opposite orientation. For example, $\mathbf{r}'_u \times \mathbf{r}'_v$ determines a specific direction, while $\mathbf{r}'_v \times \mathbf{r}'_u$ gives the opposite direction.
- * By convention the outside of a closed surface is denoted as the positive side, and so we choose a parametrization for a closed surface so that the resulting unit normal vector points to the positive side, i.e. *out* from the region contained within the surface.

For future reference, it is significant to take note that the unit normal \mathbf{N} to a (open) surface also specifies the *orientation of the boundary (and vice versa; see Definition 5.11)*! Conventionally, if \mathbf{N} satisfies the following condition then the (open) surface S is said to be *positively oriented*.

 If the little man in Figure 5.36 walks around the boundary of S so that a vector drawn from his feet to his head points in the same direction as \mathbf{N} and S is on the man's left, then S is said to be *positively oriented*.

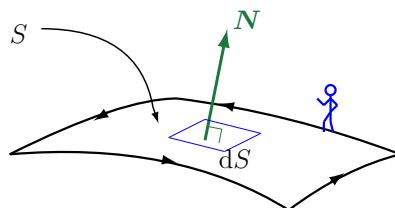


Figure 5.36 Mnemonic for determining surface orientation.

Remarks

- * If S is closed we write

$$\iint_S \mathbf{f} \cdot d\mathbf{S} = \iint_S (\mathbf{f} \cdot \mathbf{N}) dS$$

and, as mentioned already, adopt the convention that \mathbf{N} is the *outward*-pointing unit normal.

- * Why is orientation so important? Consider the uniform field \mathbf{f} and the open surfaces in Figure 5.37. Of course, for practical purposes we need to choose an \mathbf{N} for each surface, and once a choice is made we stick with it. However, the choice of \mathbf{N} determines how one interprets the travel of \mathbf{f} :

If $\mathbf{f} \cdot \mathbf{N} > 0$ we say \mathbf{f} travels *out of* S ,
 if $\mathbf{f} \cdot \mathbf{N} < 0$ we say \mathbf{f} travels *into* S .

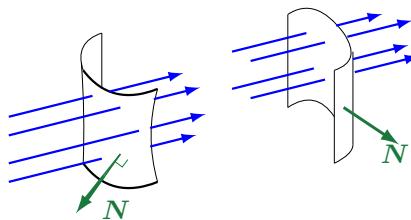


Figure 5.37 The relation between surface normal \mathbf{N} and a vector field.

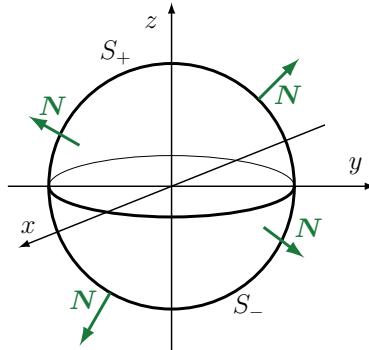


Figure 5.38 The sphere of radius a in Example 5.9.

■ **Example 5.9:**

Determine the flux of the vector field $\mathbf{f} = (x, y, 2z)$ through the surface S , defined by $x^2 + y^2 + z^2 = a^2$.

Note that this surface shown in Figure 5.38 is a closed surface. Hence, we assume the usual convention of taking the *outward*-pointing normal.

We divide the surface into the upper hemisphere,

$$S_+ = \{(x, y, z) : z = g(x, y), 0 \leq x^2 + y^2 \leq a^2\}$$

with an upward pointing normal, and the lower hemisphere,

$$S_- = \{(x, y, z) : z = -g(x, y), 0 \leq x^2 + y^2 \leq a^2\}$$

with a downward pointing normal.

Here, $g(x, y) = \sqrt{a^2 - x^2 - y^2}$.

The flux integral over the closed sphere can then be written as the sum of two flux integrals:

$$\iint_S \mathbf{f} \cdot d\mathbf{S} = \iint_{S_+} \mathbf{f} \cdot d\mathbf{S} + \iint_{S_-} \mathbf{f} \cdot d\mathbf{S}.$$

With the Cartesian variables x and y as parameters, both S_+ and S_- are defined through a one-to-one relationship with points in the planar domain $D = \{(x, y) : 0 \leq x^2 + y^2 \leq a^2\}$.

Hence, the flux integral through S can be rewritten (see Equation (5.2)) as

$$\begin{aligned} \iint_S \mathbf{f} \cdot d\mathbf{S} &= \iint_D \mathbf{f}(x, y, g(x, y)) \cdot \left(\frac{\partial \mathbf{r}_+}{\partial x} \times \frac{\partial \mathbf{r}_+}{\partial y} \right) dx dy \\ &+ \iint_D \mathbf{f}(x, y, -g(x, y)) \cdot \left(-\frac{\partial \mathbf{r}_-}{\partial x} \times \frac{\partial \mathbf{r}_-}{\partial y} \right) dx dy \end{aligned}$$

where $\mathbf{r}_\pm = (x, y, \pm g(x, y))$, and a minus sign is introduced in the second flux integral to provide the correct surface normal direction.

Now, $\frac{\partial \mathbf{r}_\pm}{\partial x} = \left(1, 0, \pm \frac{\partial g}{\partial x}\right)$ and $\frac{\partial \mathbf{r}_\pm}{\partial y} = \left(0, 1, \pm \frac{\partial g}{\partial y}\right)$.

Thus,

$$\left(\frac{\partial \mathbf{r}_\pm}{\partial x} \times \frac{\partial \mathbf{r}_\pm}{\partial y}\right) = \left(\mp \frac{\partial g}{\partial x}, \mp \frac{\partial g}{\partial y}, 1\right),$$

where $\frac{\partial g}{\partial x} = \frac{1}{2}(a^2 - x^2 - y^2)^{-1/2} \cdot (-2x) = \frac{-x}{\sqrt{a^2 - x^2 - y^2}} = \frac{-x}{g(x, y)}$, etc.

This gives the two normal vectors

$$\pm \left(\frac{\partial \mathbf{r}_\pm}{\partial x} \times \frac{\partial \mathbf{r}_\pm}{\partial y}\right) = \left(\frac{x}{g}, \frac{y}{g}, \pm 1\right),$$

and so

$$\begin{aligned} \iint_S \mathbf{f} \cdot d\mathbf{S} &= \iint_D (x, y, 2g) \cdot \left(\frac{x}{g}, \frac{y}{g}, 1\right) dx dy \\ &\quad + \iint_D (x, y, -2g) \cdot \left(\frac{x}{g}, \frac{y}{g}, -1\right) dx dy \\ &= 2 \iint_D \left(\frac{x^2}{g} + \frac{y^2}{g} + 2g\right) dx dy = 2 \iint_D \left(\frac{x^2 + y^2 + 2g^2}{g}\right) dx dy \\ &= 2 \iint_D \left(\frac{2a^2 - x^2 - y^2}{\sqrt{a^2 - x^2 - y^2}}\right) dx dy. \end{aligned}$$

Since D is a disc we can use polar coordinates:

$$x = r \cos \theta, \quad y = r \sin \theta, \quad dx dy = r dr d\theta, \quad x^2 + y^2 = r^2.$$

$$\begin{aligned} \iint_S \mathbf{f} \cdot d\mathbf{S} &= 2 \int_0^{2\pi} d\theta \int_0^a \left(\frac{2a^2 - r^2}{\sqrt{a^2 - r^2}}\right) r dr \\ &= 4\pi \left\{ a^2 \int_0^a \frac{r dr}{\sqrt{a^2 - r^2}} + \int_0^a \sqrt{a^2 - r^2} r dr \right\}. \end{aligned}$$

The final integrals can be evaluated easily using the substitution $u = a^2 - r^2 \implies du = -2r dr$:

$$\begin{aligned} \iint_S \mathbf{f} \cdot d\mathbf{S} &= 4\pi \left\{ a^2 \left[-\frac{1}{2} u^{1/2} \cdot 2 \right]_{a^2}^0 + \left[-\frac{1}{2} u^{3/2} \cdot \frac{2}{3} \right]_{a^2}^0 \right\} \\ &= 4\pi \left\{ a^2 \cdot a + a^3 \frac{1}{3} \right\} = \frac{16}{3} \pi a^3. \end{aligned}$$



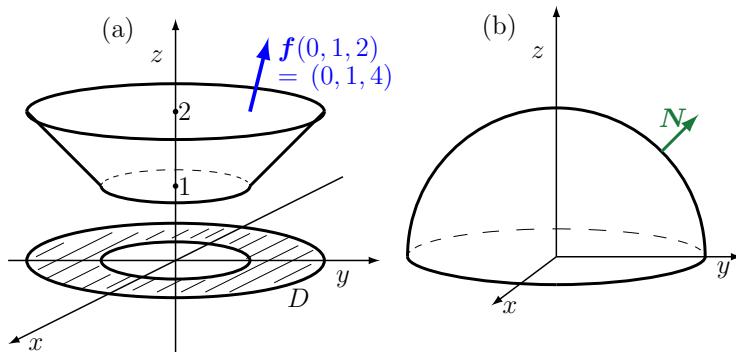


Figure 5.39 (a) The cone S and projection D in MC 5.16;
 (b) The hemisphere in MC 5.17.

Mastery Check 5.16:

Calculate the flux of $\mathbf{f} = (x^2, y^2, z^2)$ through the surface

$$S = \{\mathbf{r} : z^2 = x^2 + y^2, 1 \leq z \leq 2\}$$

in the direction $\mathbf{N} \cdot \mathbf{e}_3 > 0$ (Figure 5.39(a)).



Mastery Check 5.17:

Determine the flux of

$\mathbf{f} = (-y, x, x^2 + z)$ through the surface $S = \{(x, y, z) : x^2 + y^2 + z^2 = 1, z \geq 0\}$ in Figure 5.39(b).



5.E Gauss's theorem

In this section we consider the topic of the net fluxes through closed surfaces, Figure 5.40.

Consider a C^1 field $\mathbf{f} : D_f \subset \mathbb{R}^3 \rightarrow \mathbb{R}^3$ and a *closed* surface S which is the boundary of a *bounded* volume $V \subset D_f$.

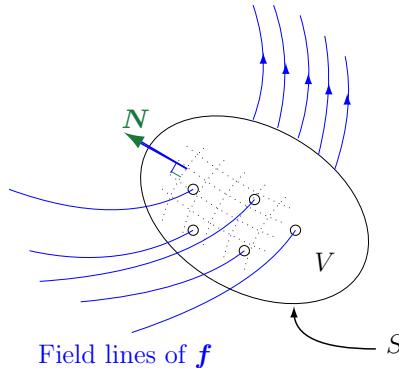


Figure 5.40 The field \mathbf{f} passing into and out of a closed region S .

With S being closed there will be some flux of \mathbf{f} into S and some flux out of S . Consequently, if we take \mathbf{N} to be the outward normal to S , then $\iint_S \mathbf{f} \cdot d\mathbf{S}$ will be a measure of the *net flux of \mathbf{f} out of S* :

$$q = \iint_S \mathbf{f} \cdot \mathbf{N} dS$$

Let's now place this q in a physical setting. Suppose \mathbf{f} describes flow of water. Then

- if $q > 0$: this says that more water flows out of S than in, meaning that there is a *production* of water *inside* V ;
- if $q < 0$: this says that less water flows out of S than in, meaning that there is a *destruction* of water *inside* V ;
- if $q = 0$: this says that what flows in flows out of S , meaning that the amount of water in V is conserved.

From this interpretation one would naturally suspect that q contains information about what occurs *inside* V . That is, there is reason to suspect a relationship of the form

$$\iiint_V (\text{production of } \mathbf{f}) dV = \iint_S \mathbf{f} \cdot d\mathbf{S}.$$

This is in fact exactly what Gauss's¹ theorem states:

¹This book uses this form of the possessive for proper names ending in “s”, such as Gauss and Stokes [21, 22]. Some texts use the form “Gauss' theorem”. The theorems are the same however they are described.

Theorem 5.3

Suppose $\mathbf{f} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ is a C^1 vector field defined on and within a domain V which is bounded by a piecewise smooth closed surface S which in turn has a continuously varying outward unit normal \mathbf{N} . Then

$$\iint_S \mathbf{f} \cdot d\mathbf{S} = \iiint_V \left(\frac{\partial f_1}{\partial x} + \frac{\partial f_2}{\partial y} + \frac{\partial f_3}{\partial z} \right) dV = \iiint_V (\nabla \cdot \mathbf{f}) dV.$$

Remarks

- * The theorem implies that the divergence $\nabla \cdot \mathbf{f}$ is a measure of the *local production* or *local destruction* of \mathbf{f} inside V : $\nabla \cdot \mathbf{f}(\mathbf{x})$ is the *source or sink strength* per unit volume of \mathbf{f} at $\mathbf{x} \in V$.
- * The theorem holds true *only* if \mathbf{N} is the unit normal pointing *away* from region V , as per the examples in Figure 5.41.

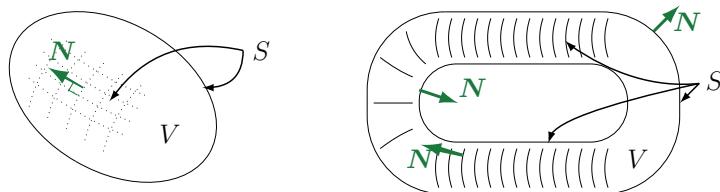


Figure 5.41 Two closed surfaces and their outward normals.

- * The reference to “piecewise smooth” means that S can have edges, just as long as S is closed.
- * Gauss's theorem is useful in rewriting relations involving surface and volume integrals so that all terms can be combined under one integral sign ([1] Chapter 18).
- * If we apply the mean value theorem for multiple integrals to the volume integral in Gauss's theorem we get

$$\iiint_V (\nabla \cdot \mathbf{f}) dV = \nabla \cdot \mathbf{f}(P_0) \cdot V,$$

where $\mathbf{f} \in C^1$ and is bounded, and where P_0 is some point in V . If we now take the limit of this result as $V \rightarrow 0$ and $S \rightarrow 0$ so as to converge

to the single point \mathbf{x} , which will coincide with P_0 , then

$$\nabla \cdot \mathbf{f}(\mathbf{x}) = \lim_{V \rightarrow 0} \frac{1}{V} \iint_S \mathbf{f} \cdot d\mathbf{S},$$

where \mathbf{x} is common to all V and S in this limit.

This result says that the divergence of a vector field \mathbf{f} is the *flux per unit volume* of \mathbf{f} out of a region of vanishing volume.

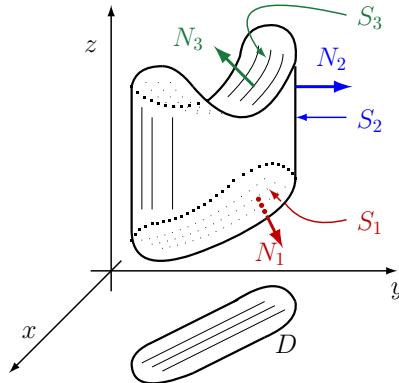


Figure 5.42 V as a z -simple region.

Sketch proof of Gauss's theorem

We start by splitting the surface and volume integrals into their component terms:

$$\begin{aligned} \iint_S \mathbf{f} \cdot \mathbf{N} dS &= \iint_S (f_1 \mathbf{e}_1) \cdot d\mathbf{S} + \iint_S (f_2 \mathbf{e}_2) \cdot d\mathbf{S} + \iint_S (f_3 \mathbf{e}_3) \cdot d\mathbf{S} \\ \iiint_V \nabla \cdot \mathbf{f} dV &= \iiint_V \frac{\partial f_1}{\partial x} dV + \iiint_V \frac{\partial f_2}{\partial y} dV + \iiint_V \frac{\partial f_3}{\partial z} dV. \end{aligned}$$

involves f_1 involves f_2 involves f_3

It is always possible to treat S as the union of piecewise smooth surfaces and V as the union of simple domains (x -simple, y -simple, and z -simple).

Suppose now that V is one of these cases, specifically a z -simple domain, as in Figure 5.42: $V = \{\mathbf{x} : h(x, y) \leq z \leq g(x, y), (x, y) \in D\}$.

$$\text{On } S_1: S_1 = \{ \mathbf{x} : (x, y) \in D, z = h(x, y) \}, \quad \mathbf{N}_1 = \frac{\left(\frac{\partial h}{\partial x}, \frac{\partial h}{\partial y}, -1 \right)}{\sqrt{\dots}}.$$

On S_2 : \mathbf{N}_2 is orthogonal to \mathbf{e}_3 .

$$\text{On } S_3: S_3 = \{ \mathbf{x} : (x, y) \in D, z = g(x, y) \}, \quad \mathbf{N}_3 = \frac{\left(\frac{\partial g}{\partial x}, \frac{\partial g}{\partial y}, 1 \right)}{\sqrt{\dots}}.$$

Now let's work with the f_3 component of the flux integral.

$$\begin{aligned} \iint_S (f_3 \mathbf{e}_3) \cdot d\mathbf{S} &= \iint_{S_1 \cup S_2 \cup S_3} (f_3 \mathbf{e}_3 \cdot \mathbf{N}) dS \\ &= \iint_{S_1} (f_3 \mathbf{e}_3 \cdot \mathbf{N}_1) dS + \iint_{S_2} (f_3 \mathbf{e}_3 \cdot \mathbf{N}_2) dS + \iint_{S_3} (f_3 \mathbf{e}_3 \cdot \mathbf{N}_3) dS \\ &= - \iint_D f_3(x, y, h) dA + \iint_D f_3(x, y, g) dA \\ &= \iint_D \left(\int_{h(x, y)}^{g(x, y)} \frac{df_3}{dz} dz \right) dA = \iiint_V \frac{\partial f_3}{\partial z} dV. \end{aligned}$$

We can manipulate the other cases in analogous ways. Adding these contributions we get the desired result. ■

Example 5.10:

Determine the flux of $\mathbf{f} = (x^3, y^3, z^3)$ out of the sphere in Figure 5.38: $S = \{(x, y, z) : x^2 + y^2 + z^2 = a^2\}$.

The vector field \mathbf{f} is clearly C^1 , and S is closed with a continuously varying normal, \mathbf{N} . So we can use Gauss's theorem.

$$\begin{aligned} \iint_S \mathbf{f} \cdot d\mathbf{S} &= \iiint_V \nabla \cdot \mathbf{f} dV \\ &= \iiint_V (3x^2 + 3y^2 + 3z^2) dV = 3 \iiint_V (x^2 + y^2 + z^2) dV. \end{aligned}$$

We naturally change to spherical coordinates. (See Mastery Check 4.7.)

$$x = \rho \sin \phi \cos \theta, \quad y = \rho \sin \phi \sin \theta, \quad z = \rho \cos \phi, \quad x^2 + y^2 + z^2 = \rho^2.$$

The values $0 \leq \rho \leq a$, $0 \leq \theta \leq 2\pi$, $0 \leq \phi \leq \pi$ cover all points in V , the

interior of S . The differential $dV = \rho^2 \sin \phi d\rho d\theta d\phi$. Thus

$$\begin{aligned}\oint_S \mathbf{f} \cdot d\mathbf{S} &= 3 \iiint_V (x^2 + y^2 + z^2) dV \\ &= 3 \int_0^{2\pi} d\theta \int_0^\pi d\phi \int_0^a \rho^2 \cdot \rho^2 \sin \phi d\rho = 3 \int_0^{2\pi} d\theta \int_0^\pi \frac{1}{5} a^5 \sin \phi d\phi \\ &= \frac{3}{5} a^5 \int_0^{2\pi} d\theta \left[-\cos \phi \right]_0^\pi = \frac{12\pi a^5}{5}.\end{aligned}$$

■

What is good about Gauss's theorem:

Assuming that the conditions of the theorem are satisfied, if we were asked to evaluate $\iiint_V g(x, y, z) dV$, we could instead evaluate $\oint_S \mathbf{f} \cdot d\mathbf{S}$, where $\nabla \cdot \mathbf{f} = g(x, y, z)$. On the other hand, if we were asked to evaluate $\oint_S \mathbf{f} \cdot d\mathbf{S}$, we could instead evaluate $\iiint_V g(x, y, z) dV$.

In these cases we choose whichever integral is easiest to evaluate.

■ Example 5.11:

According to Corollary 4.1.3, volumes of regions are found when $g = 1$. To get $g = 1$, use $\mathbf{f} = \frac{1}{3}\mathbf{r} = \frac{1}{3}(x, y, z)$, or $= (x, 0, 0)$, or $= (0, y, 0)$.

■

What is NOT so good about Gauss's theorem:

The conditions on the theorem are strict. As given, the theorem states that

$$\oint_S \mathbf{f} \cdot d\mathbf{S} = \iiint_V \nabla \cdot \mathbf{f} dV$$

is true *only if*

- a) S is a closed surface. But not every problem that is posed involves a closed surface.

However, if the given S in a problem is not closed we can create a convenient closed surface by complementing S :

$$S \longrightarrow S_c = S \cup S_{\text{xtra}}$$

closed ↗ extra bit ↗

and then apply Gauss's theorem on S_c since, by Corollary 5.2.1, we can argue that

$$\iint_S \mathbf{f} \cdot d\mathbf{S} = \iint_{S_c} \mathbf{f} \cdot d\mathbf{S} - \iint_{S_{\text{xtra}}} \mathbf{f} \cdot d\mathbf{S} = \iiint_V \nabla \cdot \mathbf{f} dV - \iint_{S_{\text{xtra}}} \mathbf{f} \cdot d\mathbf{S}.$$

— by Gauss's theorem on S_c
at the middle step

(See Mastery Check 5.18.)

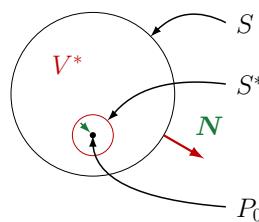


Figure 5.43 Exclusion of the singular point P_0 .

- b) \mathbf{f} is C^1 on S and in V . However, in some problems, \mathbf{f} may be defined on S so that we can evaluate a flux integral, but not defined at some point P_0 within S , so we cannot use Gauss's theorem.

However, we can always *exclude* P_0 by enclosing it with another surface. For example, suppose S is closed and is the boundary of V containing the singular point P_0 . Enclose P_0 in a new closed surface S^* , as shown in Figure 5.43.

Now $S \cup S^*$ is the boundary of a volume V^* , and $P_0 \notin V^*$. Hence, with the singular point now removed we can argue, again by appealing to Corollary 5.2.1, that

$$\iint_S \mathbf{f} \cdot d\mathbf{S} = \iint_{S \cup S^*} \mathbf{f} \cdot d\mathbf{S} - \iint_{S^*} \mathbf{f} \cdot d\mathbf{S} = \iiint_{V^*} \nabla \cdot \mathbf{f} dV - \iint_{S^*} \mathbf{f} \cdot d\mathbf{S}.$$

— by Gauss's theorem on $S \cup S^*$
at the middle step

(See Mastery Check 5.20.)

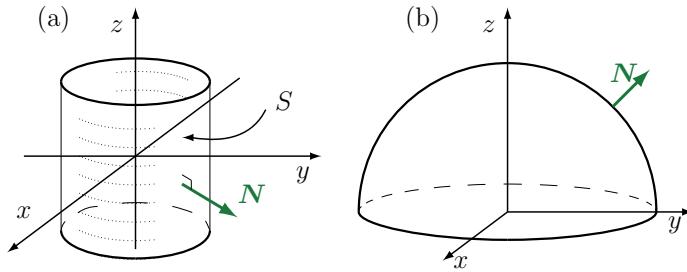


Figure 5.44 (a) The cylinder $x^2 + y^2 = 1$ in MC 5.18;
 (b) The unit hemisphere in MC 5.19.

✍ Mastery Check 5.18:

How to use Gauss's theorem when f is C^1 but S is not closed.

Determine the flux of $\mathbf{f} = (xz^2, x, x^2)$ through the cylinder $S = \{\mathbf{x} : x^2 + y^2 = 1, |z| \leq 1\}$ with outward-pointing normal (Figure 5.44(a)).



✍ Mastery Check 5.19:

Use Gauss's theorem to find the flux of the vector field $\mathbf{f} = (-y, x, x^2 + z)$ through the surface $S = \{(x, y, z) : x^2 + y^2 + z^2 = 1, z \geq 0\}$ (Figure 5.44(b)).

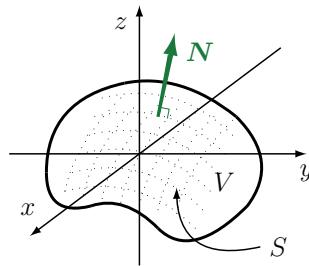


Figure 5.45 An arbitrary closed S enclosing $(0, 0, 0)$.

 **Mastery Check 5.20:**

How to use Gauss's theorem if S is closed but f is not C^1 at some point.

Determine the flux of $\mathbf{f} = -k \frac{\mathbf{r}}{|\mathbf{r}|^3}$ with $\mathbf{r} = (x, y, z)$, through *any* closed surface, Figure 5.45, enclosing a volume V which contains the point $(0, 0, 0)$.



 **Mastery Check 5.21:**

Show that for the conditions of Gauss's theorem and for two C^2 functions $\phi, \psi : \mathbb{R}^3 \rightarrow \mathbb{R}$, the following integral identity holds:

$$\iint_V (\psi \nabla^2 \phi - \phi \nabla^2 \psi) dV = \iint_S (\psi \nabla \phi - \phi \nabla \psi) \cdot \mathbf{N} dS.$$

This is known as Green's second identity ([5]) or Green's formula ([13]).



5.F Green's and Stokes's theorems

Returning to the subject of line integrals of vector fields, but focusing interest on closed curves, we come to discuss Green's theorem and Stokes's theorem, which do for circulation of vector fields (Page 255) what Gauss's theorem does for net fluxes (Page 273). However, before introducing the theorems we shall first cover a few additional curve concepts.

I. Additional notes on curves

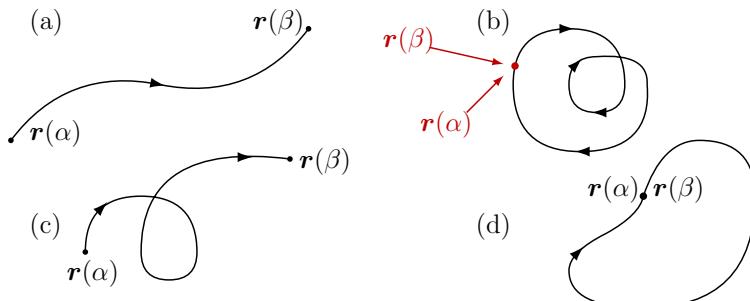


Figure 5.46 Curve types.

- (i) For any smooth curve (Figure 5.46(a)) we can find or construct a *parametrization* $\mathbf{r} = \mathbf{r}(t)$, $\alpha \leq t \leq \beta$.
- (ii) If $\mathbf{r}(\alpha) = \mathbf{r}(\beta)$ (Figure 5.46(b)), then the curve is *closed*.
- (iii) If $\mathbf{r}(t_1) = \mathbf{r}(t_2)$ for some $\alpha < t_1 < t_2 < \beta$ (Figure 5.46(c)), the curve is *self-intersecting*.
- (iv) If $\mathbf{r}(t_1) = \mathbf{r}(t_2)$ for some $\alpha \leq t_1 < t_2 \leq \beta \implies t_1 = \alpha, t_2 = \beta$ (Figure 5.46(d)), the curve is closed and non-intersecting. Such a curve is called a *simple closed curve*.
- (v) Of crucial importance is the obvious fact that a simple closed curve in 3D can be the boundary of many open surfaces, in 3D, as indicated by the three example surfaces in Figure 5.47.

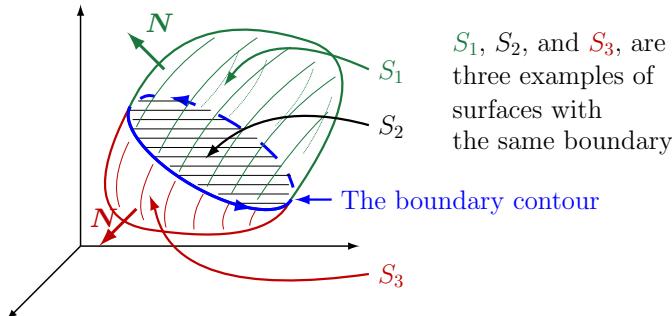


Figure 5.47 One contour, multiple surfaces.

So, in 3D, while an *open* surface has a unique boundary, which is a *simple closed curve*, the converse is not true. In 3D, a *closed* curve can be the boundary to an infinite number of surfaces. This becomes important in the context of Stokes's theorem (Page 287).

- (vi) A parametrization of a curve defines its orientation:

$$\begin{aligned} t : a &\longrightarrow b \implies \mathbf{r}_0 = \mathbf{r}(a) \longrightarrow \mathbf{r}(b) = \mathbf{r}_1 \\ t^* : \alpha &\longrightarrow \beta \implies \mathbf{r}_1 = \mathbf{r}(\alpha) \longrightarrow \mathbf{r}(\beta) = \mathbf{r}_0. \end{aligned}$$

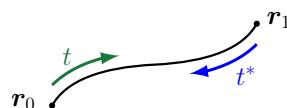


Figure 5.48 One non-closed curve, two alternative directions.

There are two alternatives to specifying orientation. The choice of alternative is obvious for a non-closed curve: Figure 5.48. For a closed curve we need to be more precise, which the next definition and Figure 5.49 attempt to address.

Definition 5.11

A closed curve Γ , the boundary of an oriented surface S , has positive orientation if

(a)  where N is a unit surface normal, and

(b) as we walk around Γ the surface is on our left.

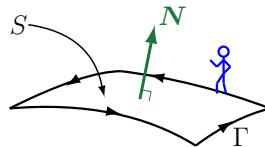


Figure 5.49 A positively oriented boundary.

■ **Example 5.12:**

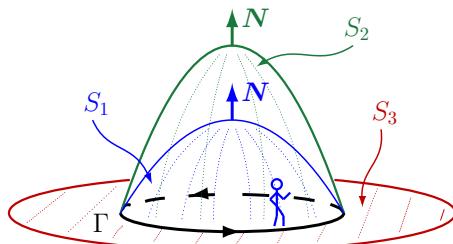


Figure 5.50 Surface and boundary orientation.

Γ in Figure 5.50 is a positively oriented boundary to both S_1 and S_2 , but not S_3 ! ■

II. Green's theorem

In the case of \mathbb{R}^2 , the ambiguity between a closed contour and the surface

it encloses vanishes; a closed positively oriented curve defines and encloses a unique region (Figure 5.51). We take advantage of that fact in the next theorem which is valid for plane vector fields and plane regions (in 2D life is much simpler).

We state and give a sketch proof of this important result, which was discovered by a gentleman having no formal education at the time [18].

Theorem 5.4

Green's theorem:

Let $\mathbf{f}(x, y) = f_1(x, y)\mathbf{e}_1 + f_2(x, y)\mathbf{e}_2$ be a smooth 2D field defined on and within the positively oriented simple closed boundary Γ of a closed and bounded region $D \subset \mathbb{R}^2$. Then

$$\underbrace{\oint_{\Gamma} \mathbf{f} \cdot d\mathbf{r}}_{\text{line integral}} = \oint_{\Gamma} f_1(x, y) dx + f_2(x, y) dy = \underbrace{\iint_D \left(\frac{\partial f_2}{\partial x} - \frac{\partial f_1}{\partial y} \right) dA}_{\text{double integral}}$$

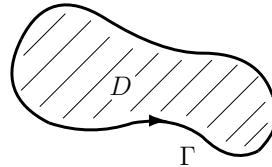


Figure 5.51 Domain D and its positively oriented boundary Γ .

Sketch proof of Green's theorem (for a simply-connected domain)

Suppose D is x -simple and y -simple. As indicated in Figure 5.52 functions ϕ_1 , ϕ_2 , ψ_1 , and ψ_2 can be found such that

$$D = \{(x, y) : a \leq x \leq b, \phi_1(x) \leq y \leq \phi_2(x)\}$$

and

$$D = \{(x, y) : c \leq y \leq d, \psi_1(y) \leq x \leq \psi_2(y)\}.$$

Suppose also that all integrals exist.

We shall show that $\iint_D \frac{\partial f_2}{\partial x} dA = \oint_{\Gamma} f_2 dy$:

Treating D as a y -simple domain we get

$$\begin{aligned}\iint_D \frac{\partial f_2}{\partial x} dA &= \int_c^d dy \int_{\psi_1(y)}^{\psi_2(y)} \frac{\partial f_2}{\partial x} dx \\ &= \int_c^d f_2(\psi_2(y), y) dy - \int_c^d f_2(\psi_1(y), y) dy \\ &= \int_{\psi_2} f_2 dy + \int_{-\psi_1} f_2 dy = \oint_{\Gamma} f_2 dy.\end{aligned}$$

An analogous argument treating D as an x -simple domain will show that $\iint_D \frac{\partial f_1}{\partial y} dA = -\oint_{\Gamma} f_1 dx$, and the theorem is proved. ■

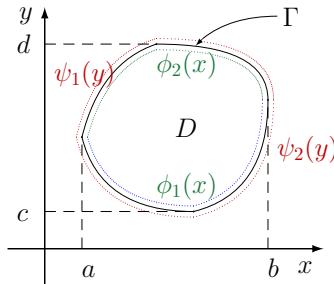


Figure 5.52 D as an x -simple and y -simple domain.

■ Example 5.13:

Use Green's theorem to evaluate a line integral.

Calculate $\oint_{\Gamma} \mathbf{f}(r) \cdot d\mathbf{r}$ where Γ is the boundary of the rectangle $D = \{(x, y) : x \in [0, 4], y \in [0, 1]\}$ when $\mathbf{f} = 3x^2\mathbf{i} - 4xy\mathbf{j}$.

We have

$$\begin{aligned}\oint_{\Gamma} \mathbf{f}(r) \cdot d\mathbf{r} &= \iint_D \left(\frac{\partial f_2}{\partial x} - \frac{\partial f_1}{\partial y} \right) dA \\ &= \int_0^4 dx \int_0^1 (-4y - 0) dy = \int_0^4 (-2) dx = -8.\end{aligned}$$
■

☞ Mastery Check 5.22:

For the function \mathbf{f} and curve Γ as defined in Example 5.13, find $\oint_{\Gamma} \mathbf{f}(r) \cdot d\mathbf{r}$

as a line integral, i.e. without using Green's theorem.



☞ Mastery Check 5.23:

Verify Green's theorem for the curve integral

$$\oint_{\Gamma} (3x^2 - 8y^2) dx + (4y - 6xy) dy$$

where Γ is the boundary of the region D bounded by curves $y = \sqrt{x}$ and $y = x^2$ shown in Figure 5.53.

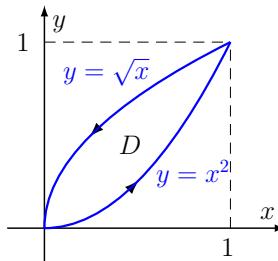


Figure 5.53 D and its boundary: $y = \sqrt{x}$ and $y = x^2$.

Remarks

- * As with Gauss's theorem, this theorem asserts a relation between properties of \mathbf{f} within a region and other properties on the boundary!
- * Special choices of 2D fields allow for the calculation of areas of planar regions using curve integrals, as in the next Mastery Check.

☞ Mastery Check 5.24:

Prove the preceding remark for each of the following cases. $\mathbf{F} = x\mathbf{e}_2$; $\mathbf{F} = -y\mathbf{e}_1$; $\mathbf{F} = \frac{1}{2}(-y\mathbf{e}_1 + x\mathbf{e}_2)$.

Hint: Recall Corollary 4.1.3 which gives a formula for area, and compare this with the statement of Green's theorem.



- * Green's theorem is valid for more complicated regions. For example, it works in the case of annular domains such as shown in Figure 5.54. Take good note of the indicated orientations of the inner and outer boundaries of the annular domain. In each case the curve complies with Definition 5.11 so that Green's theorem remains valid.

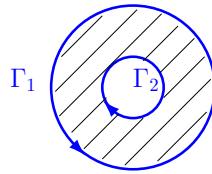


Figure 5.54 The composite boundary $\Gamma = \Gamma_1 \cup \Gamma_2$.

- * If the plane field \mathbf{f} is conservative then

$$\oint_{\Gamma} \mathbf{f} \cdot d\mathbf{r} = 0 \implies \iint_D \left(\frac{\partial f_2}{\partial x} - \frac{\partial f_1}{\partial y} \right) dA = 0,$$

which motivates the following important theorem about conservative fields.

Theorem 5.5

If the plane field $\mathbf{f} = (f_1, f_2)$ satisfies $\frac{\partial f_2}{\partial x} = \frac{\partial f_1}{\partial y}$ in a simply-connected domain D then \mathbf{f} is conservative in D .

III. Stokes's theorem

Generalizing Green's theorem to 3D we get what is commonly referred to as Stokes's theorem, which analogously relates properties of a vector field around a boundary of an (unclosed) surface to properties of the field (in this case the curl of the field) all over that surface.

Theorem 5.6

Suppose S is an oriented surface, piecewise smooth in \mathbb{R}^3 , with unit normal \mathbf{N} , and suppose that S is bounded by a piecewise smooth closed curve Γ with positive orientation. If $\mathbf{f} = (f_1(x, y, z), f_2(x, y, z), f_3(x, y, z))$ is a 3D smooth field defined on S , then

$$\oint_{\Gamma} \mathbf{f} \cdot d\mathbf{r} = \iint_S (\nabla \times \mathbf{f}) \cdot \mathbf{N} dS.$$

Remarks

- * Suppose we apply the mean value theorem to the surface integral in Stokes's theorem:

$$\iint_S (\nabla \times \mathbf{f}) \cdot d\mathbf{S} = (\nabla \times \mathbf{f}) \Big|_{\mathbf{r}_0} \cdot \mathbf{N}(\mathbf{r}_0) \iint_S dS$$

for $\mathbf{f} \in C^1$ and some point $\mathbf{r}_0 \in S$. In the limit, as $S, \Gamma \rightarrow 0$

$$(\nabla \times \mathbf{f}) \Big|_{\mathbf{r}} \cdot \mathbf{N}(\mathbf{x}) = \lim_{S \rightarrow 0} \frac{1}{S} \oint_{\Gamma} \mathbf{f} \cdot d\mathbf{r} \quad (\mathbf{x} \text{ common to all } S \text{ in this limit})$$

— the component of $\nabla \times \mathbf{f}$ normal to S at $\mathbf{x} \in S$
 is the work done per unit area in traversing
 an oriented contour Γ of
 vanishing length

- * If $\mathbf{f} \in C^1(\mathbb{R}^3)$ is conservative, then by Theorem 5.2 $\oint_{\Gamma} \mathbf{f} \cdot d\mathbf{r} = 0$, and
 Stokes's theorem implies that $\iint_S (\nabla \times \mathbf{f}) \cdot d\mathbf{S} = 0$.

- * Conversely, if for $\mathbf{f} \in C^1(\mathbb{R}^3)$ we find that

$$\frac{\partial f_1}{\partial y} = \frac{\partial f_2}{\partial x}, \quad \frac{\partial f_3}{\partial x} = \frac{\partial f_1}{\partial z}, \quad \frac{\partial f_3}{\partial y} = \frac{\partial f_2}{\partial z}$$

in some simply-connected domain, then

$$\operatorname{curl} \mathbf{f} = \nabla \times \mathbf{f} = 0.$$

— the vector field \mathbf{f} is then said to be irrotational
 and Stokes's theorem implies that $\oint_{\Gamma} \mathbf{f} \cdot d\mathbf{r} = 0$ for all closed curves Γ
 in that domain. That is, \mathbf{f} is *conservative*. This important result is an
 extension to 3D of Theorem 5.5 and warrants its own theorem status.

Theorem 5.7

If the C^1 vector field $\mathbf{f} = (f_1, f_2, f_3)$ satisfies $\nabla \times \mathbf{f} = \mathbf{0}$ in a simply-connected domain $D \subset \mathbb{R}^3$, then \mathbf{f} is conservative in D .

Example 5.14:

Verify Stokes's theorem for the C^1 vector field $\mathbf{f} = (xz, xy^2 + 2z, xy + z)$,

where the contour Γ is created by the intersection of the cone $x = 1 - \sqrt{y^2 + z^2}$ and the xy - and yz -planes, such that $0 \leq x \leq 1$, $-1 \leq y \leq 1$, $z \geq 0$. Γ is oriented counterclockwise as seen from $(0, 0, 10)$.

Solution: The geometrical situation is shown in Figure 5.55.

The piecewise curve Γ comprises three part curves: Γ_1 , Γ_2 , Γ_3 , which are oriented as shown.

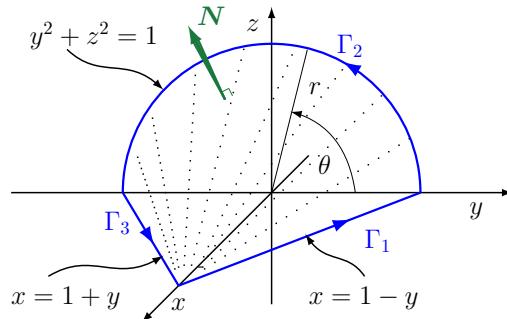


Figure 5.55 The cone $x = 1 - \sqrt{y^2 + z^2}$.

Note that no surface with which to verify the theorem is specified. However, two possibilities come to mind: The curved surface of the cone lying above the xy -plane, and the piecewise combination of the two planar pieces of the xy -plane and the yz -plane within the bottom triangle and semicircle, respectively.

First, the curve integral, $\oint_{\Gamma} \mathbf{f} \cdot d\mathbf{r} = \left(\int_{\Gamma_1} + \int_{\Gamma_2} + \int_{\Gamma_3} \right) \mathbf{f} \cdot d\mathbf{r}$, where

$$\begin{aligned} \int_{\Gamma_1} \mathbf{f} \cdot d\mathbf{r} &= \int_{\Gamma_1} xz \, dx + (xy^2 + 2z) \, dy + (xy + z) \, dz = \int_{\Gamma_1} xy^2 \, dy \, z, \, dz = 0 \\ &= \int_0^1 (1-y)y^2 \, dy = \frac{1}{3} - \frac{1}{4} = \frac{1}{12}, \end{aligned}$$

and, with $x = 0$, $y = \cos \theta$, $z = \sin \theta$,

$$\begin{aligned} \int_{\Gamma_2} \mathbf{f} \cdot d\mathbf{r} &= \int_0^{\pi} \mathbf{f}(r(\theta)) \cdot \frac{dr}{d\theta} d\theta = \int_0^{\pi} \left((2 \sin \theta)(-\sin \theta) + \sin \theta \cos \theta \right) d\theta \\ &= -2 \int_0^{\pi} \sin^2 \theta \, d\theta + \frac{1}{2} \int_0^{\pi} \sin 2\theta \, d\theta = -2 \cdot \frac{1}{2} \cdot \pi = -\pi \end{aligned}$$

$$\int_{\Gamma_3} \mathbf{f} \cdot d\mathbf{r} = \int_{\Gamma_3} xy^2 \, dy = \int_{-1}^0 (1+y)y^2 \, dy = \frac{1}{3} - \frac{1}{4} = \frac{1}{12}.$$

Therefore, adding these contributions we get $\int_{\Gamma} \mathbf{f} \cdot d\mathbf{r} = \frac{1}{12} - \pi + \frac{1}{12} = \frac{1}{6} - \pi$.

Now consider the surface integral, $\iint_S (\nabla \times \mathbf{f}) \cdot \mathbf{N} dS$.

The correct orientation of the cone surface is with \mathbf{N} pointing to +ve z .

$$\nabla \times \mathbf{f} = \begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ \frac{\partial}{\partial x} & \frac{\partial}{\partial y} & \frac{\partial}{\partial z} \\ xz & xy^2 + 2z & xy + z \end{vmatrix} = (x-2)\mathbf{i} + (x-y)\mathbf{j} + y^2\mathbf{k}.$$

Consider the parametrization for a cone

$$x = 1 - r, \quad y = r \cos \theta, \quad z = r \sin \theta; \quad D = \{0 \leq \theta \leq \pi, 0 \leq r \leq 1\}.$$

$$\frac{\partial \mathbf{x}}{\partial r} \times \frac{\partial \mathbf{x}}{\partial \theta} = \begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ -1 & \cos \theta & \sin \theta \\ 0 & -r \sin \theta & r \cos \theta \end{vmatrix} = r\mathbf{i} + r \cos \theta \mathbf{j} + r \sin \theta \mathbf{k}.$$

$$\begin{aligned} \iint_S (\nabla \times \mathbf{f}) \cdot \mathbf{N} dS &= \iint_D (\nabla \times \mathbf{f}) \cdot \left(\frac{\partial \mathbf{x}}{\partial r} \times \frac{\partial \mathbf{x}}{\partial \theta} \right) dr d\theta \quad (\text{see Page 268}) \\ &= \int_0^\pi \int_0^1 (1 - r - 2, 1 - r - r \cos \theta, r^2 \cos^2 \theta) \cdot (r, r \cos \theta, r \sin \theta) dr d\theta \\ &= \int_0^\pi \int_0^1 \left(-r(1+r) + r \cos \theta(1-r-r \cos \theta) + r^3 \cos^2 \theta \sin \theta \right) dr d\theta \\ &= \int_0^\pi \left(-\frac{5}{6} + \frac{1}{6} \cos \theta - \frac{1}{3} \cos^2 \theta + \frac{1}{4} \cos^2 \theta \sin \theta \right) d\theta. \end{aligned}$$

So,

$$\iint_S (\nabla \times \mathbf{f}) \cdot \mathbf{N} dS = -\frac{5\pi}{6} + 0 - \frac{1}{3} \cdot \frac{1}{2} \cdot \pi + \frac{1}{4} \cdot \frac{1}{3} \cdot 2 = \frac{1}{6} - \pi,$$

and the theorem is verified.

The reader should redo the surface integral calculation using the piecewise planar combination alternative.

Be mindful of choosing the correct surface orientation.



What is good about Stokes's theorem:

Assuming that the conditions of the theorem are satisfied, if we were asked

to evaluate the work $\oint_{\Gamma} \mathbf{f} \cdot d\mathbf{r}$ around a closed contour, we could instead evaluate $\iint_S (\nabla \times \mathbf{f}) \cdot d\mathbf{S}$ through any convenient surface that has Γ as its boundary.

Alternatively, if we were asked to evaluate $\iint_S (\nabla \times \mathbf{f}) \cdot d\mathbf{S}$ through a given surface S , we could instead evaluate the flux integral through any surface that has the same curve Γ as its boundary. — there is an infinity of choices

- there is an infinity of choices provided f is not singular

Or, we could evaluate $\oint_{\Gamma} \mathbf{f} \cdot d\mathbf{r}$, if \mathbf{f} could be deduced from the given expression $\nabla \times \mathbf{f}$.

In these cases we choose whichever integral is easiest to evaluate. (See Mastery Check 5.25.)

What is NOT so good about Stokes's theorem:

The conditions on the theorem are strict. As given, the theorem states that

$$\oint_{\Gamma} \mathbf{f} \cdot d\mathbf{r} = \iint_S (\nabla \times \mathbf{f}) \cdot d\mathbf{S}$$

is true *only if*

- (a) Γ is a closed contour. Unfortunately, not every problem posed involves a closed curve. However, if the given curve Γ is not closed, we make up a convenient closed contour by complementing Γ :

We can then apply Stokes's theorem on Γ_c since, by Corollary 5.2.1, it is reasonable to argue that

$$\int_{\Gamma} \mathbf{f} \cdot d\mathbf{r} = \oint_{\Gamma_c} \mathbf{f} \cdot d\mathbf{r} - \int_{\Gamma_{\text{extra}}} \mathbf{f} \cdot d\mathbf{r} = \iint_S (\nabla \times \mathbf{f}) \cdot d\mathbf{S} - \int_{\Gamma_{\text{extra}}} \mathbf{f} \cdot d\mathbf{r}.$$

— by Stokes's theorem on Γ_c
at the middle step

- (b) \mathbf{f} is non-singular on S or Γ . If \mathbf{f} is singular somewhere, then it is harder to deal with, and success will depend on the problem.

On the other hand, if we are asked to evaluate $\oint_{\Gamma} \mathbf{f} \cdot d\mathbf{r}$, and \mathbf{f} is singular on the given S , we can introduce an extra closed contour Γ_ϵ of small size ϵ around the singular point on S so as to exclude it: $\Gamma \rightarrow \Gamma_c = \Gamma \cup \Gamma_\epsilon$. Introducing this new contour on S defines a new surface S_c bounded by composite boundary Γ_c . Hence, with the singular point now removed we can argue, again by appealing to Corollary 5.2.1, that

$$\oint_{\Gamma} \mathbf{f} \cdot d\mathbf{r} = \oint_{\Gamma_c} \mathbf{f} \cdot d\mathbf{r} - \oint_{\Gamma_\epsilon} \mathbf{f} \cdot d\mathbf{r} = \iint_{S_c} \nabla \times \mathbf{f} \cdot \mathbf{N} dS - \oint_{\Gamma_\epsilon} \mathbf{f} \cdot d\mathbf{r}.$$

— by Stokes's theorem on Γ_c
at the middle step

We then need to consider the contribution from the curve integral around Γ_ϵ in the limit $\epsilon \rightarrow 0$ to regain the original surface S .

☞ Mastery Check 5.25:

Determine $\oint_{\Gamma} \mathbf{f} \cdot d\mathbf{r}$, where

$\mathbf{f} = (y, x, xy)$ and Γ is the curve of intersection of the plane $x + y + z = 1$ and

planes $x = y = z = 0$ in Figure 5.56.

Orientation of Γ : counterclockwise seen from $(10, 10, 10)$.

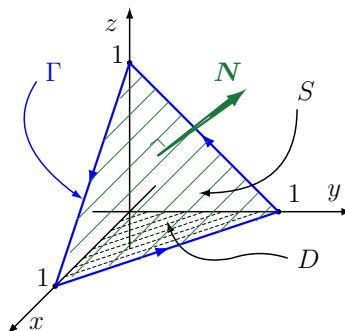


Figure 5.56 The intersection of $x + y + z = 1$ and the coordinate planes.

5.G Supplementary problems

Section 5.A

1. Given a C^2 curve in \mathbb{R}^3 representing the trajectory of a particle and described by the parametrization

$$\mathbf{r}(t) = x(t)\mathbf{i} + y(t)\mathbf{j} + z(t)\mathbf{k}, \quad 0 \leq t \leq T,$$

express the position, velocity and acceleration of the particle in terms of the local unit vectors \mathbf{T} , \mathbf{N} and \mathbf{B} , the curvature κ and torsion τ of the curve, and the speed v of the particle. Give an interpretation of your result for the particle's acceleration.

Conversely, express κ and τ in terms of \mathbf{v} and \mathbf{a} .

2. Find the curvature and torsion of the following smooth curves:

- (a) $\mathbf{r}(t) = (t^3, t^2, t)$, when $t = 1$,
- (b) $\mathbf{r}(t) = (e^t, 3t, e^{-2t})$, when $t = 1$,
- (c) $\mathbf{r}(t) = (t, t^2, \frac{2}{3}t^3)$, for general t .

3. Suppose a smooth curve is to be parameterized specifically in terms of arc length s rather than t . From the equations derived earlier, express \mathbf{T} , \mathbf{N} and \mathbf{B} in terms of s . Consequently, derive formulae for $\mathbf{T}'(s)$, $\mathbf{N}'(s)$ and $\mathbf{B}'(s)$ corresponding to those for $\mathbf{T}'(t)$, $\mathbf{N}'(t)$ and $\mathbf{B}'(t)$ appearing in Definitions 5.4 and 5.5, and Mastery Check 5.4.
- The formulae for $\mathbf{T}'(s)$, $\mathbf{N}'(s)$ and $\mathbf{B}'(s)$ are called *Frenet-Serret* formulae and are fundamental to the differential geometry of space curves.

4. Consider the 3D circle of radius a with centre at \mathbf{r}_0 described by

$$\mathbf{r}(\theta) = \mathbf{r}_0 + a(\sin \phi \cos \theta, \sin \phi \sin \theta, \cos \phi) \quad \text{with } \phi = f(\theta) \text{ and } \theta \in [0, 2\pi],$$

where θ and ϕ are spherical polar angle variables. Show that it has constant curvature, $\kappa(\theta) = 1/a$ and that the curve's torsion $\tau(\theta) = 0$.

Hint: It may help to consider a few simple cases of circles.

With this exercise we provide some detail to the claims made in Example 5.2.

5. The *catenary* curve in Figure 5.57(a) is the shape taken by a chain, of uniform density per unit length, which is allowed to hang under its

own weight.

Such a curve may be described by the vector function

$$\mathbf{r}(t) = at\mathbf{i} + a \cosh t \mathbf{j}, \quad -1 \leq t \leq 1,$$

where $a > 0$ is a constant.

Calculate the length of this curve.

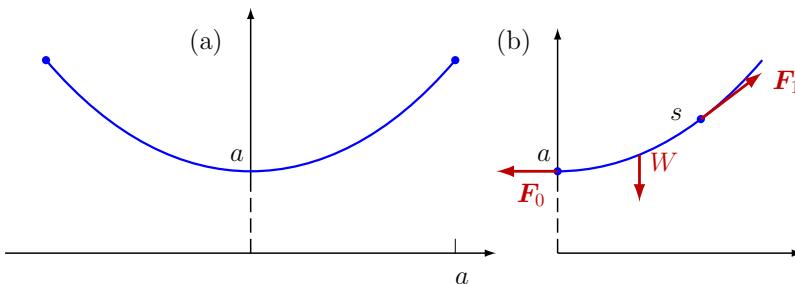


Figure 5.57 (a) The curve $\mathbf{r}(t)\mathbf{i} + a \cosh t \mathbf{j}$ in Problem 4;
 (b) The catenary curve in Problem 5.

6. Suppose a piece of chain is allowed to hang under its own weight, anchored at $(0, a)$, $a > 0$, by a force acting horizontally in the negative direction, and at some point (x, y) , $x > 0$, $y > a$, by a force tangential to the curve adopted by the chain.

Let the piece of chain have length s and let the curve be defined as $\mathbf{r}(t) = x(t)\mathbf{i} + y(t)\mathbf{j}$, $t \geq 0$. Assume the density of the chain is ρ per unit length, and that it is at rest while subject to tangential tension forces at each end.

These two tension forces \mathbf{F}_0 and \mathbf{F}_1 are shown in Figure 5.57(b).

Assume (from elementary mechanics) that the horizontal components of these forces are equal and opposite, while the sum of the vertical components is equal to the weight $\mathbf{W} = -\rho g s \mathbf{j}$ of that portion of the chain (acting through its centre of gravity).

Show that if the forces are in equilibrium then the curve adopted by the chain is the catenary curve as defined in the last problem, with a suitable choice of constants. Use the following steps:

- (i) Show that $\frac{dy}{dx} = \frac{s}{\lambda}$ for some constant λ .

- (ii) From $\frac{ds}{dt} = |\mathbf{r}'(t)|$ derive the result $\frac{ds}{dx} = \frac{\sqrt{\lambda^2 + s^2}}{\lambda}$, and a similar result for $\frac{ds}{dy}$.
- (iii) Determine x and y in terms of s .
- (iv) Show that these equations are solved by setting $x = \lambda t$, $y = \lambda \cosh t$. That is, the curve is as described with the appropriate choice of constants.

Section 5.B

7. For smooth scalar functions ψ and ϕ , verify that

$$\nabla \times (\phi \nabla \psi) = \nabla \phi \times \nabla \psi.$$

8. For the C^1 fields \mathbf{F} and \mathbf{G} , verify that

$$\nabla(\mathbf{F} \cdot \mathbf{G}) = (\mathbf{F} \cdot \nabla)\mathbf{G} + (\mathbf{G} \cdot \nabla)\mathbf{F} + \mathbf{F} \times (\nabla \times \mathbf{G}) + \mathbf{G} \times (\nabla \times \mathbf{F}).$$

9. For the C^1 field \mathbf{F} , verify that

$$\nabla \times (\mathbf{F} \times \mathbf{r}) = \mathbf{F} - (\nabla \cdot \mathbf{F})\mathbf{r} + \nabla(\mathbf{F} \cdot \mathbf{r}) - \mathbf{r} \times (\nabla \times \mathbf{F}).$$

10. For the C^2 field \mathbf{F} , verify that

$$\nabla \times (\nabla \times \mathbf{F}) = \nabla(\nabla \cdot \mathbf{F}) - \nabla^2 \mathbf{F}.$$

11. Show that in 3D a field proportional to $\frac{\mathbf{r}}{|\mathbf{r}|^3}$ is conservative.

Show also that this field is solenoidal.

12. Show that the 2D electrostatic field $\mathbf{E} = \frac{\rho}{2\pi\epsilon_0|\mathbf{r}|^2} \mathbf{r}$ (a field due to a uniformly charged wire of infinite length) is conservative.

13. Derive and solve the equations for the field lines corresponding to the 3D vector field $\frac{\mathbf{r}}{|\mathbf{r}|^3}$, $\mathbf{r} \neq \mathbf{0}$. Hence, show that these correspond to radial lines emanating from the origin.

14. Suppose the C^1 vector field \mathbf{F} satisfies $\nabla \cdot \mathbf{F} = 0$ in a domain $D \subset \mathbb{R}^3$. If

$$\mathbf{G}(x, y, z) = \int_0^1 t \mathbf{F}(\mathbf{r}(t)) \times \frac{d\mathbf{r}}{dt} dt$$

with $\mathbf{r}(t) = (xt, yt, zt)$ for $t \in [0, 1]$, show that $\nabla \times \mathbf{G} = \mathbf{F}$. This shows that \mathbf{G} is a vector potential for the solenoidal field \mathbf{F} for points $(x, y, z) \in D$.

Section 5.C

15. Evaluate

$$\int_{(0,2)}^{(4,0)} (x^2 + y^2) \, dx$$

along the path $y = \sqrt{4 - x}$.

16. Evaluate

$$\int_{(0,2)}^{(4,0)} (x^2 + y^2) \, d\mathbf{r}$$

along the path $y = \sqrt{4 - x}$.

17. Evaluate

$$\int_{(-1,0)}^{(1,0)} y(1 + x) \, dy$$

a) along the x -axis; b) along the parabola $y = 1 - x^2$.

18. Along what curve of the family $y = kx(1 - x)$ does the integral

$$\int_{(0,0)}^{(1,0)} y(x - y) \, dx$$

have the largest value?

19. Suppose a particle experiences a force directed to the origin with magnitude inversely proportional to its distance from the origin. Evaluate the work done by the force if the particle traverses the helical path $\mathbf{r}(t) = (a \cos t, a \sin t, bt)$ for $t \in [0, 2\pi]$.

20. Compute $\int_{\Gamma} \mathbf{A} \cdot d\mathbf{r}$, where $\mathbf{A} = -y\mathbf{i} + 2x\mathbf{j} + x\mathbf{k}$, and Γ is a circle in the plane $x = y$ with centre $(1, 1, 0)$ and radius 1.

Orientation: Γ is traversed counterclockwise as seen from $(1, 0, 0)$, as shown in Figure 5.58.

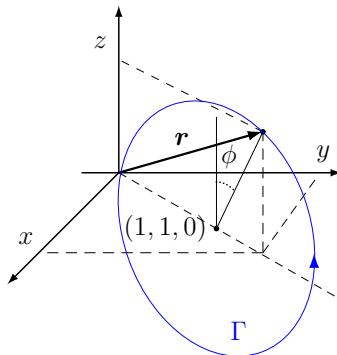


Figure 5.58 A circle intersecting the xy -plane.

21. Compute $\int_{\Gamma} \mathbf{A} \cdot d\mathbf{r}$, where

$$\mathbf{A} = x^2 \mathbf{i} + (x^2 + y^2) \mathbf{j} + (x^2 + y^2 + z^2) \mathbf{k},$$

and Γ is the boundary of $\{(x, y, z) : x + y + z = 1, x \geq 0, y \geq 0, z \geq 0\}$.

Orientation: Γ is traversed counterclockwise as seen from the point $(10, 10, 10)$, as shown in Figure 5.56.

Section 5.D

22. Evaluate $\iint_S (2x + y - 3z^2) dS$, where S is defined by
 $\mathbf{r}(u, v) = u \mathbf{i} + v \mathbf{j} + u \mathbf{k}, 0 \leq u, v \leq 1$.

23. Evaluate $\iint_S (6x + y - x^2) dS$, where S is defined by
 $\mathbf{r}(u, v) = u \mathbf{i} + u^2 \mathbf{j} + v \mathbf{k}, 0 \leq u, v \leq 1$.

24. Let S be the ellipsoid

$$\frac{x^2}{a^2} + \frac{y^2}{b^2} + \frac{z^2}{c^2} = 1,$$

and $p(x, y, z)$ be the length of the perpendicular from the plane, that is tangent to S at $(x, y, z) \in S$, to the origin. Show that the surface integral

$$\iint_S \frac{dS}{p} = \frac{4}{3} \pi abc \left(\frac{1}{a^2} + \frac{1}{b^2} + \frac{1}{c^2} \right).$$

Section 5.E

25. Use Gauss's theorem to determine the flux of the vector field $\mathbf{f} = (x^2 + y^2, y^2 + z^2, z^2 + x^2)$ through the surface of the cone $S = \{(x, y, z) : x^2 + y^2 - (z - 2)^2 = 0, 0 \leq z \leq 2, \mathbf{N} \cdot \mathbf{e}_3 > 0\}$ (Figure 5.59).

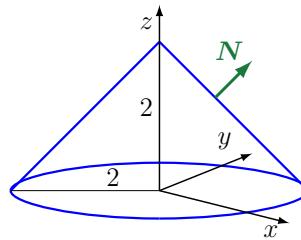


Figure 5.59 The cone $x^2 + y^2 = (z - 2)^2$.

26. Suppose the region $D \subset \mathbb{R}^3$ is bounded by a closed surface S . Using Gauss's theorem as a basis, prove the following variants of the theorem:

(a)

$$\iiint_D (\nabla \times \mathbf{F}) \, dV = \iint_S (\mathbf{N} \times \mathbf{F}) \, dS$$

(b)

$$\iiint_D (\nabla \phi) \, dV = \iint_S (\phi \mathbf{N}) \, dS$$

27. Let A be the area of a region D of the surface of a sphere of radius R centred at the origin. Let V be the volume of the solid cone comprising all rays from the origin to points on D . Show that

$$V = \frac{1}{3} A R.$$

28. Show that the electric intensity due to a uniformly charged sphere at points outside the sphere is the same as if the charge were concentrated at the centre, while at points inside the sphere it is proportional to the distance from the centre.

Section 5.F

29. Verify Green's theorem for the integral

$$\oint_{\Gamma} (x^2 + y) \, dx - xy^2 \, dy$$

where Γ is the boundary of the unit square with vertices (in order) $(0, 0)$, $(1, 0)$, $(1, 1)$, $(0, 1)$.

30. Verify Green's theorem for the integral

$$\oint_{\Gamma} (x - y) \, dx + (x + y) \, dy$$

where Γ , shown in Figure 5.60, is the boundary of the area in the first quadrant between the curves $y = x^2$ and $y = \sqrt{x}$ taken anticlockwise.

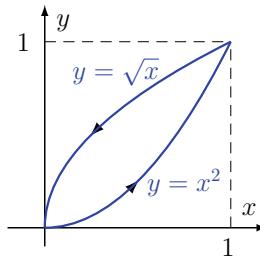


Figure 5.60 The closed contour Γ .

31. Verify Green's theorem for the integral

$$\oint_{\Gamma} (x - 2y) \, dx + x \, dy$$

where Γ is the boundary of the unit circle $x^2 + y^2 = 1$ taken anticlockwise.

32. Use Green's theorem to evaluate

$$\oint_{\Gamma} (2xy + y^2) \, dx + (x^2 + xy)x \, dy$$

where Γ is the boundary taken anticlockwise of the region cutoff from the first quadrant by the curve $y^2 = 1 - x^3$.

33. The intersection of the ellipsoid $x^2/2 + y^2 + (z - 1)^2/4 = 1$ and the plane $z + 2x = 2$ is a closed curve. Let Γ be that part of the curve

lying above the xy -plane directed from, and including, $(1, -\frac{1}{2}, 0)$ to, and including, $(1, \frac{1}{2}, 0)$. Evaluate $\oint_{\Gamma} \mathbf{F} \cdot d\mathbf{r}$ if

$$\mathbf{F}(x, y, z) = \left(z^2 + \frac{2}{3}xy^3 \right) \mathbf{i} + (x^2y^2) \mathbf{j} + \left(\frac{1}{4}z^2y^4 + \frac{1}{2}y^2 \right) \mathbf{k}.$$

34. Let Γ be the curve of intersection of the surfaces

$$\begin{aligned} S_1 &= \{(x, y, z) : 4x^2 + 4y^2 + z^2 = 40\}, \\ S_2 &= \{(x, y, z) : x^2 + y^2 - z^2 = 0, z > 0\}. \end{aligned}$$

Evaluate $\oint_{\Gamma} \mathbf{F} \cdot d\mathbf{r}$ if

$$\mathbf{F}(x, y, z) = \left(\frac{y}{x^2 + y^2} - z \right) \mathbf{i} + \left(x - \frac{x}{x^2 + y^2} \right) \mathbf{j} + (y + x^2) \mathbf{k}.$$

35. Suppose Γ is the curve of intersection of surfaces $x^2 + y^2 = x$ and $1 - x^2 - y^2 = z$, while vector field $\mathbf{f} = (y, 1, x)$. Evaluate $\oint_{\Gamma} \mathbf{f} \cdot d\mathbf{r}$.

Glossary of symbols

$a, b, c, a_1, \dots, b_1, \dots$	Scalar constants (usually).
x, y, z, u, v, w, s, t	Variables, assumed continuous.
$\mathbf{u}, \mathbf{v}, \mathbf{\underline{u}}, \mathbf{\underline{y}}, \vec{u}, \vec{v}$	Vector variables or vector constants.
f, g, h, F, G, H	Functions (usually).
$F^{(k)}(x)$	The k^{th} -order derivative of $F(x)$.
m, n, ℓ	Discrete integer variables.
\mathbb{R}	Real one-dimensional space; the real line.
$\mathbb{R}^n, n = 2, 3, \dots$	Real n -dimensional space. \mathbb{R}^3 is the 3-dimensional space we inhabit.
V, S	A volume region, and a surface embedded, in \mathbb{R}^n , $n \geq 3$.
$D_f; D, R$	Domain of a function f ; a general region and a rectangular region, respectively, (of integration) in \mathbb{R}^n .
M	A subset of \mathbb{R}^n ; n is usually specified in context.
$M^c, \partial M$	The complement, and the boundary, of a set M .
\overline{M}	The closure of a set M : $\overline{M} = M \cup \partial M$.
\mathbf{N}, \mathbf{n}	Unit normal vector and non-unit normal vector to a surface $S \subset \mathbb{R}^3$.
$ x , \mathbf{u} $	Absolute value of scalar x , and the magnitude of vector \mathbf{u} .
$ A $	Determinant of square matrix A .

$\mathbf{i}, \mathbf{j}, \mathbf{k}$	Unit vectors in the x -, y -, & z -directions, respectively.
$\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3$	Another way of writing $\mathbf{i}, \mathbf{j}, \mathbf{k}$, respectively.
θ, ϕ	Angles (usually). The symbol ϕ is also used to denote a scalar potential field in \mathbb{R}^n $n \geq 2$.
α, β, γ	Components of a fixed vector (usually). Angles (sometimes).
\equiv, \exists, \forall	“equivalent to”, “there exists”, and “for all”.
\in	“belongs to” or “is a member of”.
\subset	“is a subset of” or “is contained within”.
\cap	$A \cap B$ is the intersection of two sets A, B .
\cup	$A \cup B$ is the union of two sets A, B .
(\cdot)	An undefined, undeclared, or generic, argument of a function.
s.t.	“such that”.
:	(Inside a set-builder like $\{x : x < 0\}$) “such that”.
w.r.t.	“with respect to”.
$\frac{df}{dx} \equiv f'(x)$	Notation for the total derivative of a function f of x .
$(f \circ g)(x)$	Composite function; equivalent to $f(g(x))$.
$\frac{\partial f}{\partial x} \equiv f_x$	The partial derivative w.r.t. x of a function f of two or more variables.
$C^n(\mathbb{R}^m)$	The set of continuous functions defined on the space \mathbb{R}^m ($m \geq 1$) having continuous derivatives of order up to & including n .
∇	The vector differential operator, “del”.
∇f	The gradient of a scalar function f , “grad f ”.
$\nabla \cdot \mathbf{f}$	The divergence of a vector function \mathbf{f} , “div \mathbf{f} ”.
$\nabla \times \mathbf{f}$	The rotation vector of a vector field f , “curl \mathbf{f} ”.
J	The Jacobian determinant (usually).
L	A linear operator, a level set, a line, or a length, depending on context.

\implies, \Rightarrow	Implication: “this implies” or “this results in”.
\iff	“if and only if” or “equivalent to”.
\perp	“is orthogonal to”, “is perpendicular to”, or “is at right-angles to”.
\parallel	“is parallel to”.
\ll	“is much less than”.
\gtrless	“is greater than and also less than”.
\mapsto, \hookrightarrow	$f : x \mapsto y$ (or $f : x \hookrightarrow y$) “function f maps point x to point y ” (point mapping).
$\longrightarrow, \rightarrow$	Context dependent: $f : A \longrightarrow B$ (or $f : A \rightarrow B$) “function f maps from set A into set B ” (set mapping); $x \longrightarrow 0$ (or $x \rightarrow 0$) “ x converges to 0”; “tends to”.
Field	A scalar or vector function on \mathbb{R}^n , $n \geq 2$.
Theorem	A proposition that can be proved to be true.
Corollary	A result that follows immediately from a theorem.
1D, 2D, 3D	“one dimension” or “one-dimensional”, etc.
ODE	“ordinary differential equation”.
PDE	“partial differential equation”.
b.c.	“boundary condition”.
b.v.p.	“boundary-value problem”.

Bibliography

Any serious scientist, engineer or mathematician should be in possession of a decent personal library of reference books. Likewise, a university or college library worthy of its name should have a ready and sufficient store of text books. Although hard-copy books are becoming unfashionable (perhaps, at least, to be replaced by electronic literature — “eBooks”), it is important for the student to follow up on the material presented here by reading some of the more specialist books. The list given below is far from exhaustive, but these books do cover the areas we have discussed as well as being some of my favourites.

On multivariable and vector calculus:

1. Adams, R.A., Calculus, A.: Complete Course, 5th edn. Addison-Wesley, Boston (2003)
2. Apostol, T.M.: Mathematical Analysis: A Modern Approach to Advanced Calculus. Addison-Wesley, Boston (1957)
3. Courant, R., Hilbert, D.: Methods of Mathematical Physics, vol. 1 & 2. Wiley-Interscience, Hoboken (1962)
4. Hardy, G.H.: A Course of Pure Mathematics, 9th edn. Cambridge University Press, Cambridge (1949)
5. Kaplan, W.: Advanced Calculus, 5th edn. Addison-Wesley, Boston (2003)
6. Grossman, S.I.: Calculus, 3rd edn. Academic, New York (1984)
7. Spiegel, M.R.: Schaum’s Outline of Advanced Calculus. McGraw-Hill, New York (1974)
8. Spiegel, M.R.: Schaum’s Outline of Vector Analysis. McGraw-Hill, New York (1974)

On the approximation of functions:

9. Edwards, R.E.: Fourier Series: A Modern Introduction, vol. 1 & 2. Holt, Rinehart & Winston, New York (1967)
10. Rice, J.R.: The Approximation of Functions, vol. 1 & 2. Addison-Wesley, Boston (1964)

On partial differential equations:

11. Epstein, B.: Partial Differential Equations: An Introduction. McGraw-Hill, New York (1962)
12. Greenspan, D.: Introduction to Partial Differential Equations. McGraw-Hill, New York (1961)
13. Jeffreys, H., Jeffreys, B.: Methods of Mathematical Physics, 3rd edn. Cambridge University Press, Cambridge (1966)
14. Kreyszig, E.: Advanced: Engineering Mathematics, 7th edn. Wiley, New York (1982)
15. Morse, P.M., Feschbach, H.: Methods of Theoretical Physics, vol. 1 & 2. McGraw-Hill, New York (1953)

On linear algebra:

16. Lipschutz, S., Lipson, M.L.: Schaum's Outline of Linear Algebra, 5th edn. McGraw-Hill, New York (2013)

A little history:

17. Bell, E.T.: Men of Mathematics, vol. 1 & 2. Pelican, Kingston (1965)
18. Cannell, D.M.: George Green: Miller and Mathematician 1793–1841. City of Nottingham Arts Department (1988)
19. Dunningham, G.: Waldo Carl Friedrich Gauss: Titan of Science. Hafner (1955)
20. Newman, J.R.: The World of Mathematics, vol. 1–4. George Allen & Unwin, Crows Nest (1961)

A little grammar:

21. Fowler, H.W.: A Dictionary of Modern English Usage, 2nd edn (Revised by Sir Ernest Gowers). Oxford University Press, Oxford (1968)
22. Taggart, C., Wines, J.A.: My Grammar and I (or should that be 'me'?). Michael O'Marra Books Ltd (2011)

Index

A

- Absolute maximum, 133
- Absolute minimum, 133
- Analytic functions, 115
- Antisymmetric functions, 215
- Arc length, 231

B

- Ball-and-stick model, 86, 88, 91, 93
- Bijective transformations, 199, 201, 211, 214

C

- Cauchy-Schwarz inequality, 11
- Chain rule, 63, 84, 89, 91, 97, 237
- Class (of differentiable functions), 83
- Composite function, 84, 88, 91
- Conic sections, 34–35
- Continuous function, 21–24, 52
- Critical point, 126, 130
- Cross-derivative, 81
- Curl of a vector, 242, 287
- Curl of a vector field, 241
- Curvature, 230
- Curve integral, 246
- Cylindrical coordinates, 26

D

- Derivative, 50, 52
- Determinant, 8, 108, 109
- Differentiable function, 68, 70–74

Differential

- Differential, 236, 250
- Differential of a function, 146
- Directional derivative, 75
- Distance function, 150
- Divergence of a vector field, 240
- Domain of a function, 18

E

- Element of surface area, 262, 264, 267
- Ellipsoid, 37
- Equation of a line, 6
- Equation of a plane, 4
- Error analysis, 145
- Extension of a function, 187
- Extreme point, 125

F

- Field, vector, *see* Vector field
- Floodgate model, 84, 87, 89, 91, 95
- Flux through a surface, 266–273
- Function mapping, 18

G

- Gauss's theorem, 273
- Gradient, 235, 236
- Gradient function, 77–79, 143
- Green's theorem, 281, 284

H

- Hessian matrix, 129

I

- Implicit functions, 101–111
- Improper integral, 193, 194–197
- Integral, 180
- Integration
 - Change of variables, 198–203
 - Complex domains, 187–192
 - Double integral, 181, 184, 199
 - Iterated integral, 184–187, 207–210, 214
 - Slicing, 184, 207
 - Triple integrals, 205–212
 - Volume element, 212

J

- Jacobian, 110, 111, 112, 201, 212, 214
- Jacobian matrix, 235

L

- Lagrange multiplier, 143, 145
- Lagrangian function, 143
- Laplacian operator, 159, 162
- Least squares
 - Fitting a curve, 149
 - Fitting a straight line, 147
 - Function approximation, 150
- Leibniz's rule, 99, 100
- Level curve, 78, 102, 143
- Level sets, 31, 36, 38–42, 101, 141
- Level surface, 78, 105
- Limits, 53
 - Limit laws, 54
 - Limits in 2D, 55–62
 - Standard limits, 24
- Line integral, 246, 253
- Line of intersection, 108
- Linear approximation, 68, 102
- Local approximation, 128
- Local maximum point, 125, 130
- Local minimum point, 125, 130

M

- Maclaurin polynomial, 116
- Moving trihedral, 230

N

- Neighbourhood, 14, 53

O

- One-to-one transformations, 201
- Open sphere, 14
- Optimization, 135
 - Compact domains, 135
 - Free of constraints, 137
 - How to play, 138
 - Under constraints, 140
- Orientation of a surface, 283, 287
- Osculating plane, 230

P

- Partial derivative, 62, 64, 80, 89
- PDEs, 152
 - Boundary conditions, 154
 - Boundary-value problem, 160, 162, 163
 - Diffusion equation, 158, 160
 - Dirichlet problem, 157
 - Discriminant, 153
 - Harmonic function, 157
 - Heat equation, 158, 160
 - Initial conditions, 154
 - Laplace equation, 156
 - Laplacian operator, 157
 - Poisson equation, 157
 - Separation of variables, 162, 164
 - Wave equation, 161, 162
 - Eigenfunctions, 166
 - Eigenvalues, 166
 - Homogeneous equation, 161
 - The ideal string, 162
- Point in space, 9
- Polar coordinates, 25
- Principle of superposition, 167

Q

- Quadratic form, 129

R

- Range of a function, 18
- Regular points, 228
- Riemann integral, *see* Integral

Riemann sum, 178, 180, 205

Right-hand rule, 5

S

Saddle point, 127, 130, 132, 171

Scalar, 2

Scalar functions, 17

Sets

 Boundary point, 133

 Connected set, 259

 Open, closed, bounded, compact,
 12–16

Singularity, 52, 62, 68, 133

Smooth curve, 226

Spherical coordinates, 26

Squeeze theorems, 24, 54

Stationary point, 130

Stokes's theorem, 287

Surface integral, 232, 260

Symmetric functions, 215

T

Tangent line approximation, 111

Tangent plane, 66, 67, 71, 106, 234

Tangent vector, 236

Taylor series, 113, 115, 128

Torsion, 230

Triangle inequality, 11

Triple integral, 205

V

Vector field, 238

 Circulation, 281

 Conservative field, 244, 245, 258

 Derivative matrix, 236

 Electrostatic field, 239

 Equipotential surfaces, 245

 Gradient field, 77, 239

 Gravitational field, 238

 Irrotational field, 288

 Scalar potential, 244, 245

 Vector identities, 243

Vector-valued functions, 223, 235

Vectors

 Angle between vectors, 4

 Binormal vector, 229

 Orthogonal vectors, 4

 Principal normal vector, 229

 Rules for differentiating, 226

 Scalar product, 4, 6, 67, 76, 77

 Scalar triple product, 7–8, 212

 Tangent vector, 229, 232

 Unit vector, 2

 Vector (cross) product, 3, 5, 227

 Vector field identities, 243

 Vector triple product, 9

Volume, 181, 185

W

Work done, 251