

VOL. 6 – MATHEMATICAL METHODS

# The Undergraduate Companion to Theoretical Physics

---

**Andrea Kouta Dagnino<sup>‡</sup>**

<sup>‡</sup>*Open University, Milton Keynes, UK.*

*E-mail:* [k.y.dagnino@gmail.com](mailto:k.y.dagnino@gmail.com)

---

# Contents

<b>I Analysis</b>	<b>4</b>	
1 Unit A1: Sets, functions and vectors	5	8.2 Continuity . . . . .
2 Unit A2: Number systems	6	8.3 Properties of continuous functions
3 Unit A3: Mathematical language and proofs (why haven't you finished this yet!!)	7	8.4 Trigonometric and exponential functions . . . . .
3.1 Mathematical statements . . . . .	7	
3.2 Direct Proof . . . . .	10	
4 Unit A4: Real functions, graphs and conics(why haven't you finished this yet!!)	15	9 Unit F1: Limits
5 Unit D1: Numbers	16	9.1 Introduction to limits of functions
5.1 The set $\mathbb{N}$ of Natural Numbers . . . . .	16	9.2 Asymptotic behaviour . . . . .
5.2 The Set $\mathbb{Q}$ of Rational Numbers . . . . .	19	9.3 Continuity of functions . . . . .
5.3 The Set $\mathbb{R}$ of Real Numbers . . . . .	20	9.4 Unusual function continuity . . . . .
5.4 Absolute Value . . . . .	22	9.5 Uniform continuity . . . . .
5.5 The Completeness Axiom . . . . .	23	10 Unit F2: Differentiation
6 Unit D2: Sequences	26	10.1 Continuity and differentiability . . . . .
6.1 Introduction to sequences . . . . .	26	10.2 Rules of differentiation . . . . .
6.2 Convergence of sequences . . . . .	27	10.3 Rolle's theorem and local extrema . . . . .
6.3 Formal Proofs of Limit Theorems	28	10.4 Mean value theorem . . . . .
6.4 Null sequences . . . . .	31	10.5 L'Hopital's rule . . . . .
6.5 Limit theorems for Convergent sequences . . . . .	34	11 Unit F3: Integration
6.6 Divergent sequences . . . . .	40	11.1 The Riemann integral . . . . .
7 D3 Series	44	11.2 Inequalities and series with integrals . . . . .
7.1 Introduction to Series . . . . .	44	11.3 Series . . . . .
7.2 Telescoping series . . . . .	45	12 Unit F4: Power series
7.3 Manipulating series . . . . .	46	12.1 Taylor series . . . . .
7.4 Non-negative series . . . . .	48	12.2 Convergence . . . . .
7.5 Series with positive and negative terms . . . . .	54	12.3 The combination rules . . . . .
7.6 Monotone convergence theorem .	57	
8 Unit D4: Functions and Continuity	59	<b>II Algebra and Group Theory</b>
8.1 Real functions . . . . .	59	13 Unit B1: Symmetry and groups
		13.1 Symmetry in $\mathbb{R}^2$ . . . . .
		13.2 Representing symmetries . . . . .
		13.3 Definition of a Group . . . . .
		13.4 Properties of groups and group elements . . . . .
		13.5 Symmetry in $\mathbb{R}^3$ . . . . .
		13.6 The Dihedral group . . . . .

14	<i>Unit B2: Subgroups and isomorphisms</i>	<b>136</b>	24	<i>Modules</i>	<b>237</b>
14.1	Subgroups . . . . .	136			
14.2	Cyclic groups and subgroups . .	141			
14.3	Cyclic groups and modular arithmetic . . . . .	144			
14.4	Isomorphisms . . . . .	146			
14.5	Standard groups . . . . .	149			
14.6	Direct product of groups . . . . .	150			
15	<i>Unit B3: Permutations</i>	<b>153</b>			
15.1	Permutations . . . . .	153			
15.2	Permutation groups . . . . .	156			
15.3	Even and Odd symmetries . . . .	158			
15.4	Conjugacy of $S_n$ . . . . .	162			
15.5	Subgroups of $S_4$ . . . . .	164			
15.6	Cayley's Theorem . . . . .	166			
16	<i>Unit B4: Lagrange's Theorem and small groups</i>	<b>169</b>			
16.1	Lagrange's Theorem . . . . .	169			
16.2	Groups of small order . . . . .	171			
17	<i>Unit E1: Cosets and normal subgroups</i>	<b>175</b>			
17.1	Matrix groups . . . . .	175			
17.2	Cosets . . . . .	179			
17.3	Right cosets . . . . .	182			
17.4	Normal subgroups . . . . .	185			
18	<i>Unit E2: Quotient groups and conjugacy</i>	<b>187</b>			
18.1	Quotient groups . . . . .	187			
18.2	Quotient group of infinite groups	190			
18.3	Conjugacy . . . . .	191			
18.4	Normal subgroups and conjugacy	194			
18.5	Conjugacy in $S(\mathcal{F})$ . . . . .	199			
19	<i>Unit E3: Homomorphisms</i>	<b>205</b>			
19.1	Image and kernels . . . . .	210			
19.2	First isomorphism theorem . . . . .	212			
20	<i>Unit E4: Group actions</i>	<b>217</b>			
20.1	What are group actions? . . . . .	217			
20.2	Orbits and stabilisers . . . . .	220			
20.3	The Orbit-Stabiliser theorem . . .	225			
20.4	The Counting theorem . . . . .	229			
21	<i>Sylow theorems</i>	<b>234</b>			
22	<i>Rings</i>	<b>235</b>			
23	<i>Polynomials</i>	<b>236</b>			
			<b>III</b>	<b>Representation Theory and Lie Algebra</b>	<b>238</b>
			<b>IV</b>	<b>Differential Equations</b>	<b>239</b>
			25	<i>Fundamentals</i>	<b>240</b>
			25.1	Definitions . . . . .	240
			25.2	Integral formulation . . . . .	241
			25.3	Picard iteration . . . . .	243
			25.4	Existence and uniqueness . . . . .	245
			26	<i>First Order ODEs</i>	<b>247</b>
			26.1	Types of first order ODEs . . . . .	247
			26.2	Separable Differential Equations	248
			26.3	Exact Differential Equations . . . .	248
			26.4	Inexact Differential Equations . . .	249
			26.5	Integrating Factor Method . . . . .	250
			26.6	Bernoulli Equations . . . . .	251
			26.7	Stability and Equilibrium points	251
			27	<i>Second Order ODEs</i>	<b>252</b>
			27.1	Homogeneous equation . . . . .	252
			27.2	Non-homogeneous . . . . .	254
			27.3	Undetermined Coefficients . . . . .	255
			27.4	Variation of Constants . . . . .	255
			27.5	Reduction of Order . . . . .	256
			27.6	Euler-Cauchy equations . . . . .	256
			27.7	Intro to Green's functions . . . . .	256
			28	<i>Mechanical Vibrations and Resonance Phenomena</i>	<b>258</b>
			28.1	Homogeneous Equation . . . . .	258
			28.2	Damped Harmonic Motion . . . . .	259
			28.3	Forced Oscillations . . . . .	261
			28.4	Resonance . . . . .	262
			29	<i>General Linear ODEs</i>	<b>263</b>
			29.1	Existence and Uniqueness . . . . .	263
			29.2	Fundamental set and Wronskians	263
			29.3	Homogeneous ODE . . . . .	265
			29.4	Variation of parameters . . . . .	265
			29.5	Higher Order linear ODEs . . . . .	266
			30	<i>Sturm-Liouville theory and Green's functions</i>	<b>268</b>
			30.1	Linear differential operators . . . . .	268
			30.2	Eigenfunctions . . . . .	270

30.3 Sturm-Liouville problems . . . . .	273	40 Hyperbolic PDEs: Waves	<b>316</b>
30.4 Green's functions for BVPs . . . . .	276	40.1 The wave equation . . . . .	316
30.5 Green's functions for IVPs . . . . .	279	40.2 d'Alembert's solution . . . . .	318
<b>31 Linear systems of ODEs</b>	<b>281</b>	40.3 The 1D wave equation: strings . .	318
31.1 Non-degenerate Eigenvalues . . . . .	281	40.4 The 2D wave equation: membranes	318
31.2 Matrix exponentiation . . . . .	282	40.5 Existence and uniqueness theorem	318
31.3 Higher Order Linear Constant Coefficient Equations . . . . .	283	<b>41 Parabolic PDEs: Heat and Diffusion</b>	<b>319</b>
31.4 Triangulation . . . . .	283	41.1 Existence and uniqueness theorem	319
31.5 Jordan Form . . . . .	284	<b>42 Green's functions for PDEs</b>	<b>320</b>
<b>32 Series solutions methods</b>	<b>286</b>		
32.1 Power Series . . . . .	286	<b>V Linear Algebra</b>	<b>321</b>
32.2 Series Solutions near ordinary points . . . . .	287	<b>43 Vector spaces</b>	<b>322</b>
32.3 Euler equations . . . . .	289	43.1 Definitions . . . . .	322
32.4 Frobenius' method . . . . .	290	43.2 Basis and dimensions . . . . .	324
<b>33 Special functions</b>	<b>294</b>	43.3 Operations on subspaces . . . . .	327
33.1 Laguerre polynomials . . . . .	294	<b>44 Euclidean geometry in <math>\mathbb{R}^3</math></b>	<b>333</b>
33.2 Legendre polynomials . . . . .	294	<b>45 Matrix algebra</b>	<b>334</b>
33.3 Spherical harmonics . . . . .	294	<b>46 Linear transformations</b>	<b>335</b>
33.4 Hermite polynomials . . . . .	294	46.1 What is a map? . . . . .	335
33.5 Chebyshev polynomials . . . . .	294	46.2 What is a linear map? . . . . .	338
33.6 Bessel functions . . . . .	294	46.3 Isomorphisms . . . . .	342
<b>34 Distributions</b>	<b>295</b>	46.4 Linear maps and matrices . . . . .	342
34.1 Introducing the Dirac delta . . . . .	295	46.5 Change of basis and equivalence	349
34.2 Rigorous treatment-distributions	297	<b>47 Solving linear equations</b>	<b>354</b>
<b>35 Laplace transform methods</b>	<b>302</b>	47.1 Structure of solutions . . . . .	354
35.1 Basic definition and properties of the Laplace transform . . . . .	302	47.2 Elementary matrix operations . .	355
35.2 Solving ODEs with Laplace transforms . . . . .	305	47.3 Inverting matrices . . . . .	358
35.3 Convolutions . . . . .	306	<b>48 Determinants</b>	<b>363</b>
<b>36 Phase plane analysis</b>	<b>308</b>	48.1 The determinant of a matrix . .	363
<b>37 First order PDEs</b>	<b>309</b>	48.2 Laplace expansion . . . . .	366
37.1 Introduction . . . . .	309	48.3 Cramer's rule . . . . .	370
37.2 From ODEs to PDEs . . . . .	309	<b>49 Inner product spaces</b>	<b>372</b>
37.3 Change of variable and the chain rule . . . . .	311	49.1 Inner products . . . . .	372
37.4 The method of characteristics . .	311	49.2 Projectors . . . . .	373
<b>38 Second order PDEs: an overview</b>	<b>314</b>	49.3 Inner products and matrices . .	374
<b>39 Elliptic PDEs: Electrostatics</b>	<b>315</b>	49.4 Bilinear and Sesquilinear forms .	375
39.1 Existence and uniqueness theorem	315	<b>50 Eigen-everything</b>	<b>378</b>
		50.1 Finding eigenvalues and eigenvectors . . . . .	378
		50.2 Matrix diagonalization . . . . .	380

50.3 Orthogonal diagonalization . . . . .	<b>383</b>	56.3 de Rham cohomology and Electromagnetism . . . . .	<b>450</b>
50.4 Classifying conics . . . . .	<b>387</b>	57 <i>Connections and parallel transport</i>	<b>452</b>
50.5 Matrix exponentials and Lie algebras . . . . .	<b>390</b>	57.1 Covariant derivatives . . . . .	<b>452</b>
50.6 Schur's triangulation theorem . .	<b>391</b>	57.2 Parallel transport . . . . .	<b>454</b>
50.7 Jordan canonical form (make sure to write by October) . . . . .	<b>394</b>	57.3 Riemannian curvature . . . . .	<b>457</b>
<b>51 Multilinear algebra</b>	<b>395</b>	<b>58 Metrics and Curvature</b>	<b>460</b>
51.1 Review of linear algebra . . . . .	<b>395</b>	58.1 The metric tensor . . . . .	<b>460</b>
51.2 Multilinear maps . . . . .	<b>397</b>	58.2 Lie derivatives and symmetry .	<b>463</b>
<b>52 Tensor algebra</b>	<b>399</b>	<b>59 Integration on manifolds</b>	<b>468</b>
52.1 Einstein summation convention .	<b>399</b>		
52.2 Cartesian tensors . . . . .	<b>400</b>		
52.3 The $\delta_{ij}$ and $\epsilon_{ijk}$ tensors . . . . .	<b>406</b>		
52.4 Physical examples of cartesian tensors . . . . .	<b>406</b>		
52.5 Non-cartesian tensors . . . . .	<b>407</b>		
52.6 Covariance and contravariance .	<b>409</b>		
52.7 Application to special relativity: four-vectors . . . . .	<b>412</b>		
<b>53 Tensor calculus</b>	<b>414</b>	<b>VII Complex analysis</b>	<b>471</b>
53.1 Christoffel symbols . . . . .	<b>414</b>	<b>60 Complex numbers</b>	<b>472</b>
53.2 Differentiating tensors . . . . .	<b>415</b>	60.1 What are complex numbers . . .	<b>472</b>
53.3 Application to geometry: curvilinear coordinates . . . . .	<b>417</b>	60.2 Complex functions . . . . .	<b>477</b>
53.4 Geodesics . . . . .	<b>418</b>	60.3 Mappings under complex functions	<b>480</b>
		60.4 Special complex functions . . .	<b>480</b>
<b>VI Differential Geometry</b>	<b>419</b>	<b>61 Continuity of complex functions</b>	<b>485</b>
<b>54 Topology</b>	<b>420</b>	61.1 Complex sequences . . . . .	<b>485</b>
54.1 Why differential geometry? . . . . .	<b>420</b>	61.2 Continuity of complex functions	<b>490</b>
54.2 Topology . . . . .	<b>420</b>	61.3 Limits of complex functions . .	<b>492</b>
54.3 Topological manifolds . . . . .	<b>425</b>	61.4 Topology on $\mathbb{C}$ . . . . .	<b>494</b>
54.4 Multilinear algebra . . . . .	<b>427</b>	61.5 Extreme value theorem . . . . .	<b>497</b>
<b>55 Differentiable manifolds</b>	<b>431</b>	<b>62 Differentiating complex functions</b>	<b>500</b>
55.1 Differenti . . . . .	<b>431</b>	62.1 Derivatives of complex functions	<b>500</b>
55.2 Vectors and 1-forms . . . . .	<b>436</b>	62.2 Cauchy-Riemann equations . .	<b>502</b>
55.3 Interlude: Embeddings and immersions . . . . .	<b>439</b>	62.3 Derivative rules . . . . .	<b>504</b>
55.4 The tangent bundle . . . . .	<b>440</b>	62.4 Smooth paths . . . . .	<b>505</b>
55.5 Tensor fields . . . . .	<b>441</b>	<b>63 Integrating complex functions</b>	<b>508</b>
<b>56 Differentiable forms</b>	<b>444</b>	<b>64 Taylor and Laurent series</b>	<b>509</b>
56.1 Differentiable forms . . . . .	<b>444</b>	<b>65 Residues</b>	<b>510</b>
56.2 The exterior derivative . . . . .	<b>448</b>	<b>66 Zeros and extrema</b>	<b>511</b>
		<b>67 Conformal mappings</b>	<b>512</b>
		<b>68 Applications to fluid flows</b>	<b>513</b>
		<b>69 The Mandelbrot set and complex dynamics</b>	<b>514</b>

<b>VIII Calculus of Variations</b>	<b>515</b>	72 <i>Convolutions</i>	<b>519</b>
<b>IX Fourier Analysis</b>	<b>516</b>	<b>X Functional Analysis and Operator theory</b>	<b>520</b>
70 <i>Fourier series</i>	<b>517</b>	<i>Bibliography</i>	<b>522</b>
70.1 Dirichlet conditions . . . . .	<b>517</b>		
71 <i>Fourier transforms</i>	<b>518</b>		

---

*God used beautiful mathematics in creating  
the world.*

— P.A.M. Dirac



---

# References

Several textbooks, online courses/resources were referenced heavily (to the extend of making this text completely unoriginal, yet helpful for revision) throughout the writing of these lecture notes. We list the most relevant below:

- The Open University textbooks
- Lang, "*Undergraduate analysis*"
- Ross, "*Elementary analysis*"
- Armstrong, "*Groups and symmetries*"
- Stewart and Tall, "*Complex analysis*"
- Riley,Hobson and Bence, "*Mathematical methods for physics and engineering*"
- Lang, "*Linear algebra*"
- Kammler, "*A first course in Fourier analysis*"
- Fulton and Harris, "*Representation theory, a first course*"
- Jevanjee, "*An introduction to Tensors and Group Theory for Physicists*"
- Collins, "*Differential and Integral equations*"

# **Part I**

# **Analysis**

---

# Unit A1: Sets, functions and vectors

- (i)  $x + y = y + x, \forall x, y \in \mathbb{R}$
- (ii)  $x \cdot y = y \cdot x, \forall x, y \in \mathbb{R}$

---

## **Unit A2: Number systems**

2

# Unit A3: Mathematical language and proofs (why haven't you finished this yet!!)

## 3.1 Mathematical statements

### Types of statements

Assertions are a fundamental concept, a sort of axiom that we will not define rigorously, but can be seen as any form of some statement. For example,  $\{1, 2\}$  is greater than 2 is an assertion.

A **proposition**  $P$  is a statement that can be either true or false.  $\{1, 2\}$  is greater than 2 is not a proposition, since it is neither true nor false, it just doesn't make sense. A statement  $P(x)$  which is either true or false depending on the value of some variable(s) then it is a **variable proposition**.

A **theorem** is a mathematical statement that is (in general not obviously) true. Less important theorems are also called **propositions** (note that propositions can have two meanings depending on the context they're used in). A **lemma** is a small theorem, a result that can be used to prove other theorems. A **corollary** is a theorem that follows from another theorem.

### Logical connectives

All statements  $P$  have a related statement called **negation**  $\neg P$ , whose truth table is: For example

$P$	$\neg P$
T	F
F	T

for the following variable proposition:

$$P(x) : x \leq 0 \quad (3.1.1)$$

the negation is:

$$\neg P(x) : x > 0 \quad (3.1.2)$$

Similarly, for the proposition

$$Q : \text{there are at least 10 two-digit numbers less than 20} \quad (3.1.3)$$

$$\neg Q : \text{there are at most 9 two-digit numbers less than 20} \quad (3.1.4)$$

Given two statements  $P, Q$  we can insert the word *and*, *or* in between to give a new statement, the **conjunction**  $P \wedge Q$  ( $P$  and  $Q$ ) and **disjunction**  $P \vee Q$  ( $P$  or  $Q$ ) respectively. Their truth tables are: So  $\wedge$  is false whenever at least one of the two statements is false, whereas  $\vee$  is false only when both

$P$	$Q$	$P \wedge Q$	$P \vee Q$
T	T	T	T
T	F	F	T
F	T	F	T
F	F	F	F

statements are false. For example given  $P : 2$  is prime and  $Q : 2$  is even then:

$$P \wedge Q : 2 \text{ is prime and } 2 \text{ is even} \quad (3.1.5)$$

$$P \vee Q : 2 \text{ is prime or } 2 \text{ is even} \quad (3.1.6)$$

We can also negate conjunctions and disjunctions. For example, given the conjunction

$$p \text{ is odd and prime} \quad (3.1.7)$$

its negation is

$$p \text{ is even or not prime.} \quad (3.1.8)$$

Moreover, given the disjunction

$$\text{Either } A = B \text{ or } A \cup B = \emptyset \quad (3.1.9)$$

its negation is

$$\text{Both } A \neq B \text{ and } A \cup B \neq \emptyset. \quad (3.1.10)$$

It is therefore clear that when performing a negation *and*  $\leftrightarrow$  *or* and *both*  $\leftrightarrow$  *either*. Therefore:

$$\neg(P \wedge Q) = \neg P \vee \neg Q \quad \neg(P \vee Q) = \neg P \wedge \neg Q \quad (3.1.11)$$

## Implications

An **implication**  $P \implies Q$  (or  $P$  is sufficient for  $Q$ ,  $P$  only if  $Q$ ) is a mathematical statement of causality between two propositions  $P, Q$  in the form *if P then Q*. Its truth table is therefore:

$P$	$Q$	$P \implies Q$
T	T	T
T	F	F
F	T	T
F	F	T

so it is only false when the conclusion is false despite the hypothesis being true. For example:

$$x > 2 \text{ only if } x > 4 \quad (3.1.12)$$

can be expressed as  $x > 2 \implies x > 4$ . To negate an implication, we use a conjunction. Indeed, an implication states *if P then Q*, so its opposite would be to state that it is the case that  $P$  is true and  $Q$  is false, in other words  $P$  and not  $Q$ .

For example consider the implication:

$$(m \text{ divides } 12) \implies (m \text{ divides } 3 \text{ or } m \text{ divides } 4) \quad (3.1.13)$$

then its negation is:

$$(m \text{ divides } 12) \wedge (m \text{ does not divide neither } 3 \text{ nor } 4) \quad (3.1.14)$$

To every implication  $P \implies Q$  is an associated **converse** which states  $Q \implies P$ . It is important to remember that whether or not an implication is true tells you nothing about the truth of its converse. Indeed, consider the implication:

$$\text{If } m, n \text{ are both odd, then } m + n \text{ is even} \quad (3.1.15)$$

which is clearly true, then its converse:

$$\text{If } m + n \text{ is even, then } m, n \text{ are both odd} \quad (3.1.16)$$

which is not necessarily true.

To every implication  $P \implies Q$  is an associated **contrapositive**  $\neg Q \implies \neg P$  to which it is equivalent (therefore they state the same thing). For example the contrapositive of the previous implication 3.1.14 is:

$$\text{If } m + n \text{ is odd, then } m, n \text{ are not both odd} \quad (3.1.17)$$

## Equivalences

The statement  $P \implies Q \wedge Q \implies P$  is stated mathematically as:

$$P \iff Q \quad (3.1.18)$$

and is called an **equivalence**, and can be expressed as  $P$  if and only if (iff)  $Q$ . One common exception is in definitions, where simply *if* is used instead of *iff*. Equivalences can also be thought of as a conjunction between  $P \implies Q$  and  $\neg P \implies \neg Q$ . So for example:

$$m \text{ is even iff } m^2 \text{ is even} \quad (3.1.19)$$

is the same as stating

$$(m \text{ is even} \implies m^2 \text{ is even}) \wedge (m \text{ is odd} \implies m^2 \text{ is odd}) \quad (3.1.20)$$

## Universal and existential quantifiers

The **universal quantifier**  $\forall$  can be expressed as *for all/every*. So, the statement:

$$\forall x \in \mathbb{R}, P(x) \quad (3.1.21)$$

expresses the proposition  $P(x)$  for all  $x \in \mathbb{R}$ . Mathematical statements including the words *there are no, for all/every/each, every, any* are usually **universal statements** because they state the validity of a proposition for a certain set of values for a variable.

The **existential quantifier**  $\exists$  can be expressed as *there exists*. So the statement:

$$\exists x \in \mathbb{R}, P(x) \quad (3.1.22)$$

expresses the proposition  $P(x)$  is true for some  $x \in \mathbb{R}$ .

It is then clear that the negation of an existential statement is a universal statement and vice versa. Indeed the negation of  $\forall x \in \mathbb{R}, P(x)$  is:

$$\exists x \in \mathbb{R}, \neg P(x) \quad (3.1.23)$$

and the negation of  $\exists x \in \mathbb{R}, P(x)$  is:

$$\forall x \in \mathbb{R}, \neg P(x) \quad (3.1.24)$$

## 3.2 Direct Proof

### Proof by exhaustion

If  $n$  is an odd number between 0 and 10, then  $n^2$  is also odd.

*Proof.* The odd numbers between 0 and 10 are 1,3,5,7,9, and their respective squares are 1,9,25,49,81 which are all odd. ■

This is a **proof by exhaustion**, which consists in proving the statement for all the possible values of the variable in question, in this case  $n$ .

Another method of proof is simply algebraic verification.

[Geometric series identity] Let  $a, b \in \mathbb{R}$  and let  $n$  be a positive integer. Then:

$$a^n - b^n = (a - b)(a^{n-1} + a^{n-2}b + \dots + ab^{n-2} + b^{n-1}) \quad (3.2.1)$$

*Proof.* Expanding the RHS gives:

$$a^n + a(a^{n-2}b + a^{n-3}b^2 + \dots + ab^{n-2}) - b(a^{n-2}b + a^{n-3}b^2 + \dots + ab^{n-2}) - b^n \quad (3.2.2)$$

$$= a^n - b^n + b(a^{n-1} + a^{n-2}b + \dots + a^2b^{n-1}) - b(a^{n-2}b + a^{n-3}b^2 + \dots + ab^{n-2}) \quad (3.2.3)$$

$$= a^n - b^n \quad (3.2.4)$$

■

### Proving implications

To prove an implication  $P \implies Q$ , often we assume that  $P$  is true. We then find some true statement

$$P \implies P_1 \quad (3.2.5)$$

from which we deduce that  $P_1$  is true. We then find some other true statement:

$$P_1 \implies P_2 \quad (3.2.6)$$

from which we deduce that  $P_2$  is true. We continue this process, until we find that  $P_n$  is true and

$$P_n \implies Q \quad (3.2.7)$$

which proves that  $Q$  is true, whenever  $P$  is true.

Prove that if  $n$  is odd, then  $n^2$  is odd.

*Proof.* Let  $n$  be an odd integer, then:

$$\exists k \in \mathbb{Z} \text{ s.t. } n = 2k + 1 \implies n^2 = 4k^2 + 4k + 1 = 2(2k^2 + 2k) + 1 \quad (3.2.8)$$

which is in the form  $2m + 1$  with  $m = 2k^2 + 2k$ . Therefore  $n^2$  is odd whenever  $n$  is odd. ■

We now state an important theorem that will help us in proving other propositions.

[Fundamental Theorem of Arithmetic] Every integer greater than 1 can be written as a unique product of prime numbers.

$\forall n \in \mathbb{Z}, n^3 + 3n^2 + 2n$  is divisible by 6.

*Proof.* Let  $n \in \mathbb{Z}$  then:

$$n^3 + 3n^2 + 2n = n(n^2 + 3n + 2) = n(n+1)(n+2) \quad (3.2.9)$$

hence  $n^3 + 3n^2 + 2n$  is the product of three consecutive integers, which implies that one of them must be divisible by 2, and another must be divisible by 3. Therefore both 2 and 3 divide  $n^3 + 3n^2 + 2n$ . Hence by the Fundamental Theorem of Arithmetic we may write:

$$n^3 + 3n^2 + 2n = 2 \cdot 3 \cdot r = 6r \quad (3.2.10)$$

for some  $r$ . We may then deduce that 6 divides  $n^3 + 3n^2 + 2n$ . ■

Consider the following proof:

$$1 = -1 \implies 1^2 = (-1)^2 \implies 1 = 1 \quad (3.2.11)$$

therefore, since it is true that  $1 = 1$ , we can deduce that  $1 = -1$ .

This is clearly wrong, because we started with the statement we wanted to prove, and through a series of implications, we found a true statement. However, this does not achieve anything, since implications only work one way, we have proved that if  $1 = -1$  the  $1 = 1$ , but not the converse.

More generally, if  $P \implies Q$ , then we have no information about  $Q$  when  $P$  is false, as can be seen from the truth table.

## Proving equivalences

Recall that an equivalence  $P \iff Q$  is the same as saying  $(P \implies Q) \wedge (Q \implies P)$ . Thus, if we wish to prove  $P \iff Q$  then generally it suffices to prove both  $P \implies Q$  and  $Q \implies P$ . When dealing with equations or inequalities, often the individual implications  $P_i \implies P_{i+1}$  are

reversible, so they are actually equivalences. In that case we can simply prove  $P \iff Q$  using a series of equivalences:

$$P \iff P_1 \iff P_2 \iff \dots \iff P_n \iff Q \quad (3.2.12)$$

$$x(x - 2) = 3 \iff x = -1 \vee x = 3$$

*Proof.* We will use a series of equivalences:

$$x(x - 2) = 3 \iff x^2 - 2x - 3 = 0 \quad (3.2.13)$$

$$\iff (x + 1)(x - 3) = 0 \quad (3.2.14)$$

$$\iff (x + 1 = 0) \vee (x - 3 = 0) \quad (3.2.15)$$

$$\iff x = -1 \vee x = 3 \quad (3.2.16)$$

as required. ■

$$\text{For two arbitrary sets } A, B, A \cup B = A \iff B \subseteq A$$

*Proof.* We will use the approach of tackling the implication and its converse separately.

$\implies$  We firstly prove the rightward implication. Assume that  $A \cup B = A$ , and let  $x \in B$ , then it follows that  $x \in A \cup B$ , and thus  $x \in A$ . So,  $B \subseteq A$ , since  $x \in B \implies x \in A$ .

$\impliedby$  we prove the leftward implication. Assume that  $B \subseteq A$ .

Then let  $x \in A \cup B$  so that  $x \in A \vee x \in B$ . If  $x \in B$ , then  $x \in A$  by assumption. We can then deduce that  $x \in A$ , so  $A \cup B \subseteq A$ . Clearly we must also have  $A \subseteq A \cup B$ , so that  $A = A \cup B$  as required.

Finally, the equivalence has been proven:

$$A \cup B = A \iff B \subseteq A \quad (3.2.17)$$

■

Now consider the following (incorrect) proof:

$$n \text{ is a multiple of 5} \iff \exists k \in \mathbb{Z}, \text{s.t. } n = 5k \quad (3.2.18)$$

$$\iff \exists k \in \mathbb{Z}, \text{s.t. } n^2 = 25k^2 \quad (3.2.19)$$

$$\iff \exists k \in \mathbb{Z}, \text{s.t. } n^2 = 5(5k^2) \quad (3.2.20)$$

$$\iff \exists k \in \mathbb{Z}, \text{s.t. } n^2 \text{ is a multiple of 5.} \quad (3.2.21)$$

The problem with this proof lies in the final equivalence. Indeed, the implication is true, but its converse is not immediate and requires further justification. Indeed:

$$n^2 \text{ is a multiple of 5} \implies \exists k \in \mathbb{Z} \ n^2 = 5l \implies l = 5k^2 \implies n^2 = 5(5k^2) \quad (3.2.22)$$

for some  $k$ , since otherwise taking the square root of  $n^2 = 5l$  we would not be able to factor out the 5, and so  $n \notin \mathbb{Z}$ .

[Factor Theorem in  $\mathbb{R}$ ] Let  $p(x)$  be a real polynomial, and let  $\alpha \in \mathbb{R}$ . Then  $p(\alpha) = 0 \iff (x - \alpha)$  is a factor of  $p(x)$ .

*Proof.* Throughout the proof we assume that  $p(x)$  is real with  $\alpha \in \mathbb{R}$ .

$\implies$  Let us first prove the rightward implication. Assume that  $p(\alpha) = 0$ , and let:

$$p(x) = \sum_{i=0}^n a_i x^i, \quad a_n \neq 0 \quad (3.2.23)$$

Then since  $p(\alpha) = 0$ :

$$p(\alpha) = \sum_{i=0}^n a_i \alpha^i = 0 \implies p(x) = p(x) - p(\alpha) = \sum_{i=1}^n a_i (x^i - \alpha^i) \quad (3.2.24)$$

where the constant terms cancel out.

By the Geometric series identity, we know that  $x - \alpha$  divides the above expression. Hence we have that  $x - \alpha$  divides  $p(x)$  as required.

$\iff$  Now assume that  $(x - \alpha)$  is a factor of  $p(x)$  so that  $p(x) = (x - \alpha)q(x)$ . Then clearly:

$$p(\alpha) = 0 \cdot q(\alpha) = 0 \quad (3.2.25)$$

as required. ■

## Proving existential and universal statements

The simplest way to prove an existential statement is to provide an object that satisfies the statement.

There exists a real positive number  $x$  such that  $x \leq \sqrt{x}$ .

*Proof.* Let  $x = \frac{1}{9}$ , then  $x = \frac{1}{9} \leq \frac{1}{3} = \sqrt{x}$  as required. ■

We may also need to prove that a statement is false. For a universal statement, it suffices to find a **counterexample**. To prove that  $P(x) \implies Q(x)$  is false we need to prove that  $\exists x$  such that  $P(x)$  and not  $Q(x)$ , this value of  $x$  is a counterexample.

## Proof by induction

To prove that a mathematical statement  $P(n)$  is true for  $n = 1, 2, \dots$ :

- (i) prove that  $P(1)$  is true,
- (ii) prove that  $P(k) \implies P(k + 1)$  for  $k = 1, 2, \dots$

For all  $n \geq 7$ ,  $3^n < n!$ .

*Proof.* (i)  $P(7)$  is true since it is true that  $3^7 = 2187 < 5040 = 7!$ .

(ii) Assume that  $P(k)$  is true for some  $k \geq 7$  so that:

$$3^k < k! \quad (3.2.26)$$

Then we wish to deduce that  $P(k+1)$  is true as follows:

$$3^{k+1} = 3^k \cdot 3 < 3 \cdot k! < k \cdot k! = (k+1)! \quad (3.2.27)$$

so that  $3^{k+1} < (k+1)!$  as required. By mathematical induction, we conclude that  $P(n)$  is true for all  $n \geq 7$ . ■

For  $n \in \mathbb{N}$ ,  $2^{3n+1} + 5$  is a multiple of 7.

*Proof.*  $P(1)$  is true since  $2^{3 \cdot 1 + 1} + 5 = 21 = 3 * 7$ .

Assume that  $P(k)$  is true for some  $k \geq 1$  so that:

$$2^{3k+1} + 5 \text{ is a multiple of } 5, \quad (3.2.28)$$

Then we find that:

$$2^{3k+4} + 5 = 2^3 \cdot 2^{3k+1} + 5 = 7 \cdot 2^{3k+1} + 2^{3k+1} + 5 \quad (3.2.29)$$

which is divisible by 7 given that  $P(k)$  is true. Therefore we have shown that:

$$P(k) \implies P(k+1) \quad (3.2.30)$$

for  $k \in \mathbb{N}$ , and by mathematical induction it follows that the proposition is true. ■

Let us now try to prove a more advanced theorem, the factor theorem presented in the previous chapter.

Let  $p(x) = \sum_{i=0}^n a_i x^i$

---

# Unit A4: Real functions, graphs and conics(why haven't you finished this yet!!)

# Unit D1: Numbers

## 5.1 The set $\mathbb{N}$ of Natural Numbers

### Peano Axioms

We denote by  $\mathbb{N}$  the inductive set  $\{1, 2, 3, \dots\}$  of all positive integers, so that each positive integer  $n$  has a successor  $n + 1 = \text{succ}(n)$ . We can then state the following *Peano axioms*:

**N1.**  $1 \in \mathbb{N}$

**N2.**  $n \in \mathbb{N} \implies \text{succ}(n) \in \mathbb{N}$

**N3.**  $\forall n \in \mathbb{N}, 1 \neq \text{succ}(n)$

**N4.**  $\text{succ}(n) = \text{succ}(m) \iff n = m$

**N5.** Let  $A \subset \mathbb{N}$  which contains 1 and contains  $\text{succ}(n)$  whenever it contains  $n$ , then  $A = \mathbb{N}$

**Remark.** Assume N5 is false, then  $\mathbb{N}$  contains a set  $A$  such that:

(i)  $1 \in A$

(ii)  $n \in A \implies (n + 1) \in A$

(iii)  $A \neq \mathbb{N}$

and consider  $n_0 = \min S$ , where  $S = \{n \in \mathbb{N} \mid n \notin A\}$ . Clearly,  $n_0 \neq 1$ , so  $n_0$  is the successor of some number  $n_0 - 1$ . Since  $n_0 \in S$ , it follows that  $(n_0 - 1) \in A$ . However, by (ii)  $(n_0 - 1) \in A \implies n_0 \in A$  which is a contradiction.

### Principle of Mathematical Induction

Let  $P_1, P_2, \dots$  be a list of propositions, then the principle of mathematical induction asserts that they are true provided:

**I<sub>1</sub>**  $P_1$  is true (basis of induction)

**I<sub>2</sub>**  $(P_n \text{ is true}) \implies (P_{n+1} \text{ is true})$  (inductive step)

**Proposition 5.1 (Sum of natural numbers)**

The sum of the first  $n$  natural numbers is:

$$\sum_{i=1}^n = \frac{n(n+1)}{2}. \quad (5.1.1)$$

*Proof.* We define the  $n$ th proposition to be:

$$P_n : \sum_{i=1}^n = \frac{n(n+1)}{2} \quad (5.1.2)$$

I<sub>1</sub> For the basis for induction,  $P_1$  asserts that the sum of the first natural number is  $\frac{1 \cdot 2}{2} = 1$  which is clearly true.

I<sub>2</sub> For the inductive step, suppose  $P_n$  is true, so we assume:

$$\sum_{i=1}^n = \frac{n(n+1)}{2} \quad (5.1.3)$$

is true. Now:

$$\begin{aligned} \sum_{i=1}^{n+1} &= \frac{n(n+1)}{2} + (n+1) \\ &= \frac{n^2 + n + 2n + 2}{2} \\ &= \frac{(n+1)(n+2)}{2} \\ &= \frac{(n+1)((n+1)+1)}{2} \end{aligned}$$

so  $P_{n+1}$  is true as required. ■

**Example.** Let us prove that all numbers of the form  $5^n - 4n - 1$  are divisible by 16,  $\forall n \in \mathbb{N}$ .

*Proof.* So the  $n$ th proposition is:

$$P_n : 5^n - 4n - 1 \text{ are divisible by 16.}$$

I<sub>1</sub> The basis for induction is true, since  $5^1 - 4 - 1 = 0$  which is divisible by 16.

I<sub>2</sub> For the inductive step, suppose  $P_n$  is true, we wish to verify  $P_{n+1}$ . To do so, we write:

$$5^{n+1} - 4(n+1) - 1 = 5(5^n - 4n - 1) + 16n = 5 \cdot 16m + 16n = 16(5m + n) \quad (5.1.4)$$

where  $5^n - 4n - 1 = 16m$ , required. ■

**Example.** Let us prove that  $|\sin nx| \leq n|\sin x|, \forall n \in \mathbb{N}, \forall x \in \mathbb{R}$ .

*Proof.* Our  $n$ th proposition is:

$$P_n : |\sin nx| \leq n|\sin x|, \forall x \in \mathbb{R}$$

I<sub>1</sub> The basis for induction is clearly true, since  $|\sin x| \leq |\sin x|$ .

I<sub>2</sub> For the inductive step, assume that  $P_n$  is true, then:

$$\begin{aligned} |\sin(n+1)x| &= |\sin(nx+x)| \\ &= |\sin nx \cos x + \sin x \cos nx| \\ &\leq |\sin nx||\cos x| + |\sin x||\cos nx| \\ &\leq |\sin nx| + |\sin x| \\ &\leq n|\sin x| + |\sin x| \\ &\leq (n+1)|\sin x| \end{aligned}$$

as required,  $P_{n+1}$  holds. ■



### Theorem 5.2 (Bernoulli inequality)

Let  $x \in \mathbb{R}$  and  $n \in \mathbb{N}$  then:

$$(1+x)^n \geq 1+nx, \text{ when } x \geq -1 \quad (5.1.5)$$

*Proof.* Let  $x \geq -1$  and define:

$$P(n) : (1+x)^n \geq 1+nx \quad (5.1.6)$$

I<sub>1</sub>  $P(1)$  is obviously true, since the LHS reads  $(1+x)^1 = (1+x)$  which is equal to the RHS.

I<sub>2</sub> Now let  $P(k)$  be true for some  $k \geq 1$ , so that:

$$(1+x)^k \geq 1+kx \quad (5.1.7)$$

It follows that:

$$(1+x)^{k+1} \geq (1+kx)(1+x) \quad (5.1.8)$$

$$\geq 1+(k+1)x+kx^2 \quad (5.1.9)$$

$$\geq 1+(k+1)x \quad (5.1.10)$$

since  $kx^2 \geq 0$ . It follows that  $P(k+1)$  is true, whenever  $P(k)$  is verified.

Hence, by the principle of mathematical induction, we have that

$$(1+x)^n \geq 1+nx, \text{ when } x \geq -1 \quad (5.1.11)$$

as desired. ■

## 5.2 The Set $\mathbb{Q}$ of Rational Numbers

The set  $\mathbb{Q}$  of Rational Numbers is the set of numbers that can be written as the ratio of two integers in  $\mathbb{Z}$ .

### Definition 5.2 (Algebraic number)

A number  $x$  is called *algebraic* if it satisfies a polynomial equation:

$$c_n x^n + c_{n-1} x^{n-1} + \dots + c_1 x + c_0 = 0 \quad (5.2.1)$$

where  $c_i \in \mathbb{Z}$  and  $c_n \neq 0$ .

### Proposition 5.3 (Rational $\implies$ algebraic)

All rational numbers are algebraic numbers.

*Proof.* Consider the rational number  $x = \frac{m}{n} \in \mathbb{Q}$ , where  $m, n \in \mathbb{Z}$ . Then, clearly it satisfies the equation:

$$nx - m = 0. \quad (5.2.2)$$

■

### Theorem 5.4 (Rational Zeros Theorem)

Assume  $c_0 \dots c_n \in \mathbb{Z}$  and  $x \in \mathbb{Q}$  satisfying the equation:

$$c_n x^n + c_{n-1} x^{n-1} + \dots + c_1 x + c_0 = 0 \quad (5.2.3)$$

where  $c_n, c_0 \neq 0$ . Let  $x = \frac{c}{d}$  where  $c, d \in \mathbb{Z}$  with no common factors and  $d \neq 0$ . Then  $c$  divides  $c_0$  and  $d$  divides  $c_n$ .

*Proof.* We are given that:

$$\begin{aligned} c_n \left(\frac{c}{d}\right)^n + c_{n-1} \left(\frac{c}{d}\right)^{n-1} + \dots + c_1 \left(\frac{c}{d}\right) + c_0 &= 0 \\ c_n c^n + c_{n-1} c^{n-1} d + \dots + c_1 c d^{n-1} + c_0 d^n &= 0 \end{aligned}$$

Firstly, we solve for  $c_0 d^n$  to obtain:

$$c_0 d^n = -c [c_n c^{n-1} + c_{n-1} c^{n-2} d + \dots + c_2 c d^{n-2} + c_1 d^{n-1}] \quad (5.2.4)$$

so  $c$  divides  $c_0 d^n$ . However, since  $c$  and  $d^n$  are coprime, we must have that  $c$  divides  $c_0$ . Similarly, we solve for  $c_n c^n$ :

$$c_n c^n = -d [c_{n-1} c^{n-1} + c_{n-2} c^{n-2} d + \dots + c_1 c d^{n-2} + c_0 d^{n-1}] \quad (5.2.5)$$

so  $d$  divides  $c_n c^n$ . However, since  $c$  and  $d$  are coprime, we must have that  $d$  divides  $c_n$  as required.

■

**Corollary** Consider the equation:

$$x^n + c_{n-1}x^{n-1} + \dots + c_1x + c_0 = 0. \quad (5.2.6)$$

By applying the Rational Zeros Theorem, all rational solutions must divide  $c_0$ .

**Example.** Let us prove that  $a = \sqrt{2 + \sqrt[3]{5}}$  is irrational.

*Proof.* We firstly note that  $a$  is algebraic, since:

$$\begin{aligned} a^2 &= 2 + \sqrt[3]{5} \\ (a^2 - 2)^3 &= 5 \\ a^6 - 6a^4 + 12a^2 - 13 &= 0 \end{aligned}$$

which gives the polynomial equation:

$$x^6 - 6x^4 + 12x^2 - 13 = 0 \quad (5.2.7)$$

By corollary 1.1.1, the only possible rational solutions are  $\pm 1, \pm 13$  which clearly don't satisfy (1.5). ■

◀

### 5.3 The Set $\mathbb{R}$ of Real Numbers

The following algebraic properties hold for a field:

**A1.**  $a + (b + c) = (a + b) + c$  (addition associativity)

**A2.**  $a + b = b + a$  (addition commutativity)

**A3.**  $a + 0 = a$  (addition identity)

**A4.**  $\forall a, \exists (-a)$  s.t.  $a + (-a) = 0$  (addition inverse)

**M1.**  $a(bc) = (ab)c$  (multiplication associativity)

**M2.**  $ab = ba$  (multiplication commutativity)

**M3.**  $a \cdot 1 = a$  (multiplication identity)

**M4.**  $\forall a \neq 0, \exists a^{-1}$  s.t.  $a \cdot a^{-1} = 1$  (multiplication inverse)

**DL.**  $a(b + c) = ab + ac$  (distributivity)

The following ordering properties hold for an ordered field:

**O1.**  $\forall a, b, a \leq b \vee b \leq a$  (multiplication inverse)

**O2.**  $(a \leq b \wedge b \leq a) \implies a = b$

**O3.**  $(a \leq b \wedge b \leq c) \implies a \leq c$

**O4.**  $(a \leq b) \implies a + c \leq b + c$

**O5.**  $(a \leq b \wedge 0 \leq c) \implies ac \leq bc$

**Theorem 5.6 (Properties of fields)**

The following are consequences of the field properties:

- (i)  $a + c = b + c \implies a = b$
- (ii)  $a \cdot 0 = 0$
- (iii)  $(-a)b = -ab$
- (iv)  $(-a)(-b) = ab$
- (v)  $(ac = bc \wedge c \neq 0) \implies a = b$
- (vi)  $ab = 0 \implies (a = 0 \vee b = 0)$

*Proof.* (i)  $a + c = b + c \implies (a + c) + (-c) = (b + c) + (-c)$ , using A1 we have that  $a + [c + (-c)] = b + [c + (-c)] \implies a + 0 = b + 0$  by A4, so we finally have  $a = b$  using A3.

- (ii)  $a \cdot 0 = a \cdot (0 + 0) = a \cdot 0 + a \cdot 0$  where we used A3 and DL respectively. By (i) we conclude that  $a \cdot 0 = 0$ .
- (iii)  $a + (-a) = 0 \implies ab + (-a)b = [a + (-a)] \cdot b = 0 \cdot b = 0 = ab + (-ab)$ , so from (i) we have that  $(-a)b = -(ab)$ .
- (iv)  $(-a)(-b) + (-ab) = (-a)(-b) + (-a)b = (-a)[(-b) + b] = 0 = ab + (-ab)$ , so by (i) we have  $(-a)(-b) = ab$ .
- (v) Suppose  $ac = bc \wedge c \neq 0$ , then  $a = a \cdot 1 = a(cc^{-1}) = (ac)c^{-1} = b(cc^{-1}) = b$
- (vi) If  $ab = 0$  and  $b \neq 0$ , then  $0 = 0 \cdot b^{-1} = (ab)b^{-1} = a(bb^{-1}) = a \cdot 1 = a$

■

**Theorem 5.7 (Properties of ordered fields)**

The following are consequences of the properties of an ordered field:

- (i)  $a \leq b \implies -b \leq -a$
- (ii)  $(a \leq b \wedge c \leq 0) \implies bc \leq ac$
- (iii)  $0 \leq a \wedge 0 \leq b \implies 0 \leq ab$
- (iv)  $0 \leq a^2$
- (v)  $0 < 1$
- (vi)  $0 \leq a \implies 0 \leq a^{-1}$
- (vii)  $0 < a < b \implies 0 < b^{-1} < a^{-1}$
- (viii)  $0 \leq a, b \wedge p \in \mathbb{N} \implies (a < b \iff a^p < b^p)$

*Proof.* (i) Suppose  $a \leq b$ , then applying O4 with  $c = (-a) + (-b)$ , then  $a + [(-a) + (-b)] \leq b + [(-a) + (-b)] \implies -b \leq -a$

(ii) if  $a \leq b \wedge c \leq 0$ , then  $0 \leq -c$ . So, applying O5 gives  $a(-c) \leq -bc$ , and using (i) we get  $bc \leq ac$

(iii) This is a special case of O5 using  $a = 0$ .

(iv) For any  $a$ ,  $a \leq 0 \vee 0 \leq a$ . In the first case,  $a^2 \leq 0$  by (iii). In the latter case,  $-a \leq 0 \implies (-a)(-a) = a^2 \leq 0$  using (i).

- (v) Clearly,  $0 \neq 1$ . Indeed, consider  $x \neq 0$ , then  $x \cdot 1 = x \cdot 0 = 0$  which is a contradiction. Applying (iii) with  $a = 1$  gives the desired result.
- (vi) Suppose  $0 \leq a$  but  $a^{-1} \leq 0 \implies -a^{-1} \geq 0$ . Applying (iii) gives  $0 \leq a(-a^{-1}) = -1 \implies 1 \leq 0$  which contradicts (v).
- (vii) We multiply by  $(a^{-1})(b^{-1}) > 0$  and find  $0 < a < b \implies 0 < b^{-1} < a^{-1}$ .
- (viii) For positive integers  $p$ , we use the factor theorem:

$$b^p - a^p = (b - a) \underbrace{(b^{p-1} + b^{p-2}a + \dots + ba^{p-2} + a^{p-1})}_{>0} \quad (5.3.1)$$

but the term in brackets is positive definite, so it follows immediately that:

$$b - a > 0 \iff b^p - a^p > 0 \quad (5.3.2)$$

■

## 5.4 Absolute Value

### Definition 5.8 (Absolute value and distance)

We define the *absolute value* of  $a$  as:

$$|a| = \begin{cases} a & \text{if } a \geq 0 \\ -a & \text{if } a \leq 0 \end{cases} \quad (5.4.1)$$

For two numbers  $a, b$  we define  $\text{dist}(a, b) = |a - b|$  to be the *distance between  $a$  and  $b$* .

We present some important properties of the absolute value:

### Theorem 5.9 (Absolute value properties)

The following hold  $\forall a, b \in \mathbb{R}$ :

- (i)  $|a| \geq 0$
- (ii)  $|ab| = |a| \cdot |b|$
- (iii)  $|a + b| \leq |a| + |b|$
- (iv)  $|a - b| \geq ||a| - |b||$

*Proof.*

- (i) For  $a \in \mathbb{R}$ ,  $a \geq 0$  or  $a \leq 0$ , and since  $|a| = \pm a$ , it follows that  $|a| \geq 0$ .
- (ii) If  $a \geq 0 \wedge b \geq 0$ , then  $ab \geq 0 \implies |a| \cdot |b| = ab = |ab|$ . If  $a \leq 0 \wedge b \leq 0$ , then  $ab \geq 0 \implies |a| \cdot |b| = (-a)(-b) = ab = |ab|$ . If  $a \geq 0 \wedge b \leq 0$ , then  $ab \leq 0 \implies |a| \cdot |b| = a(-b) = -ab = |ab|$ . If  $a \leq 0 \wedge b \geq 0$ , then  $ab \leq 0 \implies |a| \cdot |b| = (-a)b = -ab = |ab|$ .
- (iii) By definition, we have  $-|a| \leq a \leq |a|$  and  $-|b| \leq b \leq |b|$ , then O4 yields:

$$-|a| - |b| \leq a + b \leq |a| + |b| \quad (5.4.2)$$

so that

$$-(|a| + |b|) \leq a + b \leq |a| + |b| \quad (5.4.3)$$

which implies  $a + b \leq |a| + |b|$  and  $-(a + b) \leq |a| + |b|$ . Since  $|a + b| = \pm(a + b)$ , it follows that  $|a + b| \leq |a| + |b|$ .

(iv) We proceed as follows:

$$|a - b| \geq ||a| - |b|| \iff (a - b)^2 \geq (|a| - |b|)^2 \quad (5.4.4)$$

$$\iff a^2 - 2ab + b^2 \geq a^2 - 2|a||b| + b^2 \quad (5.4.5)$$

$$\iff -2ab \geq -2|ab| \quad (5.4.6)$$

$$\iff ab \leq |ab| \quad (5.4.7)$$

which is true since  $x \leq |x|$  for any real number  $x$ .

■

## 5.5 The Completeness Axiom

This axiom assures us that unlike  $\mathbb{Q}$ ,  $\mathbb{R}$  has no gaps.

### Definition 5.10 (Maximum and minimum)

Let  $S$  be a nonempty subset of  $\mathbb{R}$ .

- (i) If  $S$  contains a largest element  $s_0$  s.t.  $s_0 \in S \wedge s \leq s_0, \forall s \in S$ , we write  $s = \max S$ .
- (ii) If  $S$  contains a smallest element we write it as  $\min S$ .

### Example.

- (i) The set  $\{r \in \mathbb{Q} | 0 \leq r \leq \sqrt{2}\}$  has a minimum, namely 0, but no maximum, since  $\sqrt{2} \notin \mathbb{Q}$ .
- (ii) Consider the set  $\{n^{(-1)^n} | n \in \mathbb{N}\}$ , which can be expanded into:

$$\{1, 2, \frac{1}{3}, 4, \frac{1}{5}, 6, \frac{1}{7}, \dots\} \quad (5.5.1)$$

which clearly has no maximum nor minimum.

◀

### Definition 5.11 (Upper and lower bound)

Let  $S$  be a nonempty subset of  $\mathbb{R}$ :

- (i) If a real number  $M$  satisfies  $s \leq M, \forall s \in S$ , then  $M$  is called an upper bound of  $S$ .
- (ii) If a real number  $m$  satisfies  $m \leq s, \forall s \in S$ , then  $m$  is called a lower bound of  $S$ .
- (iii) A set  $S$  is bounded if it is bounded above and below i.e.  $\exists m, M$  s.t.  $S \subseteq [m, M]$ .

**Remark.** Clearly, if a set  $S$  has a maximum, it is bounded above. Similarly, if it has a minimum, it is bounded below.

**Definition 5.12 (Supremum and infimum)**

Let  $S$  be a nonempty subset of  $\mathbb{R}$ .

- (i) If  $S$  is bounded above and has a least upper bound, then we call it the *supremum* of  $S$ , denoted by  $\sup S$ .
- (ii) If  $S$  is bounded above and has a least upper bound, then we call it the *infimum* of  $S$ , denoted by  $\inf S$ .

**Remark.** Observe that if  $S$  is bounded above, then  $M = \sup S$  iff:

- (i)  $s \leq M, \forall s \in S$
- (ii)  $\forall M_1 < M, \exists s_1 \in S \text{ s.t. } s_1 > M_1$

**Example.**

- (a) If a set  $S$  has a maximum, then  $\max S = \sup S$ . Similarly, if a set  $S$  has a minimum, then  $\min S = \inf S$ .
- (b) We have  $\inf\{n^{(-1)^n} : n \in \mathbb{N}\} = 0$ .
- (c) The set  $A = \{\frac{1}{n^2} : n \in \mathbb{N} \wedge n \geq 3\}$  is bounded. We have that  $\sup A = \max A = \frac{1}{9}$  and the minimum does not exist, however  $\inf A = 0$ .
- (d) The set  $B = \{r \in \mathbb{Q} : r^3 \leq 7\}$  is bounded above, but not below. It has no maximum, since  $\sqrt[3]{7} \notin \mathbb{Q}$ . However,  $\sup B = \sqrt[3]{7}$  and  $\inf B = -\infty$  since it has no minimum.
- (e) The set  $C = \{m + n\sqrt{2} : m, n \in \mathbb{Z}\}$  is not bounded above or below, so it has no maximum or minimum. However,  $\sup C = \infty$  and  $\inf C = -\infty$ .
- (f) The set  $D = \{x \in \mathbb{R} : x^2 < 10\}$  is the open interval  $(-\sqrt{10}, \sqrt{10})$ . So, it is bounded above and below despite not having maximum and minimum. We have  $\sup D = \sqrt{10}$  and  $\inf D = -\sqrt{10}$ .

**Theorem 5.13 (Completeness Axiom)**

Every nonempty subset  $S$  of  $\mathbb{R}$  that is bounded above has a least upper bound. In other words,  $\sup S$  exists and is a real number.

**Remark.** Note that by this definition, the set of rationals  $\mathbb{Q}$  is incomplete, that is, it contains "gaps". Indeed, consider the set  $A = \{r \in \mathbb{Q} : 0 \leq r \leq \sqrt{2}\}$ , which is bounded above by  $\frac{3}{2} \in \mathbb{Q}$  for example. If  $\mathbb{Q}$  were complete, then  $A$  would have a least upper bound that is rational, but such a number does not exist. **Corollary** Every nonempty subset  $S$  of  $\mathbb{R}$  that is bounded below has a greatest lower bound  $\inf S$ .

*Proof.* Let  $-S$  be the set  $\{-s : s \in S\}$  consisting of the negatives of  $S$ . Since  $S$  is bounded below,  $\exists m \in \mathbb{R} \text{ s.t. } m \leq s \forall s \in S$ . This implies that  $-m \geq -s, \forall s \in S$ , so since  $-S$  is bounded above by  $-m$ , by the Completeness Axiom it must have a supremum. Let us now prove that  $\inf S = -\sup(-S)$ . Let  $s_0 = \sup(-S)$ , then by definition  $-s \leq s_0 \implies s \geq -s_0$ , and since  $s_0$  is the least upper bound of  $-S$ , then if  $t \geq -s, \forall s \in S$ , then  $t \geq s_0$ . So, if  $-t \leq s, \forall s \in S$ , then  $-t \leq -s_0$ . These two conditions show that  $-s_0$  is the greatest lower bound of  $S$ , so  $\inf S = -\sup(-S)$ , as required. ■

**Theorem 5.14 (Archimedean Property)**

If  $a > 0$  and  $b > 0$ , then for some positive integer  $n$ , we have  $na > b$ .

*Proof.* Assume that Archimedean property fails, so there exists  $a > 0$  and  $b > 0$  such that  $na \leq b$ ,  $\forall n \in \mathbb{N} \implies b$  is an upper bound of  $S = \{na : n \in \mathbb{N}\}$ . Let  $s_0 = \sup S$ , and  $a > 0 \implies s_0 - a < s_0$ . Since  $s_0$  is the supremum of  $S$ ,  $s_0 - a$  can't be an upper bound since it is smaller than  $s_0$ , it follows that  $\exists n_0 \in \mathbb{N}$  s.t.  $s_0 - a < n_0 a \implies s_0 < (n_0 + 1)a \in S$  so  $s_0$  is not an upper bound of  $S$ , which is a contradiction. ■

**Theorem 5.15 (Dense ness of  $\mathbb{Q}$ )**

If  $a, b \in \mathbb{R}$  and  $a < b$ , then  $\exists r \in \mathbb{Q}$  s.t.  $a < r < b$ .

*Proof.* We wish to prove that:

$$a < r = \frac{m}{n} < b \implies an < m < bn \quad (5.5.2)$$

for some integers  $m, n$ . Since  $b - a > 0$ , by the Archimedean property,  $\exists n \in \mathbb{N}$  s.t.  $n(n - a) > 1$  ■

so it is evident that there is an integer  $m$  between  $an$  and  $bn$  since their difference is greater than 1.

**Proposition 5.16 (Linearity of sup and inf)**

For non empty bounded subsets  $A$  and  $B$  of  $\mathbb{R}$ , we have:

$$\sup(A + B) = \sup A + \sup B, \inf(A + B) = \inf A + \inf B \quad (5.5.3)$$

*Proof.* Consider  $x \in A + B \implies x = a + b$  for some  $a \in A, b \in B$ . It follows that  $x \leq \sup A + \sup B \implies \sup(A + B) \leq \sup A + \sup B$ . It remains to prove that  $\sup(A + B) \geq \sup A + \sup B$ . If one of the suprema is  $+\infty$  (without loss of generality assume it is  $B$ ), then taking some  $a_0 \in A$ , we have  $\sup(A + B) \geq \sup(a_0 + B) = a_0 + \sup B = \infty = \sup A + \sup B$ . If the sum of the suprema is finite, then we consider  $\epsilon > 0$ . Then  $\exists a \in A, b \in B$ , s.t.  $a > \sup A - \frac{\epsilon}{2}$  and  $b > \sup B - \frac{\epsilon}{2}$ . It follows that  $\sup(A + B) \geq a + b > \sup A + \sup B - \epsilon$  from which it follows that  $\sup(A + B) \geq \sup A + \sup B$ . ■

# Unit D2: Sequences

## 6.1 Introduction to sequences

A sequence is a function whose domain a set of the form  $\{n \in \mathbb{Z} : n \geq m\}$ . The sequence is denoted by  $(s)_{n=m}^{\infty}$  and the  $n$ th term in a sequence is denoted by  $s_n$ .

**Example.** Consider the sequence  $(a_n)_{n=0}^{\infty}$  where  $a_n = (-1)^n$ ,  $n \geq 0$ . We can then write the sequence as  $(1, -1, 1, -1\dots)$ , and has a set of values  $\{-1, 1\}$ . Note that the sequence contains infinite terms, but the set contains only two terms. 

### Definition (*Monotonic sequence*)

A sequence  $(a_n)$  is said to be:

- (i) **constant** if  $a_{n+1} = a_n$ , for  $n = 1, 2, 3\dots$
- (ii) **increasing (decreasing)** if  $a_{n+1} \geq a_n$  ( $a_{n+1} \leq a_n$ ) for  $n = 1, 2, 3\dots$
- (iii) **strictly increasing** if  $a_{n+1} > a_n$  ( $a_{n+1} < a_n$ ) for  $n = 1, 2, 3\dots$

If any of the above hold for  $(a_n)$ , then it is said to be **monotonic**.

Given a general sequence  $(a_n)$ , it is generally easier to show that:

- (i)  $a_{n+1} - a_n \geq (\leq)0 \implies (a_n)$  is increasing (decreasing)
- (ii)  $a_{n+1} - a_n > (<)0 \implies (a_n)$  is strictly increasing (strictly decreasing)
- (iii)  $a_{n+1} - a_n = 0 \implies (a_n)$  is constant

**Example.** Consider the sequence  $a_n = (n - 1)(n - 2)$ ,  $n = 1, 2\dots$ , which is monotonic increasing. Indeed, note that:

$$a_{n+1} - a_n = 2n - 2 \geq 0 \quad (6.1.1)$$

since  $n \geq 1$ . Note moreover that if  $n \geq 2$ , then  $(a_n)$  would be monotonic strictly increasing, since  $2n - 2 > 0$ . 

Alternatively, it is also convenient to examine the quotient  $\frac{a_{n+1}}{a_n}$ .

**Example.** Consider the sequence  $a_n = n + \frac{1}{n}$ . Then:

$$\frac{a_{n+1}}{a_n} = \frac{n+1 + \frac{1}{n+1}}{n + \frac{1}{n}} = \frac{(n+1)^2}{n^2} \cdot \frac{n}{n+1} = \frac{n+1}{n} > 1 \quad (6.1.2)$$

since  $n$  is positive. It follows that  $a_n$  is monotonic strictly increasing.  $\blacktriangleleft$

### Definition (Eventual properties)

A sequence  $(a_n)$  eventually has a property if it satisfies the property for  $n \geq n_0$  for some  $n_0 \geq 1$ .

**Example.** The sequence defined by  $a_n = \frac{n^4}{4^n}$  is eventually decreasing. Indeed:

$$\frac{a_{n+1}}{a_n} = \frac{(n+1)^4}{n^4} \cdot \frac{4^n}{4^{n+1}} = \frac{1}{4} \left( \frac{n+1}{n} \right)^4 < 1 \implies 1 + \frac{1}{n} < \sqrt{2} \implies \frac{1}{\sqrt{2}-1} < n \quad (6.1.3)$$

so the sequence eventually decreases, more specifically for  $n \geq 3$ .  $\blacktriangleleft$

## 6.2 Convergence of sequences

### Definition 6.1 (Sequence convergence)

A sequence  $(s_n)$  of real numbers converges to  $s$  (i.e.  $\lim_{n \rightarrow \infty} s_n = s$ ) provided that:

$$\forall \epsilon > 0, \exists N \in \mathbb{R} \text{ s.t. } n > N \implies |s_n - s| < \epsilon \quad (6.2.1)$$

A sequence that does not converge to a real number is divergent.

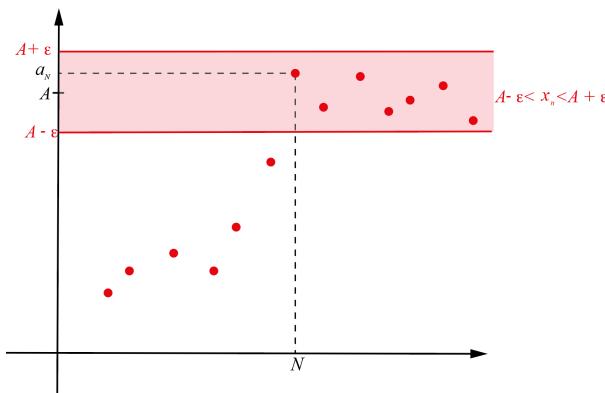


Figure 6.1. Geometrical interpretation of the epsilon-delta definition

**Example.** Consider the sequence  $s_n = \frac{3n+1}{7n-4} = \frac{3+\frac{1}{n}}{7-\frac{4}{n}}$  then clearly for large values of  $n$ , the

series should converge to  $\frac{3}{7}$ . Indeed, by definition 7.1,  $\lim_{n \rightarrow \infty} s_n = \frac{3}{7}$  means that:

$$\forall \epsilon > 0, \exists N \text{ s.t. } n > N \implies \left| \frac{3n+1}{7n-4} - \frac{3}{7} \right| < \epsilon \quad (6.2.2)$$

As  $\epsilon$  varies, getting smaller and smaller,  $N$  gets bigger and bigger, so in the end for  $n > N$ , so a very large value of  $n$ , the difference between  $s_n$  and  $s$  becomes very very small, which intuitively makes sense.  $\blacktriangleleft$

**Remark.** Finally, it must be noted that limits are unique, so:

$$\lim_{n \rightarrow \infty} s_n = s \wedge \lim_{n \rightarrow \infty} s_n = t \implies s = t \quad (6.2.3)$$

Indeed, the first implies that:

$$\exists N_1 \text{ s.t. } n > N_1 \implies |s_n - s| < \frac{\epsilon}{2} \quad (6.2.4)$$

and the second implies that:

$$\exists N_2 \text{ s.t. } n > N_2 \implies |s_n - t| < \frac{\epsilon}{2} \quad (6.2.5)$$

for some  $\epsilon > 0$ . For  $n > \max\{N_1, N_2\}$  (this allows us to use both conditions of convergence) the triangle inequality shows:

$$|s - t| = |(s - s_n) + (s_n - t)| \leq |s - s_n| + |s_n - t| \leq \frac{\epsilon}{2} + \frac{\epsilon}{2} = \epsilon \quad (6.2.6)$$

for all  $\epsilon > 0$  and thus  $|s - t| = 0 \implies s = t$  as required. Geometrically, this argument corresponds to showing that for  $n > \max\{N_1, N_2\}$ , the terms of the sequence can't belong to both  $\{y_1 : |y_1 - s| < \frac{\epsilon}{2}\}$  and  $\{y_2 : |y_2 - t| < \frac{\epsilon}{2}\}$  since the two don't intersect for sufficiently small  $\epsilon$ .  $\blacksquare$

$n$	$s_n = \frac{3n+1}{7n-4}$	$ s_n - \frac{3}{7} $
2	0.7000	0.2714
3	0.5882	0.1597
5	0.5161	0.0876
40	0.4384	0.0098
400	0.4295	0.0010

## 6.3 Formal Proofs of Limit Theorems

**Example.** Let us prove that  $\lim_{n \rightarrow \infty} \frac{3n+1}{7n-4} = \frac{3}{7}$ .

Discussion: we consider an arbitrary  $\epsilon > 0$  and show that  $\exists N$  such that  $n > N \implies \left| \frac{3n+1}{7n-4} - \frac{3}{7} \right| < \epsilon$ . Thus, we want:

$$\left| \frac{19}{7(7n-4)} \right| < \epsilon \quad (6.3.1)$$

and since  $7n - 4 > 0$  since  $n$  is positive, we drop the absolute value and write:

$$\frac{19}{7(7n-4)} < \epsilon \implies \frac{19}{7\epsilon} < 7n - 4 \implies \frac{19}{49\epsilon} + \frac{4}{7} < n \quad (6.3.2)$$

so we put  $N = \frac{19}{49\epsilon} + \frac{4}{7}$ .

Proof: let  $\epsilon > 0$  and  $N = \frac{19}{49\epsilon} + \frac{4}{7}$ . Then

$$n > N \implies n > \frac{19}{49\epsilon} + \frac{4}{7} \implies 7n - 4 > \frac{19}{7\epsilon} \implies \epsilon > \frac{19}{7(7n-4)} \quad (6.3.3)$$

Thus  $n > N \implies |\frac{3n+1}{7n-4} - \frac{3}{7}| < \epsilon$  which proves that  $\lim_{n \rightarrow \infty} \frac{3n+1}{7n-4} = \frac{3}{7}$  as required.  $\blacktriangleleft$

**Example.** Let us prove that  $\lim_{n \rightarrow \infty} \frac{4n^3+3n}{n^3-6} = 4$ .

Discussion: For each  $\epsilon > 0$ , we need the following inequality to hold:

$$|\frac{4n^3+3n}{n^3-6} - 4| < \epsilon \implies |\frac{3n+24}{n^3-6}| < \epsilon \quad (6.3.4)$$

Note that it is very hard in this case to solve for  $n$ , so instead of finding the smallest  $N$  such that  $n > N$  implies that  $|\frac{3n+24}{n^3-6}| < \epsilon$ , we will use estimates. Note that  $3n+24 \leq 27n$  for  $n > 1$  and  $n^3 - 6 \geq \frac{n^3}{2}$  for  $n > 2$ . So:

$$|\frac{3n+24}{n^3-6}| \leq \frac{27n}{\frac{1}{2}n^3} < \epsilon \implies n > \sqrt{\frac{54}{\epsilon}} \quad (6.3.5)$$

if  $n > 2$ .

Proof: let  $\epsilon > 0$  and  $N = \max\{2, \sqrt{\frac{54}{\epsilon}}\}$ . Then:

$$n > N \implies n > \sqrt{\frac{54}{\epsilon}} \implies \frac{27n}{n^3/2} < \epsilon \quad (6.3.6)$$

Since for  $n > 2$ ,  $3n+24 \leq 27n$  and  $n^3 - 6 \geq \frac{n^3}{2}$ , we find that:

$$|\frac{3n+24}{n^3-6}| \leq \frac{27n}{n^3/2} < \epsilon \quad (6.3.7)$$

and hence

$$|\frac{4n^3+3n}{n^3-6} - 4| < \epsilon \quad (6.3.8)$$

as required.  $\blacktriangleleft$

**Example.** Let  $(s_n)$  be a sequence of non-negative real numbers such that  $\lim_{n \rightarrow \infty} s_n = s$ . Then,  $\lim_{n \rightarrow \infty} \sqrt{s_n} = \sqrt{s}$ .

*Proof.*

Case I: let  $\epsilon > 0$  and  $s > 0$ , since  $\lim_{n \rightarrow \infty} s_n = s$ ,

$$\exists N \text{ s.t. } n > N \implies |s_n - s| < \sqrt{s}\epsilon \quad (6.3.9)$$

so:

$$\exists N \text{ s.t. } n > N \implies |\sqrt{s_n} - \sqrt{s}| = \frac{|s_n - s|}{\sqrt{s_n} + \sqrt{s}} \leq \frac{|s_n - s|}{\sqrt{s}} < \frac{\sqrt{s}}{\sqrt{s}}\epsilon = \epsilon \quad (6.3.10)$$

as desired.  $\blacktriangleleft$

Case II: if  $s = 0$ , let  $\epsilon > 0$  so that:

$$\exists N \text{ s.t. } n > N \implies |s_n| < \epsilon^2 \quad (6.3.11)$$

Hence,  $\sqrt{s_n} < \epsilon$  for  $n > N$  and thus:

$$|\sqrt{s_n} - 0| < \epsilon \implies \lim_{n \rightarrow \infty} \sqrt{s_n} = s = 0 \quad (6.3.12)$$

as desired.  $\blacktriangleleft$

**Example.** Let  $(s_n)$  be a convergent sequence of non-zero real numbers such that  $\lim_{n \rightarrow \infty} s_n = s \neq 0$ . Then  $\inf\{|s_n| : n \in \mathbb{N}\} > 0$ .

Discussion: The result has the geometric interpretation that all terms of the sequence are "close" to  $s$  and therefore not "close" to 0. The proof will involve three steps. First, we show that there exists  $N$  such that all terms of the sequence after  $s_N$  are all greater than  $\frac{|s|}{2}$  by using the triangle inequality. This result shows that the terms of the sequence after  $s_N$  are all "close" to  $s$ , with a maximum distance of  $\frac{|s|}{2}$ . We then take the minimum of  $\frac{|s|}{2}$ , and  $|s|_n$ , and show that it is positive so that  $\inf\{|s_n| : n \in \mathbb{N}\} > 0$  follows directly.

*Proof.*

Let  $\epsilon = \frac{|s|}{2} > 0$ , since  $\lim_{n \rightarrow \infty} s_n = s$ :

$$\exists N \text{ s.t. } n > N \implies |s_n - s| \leq \frac{|s|}{2} \quad (6.3.13)$$

We must have that:

$$\exists N \text{ s.t. } n > N \implies |s_n| \geq \frac{|s|}{2} \quad (6.3.14)$$

since:

$$|s| = |s - s_n + s_n| \leq |s - s_n| + |s_n| \leq \frac{|s|}{2} + |s_n| \implies |s_n| \geq \frac{|s|}{2}. \quad (6.3.15)$$

Setting:

$$m = \min\left\{\frac{|s|}{2}, |s_1|, \dots, |s_N|\right\} \quad (6.3.16)$$

In view of this result,  $|s_n| \geq m$ , so  $\inf\{|s_n| : n \in \mathbb{N}\} > 0$  as required.  $\blacktriangleleft$

## 6.4 Null sequences

### Definition (Null sequence)

A sequence  $(a_n)$  is said to be null if it converges to zero, that is, if:

$$\forall \epsilon > 0, \exists N \in \mathbb{R} \text{ s.t. } n > N \implies |a_n| < \epsilon \quad (6.4.1)$$

**Example.** Consider the sequence

$$a_n = \frac{(-1)^n}{n^4 + 1}, \quad n = 1, 2, \dots \quad (6.4.2)$$

Suppose  $0 < \epsilon < 1$ , then:

$$|a_n| = \frac{1}{n^4 + 1} < \epsilon \iff n^4 > \frac{1}{\epsilon} - 1 \iff n > \left(\frac{1}{\epsilon} - 1\right)^{\frac{1}{4}} = N \quad (6.4.3)$$

so for a given  $0 < \epsilon < 1$ ,  $|a_n| < \epsilon$  provided that  $n > N$  where  $N = \left(\frac{1}{\epsilon} - 1\right)^{\frac{1}{4}}$ .

If instead  $\epsilon \geq 1$ , then:

$$|a_n| = \frac{1}{n^4 + 1} < \epsilon \iff n^4 > \frac{1}{\epsilon} - 1 \quad (6.4.4)$$

but since  $\epsilon \geq 1 \implies \frac{1}{\epsilon} \leq 1$ , the LHS is non-positive, so the above inequality is true for all  $n > N$  where  $N = 1$ .  $\blacktriangleleft$

### Theorem (Power rule of null sequences)

If  $(a_n)$  is a null, non-negative sequence for  $n = 1, 2, \dots$ , then the sequence  $(a_n^p)$  given by  $a_n^p$ , with  $p \in \mathbb{R}$  for all  $n = 1, 2, \dots$  is also null.

*Proof.* The sequence  $(a_n)$  is null, therefore for each positive  $\epsilon^{1/p}$  there exists  $N$  such that:

$$a_n < \epsilon^{1/p}, \quad \forall n > N \quad (6.4.5)$$

so that for all positive  $\epsilon$ :

$$a_n^p < \epsilon, \quad \forall n > N \quad (6.4.6)$$

so  $\lim_{n \rightarrow \infty} a_n^p = 0$  as desired.  $\blacksquare$

**Example.** Consider the sequence  $a_n = \frac{1}{n}$  for  $n = 1, 2, \dots$ . Then  $a_n$  is null, since for all  $\epsilon > 0$ :

$$|a_n| = \frac{1}{n} < \epsilon \iff \epsilon < n, \quad (6.4.7)$$

so that  $|a_n| \leq \epsilon$  for all  $n > N = \epsilon$ . Hence  $\lim_{n \rightarrow \infty} a_n = 0$ . Applying the power rule to  $a_n$  we then find that  $\frac{1}{n^p}$  is also null.  $\blacktriangleleft$

**Proposition (Limit theorems for null sequences)**

Let  $(a_n)$ ,  $(b_n)$  and  $(c_n)$  be null sequences, and let  $\alpha \in \mathbb{R}$ . Then:

- (i)  $(\alpha a_n)$  is null
- (ii)  $(a_n + b_n)$  is null
- (iii)  $(a_n b_n)$  is null
- (iv)  $\left(\frac{a_n}{b_n}\right)$  is null

*Proof.* Immediate application of the Limit theorems (which will be shown in the next section) will give the desired results. ■

**Proposition (Standard null sequences)**

The following sequences are all null:

- (i)  $\left(\frac{1}{n^p}\right)$  for  $p > 0$
- (ii)  $(c^n)$  for  $|c| < 1$
- (iii)  $(n^p c^n)$  for  $p > 0$ ,  $|c| < 1$
- (iv)  $\left(\frac{c^n}{n!}\right)$  for  $c \in \mathbb{R}$
- (v)  $\left(\frac{n^p}{n!}\right)$  for  $p > 0$ .

*Proof.* (i) It has been proven in the example preceding the limit theorems.

(ii) Note that  $a_n$  is null  $\iff |a_n|$  is also null, so it suffices to prove that  $|c^n| = |c|^n$  is null (see the Theorem on the convergence of absolute values). In other words, we can limit ourselves to  $c$  non-negative.

Suppose that  $c = 0$ , then the nullity is trivial. Suppose that  $0 < c < 1$ , then:

$$c = \frac{1}{a+1}, \quad a > 0 \quad (6.4.8)$$

We may apply Bernoulli's inequality:

$$c^n = \frac{1}{(a+1)^n} \leq \frac{1}{na+1} \leq \frac{1}{na} \quad (6.4.9)$$

Since  $\frac{1}{n}$  is null, it follows from the limit theorem (ii) that  $\frac{1}{na}$  is also null. Applying the squeeze theorem then, we finally find that  $c^n$  is also null.

(iii) Once again assume that  $0 < c < 1$ , then for some  $a > 0$ :

$$c = \frac{1}{a+1} \quad (6.4.10)$$

We begin by proving that  $(nc^n)$  is null. By the binomial theorem:

$$(1+a)^n \geq 1 + na + \frac{1}{2}n(n-1)a^2 \geq \frac{1}{2}n(n-1)a^2, \quad n = 2, 3, \dots \quad (6.4.11)$$

it follows that

$$nc^n = \frac{n}{(a+1)^n} \leq \frac{n}{\frac{1}{2}n(n-1)a^2} = \frac{2}{a^2(n-1)}, \quad n = 2, 3, \dots \quad (6.4.12)$$

The sequence defined by:

$$b_n = \frac{2}{a^2(n-1)}, \quad n = 2, 3, \dots \quad (6.4.13)$$

is equivalent to:

$$b_n = \frac{2}{a^2 n}, \quad n = 1, 2, \dots \quad (6.4.14)$$

which is null by the limit theorems. By the squeeze rule, it follows that  $(nc^n)$  is also null.

For  $(n^p c^n)$ , where  $p > 0$  and  $0 < c < 1$  we write:

$$n^p c^n = (nd^n)^p, \quad n = 1, 2, \dots \quad (6.4.15)$$

where  $0 < d = c^{1/p} < 1$ . We have shown that  $(nd^n)$  is null, so that by the power rule  $(n^p c^n)$  is also null.

- (iv) Again, as in the case of (ii) we may consider only positive values of  $c$ . Let us choose an integer  $m$  such that  $m + 1 > c$ , so that for  $n > m + 1$ :

$$\frac{c^n}{n!} = \prod_{k=1}^n \frac{c}{k} \leq \frac{c}{n} \prod_{k=1}^m \frac{c}{k} = \frac{c^m}{m!} \cdot \frac{c}{n} \quad (6.4.16)$$

Since  $\frac{1}{n}$  is null, we have that  $\frac{c^m}{m!} \cdot \frac{c}{n}$  too is null, since  $\frac{c^{m+1}}{m!}$  is just a constant. By the squeeze rule, it follows that  $\frac{c^n}{n!}$  is null, as desired.

- (v) We can write:

$$\frac{n^p}{n!} = \frac{n^p}{2^n} \cdot \frac{2^n}{n!} \quad (6.4.17)$$

so that  $\left(\frac{n^p}{n!}\right)$  is null by the limit theorem (iii). ■

**Example.** Consider the sequence described by

$$a_n = \frac{1}{3n^4(2n-1)^{1/3}}, \quad n = 1, 2, \dots \quad (6.4.18)$$

We can rewrite the terms of this sequence as:

$$a_n = \frac{1}{3} \cdot \frac{1}{n^4} \cdot \left(\frac{(-1)^n}{n^4+1}\right)^{1/3} \quad (6.4.19)$$

We know that  $\frac{1}{n^4}$  is null, since it is the first standard null sequence with  $p = 1$ . Moreover, we have determined previously that  $\frac{(-1)^n}{n^4+1}$  is null, so that by the power rule,  $\left(\frac{(-1)^n}{n^4+1}\right)^{1/3}$  is also null. Finally, we exploit the product rule to conclude that  $a_n$  is indeed null. ◀

**Theorem (Squeeze rule for null sequences)**

If  $(b_n)$  is a null sequence of non-negative terms and:

$$|a_n| \leq b_n, \forall n = 1, 2, \dots \quad (6.4.20)$$

then  $(a_n)$  is null.

*Proof.* We apply the squeeze rule (which will be proven in the next section) with:

$$-b_n \leq a_n \leq b_n, \forall n = 1, 2, \dots \quad (6.4.21)$$

then since  $(-b_n)$  and  $(b_n)$  both converge to 0, then  $\lim_{n \rightarrow \infty} a_n = 0$  as well. ■

**Example.** Consider the sequence:

$$a_n = \frac{\sin(n^2)}{n^2 + 2^n}, n = 1, 2, \dots \quad (6.4.22)$$

Then:

$$\left| \frac{\sin(n^2)}{n^2 + 2^n} \right| \leq \left| \frac{1}{n^2 + 2^n} \right| \quad (6.4.23)$$

Since both  $n^2$  and  $2^n$  are positive:

$$\frac{1}{n^2 + 2^n} < \frac{1}{n^2} \quad (6.4.24)$$

so:

$$\left| \frac{\sin(n^2)}{n^2 + 2^n} \right| \leq \frac{1}{n^2} \quad (6.4.25)$$

However,  $\frac{1}{n^2}$  is one of the standard null sequences, so we find by the Squeeze rule that  $a_n$  is also null.



## 6.5 Limit theorems for Convergent sequences

**Definition (Bounded and unbounded sequence)**

A sequence  $(a_n)$  is **bounded** if  $\exists M$  such that:

$$|a_n| \leq M, n = 1, 2, \dots \quad (6.5.1)$$

and is **unbounded** otherwise

**Example.** The sequence  $a_n = \frac{2n+1}{n}$  for  $n = 1, 2, \dots$  is bounded. Indeed, for all natural numbers  $n$  we find that:

$$a_n = \frac{2n+1}{n} = 2 + \frac{1}{n} \leq 2 + 1 = 3 \quad (6.5.2)$$

since  $\frac{1}{n} \leq 1$ .



**Example.** Consider the sequence  $b_n = (-1)^n n$ , for  $n = 1, 2, \dots$ , then suppose there exists some  $M$  such that:

$$|b_n| = |n| = n < M \quad (6.5.3)$$

for all  $n = 1, 2, \dots$ , which is clearly a contradiction since  $\mathbb{N}$  has no maximum. Consequently,  $b_n$  is unbounded.



### Proposition (*Boundedness of convergent sequences*)

Convergent sequences are bounded, and unbounded sequences are divergent.

*Proof.* Suppose  $(s_n)$  is a convergent sequence with  $\lim_{n \rightarrow \infty} s_n = s$ . Then we can apply the Definition 6.1 with  $\epsilon = 1$  to find:

$$n > N \implies |s_n - s| < 1 \implies |s_n| < 1 + |s| \quad (6.5.4)$$

where we use the triangle inequality. Let us now define  $M = \max\{|s| + 1, |s_1|, \dots, |s_N|\}$  so that  $|s_n| \leq M$  for all  $n$ , proving the boundedness of convergent sequences.

The converse follows immediately. ■

**Remark.** Note that it does not suffice to show that  $|s_n| < 1 + |s|$ , because this only proves that the terms of the sequence after  $s_N$  are bounded. We must introduce  $M$  in order to prove boundedness for the initial  $N$  terms.

Also note that the choice  $\epsilon = 1$  was completely arbitrary.

**Example.** Consider the sequence  $a_n = n^{(-1)^n}$  for  $n = 1, 2, \dots$ . Then, for  $n$  even:

$$|a_n| = n^{(-1)^n} = n \quad (6.5.5)$$

which is unbounded, since  $\mathbb{N}$  is also unbounded. Instead for  $n$  odd:

$$|a_n| = \frac{1}{n} < 1 \quad (6.5.6)$$

which is bounded. Consequently,  $a_n$  overall is unbounded, and thus divergent.



**Example.** Consider the sequence  $a_n = \frac{n^2+n}{n^2+1}$  for  $n = 1, 2, \dots$ . Then

$$\frac{n^2+n}{n^2+1} \leq \frac{n^2+n^2}{n^2} = 2 \quad (6.5.7)$$

so the sequence is bounded. Moreover,

$$\lim_{n \rightarrow \infty} \frac{n^2+n}{n^2+1} = \lim_{n \rightarrow \infty} \frac{1 + \frac{1}{n}}{1 + \frac{1}{n^2}} = 1 \quad (6.5.8)$$

so it also converges. ◀

### Theorem 6.3 (*Limit theorems*)

Let  $(s_n), (t_n)$  sequences converging to  $s, t$  respectively, and let  $\alpha \in \mathbb{R}$ . Then:

- (i)  $\lim_{n \rightarrow \infty} (ks_n) = k \lim_{n \rightarrow \infty} s_n$
- (ii)  $\lim_{n \rightarrow \infty} (s_n + t_n) = \lim_{n \rightarrow \infty} s_n + \lim_{n \rightarrow \infty} t_n$
- (iii)  $\lim_{n \rightarrow \infty} (s_n t_n) = (\lim_{n \rightarrow \infty} s_n)(\lim_{n \rightarrow \infty} t_n)$
- (iv)  $\lim_{n \rightarrow \infty} \left( \frac{t_n}{s_n} \right) = \frac{\lim_{n \rightarrow \infty} t_n}{\lim_{n \rightarrow \infty} s_n}$  if  $\lim_{n \rightarrow \infty} s_n \neq 0$  and  $s_n \neq 0$  for all  $n$

*Proof.*

- (i) If  $k = 0$  then the result is trivial. If we assume that  $k \neq 0$ , and we let  $\epsilon > 0$ , then since  $\lim_{n \rightarrow \infty} s_n = s$  for some  $s \in \mathbb{R}$  there exists  $N$  such that:

$$n > N \implies |s_n - s| < \frac{\epsilon}{|k|} \quad (6.5.9)$$

Then:

$$n > N \implies |ks_n - ks| < \epsilon \quad (6.5.10)$$

showing that  $\lim_{n \rightarrow \infty} (ks_n) = k \lim_{n \rightarrow \infty} s_n$  as desired.

- (ii) Let  $\epsilon > 0$ , we need to show that for appropriately large  $n$ :

$$|s_n + t_n - (s + t)| < \epsilon \quad (6.5.11)$$

Note however that since  $\lim_{n \rightarrow \infty} s_n = s$ , there exists  $N_1$  such that:

$$n > N_1 \implies |s_n - s| < \frac{\epsilon}{2} \quad (6.5.12)$$

and similarly for  $t_n$ :

$$n > N_2 \implies |t_n - t| < \frac{\epsilon}{2} \quad (6.5.13)$$

so that, letting  $N = \max\{N_1, N_2\}$

$$n > N \implies |s_n - s| + |t_n - t| < \epsilon \quad (6.5.14)$$

Finally, we make use of the triangle inequality:

$$n > N \implies |s_n + t_n - (s + t)| \leq |s_n - s| + |t_n - t| < \epsilon \quad (6.5.15)$$

as desired.

(iii) Firstly note that:

$$|s_n t_n - st| = |s_n t_n - s_n t + s_n t - st| \leq |s_n| \cdot |t_n - t| + |t| \cdot |s_n - s| \quad (6.5.16)$$

Let  $\epsilon > 0$ , by theorem 6.2  $s_n$  must be bounded, so we can find  $M > 0$  such that  $|s_n| \leq M$  for all  $n$ . Since  $t_n$  converges:

$$n > N_1 \implies |t_n - t| < \frac{\epsilon}{2M} \quad (6.5.17)$$

and similarly for  $s_n$ :

$$n > N_2 \implies |s_n - s| < \frac{\epsilon}{2(|t| + 1)} \quad (6.5.18)$$

where we used  $|t| + 1$  since it could be the case that  $|t| = 0$ . Suppose  $N = \max\{N_1, N_2\}$  then:

$$n > N \implies |s_n t_n - st| \leq M \frac{\epsilon}{2M} + |t| \frac{\epsilon}{2(|t| + 1)} < \epsilon \quad (6.5.19)$$

as desired,  $\lim_{n \rightarrow \infty} (s_n t_n) = (\lim_{n \rightarrow \infty} s_n)(\lim_{n \rightarrow \infty} t_n)$ .

(iv) We begin by proving the following lemma:

**Lemma.** If  $s_n$  converges to  $s \neq 0$  and  $s_n \neq 0$  for all  $n$ , then  $\frac{1}{s_n}$  converges to  $\frac{1}{s}$ .

To prove this lemma, let  $\epsilon > 0$ . In the last example of the previous section, we showed that there exists  $m > 0$  such that  $|s_n| \geq m$  for all  $n$ . The convergence of  $s_n$  means that there exists  $N$  so that:

$$n > N \implies |s - s_n| < \epsilon m |s| \quad (6.5.20)$$

Then  $n > N$  implies:

$$\left| \frac{1}{s_n} - \frac{1}{s} \right| = \frac{|s - s_n|}{|s_n s|} \leq \frac{|s - s_n|}{m |s|} \leq \epsilon \quad (6.5.21)$$

as desired.

We may finally use this lemma to prove the main result:

$$\lim_{n \rightarrow \infty} \frac{t_n}{s_n} = \lim_{n \rightarrow \infty} \frac{1}{s_n} t_n = \frac{t}{s} \quad (6.5.22)$$

as desired. ■

**Remark.** Note that we also impose the condition  $s_n \neq 0$  in order for the reciprocal sequence  $\frac{1}{s_n}$  to be well-defined.

**Example.** Consider the sequence:

$$a_n = \frac{n^2 + 2^n}{3^n + n^3} \quad (6.5.23)$$

We can factorize out  $2^n$  on the numerator and  $3^n$  on the denominator to find that:

$$a_n = \frac{2^n}{3^n} \cdot \frac{\frac{n^2}{2^n} + 1}{\frac{n^3}{3^n} + 1} = \left(\frac{2}{3}\right)^n \cdot \frac{\frac{n^2}{2^n} + 1}{\frac{n^3}{3^n} + 1} \quad (6.5.24)$$

Hence:

$$\lim_{n \rightarrow \infty} a_n = \lim_{n \rightarrow \infty} \left[ \left(\frac{2}{3}\right)^n \cdot \frac{\frac{n^2}{2^n} + 1}{\frac{n^3}{3^n} + 1} \right] \quad (6.5.25)$$

$$= \lim_{n \rightarrow \infty} \left(\frac{2}{3}\right)^n \cdot \frac{\lim_{n \rightarrow \infty} \frac{n^2}{2^n} + 1}{\lim_{n \rightarrow \infty} \frac{n^3}{3^n} + 1} \quad (6.5.26)$$

$$= \lim_{n \rightarrow \infty} \left(\frac{2}{3}\right)^n \quad (6.5.27)$$

$$= 0 \quad (6.5.28)$$

so  $a_n$  is a null sequence.



### Theorem (Squeeze rule)

Let  $(a_n)$ ,  $(b_n)$  and  $(c_n)$  be convergent sequences such that:

- (i)  $b_n \leq a_n \leq c_n$  for  $n = 1, 2, \dots$
- (ii)  $\lim_{n \rightarrow \infty} b_n = \lim_{n \rightarrow \infty} c_n = l$

Then  $\lim_{n \rightarrow \infty} a_n = l$ .

*Proof.* From condition (ii) we have that for all  $\epsilon > 0$  then there exists  $N_1$  and  $N_2$  such that

$$l - \epsilon < b_n < l + \epsilon, \forall n > N_1 \quad (6.5.29)$$

$$l - \epsilon < c_n < l + \epsilon, \forall n > N_2 \quad (6.5.30)$$

Then, it follows from condition (i) that:

$$l - \epsilon < b_n \leq a_n \leq c_n < l + \epsilon, \forall n > N \quad (6.5.31)$$

where  $N = \max\{N_1, N_2\}$ , so that  $\lim_{n \rightarrow \infty} a_n = l$  as desired. ■

**Example.** Consider the sequence  $a_n = n^{1/n}$ . Using the binomial theorem for  $n \geq 2$  and  $x \geq 0$ :

$$(1 + x)^n \geq \frac{n(n-1)}{2} x^2 \quad (6.5.32)$$

since all the other terms in the binomial expansion are positive. Substituting  $x = \sqrt{\frac{2}{n-1}}$  we find that:

$$\left(1 + \sqrt{\frac{2}{n-1}}\right)^n \geq \frac{n(n+1)}{2} \frac{2}{n-1} = n \quad (6.5.33)$$

implying that

$$n^{1/n} \leq 1 + \sqrt{\frac{2}{n-1}}, \quad n = 2, 3, \dots \quad (6.5.34)$$

If we define  $b_n$  to be:

$$b_1 = 1 \quad (6.5.35)$$

$$b_n = 1 + \sqrt{\frac{2}{n-1}}, \quad n = 2, 3, \dots \quad (6.5.36)$$

then we see that  $a_n \leq b_n$  for  $n = 1, 2, 3, \dots$ . Moreover,  $\lim_{n \rightarrow \infty} b_n = 1$ . Indeed for all  $\epsilon > 0$ :

$$|b_n - 1| < \epsilon \iff n > \frac{2}{\epsilon^2} + 1 = N > 1 \quad (6.5.37)$$

so there exists some  $N$  such that  $|b_n - 1| < \epsilon$ . It is important that  $N > 1$ , since otherwise it can be the case for some  $n > N$  that  $b_n - 1 = 0$ .

It follows from the squeeze rule that:

$$\lim_{n \rightarrow \infty} n^{1/n} = 1 \quad (6.5.38)$$

◀

### Proposition (Limit inequality)

If  $\lim_{n \rightarrow \infty} a_n = l$ ,  $\lim_{n \rightarrow \infty} b_n = m$  and:

$$a_n \leq b_n, \quad \forall n = 1, 2, \dots \quad (6.5.39)$$

then  $l \leq m$ .

*Proof.* Suppose that  $\lim_{n \rightarrow \infty} a_n = l$ ,  $\lim_{n \rightarrow \infty} b_n = m$  and  $a_n \leq b_n$  for  $n = 1, 2, \dots$ . If  $l > m$  then:

$$\lim_{n \rightarrow \infty} (a_n - b_n) = l - m > 0 \quad (6.5.40)$$

so that:

$$a_n - b_n > \frac{1}{2}(l - m) \quad (6.5.41)$$

Indeed, if we take  $\epsilon = \frac{l}{2}$  then it follows that:

$$|(a_n - b_n) - (l - m)| < \frac{l}{2} \implies -\frac{l-m}{2} < (a_n - b_n) - (l - m) < \frac{l-m}{2} \quad (6.5.42)$$

$$\implies \frac{l-m}{2} < a_n, \quad \forall n > N \quad (6.5.43)$$

for some  $N$ . However, we also have that  $a_n - b_n \leq 0$  so that:

$$l - m \leq 0 \implies l \leq m \quad (6.5.44)$$

a contradiction.

**Proposition (Limit uniqueness)**

The limit of a sequence is unique.

Suppose  $\lim_{n \rightarrow \infty} a_n = l$  and  $\lim_{n \rightarrow \infty} a_n = m$ , then using the Limit inequality rule we get that  $l \leq m$  and  $m \leq l$ , giving  $l = m$  as desired. ■

**Theorem (Convergence of absolute value)**

If  $\lim_{n \rightarrow \infty} a_n = l$  then  $\lim_{n \rightarrow \infty} |a_n| = |l|$ .

*Proof.* Using the triangle inequality:

$$||a_n| - |l|| \leq |a_n - l| \quad (6.5.45)$$

Therefore, for all  $\epsilon > 0$ , there exists  $N$  such that:

$$|a_n - l| \leq \epsilon, \forall n > N \quad (6.5.46)$$

implying that:

$$||a_n| - |l|| \leq \epsilon, \forall n > N \quad (6.5.47)$$

It follows that  $\lim_{n \rightarrow \infty} |a_n| = |l|$  as expected. ■

## 6.6 Divergent sequences

**Definition (Infinite limit)**

A sequence  $(a_n)$  tends to infinity if:

$$\forall M > 0, \exists N \text{ s.t. } a_n > M, \forall n > N \quad (6.6.1)$$

We write that  $a_n \rightarrow \infty$ .

The sequence  $(a_n)$  tends to minus infinity if:

$$-a_n \rightarrow \infty \quad (6.6.2)$$

**Theorem (Reciprocal rule)**

If  $(a_n)$  is eventually positive, and  $\left(\frac{1}{a_n}\right)$  is a null sequence then  $a_n \rightarrow \infty$ .

*Proof.* Let  $M > 0$ , then since  $a_n$  is eventually positive:

$$a_n > 0, \forall n > N_1 \quad (6.6.3)$$

Moreover,  $\frac{1}{a_n}$  is null so that taking  $\epsilon = \frac{1}{M}$ :

$$\left| \frac{1}{a_n} \right| < \frac{1}{M}, \forall n > N_2 \quad (6.6.4)$$

Taking  $N = \max\{N_1, N_2\}$  then:

$$0 < \frac{1}{a_n} < \frac{1}{M}, \forall n > N \quad (6.6.5)$$

implying:

$$a_n > M, \forall n > N \quad (6.6.6)$$

as desired. ■

**Example.** Consider the sequence  $a_n = n! - 10^n$  for  $n = 1, 2, \dots$ . The dominant term is  $n!$ , so let us write:

$$a_n = n! \left(1 - \frac{10^n}{n!}\right), \quad n = 1, 2, \dots \quad (6.6.7)$$

Therefore,  $a_n$  is indeed eventually positive, since  $10^n < n!$  is true for  $n \geq 25$ . Then:

$$\lim_{n \rightarrow \infty} \frac{1}{a_n} = \lim_{n \rightarrow \infty} \frac{1}{n!} \frac{1}{1 - \frac{10^n}{n!}} = 0 \cdot \frac{0}{1 - 0} = 0 \quad (6.6.8)$$

By the reciprocal rule we see that  $a_n \rightarrow \infty$ .



### Proposition (*Properties of diverging sequences*)

If  $(a_n)$  and  $(b_n)$  both tend to infinity, then:

- (i)  $(a_n + b_n)$  tends to infinity
- (ii)  $(\alpha a_n)$  tends to infinity
- (iii)  $(a_n b_n)$  tends to infinity

### Theorem (*Squeeze theorem for sequences tending to $\infty$* )

If  $(b_n)$  tends to infinity and  $a_n \geq b_n$  for  $n = 1, 2, \dots$  then  $(a_n)$  tends to infinity.

**Example.** Consider the sequence  $a_n = \frac{2^n}{n} + 5n^9$ ,  $n = 1, 2, \dots$ . Then let us define  $b_n = \frac{2^n}{n}$  and  $c_n = n^9$  so that:

$$a_n = b_n + 5c_n \quad (6.6.9)$$

Now  $b_n \rightarrow \infty$ , since it is eventually positive for  $n \geq 1$  and:

$$\lim_{n \rightarrow \infty} \frac{1}{b_n} = \lim_{n \rightarrow \infty} n \left(\frac{1}{2}\right)^n = 0 \quad (6.6.10)$$

where we used the standard null sequence (iii) with  $p = 1$ .

Similarly,  $c_n \rightarrow \infty$  since:

$$\lim_{n \rightarrow \infty} \frac{1}{c_n} = \lim_{n \rightarrow \infty} \frac{1}{n^9} = 0 \quad (6.6.11)$$

where we used the standard null sequence (i) with  $p = 9$ . It then follows from the properties of sequences tending to infinity that  $a_n \rightarrow \infty$ .

Alternatively, we could also note that:

$$a_n \geq \frac{2^n}{n} \quad (6.6.12)$$

but  $\frac{2^n}{n} \rightarrow \infty$  so  $a_n \rightarrow \infty$  by the squeeze rule.



**Definition (Subsequence)** The sequence  $(a_{n_k})$  is a subsequence of the sequence  $(a_n)$  if  $(n_k)$  is a strictly increasing sequence of positive integers.

**Example.** Consider the sequence  $a_n = n^{(-1)^n}$  then the odd subsequence is given by the terms:

$$a_{2n+1} = \frac{1}{2n+1} \quad (6.6.13)$$

whereas the even subsequence is given by:

$$a_{2n} = 2n \quad (6.6.14)$$



### Theorem (Subsequence divergence)

For any subsequence  $(a_{n_k})$  of a sequence  $(a_n)$ :

- (i) if  $\lim_{n \rightarrow \infty} a_n = l$  then  $\lim_{n \rightarrow \infty} a_{n_k} = l$ .
- (ii) if  $a_n \rightarrow \infty$  then  $a_{n_k} \rightarrow \infty$

*Proof.*

- (i) Let  $\epsilon > 0$ , then there exists  $N$  such that:

$$|a_n - l| < \epsilon, \forall n > N \quad (6.6.15)$$

Taking  $K$  such that  $n_K \geq N$  then:

$$n_k \geq n_K \geq N, \forall k > K \quad (6.6.16)$$

so that:

$$|a_{n_k} - l| < \epsilon, \forall n_k > N \quad (6.6.17)$$

implying that  $\lim_{n \rightarrow \infty} a_{n_k} = \lim_{n \rightarrow \infty} a_n = l$  as desired.

- (ii) Let  $M > 0$ , then there exists  $N$  such that:

$$a_n > M, \forall n > N \quad (6.6.18)$$

Taking  $K$  such that  $n_K \geq N$  then:

$$n_k \geq n_K \geq N, \forall k > K \quad (6.6.19)$$

so that:

$$a_{n_k} > M, \forall n_k > N \quad (6.6.20)$$

implying that  $a_{n_k} \rightarrow \infty$  as desired. ■

**Proposition (Subsequence rules)**

The sequence  $(a_n)$  is divergent if  $(a_n)$  has two convergent subsequences with different limits.

# D3 Series

## 7.1 Introduction to Series

### Definition (Series)

Let  $(a_n)$  be a sequence. Then the expression:

$$\sum_{i=m}^{\infty} a_i \quad (7.1.1)$$

an infinite series, and define by  $n$ th partial sum:

$$s_n = \sum_{i=m}^n a_i \quad (7.1.2)$$

The series is said to be convergent to  $s$  if its sequence converges to  $s$ , that is if:

$$\lim_{n \rightarrow \infty} \left( \sum_{i=m}^n a_i \right) = s \quad (7.1.3)$$

Otherwise it is said to be divergent.

### Proposition (Geometric series)

The series of the form:

$$\sum_{k=0}^{\infty} ar^n \quad (7.1.4)$$

for some constants  $a, r \in \mathbb{R}$  are called geometric series, and converge to:

$$\sum_{k=0}^{\infty} ar^n = \frac{a}{1-r} \text{ if } |r| < 1 \quad (7.1.5)$$

If  $a \neq 0$  and  $|r| \geq 1$ , then the geometric series diverges.

*Proof.* **Lemma.** For  $r \neq 1$ , the partial sums are given by:

$$s_n = \sum_{k=0}^n ar^k = a \frac{1 - r^{n+1}}{1 - r} \quad (7.1.6)$$

Indeed, note that if  $r \neq 1$  then:

$$(1 - r) \sum_{k=0}^n ar^k = \sum_{k=0}^n ar^k - \sum_{k=0}^n ar^{k+1} \quad (7.1.7)$$

$$= (a + ar + ar^2 + \dots + ar^n) - (ar + ar^2 + \dots + ar^n + ar^{n+1}) \quad (7.1.8)$$

$$= a + ar^{n+1} \quad (7.1.9)$$

$$\Rightarrow \sum_{k=0}^n ar^k = a \frac{1 - r^{n+1}}{1 - r} \quad (7.1.10)$$

as desired.

Hence, for  $|r| < 1$ ,  $\lim_{n \rightarrow \infty} r^{n+1} = 0$  so that  $\lim_{n \rightarrow \infty} s_n = \frac{a}{1-r}$ . ■

### Definition (Cauchy criterion)

A series  $\sum a_n$  satisfies the Cauchy criterion if its partial sums form a Cauchy sequence:

$$\forall \epsilon > 0, \exists N \text{ such that } m, n > N \implies |s_n - s_m| < \epsilon \quad (7.1.11)$$

## 7.2 Telescoping series

Telescoping series are series with terms of the form:

$$a_n = b_n - b_{n+i} \quad (7.2.1)$$

for some natural number  $i$ . The partial sums of  $a_n$  are:

$$s_n = \sum_{k=1}^n (b_k - b_{n+i}) \quad (7.2.2)$$

$$= (b_1 + \dots + b_{i+1} + b_{i+2} + \dots + b_n) - (b_{i+1} + b_{i+2} + \dots + b_{n+i}) \quad (7.2.3)$$

$$= (b_1 + b_2 + \dots + b_i + b_n + b_{n+1} + \dots + b_{n+i}) \quad (7.2.4)$$

Telescoping series most naturally occur with rational series. We provide an example below.

**Example.** Consider the series:

$$\sum_{n=1}^{\infty} \frac{1}{2n(2n+2)} \quad (7.2.5)$$

We can expand the partial fraction as:

$$\frac{1}{2n(2n+2)} = \frac{1}{2} \left( \frac{1}{2n} - \frac{1}{2n+2} \right) \quad (7.2.6)$$

so that the partial sums turn into:

$$s_n = \frac{1}{2} \sum_{k=1}^n \left( \frac{1}{2k} - \frac{1}{2k+2} \right) \quad (7.2.7)$$

$$= \frac{1}{2} \left[ \left( \frac{1}{2} - \frac{1}{4} \right) + \left( \frac{1}{4} - \frac{1}{6} \right) + \dots + \left( \frac{1}{2n-2} - \frac{1}{2n} \right) + \left( \frac{1}{2n} - \frac{1}{2n+2} \right) \right] \quad (7.2.8)$$

$$= \frac{1}{2} \left( \frac{1}{2} - \frac{1}{2n+2} \right) = \frac{n}{4(n+1)} \quad (7.2.9)$$

We quickly see that the series does indeed converge:

$$\lim_{n \rightarrow \infty} s_n = \frac{1}{4} \quad (7.2.10)$$

◀

## 7.3 Manipulating series

### Proposition (Linearity of convergent series)

Suppose that  $\sum_{n=1}^{\infty} a_n = s$  and  $\sum_{n=1}^{\infty} b_n = t$  then:

$$\sum_{n=1}^{\infty} (\lambda a_n + b_n) = \lambda s + t \quad \forall \lambda \in \mathbb{R} \quad (7.3.1)$$

*Proof.* The partial sums associated to the two series are:

$$s_n = \sum_{k=1}^n a_k \quad \text{and} \quad t_n = \sum_{k=1}^n b_k \quad (7.3.2)$$

Then:

$$\sum_{k=1}^{\infty} (\lambda a_n + b_n) = (\lambda a_1 + b_1) + (\lambda a_2 + b_2) + \dots + (\lambda a_n + b_n) \quad (7.3.3)$$

$$= \lambda(a_1 + a_2 + \dots + a_n) + (b_1 + b_2 + \dots + b_n) \quad (7.3.4)$$

$$= \lambda s_n + t_n \quad (7.3.5)$$

By theorem 6.3:

$$\lim_{n \rightarrow \infty} (\lambda s_n + t_n) = \lambda \lim_{n \rightarrow \infty} s_n + \lim_{n \rightarrow \infty} t_n = \lambda s + t \quad (7.3.6)$$

as desired, the series is convergent to:

$$\sum_{n=1}^{\infty} (\lambda a_n + b_n) = \lambda s + t \quad (7.3.7)$$

■

**Example.** Consider the series:

$$\sum_{n=1}^{\infty} \left( \left( -\frac{3}{4} \right)^n - \frac{2}{n(n+1)} \right) = \sum_{n=1}^{\infty} \left( -\frac{3}{4} \right)^n - 2 \sum_{n=1}^{\infty} \frac{1}{n(n+1)} \quad (7.3.8)$$

The first is a convergent geometric series since  $|r| = \frac{3}{4} < 1$ . It converges to:

$$\sum_{n=1}^{\infty} \left( -\frac{3}{4} \right)^n = \frac{-\frac{3}{4}}{1 + \frac{3}{4}} = -\frac{3}{7} \quad (7.3.9)$$

Instead, the second series is a telescoping series whose partial sums can be expanded as:

$$s_n = \sum_{n=1}^n \frac{1}{n(n+1)} = \sum_{n=1}^n \left( \frac{1}{n} - \frac{1}{n+1} \right) \quad (7.3.10)$$

$$= \left( 1 - \frac{1}{2} \right) + \left( \frac{1}{2} - \frac{1}{3} \right) + \dots + \left( \frac{1}{n-1} - \frac{1}{n} \right) + \left( \frac{1}{n} - \frac{1}{n+1} \right) \quad (7.3.11)$$

$$= 1 - \frac{1}{n+1} = \frac{n}{n+1} \quad (7.3.12)$$

so that:

$$\sum_{n=1}^{\infty} \frac{1}{n(n+1)} = \lim_{n \rightarrow \infty} s_n = 1 \quad (7.3.13)$$

Hence the original series converges to:

$$\sum_{n=1}^{\infty} \left( \left( -\frac{3}{4} \right)^n - \frac{2}{n(n+1)} \right) = -\frac{3}{7} - 2 = -\frac{17}{7} \quad (7.3.14)$$

◀

### Theorem (Non-null test)

If  $\sum_{n=1}^{\infty} a_n$  is convergent, then  $(a_n)$  is a null sequence.

If  $(a_n)$  is not a null sequence, then  $\sum_{n=1}^{\infty} a_n$  is divergent.

*Proof.* Note that the second line is true provided the first line is true.

Let  $s_n$  be the  $n$ th partial sum of  $a_n$ . Since  $\sum a_n$  converges,  $s_n$  must also converge to some limit  $s$ . Note that:

$$a_n = s_n - s_{n-1} \implies \lim_{n \rightarrow \infty} a_n = \lim_{n \rightarrow \infty} s_n - \lim_{n \rightarrow \infty} s_{n-1} = s - s = 0 \quad (7.3.15)$$

so that  $a_n$  is indeed a null sequence. ■

**Example.** Consider the series:

$$\sum_{n=1}^{\infty} \frac{(-1)^{n+1} n^2}{2n^2 + 1} \quad (7.3.16)$$

and let us examine the sequence  $(a_n)$  given by:

$$a_n = \frac{(-1)^{n+1} n^2}{2n^2 + 1} \quad (7.3.17)$$

Clearly, we can define a subsequence  $a_{2n}$  as consisting of all even terms:

$$a_{2n} = -\frac{n^2}{2n^2 + 1} \quad (7.3.18)$$

which is convergent to:

$$\lim_{n \rightarrow \infty} a_{2n} = -\lim_{n \rightarrow \infty} \frac{1}{2 + \frac{1}{n^2}} = \frac{1}{2 + \lim_{n \rightarrow \infty} \frac{1}{n^2}} = \frac{1}{2} \quad (7.3.19)$$

Therefore  $a_n$  is not a null sequence since it has a non-null subsequence, and thus the sum  $\sum_{n=1}^{\infty} \frac{(-1)^{n+1} n^2}{2n^2 + 1}$  diverges.



## 7.4 Non-negative series

We continue our study of series by examining those containing only positive terms.

**Example.** Consider the series

$$\sum_{n=1}^{\infty} \frac{1}{n} \quad (7.4.1)$$

known as the harmonic series. Let us write the first few terms as:

$$\sum_{n=1}^{\infty} \frac{1}{n} = 1 + \frac{1}{2} + \frac{1}{3} + \frac{1}{4} + \frac{1}{5} + \frac{1}{6} + \frac{1}{7} + \frac{1}{8} + \dots \quad (7.4.2)$$

$$= 1 + \frac{1}{2} + \left( \frac{1}{3} + \frac{1}{4} \right) + \left( \frac{1}{5} + \frac{1}{6} + \frac{1}{7} + \frac{1}{8} \right) + \dots \quad (7.4.3)$$

Let  $(s_n)$  be the sequence of partial sums of the harmonic series, and consider the subsequence  $(s_{2^n})$

$$s_2 = 1 + \frac{1}{2} \quad (7.4.4)$$

$$s_4 = 1 + \frac{1}{2} + \left( \frac{1}{3} + \frac{1}{4} \right) \geq 1 + \frac{1}{2} \quad (7.4.5)$$

$$s_8 = 1 + \frac{1}{2} + \left( \frac{1}{3} + \frac{1}{4} \right) + \left( \frac{1}{5} + \frac{1}{6} + \frac{1}{7} + \frac{1}{8} \right) \quad (7.4.6)$$

$$s_{2^k} = \sum_{n=1}^{2^k} \frac{1}{n} \geq 1 + \frac{1}{2} + \frac{1}{2} + \dots + \frac{1}{2} = 1 + \frac{1}{2} k \quad (7.4.7)$$

Since  $1 + \frac{1}{2} k \rightarrow \infty$  as  $k \rightarrow \infty$  it follows from the Squeeze rule that  $s_{2^k}$  is a divergent sequence, and therefore non-null. The harmonic series therefore diverges.

**Example.** Consider instead the series:

$$\sum_{n=1}^{\infty} \frac{1}{n^2} \quad (7.4.8)$$

Consider the term:

$$\frac{1}{k^2} < \frac{1}{k(k-1)} = \frac{1}{k-1} - \frac{1}{k} \quad k > 1 \quad (7.4.9)$$

Therefore the  $k$ th partial sum of the series is:

$$s_n = \sum_{k=1}^n \frac{1}{k^2} < 1 + \sum_{k=2}^n \left( \frac{1}{k-1} - \frac{1}{k} \right) = 2 - \frac{1}{n} < 2 \quad (7.4.10)$$

Therefore,  $(s_n)$  is an increasing monotonic series that is bounded above. By the monotone convergence theorem it converges, so that:

$$\sum_{n=1}^{\infty} \frac{1}{n^2} = \lim_{n \rightarrow \infty} s_n \text{ converges} \quad (7.4.11)$$

### Theorem (Comparison test)

If  $0 \leq a_n \leq b_n$  for  $n \in \mathbb{N}$  and  $\sum b_n$  converges, then  $\sum a_n$  converges too.  
If instead  $\sum a_n$  diverges then  $\sum b_n$  diverges too.

**Example.** Consider the series:

$$\sum_{n=1}^{\infty} \frac{1}{n^3} \quad (7.4.12)$$

Then, we may write that:

$$0 \leq n^2 \leq n^3 \quad (7.4.13)$$

so that

$$0 \leq \frac{1}{n^3} \leq \frac{1}{n^2} \quad (7.4.14)$$

Since  $\sum_{n=1}^{\infty} \frac{1}{n^2}$  is convergent, it follows from the Comparison test that  $\sum_{n=1}^{\infty} \frac{1}{n^3}$  also converges.

**Example.** Consider the series:

$$\sum_{n=1}^{\infty} \frac{\cos^2 2n}{n^3} \quad (7.4.15)$$

Since  $0 \leq \cos^2(2n) \leq 1$  we find that:

$$0 \leq \frac{\cos^2 2n}{n^3} \leq \frac{1}{n^3} \quad (7.4.16)$$

Since  $\sum_{n=1}^{\infty} \frac{1}{n^3}$  converges, it follows from the Comparison test that  $\sum_{n=1}^{\infty} \frac{\cos^2 2n}{n^3}$  also converges.

◀

### Theorem (Limit comparison test)

Suppose  $\sum_{n=1}^{\infty} b_n$  converges. Suppose that  $\sum_{n=1}^{\infty} a_n$  and  $\sum_{n=1}^{\infty} b_n$  are positive term series, such that

$$\lim_{n \rightarrow \infty} \frac{a_n}{b_n} = L \neq 0 \quad (7.4.17)$$

If  $\sum_{n=1}^{\infty} b_n$  is convergent, then  $\sum_{n=1}^{\infty} a_n$  is convergent.

If  $\sum_{n=1}^{\infty} b_n$  is divergent, then  $\sum_{n=1}^{\infty} a_n$  is divergent.

*Proof.* Since  $\frac{a_n}{b_n}$  is convergent, it must be bounded:

$$\frac{a_n}{b_n} \leq K \implies a_n \leq Kb_n, \quad n = 1, 2, \dots \quad (7.4.18)$$

Since  $\sum_{n=1}^{\infty} b_n$  converges, it follows that  $\sum_{n=1}^{\infty} Kb_n$  also converges, by the linearity of series. Hence, by the comparison test,  $\sum_{n=1}^{\infty} a_n$  converges.

If instead  $\sum_{n=1}^{\infty} b_n$  diverges, then note that:

$$\lim_{n \rightarrow \infty} \frac{b_n}{a_n} = \frac{1}{L} \neq 0 \quad (7.4.19)$$

since  $L \neq 0$ . Then:

$$\frac{b_n}{a_n} \leq K \implies b_n \leq Ka_n, \quad n = 1, 2, \dots \quad (7.4.20)$$

If  $\sum_{n=1}^{\infty} a_n$  converges, it follows that  $\sum_{n=1}^{\infty} Ka_n$  also converges, by the linearity of series. Hence:

$$\sum_{n=1}^{\infty} a_n \text{ converges} \implies \sum_{n=1}^{\infty} b_n \text{ converges} \quad (7.4.21)$$

which is equivalent to its converse:

$$\sum_{n=1}^{\infty} b_n \text{ diverges} \implies \sum_{n=1}^{\infty} a_n \text{ diverges} \quad (7.4.22)$$

as desired. ■

**Example.** Consider the series:

$$\sum_{n=1}^{\infty} \frac{1}{n^3 + n} = \sum_{n=1}^{\infty} a_n \quad (7.4.23)$$

Then, define:

$$b_n = \frac{1}{n^3} \quad (7.4.24)$$

so that:

$$\lim_{n \rightarrow \infty} \frac{a_n}{b_n} = \lim_{n \rightarrow \infty} \frac{n^3}{n^3 + n} = \lim_{n \rightarrow \infty} \frac{1}{1 + \frac{1}{n^2}} = 1 \neq 0 \quad (7.4.25)$$

Therefore, by the Limit comparison test, it follows that since  $\sum_{n=1}^{\infty} b_n$  converges as was shown earlier,  $\sum_{n=1}^{\infty} \frac{1}{n^3 + n}$  must also converge. ◀

**Example.** Consider the series:

$$\sum_{n=1}^{\infty} \frac{n+4}{2n^3 - n + 1} = \sum_{n=1}^{\infty} a_n \quad (7.4.26)$$

Then we may define

$$b_n = \frac{1}{n^2} \quad (7.4.27)$$

so that:

$$\lim_{n \rightarrow \infty} \frac{a_n}{b_n} = \lim_{n \rightarrow \infty} \frac{n^3 + 4n^2}{2n^3 - n + 1} = \lim_{n \rightarrow \infty} \frac{1 + \frac{4}{n}}{2 - \frac{1}{n^3} + \frac{1}{n^2}} = \frac{1}{2} \neq 0 \quad (7.4.28)$$

However, it was shown that  $\sum_{n=1}^{\infty} \frac{1}{n^2}$  is convergent, so that  $\sum_{n=1}^{\infty} \frac{n+4}{2n^3 - n + 1}$  also converges. ◀

**Theorem (Ratio test)** Suppose  $\sum_{n=1}^{\infty} a_n$  is a series with positive terms. Then:

- (i) if  $\lim_{n \rightarrow \infty} \frac{a_{n+1}}{a_n} = l$  with  $0 \leq l < 1$  then  $\sum_{n=1}^{\infty} a_n$  converges
- (ii) if  $\lim_{n \rightarrow \infty} \frac{a_{n+1}}{a_n} = l$  with  $l > 1$  then  $\sum_{n=1}^{\infty} a_n$  diverges
- (iii) if  $\lim_{n \rightarrow \infty} \frac{a_{n+1}}{a_n} = \infty$  then  $\sum_{n=1}^{\infty} a_n$  diverges.

*Proof.*

- (i) Since  $0 \leq l < 1$  we can choose  $\epsilon > 0$  such that

$$l + \epsilon < 1 \quad (7.4.29)$$

If we let  $r = l + \epsilon$ , then since  $r > l$  there exists  $N$  such that:

$$\frac{a_{n+1}}{a_n} \leq r, \forall n \geq N \quad (7.4.30)$$

Then:

$$\frac{a_n}{a_N} = \prod_{k=n-1}^N \frac{a_{k+1}}{a_k} \leq \prod_{k=n-1}^N r = r^{n-N} \quad (7.4.31)$$

Hence:

$$a_n \leq a_N r^{n-N} \quad (7.4.32)$$

Note however that:

$$\sum_{n=1}^{\infty} a_N r^{n-N} \quad (7.4.33)$$

is a geometric series, and therefore converges. From the comparison test, it follows that  $\sum_{n=1}^{\infty} a_n$  also converges.

(ii) and (iii) Suppose that:

$$\frac{a_{n+1}}{a_n} \rightarrow \infty \text{ or } \frac{a_{n+1}}{a_n} \rightarrow l \quad (7.4.34)$$

with  $l > 1$  then there exists  $N$  such that:

$$\frac{a_{n+1}}{a_n} \geq 1, \forall n \geq N \quad (7.4.35)$$

Therefore:

$$\frac{a_n}{a_N} = \prod_{k=n-1}^N \frac{a_{k+1}}{a_k} \leq 1 \quad (7.4.36)$$

implying that:

$$a_n \geq a_N > 0, \forall n \geq N \quad (7.4.37)$$

$(a_n)$  therefore cannot be a null sequence, and hence  $\sum_{n=1}^{\infty} a_n$  must diverge. ■

**Example.** Consider the series:

$$\sum_{n=1}^{\infty} \frac{(2n)!}{n^n} = \sum_{n=1}^{\infty} a_n \quad (7.4.38)$$

Then:

$$\frac{a_{n+1}}{a_n} = \frac{(2n+2)!}{(n+1)^{n+1}} \cdot \frac{n^n}{(2n)!} \quad (7.4.39)$$

$$= (2n+2)(2n+1) \frac{n^n}{(n+1)^{n+1}} \quad (7.4.40)$$

$$= \frac{(2n+2)(2n+1)}{n+1} \frac{1}{(1+\frac{1}{n})^n} \quad (7.4.41)$$

Therefore:

$$\lim_{n \rightarrow \infty} \frac{a_{n+1}}{a_n} = \frac{2}{e} \lim_{n \rightarrow \infty} (2n+1) = \infty \quad (7.4.42)$$

so  $\sum_{n=1}^{\infty} \frac{(2n)!}{n^n}$  diverges. ◀

**Proposition (Standard series)**

The following series converge:

- (i)  $\sum_{n=1}^{\infty} \frac{1}{n^p}$  for  $p \geq 2$
- (ii)  $\sum_{n=1}^{\infty} c^n$  for  $0 \leq c < 1$
- (iii)  $\sum_{n=1}^{\infty} n^p c^n$  for  $p > 0$  and  $0 \leq c < 1$
- (iv)  $\sum_{n=1}^{\infty} \frac{c^n}{n!}$  for  $c \geq 0$

The following series is divergent:

- (v)  $\sum_{n=1}^{\infty} \frac{1}{n^p}$  for  $0 < p \leq 1$

*Proof.*

- (i) Note that if  $p \geq 2$  then:

$$\frac{1}{n^p} \leq \frac{1}{n^2}, \quad n = 1, 2, \dots \quad (7.4.43)$$

and since  $\sum_{n=1}^{\infty} \frac{1}{n^2}$  converges, by the Comparison test the series  $\sum_{n=1}^{\infty} \frac{1}{n^p}$  must also converge.

- (ii) This is the standard geometric series with common ratio  $r = c$ , which converges provided  $0 \leq c < 1$ .

- (iii) Let  $\sqrt{c} = b$ , then:

$$a_n = (n^p b^n) b^n, \quad n = 1, 2, \dots \quad (7.4.44)$$

Since  $0 \leq b < 1$ ,  $n^p b^n$  is a standard null sequence. Setting  $\epsilon = 1$ , there exists  $N$  such that we have:

$$n^p b^n < 1, \quad \forall n > N \quad (7.4.45)$$

and thus:

$$a_b < b^n, \quad \forall n > N \quad (7.4.46)$$

However,  $\sum_{n=1}^{\infty} b^n$  is a convergent geometric series, so by the Comparison test  $\sum_{n=1}^{\infty} a_n$  converges.

- (iv) If  $c = 0$ , then convergence is trivial. Suppose  $c \neq 0$ , then:

$$\frac{a_{n+1}}{a_n} = \frac{c^{n+1}}{(n+1)!} \cdot \frac{c^n}{n!} = \frac{c}{n+1} \quad (7.4.47)$$

which as  $n \rightarrow \infty$  converges to 0. We then deduce from the ratio test that  $\sum_{n=1}^{\infty} \frac{c^n}{n!}$  converges.

- (v) Note that if  $p \leq 1$  then:

$$\frac{1}{n^p} \geq \frac{1}{n^2}, \quad n = 1, 2, \dots \quad (7.4.48)$$

and since  $\sum_{n=1}^{\infty} \frac{1}{n^2}$  diverges, by the Comparison test the series  $\sum_{n=1}^{\infty} \frac{1}{n^p}$  must also diverge.

■

## 7.5 Series with positive and negative terms

### Definition (*Absolute convergence*)

The series  $\sum_{n=1}^{\infty} a_n$  is **absolutely convergent** if  $\sum_{n=1}^{\infty} |a_n|$  is convergent.

### Theorem (*Absolute convergence test*)

If  $\sum_{n=1}^{\infty} |a_n|$  is absolutely convergent, then  $\sum_{n=1}^{\infty} a_n$  is convergent.

*Proof.* Suppose that  $\sum_{n=1}^{\infty} |a_n|$  converges, and let us define two new series  $\sum_{n=1}^{\infty} b_n$  and  $\sum_{n=1}^{\infty} c_n$  such that:

$$b_n = \begin{cases} a_n, & \text{if } a_n \geq 0 \\ 0, & \text{if } a_n < 0 \end{cases}, \quad \text{and } c_n = \begin{cases} 0, & \text{if } a_n \geq 0 \\ -a_n, & \text{if } a_n < 0 \end{cases} \quad (7.5.1)$$

then both  $b_n$  and  $c_n$  are non-negative. Moreover,  $b_n \leq |a_n|$  and  $c_n \leq |a_n|$  for  $n = 1, 2, \dots$ , so by the Comparison theorem, since  $\sum_{n=1}^{\infty} |a_n|$  converges, it follows that  $\sum_{n=1}^{\infty} b_n$  and  $\sum_{n=1}^{\infty} c_n$  converge. Thus:

$$\sum_{n=1}^{\infty} a_n = \sum_{n=1}^{\infty} (b_n - c_n) = \sum_{n=1}^{\infty} b_n - \sum_{n=1}^{\infty} c_n \quad (7.5.2)$$

also converges. ■

**Example.** Consider the series:

$$\sum_{n=1}^{\infty} \frac{(-1)^{n+1} n}{n^3 + 1} \quad (7.5.3)$$

and its absolute analogue:

$$\sum_{n=1}^{\infty} \frac{n}{n^3 + 1} \quad (7.5.4)$$

Now note that if we define  $a_n = \frac{n}{n^3 + 1}$  and  $b_n = \frac{1}{n^2}$  then:

$$\lim_{n \rightarrow \infty} \frac{a_n}{b_n} = \lim_{n \rightarrow \infty} \frac{n^3}{n^3 + 1} = 1 \neq 0 \quad (7.5.5)$$

By the limit comparison theorem, since  $\sum_{n=1}^{\infty} \frac{1}{n^2}$  converges, it follows that  $\sum_{n=1}^{\infty} \frac{n}{n^3 + 1}$  converges. Hence by the absolute convergence test  $\sum_{n=1}^{\infty} \frac{(-1)^{n+1} n}{n^3 + 1}$  must also converge. ◀

**Proposition (*Triangle inequality*)** If  $\sum_{n=1}^{\infty} a_n$  is absolutely convergent:

$$\left| \sum_{n=1}^{\infty} a_n \right| \leq \sum_{n=1}^{\infty} |a_n| \quad (7.5.6)$$

*Proof.* Suppose that  $\sum_{n=1}^{\infty} |a_n|$  converges, and let us define two new series  $\sum_{n=1}^{\infty} b_n$  and  $\sum_{n=1}^{\infty} c_n$  such that:

$$b_n = \begin{cases} a_n, & \text{if } a_n \geq 0 \\ 0, & \text{if } a_n < 0 \end{cases}, \quad \text{and } c_n = \begin{cases} 0, & \text{if } a_n \geq 0 \\ -a_n, & \text{if } a_n < 0 \end{cases} \quad (7.5.7)$$

then both  $b_n$  and  $c_n$  are non-negative. Then, since:

$$\sum_{n=1}^{\infty} a_n = \sum_{n=1}^{\infty} b_n - \sum_{n=1}^{\infty} c_n \quad (7.5.8)$$

we find that

$$-\sum_{n=1}^{\infty} c_n \leq \sum_{n=1}^{\infty} a_n \leq \sum_{n=1}^{\infty} b_n \quad (7.5.9)$$

Thus, since  $c_n \leq |a_n|$  and  $b_n \leq |a_n|$  we find

$$-\sum_{n=1}^{\infty} |a_n| \leq \sum_{n=1}^{\infty} a_n \leq \sum_{n=1}^{\infty} |a_n| \implies \left| \sum_{n=1}^{\infty} a_n \right| \leq \sum_{n=1}^{\infty} |a_n| \quad (7.5.10)$$

as required. ■

**Example.** Using the absolute convergence test it is immediate that:

$$\sum_{n=1}^{\infty} \frac{(-1)^{n+1}}{2^n} \quad (7.5.11)$$

converges. Note that:

$$\left| \sum_{n=1}^{\infty} \frac{(-1)^{n+1}}{2^n} \right| \leq \sum_{n=1}^{\infty} \frac{1}{2^n} = 1 \quad (7.5.12)$$

so it follows that the value at which the sum converges must lie in the interval  $[-1, 1]$ . ◀

### Theorem (Alternating series test)

Let  $a_n = (-1)^{n+1}b_n$  for  $n = 1, 2, \dots$  then if  $(b_n)$  is a decreasing null sequence with positive terms:

$$\sum_{n=1}^{\infty} a_n \text{ converges} \quad (7.5.13)$$

*Proof.* Let us write the partial sum  $s_{2k}$  of the series as:

$$s_{2k} = (b_1 - b_2) + (b_3 - b_4) + \dots + (b_{2k-1} - b_{2k}) \quad (7.5.14)$$

Since  $(b_n)$  is decreasing, each bracket evaluates to a non-negative value. Consequently  $(s_{2n})$  is an increasing sequence. Moreover:

$$s_{2k} = b_1 - (b_2 - b_3) - (b_4 - b_5) - \dots - (b_{2k-2} - b_{2k-1}) - b_{2k} \leq b_1 \quad (7.5.15)$$

By the monotone convergence theorem it follows that

$$\lim_{n \rightarrow \infty} s_{2n} = s \quad (7.5.16)$$

for some  $s$ . Also:

$$s_{2k-1} = b_1 - (b_2 - b_3) - (b_4 - b_5) - \dots - (b_{2k-2} - b_{2k-1}) = s_{2k} + b_{2k} \quad (7.5.17)$$

so that:

$$\lim_{n \rightarrow \infty} s_{2k+1} = \lim_{n \rightarrow \infty} s_{2k} + \lim_{n \rightarrow \infty} b_{2k} = s \quad (7.5.18)$$

where since  $b_n$  is null, all its subsequences are null. Since both the odd and even subsequences converge to  $s$ , we find that  $s_n \rightarrow s \implies \sum_{n=1}^{\infty} a_n = s$ .  $\blacksquare$

**Example.** Consider the sequence:

$$\sum_{n=1}^{\infty} \frac{(-1)^{n+1}}{n + \sqrt{n}} \quad (7.5.19)$$

We can write its terms as:

$$a_n = \frac{(-1)^{n+1}}{n + \sqrt{n}} = (-1)^{n+1} \frac{1}{n + \sqrt{n}} = (-1)^{n+1} b_n \quad (7.5.20)$$

Now:

- (i)  $b_n = \frac{1}{n + \sqrt{n}} \geq 0$  for  $n = 1, 2, \dots$
- (ii) Since:

$$\frac{1}{n + \sqrt{n}} \leq \frac{1}{n} \quad (7.5.21)$$

and  $\frac{1}{n}$  is a null sequence, by the Squeeze theorem  $(b_n)$  is a null sequence.

- (iii)  $(b_n)$  is decreasing, since  $\left(\frac{1}{b_n}\right) = n + \sqrt{n}$  is increasing. Hence, by the alternating test:

$$\sum_{n=1}^{\infty} \frac{(-1)^{n+1}}{n + \sqrt{n}} \quad (7.5.22)$$

converges.  $\blacktriangleleft$

The sequence  $(a_n)$  is divergent if at least one of its subsequences tends to infinity or minus infinity.

The above proposition follows immediately from writing the converse of the subsequence divergence theorem.

**Example.** Consider the sequence  $a_n = \frac{n}{3} - \lfloor \frac{n}{3} \rfloor$  for  $n = 1, 2, \dots$ . Then the subsequence:

$$a_{3n} = n - \lfloor n \rfloor = 0 \quad (7.5.23)$$

so  $\lim_{n \rightarrow \infty} a_{3n} = 0$ . Instead:

$$a_{3n+1} = n + \frac{1}{3} - n = \frac{1}{3} \quad (7.5.24)$$

so  $\lim_{n \rightarrow \infty} a_{3n+1} = \frac{1}{3}$ . Since  $a_n$  has two subsequences converging to different values, we can conclude that  $a_n$  is divergent.  $\blacktriangleleft$

**Theorem (Convergent subsequence theorem)** Let  $(a_n)$  be made up of two subsequences  $(a_{m_k})$  and  $(a_{n_k})$  which tend to the same limit  $l$ . Then:

$$\lim_{n \rightarrow \infty} a_n = l \quad (7.5.25)$$

*Proof.* Let  $\epsilon > 0$ , then there exists  $K_1$  and  $K_2$  such that:

$$|a_{m_k} - l| < \epsilon, \forall k > K_1 \quad (7.5.26)$$

$$|a_{n_k} - l| < \epsilon, \forall k > K_2 \quad (7.5.27)$$

If we let  $N = \max\{K_1, K_2\}$  then:

$$|a_n - l| < \epsilon, \forall n > N \quad (7.5.28)$$

since for all  $n > N$ ,  $a_n = a_{m_k}$  or  $a_n = a_{n_k}$ , in which case the inequality is satisfied since  $N \geq K_1$  and  $N \geq K_2$ . ■

## 7.6 Monotone convergence theorem

**Theorem (Monotone convergence theorem)**

If the sequence  $(a_n)$  is either:

- (i) increasing and bounded above
- (ii) decreasing and bounded below

then  $(a_n)$  is convergent.

*Proof.* Suppose  $(a_n)$  is bounded above, so that  $\max\{a_n : n = 1, 2, \dots\} = l$ , and let  $\epsilon > 0$ . Then there exists an integer  $N$  such that

$$a_N > l - \epsilon, \forall n > N \quad (7.6.1)$$

since otherwise  $l - \epsilon$  would be an upper bound of  $a_n$ . Since  $a_n$  is increasing,  $a_n \geq a_N$  for  $n > N$  so that:

$$a_n > l - \epsilon \iff l - a_n < \epsilon, \forall n > N \quad (7.6.2)$$

Hence:

$$|a_n - l| = l - a_n < \epsilon, \forall n > N \quad (7.6.3)$$

proving that  $(a_n)$  converges to  $l$ . ■

Moreover, note that if  $(a_n)$  is increasing but not bounded above, then  $a_n \rightarrow \infty$ . Indeed, if it did converge to some value, then for any  $\epsilon > 0$  there exists  $N$  such that:

$$|a_n - l| < \epsilon, \forall n > N \quad (7.6.4)$$

but since  $(a_n)$  is increasing but not bounded above, we can use the triangle inequality to write that:

$$|a_n| < \epsilon + |l|, \forall n > N \quad (7.6.5)$$

proving boundedness for all terms after  $N$ . If we define  $M = \max\{|a_1|, |a_2|, \dots, |a_N|, 1 + |l|\}$  then:

$$|a_n| \leq M \quad (7.6.6)$$

which is a contradiction, since it was assumed  $a_n$  is unbounded.

We may restate this theorem as follows:

**Theorem (*Monotonic sequence theorem*)** A monotonic sequence is either convergent or diverges to  $\infty$ .

# Unit D4: Functions and Continuity

## 8.1 Real functions

### Definition (Real function)

Let  $A \subseteq \mathbb{R}$  and  $f : A \rightarrow \mathbb{R}$  is a bijective function. Then the inverse function  $f^{-1} : \text{Im}(A) \rightarrow A$  and has rule:

$$f^{-1}(f(x)) = x, \forall x \in A \quad (8.1.1)$$

**Example.** Consider the function:

$$f(x) = \frac{1}{1-x}, \forall x \in (-\infty, 1) \quad (8.1.2)$$

Let's solve the equation  $y = f(x)$ .

$$y = \frac{1}{1-x} \iff \frac{1}{y} = 1-x \iff x = 1 - \frac{1}{y} \quad (8.1.3)$$

so every value of  $y \in \text{Im}(f)$  is the image of one  $x \in A$ , showing that  $f$  is bijective, and thus invertible. Its inverse is clearly:

$$f^{-1}(x) = 1 - \frac{1}{x} \quad (8.1.4)$$

To determine the domain, we firstly find the image of  $f$ . Since  $x < -1$ , we must have that:

$$\text{Im}(f) = f(A) = \left\{ \frac{1}{1-x} : \forall x < -1 \right\} = (0, \infty) \quad (8.1.5)$$

so the domain of  $f^{-1}$  must be  $(0, \infty)$ .

Unfortunately, often times it is hard to solve  $y = f(x)$  for  $x$ . In such instances, one can use the following theorem to prove that a function is bijective.

### Theorem (Invertibility of monotonic functions)

If a function  $f$  is strictly increasing or strictly decreasing on some interval  $A$ , then it is invertible on  $A$ .

**Example.** Consider the function:

$$f(x) = x^5 + x - 1, \forall x \in \mathbb{R} \quad (8.1.6)$$

If  $x_1 < x_2$ , then  $x_1^5 < x_2^5$  so:

$$x_1^5 + x_1 < x_2^5 + x_2 \implies x_1^5 + x_1 - 1 < x_2^5 + x_2 \implies f(x_1) < f(x_2) \quad (8.1.7)$$

Thus,  $f(x)$  is strictly increasing, and thus invertible.  $\blacktriangleleft$

## 8.2 Continuity

Consider some real function  $f$ , one important question to ask about this function is whether or not it has any weird gaps, jumps, whether its graph can be drawn without lifting the pen from the paper.

Intuitively, we can define continuity as a property of a function  $f$  such that if  $(x_n)$  is any sequence on its domain that tends to  $a$ , then  $f(x_n)$  tends to  $f(a)$ .

In the case of a jump for example,  $x_n \rightarrow a$ , yet  $f(x_n)$  has no limit, because the subsequence  $f(x_n)$  with  $x_n < a$  and the subsequence  $f(x_n)$  with  $x_n > a$  do not converge to the same value.

### Definition (*Continuity*)

A function  $f : A \rightarrow \mathbb{R}$  is continuous at  $a \in A$  if  $\forall (x_n) \in A$  such that  $x_n \rightarrow a$ , we have  $f(x_n) \rightarrow f(a)$ .

If  $f$  is not continuous at  $a$ , it is discontinuous at  $a$ .

We say that  $f$  is continuous on  $A$  if it is continuous for all points in  $A$ .

**Example.** Consider the function:

$$f(x) = x^3 - 2x^2 \quad (8.2.1)$$

at the point  $a = 2$ . Consider any sequence  $x_n \rightarrow 2$ , then:

$$f(x_n) = x_n^3 - 2x_n^2 \rightarrow 2^3 - 2 \cdot 2^2 = 0 \quad (8.2.2)$$

using the limit combination properties. Moreover,  $f(2) = 0$  thus  $f$  is indeed continuous at  $a = 2$ .  $\blacktriangleleft$

**Example.** Consider the function:

$$f(x) = \lfloor x \rfloor \quad (8.2.3)$$

at the point  $a = 1$ . Consider the sequence  $x_n = 1 - \frac{1}{n}$ , so that  $x_n \rightarrow 1$ . Then:

$$f(x_n) = \left\lfloor 1 - \frac{1}{n} \right\rfloor = 0 \quad (8.2.4)$$

since for any  $n \geq 1$ ,  $1 - \frac{1}{n} < 1$  and thus  $\lfloor 1 - \frac{1}{n} \rfloor = 0$ . However,  $f(1) = 1 \neq 0$ , thus the

function is discontinuous at  $a = 1$ . ◀

**Example.** Consider the function:

$$f(x) = \begin{cases} \sin \frac{1}{x}, & x \neq 0 \\ 0, & x = 0 \end{cases} \quad (8.2.5)$$

Note that:

$$\sin\left(2n + \frac{1}{2}\right)\pi = 1 \quad (8.2.6)$$

so if we define  $x_n = \frac{1}{(2n + \frac{1}{2})\pi}$  then:

$$\sin x_n \rightarrow 1 \neq 0 = f(0) \quad (8.2.7)$$

so we see that  $f$  is discontinuous at  $x = 0$ . ◀

**Example.** Consider the function:

$$f(x) = |x|, \forall x \in \mathbb{R} \quad (8.2.8)$$

Consider any  $a \in \mathbb{R}$  and let  $(x_n)$  be any sequence in  $\mathbb{R}$  such that  $x_n \rightarrow a$  as  $n \rightarrow \infty$ . Now using the triangle inequality:

$$|x_n - a| \geq ||x_n| - |a||, n = 1, 2, \dots \quad (8.2.9)$$

and since  $x_n - a$  is a null sequence, we must have that  $|x_n| - |a|$  is also a null sequence. Thus,  $|x_n| \rightarrow |a| = f(a)$  as desired. Thus  $f$  is continuous everywhere on  $\mathbb{R}$ . ◀

**Example.** Consider the function:

$$f(x) = \sqrt{x}, \forall x \in [0, \infty) \quad (8.2.10)$$

Consider any  $a \in [0, \infty)$  and let  $(x_n)$  be any sequence in  $\mathbb{R}$  such that  $x_n \rightarrow a$  as  $n \rightarrow \infty$ . Now since  $x_n - a$  is a null sequence,  $|x_n - a|$  is also a null sequence, and by the power rule  $\sqrt{|x_n - a|}$  must also be a null sequence.

Also note that in unit D1 we derived:

$$\sqrt{|x_n - a|} \geq |\sqrt{x_n} - \sqrt{a}| \quad (8.2.11)$$

and thus  $\sqrt{x_n} - \sqrt{a}$  is also a null sequence. Thus,  $\sqrt{x_n} \rightarrow \sqrt{a} = f(a)$  as desired. Thus  $f$  is continuous everywhere on  $\mathbb{R}$ . ◀

## 8.3 Properties of continuous functions

### Proposition (*Combination of continuous functions*)

Suppose  $f, g$  are continuous functions at  $a$ , then:

- (i)  $f + g$
- (ii)  $\alpha f$ ,  $\forall \alpha \in \mathbb{R}$
- (iii)  $fg$
- (iv)  $f/g$  is  $g(a) \neq 0$
- (v)  $f \circ g$

*Proof.* We prove (v). Suppose  $f$  is continuous at  $a$  and  $g$  is continuous at  $f(a)$ . If  $f$  has domain  $A$  and  $g$  has domain  $B$ , so that the domain of  $g \circ f$  is:

$$C = \{x \in A : f(x) \in B\} \ni a \quad (8.3.1)$$

We know that  $f(x_n) \rightarrow a$  for all sequences  $(x_n) \in A$ , implying that  $(f(x_n)) \in B$ . Moreover, we know that  $g$  is continuous at  $f(a)$  so that  $g(f(x_n)) \rightarrow g(f(a))$ , as desired. ■

The following theorem results immediately

### Theorem (*Continuity of polynomials and their rationals*)

The following are continuous:

- (i) any polynomial  $p(x) = a_0 + a_1x + \dots + a_nx^n$
- (ii) any rational function  $r(x) = \frac{p(x)}{q(x)}$  where  $p, q$  are polynomials (over  $\mathbb{R} - \{x : q(x) = 0\}$ ).

**Example.** Let's prove that:

$$f(x) = \sqrt{x^2 + 2x + 2} - \frac{3x}{x^4 + 4}, \forall x \in \mathbb{R} \quad (8.3.2)$$

is continuous. To do so, first note that  $x^2 + 2x + 2$  is a polynomial, and thus continuous. Moreover, it is always positive, since it has no real roots. Therefore, if we define  $h(x) = x^2 + 2x + 2$ ,  $h$  is continuous on  $\mathbb{R}$ , and  $g(x) = \sqrt{x}$  is continuous on  $[0, \infty)$ . Thus,  $g(h(x))$  is also continuous on  $\mathbb{R}$ .

Similarly,  $e(x) = x^4 + 4$  is continuous everywhere on  $\mathbb{R}$  since it is a polynomial, and has non-zero values, since  $x^4 = -4$  has no real roots. Therefore, if we define  $d(x) = \frac{3x}{x^4 + 4}$  must also be continuous on  $\mathbb{R}$ , since  $3x$  and  $x^4 + 4$  are (non-zero) polynomials. Hence,  $f(x) = g(h(x)) + d(x)$  is continuous on  $\mathbb{R}$  by the combination rules. ■

### Theorem (*Squeeze rule*)

Let  $f, g, h$  be defined on an open interval  $I$  and let  $a \in I$ . If:

- (i)  $g(x) \leq f(x) \leq h(x)$  for  $x \in I$
- (ii)  $g(a) = f(a) = h(a)$
- (iii)  $g, h$  are continuous at  $a$

then  $f$  is continuous at  $a$ .

*Proof.* Suppose that  $f, g, h$  satisfy these conditions. Since  $x_n \rightarrow a$ , there must exist  $N$  such that:

$$x_n \in (a - \epsilon, a + \epsilon) \subseteq I, \forall n > N \quad (8.3.3)$$

Hence by condition 1:

$$g(x_n) \leq f(x_n) \leq h(x_n) \quad (8.3.4)$$

Condition 2,3 imply that:

$$\lim_{n \rightarrow \infty} g(x_n) = \lim_{n \rightarrow \infty} h(x_n) = f(a) \quad (8.3.5)$$

so by the Squeeze rule of sequences,  $\lim_{n \rightarrow \infty} f(x_n) = f(a)$ , and  $f$  is continuous at  $a$ . ■

**Example.** Consider the function:

$$f(x) = \begin{cases} x^2 \cos \frac{1}{x^2}, & x \neq 0 \\ 0, & x = 0 \end{cases} \quad (8.3.6)$$

Then, we know that:

$$-1 \leq \cos \frac{1}{x^2} \leq 1, x \neq 0 \quad (8.3.7)$$

so that:

$$-x^2 \leq x^2 \cos \frac{1}{x^2} \leq x^2, x \neq 0 \quad (8.3.8)$$

Now, since  $-x^2 = 0 \leq f(0) = 0 \leq 0 = x^2$ , we may assert that:

$$g(x) \leq f(x) \leq h(x) \quad (8.3.9)$$

where  $g(x) = -x^2$  and  $h(x) = x^2$ . Moreover,  $g(0) = h(0) = f(0) = 0$ , and since  $g, h$  are polynomials they are continuous. Thus by the squeeze rule  $f$  must be continuous at  $x = 0$ . ■

### Theorem (Glue rule)

Let  $f$  be defined on an open interval  $I$  and let  $a \in I$ . If  $h, g$  are functions satisfying:

- (i)  $f(x) = g(x)$  for  $x \in I, x < a$ , and  $f(x) = h(x)$  for  $x \in I, x > a$
  - (ii)  $f(a) = g(a) = h(a)$
  - (iii)  $g, h$  are continuous at  $a$
- then  $f$  is continuous at  $a$ .

*Proof.* Suppose  $f, g, h$  satisfy the above conditions. Then:

$$x_n \in (a - \epsilon, a + \epsilon) \subseteq I, \forall n > N \quad (8.3.10)$$

since  $x_n \rightarrow a$ . We define  $(x_n)_N^\infty$ , consists of two subsequences  $(x_{m_k})$  and  $(x_{n_k})$  satisfying:

$$x_{m_k} < a, \text{ and } x_{n_k} \geq a \quad (8.3.11)$$

The conditions give:

$$g(x_{m_k}) \rightarrow g(a) = f(a), \text{ and } h(x_{n_k}) \rightarrow h(a) = g(a) \quad (8.3.12)$$

so that:

$$f(x_{m_k}) \rightarrow f(a), \text{ and } f(x_{n_k}) \rightarrow f(a) \quad (8.3.13)$$

Therefore,  $f(x_n)$  consists of two subsequences convergent to  $f(a)$ , and thus  $f(x_n) \rightarrow f(a)$ , as desired.  $\blacksquare$

**Example.** Consider:

$$f(x) = \begin{cases} x^3 - 3x + 5, & x < 1 \\ \frac{2x+1}{3x-2}, & x \geq 1 \end{cases} \quad (8.3.14)$$

Let's define  $g(x) = x^3 - 3x + 5$  and  $h(x) = \frac{2x+1}{3x-2}$ , then:

$$f(x) = g(x), \text{ for } x < 1 \quad (8.3.15)$$

and

$$f(x) = h(x), \text{ for } x > 1 \quad (8.3.16)$$

Moreover,  $f(1) = 3 = g(1) = h(1)$ , and  $g, h$  are both continuous at  $x = 1$  since the first is a polynomial and the second is the ratio of two polynomials, with non-zero determinant. Hence, by the Glue theorem, we have that  $f$  is continuous.  $\blacktriangleleft$

## 8.4 Trigonometric and exponential functions

We will now prove that the function  $\sin$ ,  $\cos$ ,  $\tan$  and  $\exp$  are continuous.

**Proposition (Sine inequality)**

We have that:

$$\sin x \leq x, \text{ for } 0 \leq x \leq \frac{\pi}{2} \quad (8.4.1)$$

*Proof.* Consider the function  $f(x) = x - \sin x$ . If  $x = 0$  then  $f(x) = 0$ . Moreover, for  $0 \leq x \leq \frac{\pi}{2}$  we have that:

$$f'(x) = 1 - \cos x \quad (8.4.2)$$

and since  $0 \leq \cos x \leq 1$  over this interval, we have that:

$$0 \leq f'(x) \leq 1 \quad (8.4.3)$$

so  $f(x)$  is increasing over  $0 \leq x \leq \frac{\pi}{2}$ , and since  $f(0) = 0$ ,  $f(x) \geq 0$  implying that  $\sin x \leq x$  as desired.  $\blacksquare$

An important consequence of this inequality is the following.

**Corollary.**  $|\sin x| \leq |x|$

*Proof.* The sine inequality proves this result for  $0 \leq x \leq \frac{\pi}{2}$ . For  $x > \frac{\pi}{2}$ :

$$|\sin x| \leq 1 < \frac{\pi}{2} < x = |x| \quad (8.4.4)$$

Finally, if  $x < 0$  then:

$$|\sin(x)| = |\sin(-x)| \leq |-x| = |x| \quad (8.4.5)$$

as desired. ■

**Theorem (Continuity of trigonometric functions)**

The trigonometric functions  $\sin, \cos, \tan$  are continuous.

*Proof.* We need to show that:

$$\sin x_n \rightarrow \sin a \quad (8.4.6)$$

for all null sequences  $x_n - a$ .

We can use the property:

$$\sin x - \sin a = 2 \cos\left(\frac{1}{2}(x + a)\right) 2 \sin\left(\frac{1}{2}(x - a)\right) \quad (8.4.7)$$

So:

$$|\sin x_n - \sin a| = 2 \left| \cos\left(\frac{1}{2}(x_n + a)\right) 2 \sin\left(\frac{1}{2}(x_n - a)\right) \right| \quad (8.4.8)$$

$$\leq 2 \left| \sin\left(\frac{1}{2}(x_n - a)\right) \right| \quad (8.4.9)$$

$$\leq |x_n - a| \quad (8.4.10)$$

Since  $x_n - a$  is a null sequence, we must have that  $\sin x_n \rightarrow \sin a$  as desired. ■

Since  $\cos x = \sin(x + \frac{\pi}{2})$ , we can use the continuity of composite functions to state that it too must be continuous. The same goes for  $\tan x = \frac{\sin x}{\cos x}$ . ■

**Proposition (Exponential inequalities)**

- (i)  $e^x \geq 1 + x$  for  $x \geq 0$
- (ii)  $e^x \leq \frac{1}{1-x}$  for  $0 \leq x < 1$ .

*Proof.* These follow immediately from the power series representation of  $e^x$  (for  $x \geq 0$ ) and  $\frac{1}{1-x}$  (for  $0 \leq x < 1$ ). Indeed:

$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots \leq 1 + x \quad (8.4.11)$$

for  $x \geq 0$ . Similarly:

$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots \leq 1 + x + x^2 + x^3 + \dots = \frac{1}{1-x} \quad (8.4.12)$$

for  $0 \leq x < 1$ . ■

**Corollary.**  $1 + x \leq e^x \leq \frac{1}{1-x}$  for  $|x| < 1$ .

*Proof.* We have already proven the case for  $0 \leq x < 1$ . Consider now the case for  $-1 < x < 0 \implies 0 < -x < 1$ . Then:

$$1 - x \leq e^{-x} \leq \frac{1}{1+x} \quad (8.4.13)$$

Since all terms in the inequality are non-zero on  $0 < -x < 1$ , we can take the reciprocal:

$$1 + x \leq e^x \leq 1 - x \quad (8.4.14)$$

Thus, the inequality has been proven. ■

**Theorem (Continuity of the exponential function)** The exponential function  $\exp$  is continuous.

*Proof.* We need to prove that for all sequences  $x_n \rightarrow a$  we have  $e^{x_n} \rightarrow e^a$ .

If  $x_n - a$  is a null sequence, then there exists  $N$  such that  $|x_n - a| \leq 1$  for  $n > N$  (we use  $\epsilon = 1$ ). Therefore:

$$1 + (x_n + a) \leq e^{x_n - a} \leq \frac{1}{1 - (x_n - a)}, \quad \forall n > N \quad (8.4.15)$$

By the squeeze rule, we find that  $e^{x_n - a} \rightarrow 1$  and thus  $e^{x_n} \rightarrow e^a$  as desired. ■

We summarize the main continuous functions we have found below:

**Proposition (Standard continuous functions)** The following are all continuous:

- (i) polynomials and rational functions of polynomials
- (ii)  $f(x) = |x|$
- (ii)  $f(x) = \sqrt{x}$
- (iv) the trigonometric functions
- (v) the exponential function

# Unit F1: Limits

## 9.1 Introduction to limits of functions

**Definition 9.1 (Punctured neighbourhood)** The **punctured neighbourhood** of a point  $c \in \mathbb{R}$  is a bounded open interval whose midpoint is  $c$  and has been removed. If we define the width of the neighborhood to be  $2\epsilon$  then:

$$N_\epsilon(c) = (c - \epsilon, c) \cup (c, c + \epsilon) \quad (9.1.1)$$

The concept of a punctured neighborhood is essential in defining the limit of a function.

**Definition 9.2 (Limit of a function)**

Let  $f$  be a function defined on  $N_\epsilon(c)$ . Then  $f(x)$  tends to  $l$  as  $x$  tends to  $c$  if  $l \in \mathbb{R}$  and for all sequences  $(x_n)$  in  $N_\epsilon(c)$  such that  $x_n \rightarrow c$ , we have that  $f(x_n) \rightarrow l$ . We write this as:

$$\lim_{x \rightarrow c} f(x) = l \quad (9.1.2)$$

**Example.** Let us try to prove that  $\lim_{x \rightarrow 0} \frac{\sin x}{x} = 1$ .

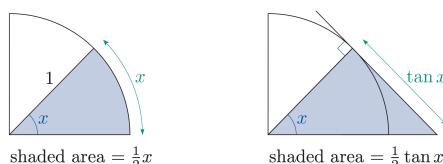
Firstly, in the previous unit we established the inequality:

$$\sin x \leq x, \text{ for } 0 < x \leq \frac{\pi}{2} \quad (9.1.3)$$

We may also deduce that:

$$x \leq \tan x, \text{ for } 0 < x \leq \frac{\pi}{2} \quad (9.1.4)$$

as can be seen from the figures below:



The first inequality (9.1.3) may be rearranged into the more useful form:

$$\frac{\sin x}{x} \leq 1, \text{ for } 0 < x \leq \frac{\pi}{2} \quad (9.1.5)$$

and similarly for the second inequality (9.1.4) may be written as:

$$\cos x \leq \frac{\sin x}{x}, \text{ for } 0 < x \leq \frac{\pi}{2} \quad (9.1.6)$$

Hence, these two may be combined into a single inequality providing upper and lower bounds for  $\frac{\sin x}{x}$ :

$$\cos x \leq \frac{\sin x}{x} \leq 1, \text{ for } 0 < x \leq \frac{\pi}{2} \quad (9.1.7)$$

Now since both  $\cos x$  and  $\frac{\sin x}{x}$  are both even functions, we may substitute  $x' = -x$  into (9.1.7) to find that:

$$\cos x \leq \frac{\sin x}{x} \leq 1, \text{ for } 0 < |x| \leq \frac{\pi}{2} \quad (9.1.8)$$

Suppose  $(x_n)$  is a null sequence in the neighborhood  $N_{\frac{\pi}{2}}(0)$  so that:

$$\cos x_n \leq \frac{\sin x_n}{x_n} \leq 1, n = 1, 2, \dots \quad (9.1.9)$$

Using the squeeze rule, we see that since  $\cos x_n \rightarrow 1$  it must be that  $\frac{\sin x_n}{x_n} \rightarrow 1$ . This proves that:

$$\lim_{x \rightarrow 0} \frac{\sin x}{x} = 1 \quad (9.1.10)$$

◀

**Example.** Consider  $\lim_{x \rightarrow 0} \lfloor x \rfloor$ , we will show that this limit does not exist. Indeed, consider the neighbourhood  $N_{\frac{1}{2}}(1)$  and the sequences  $x_n = 1 - \left(\frac{1}{2}\right)^n$  and  $y_n = 1 + \left(\frac{1}{2}\right)^n$  defined on this neighbourhood. We see that:

$$\lim_{n \rightarrow \infty} f(x_n) = \lim_{n \rightarrow \infty} \left[ 1 - \left(\frac{1}{2}\right)^n \right] = 0 \quad (9.1.11)$$

whereas:

$$\lim_{n \rightarrow \infty} f(y_n) = \lim_{n \rightarrow \infty} \left[ 1 + \left(\frac{1}{2}\right)^n \right] = 1 \quad (9.1.12)$$

Therefore, we have found two different sequences on  $N_{\frac{1}{2}}(1)$  which converge to different values, showing that  $f(x) = \lfloor x \rfloor$  does not tend to a limit as  $x \rightarrow 1$ .

◀

### Theorem 9.3 (Continuity $\iff$ limit)

Let  $f$  be a function defined on an open interval  $I$  and let  $c \in I$ . Then:

$$f \text{ is continuous at } c \iff \lim_{x \rightarrow c} f(x) = f(c) \quad (9.1.13)$$

This is particularly useful when trying to determine the limit of a function that falls within the class of standard continuous functions, such as polynomials, exponential and logarithmic functions etc...

**Proposition 9.4 (Composition rule)**

If  $\lim_{x \rightarrow c} f(x) = l$  and  $\lim_{x \rightarrow l} g(x) = L$ , then  $\lim_{x \rightarrow c} g(f(x)) = L$  provided:

$$g \text{ is defined and continuous at } l \quad (9.1.14)$$

or

$$f(x) \neq l, \forall x \in N_\varepsilon(c) \quad (9.1.15)$$

**Example.** Let us try to evaluate  $\lim_{x \rightarrow 0} \sqrt{\frac{x}{\sin x}}$ . Let us define  $f(x) = \frac{\sin x}{x}$  and  $g(x) = \frac{1}{\sqrt{x}}$ . Then, we have that:

$$\lim_{x \rightarrow 0} f(x) = 1 \quad (9.1.16)$$

as was found previously. Moreover

$$\lim_{x \rightarrow 1} g(x) = 1 \quad (9.1.17)$$

so that, since  $g$  is defined and continuous at 1 (using the composition rules of continuous functions):

$$\lim_{x \rightarrow 0} g(f(x)) = \lim_{x \rightarrow 0} \sqrt{\frac{x}{\sin x}} = 1 \quad (9.1.18)$$



**Example.** We consider the limit:

$$\lim_{x \rightarrow 0} \frac{1 - \cos x}{x} \quad (9.1.19)$$

We use the identity:

$$\cos x = 1 - 2 \sin^2 \frac{x}{2} \quad (9.1.20)$$

to find that:

$$\lim_{x \rightarrow 0} \frac{2 \sin^2 \frac{x}{2}}{x} = \lim_{x \rightarrow 0} \frac{\sin \frac{x}{2}}{\frac{x}{2}} \cdot \lim_{x \rightarrow 0} \sin \frac{x}{2} = 0 \quad (9.1.21)$$

where we used the substitution  $u = \frac{x}{2}$  to evaluate  $\lim_{x \rightarrow 0} \frac{\sin \frac{x}{2}}{\frac{x}{2}}$ .


**Theorem 9.5 (Squeeze rule for limits)**

Let  $f, g, h$  be functions defined on  $N_\varepsilon(c)$  for some  $r > 0$ . If:

- (i)  $g(x) \leq f(x) \leq h(x), \forall x \in N_\varepsilon(c)$
- (ii)  $\lim_{x \rightarrow c} g(x) = \lim_{x \rightarrow c} h(x) = l$  then  $\lim_{x \rightarrow c} f(x) = l$ .

**Example.** We have shown previously that:

$$1 + x \leq e^x \leq \frac{1}{1-x}, \text{ for } |x| < 1 \quad (9.1.22)$$

This may rearranged into:

$$1 + x \leq e^x \leq \frac{1}{1-x}, \text{ for } |x| < 1 \quad (9.1.23)$$

$$1 \leq \frac{e^x - 1}{x} \leq \frac{e^x - 1}{x} \leq \frac{1}{1-x}, \text{ for } 0 < |x| < 1 \quad (9.1.24)$$

Now we have that:

$$1 - \frac{|x|}{1-x} = \frac{x - |x|}{1-x} \leq 1 \quad (9.1.25)$$

since if  $x < 0$  then  $|x| - x = 2|x| > 0$  whereas if  $x > 0$  then  $|x| - x = 0$ . Similarly:

$$\frac{1}{1-x} \leq 1 + \frac{|x|}{1-x} = \frac{1 + (|x| - x)}{1-x} \quad (9.1.26)$$

Consequently, we have the following inequality:

$$1 - \frac{|x|}{1-x} \leq \frac{e^x - 1}{x} \leq 1 + \frac{|x|}{1-x} \quad (9.1.27)$$

Define  $g(x) = 1 - \frac{|x|}{1-x}$ ,  $h(x) = 1 + \frac{|x|}{1-x}$  and  $f(x) = \frac{e^x - 1}{x}$  on  $N_\epsilon(c)$  for some  $\epsilon > 0$ . Then the first condition of the squeeze rule is clearly satisfied:

$$g(x) \leq f(x) \leq h(x) \quad (9.1.28)$$

Next, we also have that:

$$\lim_{x \rightarrow 0} g(x) = 1 = \lim_{x \rightarrow 0} h(x) \quad (9.1.29)$$

Therefore we may conclude that

$$\lim_{x \rightarrow 0} \frac{e^x - 1}{x} = 1 \quad (9.1.30)$$

using the Squeeze rule. ◀

### Definition 9.6 (One-sided limit)

Let  $f(x)$  be a function defined on  $(c, c+r)$  for  $r > 0$ . Then we say that  $f(x)$  tends to  $l$  as  $x$  tends to  $c$  from the right:

$$\lim_{x \rightarrow c^+} f(x) = l \quad (9.1.31)$$

provided that for each sequence  $(x_n)$  in  $(c, c+r)$  such that  $x_n \rightarrow c$ ,  $f(x_n) \rightarrow l$ .

Let  $f(x)$  be a function defined on  $(c-r, c)$  for  $r > 0$ . Then we say that  $f(x)$  tends to  $l$  as  $x$  tends to  $c$  from the left:

$$\lim_{x \rightarrow c^-} f(x) = l \quad (9.1.32)$$

provided that for each sequence  $(x_n)$  in  $(c-r, c)$  such that  $x_n \rightarrow c$ ,  $f(x_n) \rightarrow l$ .

The following result follows immediately from the fact that  $(c-\epsilon, c) \cup (c, c+\epsilon) = N_\epsilon(c)$ .

### Theorem 9.6 (Ordinary limits and one sided-limits)

Let  $f$  be defined on  $N_\epsilon(c)$  for  $r > 0$ . Then:

$$\lim_{x \rightarrow c} f(x) = l \iff \lim_{x \rightarrow c^+} f(x) = \lim_{x \rightarrow c^-} f(x) = l \quad (9.1.33)$$

Finally, we also have an analogue of Theorem 9.3 for one-sided limits:

**Proposition 9.7 (Continuity  $\iff$  one-sided limit)**

Let  $f$  be a function whose domain  $I$  is an interval with left-hand endpoint  $c$  included (so either  $[c, \infty)$ ,  $[c, b)$  or  $[c, b]$  where  $b > c$ ). Then:

$$f \text{ is continuous at } c \iff \lim_{x \rightarrow c^\pm} f(x) = f(c) \quad (9.1.34)$$

**Example.** Let us evaluate  $\lim_{x \rightarrow 0^+} \left( \frac{\sin x}{x} + \sqrt{x} \right)$ .

We have from typical combination rules that:

$$\lim_{x \rightarrow 0^+} \left( \frac{\sin x}{x} + \sqrt{x} \right) = \lim_{x \rightarrow 0^+} \frac{\sin x}{x} + \lim_{x \rightarrow 0^+} \sqrt{x} = 1 + 0 = 0 \quad (9.1.35)$$

where we used Theorem 9.6 to get:

$$\lim_{x \rightarrow 0} \frac{\sin x}{x} = 1 \implies \lim_{x \rightarrow 0^+} \frac{\sin x}{x} = 1 \quad (9.1.36)$$

and Proposition 9.7 to get

$$\lim_{x \rightarrow 0} \sqrt{x} = 1 \implies \lim_{x \rightarrow 0^+} \sqrt{x} = 1 \quad (9.1.37)$$

◀

## 9.2 Asymptotic behaviour

**Definition 9.8 (Infinite limit)**

Let  $f$  be defined on  $N_\epsilon(C)$ . Then  $f(x)$  tends to  $\infty$  as  $x$  tends to  $c$  if for each sequence  $(x_n)$  in  $N_\epsilon(c)$  such that  $x_n \rightarrow c$ ,  $f(x_n) \rightarrow \infty$ . We write that:

$$f(x) \rightarrow \infty \text{ as } x \rightarrow c \quad (9.2.1)$$

**Theorem 9.9 (Reciprocal rule for limits)**

If  $f$  is a function satisfying:

(i)  $f(x) > 0$  for  $x \in N_\epsilon(c)$ , where  $\epsilon > 0$

(ii)  $f(x) \rightarrow 0$  as  $x \rightarrow c$

then  $\frac{1}{f(x)} \rightarrow \infty$  as  $x \rightarrow c$ .

**Example.** Consider the asymptotic behaviour of  $\frac{1}{x^3 - 1}$  as  $x \rightarrow 1^+$ . Define  $f(x) = x^3 - 1$ , then we have that:

$$f(x) = x^3 - 1 > 0 \quad \forall x \in (1, \infty) = (1, 1 + r) \quad (9.2.2)$$

for some  $r > 0$ . Moreover

$$\lim_{x \rightarrow 1^+} f(x) = 0 \quad (9.2.3)$$

Therefore:

$$\frac{1}{x^3 - 1} \rightarrow \infty \text{ as } x \rightarrow 1^+ \quad (9.2.4)$$



### Theorem 9.10 (Squeeze rule for $x \rightarrow \infty$ )

Let  $f, g, h$  be defined on  $(R, \infty)$ . Then:

(a) if

- (i)  $g(x) \leq f(x) \leq h(x), \forall x \in (R, \infty)$
- (ii)  $\lim_{x \rightarrow \infty} g(x) = \lim_{x \rightarrow \infty} h(x) = l$   
then  $\lim_{x \rightarrow \infty} f(x) = l$ .

(b) if

- (i)  $g(x) \leq f(x), \forall x \in (R, \infty)$
- (ii)  $g(x) \rightarrow \infty$  as  $x \rightarrow \infty$   
then  $f(x) \rightarrow \infty$  as  $x \rightarrow \infty$ .

**Example.** Let's examine the asymptotic behaviour as  $x \rightarrow \infty$  of:

$$f(x) = \frac{\sin(1/x)}{x} \quad (9.2.5)$$

We have that for  $x \neq 0$ :

$$-1 \leq \sin(1/x) \leq 1 \quad (9.2.6)$$

so that:

$$-\frac{1}{x} \leq \frac{\sin(1/x)}{x} \leq \frac{1}{x}, \quad \forall x \neq 0 \quad (9.2.7)$$

Then, since:

$$\lim_{x \rightarrow \infty} -\frac{1}{x} = \lim_{x \rightarrow \infty} \frac{1}{x} = 0 \quad (9.2.8)$$

we must have that

$$\lim_{x \rightarrow \infty} \frac{\sin(1/x)}{x} = 0 \quad (9.2.9)$$



### Proposition 9.11 (Asymptotic behaviour of standard functions)

(a) let  $a_0, a_1, \dots, a_{n-1} \in \mathbb{R}$ , and let

$$p(x) = x^n + a_{n-1}x^{n-1} + \dots + a_1x + a_0 \quad (9.2.10)$$

Then

$$p(x) \rightarrow \infty, \text{ and } \frac{1}{p(x)} \rightarrow 0 \text{ as } x \rightarrow \infty \quad (9.2.11)$$

(b) For  $n \in \mathbb{N}$ , then:

$$\frac{e^x}{x^n} \rightarrow \infty, \text{ and } \frac{x^n}{e^x} \rightarrow 0 \text{ as } x \rightarrow \infty \quad (9.2.12)$$

(c) we have  $\log x \rightarrow \infty$  as  $x \rightarrow \infty$ , and for  $a > 0$  we have:

$$\lim_{x \rightarrow \infty} \frac{\log x}{x^a} = 0 \quad (9.2.13)$$

*Proof.* (a) Firstly, we have that the zeros of the polynomial  $p$  must lie in some interval  $(-M, M)$ , so we have that:

$$p(x) > 0, \forall x \in (M, \infty) \quad (9.2.14)$$

Now if  $x \neq 0$  then:

$$p(x) = x^n \left( 1 + \frac{a_{n-1}}{x} + \dots + \frac{a_0}{x^n} \right) \quad (9.2.15)$$

Now we have that:

$$1 + \frac{a_{n-1}}{x} + \dots + \frac{a_0}{x^n} \rightarrow 1 + 0 + \dots + 0 = 1 \text{ as } x \rightarrow \infty \quad (9.2.16)$$

Therefore

$$\frac{1}{p(x)} = \frac{\frac{1}{x^n}}{1 + \frac{a_{n-1}}{x} + \dots + \frac{a_0}{x^n}} \rightarrow 0 \text{ as } x \rightarrow \infty \quad (9.2.17)$$

It follows from the reciprocal rule that, since  $p(x) \rightarrow \infty$  as  $x \rightarrow \infty$ .

(b) We use the series expansion:

$$e^x = \sum_{i=0}^{\infty} \frac{x^i}{i!} \geq \frac{x^{n+1}}{(n+1)!}, \forall x \geq 0 \quad (9.2.18)$$

Hence, if  $x > 0$  then:

$$\frac{e^x}{x^n} \geq \frac{x}{(n+1)!}, \text{ and } 0 \leq \frac{x^n}{e^x} \leq \frac{(n+1)!}{x} \quad (9.2.19)$$

Now  $\frac{x}{(n+1)!} \rightarrow \infty$  so that:

$$\frac{e^x}{x^n} \rightarrow \infty, \text{ as } x \rightarrow \infty \quad (9.2.20)$$

Similarly, using the Squeeze rule

$$\frac{x^n}{e^x} \rightarrow 0 \text{ as } x \rightarrow \infty \quad (9.2.21)$$

(c) We know that  $\log x$  is the strictly increasing inverse of the exponential, so that  $\log x \rightarrow \infty$ .

Now let  $a > 0$ , if we define  $t(x) = a \log x$  then  $x^a = e^{a \log x} = e^t$ , giving:

$$\frac{\log x}{x^a} = \frac{t}{ae^t} \quad (9.2.22)$$

We have shown that  $t \rightarrow \infty$  since  $a$  is positive. Hence, using part (b)

$$\frac{t}{ae^t} \rightarrow 0 \quad (9.2.23)$$

Using the composition rule of limits:

$$\lim_{x \rightarrow \infty} \frac{\log x}{x^a} = 0 \quad (9.2.24)$$

as desired. ■

**Example.** Let us examine the behaviour of

$$f(x) = \frac{2e^x - x^2}{e^x + \log x} \quad (9.2.25)$$

Then:

$$f(x) = \frac{\frac{2e^x}{x^2} - 1}{\frac{e^x}{x^2} + \frac{\log x}{x^2}} \quad (9.2.26)$$

Using the combination rules of limits:

$$\lim_{x \rightarrow \infty} f(x) = \frac{\lim_{x \rightarrow \infty} \frac{2e^x}{x^2} - 1}{\lim_{x \rightarrow \infty} \left( \frac{e^x}{x^2} + \frac{\log x}{x^2} \right)} \quad (9.2.27)$$

$$= \frac{\lim_{x \rightarrow \infty} \frac{2e^x}{x^2} - 1}{\lim_{x \rightarrow \infty} \frac{e^x}{x^2}} \quad (9.2.28)$$

$$= \lim_{x \rightarrow \infty} \left( 2 - \frac{x^2}{e^x} \right) \quad (9.2.29)$$

$$= 2 \quad (9.2.30)$$



**Example.** Let's prove that  $\lim_{x \rightarrow \infty} x \sin \frac{1}{x} = 1$ . We use the substitution  $u(x) = \frac{1}{x}$  so that:

$$\lim_{x \rightarrow \infty} x \sin \frac{1}{x} = \lim_{u \rightarrow 0} \frac{\sin u}{u} = 1 \quad (9.2.31)$$

as desired. ◀

## 9.3 Continuity of functions

**Definition 9.12 (Continuity of functions)**

Let  $f$  have a domain  $A$  and let  $c \in A$ . Then  $f$  is **continuous** at  $c$  if  $\forall \epsilon > 0$ , there exists  $\delta > 0$  such that:

$$|f(x) - f(c)| < \epsilon, \forall x \in A \text{ with } |x - c| < \delta \quad (9.3.1)$$

Much like in the definition of limits for series, we can view definition 9.12 as an  $\epsilon - \delta$  game. One player chooses a small positive  $\epsilon$ , and challenges the other player to find  $\delta$  suitable small such that

$$|f(x) - f(c)| < \epsilon, \forall x \in A \text{ with } |x - c| < \delta \quad (9.3.2)$$

is satisfied.

The general strategy to prove the continuity of polynomial functions  $f$  with domain  $A$  at some point  $c \in A$  is:

- (i) express  $f(x) - f(c) = (x - c)g(x)$
- (ii) obtain a bound  $|g(x)| \leq M$  for  $|x - c| \leq r$  where  $r > 0$  such that  $[c - r, c + r] \subset A$ .
- (iii) use  $|f(x) - f(c)| \leq M|x - c|$  for  $|x - c| \leq r$  and set  $\epsilon = M|x - c|$  to choose  $\delta$  such that:

$$|f(x) - f(c)| < \epsilon, \forall x \in A \text{ with } |x - c| < \delta \quad (9.3.3)$$

**Example.** Let us prove that  $f(x) = x^3$  is continuous at  $c = 1$ .

Firstly we note that the domain of  $f$  is  $\mathbb{R}$ . Suppose we are given  $\epsilon > 0$ , our goal is to choose  $\delta > 0$  such that:

$$|x^3 - 1| < \epsilon, \forall x \text{ with } |x - 1| < \delta \quad (9.3.4)$$

We can write that:

$$x^3 - 1 = (x - 1)(x^2 + x + 1) \quad (9.3.5)$$

Now we find an upper bound for  $|x^2 + x + 1|$ . When  $|x - 1| \leq 1$  then  $x \in [0, 2]$  so that:

$$|x^2 + x + 1| \leq |x^2| + |x + 1| \leq 4 + 3 = 7, \forall |x - 1| \leq 1 \quad (9.3.6)$$

Therefore:

$$|f(x) - f(1)| \leq 7|x - 1|, \forall |x - 1| \leq 1 \quad (9.3.7)$$

Therefore, if  $|x - 1| < \delta$ , then  $|f(x) - f(1)| \leq 7\delta|x - 1|$ . Now we need  $5\delta < \epsilon$  so that  $\delta \leq \frac{1}{7}\epsilon$ . Hence, we must have that if  $\delta = \min\{1, \frac{\epsilon}{7}\}$  then:

$$|f(x) - f(1)| < \epsilon, \forall x \text{ with } |x - 1| < \delta \quad (9.3.8)$$

as desired. ◀

### Theorem 9.13 (Equivalence of continuity definitions)

The  $\epsilon - \delta$  definition and the sequence definition of continuity are equivalent.

*Proof.* Let  $f$  have domain  $A$  with  $c \in A$ .

Assume that continuity according to  $\epsilon - \delta$  is satisfied, so that for  $\epsilon > 0$ ,  $\exists \delta > 0$  such that:

$$|f(x) - f(c)| < \epsilon, \text{ for } |x - c| < \delta \quad (9.3.9)$$

Now consider a sequence  $x_n \in A$  such that  $x_n \rightarrow c$ . Then, there exists  $N$  such that:

$$|x_n - c| < \delta, \forall n > N \quad (9.3.10)$$

so that:

$$|f(x_n) - f(c)| < \epsilon, \forall n > N \quad (9.3.11)$$

Consequently,  $f(x_n) \rightarrow f(c)$ , which the sequence definition of continuity.

Now suppose that  $f$  is continuous according to the sequence definition, we argue by contradiction that for some  $\epsilon > 0$ , there is no  $\delta > 0$  such that

$$|f(x) - f(c)| < \epsilon, \text{ for } |x - c| < \delta \quad (9.3.12)$$

Hence, for all  $n$ ,  $\exists x_n \in A$  with  $|x_n - c| < \frac{1}{n}$  such that

$$|f(x_n) - f(c)| \geq \epsilon \quad (9.3.13)$$

By the sequential definition  $\lim_{n \rightarrow \infty} f(x_n) = f(c)$  contradicting the above inequality. Hence the  $\epsilon - \delta$  definition must also be satisfied. ■

## 9.4 Unusual function continuity

### Proposition 9.14 (Dirichlet function)

The Dirichlet function defined as:

$$f(x) = \begin{cases} 1, & \text{if } x \text{ is rational} \\ 0, & \text{if } x \text{ is irrational} \end{cases} \quad (9.4.1)$$

is discontinuous everywhere on  $\mathbb{R}$ .

*Proof.* Let  $c \in \mathbb{R}$ . By the density of  $\mathbb{R}$ , each interval  $(c - \frac{1}{n}, c + \frac{1}{n})$  with  $n$  natural contains a rational  $x_n$  and irrational  $y_n$ . Then  $x_n \rightarrow c$  and  $y_n \rightarrow c$ , yet  $f(x_n) = 1$  and  $f(y_n) = 0$  so  $f$  is discontinuous at  $c$ . ■

### Proposition 9.15 (Blancmange function)

The sawtooth function defined as:

$$s(x) = \begin{cases} x - \lfloor x \rfloor, & \text{if } 0 \leq x - \lceil x \rceil \leq \frac{1}{2} \\ 1 - (x - \lfloor x \rfloor), & \text{if } \frac{1}{2} < x - \lfloor x \rfloor < 1 \end{cases} \quad (9.4.2)$$

and the Blancmange function  $B$  is defined as:

$$B(x) = \sum_{n=0}^{\infty} \frac{1}{2^n} s(2^n x) \quad (9.4.3)$$

is continuous everywhere on  $\mathbb{R}$ .

*Proof.* Let  $c \in \mathbb{R}$ , and let  $\epsilon > 0$ . Then:

$$B(x) - B(c) = \sum_{n=0}^{\infty} \frac{1}{2^n} (s(2^n x) - s(2^n c)) \quad (9.4.4)$$

Using the triangle inequality:

$$|B(x) - B(c)| = \left| \sum_{n=0}^{\infty} \frac{1}{2^n} (s(2^n x) - s(2^n c)) \right| \leq \sum_{n=0}^{\infty} \frac{1}{2^n} |s(2^n x) - s(2^n c)| \quad (9.4.5)$$

Now since  $s(2^n x) \in \left[0, \frac{1}{2}\right]$  and  $s(2^n c) \in \left[0, \frac{1}{2}\right]$  for all  $x, c$  and natural  $n$ , we may write:

$$|s(2^n x) - s(2^n c)| \leq \frac{1}{2} \implies \sum_{n=N}^{\infty} \frac{1}{2^n} |s(2^n x) - s(2^n c)| \leq \frac{1}{2} \sum_{n=N}^{\infty} \frac{1}{2^n} \quad (9.4.6)$$

and using our standard results for geometric series:

$$\sum_{n=N}^{\infty} \frac{1}{2^n} |s(2^n x) - s(2^n c)| \leq \frac{1}{2^N} \quad (9.4.7)$$

Now consider:

$$x \mapsto s(2^n x), n = 0, 1, 2, \dots \quad (9.4.8)$$

which is a continuous function. Consequently, for all  $n$  there is a positive  $\delta_n$  such that:

$$|s(2^n x) - s(2^n c)| < \frac{\epsilon}{4}, \forall |x - c| < \delta_n \quad (9.4.9)$$

Choosing  $\delta = \min_{n \in \mathbb{N}} \delta_n$  we get that for  $|x - c| < \delta$

$$\sum_{n=0}^{N-1} \frac{1}{2^n} |s(2^n x) - s(2^n c)| \leq \sum_{n=0}^{N-1} \frac{1}{2^N} \frac{\epsilon}{4} < 2 \cdot \frac{\epsilon}{4} = \frac{\epsilon}{2} \quad (9.4.10)$$

Consequently

$$\sum_{n=0}^{\infty} \frac{1}{2^n} |s(2^n x) - s(2^n c)| \leq \frac{\epsilon}{2} + \frac{1}{2^N} \quad (9.4.11)$$

so we need to choose  $N$  such that  $\frac{1}{2^N} < \frac{1}{2}\epsilon$  for the condition of continuity to be satisfied. We can always do so because  $\frac{1}{2^n}$  is a basic null sequence.

The blancmange function is therefore continuous ■

### Definition 9.16 ( $\epsilon - \delta$ definition of limit)

Let  $f$  be a function defined on  $N_\epsilon(c)$  of  $c$ . Then  $f(x)$  tends to  $l$  as  $x$  tends to  $c$  if  $\forall \epsilon > 0, \exists \delta > 0$  such that

$$|f(x) - l| < \epsilon, \forall x \text{ such that } 0 < |x - c| < \delta \quad (9.4.12)$$

We then write that:

$$\lim_{x \rightarrow c} f(x) = l \quad (9.4.13)$$

**Example.** We evaluate

$$\lim_{x \rightarrow 1} \frac{2x^3 + 3x - 5}{x - 1} \quad (9.4.14)$$

Since  $2x^3 + 3x - 5 = (x - 1)(2x^2 + 2x + 5)$  we guess that the limit tends to 9. Indeed, we need to show that for each  $\epsilon > 0$ , there exists  $\delta > 0$  such that:

$$|f(x) - 9| < \epsilon, \forall x \text{ with } 0 < |x - 1| < \delta \quad (9.4.15)$$

But we have that for  $0 < |x - 1| < 1$  then  $x \in (-1, 2)$  so that:

$$|2x^2 + 2x + 5| < 17, \forall 0 < |x - 1| < 1 \quad (9.4.16)$$

so that

$$|f(x) - 9| < 17|x - 1|, \forall 0 < |x - 1| < 1 \quad (9.4.17)$$

Consequently if  $|x - 1| < \delta$  then:

$$|f(x) - 9| < 17\delta, \forall 0 < |x - 1| < \delta \quad (9.4.18)$$

To have that  $|f(x) - 9| < \epsilon$ , we need  $\delta \leq \frac{\epsilon}{17}$ . So given  $\epsilon > 0$  we need to find  $\delta \leq \frac{\epsilon}{17}$ , which can always be done. Consequently  $f$  is continuous.  $\blacktriangleleft$

## 9.5 Uniform continuity

**Definition 9.17 (Uniform continuity)**

A function  $f$  defined on the interval  $I$  is **uniformly continuous** on  $I$  if for all  $\epsilon > 0$  there exists  $\delta > 0$  such that:

$$|f(x) - f(y)| < \epsilon, \forall x, y \in I \text{ such that } |x - y| < \delta \quad (9.5.1)$$

**Example.** Let us prove that  $f(x) = x^2$  is uniformly continuous on  $I = [-4, 4]$ .

Let  $\epsilon > 0$ , we have:

$$f(x) - f(y) = x^2 - y^2 = (x + y)(x - y) \quad (9.5.2)$$

Therefore  $x, y \in [-4, 4]$  implies that  $|x| \leq 4$  and  $|y| \leq 4$  so that:

$$|f(x) - f(y)| = |x + y||x - y| \quad (9.5.3)$$

$$\leq (|x| + |y|)|x - y| \quad (9.5.4)$$

$$\leq 8|x - y| \quad (9.5.5)$$

Thus if we choose  $\delta = \frac{\epsilon}{8}$  and  $|x - y| < \delta$  then:

$$|f(x) - f(y)| < \epsilon \quad (9.5.6)$$

as desired.  $\blacktriangleleft$

**Theorem 9.17 (Sequential definition of uniform discontinuity)**

Let  $f$  be defined on  $I$ , then  $f$  is not uniformly continuous on  $I$  iff  $\exists(x_n), (y_n) \in I$  and  $\exists\epsilon > 0$  such that:

- (i)  $|x_n - y_n| \rightarrow 0$  as  $n \rightarrow \infty$
- (ii)  $|f(x_n) - f(y_n)| \geq \epsilon$ , for  $n = 1, 2, \dots$

*Proof.* Suppose  $f$  is not uniformly continuous on  $I$ . Then  $\exists\epsilon > 0$  such that for all  $\delta > 0$  there are  $x, y \in I$  such that:

$$|x - y| < \delta \text{ and } |f(x) - f(y)| \geq \epsilon \quad (9.5.7)$$

Setting  $\delta = 1, \frac{1}{2}, \frac{1}{3}, \dots$  we obtain

$$|x_n - y_n| < \frac{1}{n} \text{ and } |f(x_n) - f(y_n)| \geq \epsilon, \quad n = 1, 2, \dots \quad (9.5.8)$$

Therefore  $|x_n - y_n| \rightarrow 0$  and  $|f(x_n) - f(y_n)| \geq \epsilon$  as desired.

Now suppose that  $|x_n - y_n| \rightarrow 0$  and  $|f(x_n) - f(y_n)| \geq \epsilon$  are satisfied. Furthermore, suppose that  $f$  is uniformly continuous such that there exists  $\delta > 0$  satisfying:

$$|f(x) - f(y)| < \epsilon, \quad \forall x, y \in I \text{ with } |x - y| < \delta \quad (9.5.9)$$

But by statement 1,  $|x_n - y_n| < \delta$  for  $n > N$ , so

$$|f(x_n) - f(y_n)| < \epsilon, \quad \forall n > N \quad (9.5.10)$$

contradicting the second statement. Therefore  $f$  is not uniformly continuous on  $I$ . ■

**Example.** We show that  $f(x) = x^2$  is not uniformly continuous on  $\mathbb{R}$ .

Indeed, taking  $x_n = n + \frac{1}{n}$  and  $y_n = n$  then we see that:

$$|x_n - y_n| = \left| \frac{1}{n} \right| = \frac{1}{n} \rightarrow 0 \quad (9.5.11)$$

Moreover:

$$|f(x_n) - f(y_n)| = \left| \left( n + \frac{1}{n} \right)^2 - n^2 \right| \quad (9.5.12)$$

$$= 2 + \frac{1}{n^2} > \epsilon \quad (9.5.13)$$

where  $\epsilon = 2$ . ◀

# Unit F2: Differentiation

**Definition 10.1 (Differentiability)** Let  $f$  be defined on the open interval  $I$ , and let  $c \in I$ . Then the **derivative of  $f$  at  $c$**  is defined as

$$f'(c) = \lim_{x \rightarrow c} \frac{f(x) - f(c)}{x - c} = \lim_{x \rightarrow c} Q(x) \quad (10.0.1)$$

where  $Q(x)$  is the difference quotient. If this limit exists, then  $f$  is differentiable at  $c$ . If  $f$  is differentiable  $\forall c \in I$  then we say that  $f$  is **differentiable**.

**Example.** Let us prove that  $f(x) = \frac{1}{x}$  is differentiable on  $I = \mathbb{R}^*$ . Indeed the difference quotient reads:

$$Q(x) = \frac{f(x) - f(c)}{x - c} = \frac{\frac{1}{x} - \frac{1}{c}}{x - c} = -\frac{1}{cx} \quad (10.0.2)$$

so that:

$$\lim_{x \rightarrow c} Q(x) = \lim_{x \rightarrow c} \frac{1}{cx} = \frac{1}{c^2} \quad (10.0.3)$$

where  $c \in \mathbb{R}^*$ . We may therefore conclude that:

$$f'(x) = -\frac{1}{x^2} \quad (10.0.4)$$

◀

**Example.** Let us prove that the function:

$$f(x) = \begin{cases} x^2 \cos \frac{1}{x}, & x \neq 0 \\ 0, & x = 0 \end{cases} \quad (10.0.5)$$

is differentiable at  $x = 0$ . Indeed:

$$Q(x) = \frac{f(x) - f(0)}{x} = x \cos \frac{1}{x} \quad (10.0.6)$$

so we must prove that the following limit exists:

$$\lim_{x \rightarrow 0} x \cos \frac{1}{x} \quad (10.0.7)$$

Let us define

$$f_{\text{new}}(x) = \begin{cases} x \cos \frac{1}{x}, & x \neq 0 \\ 0, & x = 0 \end{cases} \quad (10.0.8)$$

To do so, we know that:

$$-x \leq x \cos \frac{1}{x} \leq x \quad (10.0.9)$$

Now by the squeeze rule, we see that since  $h(x) = x$  and  $g(x) = -x$  are continuous at 0 and  $h(0) = g(0) = f_{\text{new}}(0) = 0$  then

$$\lim_{x \rightarrow 0} x \cos \frac{1}{x} = 0 \quad (10.0.10)$$

so that  $f$  is indeed differentiable at  $x = 0$  with  $f'(0) = 0$ .  $\blacktriangleleft$

**Example.** Let us examine the differentiability of

$$f(x) = \begin{cases} \sqrt{|x|} \sin \frac{1}{x}, & x \neq 0 \\ 0, & x = 0 \end{cases} \quad (10.0.11)$$

The difference quotient is:

$$Q(x) = \frac{\sqrt{|x|} \sin \frac{1}{x}}{x} \quad (10.0.12)$$

Now consider the sequence  $h_n = \frac{1}{n\pi}$  so that

$$Q(h_n) = n\pi \sqrt{\frac{1}{n\pi}} \sin n\pi \quad (10.0.13)$$

which does not exist. Hence we must have that  $\lim_{x \rightarrow 0} Q(x)$  does not exist.  $\blacktriangleleft$

### Proposition 10.2 (Standard derivatives)

- (i)  $f(x) = k$  with  $k \in \mathbb{R}$  then  $f'(x) = 0$
- (ii)  $f(x) = x^n$  with  $n \in \mathbb{N}$  then  $f'(x) = nx^{n-1}$
- (iii)  $f(x) = \sin x$  then  $f'(x) = \cos x$
- (iv)  $f(x) = \cos x$  then  $f'(x) = -\sin x$
- (v)  $f(x) = e^x$  then  $f'(x) = e^x$

Alternatively, in some situations it may be easier to define the difference quotient as:

$$Q(h) = \frac{f(h+c) - f(h)}{h} \quad (10.0.14)$$

so that differentiability at  $c$  is given when the following limit exists:

$$\lim_{h \rightarrow 0} Q(h) \quad (10.0.15)$$

*Proof.* (a) we have that

$$Q(x) = \frac{k - k}{x - c} = 0 \implies f'(x) = 0 \quad (10.0.16)$$

(b) we have that:

$$Q(h) = \frac{(c+h)^n - c^n}{h} = nc^{n-1} + \frac{n(n-1)}{2}c^{n-2}h + \dots + h^{n-1} \quad (10.0.17)$$

so that  $Q(h) \rightarrow nc^{n-1}$  as  $h \rightarrow 0$ .

(c) The difference quotient is:

$$Q(h) = \frac{\sin(c+h) - \sin c}{h} = \frac{\sin c \cos h + \sin h \cos c - \sin c}{h} \quad (10.0.18)$$

$$= \cos c \frac{\sin h}{h} + \sin c \left( \frac{\cos h - 1}{h} \right) \quad (10.0.19)$$

so using some standard trigonometric limits:

$$\lim_{h \rightarrow 0} Q(h) = \cos c \implies f'(x) = \cos x \quad (10.0.20)$$

(d) The difference quotient is:

$$Q(h) = \frac{\cos(c+h) - \cos c}{h} = \frac{\cos c \cos h - \sin h \sin c - \cos c}{h} \quad (10.0.21)$$

$$= -\sin c \frac{\sin h}{h} + \cos c \left( \frac{\cos h - 1}{h} \right) \quad (10.0.22)$$

so using some standard trigonometric limits:

$$\lim_{h \rightarrow 0} Q(h) = -\sin c \implies f'(x) = -\sin x \quad (10.0.23)$$

(e) The difference quotient is:

$$Q(h) = \frac{e^c e^h - e^c}{h} = e^c \frac{e^h - 1}{h} \quad (10.0.24)$$

so  $Q(h) \rightarrow e^c$  as  $h \rightarrow 0$  proving that  $f'(x) = e^x$ . ■

### Definition 10.3 (One-sided derivative)

Let  $f$  be defined on  $I$  and let  $c \in I$ . Then the **left derivative of  $f$  at  $c$**  is:

$$f'_L(c) = \lim_{x \rightarrow c^-} \frac{f(x) - f(c)}{x - c} = \lim_{h \rightarrow 0^-} Q(h) \quad (10.0.25)$$

If this limit exists, then  $f$  is left-differentiable at  $c$ . Similarly for the right derivative:

$$f'_R(c) = \lim_{x \rightarrow c^+} \frac{f(x) - f(c)}{x - c} = \lim_{h \rightarrow 0^+} Q(h) \quad (10.0.26)$$

### Proposition 10.4 (Differentiability and one-sided differentiability)

Let  $f$  be defined on  $I$  and let  $c \in I$ .

$$f \text{ is differentiable at } c \iff f'_R(c) = f'_L(c) = f'(c). \quad (10.0.27)$$

### Theorem 10.5 (Glue rule for differentiation)

Let  $f$  be defined on  $I$  and let  $c \in I$ . If there are functions  $g, h$  defined on  $I$  such that

- (1)  $f(x) = g(x)$  for  $x \in I, x < c$  and  $f(x) = h(x)$  for  $x \in I, x > c$
- (2)  $f(c) = g(c) = h(c)$
- (3)  $g, h$  are differentiable at  $c$

then  $f$  is differentiable at  $c$  iff  $g'(c) = h'(c)$ . In this case then  $f'(c) = g'(c) = h'(c)$ .

Note also that since differentiability is a local property, if we define a piece-wise function such as:

$$f(x) = \begin{cases} g(x), & x > c \\ h(x), & x \leq c \end{cases} \quad (10.0.28)$$

then we will have that:

$$f'(x) = \begin{cases} g'(x), & x > c \\ h'(x), & x < c \end{cases} \quad (10.0.29)$$

**Example.** Let us prove that:

$$f(x) = \begin{cases} -x^2, & x < 0 \\ x^2, & x \geq 0 \end{cases} \quad (10.0.30)$$

is differentiable on  $\mathbb{R}$ . Indeed define  $g(x) = -x^2$  and  $h(x) = x^2$  so that  $f(x) = g(x)$  for  $x \in \mathbb{R}, x < 0$  and  $f(x) = h(x)$  for  $x \in \mathbb{R}, x > 0$ . Moreover, we also have that:

$$f(0) = g(0) = h(0) = 0 \quad (10.0.31)$$

Finally, we also know that  $g'(0) = 0$  and  $h'(0) = 0$ , so that by the Glue rule  $f'(0) = 0$ .

$$f'(x) = \begin{cases} g'(x) = -2x, & x < 0 \\ h'(x) = 2x, & x > 0 \\ 0, & x = 0 \end{cases} = 2|x|, \quad x \in \mathbb{R} \quad (10.0.32)$$

as desired. ◀

## 10.1 Continuity and differentiability

### Theorem 10.6 (Differentiability implies continuity)

Let  $f$  be defined on  $I$ , and let  $c$ . If  $f$  is differentiable at  $c$  then it is continuous at  $c$ .

*Proof.* Suppose  $f$  is differentiable at  $c$  so that

$$f'(c) = \lim_{x \rightarrow c} \frac{f(x) - f(c)}{x - c} \quad (10.1.1)$$

For  $x \in I$  and  $x \neq c$  we have that:

$$f(x) - f(c) = f'(c)(x - c) \quad (10.1.2)$$

Hence

$$\lim_{x \rightarrow c} f(x) - f(c) = f'(c) \cdot 0 = 0 \quad (10.1.3)$$

so that  $f(x) \rightarrow f(c)$  implying continuity. ■

## 10.2 Rules of differentiation

### Proposition 10.7 (Combination rules)

Let  $f, g$  be defined on  $I$  and let  $c \in I$ . If  $f, g$  are differentiable at  $c$  then:

- (i)  $(f + g)'(c) = f'(c) + g'(c)$
- (ii)  $(\lambda f)'(c) = \lambda f'(c)$  with  $\lambda \in \mathbb{R}$
- (iii)  $(fg)'(c) = f'(c)g(c) + f(c)g'(c)$
- (iv) if  $g(c) \neq 0$  then  $(\frac{f}{g})'(c) = \frac{g(c)f'(c) - f(c)g'(c)}{(g(c))^2}$

*Proof.* (i) Let  $F = f + g$ , then:

$$\frac{F(x) - F(c)}{x - c} = \frac{f(x) - f(c)}{x - c} + \frac{g(x) - g(c)}{x - c} \quad (10.2.1)$$

$$\rightarrow f'(c) + g'(c) \quad (10.2.2)$$

as required.

(ii) Use product rule with  $f = \lambda$ .

(iii) Let  $F = fg$  then:

$$\frac{F(x) - F(c)}{x - c} = \frac{f(x)g(x) - f(c)g(c)}{x - c} + \frac{g(x) - g(c)}{x - c} \quad (10.2.3)$$

$$= \frac{f(x) - f(c)}{x - c}g(c) + f(c)\frac{g(x) - g(c)}{x - c} \quad (10.2.4)$$

$$\rightarrow f'(c)g(c) + f(c)g'(c) \quad (10.2.5)$$

as required.

(iv) Let  $F = \frac{f}{g}$ , since  $g$  is continuous at  $c$  and  $g(c) \neq 0$ , there exists  $\delta > 0$  such that  $J = (c - \delta, c + \delta) \subseteq I$  and:

$$|g(x) - g(c)| < \frac{1}{2}|g(c)|, \forall x \text{ with } |x - c| < \delta \quad (10.2.6)$$

This shows that there is some  $J$  such that  $g(x) \neq 0$  for  $x \in J$ . Therefore, for  $x \in J$  we have

that:

$$\frac{F(x) - F(c)}{x - c} = \frac{f(x)/g(x) - f(c)/g(c)}{x - c} + \frac{g(x) - g(c)}{x - c} \quad (10.2.7)$$

$$= \frac{f(x)g(c) - f(c)g(x)}{(x - c)g(x)g(c)} \quad (10.2.8)$$

$$= \frac{g(c)(f(x) - f(c)) - f(c)(g(x) - g(c))}{(x - c)g(x)g(c)} \quad (10.2.9)$$

$$= \frac{1}{g(x)g(c)} \left( g(c) \frac{f(x) - f(c)}{x - c} - f(c) \frac{g(x) - g(c)}{x - c} \right) \quad (10.2.10)$$

$$\rightarrow \frac{g(c)f'(c) - f(c)g'(c)}{(g(c))^2} \quad (10.2.11)$$

as desired. ■

**Theorem 10.8 (Composition rule)** Let  $f$  be defined on  $I$  and let  $g$  be defined on  $J$  so that  $f(I) \subset J$  and let  $c \in I$ .

If  $f$  is differentiable at  $c$  and  $g$  is differentiable at  $f(c)$  then

$$(g \circ f)'(c) = g'(f(c))f'(c) \quad (10.2.12)$$

*Proof.* Let  $F = g \circ f$  then:

$$\frac{F(x) - F(c)}{x - c} = \frac{g(f(x)) - g(f(c))}{x - c} \quad (10.2.13)$$

Let  $y = f(x)$  with  $x \in I$  and let  $d = f(c)$ . Then the right hand side of the above equation is:

$$\left( \frac{g(y) - g(d)}{y - d} \right) \left( \frac{f(x) - f(c)}{x - c} \right), \quad y \neq d \quad (10.2.14)$$

We may introduce the function  $h(y)$  to deal with the discontinuity at  $y \neq d$ :

$$h(y) = \begin{cases} \frac{g(y) - g(d)}{y - d}, & y \neq d \\ g'(f(c)), & f(x) = f(c) \end{cases} \quad (10.2.15)$$

Since  $g$  is differentiable at  $d$ , we have  $h(y) \rightarrow g'(d)$  as  $y \rightarrow d$ . Since  $h(d) = g'(d)$  we have that  $h$  is continuous at  $d$ . We deduce that:

$$(h \circ f)(x) = \begin{cases} \frac{g(f(x)) - g(f(c))}{f(x) - f(c)}, & f(x) \neq f(c) \\ g'(f(c)), & f(x) = f(c) \end{cases} \quad (10.2.16)$$

is continuous at  $c$ .

Therefore:

$$\frac{F(x) - F(c)}{x - c} = (h \circ f)(x) \left( \frac{f(x) - f(c)}{x - c} \right) \quad (10.2.17)$$

so that as  $x \rightarrow c$  then:

$$\frac{F(x) - F(c)}{x - c} \rightarrow g'(f(c))f'(c) \quad (10.2.18)$$

as desired. ■

**Example.** Let's find the derivative of  $f(x) = \cos\left(\frac{\cos 2x}{x^2}\right)$ .

Let us define

$$g(x) = \cos x, x \in I \quad (10.2.19)$$

and

$$h(x) = \frac{\cos 2x}{x^2}, x \in I = (0, \infty) \quad (10.2.20)$$

Then  $h(I) \subseteq (0, \infty) = I$  so we may apply the composition rule:

$$(g \circ f)'(x) = -\sin\left(\frac{\cos 2x}{x^2}\right) \cdot \frac{-2x^2 \sin 2x - 2x \cos 2x}{x^4} \quad (10.2.21)$$

$$= 2 \sin\left(\frac{\cos 2x}{x^2}\right) \frac{x \sin 2x + \cos 2x}{x^3} \quad (10.2.22)$$

◀

### Proposition 10.9 (*Inverse function rule*)

Let  $f$  be a function with domain  $I$  on which it is continuous and strictly monotonic. If it is differentiable on  $I$  and  $f'(x) \neq 0, \forall x \in I$ , then  $f^{-1}$  is differentiable on  $J$ . For  $c \in I$  and  $d = f(c)$  then:

$$(f^{-1})'(d) = \frac{1}{f'(c)} \quad (10.2.23)$$

*Proof.* We have that  $f$  is invertible on  $I$  with inverse  $f^{-1}$  whose domain is  $J = f(I)$ . Let  $y \in J \setminus \{d\}$ , it follows that  $f^{-1}(y) = x \in I \setminus c$  due to the strict monotonicity of  $f$ . Therefore we find that:

$$\frac{f^{-1}(y) - f^{-1}(d)}{y - d} = \frac{x - c}{f(x) - f(c)} = \frac{1}{\frac{f(x) - f(c)}{x - c}} \quad (10.2.24)$$

Taking the limit as  $y \rightarrow d \implies x = f^{-1}(y) \rightarrow c$  due to the continuity of  $f^{-1}$ . Hence:

$$\lim_{y \rightarrow d} \frac{f^{-1}(y) - f^{-1}(d)}{y - d} = \frac{1}{f'(c)} \quad (10.2.25)$$

proving that  $f^{-1}$  is differentiable at  $d$  with derivative  $(f^{-1})'(d) = \frac{1}{f'(c)}$ . ■

**Example.** Let us consider  $f(x) = \tan x, x \in (-\pi/2, \pi/2)$ . The domain of this function is  $I = (-\pi/2, \pi/2)$ , over which it is continuous and strictly increasing. Hence  $f$  will have an inverse  $f^{-1}$  with domain  $f(I) = \mathbb{R}$ .

Furthermore,  $f$  is differentiable on  $I$ , and its derivative is  $f'(x) = \sec^2 x$ , which is non-zero  $\forall x \in I$ . Therefore  $f^{-1}$  must be differentiable on  $\mathbb{R}$  by the Inverse function rule, with derivative at  $y \in f(c), \forall x \in (-\pi/2, \pi/2)$  given by:

$$(f^{-1})'(y) = \frac{1}{\sec^2 x} = \frac{1}{1 + \tan^2 x} = \frac{1}{1 + y^2} \quad (10.2.26)$$

implying that:

$$(\tan^{-1})'(x) = \frac{1}{1+x^2}, \forall x \in \mathbb{R} \quad (10.2.27)$$

◀

## 10.3 Rolle's theorem and local extrema

### Definition (Local extrema)

The function  $f$  with domain  $J$  is said to have a:

1. **local maximum**  $f(c)$  at  $x = c$  if there exists  $I = (c - r, c + r) \subseteq J$  where  $r > 0$  such that:

$$f(x) \leq f(c), \forall x \in I \quad (10.3.1)$$

2. **local minimum**  $f(c)$  at  $x = c$  if there exists  $I = (c - r, c + r) \subseteq J$  where  $r > 0$  such that:

$$f(c) \leq f(x), \forall x \in I \quad (10.3.2)$$

3. **local extremum**  $f(c)$  at  $x = c$  if  $f(c)$  is either a local maximum or minimum

Therefore, a local extremum is a value of  $f$  at some point such that  $f$  has lower values in some neighborhood of  $c$ .

### Theorem (Local extreme value theorem)

If  $f$  has a local extremum at  $c$  and is differentiable at  $c$  then  $f'(c) = 0$

*Proof.* Suppose that  $f$  has a local maximum at  $c$  so that  $\exists r > 0$  such that:

$$f(x) \leq f(c), \text{ for } c - r < x < c + r \quad (10.3.3)$$

Let  $x_n = c + \frac{r}{n}$  and  $x'_n = c - \frac{r}{n}$  for  $n = 2, 3, \dots$ . Then  $c < x_n < c + r$  so that  $f(x_n) \leq f(c)$  and  $x_n > c$ . Hence:

$$\frac{f(x_n) - f(c)}{x_n - c} \leq 0 \implies f'(c) \leq 0 \quad (10.3.4)$$

We also have that  $c - r < x'_n < c$  so that  $f(x'_n) \leq f(c)$  and  $x'_n < c$  and hence:

$$\frac{f(x'_n) - f(c)}{x'_n - c} \geq 0 \implies f'(c) \geq 0 \quad (10.3.5)$$

We deduce that  $f'(c) = 0$  as desired. ■

It is important to note that the converse of the local extreme value theorem is not necessarily true. All we know is that local extreme values of a differentiable function  $f$  on  $[a, b]$  occur either at  $x = a, x = b$  or at points  $x \in (a, b)$  where  $f'(x) = 0$ .

**Example.** Let's find the local extrema of  $f(x) = \sin^2 x + \cos x$  on  $[0, \pi/2]$ . Firstly  $f$  is continuous on  $[0, \pi/2]$ , so we have that  $f(0) = 1$  and  $f(\pi/2) = 1$ . Moreover  $f$  is differentiable

on  $(0, \pi/2)$  with:

$$f'(x) = 2 \sin x \cos x - \sin x \quad (10.3.6)$$

which vanishes when:

$$2 \sin x \cos x = \sin x \implies \sin x = 0 \text{ or } \cos x = \frac{1}{2} \quad (10.3.7)$$

We see that  $\sin x = 0 \implies x = \pi n, \forall n \in \mathbb{Z}$ . Instead  $\cos x = \frac{1}{2} \implies x = 2\pi n + \frac{\pi}{3}, \forall n \in \mathbb{Z}$ . Since we're restricted to the interval  $(0, \pi/2)$  we only consider  $x = \frac{\pi}{3}$  where  $f(x) = \frac{3}{4} + \frac{1}{2} = \frac{5}{4}$ . We see that this provides the largest value of  $f$  compared to  $x = 0, x = \frac{\pi}{2}$  and is therefore the local maximum. The local minimum, on the other hand, occurs at the endpoints  $x = 0, x = \frac{\pi}{2}$  where  $f(x) = 1$ .  $\blacktriangleleft$

### Theorem (Rolle's theorem)

Let  $f$  be continuous on  $[a, b]$  and differentiable on  $(a, b)$ . If  $f(a) = f(b)$  then there exists  $c$  with  $c \in (a, b)$  such that  $f'(c) = 0$ .

*Proof.* Suppose  $f$  is constant on  $[a, b]$ , then clearly  $f'(x) = 0$  everywhere on  $(a, b)$ .

Suppose  $f$  is not constant on  $[a, b]$ . Since  $f$  is continuous it must have both a maximum and minimum on  $[a, b]$ . At least one of these must be different from  $f(a) = f(b)$  since  $f$  is non-constant. Hence  $f$  has an extreme value for some  $c \in (a, b)$ . The extreme value theorem then shows that  $f'(c) = 0$  as desired.  $\blacksquare$

## 10.4 Mean value theorem

### Theorem (Mean value theorem)

Let  $f$  be continuous on  $[a, b]$  and differentiable on  $(a, b)$ . Then  $\exists c \in (a, b)$  such that:

$$f'(c) = \frac{f(b) - f(a)}{b - a} \quad (10.4.1)$$

*Proof.* The gradient of the chord joining the points  $(a, f(a))$  and  $(b, f(b))$  is:

$$m = \frac{f(b) - f(a)}{b - a} \quad (10.4.2)$$

Consequently its equation will be  $y = \frac{f(b) - f(a)}{b - a}(x - a) + f(a)$ . Let us then define:

$$h(x) = f(x) - \frac{f(b) - f(a)}{b - a}(x - a) - f(a) \quad (10.4.3)$$

We then have that  $h(a) = h(b) = 0$  and that  $h$  is continuous on  $[a, b]$ , differentiable on  $(a, b)$ . Consequently Rolle's theorem tells us that there exists some  $c \in (a, b)$  for which:

$$h'(c) = f'(c) - \frac{f(b) - f(a)}{b - a} = 0 \implies f'(c) = \frac{f(b) - f(a)}{b - a} \quad (10.4.4)$$

as desired. ■

**Example.** Let's consider  $f(x) = xe^x$  over the interval  $(0, 2)$ . Clearly  $f$  is differentiable on this interval, and continuous on  $[0, 2]$  by the product rule. We also have that:

$$m = \frac{2e^2}{2} = e^2 \implies \exists c \in (0, 2) \text{ s.t. } f'(c) = e^2 \quad (10.4.5)$$

by the mean value theorem. ◀

**Proposition (Increasing-Decreasing)** Let  $f$  be continuous on  $I$  and differentiable on the interior  $J$  of  $I$ . Then if:

- (i)  $f'(x) \leq 0, \forall x \in J$  then  $f$  is decreasing on  $I$ .
- (ii)  $f'(x) \geq 0, \forall x \in J$  then  $f$  is increasing on  $I$ .

*Proof.* Let us take  $x_1, x_2 \in I$  such that  $x_1 < x_2$ . Then since  $f$  satisfies the condition for the mean value theorem there exists  $c \in (x_1, x_2)$  such that:

$$f'(c) = \frac{f(x_2) - f(x_1)}{x_2 - x_1} \quad (10.4.6)$$

If  $f'(x) \leq 0 \forall x \in J$  then  $f'(c) \leq 0$  and hence  $f(x_2) - f(x_1) \leq 0$  proving that  $f$  is decreasing on  $I$ .

If  $f'(x) \geq 0 \forall x \in J$  then  $f'(c) \geq 0$  and hence  $f(x_2) - f(x_1) \geq 0$  proving that  $f$  is increasing on  $I$ . ■

This gives us an efficient way to prove inequalities. Indeed, suppose we wished to prove that

$$g(x) \geq h(x), \forall x \in [a, b] \quad (10.4.7)$$

Then we let  $f(x) = g(x) - h(x)$ , and if it is continuous on  $[a, b]$  and differentiable on  $(a, b)$  we show that either:

$$f(a) \geq 0 \text{ and } f'(x) \geq 0 \forall x \in (a, b) \quad (10.4.8)$$

which shows that  $f$  is smallest at  $a$  in  $[a, b]$ , or:

$$f(b) \geq 0 \text{ and } f'(x) \leq 0 \forall x \in (a, b) \quad (10.4.9)$$

which shows that  $f$  is smallest at  $b$  in  $[a, b]$ . In both cases we have that  $f$  will always be positive, and hence that  $g(x) \geq h(x)$  on  $[a, b]$ .

**Example.** Let's prove that for  $\alpha \geq 1$  and  $x \geq -1$ :

$$(1+x)^\alpha \geq 1 + \alpha x \quad (10.4.10)$$

The case where  $\alpha = 1$  is clearly true, so we assume that  $\alpha > 1$ .

Let us define  $f(x) = (1+x)^\alpha - 1 - \alpha x$  for  $x \in I = [-1, \infty)$  which is continuous on  $I$  and differentiable on its interior. The derivative of  $f$  over  $I$  is:

$$f'(x) = \alpha(1+x)^{\alpha-1} - \alpha = \alpha((1+x)^{\alpha-1} - 1) \quad (10.4.11)$$

We see that for  $-1 < x < 0$  then  $0 < 1 + x < 1$  so that  $0 < (1 + x)^{\alpha-1} < 1$  and hence  $f'(x) < 0$ ,  $f$  is decreasing on  $(-1, 0)$ .

Similarly for  $0 < x$  then  $1 < 1 + x$  so that  $1 < (1 + x)^{\alpha-1}$  and hence  $f'(x) > 0$ ,  $f$  is increasing on  $(0, \infty)$ .

Finally,  $f(0) = 0 \geq 0$ , from which it follows that  $f(x) \geq 0$  for all  $x \in [-1, \infty)$ , as desired.  $\blacktriangleleft$

### Theorem (Second derivative test)

Let  $f$  be a twice-differentiable function defined on the open interval  $I$  containing  $c$ , such that  $f'(c) = 0$  and  $f''$  is continuous at  $c$ .

- (i) if  $f''(c) > 0$  then  $f(c)$  is a local minimum of  $f$ .
- (ii) if  $f''(c) < 0$  then  $f(c)$  is a local maximum of  $f$ .

*Proof.* Suppose that  $f''(c) > 0$ . Since  $f''$  is continuous at  $c$ , we have that  $\exists \delta > 0$  such that  $(c - \delta, c + \delta) \subseteq I$  and:

$$|f''(x) - f''(c)| < \epsilon = \frac{1}{2}f''(c), \forall x \in (c - \delta, c + \delta) \quad (10.4.12)$$

implying that  $f''(x) > \frac{1}{2}f''(c) > 0$ . Hence  $f'$  is strictly increasing on  $(c - \delta, c + \delta)$ . Furthermore we have that  $f'(c) = 0$  so:

$$f'(x) < 0, \forall x \in (c - \delta, c) \quad f'(x) > 0, \forall x \in (c, c + \delta) \quad (10.4.13)$$

This implies that

$$f(x) \text{ is decreasing, } \forall x \in (c - \delta, c) \quad f(x) \text{ is increasing, } \forall x \in (c, c + \delta) \quad (10.4.14)$$

proving that  $f$  has a local minimum at  $c$ .  $\blacksquare$

## 10.5 L'Hopital's rule

### Theorem (Cauchy's mean value theorem)

Let  $f, g$  be continuous on  $[a, b]$  and differentiable on  $(a, b)$ . Then  $\exists c \in (a, b)$  such that:

$$f'(c)(g(b) - g(a)) = g'(c)(f(b) - f(a)) \quad (10.5.1)$$

*Proof.* Let us define:

$$h(x) = f(x)(g(b) - g(a)) - g(x)(f(b) - f(a)) \quad (10.5.2)$$

Clearly by the combination rules of continuity and differentiability,  $h$  must be continuous on  $[a, b]$  and differentiable on  $(a, b)$ . Note also that:

$$h(a) = f(a)g(b) - g(a)f(b) = h(b) \quad (10.5.3)$$

We may therefore apply Rolle's Theorem:

$$\exists c \in (a, b) \text{ s.t. } h'(c) = f'(c)(g(b) - g(a)) - g'(c)(f(b) - f(a)) = 0 \quad (10.5.4)$$

thus implying that:

$$f'(c)(g(b) - g(a)) = g'(c)(f(b) - f(a)) \quad (10.5.5)$$

as desired. ■

### Theorem (L'Hopital's rule)

Let  $f, g$  be differentiable on  $I$  which contains  $c$ , and suppose  $f(c) = g(c) = 0$ . Then:

$$\lim_{x \rightarrow c} \frac{f(x)}{g(x)} = \lim_{x \rightarrow c} \frac{f'(x)}{g'(x)} \quad (10.5.6)$$

provided the latter limit exists.

*Proof.* Assume that:

$$\lim_{x \rightarrow c} \frac{f'(x)}{g'(x)} = l \quad (10.5.7)$$

Let  $\epsilon > 0$ , then it follows from the existence of the above limit that  $\exists \delta > 0$  such that:

$$\left| \frac{f'(x)}{g'(x)} - l \right| < \epsilon, \quad 0 < |x - c| < \delta \quad (10.5.8)$$

so we cannot have that  $g'(x) = 0$  within  $0 < |x - c| < \delta$ . Note however that if  $g(x_0) = g(c)$  for some  $x_0$  (wlog  $x_0 > c$ ), then by Rolle's theorem  $g'(x) = 0$  for some  $x \in (c, x_0)$ , which is impossible. Hence we must have that  $g(x) \neq g(c)$  for all  $0 < |x - c| < \delta$ . Applying Cauchy's mean value theorem, we have some  $d \in (c, x)$  such that:

$$f'(d)(g(x) - g(c)) = g'(d)(f(x) - f(c)) \implies \frac{f'(d)}{g'(d)} = \frac{f(x) - f(c)}{g(x) - g(c)} \quad (10.5.9)$$

and since  $f(c) = g(c) = 0$  we get that:

$$\frac{f'(d)}{g'(d)} = \frac{f(x)}{g(x)} \quad (10.5.10)$$

Therefore:

$$\left| \frac{f(x)}{g(x)} - l \right| = \left| \frac{f'(d)}{g'(d)} - l \right| < \epsilon, \quad 0 < |x - c| < \delta \quad (10.5.11)$$

proving that

$$\lim_{x \rightarrow c} \frac{f(x)}{g(x)} = l \quad (10.5.12)$$

as desired. ■

**Example.** Consider

$$\lim_{x \rightarrow 0} \frac{\sin x - x \cos x}{x^3} \quad (10.5.13)$$

We see that  $f(x) = \sin x - x \cos x$  is continuous and differentiable in  $\mathbb{R}$  by the combination rules, and so is  $g(x) = x^3$ . Moreover, we also have that  $f(0) = g(0) = 0$ , hence we may apply l'Hopital's theorem:

$$\lim_{x \rightarrow 0} \frac{\sin x - x \cos x}{x^3} = \lim_{x \rightarrow 0} \frac{\cos x - \cos x + x \sin x}{3x^2} = \lim_{x \rightarrow 0} \frac{\sin x}{3x} = \frac{1}{3} \quad (10.5.14)$$



# Unit F3: Integration

## 11.1 The Riemann integral

### Definition (*Partition*)

A partition  $P$  of a closed interval  $[a, b]$  is a collection of closed non-intersecting subintervals whose union gives  $[a, b]$ :

$$P = \{[x_0, x_1], [x_1, x_2], \dots, [x_{i-1}, x_i], \dots, [x_{n-1}, x_n]\} \quad (11.1.1)$$

with:

$$a = x_0 < x_1 < \dots < x_i < \dots < x_n = b \quad (11.1.2)$$

The points  $x_i$  are known as partition points, and the  $i$ th subinterval is  $I_i = [x_{i-1}, x_i]$  whose length is  $\delta x_i = x_i - x_{i-1}$ . Instead the mesh of  $P$  is defined as  $\|P\| = \max_{1 \leq i \leq n} \{\delta x_i\}$ .

**Example.** Consider the following partition  $P$  of  $[0, 1]$ :

$$P = \left\{ \left[0, \frac{1}{2}\right], \left[\frac{1}{2}, \frac{3}{5}\right], \left[\frac{3}{5}, \frac{3}{4}\right], \left[\frac{3}{4}, 1\right] \right\} \quad (11.1.3)$$

We find that:

$$\delta x_1 = \frac{1}{2}, \delta x_2 = \frac{1}{10}, \delta x_3 = \frac{3}{20}, \delta x_4 = \frac{1}{4} \quad (11.1.4)$$

so that:

$$\|P\| = \max\{\delta x_1, \delta x_2, \delta x_3, \delta x_4\} = \frac{1}{2} \quad (11.1.5)$$

◀

### Definition (*Riemann sums*)

Let  $f$  be bounded on  $[a, b]$  and let  $P = \{[a, x_1], [x_1, x_2], \dots, [x_{i-1}, x_i], \dots, [x_{n-1}, b]\}$ . Define:

$$m_i = \inf_{x \in I_i} f, M_i = \sup_{x \in I_i} f \quad (11.1.6)$$

Then the lower Riemann sum for  $f$  on  $[a, b]$  is:

$$L(f, P) = \sum_{i=1}^n m_i \delta x_i \quad (11.1.7)$$

while the upper Riemann sum for  $f$  on  $[a, b]$  is:

$$U(f, P) = \sum_{i=1}^n M_i \delta x_i \quad (11.1.8)$$

Geometrically, the upper Riemann sum represents an upper bound for the area under  $f$ , whereas the lower Riemann sum represents a lower bound for the area under  $f$ .

At the essence of Riemann integration is that we may approximate a function as constant over a sufficiently small interval  $I_i = [x_{i-1}, x_i]$ . Doing this for the entire partition of  $P$ , we get a series of intervals over which  $f$  is taken to be constant. In other words, we may consider the area under  $f$  as a series of vertical strips of width  $\delta x_i$ . But what constant value should we take for the height of the columns? Well,  $m_i$  gives the largest lower bound of  $f$  over  $I_i$ , while  $M_i$  gives the smallest upper bound of  $f$ . Consequently  $m_i \delta x_i$  will underestimate the area of the  $i$ th vertical strip, while  $M_i \delta x_i$  will overestimate it.

**Example.** Consider the function

$$f(x) = \begin{cases} 2x, & 0 < x < 1 \\ 1, & x = 0, 1 \end{cases} \quad (11.1.9)$$

and the partition  $P = \left\{ [0, \frac{1}{4}], [\frac{1}{4}, \frac{1}{2}], [\frac{1}{2}, \frac{3}{4}], [\frac{3}{4}, 1] \right\}$ . Then we see that:

$$m_1 = 0, M_1 = 1, \delta x_1 = \frac{1}{4} \quad (11.1.10)$$

$$m_2 = \frac{1}{2}, M_2 = 1, \delta x_2 = \frac{1}{4} \quad (11.1.11)$$

$$m_3 = 1, M_3 = \frac{3}{2}, \delta x_3 = \frac{1}{4} \quad (11.1.12)$$

$$m_4 = 1, M_4 = 2, \delta x_4 = \frac{1}{4} \quad (11.1.13)$$

$$(11.1.14)$$

and therefore:

$$L(f, P) = \frac{1}{4}(0 + \frac{1}{2} + 1 + 1) = \frac{5}{8} \quad (11.1.15)$$

$$U(f, P) = \frac{1}{4}(1 + 1 + \frac{3}{2} + 2) = \frac{11}{8} \quad (11.1.16)$$

◀

**Example.** Let

$$f(x) = \begin{cases} x^2, & 0 \leq x \leq 1 \\ 2, & 1 < x \leq 2 \end{cases} \quad (11.1.17)$$

and the partition  $P = \left\{ [0, \frac{1}{n}], [\frac{1}{n}, \frac{2}{n}], \dots, [2 - \frac{1}{n}, 2] \right\}$ .

Now we see that the  $i$ th interval is  $[\frac{i-1}{n}, \frac{i}{n}]$ , while the interval width is  $\delta x_i = \frac{1}{n}$ . Also, note that  $f$  is increasing over  $[0, 2]$ , so  $m_i$  will be the value of  $f$  at the left endpoint of the  $i$ th inter-

val, while  $M_i$  will be the value of  $f$  at the right endpoint of the  $i$ th interval. Consequently, for  $1 \leq i \leq n$  we find that  $m_i = \frac{(i-1)^2}{n^2}$ , while  $M_i = \frac{i^2}{n^2}$ . Instead for  $i = n+1$ ,  $m_i = 1$  (which is coherent with  $m_i = \frac{(i-1)^2}{n^2}$ ) whereas  $M_i = 2$ . Finally for  $n+2 \leq i \leq 2n$  we get that  $m_i = M_i = 2$ . Hence:

$$L(f, P) = \sum_{i=1}^{n+1} \frac{(i-1)^2}{n^2} \frac{1}{n} + \sum_{i=n+2}^{2n} 2 \frac{1}{n} \quad (11.1.18)$$

$$= \frac{1}{n^3} \sum_{i=0}^n i^2 + \sum_{i=n+2}^{2n} 2 \frac{1}{n} \quad (11.1.19)$$

$$= \frac{1}{n^3} \frac{n(n+1)(2n+1)}{6} + \frac{2}{n}(2n-n-1) \quad (11.1.20)$$

$$= \frac{n(n+1)(2n+1)}{6n^3} + 2 - \frac{2}{n} \quad (11.1.21)$$

Instead:

$$U(f, P) = \sum_{i=1}^n \frac{i^2}{n^2} \frac{1}{n} + \sum_{i=n+1}^{2n} 2 \frac{1}{n} \quad (11.1.22)$$

$$= \frac{n(n+1)(2n+1)}{6n^3} + 2 \frac{1}{n}(2n-n) = \frac{n(n+1)(2n+1)}{6n^3} + 2 \quad (11.1.23)$$

Note that taking the limit as  $n \rightarrow \infty$ :

$$\lim_{n \rightarrow \infty} L(f, P) = \frac{1}{3} + 2 = \frac{7}{3} \quad (11.1.24)$$

and similarly:

$$\lim_{n \rightarrow \infty} U(f, P) = \frac{1}{3} + 2 = \frac{7}{3} \quad (11.1.25)$$

We define this number as the integral of  $f$  on  $[a, b]$ . ◀

### Theorem (Lower and upper Riemann sum inequality)

Let  $f$  be bounded on  $[a, b]$  and let  $P, P'$  be partitions of  $[a, b]$ . Then  $L(f, P) \leq U(f, P')$ .

*Proof.* Suppose  $f$  is non-negative on  $[a, b]$ .

Since for any given partition  $P$ ,  $m_i \leq M_i$  we have that  $L(f, P) \leq U(f, P)$ . Also, let  $P''$  be the union of the  $P$  and  $P'$  partitions. For example, if  $P = \left\{ [0, \frac{1}{4}], [\frac{1}{4}, \frac{1}{2}] \right\}$  and  $P' = \left\{ [0, \frac{1}{3}], [\frac{1}{3}, \frac{1}{2}] \right\}$ , then  $P'' = \left\{ [0, \frac{1}{4}], [\frac{1}{4}, \frac{1}{3}], [\frac{1}{3}, \frac{1}{2}] \right\}$ .

Now consider what happens when we add a new partition  $x'$  to  $P = \{[x_0, x_1], \dots, [x_{n-1}, x_n]\}$ . Since  $f$  is positive adding this partition will not increase the upper Riemann sum. Similarly, adding the partition will not decrease the lower Riemann sum.

In other words, if we add a refinement from the  $P$  partition onto the  $P'$  partition to form the  $P''$  partition, we will find that  $U(f, P'') \leq U(f, P')$ . Similarly if we add a refinement from the  $P'$  partition onto the  $P$  partition to form the  $P''$  partition, we will find that  $L(f, P) \leq L(f, P'')$ .

Therefore, we may write that:

$$L(f, P) \leq L(f, P'') \leq U(f, P'') \leq U(f, P') \quad (11.1.26)$$

as desired.

If instead  $f$  is negative over some interval, then because it is bounded we may still form the function  $g = f + c$  where  $c$  is some constant. Applying the same reasoning as before we will find that  $L(g, P) \leq U(g, P')$  and thus  $L(f, P) \leq U(f, P')$  since the upper and lower riemann sums of a constant function  $c$  are identical. ■

### Definition (*Integral*)

Let  $f$  be a bounded function on  $[a, b]$ , and let  $P$  be a partition of  $[a, b]$ . Then the lower integral of  $f$  on  $[a, b]$  is:

$$\underline{\int_a^b} f = \sup_P L(f, P) \quad (11.1.27)$$

while the upper integral of  $f$  is:

$$\overline{\int_a^b} f = \inf_P U(f, P) \quad (11.1.28)$$

If these two are equal, then we say that  $f$  is integrable on  $[a, b]$ . The values they are equal to is the integral of  $f$  on  $[a, b]$ .

In the previous example, we would write that:

$$\underline{\int_0^2} f = \overline{\int_0^2} f = \frac{7}{3} \implies \int_0^2 f = \frac{7}{3} \quad (11.1.29)$$

**Example.** Let  $f$  be the Dirichlet function defined by

$$f(x) = \begin{cases} 1, & 0 \leq x \leq 1, x \text{ is irrational} \\ 0, & 0 \leq x \leq 1, x \text{ is rational} \end{cases} \quad (11.1.30)$$

with a partition  $P = \{[0, x_1], [x_1, x_2], \dots, [x_{n-1}, 1]\}$ . Due to the density of real numbers, we can always find a rational and irrational number within each interval, so that  $m_i = 0$  and  $M_i = 1$ . Hence we find that:

$$L(f, P) = \sum_i m_i \delta x_i = 0 \quad (11.1.31)$$

while

$$U(f, P) = \sum_i M_i \delta x_i = \sum_i \delta x_i = 1 \quad (11.1.32)$$

Therefore:

$$\underline{\int_0^1} f = 0 \neq \overline{\int_0^1} f = 1 \quad (11.1.33)$$

from which it follows that  $f$  is not integrable on  $[0, 1]$ . ■

**Theorem (Integrability)**

Let  $f$  be bounded on  $[a, b]$ , if there exists a sequence of partitions  $(P_n)$  of  $[a, b]$  such that  $\|P_n\| \rightarrow 0$  then:

$$\lim_{n \rightarrow \infty} L(f, P_n) = \lim_{n \rightarrow \infty} U(f, P_n) = A \in \mathbb{R} \quad (11.1.34)$$

then  $\int_a^b f = A$ .

*Proof.* Let  $\epsilon > 0$ , then there exists  $n$  such that:

$$|L(f, P_n) - A| < \frac{1}{2}\epsilon \implies L(f, P_n) > A - \frac{1}{2}\epsilon \quad (11.1.35)$$

and similarly:

$$U(f, P_n) < A + \frac{1}{2}\epsilon \quad (11.1.36)$$

Also, by definition we must have that:

$$L(f, P_n) \leq \underline{\int_a^b f} \leq \overline{\int_a^b f} \leq U(f, P_n) \quad (11.1.37)$$

so that:

$$A - \frac{1}{2}\epsilon \leq \underline{\int_a^b f} \leq \overline{\int_a^b f} \leq A + \frac{1}{2}\epsilon \quad (11.1.38)$$

Therefore, we find that  $f$  is integrable on  $[a, b]$  with:

$$\int_a^b f = A \quad (11.1.39)$$

as desired. ■

**Proposition (Integrability)**

A function  $f$  which is

- (a) bounded and monotonic on  $[a, b]$  or
- (b) continuous on  $[a, b]$

is integrable on  $[a, b]$ .

*Proof.* (a) Consider the following partition of  $[a, b]$  into equal-sized intervals:

$$P_n = \{[a, x_1], [x_1, x_2], \dots, [x_{n-1}, n]\} \quad (11.1.40)$$

with  $x_i = a + \frac{b-a}{n}i$  and thus  $\delta x_i = \frac{b-a}{n}$ . Since  $f$  is increasing, we must have that  $m_i = f(x_{i-1})$  and  $M_i = f(x_i)$ . Hence:

$$U(f, P_n) - L(f, P_n) = \sum_i (f(x_i) - f(x_{i-1})) \frac{b-a}{n} = (f(b) - f(a)) \frac{b-a}{n} \quad (11.1.41)$$

This sequence is clearly null, so  $f$  is integrable on  $[a, b]$ . ■

**Proposition (Properties of Riemann integral)**

Let  $f$  be integrable on an interval containing  $a, b, c$ , then:

$$\int_a^c f + \int_c^b f = \int_a^b f \quad (11.1.42)$$

Suppose that  $f$  is integrable on  $[a, b]$ , then  $|f|$  is also integrable on  $[a, b]$ . Also:

$$\int_a^b (f + g) = \int_a^b f + \int_a^b g \quad (11.1.43)$$

and

$$\int_a^b \lambda f = \lambda \int_a^b f \quad (11.1.44)$$

Finally, both  $fg$  and  $f/g$  are integrable, provided that  $\frac{1}{g}$  is bounded on  $[a, b]$  in the latter case.

## 11.2 Inequalities and series with integrals

### Series

**Proposition (Inequality rules)**

Let  $f$  and  $g$  be integrable over  $[a, b]$ . then:

(i) if  $f(x) \leq g(x)$ ,  $\forall x \in [a, b]$  then:

$$\int_a^b f \leq \int_a^b g \quad (11.2.1)$$

(ii) if  $m \leq f(x) \leq M$ ,  $\forall x \in [a, b]$  then:

$$m(b-a) \leq \int_a^b f \leq M(b-a) \quad (11.2.2)$$

(iii)

$$\left| \int_a^b f \right| \leq \int_a^b |f| \quad (11.2.3)$$

*Proof.* (i) Suppose that  $f(x) \leq g(x)$ ,  $\forall x \in [a, b]$ , and let  $P$  be any partition of  $[a, b]$ . Then:

$$\inf_{[x_i, x_{i+1}]} f \leq \inf_{[x_i, x_{i+1}]} g \quad (11.2.4)$$

implying that:

$$\int_a^b f = \sup_P L(f, P) \leq \sup_P L(g, P) = \int_a^b g \quad (11.2.5)$$

as desired.

(ii) Suppose  $m \leq f(x) \leq M$  over  $[a, b]$ . Then from the results of part (a)

$$\int_a^b m \leq \int_a^b f \leq \int_a^b M \implies m(b-a) \leq \int_a^b f \leq M(b-a) \quad (11.2.6)$$

(iii) Note that  $-f(x) \leq |f(x)| \leq f(x)$  for all  $x \in [a, b]$ , so that:

$$-\int_a^b |f| \leq \int_a^b f \leq \int_a^b |f| \implies \left| \int_a^b f \right| \leq \int_a^b |f| \quad (11.2.7)$$

as desired. ■

### Example.

(a) Let us prove that:

$$\int_1^3 x \sin \frac{1}{x^{10}} dx \leq 4 \quad (11.2.8)$$

We have that for all  $x \in [1, 3]$ :

$$-1 \leq \sin \frac{1}{x^{10}} \leq 1 \quad (11.2.9)$$

implying that:

$$\int_1^3 x \sin \frac{1}{x^{10}} dx \leq \int_1^3 x = \frac{1}{2}[x^2]_1^3 = 4 \quad (11.2.10)$$

as desired.

(b) Let us prove that:

$$\frac{1}{2} \leq \int_0^{\frac{1}{2}} e^{x^2} dx \leq \frac{1}{2} e^{1/4} \quad (11.2.11)$$

Indeed, note that:

$$\frac{d}{dx}(e^{x^2}) = 2xe^{x^2} \geq 0, \forall x \in \left[0, \frac{1}{2}\right] \quad (11.2.12)$$

showing that  $e^{x^2}$  is increasing on  $[0, \frac{1}{2}]$ . It then follows that:

$$1 \leq e^{x^2} \leq e^{\frac{1}{4}}, \forall x \in \left[0, \frac{1}{2}\right] \quad (11.2.13)$$

and hence:

$$\frac{1}{2} \leq \int_0^{\frac{1}{2}} e^{x^2} dx \leq \frac{1}{2} e^{1/4} \quad (11.2.14)$$

as desired.

(c) Finally, let us prove that:

$$\left| \int_0^{\frac{\pi}{4}} \frac{\tan x}{3 - \sin x^2} dx \right| \leq \frac{1}{4} \log 2 \quad (11.2.15)$$

Indeed, we have that:

$$-1 \leq \sin(x^2) \leq 1 \implies 2 \leq 3 - \sin x^2 \leq 4 \quad (11.2.16)$$

Thus:

$$\frac{1}{4} \tan x \leq \frac{\tan x}{3 - \sin x^2} \leq \frac{1}{2} \tan x \quad (11.2.17)$$

implying that:

$$\left| \frac{\tan x}{3 - \sin x^2} \right| \leq \frac{1}{2} |\tan x| \quad (11.2.18)$$

Also  $\tan x \geq 0, \forall x \in [0, \frac{\pi}{4}]$ , so that:

$$\left| \frac{\tan x}{3 - \sin x^2} \right| \leq \frac{1}{2} \tan x \quad (11.2.19)$$

Using the limit inequality:

$$\left| \int_0^{\frac{\pi}{4}} \frac{\tan x}{3 - \sin x^2} dx \right| \leq \int_0^{\frac{\pi}{4}} \frac{1}{2} \tan x dx = \frac{1}{2} [\ln |\sec x|]_0^{\frac{\pi}{4}} \quad (11.2.20)$$

$$= \frac{1}{2} \ln \left| \sqrt{2} \right| = \frac{1}{4} \ln 2 \quad (11.2.21)$$

as desired.

◀

### Wallis' formula

**Lemma.** Let

$$I_n = \int_0^{\frac{\pi}{2}} \sin^n x dx, n = 0, 1, 2, \dots \quad (11.2.22)$$

Then  $I_n = \frac{n-1}{n} I_{n-2}$  for  $n \geq 2$ .

*Proof.* It is easy to see that:

$$I_0 = \frac{\pi}{2}, I_1 = 1 \quad (11.2.23)$$

Also, for  $n \geq 2$ :

$$I_n = \int_0^{\pi/2} \sin^n x dx = \int_0^{\pi/2} \sin x \sin^{n-1} x dx \quad (11.2.24)$$

$$= [-\cos x \sin^{n-1} x]_0^{\pi/2} + \int_0^{\pi/2} (n-1) \cos^2 x \sin^{n-2} x dx \quad (11.2.25)$$

$$= (n-1) \left( \int_0^{\pi/2} \sin^{n-2} x dx - \int_0^{\pi/2} \sin^n x dx \right) \quad (11.2.26)$$

$$\implies nI_n = (n-1)I_{n-2} \implies I_n = \frac{n-1}{n} I_{n-2} \quad (11.2.27)$$

as desired.

■

For example, we have that even powers of  $\sin x$  integrate to:

$$I_4 = \frac{3}{4} \frac{1}{2} \frac{\pi}{2}, I_6 = \frac{5}{6} \frac{3}{4} \frac{1}{2} \frac{\pi}{2} \quad (11.2.28)$$

or more generally:

$$I_{2n} = \frac{2n-1}{2n} \frac{2n-3}{2n-2} \cdots \frac{3}{4} \frac{1}{2} \frac{\pi}{2} \quad (11.2.29)$$

Similarly, we have that odd powers of  $\sin x$  integrate to:

$$I_5 = \frac{4}{5} \frac{2}{3}, \quad I_7 = \frac{6}{7} \frac{4}{5} \frac{2}{3} \quad (11.2.30)$$

or more generally:

$$I_{2n+1} = \frac{2n}{2n+1} \frac{2n-2}{2n-1} \cdots \frac{4}{5} \frac{2}{3} \quad (11.2.31)$$

**Lemma.** Let:

$$a_n = \frac{2}{1} \frac{2}{3} \frac{4}{3} \frac{4}{5} \frac{6}{5} \frac{6}{7} \cdots \frac{2n}{2n-1} \frac{2n}{2n+1} \quad (11.2.32)$$

$$b_n = \frac{(n!)^2 2^{2n}}{(2n)! \sqrt{n}} \quad (11.2.33)$$

then:

$$b_n^2 = \frac{2n+1}{n} a_n, \quad n = 1, 2, \dots \quad (11.2.34)$$

*Proof.* Firstly note that:

$$b_1^2 = \frac{2^4}{4} = 4, \quad a_1 = \frac{2}{1} \frac{2}{3} = \frac{4}{3} \implies b_1^2 = 3a_1 \quad (11.2.35)$$

More generally, suppose that for some  $n$ :

$$b_n^2 = \frac{2n+1}{n} a_n \quad (11.2.36)$$

then:

$$b_{n+1}^2 = \frac{((n+1)!)^4 2^{4(n+1)}}{((2n+2)!)^2 (n+1)} \quad (11.2.37)$$

$$= \frac{(n!)^4 (n+1)^4 2^{4n} 2^4}{(2n!)^2 (2n+2)^2 (2n+1)^2 (n+1)} \quad (11.2.38)$$

$$= \frac{(n!)^4 2^{4n}}{((2n)!)^2 n} \frac{n}{n+1} \frac{(n+1)^4 2^4}{2^2 (n+1)^2 (2n+1)^2} \quad (11.2.39)$$

$$= \frac{2n+1}{n} \frac{n}{n+1} \frac{(n+1)^2 2^2}{(2n+1)^2} a_n \quad (11.2.40)$$

$$= \frac{4(n+1)}{2n+1} a_n \quad (11.2.41)$$

Now note that:

$$a_{n+1} = \frac{2n+2}{2n+1} \frac{2n+2}{2n+3} a_n \implies \frac{4(n+1)}{2n+1} a_n = \frac{2n+3}{n+1} a_{n+1} \quad (11.2.42)$$

and hence:

$$b_{n+1}^2 = \frac{2n+3}{n+1} a_{n+1} \quad (11.2.43)$$

as desired. ■

**Theorem (Wallis' Formula)** Wallis' formula:

$$\lim_{n \rightarrow \infty} \frac{(n!)^2 2^{2n}}{(2n)! \sqrt{n}} = \sqrt{\pi} \quad (11.2.44)$$

*Proof.* We begin by proving that:

$$\lim_{n \rightarrow \infty} a_n = \frac{\pi}{2} \quad (11.2.45)$$

Indeed, note that:

$$\frac{\pi}{2} a_n = \frac{I_{2n+1}}{I_{2n}} \quad (11.2.46)$$

so we need to prove that:

$$\lim_{n \rightarrow \infty} \frac{I_{2n+1}}{I_{2n}} = 1 \quad (11.2.47)$$

Furthermore, for  $x \in [0, \pi/2]$  then:

$$\sin^{2n+2} x \leq \sin^{2n+1} x \leq \sin^{2n} x \implies I_{2n+2} \leq I_{2n+1} \leq I_{2n} \quad (11.2.48)$$

using the inequality rules for integrals. Therefore:

$$\frac{I_{2n+2}}{I_{2n}} = \frac{2n+1}{2n+2} \leq \frac{I_{2n+1}}{I_{2n}} \leq 1 \quad (11.2.49)$$

Hence, using the squeeze rule:

$$\lim_{n \rightarrow \infty} \frac{I_{2n+1}}{I_{2n}} = 1 \implies \lim_{n \rightarrow \infty} a_n = \frac{\pi}{2} \quad (11.2.50)$$

as desired.

We then have that:

$$\lim_{n \rightarrow \infty} b_n^2 = \lim_{n \rightarrow \infty} \frac{2n+1}{n} a_n = \pi \implies \lim_{n \rightarrow \infty} b_n = \sqrt{\pi} \quad (11.2.51)$$

■

## 11.3 Series

### Theorem (Integral test)

Let  $f$  be positive and decreasing on  $[1, \infty)$ , and suppose  $\lim_{x \rightarrow \infty} f(x) = 0$ . Define:

$$I_n = \int_1^n f \quad (11.3.1)$$

Then:

- (i)  $\sum_{n=1}^{\infty} f(n)$  converges if  $(I_n)$  is bounded above.
- (ii)  $\sum_{n=1}^{\infty} f(n)$  diverges if  $(I_n)$  diverges.

*Proof.* Let  $s_n = \sum_{k=1}^n f(k)$  be the  $n$ th partial sum of  $\sum_{n=1}^{\infty} f(n)$ , and let  $P_{n-1}$  be the partition of  $[1, n]$ :

$$P_n = \{[1, 2], \dots, [i, i+1], \dots, [n-1, n]\} \quad (11.3.2)$$

Now since  $f(x)$  is decreasing on  $[0, \infty)$ , we must have that:

$$m_i = f(i+1) \implies L(f, P_{n-1}) = f(2) + f(3) + \dots + f(n) = s_n - f(1) \quad (11.3.3)$$

$$M_i = f(i) \implies U(f, P_{n-1}) = f(1) + f(2) + \dots + f(n-1) = s_n - f(n) \quad (11.3.4)$$

Hence:

$$s_n - f(1) \leq \int_1^n f \leq s_n - f(n) \quad (11.3.5)$$

(i) Suppose  $I_n = \int_1^n f$  is bounded above by some  $M$ :

$$s_n - f(1) \leq I_n \leq M \implies s_n \leq M + f(1) \quad (11.3.6)$$

Since  $(s_n)$  is an increasing bounded sequence, it follows from the Monotone convergence theorem that  $s_n$  converges. Hence  $\sum_{n=1}^{\infty} f(n)$  converges.

(b) Suppose  $I_n$  is not bounded above, since  $f$  is positive:

$$I_{n+1} - I_n = \int_n^{n+1} f \geq 0 \implies I_n \text{ is increasing} \quad (11.3.7)$$

So we have that  $I_n$  diverges. Note also that:

$$s_n \geq I_n \quad (11.3.8)$$

so using the Squeeze rule,  $s_n$  diverges, and hence so does  $\sum_{n=1}^{\infty} f(n)$ . ■

**Example.** Consider:

$$\int \frac{dx}{x(\log x)^2} = \int \frac{du}{u^2} = -\frac{1}{\log x} \quad (11.3.9)$$

where we used  $u = \log x$ ,  $du = \frac{1}{x}$ . Hence:

$$I_n = \int_2^n \frac{dx}{x(\log x)^2} = \left[ \frac{1}{\log x} \right]_n^2 = \frac{1}{\log 2} - \frac{1}{\log n} \quad (11.3.10)$$

Note that  $x(\log x)^2$  is positive and increasing on  $[2, \infty)$ , implying that  $f = \frac{1}{x(\log x)^2}$  is positive and decreasing on  $[2, \infty)$ . Furthermore:

$$I_n = \frac{1}{\log 2} - \frac{1}{\log n} \leq \frac{1}{\log 2} \quad (11.3.11)$$

so  $I_n$  is bounded above. It follows that:

$$\sum_{n=2}^{\infty} \frac{1}{n(\log n)^2} \text{ is convergent} \quad (11.3.12)$$

# Unit F4: Power series

## 12.1 Taylor series

### Definition (*Taylor polynomial*)

Let  $f \in C^n(I)$  be a  $n$ -times differentiable function defined on an open interval  $I$  containing  $a$ . Then the Taylor polynomial of degree  $n$  at  $a$  is the polynomial:

$$T_n(x) = \sum_{k=0}^n \frac{f(k)(x)}{k!}(x-a)^k = f(a) + f'(a)(x-a) + \frac{f''(a)}{2}(x-a)^2 + \dots + \frac{f^{(n)}(a)}{n!}(x-a)^n \quad (12.1.1)$$

**Example.** Consider  $f(x) = \cos x$ , let us find its  $n$ th order Taylor polynomial about  $a \in \mathbb{R}$ . Firstly we evaluate the derivatives:

$$\begin{aligned} f(x) &= \cos x \implies f(a) = \cos a \\ f'(x) &= -\sin x \implies f'(a) = -\sin a \\ f''(x) &= -\cos x \implies f''(a) = -\cos a \\ f^{(3)}(x) &= \sin x \implies f^{(3)}(a) = \sin a \\ f^{(4)}(x) &= \cos x \implies f^{(4)}(a) = \cos a \\ &\vdots \\ f^{(2n)}(x) &= (-1)^n \cos x \implies f^{(2n)}(a) = (-1)^n \cos a \\ f^{(2n+1)}(x) &= (-1)^{n+1} \sin x \implies f^{(2n+1)}(a) = (-1)^{n+1} \sin a \end{aligned}$$

Consequently, even order taylor polynomials are:

$$T_{2n}(a) = \sum_{k=0}^n \frac{(-1)^k \cos a}{(2k)!} (x-a)^{2k} + \sum_{k=0}^{n-1} \frac{(-1)^{k+1} \sin a}{(2k+1)!} (x-a)^{2k+1} \quad (12.1.2)$$

while odd order taylor polynomials are:

$$T_{2n+1}(a) = \sum_{k=0}^n \frac{(-1)^k \cos a}{(2k)!} (x-a)^{2k} + \sum_{k=0}^n \frac{(-1)^{k+1} \sin a}{(2k+1)!} (x-a)^{2k+1} \quad (12.1.3)$$

Taking  $a = 0$  we find that:

$$T_{2n}(0) = T_{2n+1}(0) = \sum_{k=0}^n \frac{(-1)^k}{(2k)!} (x-a)^{2k} \quad (12.1.4)$$



**Example.** Consider  $f(x) = \sin x$ , let us find its  $n$ th order Taylor polynomial about  $a \in \mathbb{R}$ . Firstly we evaluate the derivatives:

$$\begin{aligned} f(x) &= \sin x \implies f(a) = \sin a \\ f'(x) &= \cos x \implies f'(a) = \cos a \\ f''(x) &= -\sin x \implies f''(a) = -\sin a \\ f^{(3)}(x) &= -\cos x \implies f^{(3)}(a) = -\cos a \\ f^{(4)}(x) &= \sin x \implies f^{(4)}(a) = \sin a \\ &\vdots \\ f^{(2n)}(x) &= (-1)^n \sin x \implies f^{(2n)}(a) = (-1)^n \sin a \\ f^{(2n+1)}(x) &= (-1)^n \cos x \implies f^{(2n+1)}(a) = (-1)^n \cos a \end{aligned}$$

Consequently, even order taylor polynomials are:

$$T_{2n+2}(a) = \sum_{k=0}^{n+1} \frac{(-1)^k \sin a}{(2k)!} (x-a)^{2k} + \sum_{k=0}^n \frac{(-1)^k \cos a}{(2k+1)!} (x-a)^{2k+1} \quad (12.1.5)$$

while odd order taylor polynomials are:

$$T_{2n+1}(a) = \sum_{k=0}^n \frac{(-1)^k \sin a}{(2k)!} (x-a)^{2k} + \sum_{k=0}^n \frac{(-1)^k \cos a}{(2k+1)!} (x-a)^{2k+1} \quad (12.1.6)$$

Taking  $a = 0$  we find that:

$$T_{2n+1}(0) = T_{2n+2}(0) = \sum_{k=0}^n \frac{(-1)^k}{(2k+1)!} (x-a)^{2k+1} \quad (12.1.7)$$



**Example.** Consider  $f(x) = e^x$ , let us find its  $n$ th order Taylor polynomial about  $a \in \mathbb{R}$ . First we evaluate the derivatives:

$$f^{(n)}(a) = e^a \quad (12.1.8)$$

so that:

$$T_n(a) = \sum_{k=1}^n \frac{e^a}{k!} (x-a)^k \implies T_n(0) = \sum_{k=1}^n \frac{x^k}{k!} \quad (12.1.9)$$



**Theorem (Taylor's theorem)**

Let  $f \in C^{n+1}(I)$  and  $a, x \in I$ . Then:

$$f(x) = T_n(x) + R_n(x) \quad (12.1.10)$$

where  $T_n(x)$  is the  $n$ th order Taylor about  $a$ , and:

$$R_n(x) = \frac{f^{(n+1)}(c)}{(n+1)!} (x-a)^{n+1} \quad (12.1.11)$$

for some  $c \in (a, x)$  is known as the error term.

*Proof.* Let's consider:

$$h(t) = f(t) - T_n(t) - A(t-a)^{n+1} \quad (12.1.12)$$

where  $A$  is chosen so that  $h(x) = 0$ . Note also that:

$$f^{(k)}(a) = T_n^{(k)}(a) \implies h^{(k)}(a) = 0, k = 0, 1, \dots, n \quad (12.1.13)$$

from which it follows that  $h$  is continuous and  $n$ -fold differentiable on  $I$ , and that  $h^{(k)}(a) = h^{(k)}(x) = 0$ . Using Rolle's theorem applied to  $h$  on the interval  $[a, x]$ , we see that there must be some  $c_1$  such that  $h'(c_1) = 0$ .

Similarly, applying Rolle's theorem to  $h'$  on the interval  $[a, c_1]$ , we see that there must be some  $c_2$  such that  $h''(c_2) = 0$ . Repeating this reasoning to  $h'', h^{(3)}, \dots, h^{(n)}$  on the intervals:

$$[a, c_2], [a, c_3], \dots, [a, c_n], c_2 > c_3 > \dots > c_n > a \quad (12.1.14)$$

we find that there must be some  $c \in [a, c_n]$  such that:

$$h^{(n+1)}(c) = f^{(n+1)}(c) - A(n+1)! = 0 \implies A = \frac{f^{(n+1)}(c)}{(n+1)!} \quad (12.1.15)$$

Substituting this into our original expression for  $h(t)$ , and setting  $t = x$  with  $h(x) = 0$ :

$$f(x) = T_n(x) + \frac{f^{(n+1)}(c)}{(n+1)!} (x-a)^{n+1} \quad (12.1.16)$$

as desired.

**Example.** The Taylor expansion of  $f(x) = \log x$  about  $a = 0$  is:

$$T_n(x) = \sum_{k=1}^n \frac{(-1)^{k+1} x^k}{k} \quad (12.1.17)$$

for  $x \in (-1, 1]$ . Let us limit the domain to  $I = [-0.02, 0.02] = [a-r, a+r]$  where  $r = 0.02$ . The second order polynomial is therefore:

$$T_2(x) = x - \frac{x^2}{2}, x \in I \quad (12.1.18)$$

implying that:

$$|R_2(x)| = \left| \frac{f^{(3)}(c)}{3!} x^3 \right| \leq |f^{(3)}(c)| \frac{r^3}{3!} \quad (12.1.19)$$

for some  $c \in I$ . Now we have that:

$$|f^{(3)}(c)| = \left| \frac{2}{(1+c)^3} \right| \leq 2 \quad (12.1.20)$$

so that:

$$|R_2(x)| \leq 2 \cdot \frac{(0.02)^3}{3!} = 2.66 \times 10^{-6} \quad (12.1.21)$$

so we know that the error will be negligible up to 5 decimal places. Consequently:

$$\log(1.02) \approx (0.02) - \frac{(0.02)^2}{2} \approx 0.01980 \text{ (5 d.p.)} \quad (12.1.22)$$

◀

**Example.** The fourth order Taylor polynomial  $T_4(x)$  at  $\pi$  for  $f(x) = \cos x$  is:

$$T_4(x) = \sum_{k=0}^4 \frac{(-1)^{k+1}}{(2k)!} (x - \pi)^{2k} = -1 + \frac{1}{2}(x - \pi)^2 \quad (12.1.23)$$

We need to show that  $T_4(\pi)$  approximates  $f(\pi) = \cos \pi$  to at least a 0.01 error on  $[3\pi/4, 5\pi/4] = [a - r, a + r]$  where  $a = \pi$ ,  $r = \frac{\pi}{4}$ .

To do so we must find an upper bound for the remainder  $|R_4(x)| = \left| \frac{f^{(5)}(c)}{5!} (x - a)^5 \right|$ . Now  $|x - a| \leq r = \frac{\pi}{4}$ , so:

$$|R_4(x)| \leq |f^{(5)}(c)| \frac{r^5}{5!} \quad (12.1.24)$$

for some  $c \in [3\pi/4, 5\pi/4]$ . We also have that:

$$f^{(5)}(x) = -\sin x \implies |f^{(5)}(c)| = |\sin c| \leq 1 \quad (12.1.25)$$

and thus:

$$|R_4(x)| \leq \frac{(\pi/4)^5}{5!} = 0.0025 < 0.01 \quad (12.1.26)$$

as desired. ◀

### Theorem (Taylor series)

Let  $f$  be a class  $C^\infty$  on an open interval  $I$  at points  $a$  and  $x$ . If  $R_n(x) \rightarrow 0$  as  $n \rightarrow \infty$  then

$$f(x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!} (x - a)^n \quad (12.1.27)$$

which is known as the **Taylor series** at  $a$  for  $f$ .

*Proof.* This follows immediately from the fact that  $f(x) = T_n(x) + R_n(x)$ . ■

**Proposition (Important taylor series at 0)**

$$\frac{1}{1-x} = \sum_{n=0}^{\infty} x^n, |x| < 1 \quad (12.1.28)$$

$$\sin x = \sum_{n=0}^{\infty} \frac{(-1)^n x^{2n+1}}{(2n+1)!}, x \in \mathbb{R} \quad (12.1.29)$$

$$\cos x = \sum_{n=0}^{\infty} \frac{(-1)^n x^{2n}}{(2n)!}, x \in \mathbb{R} \quad (12.1.30)$$

$$e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!}, x \in \mathbb{R} \quad (12.1.31)$$

$$\log(1+x) = \sum_{n=1}^{\infty} \frac{(-1)^{n+1} x^n}{n}, -1 < x \leq 1 \quad (12.1.32)$$

*Proof.* (a) It can easily be seen that the  $n$ th order taylor polynomial of  $f(x) = \frac{1}{1-x}$  is:

$$T_n(x) = \sum_{k=0}^n x^k \quad (12.1.33)$$

which is a geometric series. It converges to  $\frac{1}{1-x}$  only for  $|x| \leq 1$ , as desired.

(b) The  $n$ th order taylor polynomial of  $f(x) = \sin x$  is:

$$T_n(x) \sum_{n=0}^{\infty} \frac{(-1)^n x^{2n+1}}{(2n+1)!} \quad (12.1.34)$$

implying that the error term  $R_n(x)$  may be expressed as:

$$|R_n(x)| = \left| \frac{f^{(n+1)}(c)}{(n+1)!} x^{n+1} \right| \leq \frac{x^{n+1}}{(n+1)!} \quad (12.1.35)$$

since  $f^{(n+1)}(c) = \pm \sin c$  or  $\pm \cos c$ . Therefore:

$$|R_n(x)| \leq \frac{|x|^{n+1}}{(n+1)!} \rightarrow 0 \text{ as } n \rightarrow \infty \quad (12.1.36)$$

where we have used the squeeze rule for null sequences. This is true for all  $x$ , and the result follows.

(c) The  $n$ th order taylor polynomial of  $f(x) = \cos x$  is:

$$T_n(x) \sum_{n=0}^{\infty} \frac{(-1)^n x^{2n}}{(2n)!} \quad (12.1.37)$$

implying that the error term  $R_n(x)$  may be expressed as:

$$|R_n(x)| = \left| \frac{f^{(n+1)}(c)}{(n+1)!} x^{n+1} \right| \leq \frac{x^{n+1}}{(n+1)!} \quad (12.1.38)$$

since  $f^{(n+1)}(c) = \pm \sin c$  or  $\pm \cos c$ . Therefore:

$$|R_n(x)| \leq \frac{|x|^{n+1}}{(n+1)!} \rightarrow 0 \text{ as } n \rightarrow \infty \quad (12.1.39)$$

where we have used the squeeze rule for null sequences. This is true for all  $x$ , and the result follows.

(d) The  $n$ th order taylor polynomial of  $f(x) = e^x$  is:

$$T_n(x) = \sum_{k=0}^n \frac{x^k}{k!} \quad (12.1.40)$$

The error term may be written as:

$$R_n(x) = \frac{f^{(n+1)}(c)}{(n+1)!} x^{n+1} = e^c \frac{x^{n+1}}{(n+1)!} \quad (12.1.41)$$

and since  $c$  lies between 0 and  $x$ :

$$|R_n(x)| \leq e^x \frac{|x|^{n+1}}{(n+1)!} \rightarrow 0 \text{ as } n \rightarrow \infty \quad (12.1.42)$$

as desired.

(e) The  $n$ th order Taylor polynomial of  $f(x) = \log(1+x)$  is:

$$T_n(x) = \sum_{k=1}^n \frac{(-1)^{k+1} x^k}{k} \quad (12.1.43)$$

Consequently, we have that for  $0 < x \leq 1$  then the error term reads:

$$|R_n(x)| = \left| \frac{f^{(n+1)}(c)}{(n+1)!} x^{n+1} \right| \quad (12.1.44)$$

$$= \left| \frac{n!}{(1+c)^{n+1}} \frac{x^{n+1}}{(n+1)!} \right| \quad (12.1.45)$$

$$= \frac{|x|^{n+1}}{(n+1)(c+1)^{n+1}} \leq \frac{1}{n} \rightarrow 0 \text{ as } n \rightarrow \infty \quad (12.1.46)$$

where we noted that  $|x|^{n+1} \leq 1$  and  $1+c > 1$ . Using the squeeze rule we readily find the desired result. ■

## 12.2 Convergence

**Definition (Power series)** Let  $a, a_n, x \in \mathbb{R}$  for  $n = 0, 1, 2, \dots$ . Then a power series at  $a$  in  $x$  is a series of the form:

$$\sum_{n=0}^{\infty} a_n(x-a)^n \quad (12.2.1)$$

**Lemma.** If the power series  $\sum_{n=0}^{\infty} a_n x^n$  converges for some  $x_0 \neq 0$ , then it converges absolutely on  $(-|x_0|, |x_0|)$ .

*Proof.* Let  $r = |x_0|$ . Note that the convergence of  $\sum_{n=0}^{\infty} a_n x_0^n$  implies that  $(a_n x_0^n)$  is a null sequence, and hence there exists some  $K$  such that:

$$|a_n|r^n = |a_n x_0^n| \leq K, \quad n = 0, 1, 2, \dots \quad (12.2.2)$$

Let  $|x| < r$ , then clearly:

$$\sum_{n=0}^{\infty} a_n x^n = \sum_{n=0}^{\infty} a_n r^n \frac{x^n}{r^n} \implies |a_n x^n| \leq K \frac{|x|^n}{r^n} \quad (12.2.3)$$

Since we assumed that  $|x| < r$ , we have that the geometric series  $\sum_{n=0}^{\infty} \left(\frac{|x|}{r}\right)^n$  is convergent. By the comparison test, it follows that  $\sum_{n=0}^{\infty} a_n x^n$  converges for  $|x| < |x_0|$ . ■

### Theorem (Radius of convergence)

The power series  $\sum_{n=0}^{\infty} a_n(x-a)^n$  exactly one of the following is true:

- (a) it converges only for  $x = a$
- (b) it converges for all  $x$
- (c) there exists some  $R > 0$  such that the series diverges if  $|x-a| > R$  and converges if  $|x-a| < R$ .

and in all cases absolute convergence follows on the same intervals.

*Proof.* Let us define:

$$E = \left\{ x \in \mathbb{R} : \sum_{n=0}^{\infty} a_n(x-a)^n \text{ converges} \right\} \quad (12.2.4)$$

If  $E = \{a\}$  then we have satisfied condition (a).

If  $E$  is unbounded, then for all  $x \in \mathbb{R}$  there is some  $x_0 \in E$  obeying  $|x| < |x_0|$ . It follows that  $\sum_{n=0}^{\infty} a_n(x-a)^n$  converges absolutely for  $(-|x_0|, |x_0|)$  by the Lemma we have just proven. Since this holds for all  $x \in \mathbb{R}$  the power series must satisfy condition (b).

The only remaining set is bounded and containing some  $x_0 \neq a$ . Consequently the power series converges absolutely over  $(-|x_0|, |x_0|) \subseteq E$ . We see that the radius of convergence is  $R = \sup E$ , so that  $R > |x_0|$ .

In the case where  $|x - a| < R$ , then there exists some  $x_1 \in E$  such that  $|x - a| < x_1$ . Therefore, the series converges absolutely.

In the case where  $|x - a| > R$ , then we can find  $x_2 > R$  such that  $|x - a| > x_2$ . If  $\sum_{n=0}^{\infty} a_n(x - a)^n$  were to converge then we would find that  $\sum_{n=0}^{\infty} a_n x_2^n$  converges, a contradiction.

So here condition (c) is satisfied. ■

### Theorem (Ratio test)

Suppose that  $\sum_{n=0}^{\infty} a_n(x - a)^n$  is a power series with radius of convergence  $R$ .

- (a) If  $\left| \frac{a_{n+1}}{a_n} \right| \rightarrow \infty$  as  $n \rightarrow \infty$  then  $R = 0$ .
- (b) If  $\left| \frac{a_{n+1}}{a_n} \right| \rightarrow 0$  as  $n \rightarrow \infty$  then  $R = \infty$ .
- (c) If  $\left| \frac{a_{n+1}}{a_n} \right| \rightarrow L$  as  $n \rightarrow \infty$  then  $R = \frac{1}{L}$  provided  $L > 0$ .

*Proof.* (a) Suppose that  $\left| \frac{a_{n+1}}{a_n} \right| \rightarrow \infty$  as  $n \rightarrow \infty$ . If  $x \neq a$  then:

$$\frac{|a_{n+1}(x - a)^{n+1}|}{|a_n(x - a)^n|} = \left| \frac{a_{n+1}}{a_n} \right| |x - a| \rightarrow \infty \text{ as } n \rightarrow \infty \quad (12.2.5)$$

proving that  $\sum_{n=0}^{\infty} |a_n(x - a)^n|$  diverges. Since absolute convergence of power series follows from normal convergence, we have that  $\sum_{n=0}^{\infty} a_n(x - a)^n$  diverges. So the series only converges when  $x = a$ , giving  $R = 0$ .

(b) Suppose that  $\left| \frac{a_{n+1}}{a_n} \right| \rightarrow 0$  as  $n \rightarrow \infty$ . If  $x \neq a$  then:

$$\frac{|a_{n+1}(x - a)^{n+1}|}{|a_n(x - a)^n|} = \left| \frac{a_{n+1}}{a_n} \right| |x - a| \rightarrow 0 \text{ as } n \rightarrow \infty \quad (12.2.6)$$

proving that  $\sum_{n=0}^{\infty} |a_n(x - a)^n|$  converges. Since absolute convergence of power series follows from normal convergence, we have that  $\sum_{n=0}^{\infty} a_n(x - a)^n$  converges. So the series only converges when  $x = a$ , giving  $R = 0$ .

(c) Suppose that  $\left| \frac{a_{n+1}}{a_n} \right| \rightarrow L$  as  $n \rightarrow \infty$ . If  $x \neq a$  then:

$$\frac{|a_{n+1}(x - a)^{n+1}|}{|a_n(x - a)^n|} = \left| \frac{a_{n+1}}{a_n} \right| |x - a| \rightarrow L|x - a| \text{ as } n \rightarrow \infty \quad (12.2.7)$$

If  $|x - a| > \frac{1}{L}$  then we find that:

$$\frac{|a_{n+1}(x - a)^{n+1}|}{|a_n(x - a)^n|} \rightarrow L|x - a| > 1 \text{ as } n \rightarrow \infty \quad (12.2.8)$$

proving that  $\sum_{n=0}^{\infty} |a_n(x - a)^n|$  diverges over this interval, and hence so does  $\sum_{n=0}^{\infty} a_n(x - a)^n$ . It follows that  $R \leq \frac{1}{L}$ . If  $|x - a| < \frac{1}{L}$ , then we find that:

$$\frac{|a_{n+1}(x - a)^{n+1}|}{|a_n(x - a)^n|} \rightarrow L|x - a| < 1 \text{ as } n \rightarrow \infty \quad (12.2.9)$$

proving that  $\sum_{n=0}^{\infty} |a_n(x-a)^n|$  converges over this interval, and hence so does  $\sum_{n=0}^{\infty} a_n(x-a)^n$ . It follows that  $R \geq \frac{1}{L}$ .

Together, these results show that  $R = \frac{1}{L}$  as desired. ■

**Example.** Consider the power series (about 0):

$$\sum_{n=1}^{\infty} \frac{(n!)^2}{(2n)!} x^n \quad (12.2.10)$$

We see that  $a_n = \frac{(n!)^2}{(2n)!}$ , and hence:

$$\left| \frac{a_{n+1}}{a_n} \right| = \frac{((n+1)!)^2}{(n!)^2} \quad (12.2.11)$$

$$\frac{(2n)!}{(2n+2)!} = \frac{(n+1)^2}{(2n+1)(2n+2)} \quad (12.2.12)$$

$$= \frac{n+1}{2(2n+1)} \rightarrow \frac{1}{4} \text{ as } n \rightarrow \infty \quad (12.2.13)$$

Using the ratio test for power series we see that  $R = 4$ .

Consider the power series (about 0):

$$\sum_{n=1}^{\infty} n^n x^n \quad (12.2.14)$$

We see that  $a_n = n^n$ , and hence:

$$\left| \frac{a_{n+1}}{a_n} \right| = \frac{(n+1)^{n+1}}{n^n} = \left( 1 + \frac{1}{n} \right)^n (n+1) \rightarrow \infty \text{ as } n \rightarrow \infty \quad (12.2.15)$$

Using the ratio test for power series we see that  $R = 0$ .

Consider the power series (about 0):

$$\sum_{n=1}^{\infty} (n + 2^{-n})(x-1)^n \quad (12.2.16)$$

We see that  $a_n = (n + 2^{-n})$ , and hence:

$$\left| \frac{a_{n+1}}{a_n} \right| = \frac{n+1 + 2^{-n-1}}{n + 2^{-n}} \quad (12.2.17)$$

$$1 + \frac{2^{-n-1} - 2^{-n}}{n + 2^{-n}} + \frac{1}{n + 2^{-n}} \quad (12.2.18)$$

$$= 1 + \frac{1}{2(n2^n + 1)} + \frac{1}{n + 2^{-n}} \rightarrow 1 \text{ as } n \rightarrow \infty \quad (12.2.19)$$

Using the ratio test for power series we see that  $R = 1$ . ■

**Example.** Let's consider the following power series:

$$\sum_{n=0}^{\infty} \frac{\alpha(\alpha-1)\dots(\alpha-n+1)}{n!} x^n \quad (12.2.20)$$

where  $\alpha$  is not an integer. Then this is a power series about 0 with coefficients:

$$a_n = \frac{\alpha(\alpha-1)\dots(\alpha-n+1)}{n!} \quad (12.2.21)$$

We see that:

$$\left| \frac{a_{n+1}}{a_n} \right| = \left| \frac{\alpha(\alpha-1)\dots(\alpha-n)}{\alpha(\alpha-1)\dots(\alpha-n-1)} \frac{n!}{(n+1)!} \right| \quad (12.2.22)$$

$$= \left| \frac{\alpha-n}{n+1} \right| \rightarrow 1 \text{ as } n \rightarrow \infty \quad (12.2.23)$$

◀

## 12.3 The combination rules

### Proposition (Combination rules)

Let  $f, g$  be a function represented by a Taylor series at  $a$ :

$$f(x) = \sum_{n=0}^{\infty} a_n (x-a)^n, |x-a| < R \quad (12.3.1)$$

$$g(x) = \sum_{n=0}^{\infty} b_n (x-a)^n, |x-a| < R' \quad (12.3.2)$$

then for  $r = \min\{R, R'\}$  and  $\lambda \in \mathbb{R}$ :

$$(f+g)(x) = \sum_{n=0}^{\infty} (a_n + b_n) (x-a)^n, |x-a| < r \quad (12.3.3)$$

$$\lambda f(x) = \sum_{n=0}^{\infty} (\lambda a_n) (x-a)^n, |x-a| < R \quad (12.3.4)$$

$$(12.3.5)$$

Note that the theorem does not state that the radius of convergence is  $r = \min\{R, R'\}$ , it may be larger.

**Example.** Let us find the taylor series at 0 for  $f(x) = \cosh x$ . We use the identity:

$$f(x) = \cosh x = \frac{e^x + e^{-x}}{2} = \frac{1}{2} \sum_{n=0}^{\infty} (1 + (-1)^n) \frac{x^n}{n!} = \sum_{n=0}^{\infty} \frac{x^{2n}}{(2n)!} \quad (12.3.6)$$

with infinite radius of convergence. ◀

**Proposition (Product rule)**

Let  $f, g$  be a function represented by a Taylor series at  $a$ :

$$f(x) = \sum_{n=0}^{\infty} a_n (x-a)^n, |x-a| < R \quad (12.3.7)$$

$$g(x) = \sum_{n=0}^{\infty} b_n (x-a)^n, |x-a| < R' \quad (12.3.8)$$

then for  $r = \min\{R, R'\}$ :

$$(fg)(x) = \sum_{n=0}^{\infty} c_n (x-a)^n, |x-a| < r \quad (12.3.9)$$

where

$$c_n = \sum_{k=0}^n a_k b_{n-k} \quad (12.3.10)$$

**Example.** Let's consider the following function:

$$f(x) = (1+x) \log(1+x) \quad (12.3.11)$$

The taylor series for  $\log(1+x)$  at 0 is:

$$\log(1+x) = \sum_{n=0}^{\infty} \frac{(-1)^{n+1} x^n}{n}, -1 < x \leq 1 \quad (12.3.12)$$

implying that:

$$(1+x) \log(1+x) = \sum_{n=1}^{\infty} \frac{(-1)^{n+1}}{n} (x^n + x^{n+1}) \quad (12.3.13)$$

$$= \sum_{n=1}^{\infty} \frac{(-1)^{n+1}}{n} x^n + \sum_{n=1}^{\infty} \frac{(-1)^{n+1}}{n} x^{n+1} \quad (12.3.14)$$

$$= \sum_{n=1}^{\infty} \frac{(-1)^{n+1}}{n} x^n + \sum_{n=2}^{\infty} \frac{(-1)^n}{n-1} x^n \quad (12.3.15)$$

$$= x + \sum_{n=2}^{\infty} \left( \frac{(-1)^{n+1}}{n} + \frac{(-1)^n}{n-1} \right) x^n \quad (12.3.16)$$

$$= x + \sum_{n=2}^{\infty} (-1)^n \frac{x^n}{n(n-1)} \quad (12.3.17)$$

for  $-1 < x \leq 1$ .

Let's consider the following function:

$$f(x) = \frac{1}{(1-x)^2} \quad (12.3.18)$$

The taylor series for  $\frac{1}{1-x}$  is:

$$\frac{1}{1-x} = \sum_{n=0}^{\infty} x^n, |x| < 1 \quad (12.3.19)$$

Therefore:

$$\frac{1}{(1-x)^2} = \sum_{n=0}^{\infty} c_n x^{2n}, |x| < 1 \quad (12.3.20)$$

where

$$c_n = \sum_{k=0}^n 1 = n + 1 \implies \frac{1}{(1-x)^2} = \sum_{n=0}^{\infty} (n+1)x^n, |x| < 1 \quad (12.3.21)$$

Consider the function:

$$f(x) = \frac{1}{1+2x^2} \quad (12.3.22)$$

The taylor series for  $\frac{1}{1+x}$  is:

$$\frac{1}{1+x} = \sum_{n=0}^{\infty} (-1)^n x^n, |x| < 1 \quad (12.3.23)$$

Consequently:

$$\frac{1}{1+2x^2} = \sum_{n=0}^{\infty} (-1)^n (2x^2)^n = \sum_{n=0}^{\infty} (-2)^n x^{2n}, |2x^2| < 1 \quad (12.3.24)$$

so the range of validity for this expansion is  $2x^2 < 1 \implies |x| \leq \frac{1}{\sqrt{2}}$ .

Consider the function

$$f(x) = \frac{e^x}{(1-x)^2} \quad (12.3.25)$$

The taylor series reads:

$$\sum_{n=0}^{\infty} c_n x^n, |x| < 1 \quad (12.3.26)$$

where

$$c_n = \sum_{k=0}^n \frac{k+1}{(n-k)!} \quad (12.3.27)$$

so:

$$c_0 = 1, c_1 = 1 + 2 = 3, c_2 = \frac{1}{2} + 2 + 3 = \frac{11}{2} \quad (12.3.28)$$

◀

**Theorem (Differentiation rule)** The following taylor series:

$$f(x) = \sum_{n=0}^{\infty} a_n (x-a)^n, \quad (12.3.29)$$

$$g(x) = \sum_{n=1}^{\infty} n a_n (x-a)^{n-1} \quad (12.3.30)$$

have the same radius of convergence, and  $f'(x) = g(x)$  for  $|x - a| < R$ .

**Theorem (Integration rule)** The following taylor series:

$$f(x) = \sum_{n=0}^{\infty} a_n (x - a)^n, \quad (12.3.31)$$

$$F(x) = \sum_{n=0}^{\infty} \frac{a_n}{n+1} (x - a)^{n+1} \quad (12.3.32)$$

have the same radius of convergence  $R$ , and if  $R > 0$  then:

$$\int f(x) dx = F(x), |x - a| < R \quad (12.3.33)$$

*Proof.* The two series have the same radius of convergence by applying the differentiation rule to  $F(x)$ . The differentiation rule also implies that  $F'(x) = f(x)$  over  $|x - a| < R$ , giving the desired integral. ■

**Example.** Let's find the taylor series at 0 for  $f(x) = \tanh^{-1} x$ . We have that:

$$f'(x) = \frac{1}{1-x^2} = \sum_{n=0}^{\infty} x^{2n}, |x| < 1 \quad (12.3.34)$$

Consequently:

$$f(x) = \sum_{n=0}^{\infty} \frac{x^{2n+1}}{2n+1}, |x| < 1 \quad (12.3.35)$$

Let's find the taylor series of  $e^{-x^2}$ :

$$e^{-x^2} = \sum_{n=0}^{\infty} \frac{(-x^2)^n}{n!} = \sum_{n=0}^{\infty} \frac{(-1)^n}{n!} x^{2n}, x \in \mathbb{R} \quad (12.3.36)$$

implying that:

$$\int_0^1 e^{-x^2} = \sum_{n=0}^{\infty} \int_0^1 (-1)^n n! x^{2n} = \sum_{n=0}^{\infty} \frac{(-1)^n}{n!} \frac{1}{2n+1} \quad (12.3.37)$$

We define the error function as:

$$\text{erf}(x) = \frac{2}{\sqrt{\pi}} \sum_{n=0}^{\infty} \frac{(-1)^n x^{2n+1}}{n!} \frac{1}{2n+1} \implies \int_0^1 e^{-x^2} = \frac{\sqrt{\pi}}{2} \text{erf}(1) \quad (12.3.38)$$

Finally, let's find the taylor series of  $f(x) = \log(1 + x)$ . We know that:

$$\frac{1}{1+x} = \sum_{n=0}^{\infty} (-1)^n x^n, |x| < 1 \quad (12.3.39)$$

implying that:

$$\log(1+x) = \int \frac{1}{1+x} dx \sum_{n=0}^{\infty} \frac{(-1)^n}{n+1} x^{n+1} = \sum_{n=1}^{\infty} \frac{(-1)^{n+1}}{n} x^n \quad (12.3.40)$$

for  $|x| < 1$ . ◀

**Theorem (General binomial theorem)** Let  $\alpha \in \mathbb{R}$ , then:

$$(1+x)^\alpha = \sum_{n=0}^{\infty} \binom{\alpha}{n} x^n, |x| < 1 \quad (12.3.41)$$

*Proof.* Let us define:

$$f(x) = \sum_{n=0}^{\infty} \binom{\alpha}{n} x^n, g(x) = f(x)(1+x)^{-\alpha} \quad (12.3.42)$$

Differentiating  $g$  we find that:

$$g'(x) = f'(x)(1+x)^{-\alpha} - \alpha f(x)(1+x)^{-\alpha-1} \quad (12.3.43)$$

$$= (1+x)^{-\alpha-1} ((1+x)f'(x) - \alpha f(x)) \quad (12.3.44)$$

$$= (1+x)^{-\alpha-1} \left( (1+x) \sum_{n=1}^{\infty} \binom{\alpha}{n} n x^{n-1} - \alpha \sum_{n=0}^{\infty} \binom{\alpha}{n} x^n \right) \quad (12.3.45)$$

We can simplify the expression in brackets:

$$(1+x) \sum_{n=1}^{\infty} \binom{\alpha}{n} n x^{n-1} - \alpha \sum_{n=0}^{\infty} \binom{\alpha}{n} x^n \quad (12.3.46)$$

$$= \sum_{n=1}^{\infty} \binom{\alpha}{n} n x^{n-1} + \sum_{n=1}^{\infty} \binom{\alpha}{n} n x^n - \alpha \sum_{n=0}^{\infty} \binom{\alpha}{n} x^n \quad (12.3.47)$$

$$= \sum_{n=0}^{\infty} \binom{\alpha}{n+1} (n+1) x^n + \sum_{n=1}^{\infty} \binom{\alpha}{n} (n-\alpha) x^n \quad (12.3.48)$$

We find that:

$$\binom{\alpha}{n+1} (n+1) = \frac{\alpha!}{(n+1)!(\alpha-n-1)!} (n+1) = \frac{\alpha!}{(n)!(\alpha-n-1)!} \quad (12.3.49)$$

and

$$\binom{\alpha}{n} (n-\alpha) = \frac{\alpha!}{n!(\alpha-n)!} (n-\alpha) = -\frac{\alpha!}{n!(\alpha-n-1)!} \quad (12.3.50)$$

implying that:

$$\binom{\alpha}{n+1} (n+1) + \binom{\alpha}{n} (n-\alpha) = \frac{\alpha!}{(n)!(\alpha-n-1)!} - \frac{\alpha!}{n!(\alpha-n-1)!} = 0 \quad (12.3.51)$$

Consequently, we see that  $g'(x) = 0$ , that is,  $g(x) = f(x)(1+x)^{-\alpha}$  takes a constant value. Evaluating

$g(0)$  we get the desired result. ■

## **Part II**

# **Algebra and Group Theory**

# Unit B1: Symmetry and groups

## 13.1 Symmetry in $\mathbb{R}^2$

### Symmetries of plane figures

A special class of transformations of the plane are called **isometries**. They are transformations that preserve distances between points, and include **reflections**, **rotations**, **translation**, and **glide reflections** (reflections followed by a translation parallel to line of reflection).

Symmetries are special isometries, that maps a bounded plane figure to itself.

#### **Definition 14.1 (Symmetry of plane figures)**

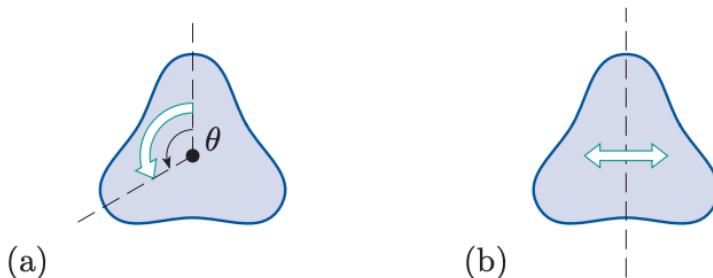
A symmetry of a plane figure  $\mathcal{F}$  is an isometry:

$$f : \mathbb{R}^2 \rightarrow \mathbb{R}^2 \quad (13.1.1)$$

$$\mathcal{F} \mapsto \mathcal{F} \quad (13.1.2)$$

For bounded plane figures, translations do not map the figure to itself, and so neither do glide-reflections. They are not symmetries, leaving only rotations and reflections. We also have the **identity transformation**, which leaves every point in  $\mathbb{R}^2$  as they are.

It is important to note that when specifying the angle through which a rotation occurs, this angle must be measured anti-clockwise by convention.



The axes of symmetry of a figure all pass through a point, called the centre of the bounded figure, is also the center of all rotational symmetries.

Consider the symmetries of a square. To keep track of its orientation and position we mark a dot on its upper left corner, and color its other side darker (this allows us to distinguish a rotation by

$\pi/2$  from a reflection about a horizontal axis of symmetry for example).

We then find that the square has four rotational symmetries including the identity symmetry) as shown below.

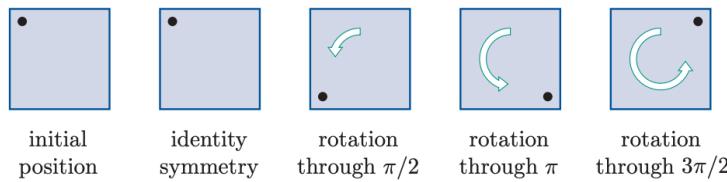


Figure 13.1. Rotational symmetry of a square

The square also has 4 reflection symmetries, 1 horizontal, 1 vertical and 2 diagonal:

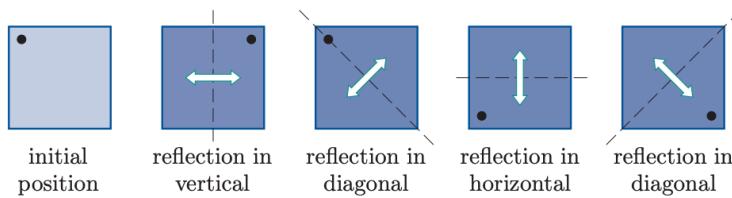


Figure 13.2. Reflection symmetries of a square

In general, it can be shown that a **regular polygon**, that is, a bounded plane figure with  $n$  equilateral sides, has  $2n$  symmetries.

### Theorem 14.2 (Symmetries of $n$ -gon)

A regular  $n$ -gon has  $2n$  symmetries, namely  $n$  rotations through  $\frac{2\pi}{k}$ ,  $k \in \mathbb{N}^*$  and  $n$  reflections. The set of all symmetries is called  $D_{2n}$ , and is called the **dihedral group**.

We shall go back to defining the dihedral group more rigorously in chapter 17.

### Identities and subsets of $S(\mathcal{F})$

**Proposition 14.3 (Properties of  $S(\mathcal{F})$ )** The set of symmetries  $S(\mathcal{F})$  of a bounded plane figure  $\mathcal{F}$  satisfies the following properties

- (i) **Closure under composition:** if  $f, g \in S(\mathcal{F})$  then  $g \circ f \in S(\mathcal{F})$
- (ii) **Associativity:** if  $f, g, h \in S(\mathcal{F})$  then  $h \circ (g \circ f) = (h \circ g) \circ f$
- (iii) **Identity existence:** for each symmetry  $f \in S(\mathcal{F})$ ,  $f \circ e = e \circ f = f$  where  $e$  is the identity symmetry.
- (iv) **Inverse existence:** for each symmetry  $f \in S(\mathcal{F})$ ,  $\exists f^{-1} \in S(\mathcal{F})$  such that  $f \circ f^{-1} = f^{-1} \circ f = e$ , where  $e$  is the identity symmetry.

Consider for example the symmetries of a square labelled below:

Then clearly, if we apply  $t \circ u$  as shown below, it is equivalent to applying  $c$ , so  $t \circ u = c \in S(\mathcal{F})$ .

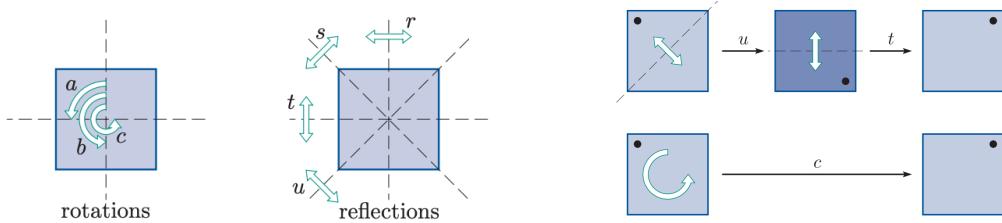


Figure 13.3. Symmetries of a square

More generally we can write that:

$\circ$	rotation	reflection
rotation	rotation	reflection
reflection	reflection	rotation

It is also interesting to note that the composition of the same reflection twice gives the identity symmetry i.e.  $r \circ r = s \circ s = t \circ t = u \circ u = e$ . We say that the reflection symmetries are **self-inverse**.

The inverses of the square symmetries can be summarised as:

Element	$e$	$a$	$b$	$c$	$r$	$s$	$t$	$u$
Inverse	$e$	$c$	$b$	$a$	$r$	$s$	$t$	$u$

#### Definition 14.4 (Direct and Indirect symmetries)

**Direct symmetries** are symmetries of a plane figure  $\mathcal{F}$  that do not require us to lift the figure out of  $\mathbb{R}^2$  and flip it. The set of direct symmetries of  $\mathcal{F}$  is denoted as  $S^+(\mathcal{F})$ .

**Indirect symmetries** are symmetries of a plane figure  $\mathcal{F}$  that require us to lift the figure out of  $\mathbb{R}^2$  and flip it.

For bounded plane figures it is immediate that rotations are direct symmetries and reflections are indirect symmetries. Therefore, for a square:

$$S^+(\square) = \{e, a, b, c\} \quad (13.1.3)$$

We then find that:

$\circ$	direct	indirect
direct	direct	indirect
indirect	indirect	direct

from which it follows that if  $f \circ f^{-1} = e$ , since  $e$  is a direct symmetry,  $f$  and its inverse must have the same nature of directness (inverse of direct is direct, inverse of indirect is indirect).

#### Theorem 14.5 (Number of direct and indirect symmetries)

If a plane figure has finite symmetries, either

- all symmetries are direct

- half the symmetries are direct and the other half indirect

*Proof.* Consider a plane figure  $\mathcal{F}$  with finite symmetries with  $n$  direct symmetries. If it has no indirect symmetries, than this case falls under category 1. Consider the case where  $\mathcal{F}$  has at least one indirect symmetry. But then, we can compose this indirect symmetry with  $n$  direct symmetries creating  $n - 1$  other indirect symmetries (since  $e$  is unique and one of the  $n$  direct symmetries, its composite with the indirect symmetry gives the same transformation, so we don't count it). So the figure also has  $n$  indirect symmetries. This algorithm is shown for  $S(\square)$ : ■

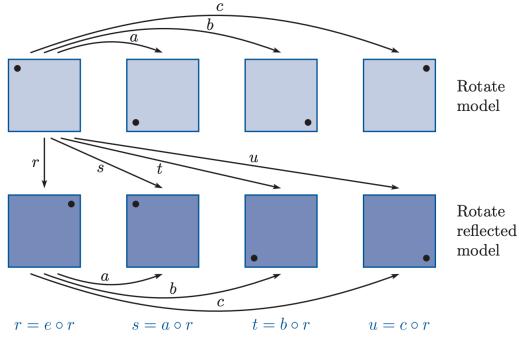


Figure 13.4. Deriving indirect symmetries from direct symmetries

An immediate consequence of this theorem is that no bounded plane figure can have solely indirect symmetry. This is easy to see applying the closure property. Indeed, if  $f, g$  are indirect symmetries then  $f \circ g$  must be a direct symmetry.

## 13.2 Representing symmetries

### Two line symbol

Because labelling each transformation by a letter may be time consuming and impractical for figures with several symmetries, we introduce the notation of **two line symbols**.

Consider for example the transformation  $r$  of the square with the vertices labelled as shown: We

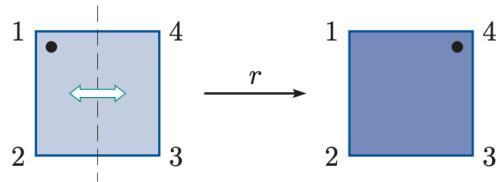


Figure 13.5. Two line symmetry notation

can then represent this transformation as:

$$r \leftrightarrow \begin{pmatrix} 1 & 2 & 3 & 4 \\ 4 & 3 & 2 & 1 \end{pmatrix} \quad (13.2.1)$$

where the top row shows the initial vertex, and the second row shows where it gets mapped.

### Definition 14.6 (*Two line symbol*)

The two line symbol representing a symmetry  $f$  of a polygon  $\mathcal{F}$  with vertices 1, 2, 3... $n$  is:

$$f \leftrightarrow \begin{pmatrix} 1 & 2 & 3 & \dots & n \\ f(1) & f(2) & f(3) & \dots & f(n) \end{pmatrix} \quad (13.2.2)$$

The inverse  $f^{-1}$  is then clearly represented by:

$$f^{-1} \leftrightarrow \begin{pmatrix} f(1) & f(2) & f(3) & \dots & f(n) \\ 1 & 2 & 3 & \dots & n \end{pmatrix} \quad (13.2.3)$$

### Cayley tables

The tables we have used so far to categorize composites of reflections and rotations, direct and indirect symmetries are called **Cayley tables**. They can be constructed by listing all the elements of  $S(\mathcal{F})$  on the top and left hand side of a square array. For any  $x, y \in S(\mathcal{F})$ , their composite  $x \circ y$  is in the  $x$ th row and  $y$ th column.

The Cayley table for the symmetries of a rectangle is shown below:

$\circ$	e	a	r	s	
e	e	a	r	s	
a	a	e	s	r	
r	r	s	e	a	
s	s	r	a	e	

Table 13.1. Rectangle symmetries and its Cayley table

Since we chose to list the direct symmetries and indirect symmetries separately, we formed four blocks each containing only direct or indirect symmetries.

The same occurs with the Cayley table for a square:

$\circ$	e	a	b	c	r	s	t	u	
e	e	a	b	c	r	s	t	u	
a	a	b	c	e	s	t	u	r	
b	b	c	e	a	t	u	r	s	
c	c	e	a	b	u	r	s	t	
r	r	u	t	s	e	c	b	a	
s	s	r	u	t	a	e	c	b	
t	t	s	r	u	b	a	e	c	
u	u	t	s	r	c	b	a	e	

$\circ$	direct	indirect
direct	direct	indirect
indirect	indirect	direct

Figure 13.6. Block patterns in Cayley table

### 13.3 Definition of a Group

#### Definition 14.7 (Binary operation)

A **binary operation**\* is a transformation mapping two members of a set  $G$  to another member of  $G$ :

$$*: G \times G \longrightarrow G \quad (13.3.1)$$

$$(f, g) \longmapsto h \quad (13.3.2)$$

where  $f, g, h \in G$ .  $G$  is therefore **closed** under  $*$ .

If we combine the set  $G$  with the binary operation  $*$ , then we get a mathematical structure known as a **group**.

#### Definition 14.8 (Group)

Let  $G$  be a set and  $*$  be a binary operation on  $G$ . Then,  $(G, *)$  is a **group** provided  $\forall f, g, h \in G$

- (i) **Associativity:**  $f * (g * h) = (f * g) * h$
- (ii) **Identity existence:**  $\exists e \in G$ , called **identity element** such that  $f * e = e * f = f$
- (iii) **Inverse existence:**  $\exists g^{-1} \in G$  called the **inverse** of  $g$  such that  $g * g^{-1} = g^{-1} * g = e$

**Example.** Let  $X = \{(a, b) \in \mathbb{R}^2 : a \neq 0\}$  and  $(a, b) * (c, d) = (ac, ad + b)$ . Show that  $(X, *)$  is a group.

- (i) **Closure:** let  $(a, b), (c, d) \in X$ , so  $a, b, c, d \in \mathbb{R}$  with  $a \neq 0$  and  $c \neq 0$ . Then:

$$(a, b) * (c, d) = (ac, ad + b) \in \mathbb{R} \quad (13.3.3)$$

since  $ac, ad + b \in \mathbb{R}$  and  $ac \neq 0$  because  $a \neq 0$  and  $c \neq 0$ .

- (ii) **Associativity:** let  $(a, b), (c, d), (e, f) \in X$  then:

$$((a, b) * (c, d)) * (e, f) = (ac, ad + b) * (e, f) = (ace, acf + ad + b) \quad (13.3.4)$$

and

$$(a, b) * ((c, d) * (e, f)) = (a, b) * (ce, cf + d) = (ace, acf + ad + b) \quad (13.3.5)$$

The two expressions are equivalent, as required, so associativity is satisfied.

- (iii) **Identity:** let  $(e_1, e_2)$  be the identity element of the group. Then, we need  $\forall (a, b) \in X$ :

$$(a, b) * (e_1, e_2) = (ae_1, ae_2 + b) = (a, b) \quad (13.3.6)$$

$$(e_1, e_2) * (a, b) = (ae_1, be_1 + e_2) = (a, b) \quad (13.3.7)$$

which upon equating terms component-wise gives:

$$\begin{cases} ae_1 = a \\ ae_2 + b = b \\ be_1 + e_2 = b \end{cases} \implies \begin{cases} e_1 = 1 \\ e_2 = 0 \end{cases} \quad (13.3.8)$$

So the identity element is  $(1, 0)$ . Let us prove this. Firstly,  $(1, 0) \in X$  since  $1 \neq 0$  and it belongs to  $\mathbb{R}^2$ . Then,  $\forall(a, b) \in X$ :

$$(a, b) * (1, 0) = (a \cdot 1, a \cdot 0 + b) = (a, b) \quad (13.3.9)$$

$$(1, 0) * (a, b) = (1 \cdot a, 1 \cdot b + 0) = (a, b) \quad (13.3.10)$$

as required,  $(1, 0)$  is the identity element of  $*$  over  $X$ .

(iv) Consider  $(a, b) \in X$ , and suppose  $(c, d)$  is its inverse. Therefore:

$$(a, b) * (c, d) = (ac, ad + b) = (1, 0) \quad (13.3.11)$$

$$(c, d) * (a, b) = (ca, cb + d) = (1, 0) \quad (13.3.12)$$

which implies:

$$\begin{cases} ac = 1 \\ ad + b = 0 \\ bc + d = 0 \end{cases} \implies \begin{cases} c = \frac{1}{a} \\ d = -\frac{b}{a} \end{cases} \quad (13.3.13)$$

where we used  $a \neq 0$  by definition since  $(a, b) \in X$ . Therefore the inverse of  $(a, b)$  is  $(\frac{1}{a}, -\frac{b}{a})$ . To prove this, consider that  $(\frac{1}{a}, -\frac{b}{a}) \in X$  because  $\frac{1}{a}, \frac{b}{a} \in \mathbb{R}$  and  $\frac{1}{a} \neq 0$  for  $a \neq 0$ . Also:

$$(a, b) * \left(\frac{1}{a}, -\frac{b}{a}\right) = \left(a \frac{1}{a}, -\frac{b}{a}a + b\right) = (1, 0) \quad (13.3.14)$$

$$(c, d) * \left(\frac{1}{a}, -\frac{b}{a}\right) = \left(a \frac{1}{a}, b \frac{1}{a} - \frac{b}{a}\right) = (1, 0) \quad (13.3.15)$$

as required. So the inverse of  $(a, b)$  is  $(\frac{1}{a}, -\frac{b}{a})$ , and every element in  $X$  has an inverse. Since all group axioms are satisfied,  $(X, *)$  is a group.  $\blacktriangleleft$

Note that the operation  $*$  need not to be commutative. In such cases we refer to the group as **Abelian**.

#### Definition 14.9 (Abelian group)

A group  $(G, *)$  where  $*$  is commutative, that is,  $\forall f, g \in G, f \circ g = g \circ f$ , is called an **Abelian group**.

#### Definition 14.10 (Finite and Infinite groups)

A group  $(G, *)$  is said to be **finite** if  $|G| = n < \infty$ .

A group is said to be **infinite** if  $G$  is an infinite set.

In Unit A1 we saw that a field  $(F, +, \times)$  is a field provided it satisfies the twelve field axioms. With our newfound knowledge of groups, we may now provide an alternative definition:

1.  $(F, +)$  is an Abelian group
2.  $(F \setminus \{0\}, \times)$  is an Abelian group
3. the distributive law holds for all  $x, y, z \in F$ , so  $x \times (y + z) = x \times y + x \times z$

### Checking group axioms with Cayley tables

Checking the group axioms 1 for a finite set can be done using a Cayley table.

Consider for example the group  $(\mathbb{Z}_4, +_4)$ , whose Cayley table is:

$+_4$	0	1	2	3
0	0	1	2	3
1	1	2	3	0
2	2	3	0	1
3	3	0	1	2

We can clearly see that every element in the body of the table belongs to  $\mathbb{Z}_4$ , so  $\mathbb{Z}_4$  is closed under  $+_4$ , and the latter is a binary operation. We also know from module A2 that  $+_4$  is associative, so the first group axiom is satisfied. The identity element is 0, because the row and column labelled 0 repeat the borders of the table. Finally, because each row and column contains the identity element 0, this means that every element of  $\mathbb{Z}_4$  has an inverse. Indeed,  $0^{-1} = 0$ ,  $1^{-1} = 3$ ,  $2^{-1} = 2$ ,  $3^{-1} = 1$ .

Hence  $(\mathbb{Z}_4, +_4)$  satisfies the group axioms, and is therefore a group.

Note that just showing that one column labelled 0 or one row labelled 0 repeats the borders of the table is not enough to show that 0 is the identity element. Both row and column must repeat the borders.

#### Proposition 14.11 (Reading $e$ from Cayley tables)

Let  $(G, *)$  be a group, then  $e$  is the identity element of the group iff both the group and column labelled  $e$  repeat the table borders.

*Proof.* In the Cayley table, the row and column labelled  $e$  contains all elements  $e * g$  and  $g * e$  respectively. So, saying that the row/column labelled  $e$  repeats the borders of the table is equivalent to saying that  $e * g = g$  and  $g * e = g$  for all  $g \in G$ . So  $e$  is the identity element of  $G$ . ■

$\circ$	... a b c ...
⋮	⋮
$e$	... a b c ...
⋮	⋮
	... $e \circ a$ $e \circ b$ $e \circ c$ ...

$\circ$	... e ...
⋮	⋮
$a$	$a$
$b$	$b$
$c$	$c$
⋮	⋮

Figure 13.7. Row and column of the identity element  $e$

**Proposition 14.12 (Reading inverses from Cayley tables)** Let  $(G, *)$  be a group, then  $h$  is an inverse of  $g$  iff  $e$  appears both in the position with row  $g$ , column  $h$  and in the position with row  $h$ , column  $g$ .

*Proof.* The element in position with row  $g$ , column  $h$  is  $g * h$  and similarly the element in position with row  $h$ , column  $g$  is  $h * g$ . Therefore, claiming that  $e$  is in both these positions is equivalent to saying  $g * h = h * g = e$ , so that  $h$  is the inverse of  $g$  as required. ■

We saw that  $\mathbb{Q}, \mathbb{R}, \mathbb{C}$  are fields under addition and multiplication so that:

$$(\mathbb{Z}, +), (\mathbb{Q}, +), (\mathbb{R}, +), (\mathbb{C}, +), (\mathbb{Q}^*, \times), (\mathbb{R}^*, \times), (\mathbb{C}^*, \times) \quad (13.3.16)$$

are all groups.

### Theorem 14.13 ( $\mathbb{Z}_n$ and $\mathbb{U}_n$ )

For  $n \geq 2$ , the set  $\mathbb{Z}_n$  is a group under  $+_n$ , and the set  $\mathbb{U}_n$  of integers in  $\mathbb{Z}_n$  coprime to  $n$  is a group under  $\times_n$ .

*Proof.* The group axioms for  $(\mathbb{Z}_n, +_n)$  hold because they are the properties of addition in  $\mathbb{Z}_n$  as explained in Unit 2.

Let us now prove that  $(\mathbb{U}_n, \times_n)$  is a group.

#### Closure

We need to prove that  $\forall a, b \in \mathbb{U}_n$ ,  $a \times_n b \in \mathbb{U}_n$ . To do so, we use the result from Unit A2 that  $a \times_n b$  is co-prime to  $n$  provided that it has a multiplicative inverse in  $\mathbb{Z}_n$ . If we denote the inverses of  $a, b$  as  $c, d$  respectively then we can write using the commutativity of  $\times_n$ :

$$(c \times_n d) \times_n (a \times_n b) = (c \times_n a) \times_n (d \times_n b) = 1 \times_n 1 = 1 \quad (13.3.17)$$

and similarly:

$$(a \times_n b) \times_n (c \times_n d) = 1 \quad (13.3.18)$$

so  $c \times_n d$  is the multiplicative inverse of  $a \times_n b$ , and the latter is therefore co-prime to  $n$ .

#### Associativity

We know that modular multiplication is associative.

#### Identity

Consider  $1 \in \mathbb{U}_n$  and  $\forall a \in \mathbb{Z}_n$ :

$$a \times_n 1 = 1 \times_n a = a \quad (13.3.19)$$

so 1 is the identity element.

#### Inverses

Let  $a \in \mathbb{U}_n$ , which implies that  $\exists b$ ,  $a \times_n b = 1$ . We have to show that  $b \in \mathbb{U}_n$ . To do so, consider:

$$a \times_n b = b \times_N a = 1 \quad (13.3.20)$$

so that  $b$  is also co-prime to  $n$ , and therefore  $b \in \mathbb{U}_n$ . Hence  $a$  has an inverse in  $\mathbb{U}_n$ .

Hence  $(\mathbb{U}_n, \times_n)$  satisfies all group axioms and is therefore a group. ■

An immediate consequence of this theorem is that  $(\mathbb{Z}_p^*, \times_p)$  is a group provided  $p$  is prime.

## 13.4 Properties of groups and group elements

### Proposition 14.14 (*Properties of groups*)

For a group  $(G, *)$  the following properties hold:

- (i) the identity element  $e$  is unique i.e.  $g * e = g$ ,  $g * e' = g \implies e = e'$
- (ii) each element has a unique inverse i.e.  $g * h = e$ ,  $g * h' = e \implies h = h'$
- (iii) the inverse of  $g^{-1}$  is  $g$
- (iv)  $(g * h)^{-1} = h^{-1} * g^{-1}$
- (v) if  $g * h = g * f$  then  $h = f$  (left cancellation law) and if  $h * g = f * g$  then  $h = f$  (right cancellation law)
- (vi)  $g^m * g^n = g^{m+n}$  for  $m, n \in \mathbb{Z}$
- (vii)  $(g^m)^n = g^{mn}$  for  $m, n \in \mathbb{Z}$ .

*Proof.*

- (i) Let  $e, e' \in (G, *)$  be both identity elements. Then  $e * e' = e$  since  $e'$  is an identity element. Similarly,  $e * e' = e'$  since  $e$  is an identity element. Therefore  $e = e'$ , and thus the identity element is unique.
- (ii) Let  $g * h = e$  and  $g * h' = e$ , then we may write  $h' * g * h = (h' * g) * h = e * h = h$  and similarly  $h' * g * h = h' * (g * h) = h' * e = h'$  so that  $h = h'$  as required.
- (iii)  $g * g^{-1} = g^{-1} * g = e$  implies that  $g$  is the inverse of  $g^{-1}$ .
- (iv) We first show that  $(g * h) * (h^{-1} * g^{-1}) = e$ :

$$(g * h) * (h^{-1} * g^{-1}) = g * (h * h^{-1}) * g^{-1} \quad (13.4.1)$$

$$= g * e * g^{-1} \quad (13.4.2)$$

$$= g * g^{-1} \quad (13.4.3)$$

$$= e \quad (13.4.4)$$

We now show that  $(h^{-1} * g^{-1}) * (g * h) = e$ :

$$(h^{-1} * g^{-1}) * (g * h) = h^{-1} * (g^{-1} * g) * h \quad (13.4.5)$$

$$= h^{-1} * e * h \quad (13.4.6)$$

$$= h^{-1} * h \quad (13.4.7)$$

$$= e \quad (13.4.8)$$

as required.

- (v) To prove the left cancellation law,  $g * h = g * f \implies g^{-1} * g * h = g^{-1} * g * f \implies (g^{-1} * g) * h = (g^{-1} * g) * f \implies e * h = e * f \implies h = f$  as required. To prove the right cancellation law  $h * g = f * g \implies h * g * g^{-1} = f * g * g^{-1} \implies h * (g * g^{-1}) = f * (g * g^{-1}) \implies h * e = f * e \implies h = f$  as required.
- (vi)  $g^m * g^n = \underbrace{g * g * \dots * g}_{m \text{ times}} * \underbrace{g * g * \dots * g}_{n \text{ times}} = \underbrace{g * g * g * \dots * g * g}_{m + n \text{ times}} = g^{m+n}$

$$(vii) (g^m)^n = \underbrace{x * x * \dots * x}_{m \text{ times}} * \dots * \underbrace{x * x * \dots * x}_{m \text{ times}} = \underbrace{g * g * g * \dots * g * g}_{mn \text{ times}} = g^{mn}$$

■

## 13.5 Symmetry in $\mathbb{R}^3$

We now adapt the study of symmetry in  $\mathbb{R}^2$  to bounded figures in  $\mathbb{R}^3$ , called **solids**. More specifically we will consider **convex polyhedra**, solids whose faces are polygons, and which have no dents or dimples, nor spikes.

### Definition 14.15 (Symmetry)

A symmetry of a figure  $\mathcal{F}$  is an isometry:

$$f : \mathbb{R}^3 \rightarrow \mathbb{R}^3 \quad (13.5.1)$$

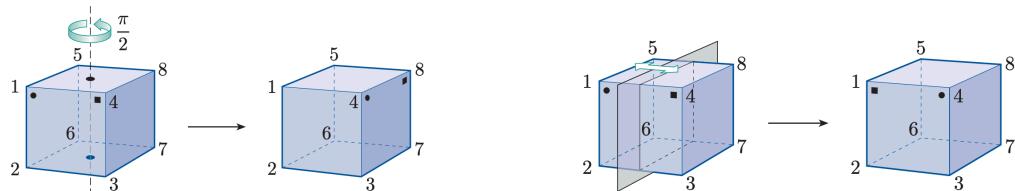
$$\mathcal{F} \mapsto f(\mathcal{F}) \quad (13.5.2)$$

Two symmetries  $f, g$  are **equal** if  $f(X) = g(X), \forall X \in \mathcal{F}$

The potential symmetries of a bounded plane figure in  $\mathbb{R}^3$  are:

1. identity transformation
2. rotation specified by an axis of symmetry, direction and angle of rotation
3. reflection in a plane
4. composite of the above

The two line symbol applies as always. For example, consider the rotation of a cube through  $\pi/2$  about its vertical axis:



(a) Rotation of a cube through  $\pi/2$  about its vertical axis      (b) Reflection of a cube in the vertical plane

Figure 13.8.

Using two-line symbols we can write:

$$g \leftrightarrow \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ 4 & 3 & 7 & 8 & 1 & 2 & 6 & 5 \end{pmatrix} \quad (13.5.3)$$

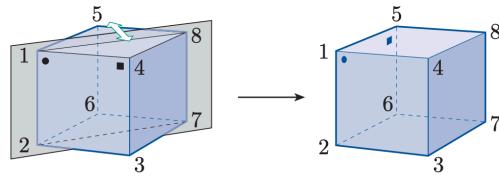
Similarly a reflection in the vertical plane is represented by:

$$f \leftrightarrow \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ 4 & 3 & 2 & 1 & 8 & 7 & 6 & 5 \end{pmatrix} \quad (13.5.4)$$

The composition of the two is performed as intuition would expect, that is, reading off  $y = g(x)$  from 14.5.3 and then  $f(y)$  from 14.5.5 for all  $x$  vertices.

$$f \circ g \leftrightarrow \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ 1 & 2 & 6 & 5 & 4 & 3 & 7 & 8 \end{pmatrix} \quad (13.5.5)$$

which is a reflection in the diagonal plane as shown below:



**Figure 13.9.** Reflection in diagonal plane of a cube

### Theorem 14.16 (Symmetry group)

$S(\mathcal{F})$  forms a group under the composition function for bounded figures in  $\mathbb{R}^n$ . The group  $(S(\mathcal{F}), \circ)$  is called the **symmetry group**.

Analogously to plane figures, the symmetries of solid figures can also be classified as direct or indirect. In this case however, it must be noted that one cannot physically demonstrate a reflection, since to do so would require accessing the fourth dimension. Therefore, we may define direct symmetries as those which can be shown physically in  $\mathbb{R}^3$ , and indirect symmetries as all the others.

### Finding the number of symmetries of a polyhedron

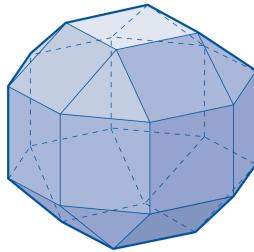
#### Theorem 14.17 (Symmetries of regular polyhedra)

For a regular polyhedron with  $F$  faces each with  $n$  symmetries, then the number of symmetries of the polyhedron is  $F \cdot n$ .

To see why this is the case, consider a tetrahedron, a polyhedron made up of 4 equilateral triangular faces. We wish to find the number of ways of replacing the figure in the space it occupied originally.

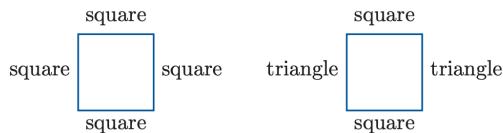
We can choose any one of the four equilateral triangles as a base, each with 6 symmetries. For each symmetry of the base, we can allow the entire tetrahedron to be transformed accordingly, resulting in a symmetry of the solid. So there are  $4 \times 6 = 24$  symmetries.

For irregular polyhedra, the process is similar. We consider for example a small rhombicuboctahedron, with 18 squares and 8 equilateral triangles as faces.

**Figure 13.10.** Rhombicuboctahedron

Again we look at all the ways of replacing the polyhedron in the space it occupied originally. To do so we first find a type of face we can use as base, such as a square face.

Now immediately we realize that not all squares of the rhombicuboctahedron can be used as a base. Indeed, we have two different types of square faces that give different symmetries. We will choose the type of square faces to the left in Figure 14.11 as our base.

**Figure 13.11.** Square faces of a rhombicuboctahedron

We then see that only 6 square faces are suitable squares of this type. Each have 8 symmetries which give symmetries of the polyhedron (this is not always the case, for some polyhedra not all symmetries of a base will correspond in symmetries of the figure). Therefore, there are 48 total symmetries.

Had we chosen the second type of square face (of which there are 12), then only 4 of the symmetries of the square would have been suitable. Indeed, rotations by  $\pi/2$  and  $3\pi/2$ , as well as the two reflections in the diagonals do not give a symmetry of the polyhedron. Hence, as was found earlier, there are 48 symmetries.

### Finding the symmetries of a polyhedron

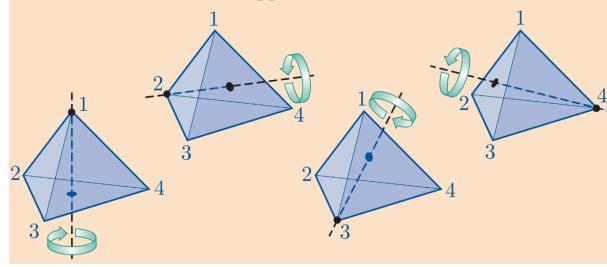
Let us try to find all the symmetries of a regular tetrahedron. Using Theorem 14.6 one can easily show that it has 12 symmetries. We start by finding all direct symmetries. As always we have the identity symmetry:

$$e = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{pmatrix} \quad (13.5.6)$$

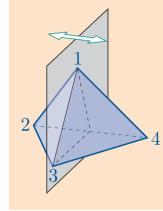
For each base, there is a rotational symmetry about a fixed axis through the opposite vertex.

One can write the two line symbol for each of these symmetries, getting to a total of 9 direct symmetries. We are therefore missing three, which can be found by composing direct symmetries with each other.

Having found all direct symmetries, we now try to find one indirect symmetry. For example, a reflection in the vertical plane as shown:



**Figure 13.12.** Rotational symmetries of the tetrahedron about axes through bases



**Figure 13.13.** Reflectional symmetry of tetrahedron

We can then compose this reflection symmetry with the 12 direct symmetries to find 12 indirect symmetries. This accounts for all 24 symmetries of the tetrahedron.

## 13.6 The Dihedral group

### Theorem 14.18 (Order of $D_{2n}$ )

The Dihedral group  $D_{2n}$ , the group of symmetries of a regular  $n$ -gon, has order  $2n$ .

*Proof.* We will consider the polygon  $\mathcal{F} \subseteq \mathbb{C}$ , with vertices at  $e^{2im\pi}n$ ,  $0 \leq m \leq n$ .

We define the following map:

$$r : \mathbb{C} \longrightarrow \mathbb{C} \quad (13.6.1)$$

$$z \mapsto z \cdot e^{\frac{2i\pi}{n}} \quad (13.6.2)$$

which is a rotation about the center of the polygon by  $\frac{2\pi}{n}$ . This is a symmetry, since it preserves distances  $\forall z, w \in \mathbb{C}$ :

$$|r(z) - r(w)| = |(z - w) \cdot e^{\frac{2i\pi}{n}}| = |z - w| \quad (13.6.3)$$

We define the reflection in the  $x$ -axis by the following map:

$$s : \mathbb{C} \longrightarrow \mathbb{C} \quad (13.6.4)$$

$$z \mapsto \bar{z} \quad (13.6.5)$$

which is again another symmetry since  $\forall z, w \in \mathbb{C}$ :

$$|s(z) - s(w)| = |\bar{z} - \bar{w}| = \sqrt{(\bar{z} - \bar{w})(z - w)} = |z - w| \quad (13.6.6)$$

We will show that:

$$D_{2n} = \underbrace{\{e, r, r^2, r^3, \dots, r^{n-1}\}}_{\text{rotations}}, \underbrace{\{s, rs, r^2s, \dots, r^{n-1}s\}}_{\text{reflections}} \quad (13.6.7)$$

so that all symmetries of  $\mathcal{F}$  are some composition of  $r$  and  $s$ .

Indeed, let  $f \in D_{2n}$  so that  $1 \in \mathcal{F} \implies f(1) \in \mathcal{F}$ , 1 gets mapped to another vertex, say  $e^{2\pi ik}n$  for some  $0 \leq k < n$ . however,  $r^k$  also maps  $1 \mapsto e^{2\pi ik}n$  so that  $g \equiv r^{-k} \circ f$  is an isometry fixing 1.

**Lemma.** The composite of two isometries over the metric space  $(\mathbb{C}, d)$  where  $d(z, w) = |z - w|$ , is an isometry over the same metric space.

*Proof.* Let  $f, g$  be two such isometries. Then:

$$d((f \circ g)(z), (f \circ g)(w)) = |(f \circ g)(z) - (f \circ g)(w)| \quad (13.6.8)$$

$$= |f(g(z)) - f(g(w))| \quad (13.6.9)$$

$$= |g(z) - g(w)| = |z - w| \quad (13.6.10)$$

$$= d(z, w) \quad (13.6.11)$$

so  $f \circ g$  is an isometry over  $(\mathbb{C}, d)$ . ■

Since  $e^{2\pi i}n$  shares an edge with 1, and since  $g$  preserves distances, we require  $g(e^{2\pi i}n)$  to also share an edge with 1.

The two possibilities are either  $e^{2\pi i}n$  or  $e^{2\pi i(n-1)}n$ .

In the first case where  $g$  fixes 1 and  $e^{2\pi i}n$ . We can repeat the same argument as before for  $g(e^{4\pi i}n)$ , which can only get mapped to itself in order to preserve distances. Suppose all vertices  $e^{2\pi ik}n$  with  $k \leq m-1$  for some  $0 < m < n$  have been mapped to themselves by  $g$ . Then,  $e^{2\pi im}n$  can only be mapped to itself or  $e^{2\pi i(m-2)}n$ . However, the latter cannot be the case, since  $g(e^{2\pi i(m-2)}n) = e^{2\pi i(m-2)}n$ , the vertex has already been "taken". Consequently,  $g(e^{2\pi im}n) = e^{2\pi i(m-2)}n$ , and by the principle of induction all vertices have been fixed. Hence,  $g = e$ , the identity transformation, implying  $f = r^k$ .

If instead  $g$  fixes 1, and  $g(e^{2\pi i}n) = e^{2\pi i(n-1)}n$ . We then have  $(s \circ g)(e^{2\pi i}n) = e^{\frac{2\pi i}{n}}$ . Also,  $(s \circ g)(1) = s(1) = 1$ , hence  $s \circ g$  fixes 1 and  $e^{\frac{2\pi i}{n}}$ . By the same argument as before then,  $s \circ g = e$ , and consequently  $f = r^k \circ s$ .

We have therefore proven that any isometry  $f$  can be expressed as either  $r^k$ , a rotation, or  $r^k \circ s$ , a reflection. In total, there are  $n$  such rotations and  $n$  such reflections<sup>1</sup>, giving  $|D_{2n}| = 2n$  as desired. ■

### Proposition 14.19 (Properties of $D_{2n}$ )

Let  $r, s \in D_{2n}$  be a rotation and reflection respectively, as defined in the previous proof.

Then:

- (i)  $sr^k s = r^{-k}$  for all  $0 \leq k < n$
- (ii)  $\text{ord}(r) = n$ ,  $\text{ord}(r^i s) = 2$ , for all  $0 \leq i < n$

<sup>1</sup>they are also all distinct. Clearly, if  $r^m = r^k$ , then 1 gets mapped to  $e^{\frac{2\pi im}{n}}$  and  $e^{\frac{2\pi ik}{n}}$ , giving  $m = k$ . In other words all the rotations  $r^i$  are distinct for all  $0 \leq i < n$ . Also,  $s \neq r^i$  for any  $i$ , since  $s$  fixes 1, whereas  $r^i$  only does so if  $r^i = e = s$ , which is a contradiction. Finally,  $r^i s \neq r^k s$  follows immediately by composing by  $s$  to the left

- (iii)  $D_{2n}$  is not abelian if  $n \geq 3$ .
- (iv) any rotation is the composition of two reflections

*Proof.*

- (i) Consider where  $r^k s$  maps the vertex  $e^{\frac{2\pi i m}{n}}$ .  $s$  sends it to  $e^{\frac{-2\pi i m}{n}}$ , followed by  $r^k$  which sends it to  $e^{\frac{2\pi i(k-m)}{n}}$ . Similarly,  $sr^{-k}$  maps  $e^{\frac{2\pi i m}{n}}$  to  $e^{\frac{2\pi i(m-k)}{n}}$ , and then to  $e^{\frac{2\pi i(k-m)}{n}}$ . Hence  $r^k s$  and  $sr^{-k}$  map  $e^{\frac{2\pi i m}{n}}$  to the same vertex, and are equivalent.  $sr^{-k} = r^k s \implies r^{-k} = sr^k s$ .
- (ii) Since  $e, r, r^2, \dots, r^{n-1}$  are all distinct, and  $r^n = e$ , it follows that  $\text{ord}(r) = n$ , since for any  $i < n$ ,  $r^i \neq r^n = e$ . Note that by definition,  $s^2(z) = s(\bar{z}) = \bar{\bar{z}} = z \implies s^2 = e$ . Also,  $(r^i s)(r^i s) = (r^i s)(sr^{-i}) = r^i s^2 r^{-i} = r^i r^{-i} = e$ , so  $\text{ord}(r^i s) = 2$ .
- (iii) We have that  $rs = sr^{-1}$ , and suppose  $rs = sr$ , then  $sr^{-1} = sr \implies r^2 = e$ , which is a contradiction since  $\text{ord}(r) = n > 2$  by assumption.
- (iv) Consider the rotation  $r^i$  for some  $0 \leq i < n$ . Then, it may be written as  $r^i = (r^{i+1}s) \circ (sr)$

■

# Unit B2: Subgroups and isomorphisms

## 14.1 Subgroups

### Definition 15.1 (Subgroup)

A **subgroup** of a group  $(G, *)$  is a group  $(H, *)$  where  $H \subseteq G$ . We write that  $(H, *) \leq (G, *)$ .

For example, the group  $S^+(\mathcal{F}, \circ)$  is a subgroup of  $S(\mathcal{F}, \circ)$  because  $S^+(\mathcal{F}) \subseteq S(\mathcal{F})$ .

Every non-empty group  $(G, *)$  has at least two subgroups,  $(e, *)$  called the **trivial subgroup** and  $(G, *)$  itself. All subgroups other than  $(H, *)$  with  $H \subsetneq H$  are called **proper subgroups**.

### Theorem 15.2 (Identity and inverses of subgroups)

Let  $(G, *)$  be a group with subgroup  $(H, *)$ . Then:

- (i) the identity element of  $(H, *)$  is the same as the identity element of  $(G, *)$
- (ii)  $\forall h \in H$ , the inverse of  $h$  in  $(G, *)$  and  $(H, *)$  is the same.

*Proof.*

- (i) Let the identity element in  $(G, *)$  and  $(H, *)$  be  $e$  and  $e_H$  respectively. We must therefore have  $e_H \circ e = e_H$  since  $e$  is the identity element of  $(G, *)$  and  $e_H \circ e = e$  since  $e_H$  is the identity element of  $(H, *)$ . It follows immediately that  $e = e_H$ .
- (ii) Let the inverses of  $h$  in  $(H, *)$  be  $x$  and  $y$ . We then have that  $h * a = h * b = e$  where  $e$  is the identity element of  $(G, *)$  and thus of  $(H, *)$ . Using the left cancellation law,  $a = b$  as required.

■

The astute reader may have noted that some of the group axioms hold for any subgroup of a group. It is therefore only necessary to prove some of the properties of a group to ascertain that it is a subgroup.

### Theorem 15.3 (Subgroup criteria)

Let  $(G, *)$  be a group with identity element  $e$  and let  $H \subseteq G$ . Then  $(H, *)$  is a subgroup of  $(G, *)$  if and only if  $\forall x, y \in H$

- (SG1)  $x * y \in H$
- (SG2)  $e_G \in H$
- (SG3)  $x_G^{-1} \in H$

where  $x_G^{-1}$  is the inverse of  $x$  in  $G$  and  $e_G$  is the identity element of  $(G, *)$ .

*Proof.* We begin by proving the  $\implies$  implication. Suppose that  $(H, *)$  is a subgroup of  $(G, *)$ , so that  $H \subseteq G$ . Then, the closure group axiom of  $H$  under  $*$  asserts that  $x * y \in H$ . Similarly, the existence of an identity element  $e_G$  was proven in Theorem 15.2 (a). Finally the existence of the inverse was proven in Theorem 15.2(b)

We now prove the  $\iff$  implication. Suppose that (SG1)-(SG3) are satisfied, we must check that the group axioms are satisfied.

**Closure** is trivial

**Associativity** since  $(G, *)$  is a group,  $\forall x, y, z \in (G, *)$  we have that  $x * (y * z) = (x * y) * z$ . Since  $H \subseteq G$  it follows that  $x, y, z \in (H, *) \implies x, y, z \in (G, H*)$  and therefore associativity holds in  $H$  as well.

**Identity** if  $e \in G$  then  $x \in H \implies x \in G$  and thus  $x * e_G = e_G * x = x$  using the identity group axiom of  $G$ .  $e_G \in H$  is trivial, since it is equivalent to (SG2).

**Inverses** if  $x \in G$  then  $x_G^{-1} \in G$  and  $x \in H$ . Thus  $x * x_G^{-1} = x_G^{-1} * x = e_G$  as required.  $x_G^{-1} \in H$  is trivial, since it is equivalent to (SG3). ■

**Example.** Show that  $(\{e, a, b, c\}, \circ)$  is a subgroup of  $(S(\square), \circ)$

We have that  $\{e, a, b, c\} \subseteq S(\square)$  and  $\circ$  is a binary operation.

The Cayley table is:

$\circ$	$e$	$a$	$b$	$c$
$e$	$e$	$a$	$b$	$c$
$a$	$a$	$b$	$c$	$e$
$b$	$b$	$c$	$e$	$a$
$c$	$c$	$e$	$a$	$b$

Closure is clearly satisfied since all elements in the table are  $\{e, a, b, c\}$ . The identity element in  $(S(\square), \circ)$  is  $e$ , which is in  $\{e, a, b, c\}$ . The elements  $e, b$  are self-inverse and  $a, c$  are inverses of each other so  $(\{e, a, b, c\}, \circ)$  contains all the inverses of its elements. ◀

### Proposition 15.4 (Integer Multiples subgroups)

The only subgroups of  $(\mathbb{Z}, +)$  are  $n\mathbb{Z} = \{kn : k \in \mathbb{Z}\}$  for  $n \in \mathbb{N}$ .

*Proof.* It is clear that  $n\mathbb{Z}$  are subgroups. Let us prove that they are the only subgroups of  $(\mathbb{Z}, +)$ .

Let  $H \leq (\mathbb{Z}, +)$ , then we must have  $0 \in H$ . If no other elements are included, then  $H = 0\mathbb{Z}$ .

If we instead have other elements, we pick the smallest positive integer  $n$  in  $H$ . Then  $H = n\mathbb{Z}$ . Otherwise, if this is not the case, then  $\exists a \in H$  such that  $n$  doesn't divide  $a$ . The division algorithm allows us to write  $a = p \cdot n + q \in H$ , with  $0 < q < n$ . By closure  $a, p \cdot n \in H$ , implying  $q \in H$ . Yet,  $q < n$  is smaller than  $n$ , the smallest element of  $H$ , giving us a contradiction. So there are no elements of  $H$  not divisible by its smallest element. Furthermore, to satisfy closure we must have all multiples of  $n$ . So,  $H = n\mathbb{Z}$ . ■

**Example.** Show that  $(A, *)$  is a subgroup of  $(X, *)$  with  $A = \{(a, b) \in X : a = 1\}$ .

(i) **Closure:** Let  $(1, b), (1, d) \in A$  then:

$$(1, b) * (1, d) = (1, d + b) \in A \quad (14.1.1)$$

since the first term is 1 and the second belongs to  $\mathbb{R}$ .

(ii) **Identity:** the identity element in  $X$  is  $(1, 0)$  which belongs to  $A$  as required.

(iii) **Inverse:** the inverse of  $(1, b) \in A$  in  $(X, *)$  is given by  $(1, -b)$ . This element belongs to  $A$  so every element in  $A$  has an inverse.

Because these subgroup properties are all satisfied,  $(A, *)$  is a subgroup of  $(X, *)$ .  $\blacktriangleleft$

### Subgroup of symmetry groups

We saw that the symmetries of a figure form a symmetry group under the composition function. We also saw how the subset of all direct symmetries of a square forms a group under composition, so that  $(S^*(\square), \circ)$  is a subgroup of  $(S(\square), \circ)$ . This is true for any figure as we shall prove now.

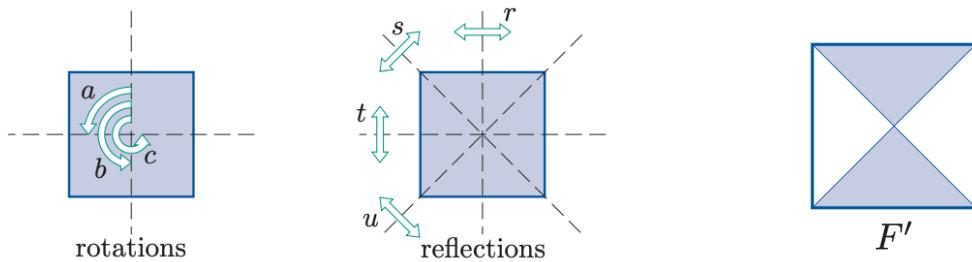
**Theorem 15.5 (Direct symmetry subgroup)** Let  $F \subsetneq \mathbb{R}^2(\mathbb{R}^3)$  be a figure. Then  $(S^+(F), \circ)$  is a subgroup of  $(S(F), \circ)$ .

*Proof.* We have  $S^+(F) \subseteq S(F)$  and  $\circ$  is the same binary operation on both sets. We now check the three subgroup properties:

- (i) **Closure:** composing any two direct symmetries gives a direct symmetry, so  $S^+(F)$  is closed.
- (ii) **Identity:** the identity element of  $(S(F), \circ)$  is  $e$ , the identity symmetry, which also belongs to  $(S^+(F), \circ)$  since it is direct.
- (iii) **Inverse:** if  $f$  is a direct symmetry, then because  $e$  is a direct symmetry  $f^{-1}$  must also be direct and hence belong to  $S^+(F)$ .

So all the subgroup properties are satisfied as required.  $\blacksquare$

Another way to produce a subgroup is to modify the figure by adding shaded patterns. For example, consider coloring the square as shown below:



**Figure 14.1.** Symmetries of  $\square$ , and the modified square  $F'$  with symmetries  $\{e, b, r, t\}$

Then clearly the symmetries of the figure are restricted. Indeed, we can see that the only symmetries are  $S(\mathcal{F}') = \{e, b, r, t\}$ , which therefore forms a subgroup under  $\circ$  of the symmetry group of a square.

Finally, a third way to find a subgroup of a symmetry group is to fix some feature of the figure (an edge or vertex usually). The resulting symmetries will still form a subgroup, as will be shown now.

**Proposition 15.6 (Fixed subset symmetries)** Let  $\mathcal{F} \subsetneq \mathbb{R}^2(\mathbb{R}^3)$ , and let  $A \subseteq \mathcal{F}$ . Then the subset of  $S(\mathcal{F})$  whose elements are all symmetries of  $\mathcal{F}$  that fix  $A$  is a subgroup under  $\circ$ .

*Proof.* Let  $H$  be a subset of  $S(\mathcal{F})$  that fixes  $A$ . Then:

- (i) **Closure:** if  $f, g \in H$ , then they both fix  $A$ . Hence, if we perform one symmetry after the other  $A$  will still remain fixed so  $f \circ g \in H$ .
- (ii) **Identity:** the identity symmetry fixes  $A$ , and so the identity element of  $S(\mathcal{F})$  belongs to  $H$ .
- (iii) **Inverses:** let  $f \in H$ , then performing the symmetry in reverse, that is,  $f^{-1}$ ,  $A$  must remain fixed. So  $H$  contains the inverse of each of its elements.

The three subgroup properties are satisfied, and thus  $H$  forms a subgroup under  $\circ$ . ■

**Example.** Find the elements of the subgroup that consists of all symmetries of the tetrahedron fixing the vertex labelled 4.

The only such symmetries are rotations about lines containing the vertex 4 and reflections about planes containing the vertex 4.

We therefore have that the only direct symmetries are rotations about the line through 4 and the centre of the opposite face, which are three (including  $e$ ).

There exists an indirect symmetry, a reflection about the plane containing 4, and the height of its opposite face. There must therefore be two other indirect symmetries using Theorem 14.5. Indeed all of the three reflections shown below (they are the only ones containing 4).

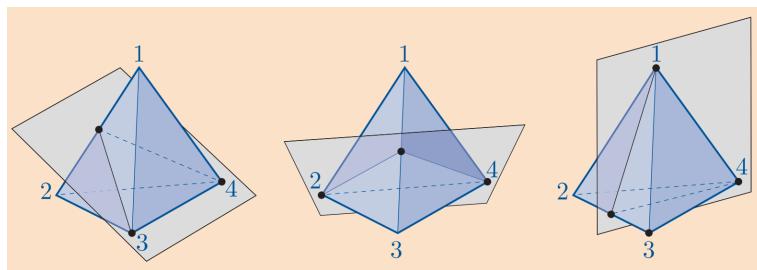


Figure 14.2. Indirect symmetries of the tetrahedron fixing a vertex

We saw in proposition 14.14 that the following properties hold for a group  $(G, *)$ :

- (i)  $g^m * g^n = g^{m+n}$
- (ii)  $(g^m)^n = g^{mn}$

Using this index notation to represent the repeated use of a binary operation works if the operation is akin to multiplication. However, it would not work well for operations using addition. For example, in the group  $(\mathbb{R}, +)$  one would not denote  $x + x + x$  as  $x^3$  but rather  $3x$ . These two types of notations are known as **multiplicative** and **additive** notation.

We can therefore write:

**Proposition 15.7 (Index laws in additive notation)**

Let  $g$  be an element of a group  $(G, *)$  then for  $m, n \in \mathbb{Z}$ :

- (i)  $mx + nx = (m + n)x$
- (ii)  $n(mx) = (nm)x$

**Definition 15.8 (Group element order)**

Let  $x$  be an element of a group  $(G, *)$ . If there exists  $n \in \mathbb{N}$  such that  $x^n = e$  then the smallest such  $n$  is the order of the element  $x$ , which has **finite order**.

If there is no such  $n$  then we say that  $x$  has **infinite order**.

For example the element  $a \in S(\square)$  which is the rotation anti-clockwise by  $\frac{\pi}{2}$  has order 4. Indeed,  $a^4 = e$  is the smallest  $n = 4$  displaying this periodicity.

Instead,  $2 \in (\mathbb{R}^*, \times)$  has infinite order, because there is no integer  $n$  such that  $2^n = 1$ .

It is also clear that the identity element has order 1. Similarly, any self inverse element  $x \neq e$  has order 2.

**Proposition 15.9 (Properties of element orders I)**

- (i) Let  $x$  be an element of a finite group  $G$ , then  $x$  has finite order.
- (ii) If  $x$  is an element of a group, then either  $x$  and  $x^{-1}$  have the same finite order, or they both have infinite order.

*Proof.*

- (i) Consider the elements  $\dots, x^{-3}, x^{-2}, x^{-1}, x^0, x, x^2, x^3 \dots$  which must belong to  $G$  by closure. Because the group is of finite order, at least one element must be repeated, or else we would have infinitely many elements. So,  $\exists s, t \in \mathbb{Z}$  with  $s < t$  such that:

$$x^s = x^t \implies x^s * (x^s)^{-1} = x^t * (x^s)^{-1} \implies e = x^{t-s} \quad (14.1.2)$$

Since  $t - s$  is positive, we have that  $x$  has finite order.

- (ii) Let  $x$  be an element of a group with identity  $e$ . First let us show that  $x^n = e \iff (x^{-1})^n = e$ . Indeed, suppose that for some  $n \in \mathbb{Z}$ :

$$x^n = e \implies (x^n)^{-1} * x^n = (x^n)^{-1} \implies (x^{-1})^n = e \quad (14.1.3)$$

so we see that  $x^{-1}$  also has the same order. To prove the converse, because the implication has been proven for any element of the group, we simply replace  $x^{-1}$  with  $x$ . So the values for which  $x^n = e$  are the same as the values for which  $(x^{-1})^n$ , so  $x$  and  $x^{-1}$  have the same order, or both have infinite order.

**Proposition 15.10 (Properties of element orders II)**

Let  $x$  be an element of  $(G, *)$  then:

- (i) if  $x$  has finite order  $n$  then:

$$e, x, x^2, \dots, x^{n-1} \quad (14.1.4)$$

are all distinct and repeat every  $n$  powers.

- (ii) If  $x$  has infinite order, then all powers of  $x$  are distinct.

*Proof.*

- (i) Suppose  $x$  has finite order  $n$ , and suppose that the powers  $e, x, x^2, \dots, x^{n-1}$  are not distinct, so that  $x^u = x^t$  for some  $0 \leq t < u \leq n - 1$ . We can then deduce that  $e = x^{u-t}$ . However, since  $0 < u - t < n$ , we see that  $x$  has order  $u - t$ , which is a contradiction. Therefore the above powers of  $x$  must all be distinct.

Now consider any integer multiple of  $n$  power:  $x^{kn}$  with  $k \in \mathbb{Z}$ . Then we have:

$$x^{kn} = (x^n)^k = e^k = e \quad (14.1.5)$$

since  $e$  has order 1. It follows that  $e$  repeats every  $n$  elements, and since all other powers are formed by composing  $x$ , it follows that all the elements listed above also repeat every  $n$  elements.

- (ii) Let  $x$  have infinite order, so that  $\nexists n \in \mathbb{N}$  such that  $x^n = e$ . Then, suppose that the powers of  $x$  are not all distinct, so that for some  $0 \leq t < u \leq n - 1$ ,  $x^u = x^t$ . But then  $x^{u-t} = e$  which implies that  $x$  has finite order, contradicting our initial assumption. ■

**Example.** Find the order of all elements in  $(\mathbb{Z}_6, +_6)$ .

The identity element 0 has order 1. For the element 1:

$$1 +_6 1 +_6 1 +_6 1 +_6 1 +_6 1 = 0 \quad (14.1.6)$$

so 1 has order 6, and  $1^{-1} = 5$  has order 6 as well.

For the element 2:

$$2 +_6 2 +_6 2 = 0 \quad (14.1.7)$$

so 2 has order 3, and  $2^{-1} = 4$  has order 3 as well.

Finally the element 3 is self-inverse and therefore has order 2. ◀

## 14.2 Cyclic groups and subgroups

We can consider the powers of a group element as forming a cycle that repeats itself after  $n$  operations.

We can see in (a) that moving around the cycle anti-clockwise is equivalent to performing  $a$  repeatedly. Moving clockwise does the opposite, so it takes  $a^{-1}$ . It follows that the element right before the identity element in a cycle of powers of  $x$  is  $x^{-1}$ .

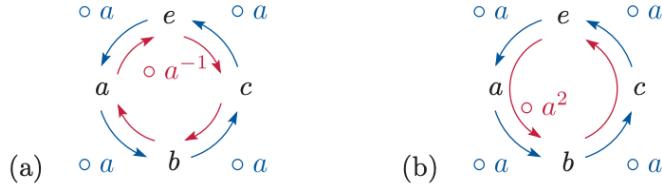


Figure 14.3. Cycle of powers of  $a$  in  $S(\square)$

In (b) we see that moving twice in the clockwise direction is the same as doing  $a \circ a = a^2 = b$  and it is then clear that  $b$  has order 2.

### Definition 15.11 (Generated subset)

Let  $x$  be an element of  $(G, *)$ , then the set of all powers of  $x$  is called the **subset of  $G$  generated by  $x$**  and denoted by  $\langle x \rangle = \{x^k : k \in \mathbb{Z}\}$  in multiplicative notation and  $\langle x \rangle = \{kx : k \in \mathbb{Z}\}$  in additive notation.

For example, the subset  $\langle a \rangle$  of  $S(\square)$  consists of the consecutive powers of  $a$  in Figure 15.3. Therefore  $\langle a \rangle = \{e, a, b, c\}$ .

### Theorem 15.12 (Cyclic subgroup)

$(\langle x \rangle, *)$  is a subgroup of  $(G, *)$  for any element  $x \in G$ . We call  $(\langle x \rangle, *)$  **cyclic subgroup of  $G$  generated by  $x$** , and  $x$  is a **generator** of  $\langle x \rangle$ .

*Proof.*

- (i) Closure: let  $g, h \in \langle x \rangle$  so that  $g = x^s$  and  $h = x^t$  for some  $s, t \in \mathbb{Z}$ . Then:

$$g * h = x^{s+t} \in \langle x \rangle \quad (14.2.1)$$

so we have closure.

- (ii) The identity element of  $(G, *)$ ,  $e$ , also belongs to  $(\langle x \rangle, *)$  since  $e = x^0$ .
- (iii) Let  $g \in \langle x \rangle$ , then  $g = x^s$ . Then  $g^{-1} = (x^s)^{-1} = x^{-s} \in \langle x \rangle$  so  $\langle x \rangle$  includes the inverses of all of its elements.

■

### Example.

Show that  $\langle x \rangle = \langle x^{-1} \rangle$ .

*Proof.*  $\langle x \rangle = \{x^k : k \in \mathbb{Z}\} = \{x^{-k} : -k \in \mathbb{Z}\} = \{x^{-k} : k \in -\mathbb{Z}\}$ . However, note that  $-\mathbb{Z} = \mathbb{Z}$  so that  $\langle x \rangle = \{x^{-k} : k \in \mathbb{Z}\} = \langle x^{-1} \rangle$  as required.

◀

Since  $\langle a \rangle = \{e, a, b, c\}$ , and  $c = a^{-1}$  we have that  $\langle c \rangle = \{e, a, b, c\}$ .

**Proposition 15.13 (Cyclic subgroup order and element order)**

If  $x$  has finite order  $n$  then the subgroup  $(\langle x \rangle, *)$  has order  $n$ .

If  $x$  has infinite order then so does the subgroup  $(\langle x \rangle, *)$ .

**Definition 15.14 (Cyclic group)**

Let  $(G, *)$  be a group with element  $x$  such that  $G = \langle x \rangle$ . Then  $(G, *)$  is called a **cyclic group**, otherwise, if there is no such  $x$  then it is **non-cyclic**.

For example, we saw that  $(\mathbb{Z}_6, +_6)$  contains two generators, 1 and 5. Note also that infinite groups can be cyclic, such as  $(\mathbb{Z}, +)$  which is generated by both 1 and  $-1$ .

**Theorem 15.15 (Abelianity of cyclic (sub)groups)**

Every cyclic group is abelian and every subgroup of a cyclic group is cyclic.

*Proof.* Let us first prove that every cyclic group is abelian. Suppose that  $(G, *)$  is a cyclic group with generator  $a$ , and let  $f, g \in G$ , so that  $f = a^s$  and  $g = a^t$ . Then:

$$f * g = a^s * a^t = a^{s+t} = a^{t+s} = a^t * a^s = g * f \quad (14.2.2)$$

so  $(G, *)$  is abelian.

Let us now prove that every subgroup  $(H, *)$  of  $(G, *)$  is cyclic. Suppose  $H = \{e\}$ , the trivial subgroup, then it is clearly cyclic (generated by  $e$ ). Suppose now that  $H$  is non-trivial, but since  $H \subseteq G$  all elements of  $H$  are powers of  $a$ . Let  $m$  be the smallest positive integer such that  $a^m \in H$  (it must exist since if  $a^m \in H$  then  $a^{-m} \in H$  and we must have at least one element in  $H$  with non-zero exponent). We will prove that  $a^m$  generates  $H$ .

Indeed, let  $h \in H \implies h = a^k$  for some  $k \in \mathbb{Z}$ . By the division theorem:  $k = qm + r$  for some  $q, r$  with  $0 \leq r < m$ . Then:

$$a^r = a^{k-qm} = a^k * (a^m)^{-q} = h * (a^m)^{-q} \quad (14.2.3)$$

but since  $H$  is a group, it must be closed under  $*$  and hence  $a^r \in H$ . But since  $m$  is the smallest positive integer such that  $a^m \in H$  and  $0 \leq r < m$ , we must have  $r = 0$  to not have a contradiction. Then  $k = qm$  and:

$$a^k = (a^m)^q \quad (14.2.4)$$

so we can generate  $H$  with  $a^m$  as required. Hence  $(H, *)$  is cyclic. ■

**Example.** Find all subgroups of  $(\mathbb{Z}_5^*, \times_5)$ .

We firstly note that  $(\mathbb{Z}_5^*, \times_5)$  is a cyclic group. Indeed looking at the generated subgroups:

$$\langle 1 \rangle = \{1\} \quad (14.2.5)$$

$$\langle 2 \rangle = \{2, 4, 3, 1\} = \mathbb{Z}_5^* \quad (14.2.6)$$

$$\langle 3 \rangle = \{3, 4, 2, 1\} = \mathbb{Z}_5^* \quad (14.2.7)$$

$$\langle 4 \rangle = \{4, 1\} \quad (14.2.8)$$

we see that  $(\mathbb{Z}_5^*, \times_5)$  is generated by 2, and is therefore cyclic. Hence, all the subgroups must be cyclic, and are included in the list above:

$$\{1\}, \{1, 4\}, \mathbb{Z}_5^* \quad (14.2.9)$$

◀

## 14.3 Cyclic groups and modular arithmetic

We have seen that the additive group  $(\mathbb{Z}_6, +_6)$  is cyclic and generated by 1 with order 6. This is true more generally for any group  $(\mathbb{Z}_n, +_n)$  with  $n \geq 2$ , which is cyclic of order  $n$ .

### Theorem 15.16 (Order of cyclic group elements)

$(\mathbb{Z}_n, +_n)$  is cyclic of order  $n$ . Any non-zero element  $m$  of the group has order  $\frac{n}{d}$  where  $d = \text{GCD}(m, n)$ .

*Proof.* We start by proving a useful lemma:

**Lemma.** Let  $m$  be a non-zero element of  $(\mathbb{Z}_n, +_n)$ , if  $m$  is a factor of  $n$  then  $m$  has order  $\frac{n}{m}$ .

Indeed, repeatedly adding  $m$  in  $(\mathbb{Z}_n, +_n)$  is the same as moving  $m$  places at a time around the cycle. Starting from 0 then and adding  $m \frac{n}{m}$  times we reach 0, so starting from  $m$  and adding  $m \frac{n}{m}$  times we reach  $m$ . Hence the order of  $m$  is  $\frac{n}{m}$ .

Now let  $m$  be a non-zero integer in  $\mathbb{Z}_n$  and let  $d = \text{GCD}(m, n)$ . Then  $\frac{m}{d}$  and  $\frac{n}{d}$  are coprime integers, and by the lemma we have proven,  $d$  has order  $\frac{n}{d}$ .

Our goal is to prove that  $\frac{n}{d} \leq \text{ord}(m) \leq \frac{n}{d}$ .

To show that  $\text{ord}(m) \leq \frac{n}{d}$ , consider the cycle of multiples of 1:

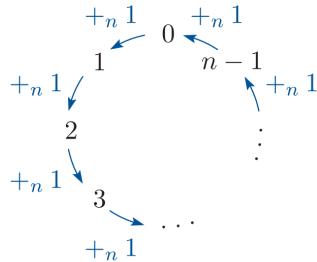


Figure 14.4. Cycle of 1 in  $\mathbb{Z}_n, +_n$

If we start from 0 and move around  $m$  places at a time  $\frac{n}{d}$  times, then we move around a total of  $\frac{mn}{d}$ . Since  $l \equiv \frac{m}{d}$  is an integer, this means that we have gone around  $l$  times, and so end up at 0. Hence,  $\text{ord}(m)$  is indeed at most  $\frac{n}{d}$ .

Now let us show that  $\frac{n}{d} \leq \text{ord}(m)$  by contradiction. Indeed, suppose that  $1 \leq \text{ord}(m) = r < \frac{n}{d}$ . Then starting from 0 and moving  $m$  places at a time  $r$  times, we end up at 0 again by definition.

Hence we must have for some  $k \in \mathbb{N}$ :

$$rm = kn \implies r \frac{m}{d} = k \frac{n}{d} \quad (14.3.1)$$

but since  $\frac{m}{d}$  and  $\frac{n}{d}$  are coprime, the only way the above equation may be true is if  $\frac{m}{d}$  is a factor of  $k \frac{n}{d}$  and hence  $\frac{m}{d}$  is a factor of  $k$ . We can then write:

$$rd = \frac{kd}{m} n \quad (14.3.2)$$

where  $\frac{kd}{m}$  is an integer since  $\frac{m}{d}$  is a factor of  $k$ . Thus  $rd$  is a multiple of  $n$ , and therefore going around the cycle  $d$  places at a time  $r$  times we end up at 0 again. In other words,  $\text{ord}(d) = r = \frac{n}{d}$ , which contradicts the assumption  $1 \leq \text{ord}(m) = r < \frac{n}{d}$ . Thus the order of  $m$  cannot be less than  $\frac{n}{d}$ , and we may conclude that:

$$\text{ord}(m) = \frac{n}{d} \quad (14.3.3)$$

as desired. ■

**Corollary 1.** For any prime  $p$ , any non-zero element  $m$  of the group  $(\mathbb{Z}_p, +_p)$  has order  $p$ .

*Proof.* Since  $p$  is prime,  $d = \text{GCD}(m, p) = 1$  and so  $\text{ord}(m) = \frac{p}{1} = p$ . ■

**Corollary 2.** A generator  $m$  of a group  $(\mathbb{Z}_n, +_n)$  must be coprime to  $n$ .

*Proof.* If  $m = 0$  then it is not a generator and it also is not coprime to  $n$ . Now assume  $m$  is non-zero. Then it is a generator iff  $\frac{n}{d} = n$  where  $d = \text{GCD}(m, n)$ . Hence  $d = 1$  and thus  $m, n$  are coprime as required. ■

This corollary also means that any element of  $(\mathbb{Z}_p, +_p)$  is a generator.

### Proposition 15.17 (Subgroups of $(\mathbb{Z}_n, +_n)$ )

The group  $(\mathbb{Z}_n, +_n)$  has exactly one cyclic subgroup of order  $q$  for each factor  $q$  of  $n$ , and no other subgroups.

- (i) the subgroup of order 1 is generated by 0
- (ii) the subgroup of order  $q$  is generated by  $d = \frac{n}{q}$

*Proof.* The cyclic subgroup of order 1 is the one generated by 0.

Let  $q$  be any factor of  $n$  that is not 1 and let  $qd = n$ . Then,  $\text{GCD}(n, d) = d$  and thus  $d$  generates a cyclic subgroup  $\langle d \rangle$  of order  $\frac{n}{d} = q$ .

Now let us prove that there are no further cyclic subgroups of order  $q$ . Let  $m \in \mathbb{Z}_n^*$  and consider  $(\langle m \rangle, *)$ . If  $d = \text{GCD}(n, m)$  then  $m \in \langle d \rangle \implies \langle m \rangle \leq \langle d \rangle$  by Theorem 15.12. However, the two subgroups must also have the same order by assumption, so they must be equal. Therefore  $\langle q \rangle$  for each factor  $q$  are all the subgroups, and they are uniquely determined. ■

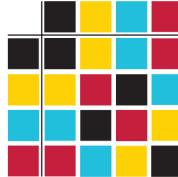
## 14.4 Isomorphisms

Let us compare two cyclic groups of order 4,  $(S^+(\square), \circ)$  and  $(\mathbb{Z}_4, +_4)$ . These two groups are structurally identical.

Indeed, to "move" from one group to another it suffices to  $e \leftrightarrow 0, a \leftrightarrow 1, b \leftrightarrow 2, c \leftrightarrow 3$  in their Cayley tables.

**Figure 14.5.** Cayley tables for  $(S^+(\square), \circ)$  and  $(\mathbb{Z}_4, +_4)$  and corresponding pattern

Indeed one could entirely remove symbols, and simply color the tiles accordingly to get a pattern. Now consider the group  $(\mathbb{Z}_5^*, \times_5)$  which is also of group 4. One can quickly see that its pattern table is:



**Figure 14.6.** Pattern table for  $(\mathbb{Z}_5^*, \times_5)$  and its relation with the pattern table for  $(\mathbb{Z}_4, +_4)$

The relation with the pattern in Figure 15.5 is then immediate. Indeed, one can switch columns and rows 3,4 to find:

$\times 5$	1	2	3	4		$\times 5$	1	2	4	3		$\times 5$	1	2	4	3
1	1	2	3	4		1	1	2	4	3		1	1	2	4	3
2	2	4	1	3	→	2	2	4	3	1	→	2	2	4	3	1
3	3	1	4	2	swap	3	3	1	2	4	swap	4	4	3	1	2
4	4	3	2	1	columns	4	4	3	1	2	rows	3	3	1	2	4
					3, 4						3, 4					
original table				intermediate table				rearranged table								

Because we have to switch columns and rows, we can claim that  $(\mathbb{Z}_4, +_4)$  and  $(\mathbb{Z}_5^*, \times_5)$  are *not* structurally identical. This concept of "moving" is important and leads to the concept of an isomorphism.

We can define a mapping  $\phi$  between  $(S^+(\square), \circ)$  and  $(\mathbb{Z}_4, +_4)$  such that:

$$\phi : (S^+(\square), \circ) \longrightarrow (\mathbb{Z}_4, +_4) \quad (14.4.1)$$

$$\{e, a, b, c\} \mapsto \{1, 2, 3, 4\} \quad (14.4.2)$$

We see that this type of mapping that transforms one Cayley table to another must be bijective, it must map every element of the first group to exactly one element of the second group.

However, we must have another property, we cannot simply map the elements randomly. Indeed if we used the rule  $\{e, a, b, c\} \mapsto \{2, 3, 4, 1\}$  then we would find an entirely different Cayley table shown in the figure below.

$\circ$	$e$	$a$	$b$	$c$		$2$	$3$	$4$	$1$	
$e$	$e$	$a$	$b$	$c$	$\longrightarrow$	$2$	$2$	$3$	$4$	$1$
$a$	$a$	$b$	$c$	$e$		$3$	$3$	$4$	$1$	$2$
$b$	$b$	$c$	$e$	$a$		$4$	$4$	$1$	$2$	$3$
$c$	$c$	$e$	$a$	$b$		$1$	$1$	$2$	$3$	$4$

$(S^+(\square), \circ)$

**Figure 14.7.** Cayley table using  $\{e, a, b, c\} \mapsto \{2, 3, 4, 1\}$

This table is clearly not the correct Cayley table for  $(\mathbb{Z}_4, +_4)$ . Although it has the same structure, the relations between different elements is no longer satisfied (e.g.  $1 +_4 3 \neq 2$ ).

So to have a coherent mapping we must have not only the border elements of the table mapped correctly, but also the body elements. In other words, for any two given elements we must have  $\phi(x * y) = \phi(x) * \phi(y)$ ,  $\forall x, y \in G$ .

$$\begin{array}{c|cccc}
 \circ & \cdots & y & \cdots \\
 \hline
 \vdots & & \vdots & \\
 x & \cdots & x \circ y & \cdots & \longrightarrow & \phi(x) & \cdots & \phi(y) & \cdots \\
 \vdots & & \vdots & \\
 \end{array}$$

**Definition 15.18** (*Isomorphic groups*)

Two groups  $(G, \circ)$ ,  $(H, *)$  are isomorphic if there exists a mapping  $\phi : G \rightarrow H$  called **isomorphism** such that:

- (i)  $\phi$  is bijective  
(ii)  $\forall x, y \in G \phi(x \circ y) = \phi(x) * \phi(y)$

We then write  $(G, \circ) \cong (H, *)$  to assert the isomorphic relation.

We can say that two groups belong to the same **isomorphism class** if they are isomorphic to each other.

### **Proposition 15.19** (*Order of isomorphic groups*)

If two groups are isomorphic than they either have both finite order, or they are both infinite.

**Example.** Let  $(G, \times)$  be a cyclic subgroup of  $(\mathbb{R}, \times)$  with  $G = \{2^k : k \in \mathbb{Z}\}$ . Show that:

$$\psi : G \longrightarrow \mathbb{Z} \quad (14.4.3)$$

$$2^k \mapsto k \quad (14.4.4)$$

We firstly show that  $\phi$  is bijective, that is, both injective and surjective.

Suppose that  $\phi(2^j) = \phi(2^k)$  for  $j, k \in \mathbb{Z}$ . Then  $j = k$  and hence  $\psi$  is injective.

Now  $\text{Im}(\phi) = \{k \in \mathbb{Z} : 2^k \in G\} = \mathbb{Z}$  and thus surjectivity is satisfied.

Finally, for  $2^j, 2^k \in G$ :

$$\phi(2^j \times 2^k) = \phi(2^{j+k}) = j + k = \phi(2^j) + \phi(2^k) \quad (14.4.5)$$

so  $\phi$  is indeed an isomorphism.  $\blacktriangleleft$

### Proposition 15.20 (Properties of isomorphisms)

Let  $(G, \circ)$  and  $(H, *)$  be groups with identities  $e_G, e_H$  respectively. Then for any isomorphism  $\phi : (G, \circ) \longrightarrow (H, *)$  and  $\forall g \in G$ :

- (i)  $\phi(e_G) = e_H$
- (ii)  $\phi(g^{-1}) = (\phi(g))^{-1}$  (this is actually true for any  $k$  in the exponent, but the case  $k = -1$  is very important)
- (iii)  $\text{ord}(g) = \text{ord}(\phi(g))$
- (iv) if  $(K, \circ) \leq (G, \circ)$  then  $(\phi(K), *) \leq (H, *)$
- (v) if  $(G, \circ)$  is abelian/cyclic then so is  $(H, *)$ .
- (vi)  $\phi(g^k) = (\phi(g))^k, \forall k \in \mathbb{Z}$ .

*Proof.*

- (i) since  $e_G \circ e_G = e_G$  we have  $\phi(e_G \circ e_G) = \phi(e_G) * \phi(e_G) = \phi(e_G)$ . We now rewrite  $\phi(e_G) = \phi(e_G) * e_H$  and use the left cancellation law to find that  $\psi(e_G) = e_H$ .
- (ii) since  $g \circ g^{-1} = g^{-1} \circ g = e_G$  we find  $\phi(g \circ g^{-1}) = \phi(g^{-1} \circ g) = e_H$  so  $\phi(g) * \phi(g^{-1}) = e_H$  thus proving that  $\phi(g^{-1}) = (\phi(g))^{-1}$ .
- (iii) Suppose  $\text{ord}(g) = k$ , since  $\phi$  is injective  $g^k = e_G \iff \phi(g^k) = (\phi(g))^k = \phi(e_G) = e_H$ . Hence  $\text{ord}(\phi(g)) \leq \text{ord}(g) = k$ . Since  $\phi$  is bijective, it has an inverse  $\phi^{-1}$ , so that repeating the argument of before using  $\text{ord}(\phi(g)) = l$  one finds:  $\phi^{-1}(\phi(g)^l) = (\phi^{-1}(\phi(g)))^l = g^l = e_G$ . Consequently,  $l = \text{ord}(\phi(g)) \geq \text{ord}(g)$ . Finally, we find that  $\text{ord}(\phi(g)) = \text{ord}(g)$ .

- (iv) we prove the three subgroup properties for  $\phi(K)$ :

**Closure:** let  $l_1 = \phi(k_1)$  and  $l_2 = \phi(k_2)$  for some  $k_1, k_2 \in K$ . Then:  $l_1 * l_2 = \phi(k_1) * \phi(k_2) = \phi(k_1 \circ k_2) \in \phi(K)$  since  $k_1 \circ k_2 \in K$  by the closure property of subgroups.

**Identity:**  $e_H = \phi(e_G) \in \phi(K)$

**Inverses:** let  $l = \phi(k)$  for some  $k \in K$ . Then  $l^{-1} = (\phi(k))^{-1} = \phi(k^{-1}) \in \phi(K)$  since  $k^{-1} \in K$  by the subgroup properties of  $K$ .

- (v) suppose that  $(G, \circ)$  is abelian, and let  $h_1 = \phi(g_1), h_2 = \phi(g_2)$  for some  $g_1, g_2 \in G$ . Then  $g_1 \circ g_2 = g_2 \circ g_1 \implies \phi(g_1 \circ g_2) = \phi(g_2 \circ g_1)$  and hence  $\phi(g_1) * \phi(g_2) = \phi(g_2) * \phi(g_1)$  and thus  $(\phi(K), *)$  is abelian as well.

Now suppose that  $(G, \circ)$  is cyclic and generated by  $a$ . Let  $h = \phi(g)$  for some  $g \in G$  then  $g = a^k \implies \phi(g) = \phi(a^k) = (\phi(a))^k$  so  $(H, *)$  is generated by  $\phi(a)$  and is thus cyclic as well.

- (vi) The case for  $k = 1$  is trivial. Now, suppose that for some  $k \geq 1$ ,  $\phi(g^k) = (\phi(g))^k$ . Then,  $\phi(g^{k+1}) = \phi(g \circ g^k) = \phi(g) \circ \phi(g^k) = (\phi(g))^{k+1}$ . Hence, by the principle of mathematical induction,  $\phi(g^k) = (\phi(g))^k, \forall k \in \mathbb{N}$ . The case for  $k = 0, -1$  have been proven in (i) and (ii) respectively. To prove this statement  $\forall k \in \mathbb{Z}$ , simply repeat the induction proof, but use  $g^{-1}$  instead of  $g$ , and apply property (ii). ■

Note that one can use the above identities to prove that two groups are not isomorphic:

1. if one group has more/less number of self-inverse elements
2. if one group has a different order than the other
3. if one group is abelian/cyclic, and the other is not

### **Proposition 15.21 (Isomorphisms of cyclic groups)**

Let  $(G, \circ)$  and  $(H, *)$  be finite cyclic groups of order  $n$  or infinite groups generated by  $a, b$  respectively. Then they are isomorphic through  $\phi : a^k \mapsto b^k, k = 0, 1, 2, \dots, n-1$  for finite ordered groups and  $k \in \mathbb{Z}$  for infinite ordered groups.

*Proof.* It is trivial to see that  $\phi$  is bijective. Also  $\forall a^j, a^k \in G$ :

$$\phi(a^j \circ a^k) = \phi(a^{j+k}) = b^{j+k} = b^j * b^k = \phi(a^j) * \phi(a^k) \quad (14.4.6)$$

as required. ■

## 14.5 Standard groups

### **Definition 15.22 (Cyclic group of order $n$ )**

We denote the standard, abstract **cyclic group of order  $n$**  as  $C_n$ .

**Example.** Find an isomorphism from  $(\mathbb{Z}_4, +_4)$  to  $(\mathbb{Z}_5, \times_5)$

We know that  $(\mathbb{Z}_4, +_4)$  is generated by 1 and  $(\mathbb{Z}_5, \times_5)$  by 2, so we can define:

$$\phi : \mathbb{Z}_4, +_4 \longrightarrow (\mathbb{Z}_5, \times_5) \quad (14.5.1)$$

$$0 \mapsto 1 \quad (14.5.2)$$

$$1 \mapsto 2 \quad (14.5.3)$$

$$2 \mapsto 4 \quad (14.5.4)$$

$$3 \mapsto 3 \quad (14.5.5)$$

**Definition 15.23 (Klein four-group)**

We denote the standard, abstract group of order 4 with all elements self inverse as the **Klein four-group**  $V$ .

There are several examples where the Klein-four group comes to use. For example, let us consider the Cayley tables of  $(S(\square), \circ)$  and  $(U_8, \times_8)$ :

$\circ$	$e$	$a$	$r$	$s$	$\times_8$	1	3	5	7
$e$	$e$	$a$	$r$	$s$	1	1	3	5	7
$a$	$a$	$e$	$s$	$r$	3	3	1	7	5
$r$	$r$	$s$	$e$	$a$	5	5	7	1	3
$s$	$s$	$r$	$a$	$e$	7	7	5	3	1

$(S(\square), \circ)$                                      $(U_8, \times_8)$

We can see that they both share the same structure, and therefore are isomorphic to each other.

## 14.6 Direct product of groups

**Definition 15.24 (Direct product of two groups)**

Given two groups  $(G, *_G)$  and  $(H, *_H)$ , then we may define their **direct product**, denoted as  $(G \times H, *)$ , where:

$$G \times H = \{(g, h) : g \in G, h \in H\} \quad (14.6.1)$$

is the **cartesian product** of  $G, H$ , equipped with the operation

$$(g_1, h_1) * (g_2, h_2) = (g_1 *_G g_2, h_1 *_H h_2) \quad (14.6.2)$$

for  $g_1, g_2 \in G$  and  $h_1, h_2 \in H$ .

We can easily see that this new group  $G \times H$  does indeed satisfy the group axioms.  $*$  is certainly binary. Indeed, for  $(g_1, h_1), (g_2, h_2) \in G \times H$ :

$$(g_1, h_1) * (g_2, h_2) = (g_1 *_G g_2, h_1 *_H h_2) \in G \times H \quad (14.6.3)$$

since  $G$  and  $H$  are closed under  $*_G$  and  $*_H$  respectively.

Associativity can be proven similarly.

Also, the identity element is  $(e_G, e_H)$ , since  $\forall (g, h) \in G \times H$ :

$$(e_G, e_H) * (g, h) = (e_G *_G g, e_H *_H h) = (g, h) \quad (14.6.4)$$

and

$$(g, h) * (e_G, e_H) = (g *_G e_G, h *_H e_H) = (g, h) \quad (14.6.5)$$

as desired.

Finally, the inverse of  $(g, h)$  is  $(g^{-1}, h^{-1})$ , since:

$$(g, h) * (g^{-1}, h^{-1}) = (g * g^{-1}, h *_H h^{-1}) = (e_G, e_H) \quad (14.6.6)$$

and

$$(g^{-1}, h^{-1}) * (g, h) = (g^{-1} *_G g, h^{-1} *_H h) = (e_G, e_H) \quad (14.6.7)$$

as desired. Taking the direct product of two groups can prove to be very useful. For example, suppose we have two independent figures  $\mathcal{F}$  and  $\mathcal{F}'$ . Then, the symmetries of this overall system form the group  $S(\mathcal{F}) \times S(\mathcal{F}')$ .

**Proposition 15.25 (Isomorphic group products)**

$$C_n \times C_m \cong C_{nm} \text{ iff } \text{GCD}(m, n) = 1.$$

*Proof.* Suppose that  $\text{GCD}(m, n) = 1$ , and let  $C_n = \langle a \rangle$ ,  $C_m = \langle b \rangle$ , and  $\text{ord}((a, b)) = k$  so that:

$$(a, b)^k = (a^k, b^k) = e \quad (14.6.8)$$

which can only be the case if  $k$  is a common multiple of  $m = \text{ord}(a)$  and  $n = \text{ord}(b)$ . Since  $k$  must be the minimum such integer, we require  $\text{LCM}(m, n) = k$ . Using the well known relation that:

$$\text{LCM}(m, n) = \frac{n \cdot m}{\text{GCD}(m, n)} = n \cdot m = k \quad (14.6.9)$$

hence the order of  $(a, b)$  is the product of the order of  $a$  and  $b$ .

Recall that  $\langle (a, b) \rangle \leq C_n \times C_m$ , so that  $|\langle (a, b) \rangle| = \text{ord}((a, b)) = n \cdot m$ . However,  $|C_n \times C_m| = n \cdot m$  as well, implying that  $\langle (a, b) \rangle = C_n \times C_m$ . It is also immediate that  $\langle (a, b) \rangle \cong C_{nm}$ , so that  $C_n \times C_m \cong C_{nm}$  as desired.

Now suppose that  $\text{GCD}(m, n) \neq 1$ , so that  $k \neq n \cdot m$ . Hence,  $C_{nm}$  is of order  $n \cdot m$ , whereas  $C_n \times C_m$  is of order  $k$ . Two groups cannot be isomorphic if they have different orders, and consequently  $C_{nm} \not\cong C_n \times C_m$ .

■

**Theorem 15.26 (Direct product theorem)**

Let  $H, F \leq (G, *)$ , and suppose  $\forall h \in H, f \in F, g \in G$ :

- (i)  $H \cap F = \{e\}$
- (ii)  $hf = fh$
- (iii)  $g = hf$

Then  $G \cong H \times F$ .

*Proof.* We will prove that:

$$\phi : H \times F \rightarrow G \quad (14.6.10)$$

$$(h, g) \mapsto h * g \quad (14.6.11)$$

is an isomorphism. Indeed,

$$\phi((h_1, f_1) * (h_2, f_2)) = \phi(h_1 * h_2, f_1 * f_2) \quad (14.6.12)$$

$$= h_1 * h_2 * f_1 * f_2 \quad (14.6.13)$$

$$= (h_1 * f_1) * (h_2 * f_2) \quad (14.6.14)$$

$$= \phi(h_1, f_1) * \phi(h_2, f_2) \quad (14.6.15)$$

where we used the commutativity of  $h, f$  in the third line.

Also,  $\phi$  is injective. Indeed:

$$\phi(h_1, f_1) = \phi(h_2, f_2) \implies h_1 * f_1 = h_2 * f_2 \implies h_1 * h_2^{-1} = f_2 * f_1^{-1} \quad (14.6.16)$$

Now by closure, the LHS belongs to  $H$ , and the RHS belongs to  $F$ , hence  $h_1 * h_2^{-1} \in H$  and  $h_1 * h_2^{-1} \in F$ , consequently  $h_1 * h_2^{-1} \in H \cap F \implies h_1 * h_2^{-1} = e \implies h_2 = h_1$ . Similar argument gives  $f_2 = f_1$ . Hence,  $(h_1, f_1) = (h_2, f_2)$  as desired.  $\blacksquare$

To prove surjectivity, note that if  $g \in G$ , then  $\exists h, f$  such that  $g = hf = \phi(h, f)$  by assumption.

In conclusion,  $\phi$  is a bijective homomorphism, hence an isomorphism.  $\blacksquare$

# Unit B3: Permutations

## 15.1 Permutations

### Definition 16.1 (*Permutation*)

A **permutation** of a finite set  $S$  is a bijective map  $\sigma : S \rightarrow S$ . The set of all permutations of this set is denoted  $\text{Sym}(S)$ .

We use the typical two-line notation to denote a permutation:

$$\sigma \leftrightarrow \begin{pmatrix} 1 & 2 & 3 & 4 \\ 4 & 1 & 2 & 3 \end{pmatrix} \quad (15.1.1)$$

A more convenient and effective notation is the **cycle form**. Indeed starting from 1 and looking at where each element gets mapped we find:

$$1 \xrightarrow{\sigma} 4 \xrightarrow{\sigma} 3 \xrightarrow{\sigma} 2 \xrightarrow{\sigma} 1 \quad (15.1.2)$$

or more concisely as:

$$\sigma = (1\ 4\ 3\ 2) \quad (15.1.3)$$

However it is not always possible to write a permutation in one single cycle. Indeed some permutations map only some symbols in each cycle:

$$\sigma_g = \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ 4 & 6 & 8 & 3 & 1 & 2 & 7 & 5 \end{pmatrix} \quad (15.1.4)$$

then we have three **disjoint** cycles (disjoint because every member appears only in one cycle):

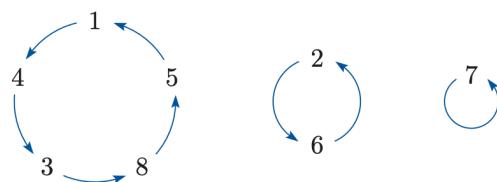


Figure 15.1. Cycles of  $\sigma_g$

We can then write:

$$\sigma_g = (1\ 4\ 3\ 8\ 5)(2\ 6)(7) \quad (15.1.5)$$

and we say that it is a product of three cycles.

### Definition 16.2 (*Permutation cycles*)

A permutation  $\sigma$  of a set  $S$  is said to be cyclic if there exist  $a_1 \dots a_k \in \{1, 2, \dots, n\}$  such that:

$$a_i\sigma = a_{i+1} \text{ for } 1 \leq i < k \quad (15.1.6)$$

$$a_k\sigma = a_1 \quad (15.1.7)$$

and is denoted in **cycle form** as  $(a_1 \ \sigma a_1 \ \sigma^2 a_1 \dots \sigma^{k-1} a_1)$ . We refer to this cycle as a  **$k$ -cycle** because it has order  $k$ . Two cycles are disjoint if they do not share any common elements.

The first line,  $a_i\sigma = a_{i+1}$  for  $1 \leq i < k$ , tells us that if some  $a_1$  gets mapped to  $a_2$  by  $\sigma$ , then  $a_2$  itself will get mapped to  $a_3$  and so forth until we reach  $a_k$ . Here, the cycle restarts, so we must have that  $a_k$  gets mapped to  $a_1$ , which is encapsulated in the second line.

For example,  $\sigma = (1\ 4\ 3\ 8\ 5)$  is a cycle. Indeed, defining  $a_1 = 1$  then:

$$a_1\sigma = 4 = a_2, \quad (15.1.8)$$

$$a_2\sigma = 3 = a_3, \quad (15.1.9)$$

$$a_3\sigma = 8 = a_4, \quad (15.1.10)$$

$$a_4\sigma = 5 = a_5 \quad (15.1.11)$$

$$a_5\sigma = 1 = a_1 \quad (15.1.12)$$

so here  $k = 5$ , thus we have a 5-cycle.

### Proposition 16.3 (*Commutativity product of disjoint cycles*)

The product of disjoint cycles is **commutative**.

*Proof.* Let  $\sigma_a = (a_1 \dots a_k)$  and  $\sigma_b = (b_1 \dots b_l)$  so that  $a_i \neq b_i$ . Then applying  $\sigma_b\sigma_a$ :

$$a_i\sigma_a\sigma_b = a_{i+1}\sigma_b = a_{i+1} \quad \text{for } i < k \quad (15.1.13)$$

$$a_k\sigma_a\sigma_b = a_1\sigma_b = a_1 \quad (15.1.14)$$

$$b_i\sigma_a\sigma_b = b_i\sigma_b = b_{i+1} \quad \text{for } i < l \quad (15.1.15)$$

$$b_l\sigma_a\sigma_b = b_l\sigma_b = b_1 \quad (15.1.16)$$

Similarly applying  $\sigma_b\sigma_a$  one finds:

$$a_i\sigma_b\sigma_a = a_i\sigma_a = a_{i+1} \quad \text{for } i < k \quad (15.1.17)$$

$$a_k\sigma_b\sigma_a = a_k\sigma_a = a_1 \quad (15.1.18)$$

$$b_i\sigma_b\sigma_a = b_{i+1}\sigma_a = b_{i+1} \quad \text{for } i < l \quad (15.1.19)$$

$$b_l\sigma_b\sigma_a = b_1\sigma_a = b_1 \quad (15.1.20)$$

■

An immediate consequence of Proposition 16.3 is that, if a permutation  $\sigma$  can be expressed as a product of disjoint cycles  $\sigma_i$ :

$$\sigma = \prod_{i=1}^n \sigma_i \implies \sigma^k = \prod_{i=1}^n \sigma_i^k \quad (15.1.21)$$

It turns out that this process of writing permutations as products of disjoint cycles can be done for any cycle, and this process is well-defined. In other words, the cycle form is uniquely determined for any permutation.

**Theorem 16.4 (Existence and uniqueness cycle form)**

Every permutation can be written in a unique cycle form (aside the choice of starting symbol and order in which the symbols are listed).

*Proof. Proof of existence*

Consider a permutaton  $\sigma \in \text{Sym}(\{1, 2, \dots, n\})$  and the infinite sequence:

$$a_1, a_1\sigma, a_1\sigma^2, a_1\sigma^3\dots \quad (15.1.22)$$

where  $a_1 \in \{1, 2, \dots, n\}$ . Because  $\{1, 2, \dots, n\}$  has finite cardinality, and the sequence continues infinitely, we must have some repetition  $a_1\sigma^i = a_1\sigma^j$  for some  $i < j$  which implies  $a_1\sigma^{j-i} = a_1$ . If we let  $k_1 = j - i$  be the smallest integer such that  $a_1\sigma^{k_1} = a_1$  then we denote  $\{a_1, a_1\sigma, \dots, a_1\sigma^{k_1-1}\}$  the **orbit** of  $a_1$  which is our first cycle.

Now if  $k_1 = n$  then the permutation  $\sigma$  is a cycle, and we are done. Otherwise, we choose  $a_1$  not in the orbit of  $a_1$  and write its orbit.

The orbits of  $a_1$  and  $a_2$  are disjoint since  $\sigma$  is bijective. Indeed, if  $\sigma^i a_1 = \sigma^j a_2$  then  $a_1 = \sigma^{j-i} a_2$ , implying that  $a_1$  belongs to the orbit of  $a_2$ , a contradiction.

Repeating this process, since the set  $S$  is finite eventually we exhaust the number of symbols and find:

$$\sigma = (a_1 a_1\sigma \dots a_1\sigma^{k_1-1})(a_2 a_2\sigma \dots a_2\sigma^{k_2-1}) \dots (a_r a_r\sigma \dots a_r\sigma^{k_r-1}) \quad (15.1.23)$$

where  $r$  is the different number of orbits.

**Proof of uniqueness** Suppose now that we can write the permutation as two distinct products of disjoint cycles:

$$\sigma = \rho_1 \rho_2 \dots \rho_k = \tau_1 \tau_2 \dots \tau_l \quad (15.1.24)$$

Then  $a_1 \in \{1, 2, \dots, n\}$  appears exactly once in  $\rho_i$  and  $\tau_j$  and as they are disjoint we reorder (through commutativity) the order of cycles so that  $i = j = 1$ . Hence WLOG assume that 1 appears in  $\tau_1$  and  $\rho_1$ .

We can then show that:

$$\rho_1 = (1 \ 1\sigma \ 1\sigma^2 \ \dots \ 1\sigma^{k-1}) = \tau_1 \quad (15.1.25)$$

an repeating this for all other  $\rho$  and  $\tau$  we get the desired result. ■

During this proof we encountered the important concept of an orbit which we shall revisit more in depth later when studying group actions:

**Definition 16.5** (*Orbit of a permutation element*)

The orbit of some element  $a_1 \in S$  in a permutation  $\sigma \in \text{Sym}(S)$  is the set  $\{a_1, \sigma a_1 \dots \sigma^{k_1-1} a_1\}$  where  $k_1$  is the smallest integer such that  $\sigma^{k_1} a_1 = a_1$ .

**Strategy.** (*Composing permutations*)

To find  $\sigma_g \circ \sigma_f$  in cycle form:

1. Start with 1 and find the symbol of 1 under  $\sigma_f$ , and then find the image of that symbol under  $\sigma_g$ , and denote it as  $x = \sigma_g \sigma_f 1$  so that  $\sigma_g \circ \sigma_f = (1 \ x \dots)$ .
2. Start with the symbol  $x$  and repeat the process.
3. Continue repeating the process until you reach the original symbol 1. The cycle is then complete.
4. Choose the smallest symbol not placed in the cycle, and repeat steps 1-3 again until the second cycle is complete.
5. Continue until every symbol has been placed.

**Example.** Write in cycle form  $(1 \ 4 \ 6)(3 \ 5) \circ (1 \ 5 \ 3 \ 2 \ 4) \circ (1 \ 2)(3 \ 5)(4 \ 6)$ .

We start with 1, which gets mapped to 2, then to 4 and finally to 6, so  $1 \mapsto 6$ .

Now we see that 6 gets mapped to 4, then to 1 and finally to 4, so  $6 \mapsto 4$ .

Now we see that 4 gets mapped to 6 and then doesn't get mapped and finally to 1, so  $4 \mapsto 1$ .

This completes the first cycle  $(1 \ 6 \ 4)$ .

We start with 2, which gets mapped to 1, then to 5 and finally to 3, so  $2 \mapsto 3$ .

Now we see that 3 gets mapped to 5, then to 3 and finally to 5 so  $3 \mapsto 5$ .

Now we see that 5 gets mapped to 3, then to 2 and then doesn't get mapped, so  $5 \mapsto 2$ .

This completes the second cycle  $(2 \ 3 \ 5)$ .

So we may write:

$$(1 \ 4 \ 6)(3 \ 5) \circ (1 \ 5 \ 3 \ 2 \ 4) \circ (1 \ 2)(3 \ 5)(4 \ 6) = (1 \ 6 \ 4)(2 \ 3 \ 5) \quad (15.1.26)$$



## 15.2 Permutation groups

**Theorem 16.6** (*Symmetric group*)

The set  $S_n$  of all permutations of  $\{1, 2, 3, \dots, n\}$  forms a group under  $\circ$  called the **symmetric group of degree  $n$**

*Proof.* We check that the group axioms hold:

(Closure) Consider  $\sigma_g \circ \sigma_f \in S_n$ . Since  $\sigma_g$  and  $\sigma_f$  are bijective maps mapping  $\{1, 2, \dots, n\}$  to itself, then  $\sigma_g \circ \sigma_f$  and  $\sigma_f \circ \sigma_g$  must necessarily also map  $S_n$  to itself and be bijective. So,  $\sigma_g \circ \sigma_f \in S_n$  as required.

(Associativity) Composition is associative.

(Identity) The identity permutation  $e$  is an identity, since its action is to map every symbol of  $\{1, 2, \dots, n\}$  to itself.

(Inverse) Since  $\sigma_f$  is bijective, it must have an inverse  $\sigma_f^{-1}$  which is also a permutation since it too maps  $\{1, 2, \dots, n\}$  to itself. ■

It is important to notice the difference between the order and degree of  $S_n$ . The order is, as always, the number of permutations  $S_n$  contains whereas its degree is how many symbols each of its permutations permute.

### Theorem 16.7 (Properties of $S_n$ )

The order  $|S_n|$  of  $S_n$  is  $n!$ , and for  $n \geq 3$ ,  $S_n$  is non-abelian.

*Proof.* Firstly, for  $\sigma \in S_n$ , there are  $n$  different possibilities for  $1\sigma$ , and since  $\sigma$  is bijective,  $1\sigma \neq 2\sigma$  and thus there are  $n - 1$  possibilities for  $2\sigma$ . We can keep doing this to find:

$$n \times (n - 1) \times (n - 2) \dots \times 2 \times 1 = n! \quad (15.2.1)$$

are all the possible permutations. Now if  $x_1, x_2, x_3 \in \{1, 2, \dots, n\}$  are distinct then we can define:

$$\sigma_1 : x_1 \mapsto x_1, x_2 \mapsto x_3, x_3 \mapsto x_2 \quad (15.2.2)$$

$$\sigma_2 : x_1 \mapsto x_2, x_2 \mapsto x_1, x_3 \mapsto x_3 \quad (15.2.3)$$

We then find that:

$$\sigma_2 \circ \sigma_1 : x_1 \mapsto x_2, x_2 \mapsto x_3, x_3 \mapsto x_1 \quad (15.2.4)$$

$$\sigma_1 \circ \sigma_2 : x_1 \mapsto x_3, x_2 \mapsto x_1, x_3 \mapsto x_2 \quad (15.2.5)$$

which are not the same. Thus for sets  $S$  of cardinality greater than or equal to 3 the corresponding  $\text{Sym}(S)$  is non-abelian. ■

### Definition 16.8 (Same cycle structure)

Two permutations in  $S_n$  have the same cycle structure if their cycle forms contain the same number of disjoint  $k$ -cycles for each  $k$ .

So for example  $(1\ 2\ 4)(3\ 8)(4\ 5)$  and  $(1\ 7)(2\ 8\ 3)(4\ 5)$  have the same cycle structure since they each consist of a 3-cycle, two **transposition** (2-cycle) and a 1-cycle (which is not shown in cycle form).

### Proposition 16.9 (Order of a permutation)

The order of a permutation  $\sigma = \rho_1 \rho_2 \dots \rho_k$  where  $\rho_i$  are disjoint cycles of order  $l_i$  is:

$$\text{ord}(\sigma) = \text{LCM}(l_1, l_2, \dots, l_k) \quad (15.2.6)$$

so that the order of a  $k$ -cycle is  $k$ .

*Proof.* Let  $n = \text{ord}(\sigma)$ . Then, since  $\rho_i$  are disjoint, we can use 16.1.21 to write:

$$\sigma^n = \rho_1^n \rho_2^n \rho_3^n \dots \rho_k^n = e \quad (15.2.7)$$

Now since  $\rho_i$  are disjoint, they permute different sets, and consequently their product can only be equal to  $e$  if each of them are equal to  $e$ :

$$\rho_i^n = e, \forall i \quad (15.2.8)$$

so that  $n|l_i$ . Since  $n$  must be the smallest such integer, it follows that  $n = \text{LCM}(l_1, l_2, \dots, l_k)$  as desired. ■

We can use cycle form to denote symmetries of figures as well. For example the symmetry:

$$\begin{pmatrix} 1 & 2 & 3 & 4 \\ 1 & 4 & 3 & 2 \end{pmatrix} \quad (15.2.9)$$

can be written as  $(2\ 4)$  in cycle form, which is much more concise.

**Example.** Write all the symmetries of an octahedron.

The octahedron has six direct symmetries: we can rotate it about the vertical line through the vertices 4,5 through  $0, 2\pi/3, 4\pi/3$ , or turn the octahedron upside down and repeat the same process. The octahedron also has at least one indirect symmetry, name a reflection in the plane through vertices 1,2,3. Hence there are 6 indirect symmetries, and 12 symmetries in total.

We first start by writing all direct symmetries of the equilateral triangle with vertices 1,2,3:

$$e, (1\ 2\ 3), (1\ 3\ 2), (1\ 2), (1\ 3), (2\ 3) \quad (15.2.10)$$

We can then compose with the reflection  $(4\ 5)$  to find:

$$(4\ 5), (1\ 2\ 3)(4\ 5), (1\ 3\ 2)(4\ 5), (1\ 3)(4\ 5), (2\ 3)(4\ 5) \quad (15.2.11)$$

These are twelve distinct symmetries, and therefore we have found all the symmetries. ◀

## 15.3 Even and Odd symmetries

**Strategy.** (*Expressing cycles as composite of transpositions*)

Consider a cycle  $(a_1\ a_2\ a_3\dots a_r)$ , then we can express:

$$(a_1\ a_2\ a_3\dots a_r) = (a_1\ a_r) \circ (a_1\ a_{r-1}) \circ \dots \circ (a_1\ a_2) \quad (15.3.1)$$

*Proof.* Firstly, consider  $a_1$ . It gets mapped to  $a_2$ .  $a_2$  then gets mapped to itself, since it does not appear in any other transposition (they are disjoint). So overall  $a_1$  gets mapped to  $a_2$ .

Now consider the  $a_s$  where  $2 \leq s \leq r - 1$ . We then see that  $(a_1\ a_2), \dots, (a_1\ a_{s-1})$  all map  $a_s$  to itself. The next transposition  $(a_1\ a_s)$  maps  $a_s$  to  $a_1$ . Then the next transposition  $(a_1\ a_{s+1})$  maps  $a_1$  to  $a_{s+1}$ .

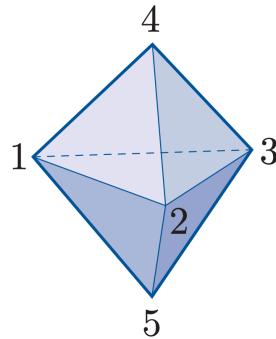


Figure 15.2. Octahedron

Finally, all the successive transpositions  $(a_1 \ a_{s+2}) \dots (a_1 \ a_r)$  map  $a_{s+1}$  to itself. So we find that overall  $a_s$  gets mapped to  $a_{s+1}$ .

Next we consider  $a_r$ . By the same argument as before all transpositions  $(a_1 \ a_2), \dots (a_1 \ a_{r-1})$  map  $a_r$  to itself. The final transpositions  $(a_1 \ a_r)$  map  $a_r$  to  $a_1$  as required.

So we find that  $a_s \mapsto a_{s+1}$  for all  $1 \leq s < r$  and  $a_r \mapsto a_1$ , which defines the cycle  $(a_1 \ a_2 \ a_3 \dots a_r)$  as in definition 16.2. ■

**Example.** We can express  $(2 \ 4 \ 3 \ 5) = (2 \ 5) \circ (2 \ 3) \circ (2 \ 4)$ . ◀

Notice that there are several ways we can express a permutation as a composite of transposition. Indeed in the previous example we could have written:

$$(2 \ 4 \ 3 \ 5) = (4 \ 3 \ 5 \ 2) = (4 \ 2) \circ (4 \ 5) \circ (4 \ 3) \quad (15.3.2)$$

$$= (2 \ 4) \circ (5 \ 4) \circ (3 \ 4) \quad (15.3.3)$$

Notice however that the number of transpositions is always even. Indeed it turns out that if a permutation can be expressed as a composition of an even number of transpositions, then it can only be expressed as a composite of even transpositions.

### Definition 16.9 (*Parity of permutation*)

A permutation is **even** if it can be expressed as a composite of an even number of transpositions.

A permutation is **odd** if it can be expressed as a composite of an odd number of transpositions.

We have that:

### Theorem 16.10 (*Parity Theorem*)

A permutation cannot be expressed as both a composite of an even number of transpositions and a composite of an odd number of transpositions.

*Proof.* Let  $x_1, \dots, x_n$  and consider the Vandermonde polynomial <sup>1</sup>

$$P = \prod_{1 \leq i < j \leq n} (x_j - x_i) \quad (15.3.4)$$

We now choose a permutation  $\sigma \in S_n$ , and define the function  $f_\sigma$  as:

$$f_\sigma(P) = \prod_{1 \leq i < j \leq n} (x_{\sigma(j)} - x_{\sigma(i)}) \quad (15.3.5)$$

which reshuffles the terms in the normal Vandermonde polynomial.

**Lemma.** Let  $\tau$  be a transposition, then  $f_\tau(P) = -P$ .

*Proof.* Let's consider the action of a transposition  $\tau = (a b)$  on  $P$ .

For  $j, i \neq a, b$  in any order, the factor  $(x_j - x_i)$  is left unchanged.

Now let us consider a pair with exactly one index equal to  $a$  or  $b$  (we assume WLOG that  $a < b$ ).

Then if the other index  $j$  is between  $a$  and  $b$ ,  $(x_j - x_a) \mapsto -(x_b - x_j)$  and  $(x_b - x_j) \mapsto -(x_j - x_a)$ . However these two signs cancel each other out, so no net change.

If the other index  $i$  is not between  $a$  and  $b$ , then it can be smaller than  $a$  or larger than  $b$ . In the first case,  $(x_a - x_i) \mapsto (x_b - x_i)$ . In the former case,  $(x_i - x_a) \mapsto (x_i - x_b)$ . In both cases the sign does not change.

Finally, if both indices are in  $\{a, b\}$  then  $x_b - x_a \mapsto -(x_b - x_a)$ .

So overall there is only one sign change due to when both indices correspond to  $a$  and  $b$ . So if  $\tau$  is a transposition, then  $f_\tau(P) = -P$  and consequently  $f_\tau(-P) = P$ . ■

Now take an arbitrary permutation  $\sigma$ , and express it as a product of transpositions  $\tau_i$  and  $\rho_i$  in two different ways:

$$\sigma = \tau_1 \cdots \tau_r = \rho_1 \cdots \rho_s \quad (15.3.6)$$

Then

$$f_\sigma(P) = f_{\tau_1 \cdots \tau_r}(P) = (f_{\tau_1} \circ \cdots \circ f_{\tau_r})(P) = (-1)^r P \quad (15.3.7)$$

$$f_\sigma(P) = f_{\rho_1 \cdots \rho_s}(P) = (f_{\rho_1} \circ \cdots \circ f_{\rho_s})(P) = (-1)^s P \quad (15.3.8)$$

Therefore,  $(-1)^r P = (-1)^s P$ , so  $r$  and  $s$  have the same parity: both odd, or both even. ■

### Proposition 16.11 (Parity of $k$ -cycles)

For  $\sigma \in S_n$  where  $\sigma$  is a  $k$ -cycle:

$$\sigma \text{ is } \begin{cases} \text{even permutation, if } k \text{ is odd} \\ \text{odd permutation, if } k \text{ is even} \end{cases} \quad (15.3.9)$$

<sup>1</sup>For  $n = 4$ , we would have  $P = (x_2 - x_1)(x_3 - x_1)(x_4 - x_1)(x_3 - x_2)(x_4 - x_2)(x_4 - x_3)$ .

*Proof.* This follows immediately from equation 16.3.1 ■

We note immediately then that the parity of the composite of two permutations can be found by simply summing their parity.

**Strategy.** (*Finding the parity of any permutation*)

1. Express the permutation as a composite of cycles
2. find the parity of each  $k$ -cycle following

$$\begin{cases} \text{even permutation, if } k \text{ is odd} \\ \text{odd permutation, if } k \text{ is even} \end{cases} \quad (15.3.10)$$

3. Combine the parities of each cycle following the Cayley table below:

+	even	odd
even	even	odd
odd	odd	even

**Proposition 16.12 (Parity of inverse)**

The inverse of a permutation has the same parity as the permutation.

*Proof.* We see that  $e$  is an even permutation, so if  $f \circ f^{-1} = e$  then we must have that  $f$  and  $f^{-1}$  have the same parity. ■

**Proposition 16.13 (Alternating group of order  $n$ )**

The group  $A_n$  of all even permutations of  $\{1, 2, \dots, n\}$  is called the **alternating group of order  $n$**  and  $A_n \leq S_n$ .

*Proof.* We check the subgroup properties:

- (i) **Closure:** the composite of two even permutations is even, so closure is satisfied.
- (ii) **Identity:** the identity permutation  $e$  is even, and thus belongs to  $A_n$ .
- (iii) **Inverses:** by theorem 16.10, the inverse of an even permutation is even.

For example, let us try to find all elements of  $A_4$ . Their structures must be:

$$e, (\underline{\quad})(\underline{\quad}), (\underline{\quad}\underline{\quad}) \quad (15.3.11)$$

and we must fill the gaps with  $\{1, 2, 3, 4\}$ .

Note that there are only 3 cycle structures of the form  $(\underline{\quad})(\underline{\quad})$ . Indeed, WLOG place 1 in the first place, then we can place 3 elements in the second place. The other transposition is then given immediately. So there are three different ways.

Instead, there are 8 cycle structures of the form  $(\_ \_ \_ \_)$ . For these, we place WLOG 1 in the first place so that there are 3 possible elements in the second place, and 2 possible in the third, the fourth is then immediate. So there are eight different ways.

In total there are then 12 different elements in  $A_4$ , which is exactly half the order of  $S_4$  interestingly.

This turns out not to be a coincidence. Indeed, we have the following general result:

**Theorem 16.14 (Order of  $A_n$ )**

The order of the alternating group is  $|A_n| = \frac{1}{2}n!$  for  $n \geq 2$ .

*Proof.* Suppose  $S_n$  has  $r$  even permutations and  $s$  odd permutations.

We first prove that  $r \leq s$ . Let  $f_1, \dots, f_r \in A_n$  then consider:

$$(1 \ 2) \circ f_1, (1 \ 2) \circ f_2, \dots, (1 \ 2) \circ f_r \quad (15.3.12)$$

these are all distinct odd permutations. So there are at least  $r$  odd permutations in  $S_n$ .

A similar argument uses  $g_1 \dots g_r$  odd permutations composed with  $(1 \ 2)$  to prove that  $s \leq r$ .

Since  $s \leq r$  and  $r \leq s$ , we have that  $s = r$ . So exactly half of the permutations in  $S_n$  are in  $A_n$ , thus  $|A_n| = \frac{1}{2}n!$ . ■

## 15.4 Conjugacy of $S_n$

Consider the permutation  $x$  of some set  $S$  with  $i, j \in S$ . Then consider another permutation  $g$  which relabels  $S$ . Our question is to know what the permutation  $x$  looks like with the relabelled set  $S'$ .

Looking at the diagram below:

$$\begin{array}{ccc} i & \xrightarrow{x} & j \\ g^{-1} \uparrow & & \downarrow g \\ g(i) & \xrightarrow{y} & g(j) \end{array}$$

we clearly see that  $(y \circ g)(i) = (g \circ x \circ g^{-1})(i)$  and by the cancellation rule we find that  $y = g \circ x \circ g^{-1}$ .

**Definition 16.15 (Conjugate permutations)**

The permutation  $\sigma$  is the conjugate of  $\rho$  in  $S_n$  if there exists a permutation  $\tau$  such that:

$$\sigma = \tau \circ \rho \circ \tau^{-1} \quad (15.4.1)$$

We then say that  $\tau$  is a **conjugating permutation** of  $\rho$  to  $\sigma$ , and that  $\sigma$  is the **conjugate** of  $\rho$  by  $\tau$ .

**Strategy.** (*Finding a conjugating permutation*)

1. Align the cycles of  $\sigma$  and  $\rho$  so that cycles of same order correspond:

$$\begin{aligned} x &= (* \ * \ \dots \ *) (* \ * \ \dots \ *) \dots (*)(*) \\ y &= (* \ * \ \dots \ *) (* \ * \ \dots \ *) \dots (*)(*) \end{aligned} \quad \downarrow \tau \quad (15.4.2)$$

2. Read off the two line form of the permutation  $\tau$ .

**Example.** Let  $\sigma = (1 \ 2 \ 4)(3 \ 5)$  and  $\rho = (1 \ 4)(2 \ 5 \ 3)$  in  $S_5$ . Find three permutations  $g \in S_5$  conjugating  $\sigma$  to  $\rho$ .

We can write that:

$$\begin{aligned} x &= (1 \ 2 \ 4)(3 \ 5) \\ y &= (2 \ 5 \ 3)(1 \ 4) \end{aligned} \quad \downarrow \tau \quad (15.4.3)$$

and read off the conjugating permutation  $\tau = (1 \ 2 \ 5 \ 4 \ 3)$ .

Alternatively, we can rewrite  $(2 \ 5 \ 3)$  as  $(3 \ 2 \ 5)$  and find:

$$\begin{aligned} x &= (1 \ 2 \ 4)(3 \ 5) \\ y &= (3 \ 2 \ 5)(1 \ 4) \end{aligned} \quad \downarrow \tau \quad (15.4.4)$$

whose corresponding conjugate permutation is  $\tau = (1 \ 3)(4 \ 5)$ .

Finally we can rewrite  $(2 \ 5 \ 3)$  as  $(3 \ 5 \ 2)$  and find:

$$\begin{aligned} x &= (1 \ 2 \ 4)(3 \ 5) \\ y &= (3 \ 5 \ 2)(1 \ 4) \end{aligned} \quad \downarrow \tau \quad (15.4.5)$$

whose corresponding conjugate permutation is  $\tau = (1 \ 3)(2 \ 5 \ 4)$  ◀

Now consider the action of a conjugating permutation not on a single permutation, but on every permutation in a subgroup.

Let  $H \leq S_n$  and let  $g \in S_b$ , then we will denote:

$$g \circ H \circ g^{-1} = \{g \circ h \circ g^{-1} : h \in H\} \quad (15.4.6)$$

So it suffices to substitute every element in  $H$  using  $g$ .

If we let for example  $H = \langle (1 \ 2 \ 4 \ 5) \rangle$  then:

$$H = \{e, (1 \ 2 \ 4 \ 5), (1 \ 2 \ 4 \ 5)^2, (1 \ 2 \ 4 \ 5)^3\} \quad (15.4.7)$$

$$= \{e, (1 \ 2 \ 4 \ 5), (1 \ 4)(2 \ 5), (1 \ 5 \ 4 \ 2)\} \quad (15.4.8)$$

Then, if we let  $g = (3 \ 5)$  we find that:

$$g \circ H \circ g^{-1} = \{e, (1 \ 2 \ 4 \ 3), (1 \ 4)(2 \ 3), (1 \ 3 \ 4 \ 2)\} \quad (15.4.9)$$

### Theorem 16.16 (Conjugate subgroups)

Let  $H \leq S_n$  and let  $g \in S_n$ . Then  $g \circ H \circ g^{-1}$  is also a subgroup of  $S_n$ .

*Proof.*

**Closure:** consider any two elements  $h, k \in g \circ H \circ g^{-1}$ . Then,  $\exists h', k'$  such that  $h = g \circ h' \circ g^{-1}$  and  $k = g \circ k' \circ g^{-1}$ . Thus:

$$h \circ k = (g \circ h' \circ g^{-1}) \circ (g \circ k' \circ g^{-1}) \quad (15.4.10)$$

$$= g \circ h' \circ g^{-1} \circ g \circ k' \circ g^{-1} \quad (15.4.11)$$

$$= g \circ (h' \circ k) \circ g^{-1} \quad (15.4.12)$$

$$= g \circ l \circ g^{-1} \quad (15.4.13)$$

where  $l = h' \circ k \in H$  since subgroups are closed. Hence  $h \circ k \in g \circ H \circ g^{-1}$  as required.

**Identity:** let the identity element in  $H$  be  $e_H$ . Then:

$$e_H = g \circ g^{-1} = g \circ e_H \circ g^{-1} \in g \circ H \circ g^{-1} \quad (15.4.14)$$

as required.

**Inverses** Let  $h \circ h^{-1} = e_H$  for all  $h \in H$ . Then, the inverse of  $g \circ h \circ g^{-1}$  is  $g \circ h^{-1} \circ g^{-1}$ . Indeed:

$$g \circ h \circ g^{-1} \circ g \circ h^{-1} \circ g^{-1} = g \circ h \circ h^{-1} \circ g^{-1} = g \circ g^{-1} = e_H \quad (15.4.15)$$

Let us now check that the inverse belongs to  $g \circ H \circ g^{-1}$ . This is clearly true, since  $h^{-1} \in H$  due to the inverse property of subgroups.

■

## 15.5 Subgroups of $S_4$

We will now tackle the problem of finding all subgroups of  $S_4$ . To do so we will find cyclic and all non-cyclic subgroups of  $S_4$ .

### Cyclic subgroups

The cyclic elements of  $S_4$  must have the structures shown in Figure 16.3. We see that each element in  $S_4$  has order 1,2,3,4 so the order of each subgroup must also be 1,2,3 or 4.

The cyclic subgroup of order 1 is obviously  $\{e\}$ .

The cyclic subgroups of order 2 are the identity permutation with one permutation of order 2:

$$\{e, (1 2)\}, \{e, (1 3)\}, \{e, (1 4)\}, \{e, (2 3)\}, \{e, (2 4)\}, \{e, (3 4)\} \quad (15.5.1)$$

and:

$$\{e, (1 2)(3 4)\}, \{e, (1 3)(2 4)\}, \{e, (1 4)(2 3)\} \quad (15.5.2)$$

which are 9 in total.

To find all cyclic subgroups of order 3, note that  $(1 2 3)$  and  $(1 3 2)$  all generate the same subgroup  $\{e, (1 2 3), (1 3 2)\}$ . Similarly, the other six 3 – cycles each couple up in a similar way. So the cyclic groups are:

$$\{\langle (1 2 3) \rangle, \langle (1 3 4) \rangle, \langle (1 4 2) \rangle, \langle (2 3 4) \rangle\} \quad (15.5.3)$$

Cycle structure	Order	Elements of $S_4$	Description
$e$	1	$e$	identity
(--)	2	(1 2), (1 3), (1 4), (2 3), (2 4), (3 4)	transpositions
(---)	3	(1 2 3), (1 3 2), (1 2 4), (1 4 2), (1 3 4), (1 4 3), (2 3 4), (2 4 3)	3-cycles
(----)	4	(1 2 3 4), (1 2 4 3), (1 3 2 4), (1 3 4 2), (1 4 2 3), (1 4 3 2)	4-cycles
(--)(--)	2	(1 2)(3 4), (1 3)(2 4), (1 4)(2 3)	products of 2-cycles

**Figure 15.3.** Cycle structures of elements in  $S_4$

The cyclic subgroups of order 4 can be found similarly.  $(1 2 3 4)$  generates the following set:

$$\langle(1 2 3 4)\rangle = \{e, (1 2 3 4), (1 3)(2 4), (1 4 3 2)\} \quad (15.5.4)$$

We now choose a permutation of order 4 that is not in this list, and finds its generator:

$$\langle(1 2 4 3)\rangle = \{e, (1 2 4 3), (1 4)(2 3), (1 3 4 2)\} \quad (15.5.5)$$

Repeat this process one final time:

$$\langle(1 3 2 4)\rangle = \{e, (1 3 2 4), (1 2)(3 4), (1 4 2 3)\} \quad (15.5.6)$$

We see that all six permutations of order 4 have been found, so we found all the cyclic subgroups of  $S_4$  of order 4.

So in conclusion:

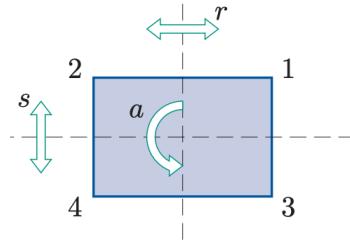
Order	Number of cyclic subgroups
1	1
2	9
3	4
4	3

with 16 total cyclic subgroups.

### Non-cyclic subgroups

We now try to find non-cyclic subgroups of  $S_4$ . We can do so by drawing a figure labels 1, 2...n. We can then find the symmetry group of the figure, which is a subgroup of  $S_n$ .

For example, if we draw and label the rectangle as shown below:



Then the symmetry group is:

$$\{e, (1 2)(3 4), (1 3)(2 4), (1 4)(2 3)\} \quad (15.5.7)$$

which is not cyclic. Indeed it has 4 elements, but its elements are of order 2.

We can now find other non-cyclic subgroup from the old subgroup by relabelling it through conjugating permutations.

## 15.6 Cayley's Theorem

We saw that the symmetry groups of most figures can be represented as permutation groups, and are therefore isomorphic.

It turns out this is true for any finite group.

### Theorem 16.17 (Cayley's theorem)

Let  $(G, *)$  be a finite group. For each  $x \in G$ , let  $p_x$  be the permutation whose two-line symbol has as its first line the column heading of the Cayley table of  $G$ , and as its second line the row labelled  $x$  in the group table.

If we let  $P = \{p_x : x \in G\}$  then  $(P, \circ)$  is a permutation group isomorphic to  $(G, *)$ .

*Proof.* Let  $G = \{g_1, \dots, g_n\}$  so that for each element  $g_i$  the table shows:

*	$g_1$	$g_2$	$g_3$	$\dots$	$g_n$
$g_1$					
$\vdots$					
$x$	$x * g_1$	$x * g_2$	$x * g_3$	$\dots$	$x * g_n$
$\vdots$					
$g_n$					

so that:

$$p_x = \begin{pmatrix} g_1 & g_2 & \dots & g_n \\ x * g_1 & x * g_2 & \dots & x * g_n \end{pmatrix} \quad (15.6.1)$$

which is a permutation since every element of  $G$  is repeated only once in the column headings and in the row labelled  $x$ .

One can then easily verify using associativity that:

$$p_x \circ p_y = \begin{pmatrix} g_1 & g_2 & \cdots & g_n \\ x * y * g_1 & x * y * g_2 & \cdots & x * y * g_n \end{pmatrix} = p_{x*y} \quad (15.6.2)$$

so that:

*	...	y	...		o	...	$p_y$	...
:		:			:		:	
x	...	$x * y$	...		$p_x$	...	$p_{x*y}$	...
:		:			:		:	

$(G, *)$                                      $(P, o)$

The two Cayley tables therefore are identical in structure, and thus  $(P, o)$  is a group. Not only that, it is a group isomorphic to  $(G, *)$ , since the map:

$$p : G \rightarrow P \quad (15.6.3)$$

$$x \mapsto p_x \quad (15.6.4)$$

Indeed, we have already verified that  $p(x*y) = p(x)o(p(y))$ . Also,  $p$  is injective, since  $p(x) = p(y) \implies x * g_i = y * g_i$  for all  $g_i \in G$ , and by the right cancellation law  $x = y$ . Surjectivity is trivial. ■

For example, consider the group  $(\mathbb{Z}_6, +_6)$ . For each  $x \in \mathbb{Z}_6$ , we associate a permutation  $p_x$  whose

$+_6$	0	1	2	3	4	5
0	0	1	2	3	4	5
1	1	2	3	4	5	0
2	2	3	4	5	0	1
3	3	4	5	0	1	2
4	4	5	0	1	2	3
5	5	0	1	2	3	4

**Figure 15.4.** Cayley table for  $(\mathbb{Z}_6, +_6)$

two line form has its first line as the column heading and its second line as the row labelled  $x$ . So:

$$p_2 = \begin{pmatrix} 0 & 1 & 2 & 3 & 4 & 5 \\ 2 & 3 & 4 & 5 & 0 & 1 \end{pmatrix} = (0 \ 2 \ 4)(1 \ 3 \ 5) \quad (15.6.5)$$

This gives us the permutations  $p_0, p_1 \dots p_6$  for each element of  $\mathbb{Z}_6$ , obtaining:

$$p_0 = e \quad (15.6.6)$$

$$p_1 = (0 \ 1 \ 2 \ 3 \ 4 \ 5) \quad (15.6.7)$$

$$p_2 = (0 \ 2 \ 4)(1 \ 3 \ 5) \quad (15.6.8)$$

$$p_3 = (0 \ 3)(1 \ 4)(2 \ 5) \quad (15.6.9)$$

$$p_4 = (0 \ 4 \ 2)(1 \ 5 \ 3) \quad (15.6.10)$$

$$p_5 = (0 \ 5 \ 4 \ 3 \ 2 \ 1) \quad (15.6.11)$$

Let  $P = \{p_0, p_1, p_2 \dots p_5\}$ . If we draw the Cayley table for  $(P, \circ)$  we find:

$\circ$	$p_0$	$p_1$	$p_2$	$p_3$	$p_4$	$p_5$
$p_0$	$p_0$	$p_1$	$p_2$	$p_3$	$p_4$	$p_5$
$p_1$	$p_1$	$p_2$	$p_3$	$p_4$	$p_5$	$p_0$
$p_2$	$p_2$	$p_3$	$p_4$	$p_5$	$p_0$	$p_1$
$p_3$	$p_3$	$p_4$	$p_5$	$p_0$	$p_1$	$p_2$
$p_4$	$p_4$	$p_5$	$p_0$	$p_1$	$p_2$	$p_3$
$p_5$	$p_5$	$p_0$	$p_1$	$p_2$	$p_3$	$p_4$

which is structurally identical to the table for  $(\mathbb{Z}_6, +_6)$ . So we can conclude that the map:

$$\phi : \mathbb{Z}_6 \rightarrow P \quad (15.6.12)$$

$$x \mapsto p_x \quad (15.6.13)$$

is an isomorphism.

# Unit B4: Lagrange's Theorem and small groups

## 16.1 Lagrange's Theorem

### Theorem 17.1 (Lagrange's Theorem)

Let  $G$  be a finite group, and let  $H \leq G$ . Then  $\text{ord}(H) | \text{ord}(G)$ , that is, the order of  $H$  divides the order of  $G$ .

*Proof.* Let  $(G, \circ)$  be a finite group and let  $H \leq G$  so that  $\text{ord}(G) = s$  and  $\text{ord}(H) = r$ .

We begin by writing down all the elements of  $H$ :

$$(h_1 \ h_2 \ \dots \ h_r) \quad (16.1.1)$$

Next we choose any element of  $G$  that is not included in the above array, such as  $g_2$ , and compose it to the left with the first row.

$$\begin{pmatrix} h_1 & h_2 & \dots & h_r \\ g_2 \circ h_1 & g_2 \circ h_2 & \dots & g_2 \circ h_r \end{pmatrix} \quad (16.1.2)$$

If there are no other elements of  $G$  excluded from the array, then we are done. Otherwise, choose another element, say  $g_3$  that is not included and compose it to the left with the first row to find:

$$\begin{pmatrix} h_1 & h_2 & \dots & h_r \\ g_2 \circ h_1 & g_2 \circ h_2 & \dots & g_2 \circ h_r \\ g_3 \circ h_1 & g_3 \circ h_2 & \dots & g_3 \circ h_r \end{pmatrix} \quad (16.1.3)$$

We repeat this process until all the elements of  $G$  have been exhausted. This must happen since  $G$  has finite order and each row introduces a new element  $g_i \in G$ .

At the end, we reach the following array:

$$\begin{pmatrix} h_1 & h_2 & \dots & h_r \\ g_2 \circ h_1 & g_2 \circ h_2 & \dots & g_2 \circ h_r \\ g_3 \circ h_1 & g_3 \circ h_2 & \dots & g_3 \circ h_r \\ \vdots & \vdots & & \vdots \\ g_k \circ h_1 & g_k \circ h_2 & \dots & g_k \circ h_r \end{pmatrix} \quad (16.1.4)$$

Next, we have to show that all the elements in the array are distinct. We start by showing that all the elements in a row are distinct. This is clearly true for the first row, since they are all distinct elements of  $H$ . For the  $k$ th row, we have that if for some  $h_i, h_j \in H$  distinct:

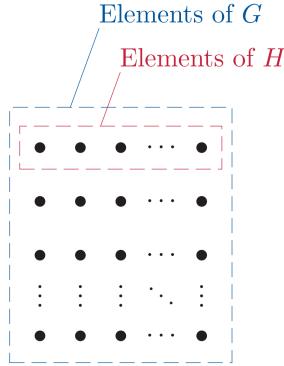
$$g_k \circ h_i = g_k \circ h_j \quad (16.1.5)$$

then by the left cancellation law  $h_i = h_j$  which is a contradiction.

Secondly, we show that elements in a row are not repeated in any other row. Again, we go by contradiction, and suppose that in some row  $l$ , the element  $g_l \circ h_i$  is repeated as  $g_k \circ h_j$  another row  $k$ :

$$g_l \circ h_i = g_k \circ h_j \implies g_l = g_k \circ h_j \circ h_i^{-1} \quad (16.1.6)$$

By the closure property of  $H$ ,  $h_j \circ h_i^{-1} \in H$ , which would imply that  $g_l = g_k \circ h_m$  for some  $m$  and that therefore  $g_l$  belongs to the  $k$ th row. This is a contradiction, since we assumed that the rows  $l$  and  $k$  are different.



**Figure 16.1.** Visualization of Lagrange's proof

Thus none of the elements in each row are repeated in other rows. We can therefore conclude that the order of  $G$  is the size of the complete matrix, that is,  $\text{ord}(G) = k \cdot r = k \cdot \text{ord}(H)$ . It follows immediately that  $\text{ord}(H) | \text{ord}(G)$ . ■

### Corollary.

- (i) let  $g \in G$ , then  $\text{ord}(g) | \text{ord}(G)$ .
- (ii) let  $G$  be a group of prime order. Then  $G$  is cyclic, with every non-identity element being a generator.
- (iii) let  $G$  be a group of prime order  $p$ , then  $(G, \circ) \cong (\mathbb{Z}_p, +_p)$ .

*Proof.*

- (i) we have that  $\langle g \rangle \leq G$  and  $\text{ord}(g) = \text{ord}(\langle g \rangle)$  so that  $\text{ord}(g)|\text{ord}(G)$ .
- (ii) if  $G$  has prime order  $p$ , then for every element  $g \in G$ ,  $\langle g \rangle$  can have order 1 or  $p$ . However, only  $\langle e \rangle$  has order 1, therefore  $\langle x \rangle$  must have order  $p$ , and therefore generate  $G$ .
- (iii) we have that  $(G, \circ)$  is a cyclic group of order  $p$ , and that  $(\mathbb{Z}_p, +_p)$  too is a cyclic group of order  $p$ . From proposition 15.21, the two groups must therefore be isomorphic.

■

## 16.2 Groups of small order

The goal of this section will be to justify the following classification of isomorphism classes for small groups:

Order	Standard group(s)	Properties	Further examples
1	$C_1$	cyclic	$(\{0\}, +), (\{1\}, \times)$
2	$C_2, (\mathbb{Z}_2, +_2)$	cyclic	$S^+(\square), (\mathbb{Z}_3^*, \times_3)$
3	$C_3, (\mathbb{Z}_3, +_3)$	cyclic	$S^+(\triangle), (\{0, 4, 8\}, +_{12}), (\{1, 4, 7\}, \times_9)$
4	$C_4, (\mathbb{Z}_4, +_4)$	cyclic	$(\mathbb{Z}_5^*, \times_5), S^+(\square), S(\diamond), (\{0, 3, 6, 9\}, +_{12}), (\{1, -1, i, -i\}, \times)$
	$V, S(\square)$	abelian, non-cyclic	$(U_8, \times_8), (U_{12}, \times_{12}), (\{1, 7, 9, 15\}, \times_{16}), (\{1, 9, 11, 19\}, \times_{20})$
5	$C_5, (\mathbb{Z}_5, +_5)$	cyclic	$S^+(\diamond)$
6	$C_6, (\mathbb{Z}_6, +_6)$	cyclic	$S^+(\square), (\mathbb{Z}_7^*, \times_7), (U_9, \times_9), (\{0, 2, 4, 6, 8, 10\}, +_{12}), (U_{14}, \times_{14})$
	$S(\triangle)$	non-abelian	$S_3, \{e, (2 3), (2 4), (3 4), (2 3 4), (2 4 3)\}$
7	$C_7, (\mathbb{Z}_7, +_7)$	cyclic	$S^+(\text{heptagon})$
8	$C_8, (\mathbb{Z}_8, +_8)$	cyclic	$S^+(\text{octagon})$
	$S(\text{cuboid})$	abelian	
	$(U_{15}, \times_{15})$	abelian	$(U_{20}, \times_{20})$
	$S(\square)$	non-abelian	
	$Q_8$	non-abelian	

### Proposition 17.2 (Useful results)

Let  $G$  be a group of finite order:

- (i) if each element except the identity has order 2, then  $G$  is abelian.
- (ii) if  $\text{ord}(G) > 2$  and each element except  $e$  has order 2, then  $4|\text{ord}(G)$ .
- (iii) if  $\text{ord}(G)$  is even, then at least one element of  $G$  has order 2.

*Proof.*

- (i) let  $x, y \in G$ , then  $xy$  is either the identity element or it has order 2 so that:

$$(xy)^2 = e \implies xyxy = e \implies xey = x^2yxy^2 \implies xy = yx \quad (16.2.1)$$

since  $x^2 = y^2 = e$ . The group  $G$  is therefore abelian.

- (ii) Firstly, by the previous point  $G$  must be abelian. Also,  $G$  has at least 3 elements,  $e, x, y$ , with  $x^2 = y^2 = e$ . Now consider  $z = xy$ , this must be distinct from  $x, y, e$  since

$$z = e \implies xy = e \implies y = x \quad (16.2.2)$$

$$z = x \implies xy = x \implies y = e \quad (16.2.3)$$

$$z = y \implies xy = y \implies x = e \quad (16.2.4)$$

It now remains to prove that  $\{e, x, y, z\} \leq G$ . We construct the following Cayley table:

	$e$	$x$	$y$	$z$
$e$	$e$	$x$	$y$	$z$
$x$	$x$	$e$	$z$	$y$
$y$	$y$	$z$	$e$	$x$
$z$	$z$	$y$	$x$	$e$

where for example  $yxy = yyx = x$ . The subgroup properties are then readily verified. Closure holds since every element in the body of the table is in  $\{e, x, y, z\}$ . The identity element of  $G$  is  $e \in \{e, x, y, z\}$ . Finally, all elements are self-inverse, and consequently their inverses belong to the same set.

We conclude that  $\{e, x, y, z\} \leq G$ , and by Lagrange's theorem,  $4|\text{ord}(G)$ .

- (iii) the elements that are not-self inverse can be paired up with their inverses, so they must be even. It follows that the number of self-inverse elements must also be even (for if they were odd then  $G$  would have odd order). The identity element is one such self-inverse element, so there must be at least one more self-inverse element, which of course has order 2.

■

## Groups of order 1,2,3,5,7

Obviously, there is only one isomorphism class for groups of order 1, and that is  $C_1$ .

For the other groups of prime order  $p$ , we have from the Corollary to Lagrange's theorem that they are isomorphic to  $(\mathbb{Z}_p, +_p)$ , and therefore belong to the same isomorphism class.

So the isomorphism class for each group of prime order  $p$  is the one containing  $C_p$ .

## Groups of order 4

If  $G$  is a group of order 4, then by the Corollary to Lagrange's theorem, we must have that each element  $g \in G$  must have order 1,2 or 4.

### G has an element of group 4

If  $G$  has an element of order 4, then  $G$  is cyclic (generated by this element) and so isomorphic to  $C_4$ . They all belong to the same isomorphism class.

### G has no element of group 4

Only the identity element has order 1, so the other three elements must have order 2. By proposition 17.2 then,  $G$  is abelian, and if we let  $G = \{e, x, y, z\}$ , then by the same logic as in the proof of (ii)  $z = xy$  and we retrieve the following Cayley table: which is the table of the Klein four-group  $V$ , so

	$e$	$x$	$y$	$z$
$e$	$e$	$x$	$y$	$z$
$x$	$x$	$e$	$z$	$y$
$y$	$y$	$z$	$e$	$x$
$z$	$z$	$y$	$x$	$e$

that  $G \cong V$ .

Therefore the two isomorphism classes for groups of order 4 are the one containing  $C_4$  and the one containing  $V$ .

### Groups of order 6

Suppose that  $G$  is a group of order 6, so that each element of  $G$  has order 1,2,3 or 6.

**G has an element of group 6** In this case,  $G$  is a cyclic group, and is therefore isomorphic to  $C_6$ . It can be classified in the isomorphism class containing  $C_6$ .

**G has no element of group 6** In this case, each non-identity element has order 2 or 3. We can assert that it must have at least one element of order 2 by proposition 17.2, and similarly there must be at least one element of order 3, for if they were all of order 2 then  $4|\text{ord}(G)$  which clearly isn't the case.

So let  $g, h \in G$  be some elements of order 2 and 3 respectively. We define:

$$H = \langle h \rangle = \{e, h, h^2\} \quad (16.2.5)$$

Obviously  $g \neq H$ , since all elements of  $H$  must have order 1 or 3 by the Corollary to Lagrange's theorem. We can then adopt the proof we used for Lagrange's theorem and write the following matrix:

$$\begin{pmatrix} e & h & h^2 \\ g & gh & gh^2 \end{pmatrix} \quad (16.2.6)$$

which contains six distinct elements, and must therefore include all elements of  $G$ .

Hence, we know that  $G = \{e, h, h^2, g, g^2, gh, gh^2\}$ . We construct the following incomplete Cayley table:

To evaluate the missing entries, we need to calculate  $hg$ , which must be equal to  $gh$  or  $gh^2$  (not  $g$  since it already appears in the same column). However, if  $hg = gh$  then:

$$hg = gh \neq e \implies (hg)^3 = (hg)(gh)(hg) = g \neq e \quad (16.2.7)$$

so that  $hg$  has order greater than 3. The only possible value for  $\text{ord}(hg)$  is then 6, a contradiction.

	$e$	$h$	$h^2$	$g$	$gh$	$gh^2$
$e$	$e$	$h$	$h^2$	$g$	$gh$	$gh^2$
$h$	$h$	$h^2$	$e$			
$h^2$	$h^2$	$e$	$h$			
$g$	$g$	$gh$	$gh^2$	$e$	$h$	$h^2$
$gh$	$gh$	$gh^2$	$g$			
$gh^2$	$gh^2$	$g$	$gh$			

Therefore we must have that  $hg = gh^2$ . We then obtain (using the fact that each element must repeat only once in each row and column):

	$e$	$h$	$h^2$	$g$	$gh$	$gh^2$
$e$	$e$	$h$	$h^2$	$g$	$gh$	$gh^2$
$h$	$h$	$h^2$	$e$	$gh^2$	$g$	$gh$
$h^2$	$h^2$	$e$	$h$	$gh$	$gh^2$	$g$
$g$	$g$	$gh$	$gh^2$	$e$	$h$	$h^2$
$gh$	$gh$	$gh^2$	$g$	$h^2$	$e$	$h$
$gh^2$	$gh^2$	$g$	$gh$	$h$	$h^2$	$e$

This Cayley table is identical in structure to the tables of the groups  $S_3$  and  $S(\Delta)$ . This means that the second isomorphism class is that containing  $S(\Delta)$ .

Therefore, the two isomorphism classes for groups of order 6 are the one containing  $C_6$  and the one containing  $S_3$ .

# Unit E1: Cosets and normal subgroups

## 17.1 Matrix groups

We will introduce some important sets of matrices which form groups and subgroups under matrix multiplication. They are especially important in physics, we will study them in much more detail when studying representation theory and Lie groups.

### Definition 18.1 (*General linear group $GL(n, \mathbb{R})$* )

Let  $M(n, \mathbb{R})$  be the set of all real invertible  $n \times n$  matrices. These form a group under matrix multiplication, called the **general linear group** denoted  $GL(n, \mathbb{R})$ .

*Proof.* We need to show that the group axioms are satisfied.

**Closure** Let  $A, B \in GL(n, \mathbb{R})$ , and let  $A^{-1}$  and  $B^{-1}$  be their inverses. Then, since  $\det A \neq 0$  and  $B \neq 0$ , we find that:

$$\det(AB) = \det A \cdot \det B \neq 0 \implies AB \in GL(n, \mathbb{R}) \quad (17.1.1)$$

**Associativity** matrix multiplication is associative.

**Identity** The identity matrix  $\mathbb{I}$  is the identity of  $GL(n, \mathbb{R})$ . Indeed  $\forall A \in GL(n, \mathbb{R})$ :

$$\mathbb{I}A = A\mathbb{I} = A \quad (17.1.2)$$

as desired. Moreover,  $\mathbb{I} \in GL(n, \mathbb{R})$ , since  $\det \mathbb{I} = 1 \neq 0$ .

**Inverses** Let  $A \in GL(n, \mathbb{R})$ , then it must have an inverse  $A^{-1}$ , such that:

$$AA^{-1} = A^{-1}A = \mathbb{I} \quad (17.1.3)$$

Moreover,  $A^{-1}$  is invertible (its inverse if  $A$ , so  $A^{-1} \in GL(n, \mathbb{R})$ ).

Therefore all the group axioms are satisfied, and  $GL(n, \mathbb{R})$  form a group under matrix multiplication. ■

Clearly,  $\dim GL(n, \mathbb{R}) = n^2$ . Indeed, the restriction that  $\det A \neq 0$  only restrict the values that the matrix elements cannot take, but it does not force some matrix elements to take a specific value. Hence they are spanned by the standard basis of  $Mat_n(\mathbb{R})$ .

**Definition 18.2 (Special linear group  $SL(n, \mathbb{R})$ )**

The subset of  $GL(n, \mathbb{R})$  comprising of all invertible  $n \times n$  matrices with unit determinant forms a subgroup of the general linear group, called the **special linear group**, and denoted  $SL(n, \mathbb{R})$ .

*Proof.* We need to show that the subgroup axioms are satisfied.

**Closure** Let  $A, B \in SL(n, \mathbb{R})$ , and let  $A^{-1}$  and  $B^{-1}$  be their inverses. Then, since  $\det A = 1$  and  $\det B = 1$ , we find that:

$$\det(AB) = \det A \cdot \det B = 1 \implies AB \in SL(n, \mathbb{R}) \quad (17.1.4)$$

**Identity** The identity matrix  $I \in SL(n, \mathbb{R})$  since  $\det I = 1$ .

**Inverses** Let  $A \in SL(n, \mathbb{R})$ , then it must have an inverse  $A^{-1}$ , and:

$$\det A^{-1} = \frac{1}{\det A} = 1 \quad (17.1.5)$$

so  $A^{-1} \in SL(n, \mathbb{R})$ .

Therefore all the subgroup axioms are satisfied, and  $SL(n, \mathbb{R})$  forms a sgroup under matrix multiplication. ■

Perhaps more difficult to see is the fact that  $\dim SL(n, \mathbb{R})$ . Indeed, the most general form of a matrix  $A \in SL(n, \mathbb{R})$ :

$$A = \begin{pmatrix} a_{11} & a_{12} & a_{13} & \dots & a_{1n} \\ a_{21} & a_{22} & a_{23} & \dots & a_{2n} \\ a_{31} & a_{32} & a_{33} & \dots & a_{3n} \\ \vdots & \vdots & \vdots & & \vdots \\ a_{n1} & a_{n2} & a_{n3} & \dots & a_{nn} \end{pmatrix} \quad (17.1.6)$$

The condition  $\det A = 1$  can be expanded as:

$$\det A = \sum_{\sigma \in S_n} \text{sgn}(\sigma) a_{\sigma(1),1} a_{\sigma(2),2} a_{\sigma(3),3} \dots a_{\sigma(n),n} \quad (17.1.7)$$

We can solve this equation for  $a_{nn}$  in terms of the other  $n^2 - 1$  components, so there are in total  $n^2 - 1$  independent elements in  $A$ , giving:

$$\dim SL(n, \mathbb{R}) = n^2 - 1 \quad (17.1.8)$$

The extension of these results for  $SL(n, \mathbb{C})$  is immediate:

$$\dim SL(n, \mathbb{C}) = 2 \cdot \dim SL(n, \mathbb{R}) = 2n^2 - 2 \quad (17.1.9)$$

**Definition 18.3 (Orthogonal group  $O(n, \mathbb{R})$ )**

The subset of  $GL(n, \mathbb{R})$  comprising of all invertible  $n \times n$  orthogonal matrices ( $A^T A = I$ ) forms a subgroup of the general linear group, called the **orthogonal group**, and denoted  $O(n, \mathbb{R})$ .

*Proof.* Firstly note that if  $A \in O(n, \mathbb{R})$ , then:

$$\det AA^T = (\det\{A\})^2 = \det \mathbb{I} = 1 \implies \det\{A\} = \pm 1 \quad (17.1.10)$$

so  $A$  must also be invertible.

We need to show that the subgroup axioms are satisfied.

**Closure** Let  $A, B \in O(n, \mathbb{R})$ . Then we find that:

$$(AB)(AB)^T = ABB^T A^T = \mathbb{I} \quad (17.1.11)$$

so  $AB \in O(n, \mathbb{R})$ .

**Identity** The identity matrix  $\mathbb{I} \in O(n, \mathbb{R})$  since  $\mathbb{I}^T = \mathbb{I}$ .

**Inverses** Let  $A \in O(n, \mathbb{R})$ . Then its inverse  $A^{-1}$  in  $GL(n, \mathbb{R})$  satisfies

$$(A^{-1})(A^{-1})^T = (A^{-1})(A^T)^{-1} = (A^T A)^{-1} = \mathbb{I}^{-1} = \mathbb{I} \quad (17.1.12)$$

so that  $A \in O(n, \mathbb{R})$  as desired. So  $A^{-1} \in SL(n, \mathbb{R})$ .

Therefore all the subgroup axioms are satisfied, and  $O(n, \mathbb{R})$  forms a group under matrix multiplication. ■

Again, finding the dimension of the orthogonal group of order  $n$  is slightly more involved. Consider the equation:

$$AA^T - \mathbb{I} = 0 \quad (17.1.13)$$

Note that  $(AA^T)^T = (A^T)^T A^T = AA^T$ , so  $AA^T$  only has  $\frac{1}{2}n(n+1)$  independent components, namely the  $n$  diagonal components, and then the lower/upper triangular components  $\frac{1}{2}(n^2 - n)$ . Hence can be expressed as:

$$\begin{pmatrix} b_{11} & \dots & \dots & B \\ \vdots & b_{22} & & \vdots \\ \vdots & & & \vdots \\ B & \dots & & b_{nn} \end{pmatrix} = 0 \quad (17.1.14)$$

where we absorbed  $\mathbb{I}$  into the diagonal  $b_{kk}$  components. As we said earlier the above matrix equation has  $\frac{1}{2}n(n+1)$  independent equations in  $a_{ij}$ , fixing  $\frac{1}{2}n(n+1)$  elements of  $A$ . Hence:

$$\dim O(n, \mathbb{R}) = n^2 - \frac{1}{2}n(n+1) = \frac{1}{2}n(n-1) \quad (17.1.15)$$

#### Definition 18.4 (Special orthogonal group $SU(n, \mathbb{R})$ )

The subset of  $O(n, \mathbb{R})$  comprising of all invertible  $n \times n$  orthogonal matrices with unit determinant forms a subgroup of the orthogonal group, called the **special orthogonal group**, and denoted  $SU(n, \mathbb{R})$ .

The proof is a combination of the proofs in definitions 18.2 and 18.3.

The dimension of  $SU(n, \mathbb{R})$  is surprisingly equal to the dimension of  $O(n, \mathbb{R})$ , despite the additional constraint that  $\det A = +1$ . Note that for  $A = O(n, \mathbb{R})$ , we have  $A = \pm 1$ , so we're only removing

the matrices with  $\det A = -1$ . The constraints however remain the same, so we still have  $\frac{1}{2}n(n-1)$  constraints.

#### Definition 18.4 (Unitary group $U(n, \mathbb{R})$ )

The subset of  $GL(n, \mathbb{C})$  comprising of all unitary  $n \times n$  matrices ( $AA^\dagger = \mathbb{I}$ ) forms a subgroup of the general linear group, called the **unitary group**, and denoted  $U(n, \mathbb{R})$ .

*Proof.* We need to show that the subgroup axioms are satisfied.

**Closure** Let  $A, B \in U(n, \mathbb{C})$ . Then we find that:

$$(AB)(AB)^\dagger = ABB^\dagger A^\dagger = \mathbb{I} \quad (17.1.16)$$

so  $AB \in O(n, \mathbb{C})$ .

**Identity** The identity matrix  $\mathbb{I} \in U(n, \mathbb{C})$  since  $\mathbb{I}\mathbb{I}^\dagger = \mathbb{I}$ .

**Inverses** Let  $A \in U(n, \mathbb{C})$ . Then its inverse  $A^{-1}$  in  $GL(n, \mathbb{C})$  satisfies

$$(A^{-1})(A^{-1})^\dagger = (A^{-1})(A^\dagger)^{-1} = (A^\dagger A)^{-1} = \mathbb{I}^{-1} = \mathbb{I} \quad (17.1.17)$$

so that  $A^{-1} \in U(n, \mathbb{C})$  as desired.

Therefore all the subgroup axioms are satisfied, and  $U(n, \mathbb{C})$  forms a group under matrix multiplication. ■

What is the dimension of  $U(n, \mathbb{C})$ ? Note that  $Mat_n(\mathbb{C})$  has dimension  $2n^2$ ,  $n^2$  free real parameters, and  $n^2$  free complex parameters. Now let's see how many constraints unitarity ( $AA^\dagger = \mathbb{I}$ ) imposes.

Note that  $AA^\dagger$  is hermitian, since  $(AA^\dagger)^\dagger = (A^\dagger)^\dagger A^\dagger = AA^\dagger$ . A hermitian matrix has a total of  $n^2$  free parameters,  $n$  diagonal (the diagonal components have to be real, since if they had an imaginary part, the hermitian conjugate would turn it negative, thus violating hermiticity), and  $n^2 - n$  off-diagonal ( $\frac{1}{2}(n^2-n)$  for the real part, and  $\frac{1}{2}(n^2-n)$  for the imaginary part). Therefore the hermiticity constraint consists of  $n^2$  independent equations, and thus fixes  $n^2$  elements. Hence:

$$\dim U(n, \mathbb{C}) = 2n^2 - n^2 = n^2 \quad (17.1.18)$$

#### Definition 18.4 (Special unitary group $SU(n, \mathbb{C})$ )

The subset of  $U(n, \mathbb{C})$  comprising of all invertible  $n \times n$  hermitian matrices with unit determinant forms a subgroup of the unitary group, called the **special unitary group**, and denoted  $SU(n, \mathbb{C})$ .

Note that the determinant of  $A \in U(n, \mathbb{R})$  is such that:

$$\det AA^\dagger = (\det A)(\det A)^* = |\det A|^2 = 1 \implies \det A = e^{i\theta}, \forall \theta \in [0, 2\pi) \quad (17.1.19)$$

whereas we are restricting  $\det A = 1 \implies \theta = 0$ . Unlike in the orthogonal and special orthogonal group case, where we were simply restricting the sign of the determinant, here we are restricting a whole continuum of values that the determinant can take, hence we have an additional constraint.

Therefore:

$$\dim \mathrm{SU}(n, \mathbb{C}) = n^2 - 1 \quad (17.1.20)$$

We summarize the dimensionalities of these matrix groups below:

**Theorem 18.5 (Dimensions of matrix groups)**

We have that:

- (i)  $\dim \mathrm{GL}(n, \mathbb{R}) = n^2$
- (ii)  $\dim \mathrm{SL}(n, \mathbb{R}) = n^2 - 1$
- (iii)  $\dim O(n, \mathbb{R}) = \frac{1}{2}n(n-1)$
- (iv)  $\dim \mathrm{SU}(n, \mathbb{R}) = \frac{1}{2}n(n-1)$
- (v)  $\dim U(n, \mathbb{C}) = n^2$
- (vi)  $\dim \mathrm{SU}(n, \mathbb{C}) = n^2 - 1$

## 17.2 Cosets

**Definition 18.6 (Left coset)**

Let  $H < G$  be a subgroup of  $G$ , and let  $g \in G$ . Then, the **left coset**  $gH$  of  $H$  in  $g$  is given by:

$$gH = \{gh : h \in H\} \quad (17.2.1)$$

and is the subset of  $G$  obtained by composing each element in  $H$  with  $g$  to the left.

**Example.** For example, let's try to find all the left cosets of the subgroup  $H = \{e, s\}$  in the group  $S(\Delta)$ .

We can use the Cayley table for  $S(\Delta)$ :

$$eH = e\{e, s\} = \{e, s\} \quad (17.2.2)$$

$$aH = a\{e, s\} = \{a, r\} \quad (17.2.3)$$

$$bH = b\{e, s\} = \{b, t\} \quad (17.2.4)$$

$$rH = r\{e, s\} = \{r, a\} = aH \quad (17.2.5)$$

$$sH = s\{e, s\} = \{s, e\} = eH \quad (17.2.6)$$

$$tH = t\{e, s\} = \{t, b\} = bH \quad (17.2.7)$$

so the distinct left cosets are  $\{e, s\}, \{a, r\}, \{b, t\}$ . ◀

**Example.** For example, let's try to find all the left cosets of the subgroup  $H = \{1, 2, 4\}$  in  $\mathbb{Z}_7^*$ .

We find that:

$$1H = 1\{1, 2, 4\} = \{1, 2, 4\} \quad (17.2.8)$$

$$2H = 2\{1, 2, 4\} = \{2, 4, 1\} \quad (17.2.9)$$

$$3H = 3\{1, 2, 4\} = \{3, 6, 5\} \quad (17.2.10)$$

$$4H = 4\{1, 2, 4\} = \{4, 1, 2\} \quad (17.2.11)$$

$$5H = 5\{1, 2, 4\} = \{5, 3, 6\} \quad (17.2.12)$$

$$6H = 6\{1, 2, 4\} = \{6, 5, 3\} \quad (17.2.13)$$

so we see that the distinct left cosets are  $\{1, 2, 4\}$  and  $\{3, 5, 6\}$ .  $\blacktriangleleft$

It is interesting to note that the distinct left cosets of  $\{e, s\}$  in  $S(\Delta)$  and the distinct left cosets of  $\{1, 2, 4\}$  in  $\mathbb{Z}_7^*$  partition their respective groups.

For example  $\{e, s\} \cup \{a, r\} \cup \{b, t\} = S(\Delta)$ , and all three sets are disjoint.

This is not a coincidence, it turns out that all distinct left cosets are partitions.

### Theorem 18.7 (Left coset partition)

Let  $H < G$  be a subgroup of a group  $G$ . Then the distinct left cosets of  $H$  in  $G$  form a partition of  $G$ .

*Proof.* Recall that the equivalence classes of an equivalence relation on some set  $X$  form a partition of the set  $X$ . Consequently, if we can show that a particular operation has left cosets of a subgroup  $H$  in a group  $G$ , then immediately we find that the left cosets form a partition.

Lemma. Let  $\sim$  be the relation defined on  $G$  by:

$$x \sim y \text{ if } x \in yH \quad (17.2.14)$$

Then  $\sim$  is an equivalence relation. Indeed:

- (i) **Reflexive property:** let  $x \in G$ , we have to show that  $x \sim x$ , that is,  $x \in xH$ . This is clearly true since  $x = xe$ , and  $e \in H$  due to the subgroup axioms.
- (ii) **Symmetric property:** let  $x, y \in G$ , and let  $x \sim y$ , so that  $x \in yH$ . Therefore,  $\exists h \in H$  such that  $x = yh \implies y = xh^{-1}$ . Now, due to the inverses property of subgroups,  $h^{-1} \in H \implies y \in xH$  so  $y \sim x$ .
- (iii) **Transitive property:** let  $x, y, z \in G$ , and suppose  $x \sim y, y \sim z$ , so  $x \in yH$  and  $y = zH$ . Therefore,  $\exists h_1, h_2 \in H$  such that:

$$x = yh_1 \text{ and } y = zh_2 \implies x = zh_1h_2 \quad (17.2.15)$$

Since  $h_1h_2 \in H$ , we find that  $x \in zH$ , and thus  $x \sim z$ .

We have therefore shown that  $\sim$  is an equivalence relation. Each element  $x \in G$  has equivalence class:

$$[x] = \{y \in G : y \sim x\} = \{y \in G : y \in xH\} = xH \quad (17.2.16)$$

which are the left cosets of  $H$  in  $G$ . It follows then that the left cosets (which are the equivalence classes) form a partition of  $G$ . ■

**Proposition 18.9 (Properties of left cosets)** Let  $H < G$ , then:

- (i)  $\forall g \in G, g \in gH$
- (ii) one of the left cosets of  $H$  in  $G$  is  $H$
- (iii) Any two left cosets  $g_1H$  and  $g_2H$  are either the same or disjoint
- (iv) If  $|H| < \infty$ , then each left coset  $gH$  has the same number of elements.

*Proof.* (i) Trivial, since  $e \in H$  due to subgroup axioms.

(ii)  $H = \{h : h \in H\} = \{eh : h \in H\} = eH$  so  $H$  is indeed a left coset.

(iii) Immediate from theorem 18.8.

(iv) Let  $|H| = m$  so that  $H = \{h_1, h_2, \dots, h_m\}$ , and let  $g \in G$ . Then:

$$gH = \{gh_1, gh_2, \dots, gh_m\} \quad (17.2.17)$$

Suppose  $gh_i = gh_j$ , then by the left cancellation law  $h_i = h_j$ , so it follows that  $|gH| = |H| = m$ . ■

**Example.** Let's try to find all the left cosets of  $H = \{e, a, b, c\}$  in  $S(\square)$ . The remaining elements are  $r, s, t, u$ , and the remaining left cosets all contain 4 elements. So the only way to make them disjoint is if we have  $\{r, s, t, u\}$  as the other coset. Hence the distinct left cosets are  $\{e, a, b, c\}$  and  $\{r, s, t, u\}$ .

To partition a finite group  $G$  into left cosets of some subgroup  $H < G$ :

- (i)  $H$  is the first cosets
- (ii) find an element in  $G$  that is not in  $H$ , and determine  $gH$ .
- (iii) repeat until all elements in  $G$  have been exhausted.

**Example.** Consider the group  $U_{20} = \{1, 3, 7, 9, 11, 13, 17\}$  and  $H = \{1, 19\}$ .

Firstly,  $H$  is one of the left cosets of  $H$  in  $U_{20}$ , so  $\{1, 19\}$ . One remaining element is 3, and its corresponding left coset is:

$$3H = 3\{1, 19\} = \{3, 17\} \quad (17.2.18)$$

An element missing from both  $H$  and  $3H$  is 7, and its corresponding left coset is:

$$7H = 7\{1, 19\} = \{7, 13\} \quad (17.2.19)$$

The final missing element from all these left cosets is 9:

$$9H = 9\{1, 19\} = \{9, 11\} \quad (17.2.20)$$

So we have found that the distinct left cosets are  $\{1, 19\}, \{3, 17\}, \{7, 13\}$  and  $\{9, 11\}$ . ◀

**Example.** Let's try to partition the alternating group of order 4:

$$A_4 = \{e, (1 2)(3 4), (1 3)(2 4), (1 4)(2 3), (1 2 3), (1 3 2), (1 2 4), (1 4 2), \quad (17.2.21)$$

$$(1 3 4), (1 4 3), (2 3 4), (2 4 3)\} \quad (17.2.22)$$

into left cosets of the subgroup  $H = \{e, (1 2)(3 4), (1 3)(2 4), (1 4)(2 3)\}$ .

Firstly, we note that  $H$  itself is a left coset in  $A_4$ . A remaining element not included is for example  $(1 2 3)$ , and its associated left coset is:

$$(1 2 3)H = \{(1 2 3), (1 3 4), (2 4 3), (1 4 2)\} \quad (17.2.23)$$

Therefore, the only remaining permutations are  $(1 3 2), (1 2 4), (1 4 3), (2 3 4)$ . Since the left cosets must contain as many elements as  $H$ , that is 4, and they must be disjoint, there is only one remaining left coset  $\{(1 3 2), (1 2 4), (1 4 3), (2 3 4)\}$ .

Hence the partition of  $A_4$  into the left cosets of  $H$  is:

$$H \cup \{(1 3 2), (1 2 4), (1 4 3), (2 3 4)\} = A_4 \quad (17.2.24)$$



We can use the left cosets to provide a more efficient prove of Lagrange's theorem.

### Theorem 17.1 (Lagrange's Theorem)

Let  $G$  be a finite group, and let  $H \leq G$ . Then  $\text{ord}(H)|\text{ord}(G)$ , that is, the order of  $H$  divides the order of  $G$ .

*Proof.* Let  $|G| = n$  and  $|H| = m$ , and let the number of left cosets of  $H$  in  $G$  be  $k$ . Since each left coset has  $m$  elements, and they partition  $G$  into disjoint sets, it follows that  $n = m \cdot k$ , that is,  $|H|$  divides  $|G|$ . ■

## 17.3 Right cosets

### Definition 18.10 (Right cosets)

Let  $H < G$  be a subgroup of  $G$ , and let  $g \in G$ . Then, the **right coset**  $Hg$  of  $H$  in  $g$  is given by:

$$Hg = \{hg : h \in H\} \quad (17.3.1)$$

and is the subset of  $G$  obtained by composing each element of  $H$  with  $g$  to the right.

**Example.** Let's find all the right cosets of  $H = \{e, r\}$  in  $S(\square)$ .

We find that:

$$He = \{e, r\}e = \{e, r\} \quad (17.3.2)$$

$$Ha = \{e, r\}a = \{a, u\} \quad (17.3.3)$$

$$Hb = \{e, r\}b = \{b, t\} \quad (17.3.4)$$

$$Hc = \{e, r\}c = \{c, s\} \quad (17.3.5)$$

$$Hr = \{e, r\}r = \{r, e\} \quad (17.3.6)$$

$$Hs = \{e, r\}s = \{s, c\} \quad (17.3.7)$$

$$Ht = \{e, r\}t = \{t, b\} \quad (17.3.8)$$

$$Hu = \{e, r\}u = \{u, a\} \quad (17.3.9)$$

so the distinct right cosets are:  $\{e, r\}, \{a, u\}, \{b, t\}, \{c, s\}$ . ◀

It is interesting to note that the right coset of  $H$  in some element  $g$  is not necessarily equal to the left coset of  $H$  in  $g$ .

All the results proven for left cosets are easily proven for right cosets as well.

### Theorem 18.11 (Right coset partition)

Let  $H < G$  be a subgroup of  $G$ , then the distinct right cosets of  $H$  in  $G$  form a partition of  $G$ .

### Proposition 18.12 (Properties of right coset)

Let  $H < G$ , then:

- (i)  $g \in Hg$ , for all  $g \in G$
- (ii)  $H$  is a right coset of  $H$  in  $G$
- (iii) two right cosets  $Hg_1$  and  $Hg_2$  are either the same set or disjoint
- (iv) if  $|H| < \infty$  then  $\forall g \in G |Hg| = |H|$ .

**Example.** Let's try to partition  $S(\Delta)$  into right cosets of  $H = \{e, s\}$ .

Firstly, we note that all right cosets must have 2 elements, and one of these right cosets is  $H$  itself. An element in  $G$  that does not belong to  $H$  is  $a$ , and its right coset is:

$$Ha = \{e, s\}a = \{a, t\} \quad (17.3.10)$$

Now another element in  $G$  that does not belong to these two cosets is  $b$ , so:

$$Hb = \{e, s\}b = \{b, r\} \quad (17.3.11)$$

We have exhausted all elements of  $S(\Delta)$ , so we may write the partition of the latter as:

$$S(\Delta) = \{e, s\} \cup \{a, t\} \cup \{b, r\} \quad (17.3.12)$$

Note that if  $G$  is abelian and  $H < G$ , then  $H$  is also abelian and thus:

$$Hg = \{hg : \forall h \in H\} = \{gh : \forall h \in H\} = gH \quad (17.3.13)$$

so the left and right cosets are the same.

**Theorem 18.13 (Transforming left coset to right coset)**

Let  $H < G$ . Then if every element in the partition of  $G$  into left cosets of  $H$  is replaced by its inverse, the result is the partition of  $G$  into right cosets of  $H$  (and viceversa).

*Proof.* Consider a pair of elements  $x, y$  in the same left coset of  $H$ . Therefore,

Suppose  $x \in gH$  for some  $g \in G$ . Then, we need to prove that  $x^{-1} \in Hg'$  for some  $g' \in G$ . Indeed:

$$gH = \{gh : h \in H\} \text{ and } \{(gh)^{-1} : h \in H\} = \{h^{-1}g^{-1} : h \in H\} = \{hg^{-1} : h \in H\} = Hg^{-1} \quad (17.3.14)$$

due to the inverse axiom of subgroups<sup>1</sup>. It follows that if we replace every element in  $gH$  by its inverse, we get a right coset  $Hg^{-1}$ . ■

Immediately, we find that since the act of replacing each element by its inverse is bijective, the number of distinct left and right cosets is the same.

**Proposition 18.14 (Number of left and right cosets)** Let  $H < G$ , then the number of distinct left cosets of  $H$  in  $G$  is equal to the number of distinct right cosets of  $H$  in  $G$ .

**Definition 18.15 (Index)** Let  $H < G$ . The **index** of  $H$  in  $G$  is the number of distinct left cosets (or equivalently distinct right cosets) of  $H$  in  $G$ .

For finite groups, we have a nice expression for the index.

**Proposition 18.16 (Finite index)**

Let  $H < G$ , then the index of  $H$  in  $G$  is  $\frac{|G|}{|H|}$ .

*Proof.* The left cosets of  $H$  partition  $G$ , and since each left coset has  $|H|$  elements, the number of left cosets (the index) must be  $\frac{|G|}{|H|}$ . ■

**Example.** Let's partition  $(2\mathbb{Z}, +)$  into cosets of  $6\mathbb{Z}$ .

Note that  $2\mathbb{Z} = \{\dots, -6, -4, -2, 0, 2, 4, 6, \dots\}$ , and  $6\mathbb{Z} = \{\dots, -18, -12, -6, 0, 6, 12, 18, \dots\}$ . We know that  $6\mathbb{Z} < 2\mathbb{Z}$  since  $6\mathbb{Z}$  is generated by  $6 \in 2\mathbb{Z}$ , and is therefore a cyclic subgroup of  $2\mathbb{Z}$ .

<sup>1</sup>indeed suppose

$$H = \{h_1, h_2, \dots, h_n, h_1^{-1}, h_2^{-1}, \dots, h_n^{-1}\}$$

then

$$\{h^{-1} : h \in H\} = \{h_1^{-1}, h_2^{-1}, \dots, h_n^{-1}, h_1, h_2, \dots, h_n\} = H$$

Now note that  $6\mathbb{Z}$  is itself a coset so:

$$0 + 6\mathbb{Z} = \{6k : k \in \mathbb{Z}\} \quad (17.3.15)$$

An element that was not included is 2:

$$2 + 6\mathbb{Z} = \{\dots, -16, -10, -4, 2, 8, 14, 20, \dots\} = \{2 + 6k : k \in \mathbb{Z}\} \quad (17.3.16)$$

An element that was not included is 4, so:

$$4 + 6\mathbb{Z} = \{\dots, -14, -8, -2, 4, 10, 16, 22, \dots\} = \{4 + 6k : k \in \mathbb{Z}\} \quad (17.3.17)$$

Note that:

$$\{2(3k+2) : k \in \mathbb{Z}\} \cup \{2(3k+1) : k \in \mathbb{Z}\} \cup \{2(3k) : k \in \mathbb{Z}\} = \{2k : k \in \mathbb{Z}\} = 2\mathbb{Z} \quad (17.3.18)$$

so we do indeed have a partition. Moreover the index of  $6\mathbb{Z}$  in  $2\mathbb{Z}$  is 3, something that we could not have found using proposition 18.16.  $\blacktriangleleft$

## 17.4 Normal subgroups

Although this doesn't generally happen, for certain special groups the left coset partition and right coset partition can be exactly the same.

### Definition 18.17 (Normal subgroup)

Let  $G$  be a group and let  $H < G$ . We say that  $H$  is a **normal subgroup** if the left coset partition of  $G$  in  $H$  and the right coset partition of  $G$  in  $H$  are the same. We say that  $H$  is normal in  $G$  as well, denoted as  $H \trianglelefteq G$ .

**Proposition 18.18 (Standard normal subgroups)** For any group  $G$ :

- (i) the trivial subgroup  $\{e\}$
  - (ii)  $G$
- are both normal subgroups.

*Proof.* In the first case the left coset partition contains one-element subsets of  $G$ , and so does the right coset partition. In the second case the left coset partition is simply  $G$ , and so is the right coset partition.  $\blacksquare$

**Example.** Consider  $A_4$  and  $H = \{e, (1 2)(3 4), (1 3)(2 4), (1 4)(2 3)\}$ . The partition of  $A_4$  into the left cosets of  $H$  was found previously to be:

$$\{e, (1 2)(3 4), (1 3)(2 4), (1 4)(2 3)\} \cup \{(1 3 2), (1 2 4), (1 4 3), (2 3 4)\} = A_4 \quad (17.4.1)$$

Now let's try to find the right coset partition. We can do this by simply replacing each element in the left coset partition by its inverse:

$$\{e, (1 2)(3 4), (1 3)(2 4), (1 4)(2 3)\} \cup \{(1 2 3), (1 4 2), (1 3 4), (2 4 3)\} = A_4 \quad (17.4.2)$$

Note that the left and right coset partition are identical, hence  $H$  is indeed a normal subgroup of  $G$ .  $\blacktriangleleft$

**Theorem 18.19 (Normal subgroups of abelian groups)** Every subgroup  $H$  of an abelian group  $G$  is normal.

*Proof.* Let  $H < G$ , and let  $g \in G$ . Note that:

$$gH = \{gh : h \in H\} = \{hg : h \in H\} = Hg \quad (17.4.3)$$

thus  $H$  is normal.  $\blacksquare$

**Proposition 18.20 (Subgroups of index 2)** Every subgroup of index 2 in a group is a normal subgroup

*Proof.* Let  $H < G$  so that it has exactly two left cosets and exactly two right cosets in  $G$ . One of these cosets is  $H$  itself, so the other coset must be  $G \setminus H$  both for the left and right coset partition. The two partitions are therefore identical, proving that  $H$  is a normal subgroup of  $H$ .  $\blacksquare$

**Example.** If  $n \geq 2$ , it follows that  $A_n$  is a normal subgroup of  $S_n$ , since  $\frac{|S_n|}{|A_n|} = 2$ , so  $A_n$  is normal to  $S_n$ . If instead  $n = 1$ , then  $A_n = S_n$ , and from proposition 18.18 we see that  $A_1$  is normal to  $S_1$ .  $\blacktriangleleft$

**Example.** Consider  $S(\square)$  and  $S^+(\square)$ . Since the number of direct symmetries is equal to the indirect symmetries, we find that the index of  $S^+(\square)$  in  $S(\square)$  is  $\frac{|S(\square)|}{|S^+(\square)|} = 2$ , and thus  $S^+(\square)$  is a normal subgroup of  $S(\square)$ .  $\blacktriangleleft$

**Proposition 18.21 (Equivalent condition for normality)**

Let  $H < g$ , then  $H$  is normal in  $G$  iff  $gH = Hg$ ,  $\forall g \in G$ .

*Proof.*  $\Rightarrow$  Suppose that  $H$  is normal in  $G$ , and let  $g \in G$ . Then  $g \in gH$  and  $g \in Hg$ . Now since the left and right coset partitions are identical, and the cosets are disjoint, we must have that since  $g$  belongs to both  $Hg$  and  $gH$ ,  $Hg = gH$  as desired.

$\Leftarrow$  Suppose  $gH = Hg$ , then it follows that the left coset and right coset partitions are the same.  $\blacksquare$

# Unit E2: Quotient groups and conjugacy

## 18.1 Quotient groups

**Definition 19.1 (Set composition)** Let  $G$  be a group, then the operation  $\cdot$ , called a **set composition** in  $G$ , is defined as:

$$X \cdot Y = \{xy : x \in X, y \in Y\} \quad (18.1.1)$$

where  $X, Y \subseteq G$ .

Note that for an arbitrary group, set composition is not necessarily commutative. Only for Abelian groups is set composition is commutative.

Consider the following Cayley table for the cosets of the normal subgroup  $\{e, b\}$  in  $S(\square)$ :

$\cdot$	$\{e, b\}$	$\{a, c\}$	$\{r, t\}$	$\{s, u\}$
$\{e, b\}$	$\{e, b\}$	$\{a, c\}$	$\{r, t\}$	$\{s, u\}$
$\{a, c\}$	$\{a, c\}$	$\{e, b\}$	$\{s, u\}$	$\{r, t\}$
$\{r, t\}$	$\{r, t\}$	$\{u, s\}$	$\{e, b\}$	$\{a, c\}$
$\{s, u\}$	$\{s, u\}$	$\{r, t\}$	$\{a, c\}$	$\{e, b\}$

Interestingly all the sets in the body of the table are also cosets of  $\{e, b\}$ . Therefore the set of these cosets are closed under set composition. It turns out that we can extend this result more generally for any cosets of a normal subgroup  $N$  of a group  $G$ .

**Theorem 19.2 (Closure of cosets under set composition)** Let  $N \trianglelefteq G$ , then:

$$xN \cdot yN = (xy)N, \forall x, y \in G \quad (18.1.2)$$

*Proof.* Let us firstly show that  $xN \cdot yN \subseteq (xy)N$ . Indeed, let  $z \in xN \cdot yN$ , so that there are some  $n_1, n_2 \in N$  such that:

$$z = xn_1yn_2 \quad (18.1.3)$$

Since  $N$  is a normal subgroup,  $n_1y \in Ny \implies n_1y \in yN$  so

$$n_1y = yn_3, n_3 \in N \quad (18.1.4)$$

Then:

$$z = xyn_3n_2 = xyn, \quad n = n_3n_2 \in N \quad (18.1.5)$$

implying that  $z \in xyN$ , as desired.

Let us now show that  $(xy)N \subseteq xN \cdot yN$ . Suppose that  $z \in (xy)N$ , so that there is some  $n_1 \in N$  such that:

$$z = xyn_1 \quad (18.1.6)$$

Since  $x \in xN$  and  $yn_1 \in yN$ , we find that  $z \in xN \cdot yN$ , as desired. ■

Due to this theorem, we see that if we perform set composition on the cosets of some normal subgroup  $N \trianglelefteq G$  in  $G$ , then the result will be another coset.

If we examine the Cayley table for the cosets of  $\{e, b\}$  in  $S(\square)$ , we find that not only is closure under  $\cdot$  satisfied, all the other group properties are also satisfied! Associativity of set composition follows from associativity of composition in  $S(\square)$ . Moreover, the identity element can be verified to be  $\{e, b\}$  so that all cosets are self-inverse.

Let us prove that the cosets of a normal subgroup form a group under set composition in the most general case.

**Theorem 19.3 (Group of cosets)** Let  $N \trianglelefteq G$ , then the set of cosets of  $N$  in  $G$  forms a group under set composition. This group is called the **quotient group** of  $G$  by  $N$ , denoted  $G \setminus N$ , and often read  $G$  mod  $N$  for brevity.

*Proof.*

**Closure** We have proven closure in Theorem 19.2

**Associativity** Let  $xN, yN, zN$  be cosets of  $N$  in  $G$ . Then:

$$xN \cdot (yN \cdot zN) = xN \cdot (yz)N = (x(yz))N = (xyz)N \quad (18.1.7)$$

since  $G$  is a group, and therefore its operation is associative. Similarly

$$(xN \cdot yN) \cdot zN = ((xy)z)N = (xyz)N \quad (18.1.8)$$

**Identity** We prove that  $eN$  is the identity element. Indeed:

$$eN \cdot xN = (ex)N = xN, \quad \forall x \in G \quad (18.1.9)$$

and similarly:

$$xN \cdot eN = (xe)N = xN, \quad \forall x \in G \quad (18.1.10)$$

as desired.

**Inverses** Suppose  $xN$  is a coset of  $N$  in  $G$ . Then:

$$x^{-1}N \cdot xN = (x^{-1}x)N = eN, \quad \forall x \in G \quad (18.1.11)$$

and similarly:

$$xN \cdot x^{-1}N = (xx^{-1})N = eN, \quad \forall x \in G \quad (18.1.12)$$

Therefore,  $x^{-1}N$  is the inverse of  $xN$ , and must belong to the set of cosets since  $x \in G \implies x^{-1} \in G$

■

Note that if  $G$  is finite, then  $|G \setminus N| = \frac{|G|}{|N|}$  is the number of cosets of  $N$  in  $G$ .

We can use our knowledge of normal subgroups and quotient groups to explain the "block" effect in the Cayley table of  $S(\mathcal{F})$ .

Indeed, note that since  $S^+(\mathcal{F})$  is a normal subgroup, we can construct its cosets in  $S(\mathcal{F})$ . Since  $S^+(\mathcal{F})$  has index 2, there will be two such cosets, with it being one of them. The remaining coset must therefore be the set of indirect symmetries  $S^-(\mathcal{F})$ . Consequently:

.	$S^+(\mathcal{F})$	$S^-(\mathcal{F})$
$S^+(\mathcal{F})$	$S^+(\mathcal{F})$	$S^-(\mathcal{F})$
$S^-(\mathcal{F})$	$S^-(\mathcal{F})$	$S^+(\mathcal{F})$

If we now expand  $S^+(\mathcal{F})$  and  $S^-(\mathcal{F})$  into its various components, then we find the blocks:

$\circ$	$e \ a \ b \ c \ r \ s \ t \ u$	$\circ$	direct	indirect
$e$	$e \ a \ b \ c \ r \ s \ t \ u$	direct	direct	indirect
$a$	$a \ b \ c \ e \ s \ t \ u \ r$			
$b$	$b \ c \ e \ a \ t \ u \ r \ s$			
$c$	$c \ e \ a \ b \ u \ r \ s \ t$			
$r$	$r \ u \ t \ s \ e \ c \ b \ a$	indirect	indirect	direct
$s$	$s \ r \ u \ t \ a \ e \ c \ b$			
$t$	$t \ s \ r \ u \ b \ a \ e \ c$			
$u$	$u \ t \ s \ r \ c \ b \ a \ e$			

Figure 18.1. Block effect in  $S(\square)$

**Example.** Consider the subgroup  $H = \langle 6 \rangle$  of  $\mathbb{Z}_{12}$ . The elements of  $H$  are:

$$\langle 6 \rangle = \{0, 6\} \implies \text{ord}(6) = 2 \quad (18.1.13)$$

Since  $\mathbb{Z}_{12}$  is an abelian group, all of its subgroups are normal, including  $\langle 6 \rangle$ , which is its cyclic subgroup of order 3.

The cosets of  $H$  in  $\mathbb{Z}_{12}$  must be:

$$H = \{0, 6\} \quad (18.1.14)$$

$$1 + H = \{1, 7\} \quad (18.1.15)$$

$$2 + H = \{2, 8\} \quad (18.1.16)$$

$$3 + H = \{3, 9\} \quad 4 + H = \{4, 10\} \quad (18.1.17)$$

$$5 + H = \{5, 11\} \quad (18.1.18)$$

The quotient group, formed by the above cosets, must then have Cayley table:

$+$	$H$	$1 + H$	$2 + H$	$3 + H$	$4 + H$	$5 + H$
$H$	$H$	$1 + H$	$2 + H$	$3 + H$	$4 + H$	$5 + H$
$1 + H$	$1 + H$	$2 + H$	$3 + H$	$4 + H$	$5 + H$	$H$
$2 + H$	$2 + H$	$3 + H$	$4 + H$	$5 + H$	$H$	$1 + H$
$3 + H$	$3 + H$	$4 + H$	$5 + H$	$H$	$1 + H$	$2 + H$
$4 + H$	$4 + H$	$5 + H$	$H$	$1 + H$	$2 + H$	$3 + H$
$5 + H$	$5 + H$	$H$	$1 + H$	$2 + H$	$3 + H$	$4 + H$

For example,  $(4 + H) + (3 + H) = (4 +_{12} 3)H = 7 + H = 1 + H$ . We clearly see that  $H$  is the identity element, and that  $H, 3 + H$  are self inverse, whereas  $2 + H$  and  $4 + H$  are inverses of each other and so are  $1 + H$  and  $5 + H$ .

We see that this Cayley table has the same structure as the sixth order cyclic group  $C_6$ .  $\blacktriangleleft$

## 18.2 Quotient group of infinite groups

Consider the group  $\mathbb{Z} \setminus 4\mathbb{Z}$ , that is, the set of cosets of  $4\mathbb{Z} = \{4n : n \in \mathbb{Z}\}$  in  $\mathbb{Z}$ . Its elements are:

$$4\mathbb{Z} \tag{18.2.1}$$

$$1 + 4\mathbb{Z} = \{1 + 4n : n \in \mathbb{Z}\} \tag{18.2.2}$$

$$2 + 4\mathbb{Z} = \{2 + 4n : n \in \mathbb{Z}\} \tag{18.2.3}$$

$$3 + 4\mathbb{Z} = \{3 + 4n : n \in \mathbb{Z}\} \tag{18.2.4}$$

$$(18.2.5)$$

We know that these are all the cosets since they partition  $\mathbb{Z}$ . We may construct the Cayley table for  $\mathbb{Z} \setminus 4\mathbb{Z}$  as:

$+$	$4\mathbb{Z}$	$1 + 4\mathbb{Z}$	$2 + 4\mathbb{Z}$	$3 + 4\mathbb{Z}$
$4\mathbb{Z}$	$4\mathbb{Z}$	$1 + 4\mathbb{Z}$	$2 + 4\mathbb{Z}$	$3 + 4\mathbb{Z}$
$1 + 4\mathbb{Z}$	$2 + 4\mathbb{Z}$	$3 + 4\mathbb{Z}$	$4\mathbb{Z}$	
$2 + 4\mathbb{Z}$	$2 + 4\mathbb{Z}$	$3 + 4\mathbb{Z}$	$4\mathbb{Z}$	$1 + 4\mathbb{Z}$
$3 + 4\mathbb{Z}$	$3 + 4\mathbb{Z}$	$4\mathbb{Z}$	$1 + 4\mathbb{Z}$	$2 + 4\mathbb{Z}$

Note that this has the exact same structure as the Cayley table for  $\mathbb{Z}_4$ , so  $\mathbb{Z}_4 \cong \mathbb{Z} \setminus 4\mathbb{Z}$  through the isomorphism:

$$\phi : \mathbb{Z}/4\mathbb{Z} \rightarrow \mathbb{Z}_4 \tag{18.2.6}$$

$$a + 4\mathbb{Z} \mapsto a \forall a \in \mathbb{Z}_4 \tag{18.2.7}$$

We can prove this result more generally.

### Proposition 19.4 ( $\mathbb{Z} \setminus n\mathbb{Z} \cong \mathbb{Z}_n$ )

For  $n \geq 2$ , then  $\mathbb{Z} \setminus n\mathbb{Z} \cong \mathbb{Z}_n$ . One isomorphism between them is:

$$\phi : \mathbb{Z}/n\mathbb{Z} \rightarrow \mathbb{Z}_n \tag{18.2.8}$$

$$a + n\mathbb{Z} \mapsto a, \forall a \in \mathbb{Z}_n \tag{18.2.9}$$

*Proof.* Let us firstly find the distinct cosets of  $n\mathbb{Z}$ . We have that:

$$a + n\mathbb{Z} = b + n\mathbb{Z} \iff a \in b + n\mathbb{Z} \iff a \equiv b \pmod{n} \quad (18.2.10)$$

so the distinct cosets must be:

$$n\mathbb{Z}, 1 + n\mathbb{Z}, \dots, (n - 1) + n\mathbb{Z} \quad (18.2.11)$$

It then follows that  $\phi$  is bijective. Indeed, it is injective since

$$\phi(a + n\mathbb{Z}) = \phi(b + n\mathbb{Z}) \implies a \equiv b \pmod{n} \implies a + n\mathbb{Z} = b + n\mathbb{Z} \quad (18.2.12)$$

It is also surjective since:

$$a \in \mathbb{Z}_n \implies a \leq n - 1 \implies \phi(a + n\mathbb{Z}) = a \quad (18.2.13)$$

Finally, for  $a, b \in \mathbb{Z}_n$  we find that:

$$\phi((a + n\mathbb{Z}) + (b + n\mathbb{Z})) = \phi((a +_n b) + n\mathbb{Z}) \quad (18.2.14)$$

$$= \phi(c + n\mathbb{Z}) \quad (18.2.15)$$

$$= c \quad (18.2.16)$$

where  $c \equiv a + b \pmod{n}$ . Moreover:

$$\phi(a + n\mathbb{Z}) +_n \phi(b + n\mathbb{Z}) = a +_n b = c \quad (18.2.17)$$

so that:

$$\phi((a + n\mathbb{Z}) + (b + n\mathbb{Z})) = \phi(a + n\mathbb{Z}) +_n \phi(b + n\mathbb{Z}) \quad (18.2.18)$$

as desired. ■

**Example.** Consider the group  $\mathbb{Z}/6\mathbb{Z}$ . We have established that:

$$\phi : \mathbb{Z}/6\mathbb{Z} \rightarrow \mathbb{Z}_6 \quad (18.2.19)$$

$$a + n\mathbb{Z} \mapsto a, \forall a \in \mathbb{Z}_6 \quad (18.2.20)$$

is an isomorphism. Since  $\mathbb{Z}_6$  is cyclic, it follows that if  $\phi(g)$  is a generator of  $\mathbb{Z}_6$  then  $g$  must be a generator of  $\mathbb{Z}/6\mathbb{Z}$ . Now the generators of  $\mathbb{Z}_6$  are 1 and 5 (integers coprime to 6), which are the images of  $1 + 6\mathbb{Z}$  and  $5 + \mathbb{Z}_6$ . The latter two must therefore be generators of  $\mathbb{Z}/6\mathbb{Z}$ . ◀

## 18.3 Conjugacy

### Definition 19.5 (Conjugacy)

Let  $x, y \in G$ . Then  $y$  is a conjugate of  $x$  in  $G$  if there exists some  $g \in G$  such that:

$$y = gxg^{-1} \quad (18.3.1)$$

**Proposition 19.6 (Powers of conjugate elements)** Let  $x, y, g \in G$  such that  $y = gxg^{-1}$ . Then  $y^n = gx^n g^{-1}$  for all positive integers  $n$ .

*Proof.* We proceed by mathematical induction. Let  $P(n) : y^n = gx^n g^{-1}$ , then  $P(1)$  is clearly true. Moreover, suppose that  $P(k)$  is true for some positive integer  $k$ . Then:

$$y^k = gx^k g^{-1} \implies y^{k+1} = gx^k g^{-1} y \quad (18.3.2)$$

$$= gx^k g^{-1} gxg^{-1} \quad (18.3.3)$$

$$= gx^k xg^{-1} = gx^{k+1} g^{-1} \quad (18.3.4)$$

so  $P(k + 1)$  must be true. Hence, by the principle of mathematical induction, we have that  $P(n)$  is true for any positive integer  $n$ . ■

### Theorem 19.7 (Order of conjugate elements)

Let  $x, y \in G$  be conjugate elements. Then either  $x, y$  have the same finite order or they both have infinite order.

*Proof.* There exists some  $g \in G$  such that  $y = gxg^{-1}$ . Suppose  $x^n = e$ , then:

$$y^n = gx^n g^{-1} = geg^{-1} = e \quad (18.3.5)$$

Similarly, suppose that  $y^n = e$ , then:

$$x^n = g^{-1}y^n g = g^{-1}eg = e \quad (18.3.6)$$

It follows that if there are positive integers  $n$  such that  $x^n = e$ , then  $y^n = e$  and vice versa. So either  $x, y$  have the same finite order or they have infinite order. ■

### Definition 19.8 (Conjugacy class)

Let  $x \in G$ , then the **conjugacy class** of  $x$  in  $G$  is the set of all elements in  $G$  that are conjugate to  $x$ :

$$\{gxg^{-1} : g \in G\} \quad (18.3.7)$$

### Theorem 19.9 (Conjugacy class partition)

Let  $G$  be a group, then the conjugation relation is an equivalence relation on  $G$ . Consequently, the distinct conjugacy classes form a partition of  $G$ .

Reflexive: let  $x \in G$ , then  $x = exe^{-1}$ , so  $x$  is conjugate to itself.

Symmetric: let  $x, y \in G$ , and suppose  $x$  is conjugate to  $y$ . That is, there exists some element  $g \in G$ :

$$x = gyg^{-1} \implies y = g^{-1}xg = g^{-1}x(g^{-1})^{-1} \quad (18.3.8)$$

so it follows that  $y$  is conjugate to  $x$ .

Transitive: let  $x, y, z \in G$ , and suppose  $x$  is conjugate to  $y$ , and  $y$  is conjugate to  $z$ , so that:

$$x = g_1 y g_1^{-1}, \quad y = g_2 z g_2^{-1} \quad (18.3.9)$$

for some  $g_1, g_2 \in G$ . It follows that:

$$x = g_1 g_2 z g_2^{-1} g_1^{-1} = g_3 z g_3^{-1} \quad (18.3.10)$$

where  $g_3 = g_1 g_2 \in G$ . Hence  $x$  is conjugate to  $z$ .

Since the equivalence classes of an equivalence class on some set partition the set, we find that the conjugacy classes of  $G$  form a partition of the group. ■

It is important to remember that two elements of different order cannot be conjugate to each other. This was proven in Theorem 19.7, and gives us a useful strategy when trying to partition a group into its conjugacy classes. We show this strategy in the example below.

**Example.** Consider the group  $S(\Delta)$ . Since the conjugacy classes of this group must contain elements of the same order, we can start by partitioning  $S(\Delta)$  into sets of all elements of the same order:

$$\{e\}, \{r, s, t\}, \{a, b\} \quad (18.3.11)$$

where  $e$  has order 1,  $r, s, t$  have order 2 and  $a, b$  have order 3.

Clearly,  $\{e\}$  must be a conjugacy class, since there are no other elements of the same order. We can also see that this must be the case by noting that if  $y$  is in the conjugacy class of  $e$ ,  $y = geg^{-1} = e$ .

Next, let's see if by conjugating  $r$  with other elements in  $S(\Delta)$  we retrieve  $s, t$ . We get that:

$$ara^{-1} = arb = at = s \quad (18.3.12)$$

$$brb^{-1} = bra = bs = t \quad (18.3.13)$$

so we see that  $\{r, s, t\}$  is indeed a conjugacy class. Since there are no other elements of order 2 we know that there are no other elements in this class.

Finally, let's see if  $\{a, b\}$  is a conjugacy class by the same method. We get that:

$$rar^{-1} = rar = rt = b \quad (18.3.14)$$

so we see that  $\{a, b\}$  is indeed another conjugacy class. ◀

For some other groups, such as Abelian groups, there are simpler ways to find the conjugacy partition.

**Example.** Consider the group  $\mathbb{Z}_7^* = \{1, 2, 3, 4, 5, 6\}$ .

Since  $\mathbb{Z}_7^*$  is an abelian group, it follows that if  $x$  is conjugate to  $y$ , then  $y = x$ :

$$g \in G, y = gxg^{-1} = gg^{-1}x = x \quad (18.3.15)$$

Hence, each conjugacy class can only contain one element:

$$\{1\}, \{2\}, \{3\}, \{4\}, \{5\}, \{6\} \quad (18.3.16)$$

This example leads us to believe that for any abelian group, the conjugacy classes must only contain one element.

**Proposition 19.10 (Conjugacy classes of abelian groups)**

The conjugacy classes of an Abelian group contain only one element each.

*Proof.* Suppose  $G$  is an abelian group, and let  $x \in G$ . Then, for  $g \in G$ :

$$gxg^{-1} = gg^{-1}x = ex = x \quad (18.3.17)$$

Therefore,  $x$  is only conjugate to itself. Hence the conjugacy class of  $x$  is  $\{x\}$ , as desired. ■

It is important when talking about conjugacy to express what group the elements are conjugate in i.e. the elements  $x$  and  $y$  are conjugate in the group  $G$ . Indeed, suppose  $H < G$  is a subgroup, then it is not necessarily true that  $x, y$  are conjugate in  $H$ . We know that  $\exists g \in G$  such that  $y = gxg^{-1}$ , but we cannot state that  $g \in H$ . However, the converse is true, that is if  $x, y$  are conjugate in  $H$ , then  $x, y$  must also be conjugate in  $G$ .

**Proposition 19.11 (Conjugacy in subgroups)** Let  $H < G$  be a subgroup of some group  $G$ , and let  $x, y \in H$ . Then:

- (i) if  $x, y$  are conjugate in  $H$  then they are also conjugate in  $G$
- (ii) if  $x, y$  are conjugate in  $G$  then they are not necessarily conjugate in  $H$

**Example.** Consider the subgroup  $H = \{e, (1 2)(3 4), (1 3)(2 4), (1 4)(2 3)\}$  of  $S_4$ .

Since  $H$  has order 4 it must be Abelian, and hence its conjugacy classes can only contain one element. So no two elements of  $H$  can be conjugate to each other in  $H$ . Yet, because all non-identity elements have the same cyclic structure, they are conjugate to each other in  $S_4$ . ■

## 18.4 Normal subgroups and conjugacy

**Theorem 19.12 (Normality criteria)** Let  $H < G$  be a subgroup, then  $H$  is normal in  $G$  iff:

- (a)  $gH = Hg$  for all  $g \in G$
- (b)  $ghg^{-1} \in H$  for each  $h \in H, g \in G$
- (c)  $gHg^{-1} = H$  for each  $g \in G$
- (d)  $H$  is a union of conjugacy classes of  $G$

*Proof.* The goal of this proof will be to show that the following equivalences and implications hold:

$$(a) \iff (c), (b) \implies (c), (c) \implies (d), (d) \implies (b) \quad (18.4.1)$$

(a)  $\implies$  (c) Suppose that  $gH = Hg$ ,  $\forall g \in G$ . Now suppose  $x \in gHg^{-1}$ , then  $\exists h \in H$  such that:

$$x \in gHg^{-1} \quad (18.4.2)$$

$$\iff x = ghg^{-1} \quad (18.4.3)$$

$$\iff x = h'gg^{-1} = h', h' \in H \implies x \in H \quad (18.4.4)$$

since  $gH = Hg$ . Therefore we have proven that  $gHg^{-1} = H$ , as desired.

(c)  $\implies$  (a) Suppose that  $gHg^{-1} = H$ ,  $\forall g \in G$ , and let  $g \in G$ . Then  $\exists h \in H$  such that:

$$x \in gH \quad (18.4.5)$$

$$\iff x = gh \quad (18.4.6)$$

$$\iff x = ghg^{-1}g \quad (18.4.7)$$

$$\iff x = h_1g \implies x \in Hg \quad (18.4.8)$$

Therefore  $gH = Hg$  as desired.

(b)  $\implies$  (c) Suppose that  $ghg^{-1} \in H$ , for each  $h \in H, g \in G$ . Then:

$$h = gg^{-1}hgg^{-1} = g \underbrace{(g^{-1}h(gg^{-1})^{-1})}_{\in H} g^{-1} \quad (18.4.9)$$

Now  $(g^{-1}h(gg^{-1})^{-1}) \in H$  by assumption, so that  $h \in gHg^{-1}$  as desired. Hence  $H \subseteq gHg^{-1}$ .

Moreover, since we assumed that  $gHg^{-1} \subseteq H$  it follows that  $gHg^{-1} = H$  as desired.

(c)  $\implies$  (d) Suppose that  $gHg^{-1} = H$  for all  $g \in G$ , and let  $h \in H$ . Then  $ghg^{-1} \in H$  implying that  $H$  contains all the conjugates in  $G$  of its elements.

(d)  $\implies$  (b) Suppose that  $H$  is a union of conjugacy classes. Suppose that  $h \in H, g \in G$ . Then,  $ghg^{-1}$  is conjugate to  $h$ , and must therefore belong to  $H$ . Hence,  $ghg^{-1} \in H$ . ■

**Example.** Suppose that  $H, K$  are normal subgroups of  $G$ . We have proven in Lagrange's theorem that since  $H, K$  are subgroups of  $G$ ,  $H \cap K < G$ , so let us also prove that  $H \cap K \trianglelefteq G$ . That is, we need to prove that  $ghg^{-1} \in H \cap K$  for all  $h \in H \cap K, g \in G$ .

Since  $H \trianglelefteq G$ , we have that  $gH = Hg$ ,  $\forall g \in G$ , and similarly  $gK = Kg$ ,  $\forall g \in G$ . Therefore, if we let  $x \in H \cap K$  then:

$$gxg^{-1} = ghg^{-1} = h'gg^{-1} = h \quad (18.4.10)$$

for some  $h, h' \in H$ . Similarly:

$$gxg^{-1} = gkg^{-1} = k'gg^{-1} = k' \quad (18.4.11)$$

for some  $k, k' \in K$ . It follows then that  $gxg^{-1} \in H \cap K$ . ◀

**Example.** Consider the group  $X = \{(a, b) \in \mathbb{R}^2 : a \neq 0\}$  equipped with the binary operation:

$$(a, b) * (c, d) = (ac, ad + b) \quad (18.4.12)$$

Consider the subset  $K = \{(1, n) : n \in \mathbb{Z}\}$ . We firstly prove that this is a subgroup of  $X$ .

**Closure:** suppose  $k_1 = (1, n_1), k_2 = (1, n_2) \in K$ , then:

$$k_1 * k_2 = (1, n_1) * (1, n_2) = (1, n_1 + n_2) \in K \quad (18.4.13)$$

due to the closure of  $\mathbb{Z}$ .

**Identity:** the identity of  $X$  was shown to be  $(1, 0)$ . This clearly belongs to  $X$ , since  $0 \in \mathbb{Z}$ .

**Inverses:** the inverse of some element  $k = (1, n) \in X$  is  $(1, -n)$ . This clearly also belongs to  $K$ , since  $-n \in \mathbb{Z}$  provided  $n \in \mathbb{Z}$ .

Since the subgroup axioms are satisfied, we have that  $K < X$ . Now let's see if  $K$  is a normal subgroup of  $X$ , that is  $xkx^{-1} \in K$  for  $x \in X, k \in K$ . Indeed:

$$xkx^{-1} = (a, b) * (1, n) * \left(\frac{1}{a}, -\frac{b}{a}\right) = (a, b) * \left(\frac{1}{a}, n - \frac{b}{a}\right) \quad (18.4.14)$$

$$= (1, an) \quad (18.4.15)$$

This element does not necessarily belong to  $K$ . Indeed, if  $a \in \mathbb{R}$  and  $n \in \mathbb{Z}$  then  $an$  need not to be necessarily an integer. Hence  $K$  is not a normal subgroup of  $X$ .  $\blacktriangleleft$

**Theorem 19.13 (Conjugate subgroup)** Let  $H < G$  and let  $g \in G$ . Then  $gHg^{-1} < G$ .

*Proof.* Let's check the subgroup axioms.

**Closure:** let  $ghg^{-1}, gkg^{-1} \in gHg^{-1}$ , then:

$$(ghg^{-1})(gkg^{-1}) = ghkg^{-1} = gxg^{-1} \quad (18.4.16)$$

where  $x = hk \in H$  due to the closure property of subgroups.

**Identity:** the identity element  $e$  in  $G$  also belongs to  $gHg^{-1}$ , since  $e = geg^{-1}$  and  $e \in H$ .

**Inverses:** let  $ghg^{-1} \in gHg^{-1}$ . The inverse of this element in  $G$  is  $gh^{-1}g^{-1}$ , which must also belong to  $gHg^{-1}$  since  $h^{-1} \in H$ .  $\blacksquare$

**Example.** Consider the subgroup  $H = \langle s \rangle$  of  $S(\square)$ . Then, the conjugate subgroup in  $a$  is:

$$aHa^{-1} = aHc = a\{e, s\}c = a\{c, t\} = \{e, u\} \quad (18.4.17)$$

We can use the definition of conjugate subgroups for two subgroups. Indeed, if some element  $g$  conjugates  $H$  to  $K$ , then we say that  $H, K$  are conjugate subgroups in  $G$ .

**Proposition 19.14 (Isomorphism of conjugate subgroups)**

If  $H, K$  are conjugate subgroups in  $G$ , then  $H, K$  are also isomorphic.

*Proof.* Suppose  $H, K$  are conjugate subgroups in  $G$ . Then,  $\exists g \in G$  such that  $K = gHg^{-1}$ . We prove that the following is an isomorphism:

$$\phi : H \rightarrow K \quad (18.4.18)$$

$$h \mapsto ghg^{-1} \quad (18.4.19)$$

This mapping is injective since  $\phi(x) = \phi(y)$  implies that  $gh_1g^{-1} = gh_2g^{-1} \implies h_1 = h_2$ . Moreover, it is surjective since every element of  $K$  must be of the form  $ghg^{-1}$  where  $h \in H$ .

Finally:

$$\phi(xy) = gxyg^{-1} = (gxg^{-1})(gyg^{-1}) = \phi(x)\phi(y) \quad (18.4.20)$$

as desired.  $\blacksquare$

**Example.** Consider the following subgroup of  $A_4$ :

$$K = \{e, (1\ 2)(3\ 4), (1\ 3)(2\ 4), (1\ 4)(2\ 3)\} \quad (18.4.21)$$

Let's find the following conjugate subgroups:

$$(1\ 2\ 4)K(1\ 2\ 4)^{-1} = (1\ 2\ 4)K(1\ 4\ 2), \text{ and } (2\ 4\ 3)K(2\ 4\ 3)^{-1} = (2\ 4\ 3)K(2\ 3\ 4) \quad (18.4.22)$$

Firstly, using the fact that conjugacy does not affect the cycle structure:

$$(1\ 2\ 4)e(1\ 4\ 2) = (1\ 2\ 4)(1\ 4\ 2) = e \quad (18.4.23)$$

$$(1\ 2\ 4)(1\ 2)(3\ 4)(1\ 4\ 2) = (1\ 3)(2\ 4) \quad (18.4.24)$$

$$(1\ 2\ 4)(1\ 3)(2\ 4)(1\ 4\ 2) = (1\ 4)(2\ 3) \quad (18.4.25)$$

$$(1\ 2\ 4)(1\ 4)(2\ 3)(1\ 4\ 2) = (1\ 2)(3\ 4) \quad (18.4.26)$$

Similarly:

$$(2\ 4\ 3)e(2\ 3\ 4) = (2\ 4\ 3)(2\ 3\ 4) = e \quad (18.4.27)$$

$$(2\ 4\ 3)(1\ 2)(3\ 4)(2\ 3\ 4) = (1\ 4)(2\ 3) \quad (18.4.28)$$

$$(2\ 4\ 3)(1\ 3)(2\ 4)(2\ 3\ 4) = (1\ 2)(3\ 4) \quad (18.4.29)$$

$$(2\ 4\ 3)(1\ 4)(2\ 3)(2\ 3\ 4) = (1\ 3)(2\ 4) \quad (18.4.30)$$

Therefore,  $(1\ 2\ 4)K(1\ 2\ 4)^{-1} = (2\ 4\ 3)K(2\ 4\ 3)^{-1} = K$ .

Since conjugating subgroups must leave the cycle structure invariant, and since there are only three permutations of structure  $(--)(- -)$  in  $A_4$ , it follows that the only conjugating subgroup of  $K$  is itself.  $\blacktriangleleft$

**Example.** Consider the subgroups of  $S(\Delta)$ :

Order	Subgroups
1	$\{e\}$
2	$\{e, r\}, \{e, s\}, \{e, t\}$
3	$\{e, a, b\}$
6	$S(\Delta)$

and its conjugacy classes which we found earlier:

$$\{e\}, \{a, b\}, \{r, s, t\} \quad (18.4.31)$$

Let's try to find all the normal subgroups of  $S(\Delta)$ . Recall that  $H \trianglelefteq S(\Delta)$  iff  $H$  is a union of conjugacy classes of  $G$ . We then see that the only subgroups which can be expressed as such unions are:

$$\{e\} = \{e\} \quad (18.4.32)$$

$$\{e, a, b\} = \{e\} \cup \{a, b\} \quad (18.4.33)$$

$$S(\Delta) = \{e\} \cup \{a, b\} \cup \{r, s, t\} \quad (18.4.34)$$

◀

Unfortunately, often times we do not have a list of all the subgroups of a group. In such cases, it is easier to find which unions of conjugacy classes are normal subgroups. These must contain the conjugacy class  $\{e\}$  and they must have order which divides the group's order, by Lagrange's theorem.

**Strategy** (*Determining normal subgroups from conjugacy classes*)

- (i) partition  $G$  into conjugacy classes
- (ii) Find all the unions of conjugacy classes which include  $\{e\}$  and whose order divides  $|G|$  as required by Lagrange's theorem.
- (iii) Determine which of these unions are subgroups, and hence normal subgroups.

We illustrate this method below:

**Example.** These are the conjugacy classes of  $A_5$ :

Conjugacy class	Description	Order
A	$\{e\}$	1
B	(- - -)	20
C	(- -)(- -)	15
D	conjugate to (1 2 3 4 5)	12
E	conjugate to (1 2 3 5 4)	12

We need to find all possible unions which contain  $A$ , whose order divides  $|A_5| = 60$ , so 1,2,3,4,5,6,10,12,15,20,30,60.

Firstly, the only union which has only one element is  $A = \{e\}$ .

Secondly, the unions which have 2,3,4,5,6,10 elements do not exist.

Thirdly, the unions which have 12 elements are two,  $D$  and  $E$ . However this does not contain  $e$ , so we scrap it.

Similarly, the only union which has 15 elements is  $C$ . However this does not contain  $e$ , so we scrap it.

Also, the only union which has 20 elements is  $B$ . However this does not contain  $e$ , so we scrap it.

There are no unions which contain 30 elements.

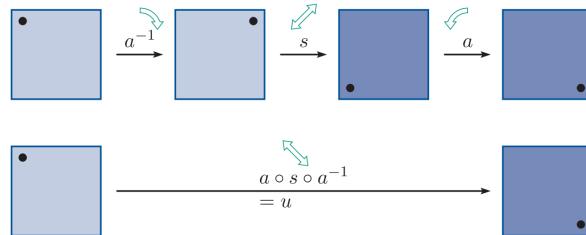
Finally, the only union which contains 60 elements is  $A \cup B \cup C \cup D \cup E = A_5$ .

So we see that the only candidates for normal subgroups are  $A$  and  $A_4$ . These are clearly

subgroups, and thus also normal subgroups. ◀

## 18.5 Conjugacy in $S(\mathcal{F})$

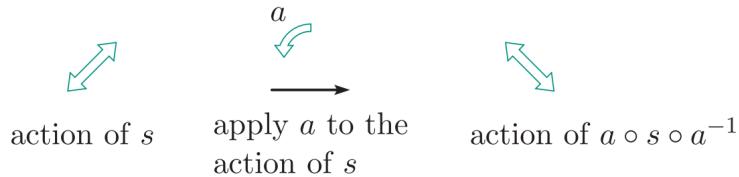
In the case of the symmetry groups, we can define conjugacy more concretely by looking at the action of conjugating some symmetry by another symmetry.



**Figure 18.2.** Application of  $rar^{-1}$

Consider for example the effect of  $rar^{-1} = c$ . We can view this symmetry as applying  $a$  on the square that has been reflected about the vertical line of symmetry.

In other words, since we are rotating through  $\frac{\pi}{2}$  anti-clockwise on the reflected square, when we reflect back to the original square we see that the overall action of  $rar^{-1}$  was to rotate clockwise. Thus, the conjugate symmetry is equivalent to the symmetry we obtain by applying  $r$  to the action of  $a$ , that is, reflect the action of  $a$  so that it goes from anti-clockwise to clockwise.



**Figure 18.3.** Effect of conjugacy on symmetries

More generally, given two symmetries  $x, g \in S(\mathcal{F})$ , then  $gxg^{-1}$  is the symmetry obtained by acting  $g$  on the visual effect of  $x$ .

It follows that two symmetries are conjugate *iff* there is a symmetry whose action on the visual effect of one of the symmetries is to give the visual effect of the other symmetry.

In other words, if there is a way to relabel the figure in such a way for the effect of the two symmetries to be identical, then they are conjugate.

For example, in the case of  $r$  and  $s$ , these two symmetries are not conjugate, since there is no possible symmetry of  $S(\square)$  that can map the effect of  $r$  to the effect of  $s$ . We cannot rename the vertices of the square in any possible way for the effect of  $r$  and  $s$  to coincide.

If instead we considered  $S(\text{octagon})$  then we would indeed have that  $r$  and  $s$  are conjugate symmetries, for example through the rotation by  $\frac{\pi}{4}$  anti-clockwise.

Similarly, it is easy to see that  $r$  and  $t$  are conjugate symmetries. For example,  $t = ara^{-1}$ , since applying  $r$ , a vertical reflection, on a square which has been rotated by  $\frac{\pi}{2}$  anti-clockwise, is equivalent to applying  $t$ , a horizontal reflection.

**Definition 19.15 (Fixed point set)**

Let  $f \in S(\mathcal{F})$ , then the **fixed point set** of  $f$  is defined as:

$$\text{Fix } f = \{P \in \mathcal{F} : f(P) = P\} \quad (18.5.1)$$

that is, the subset of  $\mathcal{F}$  that is invariant under  $f$ .

In two dimensions, for a rotation, the fixed point set consists of the center of rotation. Similarly, for a reflection, the fixed point set consists of the line of reflection.

Since applying  $g$  to the visual effect of  $x$  gives the action of  $gxg^{-1}$ , we expect that if  $\text{Fix } f$  are invariant under  $x$ , then  $g(\text{Fix } f)$  must be the fixed point set of  $gxg^{-1}$ .

**Theorem 19.16 (Fixed point set of conjugate symmetries)**

Let  $f, g \in S(\mathcal{F})$ , then  $\text{Fix } gfg^{-1} = g(\text{Fix } f)$ .

*Proof.* Firstly we prove that  $g(\text{Fix } f) \subseteq \text{Fix } gfg^{-1}$ . Indeed, suppose that  $P \in g(\text{Fix } f)$ . Then:

$$g(-1)(P) \in \text{Fix } f \implies (fg^{-1})(P) = g^{-1}(P) \implies (gfg^{-1})(P) = (P) \quad (18.5.2)$$

In other words,  $P \in \text{Fix } gfg^{-1}$ .

Now, suppose that  $P \in \text{Fix } gfg^{-1}$ . Then:

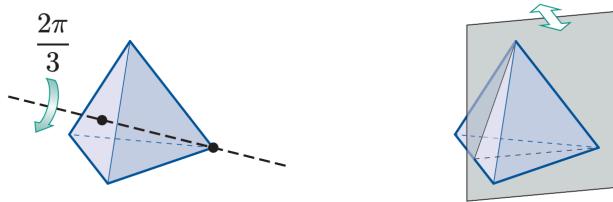
$$(gfg^{-1})(P) = P \implies (fg^{-1})(P) = g^{-1}(P) \implies f(g^{-1})(P) = (g^{-1}(P)) \quad (18.5.3)$$

proving that  $g^{-1}(P) \in \text{Fix } f$ , and thus that  $P \in g(\text{Fix } f)$ , as desired.

Thus, since  $g(\text{Fix } f) \subseteq \text{Fix } gfg^{-1}$  and  $\text{Fix } gfg^{-1} \subseteq g(\text{Fix } f)$ , we find that  $g(\text{Fix } f) = \text{Fix } gfg^{-1}$ . ■

This theorem is extremely useful when trying to tell whether or not two symmetries are conjugate. Indeed, the only candidate conjugating symmetries are those that map the fixed point set of one symmetry to the other.

For example, consider the following two symmetries of a tetrahedron:



**Figure 18.4.** Two non-conjugate symmetries of a tetrahedron

The fixed point set of the rotation is simply the axis of rotation, so it is 1-dimensional. The fixed point set of the reflection instead is the plane of reflection, which is 2-dimensional. It follows immediately that no symmetry can map these two fixed point sets to one another.

More generally, direct symmetries cannot be conjugate to indirect symmetries.

### 19.17 (Conjugacy direct and indirect symmetries)

A direct symmetry cannot be conjugate to an indirect symmetry.

*Proof.* Let  $x$  be a direct symmetry and  $y$  be any symmetry. If  $g$  is direct, then  $gxg^{-1}$  is also direct. If  $g$  is indirect, then  $gxg^{-1}$  is direct. Therefore  $x$  can only be conjugate to direct symmetries. ■

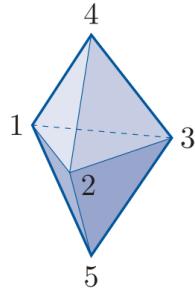
#### Strategy (Finding conjugacy classes of finite symmetry groups)

- (i) Represent  $S(\mathcal{F})$  as a subgroup of a symmetric group.
- (ii) Partition  $S(\mathcal{F})$  by cycle structure.
- (iii) For each cycle structure class, determine which of the symmetries are conjugate to each other. R

Recall that the number of elements in each conjugacy class divides  $|S(\mathcal{F})|$ , and cannot contain both a direct and indirect symmetry.

Also, remember that two symmetries whose fixed point sets cannot be mapped to each other are not conjugate.

**Example.** Consider the double tetrahedron below:



Firstly, since  $|S(\mathcal{F})| = 12$ , we see that the conjugacy classes can only contain 1, 2, 3, 4, 6, or 12 elements.

We can see that its symmetries may be categorized in terms of their cycle structure as follows:

$$\{e\} \tag{18.5.4}$$

$$\{(1 2), (1 3), (2 3), (4 5)\} \tag{18.5.5}$$

$$\{(1 2 3), (1 3 2)\} \tag{18.5.6}$$

$$\{(1 2)(4 5), (1 3)(4 5), (2 3)(4 5)\} \tag{18.5.7}$$

$$\{(1 2 3)(4 5), (1 3 2)(4 5)\} \tag{18.5.8}$$

Clearly,  $\{e\}$  is a conjugacy class.

Next, we see that  $(1 2), (1 3), (2 3)$  are all conjugate through rotations by  $\frac{\pi}{3}$  clockwise or anticlockwise, but are not conjugate to  $(4 5)$ .

To check this note that using the renaming method:

$$(2\ 3)(1\ 2)(2\ 3)^{-1} = (1\ 3) \quad (18.5.9)$$

$$(1\ 2)(1\ 3)(1\ 2)^{-1} = (2\ 3) \quad (18.5.10)$$

and since conjugacy is an equivalence relation, transitivity implies that if  $(1\ 2)$  is conjugate to  $(1\ 3)$ , and  $(1\ 3)$  is conjugate to  $(2\ 3)$ , then these must all be conjugate to each other.

To prove that  $(4\ 5)$  is not conjugate to the other three, we note that conjugacy does not affect cycle structure. Therefore, if there exists a conjugating symmetry  $g \in S(\mathcal{F})$  between  $(4\ 5)$  and, say,  $(1\ 2)$  then we would need  $g = (--)$  such that:

$$(-)(4\ 5)(--) = (1\ 2) \quad (18.5.11)$$

Such a symmetry  $g$  does not belong to  $S(\mathcal{F})$  (we can check individually all four permutations of structure  $(--)$ ).

So, we have found two other conjugacy classes,  $\{(1\ 2), (1\ 3), (2\ 3)\}$  and  $\{(4\ 5)\}$ .

Next we see immediately that  $(1\ 2\ 3)(4\ 5)$  and  $(1\ 3\ 2)(4\ 5)$  are conjugate through  $(2\ 3) \in S(\mathcal{F})$ :

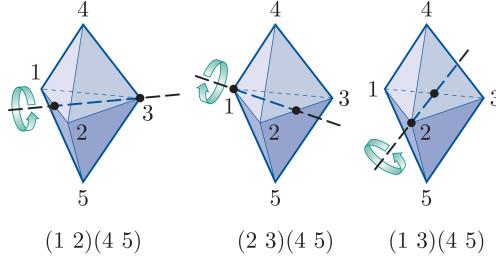
$$(2\ 3)(1\ 2\ 3)(2\ 3)^{-1} = (1\ 3\ 2) \quad (18.5.12)$$

Similarly, we also find that  $(1\ 2\ 3)(4\ 5)$  and  $(1\ 3\ 2)(4\ 5)$  are conjugate through  $(2\ 3) \in S(\mathcal{F})$ :

$$(2\ 3)(1\ 2\ 3)(4\ 5)(2\ 3)^{-1} = (1\ 3\ 2)(4\ 5) \quad (18.5.13)$$

Hence we find the conjugacy classes  $\{(1\ 2\ 3), (1\ 3\ 2)\}$  and  $\{(1\ 2\ 3)(4\ 5), (1\ 3\ 2)(4\ 5)\}$

Finally,  $\{(1\ 2)(4\ 5), (1\ 3)(4\ 5), (2\ 3)(4\ 5)\}$  is another conjugacy class. We see this intuitively since these are all rotations about axes in the plane of the triangular base, passing through a vertex and the midpoint of the opposite edge. Therefore, we can conjugate them through rotations by  $\frac{\pi}{3}$  clockwise and anti-clockwise.



More rigorously:

$$(2\ 3)(1\ 2)(4\ 5)(2\ 3)^{-1} = (1\ 3)(4\ 5) \quad (18.5.14)$$

$$(1\ 2)(1\ 3)(4\ 5)(1\ 2)^{-1} = (2\ 3)(4\ 5) \quad (18.5.15)$$

as desired. Hence the conjugacy classes of  $S(\mathcal{F})$  are:

$$\{e\} \quad (18.5.16)$$

$$\{(1 2), (1 3), (2 3)\} \quad (18.5.17)$$

$$\{(4 5)\} \quad (18.5.18)$$

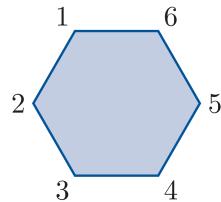
$$\{(1 2 3), (1 3 2)\} \quad (18.5.19)$$

$$\{(1 2)(4 5), (1 3)(4 5), (2 3)(4 5)\} \quad (18.5.20)$$

$$\{(1 2 3)(4 5), (1 3 2)(4 5)\} \quad (18.5.21)$$

◀

**Example.** We label the hexagon as shown:



Consider the symmetries of a hexagon:

$$\{e\} \quad (18.5.22)$$

$$\{(1 3)(4 6), (2 6)(3 5), (1 5)(2 4)\} \quad (18.5.23)$$

$$\{(1 4)(2 5)(3 6), (1 6)(2 5)(3 4), (1 2)(3 6)(4 5), (1 4)(2 3)(5 6)\} \quad (18.5.24)$$

$$\{(1 3 5)(2 4 6), (1 5 3)(2 6 4)\} \quad (18.5.25)$$

$$\{(1 2 3 4 5 6), (1 6 5 4 3 2)\} \quad (18.5.26)$$

Clearly,  $\{e\}$  is a conjugacy class.

Instead,  $\{(1 4)(2 5)(3 6), (1 6)(2 5)(3 4), (1 2)(3 6)(4 5), (1 4)(2 3)(5 6)\}$  contains one direct symmetry (the first is a rotation by  $\pi$ ) and three indirect symmetries which are reflections in axes passing through the midpoints of opposite edges. Hence, we must have the conjugacy class  $\{(1 4)(2 5)(3 6)\}$ .

Also,  $\{(1 6)(2 5)(3 4), (1 2)(3 6)(4 5), (1 4)(2 3)(5 6)\}$  are conjugate through rotations by  $\frac{\pi}{3}$  clockwise and anti-clockwise, and the proof is quite similar to the previous example.

Similarly,  $\{(1 2 3 4 5 6), (1 6 5 4 3 2)\}$  are conjugate through a vertical reflection  $(2 6)(3 5)$ .

Repeating this logical process we quickly see that the conjugacy classes are:

$$\{e\} \quad (18.5.27)$$

$$\{(1 4)(2 5)(3 6)\} \quad (18.5.28)$$

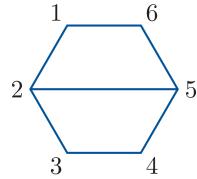
$$\{(1 3)(4 6), (2 6)(3 5), (1 5)(2 4)\} \quad (18.5.29)$$

$$\{(1 6)(2 5)(3 4), (1 2)(3 6)(4 5), (1 4)(2 3)(5 6)\} \quad (18.5.30)$$

$$\{(1 3 5)(2 4 6), (1 5 3)(2 6 4)\} \quad (18.5.31)$$

$$\{(1 2 3 4 5 6), (1 6 5 4 3 2)\} \quad (18.5.32)$$

If we now modify the hexagon as shown below:



The symmetries of the new modified picture form a subgroup of the symmetries of a hexagon. This subgroup can be partitioned into permutations of the same cyclic structure as:

$$\{e\} \tag{18.5.33}$$

$$\{(1 3)(4 6)\} \tag{18.5.34}$$

$$\{(1 6)(2 5)(3 4)\} \tag{18.5.35}$$

◀

Recall that if  $G$  is a group with subgroup  $H$ , then  $H$  is a normal subgroup of  $G$  iff  $H$  is a union of conjugacy classes of  $G$ . In our case, the symmetry of the modified hexagon definitely does not form a normal subgroup, since it is not the union of conjugacy classes of the normal hexagon's symmetry group.

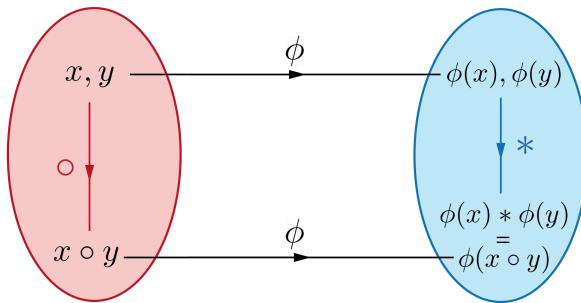
# Unit E3: Homomorphisms

## Definition 20.1 (Homomorphism)

Let  $(G, \circ)$  and  $(H, *)$  be groups. A mapping  $\phi : (G, \circ) \rightarrow (H, *)$  is a **homomorphism** if it satisfied the property:

$$\phi(x \circ y) = \phi(x) * \phi(y), \forall x, y \in G \quad (19.0.1)$$

It follows that isomorphisms are homomorphisms that are also bijective.



**Figure 19.1.** Why the homomorphism property must be satisfied for "sensible" maps

**Example.** Consider the following mapping:

$$\phi : (S_n, \circ) \longrightarrow (\mathbb{Z}_2, +_2) \quad (19.0.2)$$

$$\sigma \mapsto \begin{cases} 0 & \text{if } \sigma \text{ is an even permutation} \\ 1 & \text{if } \sigma \text{ is an odd permutation} \end{cases} \quad (19.0.3)$$

Lets show that it is a homomorphism, thus satisfying  $\phi(x \circ y) = \phi(x) +_2 \phi(y), \forall x, y \in S_n$ . Suppose  $x, y$  are both of even permutations, then  $x \circ y$  is also an even permutation, so that:

$$\phi(x \circ y) = 0 = 0 +_2 0 = \phi(x) +_2 \phi(y) \quad (19.0.4)$$

If instead  $x, y$  are both odd permutations, then  $x \circ y$  must instead be an even permutation, so that:

$$\phi(x \circ y) = 0 = 1 +_2 1 = \phi(x) +_2 \phi(y) \quad (19.0.5)$$

Finally, if  $x, y$  have opposite parity, then  $x \circ y$  is an odd permutation, so that:

$$\phi(x \circ y) = 1 = 0 +_2 1 = \phi(x) +_2 \phi(y) \quad (19.0.6)$$

Therefore,  $\phi$  is a homomorphism.  $\blacktriangleleft$

**Example.** Let us show that the following mapping is not a homomorphism:

$$\phi : (\mathrm{GL}(n, \mathbb{R}), \times) \longrightarrow (\mathrm{GL}(n, \mathbb{R}), \times) \quad (19.0.7)$$

$$A \longmapsto A^{-1} \quad (19.0.8)$$

Consider  $A, B \in \mathrm{GL}(n, \mathbb{R})$  then:

$$\phi(A) \times \phi(B) = A^{-1} \times B^{-1} \quad (19.0.9)$$

whereas

$$\phi(A \times B) = (A \times B)^{-1} = B^{-1} \times A^{-1} \quad (19.0.10)$$

However note that  $B^{-1} \times A^{-1} \neq A^{-1} \times B^{-1}$  generally.  $\blacktriangleleft$

### Proposition 20.2 ( $\mathbb{Z}$ homomorphic to $\mathbb{Z}_n$ )

For  $n \geq 2$ , the following is a homomorphism to:

$$\phi : (\mathbb{Z}, +) \longrightarrow (\mathbb{Z}_n, +_n) \quad (19.0.11)$$

$$k \longmapsto k_{(\text{mod } n)} \quad (19.0.12)$$

where  $k_{(\text{mod } n)}$  is the least residue of  $k$  modulo  $n$  (remained of  $k$  when divided by  $n$ ).

*Proof.* Suppose  $n \geq 2$  and let  $r, s \in \mathbb{Z}$ . Then:

$$\phi(r + s) = (r + s)_{(\text{mod } n)} \quad (19.0.13)$$

$$\equiv r + s \pmod{n} \quad (19.0.14)$$

$$\equiv r_{(\text{mod } n)} + s_{(\text{mod } n)} \pmod{n} \quad (19.0.15)$$

$$\equiv r_{(\text{mod } n)} +_n s_{(\text{mod } n)} \pmod{n} \quad (19.0.16)$$

$$= \phi(r) +_n \phi(s) \quad (19.0.17)$$

and since  $\phi(r + s)$  and  $\phi(r) +_n \phi(s)$  both belong to  $\mathbb{Z}_n$ , we have that the two must be equal. Thus  $\phi$  is a homomorphism.  $\blacksquare$

**Example.** Let us prove that:

$$\phi : (G, \circ) \longrightarrow (G, \circ) \quad (19.0.18)$$

$$x \longmapsto x \circ x \quad (19.0.19)$$

is a homomorphism iff  $(G, \circ)$  is abelian.

Suppose  $\phi$  is indeed a homomorphism, so that:

$$\phi(x \circ y) = (x \circ y) \circ (x \circ y) = (x \circ x) \circ (y \circ y) = \phi(x) \circ \phi(y) \quad (19.0.20)$$

$$\iff x \circ y \circ x \circ y = x \circ x \circ y \circ y \quad (19.0.21)$$

$$\iff y \circ x = x \circ y \quad (19.0.22)$$

This implies that  $(G, \circ)$  is indeed an abelian group, as desired.  $\blacktriangleleft$

### Proposition 20.3 (Trivial homomorphism)

Let  $(G, \circ)$  and  $(H, *)$  be groups, and let  $e_G, e_H$  be their respective identity elements. Then the following is a homomorphism:

$$\phi : (G, \circ) \longrightarrow (H, *) \quad (19.0.23)$$

$$x \longmapsto e_H \quad (19.0.24)$$

*Proof.* Let  $x, y \in G$ , then:

$$\phi(x \circ y) = e_H = e_H * e_H = \phi(x) * \phi(y) \quad (19.0.25)$$

as desired.  $\blacksquare$

### Proposition 20.4 (Properties of homomorphisms)

Let  $\phi : (G, \circ) \longrightarrow (H, \odot)$  be a homomorphism, then

$$(i) \text{ let } x_1, x_2, \dots, x_n \in G, \text{ then } \phi\left(\bigcirc_{k=1}^n x_k\right) = \bigodot_{k=1}^n \phi(x_k)$$

(ii)  $\phi(e_G) = e_H$  where  $e_G, e_H$  are the identity elements of  $G, H$  respectively.

(iii) for  $x \in G$ ,  $\phi(x^{-1}) = (\phi(x))^{-1}$

(iv) for  $x \in G$ ,  $\phi(x^n) = (\phi(x))^n$

*Proof.*

(i) For the case  $n = 1$ , it is clear that  $\phi(x_1) = \phi(x_1)$ . Suppose that for  $n \in \mathbb{N}$ , we have that:

$$\phi\left(\bigcirc_{k=1}^n x_k\right) = \bigodot_{k=1}^n \phi(x_k) \quad (19.0.26)$$

Then:

$$\phi\left(\bigcirc_{k=1}^{n+1} x_k\right) = \phi\left(\left(\bigcirc_{k=1}^n x_k\right) \circ x_{n+1}\right) \quad (19.0.27)$$

$$= \left(\phi\left(\bigcirc_{k=1}^n x_k\right)\right) \odot \phi(x_{n+1}) \quad (19.0.28)$$

$$= \bigodot_{k=1}^n \phi(x_k) \odot \phi(x_{n+1}) \quad (19.0.29)$$

$$= \bigodot_{k=1}^{n+1} \phi(x_k) \quad (19.0.30)$$

---

as desired. Hence by the principle of mathematical induction, we have that  $\phi\left(\bigodot_{k=1}^n x_k\right) = \bigodot_{k=1}^n \phi(x_k)$ .

(ii) We have that  $e_G \circ e_G = e_G$ , then:

$$\phi(e_G \circ e_G) = \phi(e_G) \quad (19.0.31)$$

giving:

$$\phi(e_G) \odot \phi(e_G) = \phi(e_G) = \phi(e_G) \odot e_H \quad (19.0.32)$$

from which we find that  $\phi(e_G) = e_H$ .

(ii) Let  $x \in G$ . Then:

$$x \circ x^{-1} = x^{-1} \circ x = e_G \quad (19.0.33)$$

Then, applying  $\phi$ :

$$\phi(x \circ x^{-1}) = \phi(x^{-1} \circ x) = \phi(e_G) = e_H \quad (19.0.34)$$

Therefore:

$$\phi(x) \odot \phi(x^{-1}) = \phi(x) \odot \phi(x^{-1}) = e_G \quad (19.0.35)$$

implying that  $\phi(x^{-1}) = (\phi(x))^{-1}$ , as desired.

(iv) We consider two different cases,  $n \geq 0$  and  $n < 0$ .

If  $n \geq 0$ , then the case  $n = 0$  is trivial since:

$$\phi(x^0) = \phi(e_G) = (\phi(x))^0 = e_H \quad (19.0.36)$$

which was proven previously. Now suppose that for some integer  $k \geq 0$

$$\phi(x^k) = (\phi(x))^k \quad (19.0.37)$$

Then:

$$\phi(x^{k+1}) = \phi(x^k \circ x) = \phi(x^k) \odot \phi(x) \quad (19.0.38)$$

$$= (\phi(x))^k \odot \phi(x) \quad (19.0.39)$$

$$= (\phi(x))^{k+1} \quad (19.0.40)$$

as desired.

Now suppose  $n < 0$ , and let  $x \in G$ . We can write  $n = -m$  where  $m > 0$ . Then:

$$\phi(x^n) = \phi(x^{-m}) \quad (19.0.41)$$

$$= \phi((x^{-1})^m) \quad (19.0.42)$$

$$= (\phi(x^{-1}))^m \quad (19.0.43)$$

$$= (\phi(x))^{-m} \quad (19.0.44)$$

$$= (\phi(x))^n \quad (19.0.45)$$

Therefore, both from cases 1 and 2, we find that:

$$\phi(x^n) = (\phi(x))^n \quad (19.0.46)$$

as desired. ■

### Theorem 20.5 (Order of element and homomorphism image)

Let  $\phi : (G, \circ) \rightarrow (H, *)$  be a homomorphism and let  $x \in G$  be an element of finite order. Then the order of  $\phi(x)$  is also finite and divides the order of  $x$ .

*Proof.* We begin by proving the following lemma.

**Lemma.** Let  $x \in G$ . If  $r > 0$  is a positive integer such that  $x^r = e$ , then the order of  $x$  divides  $r$ .

Suppose that  $x^r = e$ , and suppose that  $\text{ord}(x) = s$ . Then, we must have that  $r = as + b$  for some integers  $a, b$  with  $0 \leq b < s$ .

Hence:

$$e = x^r \quad (19.0.47)$$

$$= x^{as+b} \quad (19.0.48)$$

$$= (x^s)^a \circ x^b \quad (19.0.49)$$

$$= e^a \circ x^b \quad (19.0.50)$$

$$= x^b \quad (19.0.51)$$

Since  $b < s$ , we find that  $b = 0$ , or else we would have a contradiction. Therefore,  $r = as$ , in other words  $s$  divides  $r$ .

Now since the order of  $x$  is  $s$ :

$$(\phi(x))^s = \phi(x^s) = \phi(e_G) = e_H \quad (19.0.52)$$

hence the order of  $\phi$  must, by the above lemma, have order that divides  $s$ . In other words, the order of  $\phi(x)$  must be finite and divide the order of  $x$ . ■

### Theorem 20.6 (Conjugacy of element and homomorphism image)

Let  $\phi : (G, \circ) \rightarrow (H, *)$  be a homomorphism, and let  $x, y \in G$ . If  $x, y$  are conjugate then  $\phi(x), \phi(y)$  are conjugate too.

*Proof.* Suppose  $x, y$  are conjugate, so that  $\exists g \in G$  such that:

$$y = g \circ x \circ g^{-1} \quad (19.0.53)$$

Therefore:

$$\phi(y) = \phi(g \circ x \circ g^{-1}) = \phi(g) * \phi(x) * \phi(g^{-1}) = \phi(g) * \phi(x) * (\phi(g))^{-1} \quad (19.0.54)$$

so that  $\phi(x), \phi(y)$  are also conjugate. ■

## 19.1 Image and kernels

### Definition 20.7 (Image and kernel of homomorphism)

Let  $\phi : (G, \circ) \rightarrow (H, *)$  be a homomorphism. Then the **image** of  $\phi$  is:

$$\text{Im } \phi = \{\phi(g) : g \in G\} \quad (19.1.1)$$

and the **kernel** of  $\phi$  is:

$$\ker \phi = \{g \in G : \phi(g) = e_H\} \quad (19.1.2)$$

### Theorem 20.8 (Image subgroup of homomorphism)

Let  $\phi : (G, \circ) \rightarrow (H, *)$  be a homomorphism. Then  $\text{Im } \phi \leq (H, *)$ .

*Proof.* We check the three subgroup axioms:

**Closure:** let  $h_1, h_2 \in \text{Im } \phi$ , so that  $\exists g_1, g_2 \in G$  satisfying  $\phi(g_1) = h_1$  and  $\phi(g_2) = h_2$ . We find that:

$$h_1 * h_2 = \phi(g_1) * \phi(g_2) = \phi(g_1 \circ g_2) \quad (19.1.3)$$

Therefore,  $h_1 * h_2$  is the image of  $g_1 \circ g_2$ , and thus  $h_1 * h_2 \in \text{Im } \phi$ .

**Identity:** the identity of  $H$  is  $e_H$ , and also belongs to  $\text{Im } \phi$  since  $\phi(e_G) = e_H$ , in other words it is the image of the identity of  $G$ .

**Inverses:** suppose  $h \in \text{Im } \phi$ . Then, there exists  $g \in G$  such that:

$$h = \phi(g) \implies h^{-1} = (\phi(g))^{-1} = \phi(g^{-1}) \quad (19.1.4)$$

so that  $h^{-1} \in \text{Im } \phi$ .

Hence all three subgroup axioms are satisfied, and thus  $\text{Im } \phi \leq (H, *)$ . ■

### Proposition 20.9 $((G, \circ) \cong \text{Im } \phi)$

Let  $\phi : (G, \circ) \rightarrow (H, *)$  be an injective homomorphism and let  $\varphi$  be the map obtained by shrinking the codomain of  $\phi$  to  $\text{Im } \phi$ . Then  $\varphi$  is an isomorphism, so that  $(G, \circ) \cong \text{Im } \phi$ .

*Proof.*  $\varphi$  is still an injective homomorphism, since shrinking the codomain to  $\text{Im } \phi$  does not affect the homomorphism property.

However,  $\varphi$  is also surjective, since  $\text{Im } \phi = \text{Im } \varphi$  is the domain of this mapping. Consequently,  $\varphi$  is an isomorphism  $(G, \circ) \rightarrow \text{Im } \phi$ , as desired. ■

### Theorem 20.10 (Preserved structures under homomorphism)

Let  $\phi : (G, \circ) \rightarrow (H, *)$  be a homomorphism:

- (i) if  $G$  is abelian then  $(\text{Im } \phi, *)$  is abelian.
- (ii) If  $G$  is cyclic then  $(\text{Im } \phi, *)$  is cyclic.

Moreover, if  $G$  is generated by  $a$ , then  $(\text{Im } \phi, *)$  is generated by  $\phi(a)$ .

*Proof.* (i) Suppose that  $G$  is abelian, and let  $h_1, h_2 \in \text{Im } \phi$ , so that  $\exists g_1, g_2$  such that  $h_1 = \phi(g_1)$  and  $h_2 = \phi(g_2)$ . Therefore:

$$\phi(g_1 \circ g_2) = h_1 * h_2 = \phi(g_2 \circ g_1) = h_2 * h_1 \quad (19.1.5)$$

proving that  $H$  is also abelian.

(ii) Suppose that  $\langle a \rangle = G$ , and let  $h \in \text{Im } \phi$ , so that  $h = \phi(g)$  for some  $g \in G$ . Now since  $G$  is generated by  $a$  we find that:

$$g = a^k \implies \phi(g) = \phi(a^k) = (\phi(a))^k \quad (19.1.6)$$

Consequently:

$$h = (\phi(a))^k \quad (19.1.7)$$

proving that  $\text{Im } \phi$  is generated by  $\phi(a)$ , and thus also cyclic. ■

### Theorem 20.11 (Kernel normal subgroup)

Let  $\phi : (G, \circ) \rightarrow (H, *)$  be a homomorphism. Then  $\text{Ker}(\phi) \trianglelefteq (G, \circ)$ .

*Proof.* We firstly need to prove that  $\text{Ker}(\phi)$  is a subgroup of  $(G, \circ)$  by checking the subgroup properties.

**Closure:** let  $k_1, k_2 \in \text{Ker}(\phi)$ . Then  $\phi(k_1) = e_H$  and  $\phi(k_2) = e_H$ , so that:

$$\phi(k_1 \circ k_2) = \phi(k_1) * \phi(k_2) = e_H * e_H = e_H \quad (19.1.8)$$

Hence  $k_1 \circ k_2 \in \text{Ker}(\phi)$ .

**Identity:** we have  $\phi(e_G) = e_H$ , so that  $e_G \in \text{Ker}(\phi)$ .

**Inverses:** let  $k \in \text{Ker}(\phi)$ , then  $\phi(k) = e_H$  so that:

$$\phi(k^{-1}) = (\phi(k))^{-1} = e_H^{-1} = e_H \quad (19.1.9)$$

so that  $k^{-1} \in \text{Ker}(\phi)$ .

So, we have that  $(\text{Ker}(\phi), \circ) \leq (G, \circ)$ .

To prove normality, we must prove that  $g \circ k \circ g^{-1} \in \text{Ker}(\phi)$  for  $k \in \text{Ker}(\phi)$  and  $g \in G$ . Indeed:

$$\phi(g \circ k \circ g^{-1}) = \phi(g) * \phi(k) * \phi(g^{-1}) \quad (19.1.10)$$

$$= \phi(g) * e_H * (\phi(g))^{-1} \quad (19.1.11)$$

$$= e_H \quad (19.1.12)$$

as desired, we find that  $g \circ k \circ g^{-1} \in \text{Ker}(\phi)$ . Hence  $(\text{Ker}(\phi), \circ) \trianglelefteq (G, \circ)$  ■

### Theorem 20.12 (Injectivity of homomorphism)

Let  $\phi : (G, \circ) \rightarrow (H, *)$  be a homomorphism. Then  $\phi$  is injective iff  $\text{Ker}(\phi) = \{e_G\}$ .

*Proof.* We begin by proving  $\implies$ . Suppose  $\phi$  is injective, and suppose  $g_1, g_2$  are such that  $\phi(g_1) = \phi(g_2) = e_H$ . Then,  $g_1 = g_2$  due to injectivity, and since  $\phi(e_G) = e_H$  it follows that  $\text{Ker}(\phi) = \{e_G\}$ .

Let us now prove the  $\impliedby$  part. Suppose that  $\text{Ker}(\phi) = \{e_G\}$ , and let  $\phi(x) = \phi(y)$  for some  $x, y \in G$ . Then:

$$e_H = \phi(x) * (\phi(y))^{-1} \quad (19.1.13)$$

$$= \phi(x) * \phi(y^{-1}) \quad (19.1.14)$$

$$= \phi(x \circ y^{-1}) \quad (19.1.15)$$

so that  $x \circ y^{-1} \in \text{Ker}(\phi) = \{e_G\}$ . It follows that  $x \circ y^{-1} = e_G \implies x = y$ . Hence,  $\phi$  is injective. ■

### Theorem 20.13 (Normality $\iff$ kernel)

Let  $K \leq G$ , then  $K \trianglelefteq G \iff K = \text{Ker}(\phi)$  for some homomorphism  $\phi$  with domain  $G$ .

*Proof.* We begin by proving  $\implies$ . Suppose  $\phi$  is a homomorphism with domain  $G$ , then by Theorem 20.11  $\text{Ker}(\phi) = K$  is a normal subgroup of  $G$ , as desired.

Now let us prove  $\impliedby$ . Suppose that  $K$  is a normal subgroup of  $G$ . Then it is the kernel of the map  $\phi$  defined by:

$$\phi : (G, \circ) \longrightarrow (G \setminus K, \cdot) \quad (19.1.16)$$

$$x \mapsto xK \quad (19.1.17)$$

so with domain  $G$ , and codomain  $G \setminus K$ , that is the set of cosets of  $K$  in  $G$ . Indeed  $\phi$  is a homomorphism since for all  $x, y \in G$ :

$$\phi(x \circ y) = (x \circ y)K \quad (19.1.18)$$

$$= (xK) \cdot (yK) \quad (19.1.19)$$

$$= \phi(x) \cdot \phi(y) \quad (19.1.20)$$

Also:

$$\text{Ker}(\phi) = \{x \in G : \phi(x) = K\} = \{x \in G : xK = K\} = K \quad (19.1.21)$$

■

## 19.2 First isomorphism theorem

### Theorem 20.14 (Kernel cosets)

Let  $\phi : (G, \circ) \longrightarrow (H, *)$  be a homomorphism, and let  $x, y \in G$ . Then for some  $g \in G$ :

$$\phi(x) = \phi(y) \iff x, y \in g\text{Ker}(\phi) \quad (19.2.1)$$

*Proof.* Firstly, suppose that  $\phi(x) = \phi(y)$  for any  $x, y \in G$ . Then:

$$\phi(x) * \phi(y^{-1}) = \phi(x \circ y^{-1}) = e_H \quad (19.2.2)$$

so that  $x \circ y^{-1} \in \text{Ker}(\phi)$ . Since  $y^{-1} \in G$ , it follows that  $x \in y\text{Ker}(\phi)$ . However, we also have that  $y \in y\text{Ker}(\phi)$ , since  $e_G \in \text{Ker}(\phi)$ . Hence  $x, y$  both belong to the same coset, with  $g = y$ . Similarly, we could have also proven that  $y \in x\text{Ker}(\phi)$ .

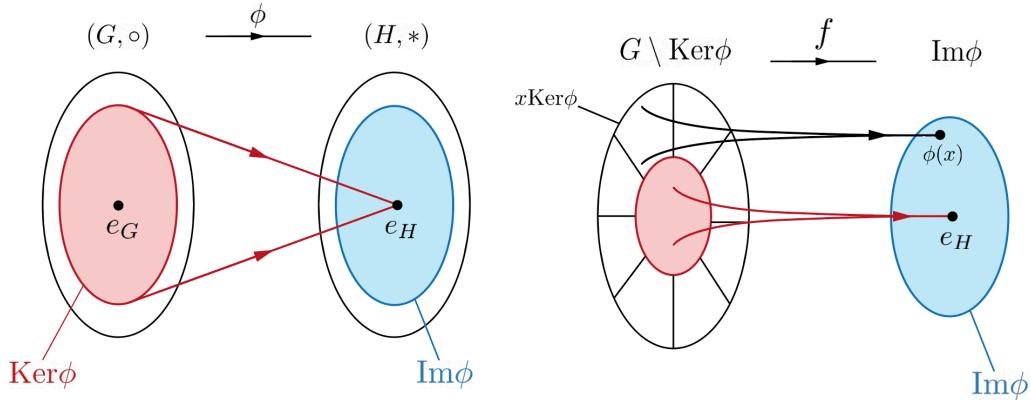
Now suppose that  $x, y$  lie in the same coset of  $\text{Ker}(\phi)$  in  $G$ , so that for some  $g \in G$  and  $k_1, k_2 \in \text{Ker}(\phi)$ :

$$x, y \in g\text{Ker}(\phi) \implies x = g \circ k_1, y = g \circ k_2 \implies x = y \circ k_2^{-1} \circ k_1 \quad (19.2.3)$$

Then:

$$\phi(x) = \phi(y \circ k_2^{-1} \circ k_1) = \phi(y) * (\phi(k_2))^{-1} * \phi(k_1) = \phi(y) \quad (19.2.4)$$

as desired.  $\blacksquare$



**Figure 19.2.** Set diagram of first isomorphism theorem

### Theorem 20.15 (First isomorphism theorem)

Let  $\phi(G, \circ) \rightarrow (H, *)$  be a homomorphism. Then:

$$f : G \setminus \text{Ker}(\phi) \rightarrow \text{Im}(\phi) \quad (19.2.5)$$

$$x\text{Ker}(\phi) \mapsto \phi(x) \quad (19.2.6)$$

is an isomorphism, so that  $G \setminus \text{Ker}(\phi) \cong \text{Im}(\phi)$ .

*Proof.* For sake of brevity, let  $K = \text{Ker}(\phi)$ .

Since elements of different cosets of  $K$  have different images under  $\phi$  (converse of theorem 20.13), we find that  $\phi$  must be injective. Indeed, suppose that  $\phi(x) = \phi(y)$ , so that:

$$x \in y\text{Ker}(\phi) = \{yk_1, yk_2, \dots\} \quad (19.2.7)$$

$$y \in x\text{Ker}(\phi) = \{xk_1, xk_2, \dots\} \quad (19.2.8)$$

from which it follows that  $x = yk_i$  and  $y = xk_i^{-1} = xk_j$  for some  $k_i \in \text{Ker}(\phi)$ . Therefore, we have that  $x\text{Ker}(\phi) \subseteq y\text{Ker}(\phi)$ , since if  $xk_n \in x\text{Ker}(\phi)$  then  $xk_n = yk_i k_n \in y\text{Ker}(\phi)$ , using the closure of  $\text{Ker}(\phi)$ . Similarly,  $y\text{Ker}(\phi) \subseteq x\text{Ker}(\phi)$ , and thus  $x\text{Ker}(\phi) = y\text{Ker}(\phi)$  as desired.

Also,  $f$  is surjective, since any element  $\phi(x) \in \text{Im}(\phi)$  is the image under  $f$  of the coset  $xK$ .

Finally, let us check the homomorphism property:

$$f(xK \cdot yK) = f((x \circ y)K) \quad (19.2.9)$$

$$= \phi(x \circ y) \quad (19.2.10)$$

$$= \phi(x) * \phi(y) \quad (19.2.11)$$

$$= f(xK) * f(yK) \quad (19.2.12)$$

as desired. It follows that  $f$  is a bijective homomorphism, hence an isomorphism, so that  $G \setminus \text{Ker}(\phi) \cong \text{Im}(\phi)$ .  $\blacksquare$

**Example.** Consider the following mapping  $\phi$ :

$$\phi : (L, \times) \longrightarrow (\mathbb{R}^*, \times) \quad (19.2.13)$$

$$\begin{pmatrix} a & 0 \\ b & c \end{pmatrix} \longmapsto ac \quad (19.2.14)$$

where  $L$  is the group of lower triangular  $2 \times 2$  matrices.

This is clearly a homomorphism, since for any  $A, B \in L$ :

$$A = \begin{pmatrix} a_1 & 0 \\ b_1 & c_1 \end{pmatrix}, B = \begin{pmatrix} a_2 & 0 \\ b_2 & c_2 \end{pmatrix} \quad (19.2.15)$$

then:

$$AB = \begin{pmatrix} a_1 a_2 & 0 \\ a_2 b_1 + c_1 b_2 & c_1 c_2 \end{pmatrix} \quad (19.2.16)$$

Hence:

$$\phi(AB) = \phi \begin{pmatrix} a_1 a_2 & 0 \\ a_2 b_1 + c_1 b_2 & c_1 c_2 \end{pmatrix} \quad (19.2.17)$$

$$= (a_1 a_2)(c_1 c_2) \quad (19.2.18)$$

$$= \phi(A)\phi(B) \quad (19.2.19)$$

as desired.

The image of  $\phi$  is:

$$\text{Im}(\phi) = \{\phi(A) : A \in L\} = \{ac : a, c \in \mathbb{R}^*\} = \mathbb{R}^* \quad (19.2.20)$$

The kernel of  $\phi$  is:

$$\text{Ker}(\phi) = \{A \in L : \det\{A\} = 1\} = \left\{ \begin{pmatrix} a & 0 \\ b & \frac{1}{a} \end{pmatrix} : a \in \mathbb{R}^*, b \in \mathbb{R} \right\} \quad (19.2.21)$$

By the first isomorphism theorem:

$$L \setminus \text{Ker}(\phi) \cong \text{Im}(\phi) = (\mathbb{R}^*, \times) \quad (19.2.22)$$



**Example.** Consider the following map:

$$\phi : (L, \times) \longrightarrow (L, \times) \quad (19.2.23)$$

$$\begin{pmatrix} a & 0 \\ b & c \end{pmatrix} \mapsto \begin{pmatrix} \frac{1}{a} & 0 \\ 0 & 1 \end{pmatrix} \quad (19.2.24)$$

This is clearly a homomorphism, since for any  $A, B \in L$ :

$$A = \begin{pmatrix} a_1 & 0 \\ b_1 & c_1 \end{pmatrix}, \quad B = \begin{pmatrix} a_2 & 0 \\ b_2 & c_2 \end{pmatrix} \quad (19.2.25)$$

then:

$$AB = \begin{pmatrix} a_1 a_2 & 0 \\ a_2 b_1 + c_1 b_2 & c_1 c_2 \end{pmatrix} \quad (19.2.26)$$

Hence:

$$\phi(AB) = \begin{pmatrix} \frac{1}{a_1 a_2} & 0 \\ 0 & 1 \end{pmatrix} \quad (19.2.27)$$

$$= \begin{pmatrix} \frac{1}{a_1} & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \frac{1}{a_2} & 0 \\ 0 & 1 \end{pmatrix} \quad (19.2.28)$$

$$= \phi(A)\phi(B) \quad (19.2.29)$$

Moreover it is easy to see that:

$$\text{Im}(\phi) = \left\{ \begin{pmatrix} \frac{1}{a} & 0 \\ 0 & 1 \end{pmatrix} : a \in \mathbb{R}^* \right\} \quad (19.2.30)$$

and:

$$\text{Ker}(\phi) = \left\{ \begin{pmatrix} 1 & 0 \\ b & c \end{pmatrix} : b \in \mathbb{R}, c \in \mathbb{R}^* \right\} \quad (19.2.31)$$

Now, by the first isomorphism theorem, we have that:

$$L \setminus \text{Ker}(\phi) \cong \left\{ \begin{pmatrix} \frac{1}{a} & 0 \\ 0 & 1 \end{pmatrix} : a \in \mathbb{R}^* \right\} \quad (19.2.32)$$

We now prove that  $\text{Im}(\phi) \in (\mathbb{R}^*, \times)$ . An example of an isomorphism between them is:

$$\varphi : \text{Im}(\phi) \longrightarrow (\mathbb{R}^*, \times) \quad (19.2.33)$$

$$\begin{pmatrix} \frac{1}{a} & 0 \\ 0 & 1 \end{pmatrix} \mapsto a \quad (19.2.34)$$

This is a homomorphism since for  $A, B \in \text{Im}(\phi)$ :

$$A = \begin{pmatrix} \frac{1}{a} & 0 \\ 0 & 1 \end{pmatrix}, \quad B = \begin{pmatrix} \frac{1}{b} & 0 \\ 0 & 1 \end{pmatrix} \quad (19.2.35)$$

with  $a \neq 0, b \neq 0$ , we have that

$$\varphi(AB) = ab = \varphi(A)\varphi(B) \quad (19.2.36)$$

Moreover,  $\varphi$  is injective, since  $\varphi(A) = \varphi(B) \implies a = b \implies A = B$ . Finally,  $\varphi$  is surjective, since every  $x \in \mathbb{R}^*$  is the map of the matrix:

$$X = \begin{pmatrix} \frac{1}{x} & 0 \\ 0 & 1 \end{pmatrix} \quad (19.2.37)$$

Hence, we may conclude that  $L \setminus \text{Ker}(\phi) \cong (\mathbb{R}^*, \times)$ .  $\blacktriangleleft$

**Proposition 20.16 (Order of kernel, image and group)** Let  $\phi : (G, \circ) \rightarrow (H, *)$  be a homomorphism with a finite group domain, then:

$$|\text{Ker}(\phi)| \cdot |\text{Im}(\phi)| = |G| \quad (19.2.38)$$

*Proof.* From the first isomorphism theorem, since  $\phi$  is a homomorphism:

$$G \setminus \text{Ker}(\phi) \cong \text{Im}(\phi) \quad (19.2.39)$$

and since isomorphic finite groups must have same order:

$$|G \setminus \text{Ker}(\phi)| = |\text{Im}(\phi)| \quad (19.2.40)$$

Now since  $G$  is finite dimensional, we must have that  $|G \setminus \text{Ker}(\phi)| = \frac{|G|}{|\text{Ker}(\phi)|}$  so that:

$$|\text{Ker}(\phi)| \cdot |\text{Im}(\phi)| = |G| \quad (19.2.41)$$

as desired.  $\blacksquare$

To summarize, we have that for finite groups  $G, H$ , and any homomorphism  $\phi$  between them:

- (i)  $|\text{Ker}(\phi)|$  divides  $|G|$  by Lagrange's theorem, since  $\text{Ker}(\phi) \leq G$
- (ii)  $|\text{Im}(\phi)|$  divides  $|H|$  by Lagrange's theorem, since  $\text{Im}(\phi) \leq H$
- (iii)  $|\text{Im}(\phi)|$  divides  $|G|$  by Proposition 20.16

**Example.** Let us try to find all homomorphisms  $\phi$  from  $S(\Delta)$  to  $(\mathbb{Z}_3, +_3)$ .

From the above considerations, we must have that  $|\text{Ker}(\phi)|$  and  $|\text{Im}(\phi)|$  divide  $|S(\Delta)| = 6$ . Hence they can have values of 1,2,3,6. Moreover, we must have that  $|\text{Im}(\phi)|$  divides 3, hence it can only take values 1,3, for which  $|\text{Ker}(\phi)|$  takes the values of 6,2 respectively.

Now we know that  $|\text{Ker}(\phi)|$  is a normal subgroup of  $S(\Delta)$ , and we found no normal subgroups of order 2. Hence  $|\text{Ker}(\phi)| = 6 = |S(\Delta)|$  and  $\text{Im}(\phi) = 1$ . Consequently, since  $\text{Ker}(\phi) \leq S(\Delta)$  and  $\text{Im}(\phi) \leq \mathbb{Z}_3$ , we have that  $\text{Ker}(\phi) = |S(\Delta)|$  and  $\text{Im}(\phi) = \{0\}$ . The only possible  $\phi$  with such structures is the trivial homomorphism.  $\blacktriangleleft$

# Unit E4: Group actions

## 20.1 What are group actions?

Several groups that we have considered consist of functions from a set to itself. For example, the elements of the group  $S(\square)$  are symmetries, or maps, of the set  $\{1, 2, 3, 4\}$  to itself.

We say that when a group element  $g$  maps an element  $x$  in some set to some other element in the set, then  $g : x \mapsto g \wedge x$ . In other words, we will denote the image of a set element  $x$  under  $g$  by  $g \wedge x$ .

In the case of  $S(\square)$ , we have for example that  $r \wedge 2 = 3$ .

For some definitions of  $\wedge$ , there are a set of interesting properties which promote it from a simple mapping to a **group action**.

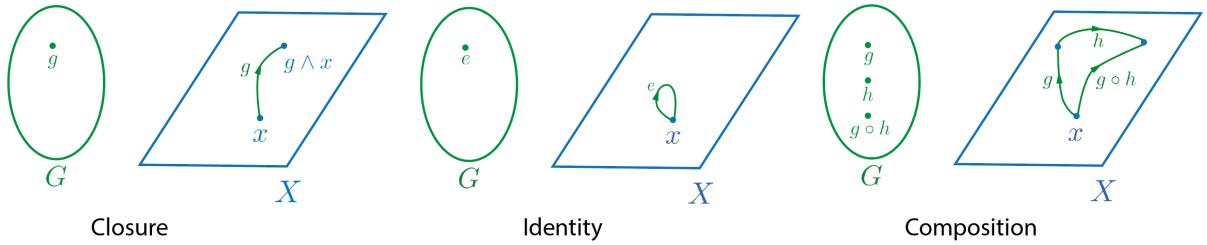


Figure 20.1. Visual illustration of group action axioms.

### Definition 21.1 (Group action)

Let  $(G, \circ)$  be a group with identity  $e$ , and let  $X$  be a set. Furthermore, suppose that we associate to each element  $g \in G$  and some element  $x \in X$  an object  $g \wedge x$ .

We then say that the effect of  $\wedge$  of  $(G, \circ)$  on  $X$  is a **group action of  $(G, \circ)$  on  $X$** , provided that the following properties:

**GA1 Closure:** for  $g \in G, x \in X$  we have that  $g \wedge x \in X$

**GA2 Identity:** for  $x \in X$ , we have that  $e \wedge x = x$

**GA3 Composition:** for  $g, h \in G$  and  $x \in X$  we have that  $g \wedge (h \wedge x) = (g \circ h) \wedge x$ .

known as the **group action axioms** are satisfied.

**Example.** Let  $G \leq S_5$  be the subgroup consisting of all permutations in  $S_5$  that fix symbols

4 or 5, or transpose them. Consider the set  $X = \{1, 2, 3\}$ , and let us define:

$$g \wedge x = g(x) \quad (20.1.1)$$

for all  $g \in G, x \in X$ . In other words,  $G$  is the set of permutations which does not map 1, 2, 3 to 4 or 5.

We check that the group action axioms are satisfied:

**GA1** for  $g \in G, x \in X$  we have that  $g \wedge x = g(x) \in X$ , since  $g$  must be a permutation of  $\{1, 2, 3, 4, 5\}$  which does not map any of  $x \in X = \{1, 2, 3\}$  to 4 or 5.

**GA2** the identity element of  $G$  must be  $\text{id}$ , the identity permutation, defined so that  $\text{id}(x) = \text{id} \wedge x = x$  for all  $x \in X$ .

**GA3** For  $g, h \in G$  and  $x \in X$  we have that:

$$g \wedge (h \wedge x) = g \wedge h(x) \quad (20.1.2)$$

$$= g(h(x)) = (g \circ h)(x) \quad (20.1.3)$$

$$= (g \circ x) \wedge x \quad (20.1.4)$$

as desired.

It follows that  $\wedge$  is indeed a group action of  $(G, \circ)$  on  $X$ . ◀

**Example.** Consider now the mapping  $\wedge$  of  $(\mathbb{R}^*, \times)$  on  $\mathbb{R}^2$  defined by:

$$g \wedge (x, y) = (x + g, y + g) \quad (20.1.5)$$

This is not a group action because it does not satisfy the composition axiom. Indeed  $\forall g, h \in (\mathbb{R}^*, \times), \forall (x, y) \in \mathbb{R}^2$  we find that:

$$g \wedge (h \wedge x) = g \wedge (x + h, y + h) = (x + h + g, y + h + g) \quad (20.1.6)$$

whereas:

$$(g \times h) \wedge x = (gh) \wedge x = (x + gh, y + gh) \neq (x + h + g, y + h + g) \quad (20.1.7)$$

as expected. ◀

**Theorem 21.2 (Properties of group actions)** Let  $\wedge$  be an action of  $G$  on  $X$ , then  $\wedge$  is a bijection, that is:

- (i)  $\forall g \in G$ , if  $x, y \in X$  such that  $g \wedge x = g \wedge y$  then  $x = y$ .
- (ii)  $\forall g \in G$ , if  $y \in X$  then  $\exists x \in X$  such that  $g(x) \in y$

*Proof.* Let  $g \in G$ , and suppose we have  $x, y \in X$  such that  $g \wedge x = g \wedge y$  then  $x = y$ . Then we may apply  $g^{-1} \in G$ :

$$g^{-1} \wedge (g \wedge x) = g^{-1} \wedge (g \wedge y) \implies x = y \quad (20.1.8)$$

using the composition axiom.

Next, let  $g \in G$  and  $y \in X$ . Then:

$$e \wedge y = y \implies g \wedge (g^{-1} \wedge y) = y \quad (20.1.9)$$

and since  $g^{-1} \wedge y \in X$  by the closure of group actions, we have that if  $x = g^{-1} \wedge y \in X$  then  $g \wedge x = y$  as desired.  $\blacksquare$

### Theorem 21.3 (Actions of group of symmetries)

Let  $G$  be the group of symmetries of some figure  $\mathcal{F} \subseteq \mathbb{R}^2$  and let  $X$  be a set of figures in  $\mathbb{R}^2$ . Then, if we define  $\wedge$  by:

$$g \wedge A = g(A), \forall g \in G, A \in X \quad (20.1.10)$$

then  $\wedge$  is a group action *iff* the closure axiom of group actions is satisfied.

*Proof.* We need to prove that GA2 and GA3 are satisfied.

Firstly let us prove that GA2 holds. Let  $e$  be the identity and let  $A \in X$ . Since  $g \in G$  are symmetries of not only  $\mathcal{F}$ , but also  $\mathbb{R}^2$ , it follows that  $e$  will satisfy  $e(P) = P$  for any  $P \in \mathbb{R}^2$ . Consequently  $e \wedge A = e(A) = A$  as desired.

Let  $g, h \in G$  and let  $A \in X$ . Then:

$$g \wedge (h \wedge A) = g(h(A)) = (g \circ h)A \quad (20.1.11)$$

by definition of the composition operation.  $\blacksquare$

**Example.** Consider now the group  $S(\square)$  and the set  $X$  whose elements are all the modified  $2 \times 2$  squares obtained by coloring each of the four small squares either blue, yellow or red.

Then  $g \wedge A = g(A)$  for all  $A \in X$  is clearly a group action since it satisfies the closure axiom. Indeed, suppose we have a square  $A \in X$  with some initial color configuration  $(c_1, c_2, c_3, c_4)$  (in clockwise direction starting from top left small square). Then, the action of each element in  $S(\square)$  are:

Element	$g \wedge A$
$e$	$(c_1, c_2, c_3, c_4)$
$a$	$(c_2, c_3, c_4, c_1)$
$b$	$(c_3, c_4, c_1, c_2)$
$c$	$(c_4, c_1, c_2, c_3)$
$r$	$(c_2, c_1, c_4, c_3)$
$s$	$(c_1, c_4, c_3, c_1)$
$t$	$(c_4, c_3, c_2, c_1)$
$u$	$(c_3, c_2, c_1, c_4)$

It is trivial to verify that  $g \wedge A \in X$  for all  $g \in S(\square)$ , thus proving that  $\wedge$  is a group action.  $\blacktriangleleft$

An effective way to show the action of a symmetry group on some set is by using cycle notation, as we investigate in the following example.

**Example.** Consider the action of  $S(\square)$  on the set  $X = \{R, S, T, U\}$  of figures shown below:



One can easily verify that this is indeed a group action, since the closure property is clearly verified.

We may write down the effect of each element in  $S(\square)$  on  $X$  using cycle notation:

Element	Permutation
$e$	$i$
$a$	$(R\ T)(S\ U)$
$b$	$i$
$c$	$(R\ T)(S\ U)$
$r$	$(S\ U)$
$s$	$(R\ T)$
$t$	$(S\ U)$
$u$	$(R\ T)$

◀

## 20.2 Orbits and stabilisers

### Definition 21.4 (Orbit)

Let  $\wedge$  be a group action of  $G$  on a set  $X$ , and let  $x \in X$ . Then, the **orbit** of  $x$  under  $\wedge$  is defined as:

$$\text{Orb } x = \{g \wedge x : g \in G\} \subseteq X \quad (20.2.1)$$

that is, the set of elements in  $X$  which are the image of  $x$  under  $g$ .

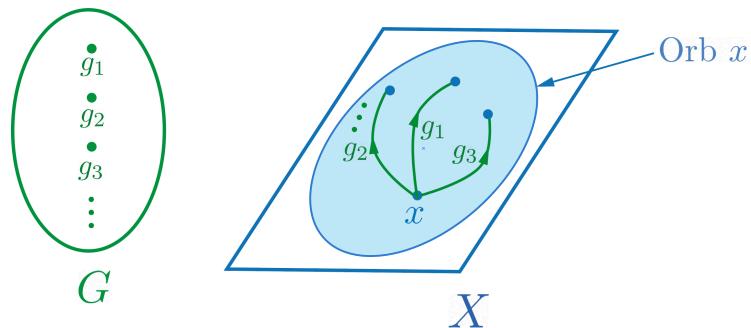


Figure 20.2. Visualization of the orbit of some  $x$  under  $\wedge$ .

**Example.** We consider the action of  $S(\square)$  on the set  $\{1, 2, 3, 4\}$  of labelled vertices of a square. Then:

$$\text{Orb } 1 = \{a \wedge 1, b \wedge 1, c \wedge 1, r \wedge 1, s \wedge 1, t \wedge 1, u \wedge 1\} \quad (20.2.2)$$

$$= \{2, 3, 4, 4, 1, 2, 3\} = \{1, 2, 3, 4\} \quad (20.2.3)$$

Similarly:

$$\text{Orb } 2 = \{a \wedge 2, b \wedge 2, c \wedge 2, r \wedge 2, s \wedge 2, t \wedge 2, u \wedge 2\} \quad (20.2.4)$$

$$= \{3, 4, 1, 3, 4, 1, 2\} = \{1, 2, 3, 4\} \quad (20.2.5)$$

It is easy to verify that  $\text{Orb } 3 = \text{Orb } 4 = \{1, 2, 3, 4\}$  also.  $\blacktriangleleft$

**Example.** Consider the action of  $S(\bigcirc)$  on the plane  $\mathbb{R}^2$ , where the disc  $\bigcirc$  is placed with its center on the origin.

Then, we find that  $\text{Orb } P$  definitely contains the circle  $\mathcal{C}_P$  centered at the origin passing through  $P$ , which is created by acting members of  $S^+(\bigcirc)$  on  $P$ .

Instead, for the reflections, note that we can obtain all reflection symmetries of the disk by reflecting about the  $y$  axis and composing with all direct symmetries. Consequently, we see that the action of the indirect symmetries on  $P$  is to create again a circle centered at the origin passing through  $P$ , which is created by acting members of  $S^+(\bigcirc)$  on  $P$ .

Consequently  $\text{Orb } P = \mathcal{C}_P$ .  $\blacktriangleleft$

As in the case of conjugacy classes and coset classes, we can prove that orbit classes partition the set  $X$  on which the group action acts.

### Theorem 21.5 (Orbit partition)

Let  $\wedge$  be a group action of  $(G, \circ)$  on a set  $X$ . Then the distinct orbits of  $X$  under  $\wedge$  partition  $X$ .

*Proof.* We define  $\sim$  on  $X$  by:

$$x \sim y \text{ if } y \in \text{Orb } x \quad (20.2.6)$$

and prove that it is an equivalence relation.

Indeed:

- (i) **Reflexivity:** let  $x \in X$ , then we see that  $x \in \text{Orb } x$  since  $x = e \wedge x$ . Hence  $x \sim x$ .
- (ii) **Symmetry:** let  $x, y \in X$  such that  $x \sim y$ . Then  $y \in \text{Orb } x$ , that is  $y = g \wedge x$  for some  $g \in G$ . Then:

$$g^{-1} \wedge y = (g^{-1} \circ g) \wedge x = x \implies x \in \text{Orb } y \quad (20.2.7)$$

since  $g^{-1} \in G$ . Consequently  $y \sim x$ .

- (iii) **Transitivity:** let  $x, y, z \in X$  such that  $x \sim y$  and  $y \sim z$ , or equivalently  $y \in \text{Orb } x$  and  $z \in \text{Orb } y$ . We can write these as:

$$y = g_1 \wedge x, \text{ and } z = g_2 \wedge y \quad (20.2.8)$$

for some  $g_1, g_2 \in G$ . Then:

$$z = g_2 \wedge (g_1 \wedge x) = (g_2 \circ g_1) \wedge x \implies z \in \text{Orb } x \quad (20.2.9)$$

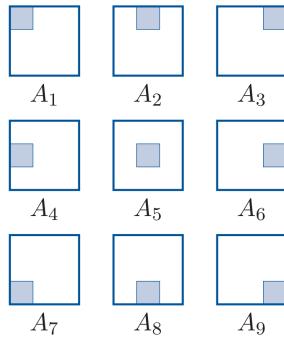
since  $g_2 \circ g_1 \in G$ . Consequently  $x \sim z$  as desired.

It follows that  $\sim$  is an equivalence relation, and that its equivalence classes thus partition  $X$ . The equivalence classes may be expressed as:

$$[x] = \{y : y \in \text{Orb } x\} = \text{Orb } x \quad (20.2.10)$$

so it follows that the distinct orbits of an element  $x \in X$  under  $\wedge$  partition  $X$ . ■

**Example.** Consider the action of  $S(\square)$  on the set  $X = \{A_i : 1 \leq i \leq 9\}$  of squares shown below:



We begin by writing down  $\text{Orb } A_1$ :

$$\text{Orb } A_1 = \{A_1, A_7, A_9, A_3\} \quad (20.2.11)$$

Next, choosing an element that was already included would have resulted in the same set. Hence, we see that  $A_2$  was not included in the above orbit, so we find its orbit, :

$$\text{Orb } A_2 = \{A_2, A_4, A_8, A_6\} \quad (20.2.12)$$

The only remaining element of  $X$  is  $A_5$ , so its orbit must only contain itself:

$$\text{Orb } A_5 = \{A_5\} \quad (20.2.13)$$

Hence the orbit partition of  $X$  is:

$$X = \{A_1, A_3, A_7, A_9\} \cup \{A_2, A_4, A_6, A_8\} \cup \{A_5\} \quad (20.2.14)$$

◀

**Example.** Consider the matrix group

$$G = \left\{ \begin{pmatrix} a & 0 \\ 0 & b \end{pmatrix} : a, b \in \mathbb{R}^+ \right\} \quad (20.2.15)$$

and the group action  $\wedge$  on  $\mathbb{R}^2$  defined by:

$$\begin{pmatrix} a & 0 \\ 0 & b \end{pmatrix} \wedge (x, y) = (ax, by) \quad (20.2.16)$$

for all  $(x, y) \in \mathbb{R}^2$ . It follows that:

$$\text{Orb } (x, y) = \{(ax, by) : a, b \in \mathbb{R}^+\} \quad (20.2.17)$$

For example:

$$\text{Orb } (1, 0) = \{(a, 0) : a \in \mathbb{R}^+\} \quad (20.2.18)$$

which is the positive  $x$ -axis excluding the origin. By symmetry,  $\text{Orb } (-1, 0)$  must be the negative  $x$ -axis excluding the origin.

Similarly:

$$\text{Orb } (0, -1) = \{(0, -b) : b \in \mathbb{R}^+\} \quad (20.2.19)$$

which is the negative  $y$ -axis excluding the origin. By symmetry,  $\text{Orb } (0, 1)$  must be the positive  $y$ -axis excluding the origin.

One element of  $\mathbb{R}^2$  which we did not include is the origin. We can guess that it is the orbit of the origin, indeed:

$$\text{Orb } (0, 0) = \{(0, 0)\} \quad (20.2.20)$$

We are missing the four quadrants of  $\mathbb{R}^2$ . Note that

$$\text{Orb } (1, 1) = \{(a, b) \in \mathbb{R}^+ : a, b \in \mathbb{R}^+\} \quad (20.2.21)$$

which is the upper right quadrant excluding the origin and the axes. By symmetry,  $\text{Orb } (1, -1)$  must be the lower right quadrant,  $\text{Orb } (-1, 1)$  must be the upper left quadrant and  $\text{Orb } (-1, -1)$  must be the lower left quadrant.  $\blacktriangleleft$

### Definition 21.6 (Stabiliser)

Let  $\wedge$  be an action of a group  $G$  on a set  $X$ , and let  $x \in X$ . Then, the **stabiliser** of  $x$  under  $\wedge$  is defined as:

$$\text{Stab } x = \{g \in G : g \wedge x = x\} \quad (20.2.22)$$

We may interpret the stabiliser geometrically as shown below:

**Example.** Let's consider the action of  $S(\bigcirc)$  on  $\mathbb{R}^2$ , where the disc  $\bigcirc$  is centered at the origin, and let  $P \in \mathbb{R}^2$ . We consider the following possibilities:

- (i)  $P$  is the origin: then  $\text{Stab } P = S(\bigcirc)$
- (ii)  $P$  is not the origin: then  $\text{Stab } P = \{e, q\}$  contains  $e$  and the reflection in the line containing  $P$  and  $O$ , which we call  $q$ .

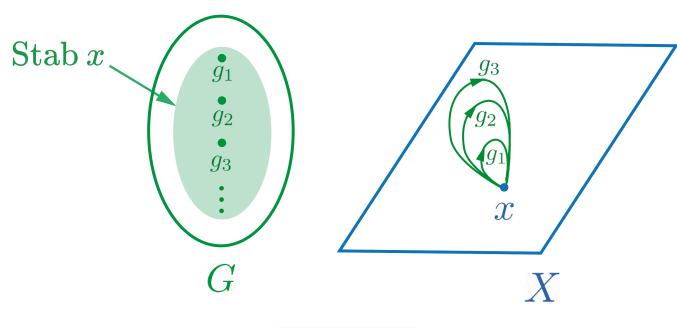


Figure 20.3. Geometrical interpretation of the stabiliser

Interestingly, the stabilisers we have found in the previous example can be verified to be subgroups of  $S(\square)$ . This is no coincidence, as the following theorem shows.

**Theorem 21.7 (Stabiliser subgroup)**

Let  $\wedge$  be an action of  $(G, \circ)$  on  $X$ . Then, for any  $x \in X$ ,  $\text{Stab } x \leq G$ .

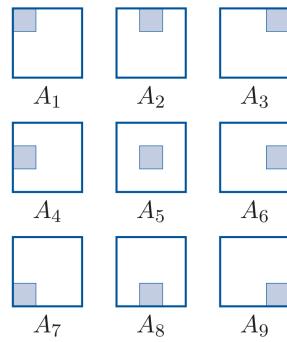
*Proof.* We show that the subgroup axioms are satisfied:

**Closure:** let  $g, h \in \text{Stab } x$ . Then  $(g \circ h) \wedge x = g \wedge (h \wedge x) = g \wedge x = x$  so  $g \circ h \in \text{Stab } x$  as desired.

**Identity:** let  $e$  be the identity of  $G$ . Then  $e \wedge x = x$  by the group action axioms, hence  $e \in \text{Stab } x$ .

**Inverses:** let  $g \in \text{Stab } x$ , then  $g \wedge x = x \implies x = g^{-1} \wedge x$  so that  $g^{-1} \in \text{Stab } x$ . ■

**Example.** Consider the action of  $S(\square)$  on the set  $X = \{A_i : 1 \leq i \leq 9\}$  of squares shown below:



Then:

$$\text{Stab } A_1 = \{e, s\} = \text{Stab } A_9 \quad (20.2.23)$$

$$\text{Stab } A_2 = \{e, r\} = \text{Stab } A_8 \quad (20.2.24)$$

$$\text{Stab } A_3 = \{e, u\} = \text{Stab } A_7 \quad (20.2.25)$$

$$\text{Stab } A_4 = \{e, t\} = \text{Stab } A_6 \quad (20.2.26)$$

$$\text{Stab } A_5 = S(\square) \quad (20.2.27)$$

These are all subgroups of  $S(\square)$ , since the first four lines contain  $e$  and a reflection (which are self inverse).  $\blacktriangleleft$

**Example.** Consider the action of  $G = \left\{ \begin{pmatrix} a & b \\ 0 & a \end{pmatrix} : a, b \in \mathbb{R}, a \neq 0 \right\}$  on  $\mathbb{R}^2$  defined by:

$$\begin{pmatrix} a & b \\ 0 & a \end{pmatrix} \wedge (x, y) = (ax, ay), \quad \forall (x, y) \in \mathbb{R}^2 \quad (20.2.28)$$

Then, we have that if  $(x, y)$  is not the origin:

$$\text{Stab } (x, y) = \left\{ \begin{pmatrix} a & b \\ 0 & a \end{pmatrix} \in G : (ax, ay) = (x, y) \right\} = \left\{ \begin{pmatrix} 1 & b \\ 0 & 1 \end{pmatrix} : b \in \mathbb{R} \right\} \quad (20.2.29)$$

If instead  $(x, y) = (0, 0)$  then:

$$\text{Stab } (0, 0) = \left\{ \begin{pmatrix} a & b \\ 0 & a \end{pmatrix} \in G : (a \cdot 0, a \cdot 0) = (0, 0) \right\} = G \quad (20.2.30)$$

$\blacktriangleleft$

## 20.3 The Orbit-Stabiliser theorem

**Theorem 21.8 (Left coset of stabiliser)** Let  $\wedge$  be an action of  $(G, \circ)$  on a set  $X$ , and let  $x \in X$ , and  $g, h \in G$ , then:

$$g \wedge x = h \wedge x \iff g, h \text{ lie in the same left coset of } \text{Stab } x \quad (20.3.1)$$

*Proof.* We firstly prove  $\implies$ . Indeed, suppose  $g, h$  lie in the same left coset of  $\text{Stab } x$ , so that  $h \in \text{Stab } x$ . Then  $\exists k \in \text{Stab } x$  such that  $h = g \circ k$  and thus:

$$h \wedge x = g \wedge k \wedge x = g \wedge x \quad (20.3.2)$$

as desired.

Now suppose that  $h \wedge x = g \wedge x$ . Then:

$$(g^{-1} \circ h) \wedge x = g^{-1} \wedge (h \wedge x) \quad (20.3.3)$$

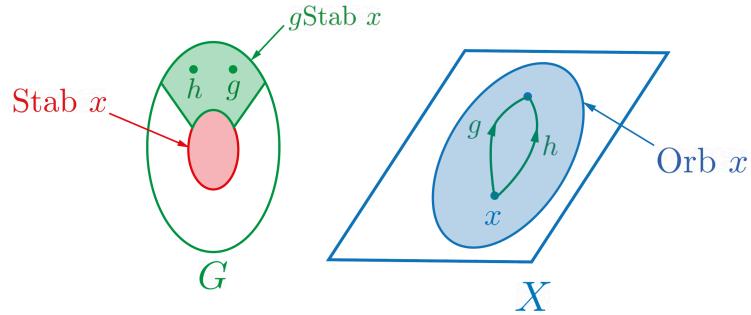
$$= g^{-1} \wedge (g \wedge x) \quad (20.3.4)$$

$$= (g^{-1} \circ g) \wedge x \quad (20.3.5)$$

$$= e \wedge x \quad (20.3.6)$$

$$= x \quad (20.3.7)$$

so that  $g^{-1} \circ h \in \text{Stab } x$ , and consequently  $h \in g \text{Stab } x$ , so that  $h, g$  both lie in the same left coset.  $\blacksquare$



**Figure 20.4.** Visual interpretation of the left coset  $g\text{Stab } x$  and its action on  $X$

**Example.** We consider the action of  $S_3$  on  $\{1, 2, 3\}$ , and find that:

$$\text{Stab } 1 = \{\sigma \in S_3 : \sigma(1) = 1\} = \{e, (2\ 3)\} \quad (20.3.8)$$

Then, the left cosets of this stabiliser are:

$$e\text{Stab } 1 = \text{Stab } 1 \quad (20.3.9)$$

$$(1\ 2)\text{Stab } 1 = \{(1\ 2), (1\ 2\ 3)\} \quad (20.3.10)$$

$$(1\ 3)\text{Stab } 1 = \{(1\ 3), (1\ 3\ 2)\} \quad (20.3.11)$$

$$(2\ 3)\text{Stab } 1 = \text{Stab } 1 \quad (20.3.12)$$

$$(1\ 2\ 3)\text{Stab } 1 = \{(1\ 2), (1\ 2\ 3)\} \quad (20.3.13)$$

$$(1\ 3\ 2)\text{Stab } 1 = \{(1\ 3), (1\ 3\ 2)\} \quad (20.3.14)$$

Now we partition  $S_3$  according to where its elements map 1:

$$\{e, (2\ 3)\} \text{ map 1 to 1} \quad (20.3.15)$$

$$\{(1\ 2), (1\ 2\ 3)\} \text{ map 1 to 2} \quad (20.3.16)$$

$$\{(1\ 3), (1\ 3\ 2)\} \text{ map 1 to 3} \quad (20.3.17)$$

which are precisely the left cosets we calculated earlier. ◀

### Proposition 21.9 (*Stabiliser coset to orbit map*)

Let  $\wedge$  be an action of  $G$  on the set  $X$  and let  $x \in X$ . Then:

$$\phi : G\text{Stab } x \longrightarrow \text{Orb } x \quad (20.3.18)$$

$$g\text{Stab } x \mapsto g \wedge x \quad (20.3.19)$$

is a bijective map.

*Proof.* Since elements of different cosets of  $\text{Stab } x$  have different images under  $\phi$ , we find that  $\phi$

must be injective. Indeed, suppose that  $\phi(g\text{Stab } x) = \phi(h \wedge x)$ , so that  $g \wedge x = h \wedge x$ :

$$g \in h\text{Stab } x = \{hs_1, hs_2, \dots\} \quad (20.3.20)$$

$$h \in g\text{Stab } x = \{gs_1, gs_2, \dots\} \quad (20.3.21)$$

from which it follows that  $g = hs_i$  and  $h = gs_i^{-1} = gs_j$  for some  $k_i, k_j \in \text{Stab } x$ . Therefore, we have that  $g\text{Stab } x \subseteq h\text{Stab } x$ , since if  $gs_n \in g\text{Stab } x$  then  $gs_n = hs_is_n \in h\text{Stab } x$ , using the closure of  $\text{Stab } x$ . Similarly,  $h\text{Stab } x \subseteq g\text{Stab } x$ , and thus  $g\text{Stab } x = h\text{Stab } x$  as desired.

Also,  $f$  is surjective, since any element  $g \wedge x \in \text{Im}(\phi)$  is the image under  $f$  of the left coset  $x\text{Stab } x$ . ■

### Theorem 21.10 (Orbit-Stabiliser theorem)

Suppose that  $G$  is a finite group acting on the set  $X$ . Then:

$$\forall x \in X, |\text{Orb } x| \times |\text{Stab } x| = |G| \quad (20.3.22)$$

*Proof.* Let  $x \in X$ , we know from Proposition 21.9 that the left cosets of  $\text{Stab } x$  in  $G$  have a bijective correspondence with the elements of  $\text{Orb } x$ . It follows that  $|G\text{Stab } x| = |\text{Orb } x|$ , the number of distinct left cosets of  $\text{Stab } x$  is equal to the number of elements in  $\text{Orb } x$ . However, since  $G$  is a finite group  $|G\text{Stab } x| = \frac{|G|}{|\text{Stab } x|}$  so that:

$$|\text{Orb } x| \cdot |\text{Stab } x| = |G| \quad (20.3.23)$$

Interestingly, our choice of  $X$  is not limited to sets. Indeed, we can consider group actions on groups themselves, in other words  $X$  can be a group.

One example of such a group action is conjugation.

### Proposition 21.11 (Conjugation group action)

Let  $G$  be a group with  $g, x \in G$  and define  $\wedge$  by:

$$g \wedge x = gxg^{-1} \quad (20.3.24)$$

Then  $\wedge$  is a group action.

*Proof.* We prove that the three group action axioms hold.

**GA1 Closure:** let  $g, x \in G$ . Then  $g \wedge x = gxg^{-1} \in G$

**GA2 Identity:** let  $x \in G$  and let  $e \in G$  be the identity element. Then  $e \wedge x = exe^{-1} = x$  as desired.

**GA3 Composition:** let  $g, h, x \in G$ . Then:

$$g \wedge (h \wedge x) = g \wedge (hxh^{-1}) \quad (20.3.25)$$

$$= g(hxh^{-1})g^{-1} \quad (20.3.26)$$

$$= (gh)x(gh)^{-1} \quad (20.3.27)$$

$$= (gh) \wedge x \quad (20.3.28)$$

as desired. ■

**Example.** We prove that  $h \wedge g = hg$  is a group action, where  $h \in H, g \in G$  and  $H \leq G$ .

Indeed:

**GA1 Closure:** let  $g \in G$  and  $h \in H \implies h \in G$ . Then  $g \wedge h = gh \in G$  as desired.

**GA2 Identity:** let  $g \in G$  and let  $e \in H$  be the identity element of  $H$ , and thus of  $G$  too. Then  $e \wedge g = eg = g$  as desired.

**GA3 Composition:** let  $g, h, f \in G$ . Then:

$$f \wedge (h \wedge g) = f \wedge (hx) \quad (20.3.29)$$

$$= fhx \quad (20.3.30)$$

$$= (fh)x \quad (20.3.31)$$

$$= (fh) \wedge x \quad (20.3.32)$$

where we used the associativity in  $G$ . ◀

### Proposition 21.12 (Cardinality of conjugacy class)

For a finite group  $G$ , the number of elements in each conjugacy class divides  $|G|$ .

*Proof.* Let  $G$  be a finite group and let  $\wedge$  be the conjugacy action  $g \wedge x = gxg^{-1}$  for  $g, x \in G$ . Then:

$$\text{Orb } x = \{gxg^{-1} : g \in G\} = [x] \quad (20.3.33)$$

so the orbit of  $x$  is the conjugacy class of  $x$ . We also have that:

$$\text{Orb } x \text{ divides } |G| \quad (20.3.34)$$

giving the desired result. ■

### Proposition 21.13 (Homomorphism group action)

Let  $\phi : (G, \circ) \rightarrow (H, *)$  be a homomorphism, and let  $\wedge$  be defined as:

$$g \wedge h = \phi(g) * h \quad (20.3.35)$$

for  $g \in G, h \in H$ . Then we have that  $\wedge$  is a group action.

*Proof.* We show that the three group action axioms are satisfied:

**GA1 Closure:** let  $g \in G, h \in H$ , then  $g \wedge h = \phi(g) * h \in H$  since  $\phi(g) \in H$ .

**GA2 Identity:** let  $e_G \in G$  be the identity of  $G$  and let  $h \in H$ . Then  $e_G \wedge h = \phi(e_G) * h = e_H * h = h$  as desired.

**GA3 Composition:** let  $g, f \in G$  and  $h \in H$ . Then:

$$g \wedge (f \wedge h) = g \wedge (\phi(f) * h) \quad (20.3.36)$$

$$= \phi(g) * (\phi(f) * h) \quad (20.3.37)$$

$$= \phi(g \circ f) * h \quad (20.3.38)$$

$$= (g \circ f) \wedge h \quad (20.3.39)$$

as desired. ■

Notice that for the homomorphism group action of  $G$  on  $H$ :

$$\text{Orb } e_H = \{\phi(g) * e_H : g \in G\} = \{\phi(g) : g \in G\} = \text{Im}(\phi) \quad (20.3.40)$$

and:

$$\text{Stab } e_H = \{g : \phi(g) * e_H = e_H\} = \{g : \phi(g) = e_H\} = \text{Ker}(\phi) \quad (20.3.41)$$

Applying the orbit stabiliser theorem:

$$|\text{Orb } e_H| \cdot |\text{Stab } e_H| = |\text{Im}(\phi)| \cdot |\text{Ker}(\phi)| = |G| \quad (20.3.42)$$

which is precisely the result proven in Proposition 20.16.

## 20.4 The Counting theorem

Consider a  $2 \times 2$  square pattern, where each of the four smaller squares is colored either blue, yellow, red, green or purple?

We see that since repetitions are allowed, each tile has 5 different possible colours. Therefore, we should have  $5^4 = 625$  different patterns.

However, note that the following two patterns, which are rotations of each other, were counted twice in our procedure:

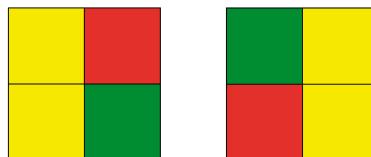


Figure 20.5. Two identical patterns which were double counted.

Surprisingly, this problem has to do with group actions. Indeed, let  $X$  be the set of  $5^4$  colored squares, where each small square is fixed in space. We can think of these squares as the vertices of a larger square, with symmetry group  $S(\square)$ . Two patterns which are double counted are therefore in the same orbit of the action of  $S(\square)$  on  $X$ .

We can reformulate this square coloring problem as finding the number of orbits of the action of  $S(\square)$  on  $X$ .

**Definition 20.17 (Fixed set)**

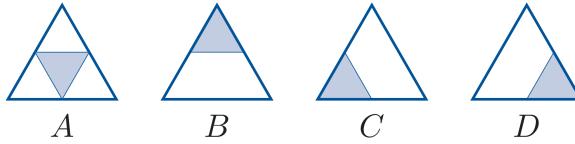
Let  $\wedge$  be a group action of  $G$  on  $X$ , and let  $g \in G$ . Then the **fixed set** of  $g$  under  $\wedge$  is defined as:

$$\text{Fix } g = \{x \in X : g \wedge x = x\} \quad (20.4.1)$$

that is, the subset of  $X$  whose elements are mapped to themselves by  $g$ .

**Example.** Consider the action of  $S(\Delta)$  on the set  $X = \{A, B, C, D\}$  of triangles shown

below:



Then, we see that:

$$\text{Fix } e = X \quad (20.4.2)$$

$$\text{Fix } a = \text{Fix } b = \{A\} \quad (20.4.3)$$

$$\text{Fix } r = \{A, B\} \quad (20.4.4)$$

$$\text{Fix } s = \{A, C\} \quad (20.4.5)$$

$$\text{Fix } t = \{A, D\} \quad (20.4.6)$$

◀

**Example.** Let

$$G = \left\{ \begin{pmatrix} a & b \\ 0 & 1 \end{pmatrix} : a, b \in \mathbb{R}, a \neq 0 \right\} \quad (20.4.7)$$

and consider the action of  $G$  on  $\mathbb{R}^2$  defined by:

$$\begin{pmatrix} a & b \\ 0 & 1 \end{pmatrix} \wedge (x, y) = (ax, y) \quad (20.4.8)$$

Then, since  $a \neq 0$  it follows that

$$\text{Fix } \begin{pmatrix} a & b \\ 0 & 1 \end{pmatrix} = \{(x, y) \in \mathbb{R}^2 : (x, y) = (ax, y)\} = \{(0, y) : y \in \mathbb{R}\} \quad (20.4.9)$$

if  $a \neq 1$  and

$$\text{Fix } \begin{pmatrix} 1 & b \\ 0 & 1 \end{pmatrix} = \mathbb{R}^2 \quad (20.4.10)$$

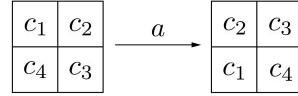
if  $a = 1$ . ▶

Let's use the notion of fixed sets in the context of the  $2 \times 2$  square pattern problem.

**Example.** Let's find the number of elements in the fixed sets of each symmetry in  $S(\square)$ . Clearly,  $\text{Fix } e = X$ , so that  $|\text{Fix } e| = 5^4$ .

Moreover,  $\text{Fix } a$  is the set of squares where all squares are colored in the same way. To see

why this must be the case, let us number the four squares in a pattern in  $\text{Fix } a$  by  $c_1, c_2, c_3, c_4$  representing their colors. If we rotate this pattern, we find that:



so that  $c_1 = c_2, c_4 = c_1, c_3 = c_4, c_2 = c_3$  implying that all colors must be the same. There are 5 such squares.

By similar arguments, we can see that  $\text{Fix } b$  must be the set of squares where the diagonal small squares are of the same color. Note that the color of small squares on different diagonals need not to be the same, hence there are  $5^2$  such squares in this fixed set.

Similarly,  $\text{Fix } c$  contains 5 elements,  $|\text{Fix } r| = |\text{Fix } t| = 5^2$ . Finally,  $|\text{Fix } s| = |\text{Fix } u| = 5^3$ . We summarize these results in the table below:

$g \in S(\square)$	$ \text{Fix } g $
$e$	$5^4$
$a$	5
$b$	$5^2$
$c$	5
$r$	$5^2$
$s$	$5^3$
$t$	$5^2$
$u$	$5^3$

◀

### Theorem 20.18 (Counting theorem)

Let  $\wedge$  be an action of a finite group  $G$  on the set  $X$ , then the number of orbits of  $\wedge$  is:

$$\frac{1}{|G|} \sum_{g \in G} |\text{Fix } g| \quad (20.4.11)$$

*Proof.* Let  $t$  be the number of orbits, and let  $B$  be one such orbit. Then:

$$\sum_{x \in B} |\text{Stab } x| = \sum_{x \in B} \frac{|G|}{|\text{Orb } x|} \quad (20.4.12)$$

$$= |G| \sum_{x \in B} \frac{1}{|\text{Orb } x|} \quad (20.4.13)$$

$$= |G| \sum_{x \in B} \frac{1}{|B|} \quad (20.4.14)$$

$$= |G| \cdot |B| \cdot \frac{1}{|B|} \quad (20.4.15)$$

$$= |G| \quad (20.4.16)$$

Therefore, since orbits partition  $X$ :

$$\sum_{x \in X} |\text{Stab } x| = t|G| \iff t = \frac{1}{|G|} \sum_{x \in X} |\text{Stab } x| \quad (20.4.17)$$

However, we have that:  $\sum_{x \in X} |\text{Stab } x| = \sum_{g \in G} |\text{Fix } g|$

Indeed, suppose we construct a table with a row heading containing  $x \in X$  and with a column heading containing  $g \in G$ . We place a  $y$  at each position where  $g \wedge x = x$ . Then  $\sum_{x \in X} |\text{Stab } x|$  corresponds in counting the number of ticks in each column labelled  $x$ , and summing them all up. This surely must be equivalent to counting the number of ticks in each row labelled  $g$  and summing them all up, which corresponds to  $\sum_{g \in G} |\text{Fix } g|$ .

	... $x$ ...
⋮	⋮
$g$	...    ✓    ✓    ...    ✓    ...    ✓    ...
⋮	⋮
	✓
	✓
	⋮

■

**Example.** Let us return to the problem of coloring a  $2 \times 2$  square. We need to be careful in defining what we mean by two squares being the same. In this particular example, we consider two squares as being the same if one can be rotated or flipped to give the other. We therefore need to find the number of orbits of the action of  $S(\square)$ :

$$\frac{1}{8}(5^4 + 5 + 5^2 + 5 + 5^2 + 5^3 + 5^2 + 5^3) = 120 \quad (20.4.18)$$

so there are 120 different patterns. ◀

Let us apply our results to one final coloring problem.

**Example.** Let's see how many different ways there are to color a cube's faces using three colors. We consider two cubes identical if one can be rotated to give the other (but not reflected, obviously).

First, we need to find the fixed sets of each element in  $S^+(\text{cube})$  (we do not consider reflections). The elements in  $S^+(\text{cube})$  are:

- (a) identity symmetry
- (b) rotations by  $\pm \frac{\pi}{2}$  about axes through centers of opposite faces
- (c) rotations by  $\pi$  about axes through centers of opposite faces
- (d) rotations by  $\pm \frac{2\pi}{3}$  about axes through opposite vertices
- (e) rotations by  $\pi$  about axes through midpoints of opposite edges.

We see that there is 1 symmetry of type a, 6 rotations of type b, 3 rotations of type c, 8

rotations of type d and 6 rotations of type e.

For each, we find through the labelling method that:

type $g \in S(\square)$	$ \text{Fix } g $
(a)	$3^6$
(b)	$3^3$
(c)	$3^4$
(d)	$3^2$
(e)	$3^3$

so that the number of orbits is:

$$\frac{1}{24}(3^6 + 6 \cdot 3^3 + 3 \cdot 3^4 + 8 \cdot 3^2 + 6 \cdot 3^3) = 57 \quad (20.4.19)$$

hence there are 57 different colored cubes. ◀

---

# Sylow theorems

21

---

# Rings

22

---

# Polynomials

23

---

# Modules

24

## **Part III**

# **Representation Theory and Lie Algebra**

## **Part IV**

# **Differential Equations**

# Fundamentals

## 25.1 Definitions

### Definition 23.1 (*n*th order ODE)

An *n*th order ordinary differential equation (ODE) in  $\mathbb{K}^N$  is an equation:

$$y^{(n)} = f(t, y, y' \dots y^{(n-1)}) \quad (25.1.1)$$

for  $(t, y, y' \dots y^{(n-1)}) \in \Omega, t \in \mathbb{R}, y \in \mathbb{K}^n$  where  $\Omega \subseteq \mathbb{R} \times (\mathbb{K}^N)^n$  and  $f : \Omega \rightarrow \mathbb{K}^n$ , with  $n, N \in \mathbb{N}^*$ .

We highlight the fundamental case where  $n = 1$ , in which case:

$$y' = f(t, y), \quad (t, y) \in \Omega, \quad t \in \mathbb{R}, \quad y \in \mathbb{K}^N \quad (25.1.2)$$

### Definition 23.2 (*Solution*)

A solution to an *n*th order ODE consists of a function  $y : I \rightarrow \mathbb{K}^N$  *n* times differentiable such that:

- (i)  $\forall t \in I, (t, y, y' \dots y^{(n-1)}) \in \Omega$
- (ii)  $\forall t \in I, y^{(n)} = f(t, y, y' \dots y^{(n-1)})$

### Definition 23.3 (*Linear and Homogeneous*)

An ODE is said to be **linear** if  $f$  is a polynomial function, that is:

$$f(t, y, y' \dots y^{(n-1)}) = \sum_{i=0}^n A_i(t) y^{(i)} \quad (25.1.3)$$

where  $A_i(t) \in \text{Mat}_N(\mathbb{K})$  for  $i = 1, 2..n$  and  $A_0(t) \in \mathbb{K}^N$ . A linear ODE is further said to be **homogeneous** if  $A_0(t) = 0$ . Homogeneous solutions are invariant under scaling  $(t, y, y' \dots y^{(n-1)}) \rightarrow (\lambda t, \lambda y, \lambda y' \dots \lambda y^{(n-1)})$  for  $\lambda \in \mathbb{K}$ .

### Proposition 23.4 (*Smoothness of Solutions*)

If  $f : \Omega \rightarrow \mathbb{K}^N$  is of class  $C^k$  then all solutions  $y$  of (23.0.2) are of class  $C^{k+1}$ .

*Proof.* We provide a proof by induction. For  $k = 0$ , then  $f$  is continuous everywhere, and let  $y$  be

a solution. Then, we have that  $y' = f(t, y)$  and so  $y$  is continuous and differentiable everywhere, hence of class  $C^1$ .

Let us now suppose that proposition 23.4 is true for some  $k \in \mathbb{N}$ , let  $f$  be of class  $C^{k+1}$  and let  $y$  be a solution. Then, since  $f$  is also of class  $C^k$ , then by hypothesis  $y$  must be of class  $C^{k+1}$ . However,  $y' = f(t, y)$  is of class  $C^{k+1}$  and  $y$  is therefore of class  $C^{k+2}$ . ■

## 25.2 Integral formulation

There is a remarkable relationship between differential equations and integral equations. In this course we will consider four main types of integral equations. Suppose that  $f : [a, b] \rightarrow \mathbb{R}$  and  $K : [a, b]^2 \rightarrow K$  are continuous, with  $t \in [a, b]$  then:

$$\text{Volterra non-homogeneous : } y(t) = f(t) + \int_a^t K(t, s)y(s)ds \quad (25.2.1)$$

$$\text{Fredholm non-homogeneous : } y(t) = f(t) + \lambda \int_a^b K(t, s)y(s)ds \quad (25.2.2)$$

and the two corresponding homogeneous equations. We call  $K(t, s)$  the **kernel** of the integral equation.

Note that for the Fredholm homogeneous equation:

$$y(t) = \lambda \int_a^b K(t, s)y(s)ds \quad (25.2.3)$$

we may consider this as an eigenfunction equation, with  $\lambda$  as an eigenvalue and  $y$  as an eigenfunction.

**Lemma 1** Suppose that  $f : [a, b] \rightarrow \mathbb{R}$  is continuous. Then:

$$\int_a^x \int_a^{x'} f(t)dt dx' = \int_a^x (x-t)f(t)dt \quad (25.2.4)$$

*Proof.* Define the integral transform  $G : [a, b] \rightarrow \mathbb{R}$  by:

$$F(x) = \int_a^x (x-t)f(t)dt \quad (25.2.5)$$

then because  $f(t)$  and  $x-t$  are continuous we may use the Leibniz integral rule to find:

$$F'(x) = \underbrace{[(x-t)f(t)]}_{t=x} \xrightarrow{d}{dx}(x) + \int_a^x \frac{\partial}{\partial x}[(x-t)f(t)]dt = \int_a^x f(t)dt \quad (25.2.6)$$

We then deduce from the fundamental theorem of Calculus that:

$$\int_a^{x'} F'(x)dx = F(x') - \cancel{F(a)}^0 = \int_a^{x'} \int_a^x f(t)dt dx \quad (25.2.7)$$

Substituting  $x \rightarrow x'$  then:

$$F(x) = \int_a^x \int_a^{x'} f(t)dt dx' \quad (25.2.8)$$

as we wished to show. ■

Consider the differential equation  $y'' + \lambda y = g(t)$  with  $t \in [0, L]$ . The reader will probably be familiar already with the solution, but here we wish to find the equivalent integral equation.

The first step in doing so is integrating from 0 to  $x$  to find (using the fact that  $y, y''$  must both be continuous):

$$y'(t) - y'(0) + \lambda \int_0^t y(s)ds = \int_0^t g(s)ds \quad (25.2.9)$$

Further integration gives:

$$y(t) - y(0) - ty'(0) + \lambda \int_0^t (t-s)f(s)ds = \int_0^t (t-s)g(s)ds \quad (25.2.10)$$

where we used Lemma 1 to simplify the two double integrals.

We must now set some conditions to solve the problem explicitly.

### Definition 23.5 (*Initial and Boundary conditions*)

An *initial condition* is a specification of  $(t, y, y' \dots y^{(n)})$  for  $n$  initial values  $t = t_i$ .

A *boundary condition* is a specification of  $y$  at the end-points of an interval.

1. **Initial condition:** suppose  $y(0) = 0$  and  $y'(0) = A$ . Then:

$$y(t) = At + \int_0^t (t-s)g(s)ds - \lambda \int_0^t (t-s)y(s)ds \quad (25.2.11)$$

which is a Volterra non-homogeneous integral equation with  $K(t, s) = \lambda(t-s)$  and  $f(t) = At + \int_0^x (t-s)g(s)ds$ .

2. **Boundary condition:** suppose  $y(0) = 0$  and  $y(L) = B$ . Then we find upon inserting  $t = L$  that:

$$y'(0) = \frac{1}{L} \left( \lambda \int_0^L (L-s)y(s)ds - \int_0^L (L-s)g(s)ds + B \right) \quad (25.2.12)$$

and substituting back into the original integral equation we find:

$$y = \frac{Bt}{L} - \int_0^L \frac{t}{L} (L-s)g(s)ds + \int_0^t (t-s)g(s)ds + \lambda \left( \int_0^L \frac{t}{L} (L-s)g(s)ds - \int_0^t (t-s)g(s)ds \right) \quad (25.2.13)$$

If we now define a function  $K(s, t)$  such that:

$$\int_0^L K(s, t)f(s)ds = \int_0^L \frac{t}{L} (L-s)f(s)ds - \int_0^t (t-s)f(s)ds \quad (25.2.14)$$

then we may write:

$$y = \underbrace{\frac{Bt}{L} - \int_0^L K(s, t)g(s)ds}_{f(x)} + \lambda \left( \int_0^L K(s, t)y(s)ds \right) \quad (25.2.15)$$

It turns out that the kernel is<sup>1</sup>:

$$K(s, t) = \begin{cases} \frac{s}{L}(L-t) & \text{when } 0 \leq s \leq t \leq L \\ \frac{t}{L}(L-s) & \text{when } 0 \leq t \leq s \leq L \end{cases} \quad (25.2.20)$$

We therefore have a non-homogeneous Fredholm equation.

It is clear that the set of conditions we impose also affects the form of the equivalent integral equation. An ODE by itself without initial/boundary conditions is not enough.

This is because an ODE by itself can't be solved exactly, that is, we cannot find a particular solution, just a general solution. An integral equation however has an exact solution, with no free parameters, and can't therefore be associated to an ODE alone.

## 25.3 Picard iteration

Consider the Volterra integral equation:

$$y(t) = f(t) + \int_a^t K(s, t)y(s)ds \quad (25.3.1)$$

with  $f$  continuous on  $[a, b]$  and  $K, \partial_x K$  continuous on  $[a, b]^2$ . Our goal will be to define an iterative sequence  $(y_n)$  which improves as  $n$  increases. We should therefore make an initial guess, and insert that into the equation to find a better solution. We will define this sequence, known as a **Picard iteration** as follows:

$$\begin{cases} y_0 = f(t) \\ y_k(t) = f(t) + \int_a^t K(s, t)y_{k-1}(s)ds \end{cases} \quad (25.3.2)$$

Because  $f(t)$  is continuous by hypothesis, so are all  $y_i$  for  $i = 0, 1, 2, \dots$ . We now conjecture that:

### Proposition 23.6 (Picard iteration convergence)

We have that:

$$|u_n(t)| = |y_n(t) - y_{n-1}(t)| \leq M_n \quad (25.3.3)$$

with  $\sum_{n=1}^{\infty} M_n$  converging.

*Proof.* Since  $K, f$  are continuous over  $[a, b]$ , they must be bounded:

$$|K(s, t)| \leq L, |f(t)| \leq M \quad \forall s, t \in [a, b] \quad (25.3.4)$$

<sup>1</sup>Indeed:

$$\int_0^L K(s, t)f(s)ds = \int_0^t K(s, t)f(s)ds + \int_t^L K(s, t)f(s)ds \quad (25.2.16)$$

$$= \int_0^t \frac{s}{L}(L-t)f(s)ds + \int_t^L \frac{t}{L}(L-s)f(s)ds \quad (25.2.17)$$

$$= \int_0^t \frac{s}{L}(L-t)f(s)ds + \int_0^L \frac{t}{L}(L-s)f(s)ds - \int_0^t \frac{t}{L}(L-s)f(s)ds \quad (25.2.18)$$

$$= \int_0^t (s-t)f(s)ds + \int_0^L \frac{t}{L}(L-s)f(s)ds \quad (25.2.19)$$

as required.

We can therefore write:

$$|y_1(t) - y_0(t)| = \left| \int_a^t K(s, t)y_0(s)ds \right| = \int_a^x |K(s, t)||f(s)|ds \leq LM(x - a) \quad (25.3.5)$$

Let us now suppose that for some  $n \geq 2$ :

$$|y_{n-1}(t) - y_{n-2}(t)| \leq L^{n-1}M \frac{(t-a)^{n-1}}{(n-1)!} \quad (25.3.6)$$

We then find:

$$|y_n(s) - y_{n-1}(s)| = \left| \int_a^t K(s, t)(y_{n-1}(s) - y_{n-2}(s))ds \right| \quad (25.3.7)$$

$$\leq \int_a^t |K(s, t)||y_{n-1}(s) - y_{n-2}(s)|ds \quad (25.3.8)$$

$$\leq \int_a^t L^n M \frac{(s-a)^{n-1}}{(n-1)!} ds \quad (25.3.9)$$

$$= \leq L^n M \frac{(t-a)^n}{n!} \quad (25.3.10)$$

as required.

We can therefore define  $M_n$  as:

$$|y_n(t) - y_{n-1}(t)| \leq L^n M \frac{(t-a)^n}{n!} \leq L^n M \frac{(b-a)^n}{n!} \equiv M_n \quad (25.3.11)$$

Consequently:

$$\sum_{n=1}^{\infty} M_n = M(e^{L(b-a)} - 1) \quad (25.3.12)$$

and we then find that  $\sum_{n=1}^{\infty} (y_n - y_{n-1})$  converges uniformly to  $u$  on  $[a, b]$  using the Weierstrass test.  $\blacksquare$

Notice however that this is a telescopic sum equal to  $y - y_0 = u$  which implies that  $y = u + y_0$ .

We can then assert that  $\forall \epsilon > 0, \exists N$  with:

$$|y(x) - y_n(x)| < \epsilon \quad \forall n \geq N \quad (25.3.13)$$

This implies that:

$$|K(s, t)y(s) - K(s, t)y_n(s)| < L\epsilon \quad \forall n \geq N \quad (25.3.14)$$

This is equivalent to saying that our iteration converges to the non-homogeneous Fredholm equation:

$$\int_a^t K(s, t)y_n(s)ds \longrightarrow \int_a^t K(s, t)y(s)ds, \quad \text{as } n \rightarrow \infty \quad (25.3.15)$$

We have therefore shown the existence of a continuous solution, but what about its uniqueness?

Suppose there is another solution  $Y$  so that:

$$|y(t) - Y(t)| \leq P \quad (25.3.16)$$

Let us suppose inductively that:

$$|y(t) - Y(t)| \leq L^{n-1} \frac{(t-a)^{n-1}}{(n-1)!} \quad (25.3.17)$$

Then:

$$|y(t) - Y(t)| = \left| \int_a^t K(s,t)(y(s) - Y(s))ds \right| \quad (25.3.18)$$

$$\leq L^n P \frac{(t-a)^n}{n!} \quad (25.3.19)$$

$$\leq L^n P \frac{(b-a)^n}{n!} \quad (25.3.20)$$

As  $n \rightarrow \infty$  the RHS tends to zero, and we therefore find that  $y = Y$  thus proving the uniqueness.

One can use a very similar process to the Fredholm equation using the iteration:

$$\begin{cases} y_0 = f(t) \\ y_k(t) = f(t) + \lambda \int_a^t K(s,t)y_{k-1}(s)ds \end{cases} \quad (25.3.21)$$

In this case we find that:

$$|y_n(t) - y_{n-1}(t)| \leq |\lambda|^n L^n M(b-a)^n \quad (25.3.22)$$

is uniformly convergent only if  $|\lambda| \leq \frac{1}{L(b-a)}$ . This is the sufficient condition that must be met for a solution to exist.

## 25.4 Existence and uniqueness

Consider the Cauchy problem:

$$y' = f(x, y), \quad y(a) = c \quad (25.4.1)$$

where  $f$  satisfies the following two conditions:

- (i)  $f$  is continuous in a region  $U$  containing  $R = \{(x, y) : |x - a| \leq h, |y - c| \leq k\} \subseteq U$ .
- (ii)  $f$  satisfies the Lipschitz condition:

$$|f(x, y_1) - f(x, y_2)| \leq A|y_1 - y_2|, \quad \forall (x, y_1), (x, y_2) \in U \quad (25.4.2)$$

- (iii) Defining:

$$M = \sup\{|f(x, y)| : (x, y) \in R\} \quad (25.4.3)$$

then we require:

$$Mh \leq k \quad (25.4.4)$$

If these three conditions are satisfied, a very important result, known as the Cauchy-Picard existence and uniqueness theorem, is established.

**Theorem (Cauchy-Picard Existence and Uniqueness theorem)**

If (i), (ii), (iii) are all satisfied then there exists for  $|x - a| \leq h$  a solution to the Cauchy problem:

$$y' = f(x, y), \quad y(a) = c \quad (25.4.5)$$

and this solution is unique in  $U$ .

We will present the proof of a more general result later.

# First Order ODEs

## 26.1 Types of first order ODEs

In this chapter we will consider first order differential equations which can be written in **standard form**:

$$\frac{dy}{dx} = f(x, y) \quad (26.1.1)$$

where  $f(x, y)$  is a function. Since  $f$  can always be written as the ratio of two functions  $M(x, y)$  and  $-N(x, y)$ , a first order ODE in standard form can also be written equivalently in **differential form**:

$$M(x, y)dx + N(x, y)dy = 0 \quad (26.1.2)$$

The main first order ODEs we will be concerned with are:

- (i) Separable ODEs: if  $M(x, y) = M(x)$  and  $N(x, y) = N(y)$ , then we can separate the  $x$  and  $y$  variables and integrate them individually.
- (ii) Exact ODEs: if  $\frac{\partial M}{\partial y} = \frac{\partial N}{\partial x}$  then the differential form of a first order ODE can be written as the total differential of a function.
- (iii) Special inexact ODEs: there are some cases where even though an ODE may not be exact, there are ways to bring it to a separable form by a change of variables.
- (iv) Linear ODEs: if  $f(x, y) = -p(x)y + q(x)$  then we can find an integrating factor that will separate variables.
- (v) Bernoulli equations: linear ODEs but with  $q(x)$  containing a  $x^r$  term for  $r \in \mathbb{R}$ .

Before investigating how these first order ODEs may be solved, we state an important theoretical result on the existence and uniqueness of solutions to first order ODEs.

### Theorem (*Existence and uniqueness*)

Let  $f(t, y)$  and  $\frac{\partial f}{\partial y}$  exist and be continuous on some domain  $\mathcal{D} \subset \mathbb{R}^2$ . Then:  $\forall (t_0, y_0) \in \mathcal{D}, \exists P t$  such that the Cauchy problem:

$$\begin{cases} \dot{y} = f(t, y) \\ y(t_0) = y_0 \end{cases} \quad (26.1.3)$$

has a unique solution in the interval  $I = [t_0 - Pt, t_0 + Pt]$ . If  $y_1(t)$  and  $y_2(t)$  are both solutions on  $I_1$  and  $I_2$  respectively, then  $y_1(t) = y_2(t)$ , that is, a the solution is unique.

## 26.2 Separable Differential Equations

Suppose that  $M(x, y) = M(x)$  and  $N(x, y) = N(y)$ :

$$M(x)dx + N(y)dy = 0 \iff \frac{dy}{dx} = -\frac{M(x)}{N(y)} \quad (26.2.1)$$

These are perhaps the simplest to solve, and luckily crop up quite often in physics, especially in classical mechanics when we relate the rate of change of some quantity with another observable. To solve (26.2.1), we simply integrate:

$$\boxed{\int M(x)dx + \int N(y)dy = 0} \quad (26.2.2)$$

For example, consider an electron of charge  $-e$ , mass  $m$  orbiting a proton of charge  $e$  at a radius  $r$  satisfies the ODE:

$$\frac{dr}{dt} = -\frac{\mu_0 e^4}{12\pi^2 \epsilon_0 m^2 c} \frac{1}{r^2} \quad (26.2.3)$$

This is clearly a separable equation with  $M(t) = -\frac{\mu_0 e^4}{12\pi^2 \epsilon_0 m^2 c}$  and  $N(r) = \frac{1}{r^2}$ . Integrating:

$$\int r^2 dr = -\int \frac{\mu_0 e^4}{12\pi^2 \epsilon_0 m^2 c} dt \implies r(t)^3 = -\frac{\mu_0 e^4}{4\pi^2 \epsilon_0 m^2 c} t + C \quad (26.2.4)$$

where  $C$  is a constant we need to determine through suitable initial conditions. Suppose for example that the electron is initially orbiting the proton at a radius  $r_0$  at  $t = 0$ . Then we find that  $r_0^3 = C$ . If we want to find the time  $T$  it takes for the electron to collapse into the proton, then we must require  $r(T) = 0$  and thus:

$$T = \frac{4\pi^2 \epsilon_0 m^2 c}{\mu_0 e^4} r_0^3 \quad (26.2.5)$$

It turns out that for  $r_i \approx 5 \times 10^{-11}$  m, it takes the electron just  $10^{-11}$  seconds to collapse! Luckily this paradox is not due to an error in our solution but rather a sign that classical physics cannot provide a coherent description of atomic phenomena.

## 26.3 Exact Differential Equations

Consider an ODE in full differentials, with solutions  $\Phi(x, y)$  such that:

$$\begin{cases} \forall (x, y) \in \mathcal{D}, \Phi_x = P(x, y), \Phi_y = Q(x, y) \\ d\Phi = 0 \end{cases} \quad (26.3.1)$$

We can then rewrite, using the chain rule, the equation as:

$$P(x, y)dx + Q(x, y)dy = 0 \quad (26.3.2)$$

### Theorem (Exactness condition)

If  $\frac{\partial P}{\partial y} = \frac{\partial Q}{\partial x}$  through a simply connected domain  $\mathcal{D}$ , then  $Pdx + Qdy = 0$  is an **exact first order ODE**.

*Proof.* Consider the solution  $\Phi(x, y) = C$ , for some constant  $C$ . It follows from the chain rule

that:

$$\frac{\partial \Phi}{\partial x} = P, \frac{\partial \Phi}{\partial y} = Q \implies \frac{\partial^2 \Phi}{\partial x \partial y} = \frac{\partial^2 \Phi}{\partial y \partial x} \implies \frac{\partial P}{\partial y} = \frac{\partial Q}{\partial x} \blacksquare$$

**Strategy (Exact)** To solve:

- Set the equation:  $\Phi_x = P(x, y)$ , and integrate directly with respect to x:

$$\Phi = \int P(x, y) dx + \phi(y) \quad (26.3.3)$$

- Substitute this into  $\Phi_y = Q(x, y)$ :

$$\frac{d}{dy} \left( \int P(x, y) dx + \phi(y) \right) = Q(x, y) \quad (26.3.4)$$

and rearrange to find:

$$\phi(y) = \int \left( Q(x, y) - \frac{d}{dy} \int P(x, y) dx \right) dy \quad (26.3.5)$$

- Set:

$$\int P(x, y) dx + \int \left( Q(x, y) - \frac{d}{dy} \int P(x, y) dx \right) dy = C \quad (26.3.6)$$

for some constant  $C$ .

## 26.4 Inexact Differential Equations

Inexact differential equations are of the form:

$$\begin{cases} P(x, y) dx + Q(x, y) dy = 0 \\ \frac{\partial P}{\partial y} \neq \frac{\partial Q}{\partial x} \end{cases} \quad (26.4.1)$$

unlike exact differential equations. Firstly, consider an ODE of the following form:

$$\frac{dy}{dx} = f(ax + by). \quad (26.4.2)$$

**Strategy (Inexact)**

To solve:

- Apply the change of variables  $z(x) = ax + by(x)$  to find:

$$\frac{dz}{dx} = a + b \frac{dy}{dx} \quad (26.4.3)$$

- Substitute the expression for  $y'$ :

$$\frac{dz}{dx} = a + bf(z) \quad (26.4.4)$$

3. Solve as a separable differential equation:

$$\int \frac{dz}{a + bf(z)} = \int dx \quad (26.4.5)$$

Next, consider the case where (2.3) is homogeneous, that is:

$$\frac{P(\lambda x, \lambda y)}{P(x, y)} = \frac{Q(\lambda x, \lambda y)}{Q(x, y)}, \forall \lambda \neq 0 \quad (26.4.6)$$

Then:

$$\frac{dy}{dx} = -\frac{P(x, y)}{Q(x, y)} = -\frac{P(\lambda x, \lambda y)}{Q(\lambda x, \lambda y)}|_{\lambda=\frac{y}{x}} = -\frac{P(1, \frac{y}{x})}{Q(1, \frac{y}{x})} = f\left(\frac{y}{x}\right) \quad (26.4.7)$$

### Strategy (Homogeneous)

To solve:

1. Apply the change of variables  $u(x)x = y(x)$ :

$$x \frac{du}{dx} + u = \frac{dy}{dx} \quad (26.4.8)$$

2. Substitute into original ODE:

$$f(u) = u + x \frac{du}{dx} \quad (26.4.9)$$

and separate variables:

$$\int \frac{du}{f(u) - u} = \int \frac{dx}{x}. \quad (26.4.10)$$

where  $f(u) \neq u$ .

## 26.5 Integrating Factor Method

Consider the non homogeneous ODE of the form:

$$\dot{y}(t) = a(t)y(t) + f(t) \quad (26.5.1)$$

To solve, we wish to multiply the whole equation by a so called **integrating factor**  $\Lambda$  such that  $\dot{\Lambda} = \Lambda a(t)$ . Then:

$$\Lambda \dot{y}(t) - \overbrace{\Lambda a(t)}^{\dot{\Lambda}} y(t) = \Lambda f(t) \implies \frac{d\Lambda y(t)}{dt} = \Lambda f(t) \quad (26.5.2)$$

$$\implies y(t) = \frac{1}{\Lambda(t)} \left[ C + \int_0^t \Lambda(t') f(t') dt' \right] \quad (26.5.3)$$

To find the integrating factor, we solve the separable equation:

$$\frac{d\Lambda}{dt} = \Lambda(t)a(t) \implies \int \frac{d\Lambda}{\Lambda} = \int a(t) dt \implies \boxed{\Lambda(t) = \exp \left( \int_0^t a(t') dt' \right)} \quad (26.5.4)$$

We have proven the following, very useful result:

**Theorem (Integrating factor)**

Let  $\dot{y}(t) = a(t)y(t) + f(t)$  be a first order linear ODE. Its general solution is then given by:

$$y(t) = \frac{1}{\Lambda(t)} \left[ C + \int_0^t \Lambda(t')f(t')dt' \right], \quad \Lambda(t) = \exp \left( \int_0^t a(t')dt' \right) \quad (26.5.5)$$

## 26.6 Bernoulli Equations

Finally, let us look at the Bernoulli Equations:

$$y' + a(x)y = b(x)y^n \quad (26.6.1)$$

which due to the  $y^n$  term, is non linear.

**Strategy (Bernoulli equation)**

To solve:

1. Divide through by  $y^n$  to find:

$$\frac{dy}{dx}y^{-n} + a(x)y^{1-n} = b(x) \quad (26.6.2)$$

2. Apply a change of variables  $u = y^{1-n}$ :

$$\frac{du}{dx} + (1-n)a(x)u = (1-n)b(x) \quad (26.6.3)$$

3. Use integrating factor method.

## 26.7 Stability and Equilibrium points

**Definition (Equilibrium points)**

An **equilibrium point** of a differential equation is a constant solution  $y' = 0, \forall t \in \mathcal{D}$ . It is:

Stable: if  $y \rightarrow c$  as  $t \rightarrow \infty$ , in other words the deviation decays.

Unstable: if  $y \rightarrow \infty$  as  $t \rightarrow \infty$ , in other words the deviation grows.

**Definition 2.1.** We can linearize differential equations by doing a perturbative analysis. Suppose  $y = a$  is an equilibrium point of  $y' = f(x, y)$ . We then induce an arbitrarily small perturbation  $y = a + \epsilon(t)$ , so that:

$$\frac{d\epsilon}{dt} = \frac{dy}{dt} = f(a, t) + \epsilon \frac{\partial f}{\partial y}(a, t) + O(\epsilon^2) \quad (26.7.1)$$

$$\approx \epsilon \frac{\partial f}{\partial t}(a, t) \quad (26.7.2)$$

If  $\dot{\epsilon} > 0$ , we have unstable equilibrium, if  $\dot{\epsilon} < 0$ , we have stable equilibrium.

# Second Order ODEs

## 27.1 Homogeneous equation

We consider the second order homogeneous differential equation with initial conditions:

$$\begin{cases} a(x)y'' + b(x)y' + c(x)y = 0, & (x \in [a, b]) \\ y(x_0) = y_0 \\ y'(x_0) = y'_0 \end{cases} \quad (27.1.1)$$

where  $a(x) > 0, \forall x \in [a, b]$ .

**Definition (Fundamental matrix and Wronskian)** Suppose  $y_1(x)$  and  $y_2(x)$  are differentiable functions on  $[a, b]$ . We then define their **Fundamental matrix** to be:

$$\Upsilon(y_1, y_2) = \begin{pmatrix} y_1(x) & y_2(x) \\ y'_1(x) & y'_2(x) \end{pmatrix} \quad (27.1.2)$$

and define their **Wronskian** to be the determinant of the fundamental matrix:

$$W(y_1, y_2)(x) = \begin{vmatrix} y_1(x) & y_2(x) \\ y'_1(x) & y'_2(x) \end{vmatrix} = y_1(x)y'_2(x) - y'_1(x)y_2(x) \quad (27.1.3)$$

### Proposition (Linear (in)dependence)

For two non-constant functions  $y_1(x), y_2(x)$  differentiable on  $[a, b]$ :

- a) If the Wronskian  $W(y_1, y_2)(x_0) \neq 0$  for some  $x_0 \in [a, b]$ , the two functions  $y_1(x)$  and  $y_2(x)$  are linearly independent on  $[a, b]$ .
- b) If they are linearly dependent then  $W(y_1, y_2)(x) = 0, \forall x \in [a, b]$ .

*Proof.*

- a. Assume that the Wronskian is non-zero for some  $x_0 \in [a, b]$ , then:

$$y_1(x_0)y'_2(x_0) - y_2(x_0)y'_1(x_0) \neq 0 \implies y_1(x_0) \neq \frac{y'_1(x_0)}{y'_2(x_0)}y_2(x_0) = cy_1(x_0) \quad (27.1.4)$$

where we set  $c = \frac{y'_1(x_0)}{y'_2(x_0)}$ . Therefore  $y_1(x_0)$  and  $y_2(x_0)$  are linearly independent. Using the result b) we see that  $y_1(x)$  and  $y_2(x)$  must be linearly independent for all  $x \in [a, b]$ , since otherwise the Wronskian would vanish identically.

- b. Assume that the solutions are linearly dependent, so that:  $c_1y_1(x) + c_2y_2(x) = 0$  for some  $c_1, c_2 \in \mathbb{R}$  not both zero. Then differentiating with respect to  $x$  one finds that:

$$\begin{cases} c_1y_1(x) + c_2y_2(x) = 0 \\ c_1y'_1(x) + c_2y'_2(x) = 0 \end{cases} \quad (27.1.5)$$

which we may consider as a system of equations in  $c_1, c_2$ . For a non-zero solution  $(c_1, c_2)$  to exist the following determinant must vanish:

$$\begin{vmatrix} y_1(x) & y_2(x) \\ y'_1(x) & y'_2(x) \end{vmatrix} = 0 \implies W(y_1, y_2)(x) = 0, \forall x \in [a, b] \quad (27.1.6)$$

thus proving the desired result. as required. ■

Suppose that by some stroke of luck we have already found two linearly independent solutions,  $y_1(x)$  and  $y_2(x)$ , of (27.1.1). Due to the linearity of (27.1.1), we may use the principle of superposition and state that:

$$y(x) = c_1y_1(x) + c_2y_2(x), \forall c_1, c_2 \in \mathbb{R} \quad (27.1.7)$$

will also be a solution. We now determine the constants:

$$\begin{pmatrix} y(x_0) \\ y'(x_0) \end{pmatrix} = \begin{pmatrix} y_1(x_0) & y_2(x_0) \\ y'_1(x_0) & y'_2(x_0) \end{pmatrix} \cdot \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} = Y \cdot \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} \quad (27.1.8)$$

Then, the coefficients are determined as:

$$\begin{pmatrix} c_1 \\ c_2 \end{pmatrix} = \begin{pmatrix} y_0 \\ y'_0 \end{pmatrix} \cdot Y^{-1} \quad (27.1.9)$$

$$= \frac{1}{W(y_1, y_2)(x_0)} \begin{pmatrix} y'_2(x_0) & -y_2(x_0) \\ -y'_1(x_0) & y_1(x_0) \end{pmatrix} \cdot \begin{pmatrix} y(x_0) \\ y'(x_0) \end{pmatrix} \quad (27.1.10)$$

Note that had we chosen linearly dependent solutions then  $Y$  would not have been invertible, and thus we would not be able to find  $c_1, c_2$  directly through this method.

### Theorem (Abel's Identity for second order ODE)

If  $y_1(x)$  and  $y_2(x)$  are solutions to the homogeneous ODE:

$$y'' + p(x)y' + q(x)y = 0 \quad (27.1.11)$$

then:

$$W(y_1, y_2)(x) = W(y_1, y_2)(x_0) \exp \left[ - \int_{t_0}^x p(s)ds \right] \quad (27.1.12)$$

*Proof.* The derivative of a determinant is given by Jacobi's formula:

$$\frac{d}{dx} \det\{\mathbf{A}(x)\} = \text{tr} \left( \text{adj } \mathbf{A}(x) \frac{d\mathbf{A}(x)}{dx} \right) \quad (27.1.13)$$

For a  $2 \times 2$  matrix  $\mathbf{A}$ :

$$\mathbf{A} = \begin{pmatrix} a(x) & b(x) \\ c(x) & d(x) \end{pmatrix} \implies \text{adj } \mathbf{A} = \begin{pmatrix} d(x) & -b(x) \\ -c(x) & a(x) \end{pmatrix} \quad (27.1.14)$$

so that:

$$\text{adj } A(x) \frac{dA(x)}{dx} = \begin{pmatrix} d(x) & -b(x) \\ -c(x) & a(x) \end{pmatrix} \begin{pmatrix} a'(x) & b'(x) \\ c'(x) & d'(x) \end{pmatrix} \quad (27.1.15)$$

$$= \begin{pmatrix} a'(x)d(x) - b(x)c'(x) & b'(x)d(x) - b(x)d'(x) \\ a(x)c'(x) - a'(x)c(x) & a(x)d'(x) - c(x)b'(x) \end{pmatrix} \quad (27.1.16)$$

and hence:

$$\frac{d}{dx} \det\{A\} = a'(x)d(x) - b(x)c'(x) + a(x)d'(x) - c(x)b'(x) \quad (27.1.17)$$

Therefore, substituting  $a(x) = y_1(x)$ ,  $c(x) = y'_1(x)$ ,  $b(x) = y_2(x)$ ,  $d(x) = y'_2(x)$  we find that:

$$W'(y_1, y_2)(x) = y'_1(x)y'_2(x) - y_2(x)y''_1(x) + y_1(x)y''_2(x) - y'_1(x)y'_2(x) \quad (27.1.18)$$

$$= y_1(x)y''_2(x) - y_2(x)y''_1(x) \quad (27.1.19)$$

We may now substitute  $y''_i(x) + p_i(x)y'_i(x) + q_i(x) = 0$  for  $i = 1, 2$  to find that

$$W'(y_1, y_2)(x) = y_1(x)(-p(x)y'_2(t) - q(x)y_2(x)) - y_2(x)(-p(x)y'_1(x) - q(x)y_1(x)) \quad (27.1.20)$$

$$= -(y_1(x)y'_2(x) - y_2(x)y'_1(x))p(x) \quad (27.1.21)$$

$$= -W(y_1, y_2)(x)p(x) \quad (27.1.22)$$

Separating variables, and integrating from  $x_0$  to  $x$ , one finally finds:

$$W(y_1, y_2)(x) = W(y_1, y_2)(x_0) \exp \left[ - \int_{x_0}^x p(s)ds \right] \quad (27.1.23)$$

as required. ■

■

## 27.2 Non-homogeneous

**Definition 3.4.** The second order non-homogeneous differential equation is:

$$\begin{cases} a(t)y'' + b(t)y' + c(t)y = f(t) \\ y(t_0) = y_0 \\ y'(t_0) = y'_0 \end{cases} \quad (27.2.1)$$

and its *associated homogeneous equation* is (3.1).

### Theorem (Complementary Function and Particular Integral)

The general solution to the non-homogeneous differential equation may be written as the sum:

$$y(t) = y_{CF}(t) + y_{PI}(t) \quad (27.2.2)$$

where  $y_{CF}$  is the complementary function, that is, the solution to the associated homogeneous equation, and  $y_{PI}$  is a particular solution.

*Proof.*

Consider the difference between the general solution and the particular integral:  $y(t) - y_{PI}(t)$ . This must be a solution to the associated homogeneous equation by the superposition principle. In-

deed:

$$a(t)(y''(t) - y_{PI}''(t)) + b(t)(y'(t) - y_{PI}'(t)) + c(t)(y(t) - y_{PI}(t)) \quad (27.2.3)$$

$$= a(t)y''(t) + b(t)y'(t) + c(t)y(t) - (a(t)y_{PI}''(t) + b(t)y_{PI}'(t) + c(t)y_{PI}(t)) \quad (27.2.4)$$

$$= f(t) - f(t) = 0 \quad (27.2.5)$$

as required. Hence, since  $y_1(t)$  and  $y_2(t)$  form a fundamental set of solutions, we may write any solution of the associated homogeneous equation, including  $y(t) - y_{PI}(t)$ , as a linear combination:

$$y(t) - y_{PI}(t) = c_1 y_1(t) + c_2 y_2(t) \implies y(t) = y_{PI}(t) + c_1 y_1(t) + c_2 y_2(t) \quad (27.2.6)$$

as required. ■

It may seem like this theorem is of very little use, since a particular solution is *ipso facto* given by the general solution. However, the following two methods may be used to determine the particular solution:

1. Undetermined Coefficients: guessing and checking, quick but works only in special cases.
2. Variation of parameters: more general, almost always works.

## 27.3 Undetermined Coefficients

This method consists in guessing the form of the particular integral leaving the coefficients indeterminate, and then plug them into the differential equation.

The following table summarizes possible particular integrals for different functions  $f(t)$ :

$g(t)$	$y_{PI}(t)$
$\alpha e^{\beta t}$	$Ae^{\beta t}$
$a \cos(\beta t) + b \sin(\beta t)$	$A \cos(\beta t) + B \sin(\beta t)$
$\sum_{i=0}^n a_i x^i$	$\sum_{i=0}^n A_i x^i$

## 27.4 Variation of Constants

Assume that we have found the complementary solution to the non-homogeneous equation:

$$y_{CF}(t) = c_1 y_1(t) + c_2 y_2(t) \quad (27.4.1)$$

and look for the particular integral of the form:

$$y_{PI}(t) = \psi_1 y_1 + \psi_2 y_2 \quad (27.4.2)$$

such that:

$$\psi'_1 y_1 + \psi'_2 y_2 = 0 \quad (27.4.3)$$

Differentiating:

$$y'_{PI} = \psi_1 y'_1 + \psi_2 y'_2 \quad (27.4.4)$$

$$y''_{PI}(t) = \psi'_1 y'_1 + \psi'_2 y'_2 + \psi_1 y''_1 + \psi_2 y''_2 \quad (27.4.5)$$

Inserting these into the differential equation and simplifying:

$$\psi'_1 y'_1 + \psi'_2 y'_2 = \frac{f(t)}{a(t)} \quad (27.4.6)$$

Let us finally assume that  $a(t) = 1$  (which corresponds to rearranging the equation so that the coefficient of  $y''$  is 1), then:

$$\begin{cases} \psi'_1 y'_1 + \psi'_2 y'_2 = f(t) \\ \psi'_1 y_1 + \psi'_2 y_2 = 0 \end{cases} \quad (27.4.7)$$

$$\Rightarrow \psi'_1 = -\frac{\psi'_2 y_2}{y_1} \Rightarrow -\frac{\psi'_2 y_2}{y_1} y'_1 + \psi'_2 y'_2 = f(t) \quad (27.4.8)$$

$$\Rightarrow \psi'_2 = \frac{y_1 f(t)}{y_1 y'_2 - y_2 y'_1}, \quad \psi'_1 = -\frac{y_2 f(t)}{y_1 y'_2 - y_2 y'_1} \quad (27.4.9)$$

Since  $y_1$  and  $y_2$  are linearly independent, their Wronskian is non-zero, and thus:

$$\psi'_1 = -\frac{y_2 f(t)}{W(y_1, y_2)}, \quad \psi'_2 = \frac{y_1 f(t)}{W(y_1, y_2)}, \quad (27.4.10)$$

which can be integrated directly in (3.4) to get:

$$y_{PI}(t) = -y_1 \int \frac{y_1 f(t)}{y_2 y'_2 - y_2 y'_1} dt + y_2 \int \frac{y_1 f(t)}{W(y_1, y_2)} \quad (27.4.11)$$

## 27.5 Reduction of Order

Finally, let us look at how we may simplify differential equations when one solution to the associated homogeneous equation,  $y_1(t)$  is known.

Then, to solve:

1. We search for solutions of the form:

$$y(t) = \phi(t)y_1(t) \quad (27.5.1)$$

and evaluate its first and second derivatives.

2. We substitute into the original differential equation.
3. Solve the resulting differential equation by substituting  $\psi = \phi'$ .

## 27.6 Euler-Cauchy equations

## 27.7 Intro to Green's functions

Up until now all our methods have relied on producing a general solution and subsequently fitting it to given boundary conditions. The method of Green's functions (which will prove to be powerful for PDEs too) already takes the boundary conditions from the beginning, building a particular solution directly.

**Theorem ()**

# Mechanical Vibrations and Resonance Phenomena

Constant coefficient second-order linear ODEs are of particular interest in the area of mechanical vibrations.

## 28.1 Homogeneous Equation

Consider the homogeneous equation:

$$\ddot{y} + a\dot{y} + by = 0 \quad (28.1.1)$$

and we guess a solution in exponential form:

$$y(t) = Ce^{\lambda t} \quad (28.1.2)$$

Plugging into (4.1) one finds:

$$\lambda^2 + a\lambda + b = 0 \quad (28.1.3)$$

which is called the *auxiliary equation*. The quadratic formula then yields:

$$\lambda_{1,2} = \frac{-a \pm \sqrt{a^2 - 4b}}{2} \quad (28.1.4)$$

which brings us to the following result.

### Proposition 1

The solutions to the second order homogeneous ODE with constant coefficients is:

$$y(t) = C_1 e^{\lambda_1 t} + C_2 e^{\lambda_2 t} \quad (28.1.5)$$

provided that the solutions  $\lambda_1, \lambda_2$  are non-degenerate. In the case where:

$a^2 > 4b$ : we have two real, distinct solutions, and by the superposition principle:

$$y(t) = e^{-at/2} (C_1 e^{t\sqrt{a^2 - 4b}/2} + C_2 e^{-t\sqrt{a^2 - 4b}/2}) \quad (28.1.6)$$

$a^2 < 4b$ : we have two complex, distinct solutions, and by the superposition principle:

$$y(t) = e^{-at/2} (C_1 e^{it\sqrt{4b-a^2}/2} + C_2 e^{-it\sqrt{4b-a^2}/2})$$

Euler's identity then gives:

$$y(t) = e^{-at/2} (A \cos \Omega t + B \sin \Omega t) = \alpha e^{-at/2} \cos(\Omega t + \phi) \quad (28.1.7)$$

$$\text{where } \Omega = \frac{\sqrt{4b-a^2}}{2}.$$

If instead we have two degenerate solutions, so that  $a^2 = 4b$ , then we have only found one solution:

$$y_1(t) = C e^{-at/2}. \quad (28.1.8)$$

We use the method of reduction of variables to find the general solution:

$$y(t) = \psi(t) e^{-at/2} \implies \psi'' + \left(b - \frac{a^2}{4}\right)\psi = 0 \quad (28.1.9)$$

$$\therefore \psi'' = 0 \implies \psi = C_1 t + C_2 \quad (28.1.10)$$

since  $b = \frac{a^2}{4}$ . Finally:

$$y(t) = (C_1 t + C_2) e^{-at/2} \quad (28.1.11)$$

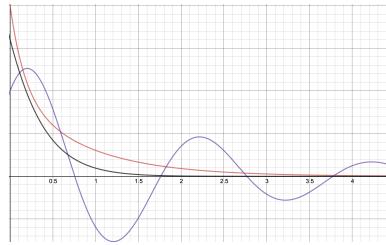
## 28.2 Damped Harmonic Motion

For a damped harmonic system, there will be two forces in action, a restoring force  $F_{res} = -m\omega_0^2 y$  and a damping force  $F_{damp} = -m\gamma\dot{y}$ . Newton's second law yields the second order ODE with constant coefficients:

$$\ddot{y} + \gamma\dot{y} + \omega_0^2 y = 0 \quad (28.2.1)$$

In this case, we may define:

$$\Omega \equiv \frac{\sqrt{\omega_0^2 - \frac{\gamma^2}{4}}}{2} \quad (28.2.2)$$



**Figure 28.1.** Plots of over damped (red), critically damped (black) and under damped (purple) solutions.

We may further impose the initial conditions  $y(0) = y_0, \dot{y}(0) = 0$ . As in the previous section, we consider three special cases:

$\omega_0 > \lambda/2$ , the oscillator is *under damped*. The solution is:

$$y(t) = y_0 e^{-\gamma t/2} \left( \cos \Omega t + \frac{\gamma}{2\Omega} \sin \Omega t \right) \quad (28.2.3)$$

Note that this can be rewritten as:

$$y(t) = A e^{-\gamma t/2} \cos(\Omega t + \phi) \quad (28.2.4)$$

which implies that the period of oscillations is:

$$\tau = \frac{2\pi}{\Omega} \quad (28.2.5)$$

and over a cycle the amplitude is multiplied by:

$$e^{-\gamma\tau/2} = \exp \left( -\frac{2\pi\gamma}{\Omega} \right) \quad (28.2.6)$$

called the *amplitude decay factor*.

$\omega_0 < \lambda/2$ , the oscillator is *over damped*. The solution is:

$$y(t) = \frac{y_0}{2\Omega} e^{-\gamma t/2} \left[ \left( \Omega + \frac{\gamma}{2} \right) e^{\Omega t} + \left( \Omega - \frac{\gamma}{2} \right) e^{-\Omega t} \right] \quad (28.2.7)$$

$\omega_0 = \lambda/2$ , the oscillator is *critically damped*, the solution is:

$$y(t) = y_0 e^{-\gamma t/2} \left( 1 + \frac{\gamma}{2} t \right) \quad (28.2.8)$$

**Definition 4.1** For a lightly damped oscillator in the regime  $\gamma \ll \omega_0$ , with initial stored energy  $E_0$  and energy lost per period of oscillation  $PE_\tau$ , the *quality factor* is defined as:

$$Q = \frac{2\pi E_0}{PE_\tau} \quad (28.2.9)$$

**Proposition 2.** *The quality factor for a damped oscillator is:*

$$Q = \frac{\omega_0}{\gamma} \quad (28.2.10)$$

*Proof.* The energy stored in the oscillator is:

$$E_0 = \frac{1}{2} k y_0^2 = \frac{1}{2} m \omega_0^2 y_0^2 \quad (28.2.11)$$

The under damped solution gives:

$$\begin{cases} y(t) = y_0 e^{-\gamma t/2} \cos(\omega_0 t - \phi) \\ \dot{y}(t) = -y_0 e^{-\gamma t/2} \left( \frac{\gamma}{2} \cos(\omega_0 t - \phi) + \omega_0 \sin(\omega_0 t - \phi) \right) \end{cases} \quad (28.2.12)$$

We then define the energy at  $t$  to be the sum of the kinetic and potential energy:

$$E(t) = \frac{1}{2}m\dot{y}^2 + \frac{1}{2}m\omega_0^2 y^2 \quad (28.2.13)$$

$$\Rightarrow \frac{dE(t)}{dt} = m\ddot{y}\dot{y} + m\omega_0^2 y\dot{y} = m\dot{y} \underbrace{(\ddot{y} + \omega_0^2 y)}_{-\omega^2 y} = -m\gamma\dot{y}^2 \quad (28.2.14)$$

$$\therefore PE_\tau = \int_t^{t+Pt} -m\gamma\dot{y}^2 dt \quad (28.2.15)$$

$$= \int_t^{t+Pt} m\gamma y_0^2 e^{-\gamma t} \left( \frac{\gamma}{2} \cos(\omega_0 t - \phi) + \omega_0 \sin(\omega_0 t - \phi) \right)^2 dt \quad (28.2.16)$$

$$(28.2.17)$$

We can now apply the substitution  $t' = \omega_0 t - \phi$ , so that the limits of integration become 0 and  $2\pi$ :

$$PE_\tau = m\gamma y_0^2 \int_0^{2\pi} e^{-\gamma \frac{t'+\phi}{\omega_0}} \left( \frac{1}{2} \cos t' + \omega_0 \sin t' \right)^2 \frac{1}{\omega_0} dt' \quad (28.2.18)$$

$$= \frac{m\gamma y_0^2}{\omega_0} \frac{\pi}{4} (4\omega_0^2 + \gamma^2) \quad (28.2.19)$$

$$= \pi m\gamma y_0^2 \omega_0 \left( 1 + \frac{\gamma^2}{4\omega_0^2} \right) \quad (28.2.20)$$

$$= \pi m\gamma y_0^2 \omega_0 \quad (28.2.21)$$

hence:

$$Q = 2\pi \frac{\frac{1}{2}m\omega_0^2 y_0^2}{\pi m\gamma y_0^2 \omega_0} = \frac{\omega_0}{\gamma} \quad (28.2.22)$$

as required. ■

## 28.3 Forced Oscillations

Finally, let us consider the damped, forced oscillations equation:

$$\ddot{y} + \gamma\dot{y} + \omega_0^2 y = F \cos \omega t \quad (28.3.1)$$

We have already found the complementary function in the previous section, we must now find a particular solution. This can be done using the method of undetermined coefficients. To do so, we solve the complex version of (4.10):

$$\ddot{z} + \gamma\dot{z} + \omega_0^2 z = F e^{i\omega t} \quad (28.3.2)$$

and use as a trial solution  $z = C e^{i\omega t}$ . Then:

$$C = \frac{F}{\omega_0^2 - \omega^2 + i\omega\gamma} = \frac{F(\omega_0^2 - \omega^2 - i\omega\gamma)}{(\omega_0^2 - \omega^2)^2 + \omega^2\gamma^2} \quad (28.3.3)$$

Next, we define:

$$\cos \phi = \frac{\omega_0^2 - \omega^2}{\sqrt{(\omega_0^2 - \omega^2)^2 + \omega^2\gamma^2}}, \sin \phi = \frac{\omega\gamma}{\sqrt{(\omega_0^2 - \omega^2)^2 + \omega^2\gamma^2}} \quad (28.3.4)$$

so that:

$$C = \frac{Fe^{-i\phi}}{\sqrt{(\omega_0^2 - \omega^2)^2 + \omega^2\gamma^2}} \quad (28.3.5)$$

This finally gives the solution:

$$z(t) = \frac{Fe^{i(\omega_0 t - \phi)}}{\sqrt{(\omega_0^2 - \omega^2)^2 + \omega^2\gamma^2}} \quad (28.3.6)$$

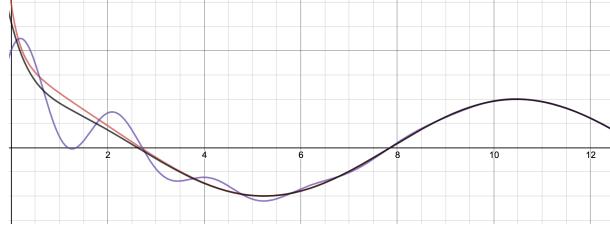
Taking the real part yields:

$$y_{PI}(t) = A \cos(\omega t - \phi), \quad A = \frac{F}{\sqrt{(\omega_0^2 - \omega^2)^2 + \omega^2\gamma^2}} \quad (28.3.7)$$

Finally, the general solution is:

$$y(t) = \underbrace{A \cos(\omega t - \phi)}_{\text{steady state}} + \underbrace{e^{-\gamma t/2} [C_1 \cos \Omega t + C_2 \sin \Omega t]}_{\text{transient}} \quad (28.3.8)$$

where the first term is the steady state solution, and the second term is the transient solution, and quickly decays after  $t \gg \gamma^{-1}$ .



**Figure 28.2.** Forced solutions for over damped, under damped and critically damped oscillators.

## 28.4 Resonance

Let us now consider the scenario in which  $\gamma \ll \omega_0$ , then  $A(\omega)$  has a peak near  $\omega_0$ , that is, when the frequency of the forced oscillations are close to the natural frequency of the oscillator. This phenomenon is known as *resonance*. To find the peak, we set  $A'(\omega) = 0$ :

$$4\omega_{res}(\omega_0 - \omega_{res}^2) + 2\omega_{res}\gamma^2 = 0 \implies \omega_{res} = \sqrt{\omega_0^2 - \frac{\gamma^2}{2}} \approx \omega_0 \quad (28.4.1)$$

The peak amplitude is then:

$$A_{res} \approx A(\omega_0) = \frac{F}{\omega_0\gamma} \quad (28.4.2)$$

# General Linear ODEs

## 29.1 Existence and Uniqueness

We now consider the general theory of linear ODEs.

**Theorem (Existence and uniqueness linear)**

If all  $A_{ij}(t)$  and  $f_i(t)$  are continuous on  $\mathcal{I} = (t_1, t_2)$ , then  $\forall t_0 \in \mathcal{I}, \forall \mathbf{y}_0 \in \mathbb{R}^n$ , the Cauchy problem:

$$\begin{cases} \dot{\mathbf{y}} = \mathbf{A}(t)\mathbf{y} + \mathbf{f}(t) \\ \mathbf{y}(t_0) = \mathbf{y}_0 \end{cases} \quad (29.1.1)$$

has a unique solution in  $\mathcal{I}$

## 29.2 Fundamental set and Wronskians

**Definition (Fundamental set of solutions)** The *fundamental system of solutions* of the homogeneous equation:

$$\dot{\mathbf{y}} = \mathbf{A}(t)\mathbf{y} \quad (29.2.1)$$

is the set of linearly independent solutions  $\{\mathbf{y}_1(t), \mathbf{y}_2(t), \dots, \mathbf{y}_n(t)\}$  to the associated homogeneous differential equation:

$$\dot{\mathbf{y}} = \mathbf{A}(t)\mathbf{y} \quad (29.2.2)$$

**Proposition (Independent solutions)**

The following are true for any Cauchy system of the form (3.1):

- a. If  $\exists t_0 \in \mathcal{I}$  such that  $\{\mathbf{y}_i(t_0)\}$  is linearly independent, then the fundamental set of solutions is linearly independent.
- b. If  $\exists t_0 \in \mathcal{I}$  such that  $\{\mathbf{y}_i(t_0)\}$  is linearly dependent, then the fundamental set of solutions is linearly dependent.

*Proof.*

- a. Suppose that  $\{\mathbf{y}_i(t_0)\}$  is linearly dependent. Then,  $\exists \{C_i\}$  not all equal to zero such that,  $\forall t \in \mathcal{I}, C_i \mathbf{y}_i = \mathbf{0}$ . However, this is not satisfied for  $t = t_0$ , thus we have a contradiction.

- b. Suppose that  $\{\mathbf{y}_i(t_0)\}$  is linearly dependent. Then,  $\forall t \in \mathcal{I}, \exists \{C_i\}$  such that:  $\mathbf{y} = C_i \mathbf{y}_i = \mathbf{0}$ . This implies that  $\mathbf{y}(t_0) = \mathbf{0}$ , and by the uniqueness theorem,  $\mathbf{y}(t) = \mathbf{0}$  identically, and therefore  $C_i \mathbf{y}_i(t) = (0), \forall t$  as required. ■

### Definition (Wronskian)

The fundamental set of solutions can be packed into a matrix, called the **fundamental matrix**:

$$\mathbf{Y}(t) = (\mathbf{y}_1(t) \ \mathbf{y}_2(t) \ \dots, \mathbf{y}_n(t)) \quad (29.2.3)$$

The **Wronskian** is then defined as the determinant of the fundamental matrix:

$$W(t) = \det \mathbf{Y}(t) \quad (29.2.4)$$

### Theorem (Liouville's formula)

To compute the Wronskian, one can use Liouville's formula:

$$W(t) = W(t_0) \exp \left[ \int_{t_0}^t \text{tr}\mathbf{A}(t') dt' \right] \quad (29.2.5)$$

*Proof.* The derivative of the Wronskian is given by differentiating row by row and then summing:

$$\dot{W}(t) = \sum_{i=1}^n \det \mathbf{Y}_i^*(t) \quad (29.2.6)$$

where we define:

$$\mathbf{Y}_i^*(t) = \begin{pmatrix} y_{11} & y_{12} & \dots & y_{1n} \\ \vdots & \vdots & & \\ \dot{y}_{i1} & \dot{y}_{i2} & \dots & \dot{y}_{in} \\ \vdots & \vdots & & \\ y_{n1} & y_{n2} & \dots & y_{nn} \end{pmatrix} \quad (29.2.7)$$

Now note that since  $\dot{\mathbf{Y}} = \mathbf{AY} \iff \dot{y}_{ik} = \sum_j A_{ij} y_{jk}$ , we may write the above as:

$$\mathbf{Y}_i^*(t) = \begin{pmatrix} y_{11} & y_{12} & \dots & y_{1n} \\ \vdots & \vdots & & \\ \sum_j A_{ij} y_{j1} & \sum_j A_{ij} y_{j2} & \dots & \sum_j A_{ij} y_{jn} \\ \vdots & \vdots & & \\ y_{n1} & y_{n2} & \dots & y_{nn} \end{pmatrix} \quad (29.2.8)$$

Now since the determinant is unchanged if we subtract from one row a linear combination of all the others, we can subtract from the  $i$ th row the following linear combination of all other rows:

$$\sum_{j \neq i}^n A_{ij} (y_{j1} \dots y_{jn}) \quad (29.2.9)$$

which will leave only the  $a_{ii}$  coefficients in the matrix:

$$\det(\mathbf{Y}_i^*(t)) = \det \begin{pmatrix} y_{11} & y_{12} & \dots & y_{1n} \\ \vdots & \vdots & & \\ A_{ii}y_{j1} & A_{ii}y_{j2} & \dots & A_{ii}y_{jn} \\ \vdots & \vdots & & \\ y_{n1} & y_{n2} & \dots & y_{nn} \end{pmatrix} = A_{ii} \det \mathbf{Y} \quad (29.2.10)$$

Therefore, summing over all  $i$  we get:

$$\det \dot{\mathbf{W}}(t) = (\text{tr}\mathbf{A}) \det \mathbf{Y}(t) \implies \dot{\mathbf{W}}(t) = (\text{tr}\mathbf{A})\mathbf{Y}(t) \quad (29.2.11)$$

which can be solved to yield the required formula. ■

## 29.3 Homogeneous ODE

Let us now look at how to solve the homogeneous ODE:

$$\begin{cases} \dot{\mathbf{y}}(t) = \mathbf{A}(t)\mathbf{y} \\ \mathbf{y}(t_0) = \mathbf{y}_0 \end{cases} \quad (29.3.1)$$

If we know a fundamental system of solutions  $\mathbf{y}_i(t)$ , then by the principle of superposition, seeing as the elements of this set form a basis for all solutions, one finds that:

$$\mathbf{y}(t) = \mathbf{Y}(t)\mathbf{C} \quad (29.3.2)$$

We impose the condition  $\mathbf{u}(t_0) = \mathbf{y}_0$  to get:

$$\mathbf{C} = \mathbf{Y}^{-1}(t_0)\mathbf{y}_0 \quad (29.3.3)$$

where  $\exists \mathbf{Y}^{-1}(t_0)$  since the columns of the Wronskian are all linearly independent (linearly independent solutions). The solution to the Cauchy problem is thus:

$$\boxed{\mathbf{y}(t) = \mathbf{Y}(t)\mathbf{Y}^{-1}(t_0)\mathbf{y}_0} \quad (29.3.4)$$

## 29.4 Variation of parameters

Once we have solved the homogeneous equation, we may generalise our results to non-homogeneous differential equations. Consider the following solution

$$\mathbf{y}(t) = \mathbf{Y}(t)\mathbf{C} + \mathbf{y}_{PI}(t) \quad (29.4.1)$$

which subject to the initial condition  $\mathbf{y}(t_0) = \mathbf{y}_0$  yields:

$$\mathbf{C} = \mathbf{Y}^{-1}(t_0)[\mathbf{y}_0 - \mathbf{y}_{PI}(t_0)] \quad (29.4.2)$$

So how can we find the particular integral? It suffices to use the method of variation of constants,  $C_i \rightarrow \psi_i(t)$ :

$$\mathbf{y}(t) = \sum_i \mathbf{y}_i(t)\psi_i(t) = \mathbf{Y}(t)\psi \quad (29.4.3)$$

We substitute this into the non homogeneous equation:

$$\dot{\psi}(t) = \dot{Y}\psi + Y\dot{\psi} = AY\psi + f \implies Y\dot{\psi} = f \quad (29.4.4)$$

since  $\dot{Y} = AY$ . Noting that  $\det Y(t) = W(t) \neq 0$  then:

$$\dot{\psi} = Y^{-1}f \implies \psi = C + \int_{t_0}^t Y^{-1}(t')f(t')dt' \quad (29.4.5)$$

so that:

$$y(t) = Y(t) \underbrace{\left[ C + \int_{t_0}^t Y^{-1}(t')f(t')dt' \right]}_{CF} \quad (29.4.6)$$

We substitute back the expression for  $C$  and find:

$$y(t) = Y(t) \left[ Y^{-1}(t_0)y_0 + \int_{t_0}^t Y^{-1}(t)f(t')dt' \right] \quad (29.4.7)$$

## 29.5 Higher Order linear ODEs

Consider the general form of a linear ODE:

$$y^{(n)} + p_{n-1}(t)y^{(n-1)} + \dots + p_1(t)y' + p_0(t)y = f(t) \quad (29.5.1)$$

and assume we have found a fundamental set of solutions  $\{y_i(t)\}$  for the associated homogeneous differential equation. Then, by the principle of superposition, the complementary function is:

$$y_{CF}(t) = \sum_{i=1}^n c_i y_i(t) \quad (29.5.2)$$

We now perform a variation of parameters and write the particular integral as

$$y_{PI}(t) = \sum_{i=1}^n \psi_i y_i \quad (29.5.3)$$

and assuming that  $\sum_{i=1}^n \psi'_i y_i = 0$ , differentiating yields:

$$y'_{PI}(t) = \sum_{i=1}^n \psi_i y'_i \quad (29.5.4)$$

In general, we will set:

$$\begin{cases} y_{PI}^{(k)}(t) = \sum_{i=1}^n \psi_i y_i^{(k)}, & k = 1, \dots, n-1 \\ \sum_{i=1}^n \psi'_i y_i^{(k)} = 0, & k = 0, \dots, n-2 \end{cases} \quad (29.5.5)$$

Finally, we evaluate the  $n$ th derivative as usual without making any further special assumptions:

$$y_{PI}^{(n)}(t) = \sum_{i=1}^n (\psi_i y_i^{(n)} + \psi'_i y_i^{(n-1)}) \quad (29.5.6)$$

We are now ready to substitute everything into (5.5):

$$\sum_{i=1}^n (\psi_i y_i^{(n)} + \psi'_i y_i^{(n-1)}) + p_{n-1}(t) \sum_{i=1}^n \psi_i y_i^{(n-1)} + \dots + p_1(t) \sum_{i=1}^n \psi_i y'_i + p_0(t) \sum_{i=1}^n \psi_i y_i = f(t) \quad (29.5.7)$$

which rearranging gives:

$$\sum_{i=1}^n \left( \psi_i \left[ \sum_{j=0}^n p_j y_i^{(j)} \right] \right) + \sum_{i=1}^n \psi'_i y_i^{(n-1)} = f(t) \quad (29.5.8)$$

Note that since  $\{y_i(t)\}$  are all solutions to the associated homogeneous equation, then  $\sum_{j=0}^n p_j y_i^{(j)} = 0$ , so that:

$$\sum_{i=1}^n \psi'_i y_i^{(n-1)} = f(t) \quad (29.5.9)$$

We therefore have the following system of equations:

$$\begin{cases} \sum_{i=1}^n \psi'_i y_i^{(k)} = 0, k = 0, \dots, n-2 \\ \sum_{i=1}^n \psi'_i y_i^{(n-1)} = f(t) \end{cases} \quad (29.5.10)$$

To solve this system of equations, we will use Cramer's rule. As always, the Wronskian is:

$$W(t) = \begin{vmatrix} y_1 & y_2 & \dots & y_n \\ y'_1 & y'_2 & \dots & y'_n \\ \vdots & \vdots & & \vdots \\ y_1^{(n-1)} & y_2^{(n-1)} & \dots & y_n^{(n-1)} \end{vmatrix} \quad (29.5.11)$$

To use Cramer's rule, we must successively substitute each column of the Wronskian with  $(0 \ 0 \ \dots \ f(t))^T$ . We will thus denote:

$$W_i = \begin{vmatrix} y_1 & \dots & y_{i-1} & 0 & \dots & y_n \\ y'_1 & \dots & y'_{i-1} & 0 & \dots & y'_n \\ \vdots & & \vdots & & & \vdots \\ y_1^{(n-1)} & \dots & y_{i-1}^{(n)} & f(t) & \dots & y_n^{(n-1)} \end{vmatrix} = f(t) \begin{vmatrix} y_1 & \dots & y_{i-1} & 0 & \dots & y_n \\ y'_1 & \dots & y'_{i-1} & 0 & \dots & y'_n \\ \vdots & & \vdots & & & \vdots \\ y_1^{(n-1)} & \dots & y_{i-1}^{(n)} & 1 & \dots & y_n^{(n-1)} \end{vmatrix} \quad (29.5.12)$$

where the  $i$ th column was altered. Cramer's rule finally gives:

$$\psi'_i = \frac{f(t)W_i(t)}{W(t)} \implies u_i = \int \frac{f(t)W_i(t)}{W(t)} dt \quad (29.5.13)$$

which substituting back into (5.6) gives:

$$y_{PI}(t) = \sum_{i=1}^n \left( y_i(t) \int \frac{f(t)W_i(t)}{W(t)} dt \right)$$

(29.5.14)

# Sturm-Liouville theory and Green's functions

## 30.1 Linear differential operators

We will try to study differential equations from a more general point of view. To do so we must introduce the concept of differential operators.

**Definition (Differential operator)** A **differential operator** is a linear operator  $L$  on a function space  $V$ , such that:

$$L[y] = f(y, y', \dots, y^{(n)}, \dots), \quad \forall y \in V \quad (30.1.1)$$

where  $f$  is a linear function in its arguments, so that:

- (i)  $L[y_1 + y_2] = L[y_1] + L[y_2]$
- (ii)  $L[\alpha y_1] = \alpha L[y_1]$

This notation allow us to write differential equations more succinctly. For example, we may write:

$$\frac{d^2y}{dx^2} + x \frac{dy}{dx} + x^2 y = 0 \quad (30.1.2)$$

as

$$L[y] = 0, \quad \text{where } L \equiv \frac{d^2}{dx^2} + x \frac{d}{dx} + x^2 \quad (30.1.3)$$

**Definition (Boundary value problem)** A **boundary value problem** is an equation  $L[y] = f$  defined on an interval  $(a, b)$  together with a boundary constraints  $B[y] = g$ , which can be:

- (i) Dirichlet: if  $y(a) = c$ , that is we define the value of the solution on a boundary
- (ii) Neumann: if  $y'(a) = c$ , that is we define the value of the derivative of the solution on a boundary.
- (iii) Robin: if  $c_1 y(a) + y'(a) = c_2$ , that is a combination of Dirichlet and Neumann conditions.
- (iv) Periodic: if  $y(a) = y(b)$ ,  $y'(a) = y'(b)$  if we require the solution and its derivative to match at two boundaries.

If  $f$  is zero the BVP is said to be homogeneous, and if  $g$  is zero the BVP is said to have homogeneous boundary constraints.

In most cases solving a BVP is a very difficult, if not impossible task. Often one can exploit the linearity of the differential operator to simplify matters. One way is to express the solution as a superposition of a special set of functions and find the coefficients for this superposition by substitution. This procedure is particularly reminiscent of linear algebra, so it is important to look at the properties of infinite dimensional vector spaces (more detailed discussions can be found in the Functional analysis part).

### Definition (*Weight function*)

A **weight function** on  $[a, b]$  is a real, non-negative function on the interval with a finite number of zeros.

### Definition (*Inner product*)

We define an inner product on the function space  $L^2[a, b]$  by:

$$\langle f, g \rangle_w = \int_a^b f^*(x)g(x)\rho(x)dx \quad (30.1.4)$$

where  $\rho(x)$  is a weight function. This defines a **Hilbert space**. By convention  $\langle f, g \rangle_1 \equiv \langle f, g \rangle$ .

Much like on normal vector spaces, we can construct a basis on Hilbert spaces too. Any set of linearly independent functions  $\{u_n(x)\}$  that spans the Hilbert space is a basis. One can use the Gram-Schmidt procedure to then produce an orthonormal basis  $\{\phi_n(x)\}$ . Given any function  $f(x)$ , we can express it as a superposition of this orthonormal basis:

$$f(x) = \sum_{n=0}^{\infty} c_n \phi_n, \quad c_n = \langle \phi_n, f \rangle = \int_a^b \phi_n^*(x)f(x)\rho(x)dx \quad (30.1.5)$$

Returning to linear operators, we can define the adjoint of an operator as usual

### Definition (*Adjoint of a differential operator*)

Let  $L$  be a linear operator on a Hilbert space. We define its adjoint  $L^\dagger$  so that:

$$\int_a^b (L[f])^*(x)g(x)dx = \int_a^b f^*(x)(L^\dagger[g])(x)dx + \text{boundary terms} \quad (30.1.6)$$

If  $L^\dagger = L$  then the operator is **self-adjoint**, and if in addition the boundary terms vanish then it is **hermitian**.

**Example.** Consider for example  $L[y] = y''(x) + p(x)y'(x)$  on  $[0, 1]$ . We have that:

$$\int_a^b f^*(x)(L^\dagger[g])(x)dx = \int_0^1 (f'')^*g dx + \int_0^1 (pf')^*g dx \quad (30.1.7)$$

$$= [f'g]_0^1 - \int_0^1 (f')^*g' dx + [p^*fg] - \int_0^1 (pf)^*g' dx \quad (30.1.8)$$

$$= [f'g - fg' + p^*fg]_0^1 + \int_0^1 f^*(g'' - p^*g') dx \quad (30.1.9)$$

so we see that:

$$L^\dagger[y] = y''(x) - p^*(x)y(x) \quad (30.1.10)$$

If  $p(x)$  is purely imaginary then we see that  $L$  is self-adjoint. Furthermore,  $L$  is hermitian if in addition

$$[f'g - fg' + p^*fg]_0^1 = 0 \quad (30.1.11)$$

◀

If we are also given a set of boundary conditions together with  $L$ , then we can find their adjoint by looking at how one can make the boundary terms in (30.1.6) vanish (we need this or else retaking the adjoint of  $L^\dagger$  would yield several boundary terms +  $L$ ). Indeed for the boundary terms to vanish for any suitable  $f$ , we need both the term at  $x = a$  and  $x = b$  to be zero.

**Example.** Suppose we add the BCs:  $2y(0) - y'(0) = 0$  and  $y'(1) = 0$  to the previous example's linear operator. The boundary condition at  $x = 0$  reads:

$$2f(0)g(0) - f(0)g'(0) + p^*(0)f(0)g(0) = 0 \implies g'(0) = (2 + p^*(0))g(0) \quad (30.1.12)$$

Similarly, the boundary condition at  $x = 1$  reads:

$$-f(1)g'(1) + p^*(1)f(1)g(1) \implies g'(1) = p^*(1)g(1) \quad (30.1.13)$$

These are the adjoint boundary conditions. ◀

## 30.2 Eigenfunctions

Continuing with our analogy between finite and infinite dimensional vector spaces, we now seek to find eigenfunctions and eigenvalues of  $L$ .

### Definition (*Eigenfunctions of differential operators*)

Let  $L$  be a linear operator on a Hilbert space, and suppose there is a function  $\phi$  such that:

$$L[\phi](x) = \lambda\rho(x)\phi(x) \quad (30.2.1)$$

where  $\lambda$  is a non-zero constant. Then  $\phi(x)$  is a **eigenfunction** of  $L$  with eigenvalue  $\lambda$ . We can in addition require the function to satisfy a set of boundary conditions.

**Example.** Consider for example  $L[y] = -y''$  on  $[0, \pi]$  and the boundary conditions  $y(0) = 0, y'(\pi) = 0$ . Let us find the eigenfunctions and eigenvalues of  $L$  with unit weight  $\rho(x) = 1$ . We need to solve

$$-\phi''(x) = \lambda\phi(x) \quad (30.2.2)$$

(i) if  $\lambda = 0$  then the general solution is

$$\phi(x) = c_1x + c_2 \quad (30.2.3)$$

and applying the boundary conditions we see that the only possible eigenfunction is

the trivial one.

- (ii) if  $\lambda > 0$  then the general solution is

$$\phi(x) = c_1 e^{i\omega x} + c_2 e^{-i\omega x}, \omega = \sqrt{\lambda} \quad (30.2.4)$$

and applying the boundary conditions:

$$\begin{cases} c_1 + c_2 = 0 \\ i\omega(c_1 e^{i\omega\pi} - c_2 e^{-i\omega\pi}) = 0 \end{cases} \implies \begin{pmatrix} 1 & 1 \\ e^{i\omega\pi} & -e^{-i\omega\pi} \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad (30.2.5)$$

we see that the only non-trivial solutions arise when:

$$\det \begin{pmatrix} 1 & 1 \\ e^{i\omega\pi} & -e^{-i\omega\pi} \end{pmatrix} = -e^{-i\omega\pi} - e^{i\omega\pi} = 0 \implies \omega = \left(n + \frac{1}{2}\right) \quad (30.2.6)$$

in which case  $c_1 = -c_2$ . Consequently, the eigenfunctions are:

$$\phi_n(x) = A \sin \left(n + \frac{1}{2}\right) x, n = 0, 1, \dots \quad (30.2.7)$$

with eigenvalues:

$$\lambda_n = \left(n + \frac{1}{2}\right)^2, n = 0, 1, \dots \quad (30.2.8)$$

- (iii) if  $\lambda < 0$  then the general solution is

$$\phi(x) = c_1 e^{-\omega x} + c_2 e^{\omega x}, \omega = \sqrt{-\lambda} \quad (30.2.9)$$

and applying the boundary conditions:

$$\begin{cases} c_1 + c_2 = 0 \\ \omega(-c_1 e^{-\omega\pi} + c_2 e^{\omega\pi}) = 0 \end{cases} \implies \begin{pmatrix} 1 & 1 \\ -e^{-\omega\pi} & e^{\omega\pi} \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad (30.2.10)$$

we see that the only non-trivial solutions arise when:

$$\det \begin{pmatrix} 1 & 1 \\ -e^{-\omega\pi} & e^{\omega\pi} \end{pmatrix} = e^{\omega\pi} - e^{-\omega\pi} = 0 \quad (30.2.11)$$

which has no solutions. Thus there are no non-trivial eigenfunctions with negative eigenvalues.

Finally, we should normalize the eigenfunctions we have found for  $\lambda > 0$ :

$$\langle \phi_n, \phi_n \rangle = \int_0^\pi \sin^2 \left(n + \frac{1}{2}\right) x dx = \frac{\pi}{2} \implies \phi_n(x) \equiv \frac{2}{\pi} \sin \left(n + \frac{1}{2}\right) x \quad (30.2.12)$$

which of course does not affect the corresponding eigenvalues. ◀

There are three properties regarding the eigenfunctions of Hermitian differential operator that are particularly important.

**Theorem (Spectral properties of Hermitian differential operators)**

Let  $L$  be a hermitian differential operator over an interval  $[a, b]$  coupled with a set of BCs, and let  $L_w^2[a, b]$  be the space of square-integrable functions on  $[a, b]$  with weight  $w(x)$ , satisfying the relevant BCs. Then:

- (i) **Real-valued:**  $L$  has real eigenvalues,
- (ii) **Orthonormality:** for a given weight function, eigenfunctions corresponding to distinct eigenvalues are orthogonal.
- (iii) **Completeness:** the normalized eigenfunctions form an orthonormal basis of the Hilbert space  $L_w^2[a, b]$ .

*Proof.* We do not prove (iii) as it is rather involved, but it is a fundamental result that is important to remember.

- (i) let  $\phi$  be an eigenfunction of  $L$  with eigenvalue  $\lambda$  and weight  $\rho$ :

$$L[\phi](x) = \lambda\phi(x)\rho(x) \quad (30.2.13)$$

Then we get that:

$$\int_a^b \phi^*(x)L[\phi](x)dx = \lambda \int_a^b \phi^*(x)\phi(x)\rho(x)dx \quad (30.2.14)$$

$$\Rightarrow \int_a^b \phi(x)(L[\phi](x))^*dx = \lambda^* \int_a^b \phi^*(x)\phi(x)\rho(x)dx \quad (30.2.15)$$

$$(30.2.16)$$

Since  $L$  is hermitian, we may write that:

$$\int_a^b \phi(x)(L[\phi](x))^*dx = \int -a^b L^\dagger[\phi](x)\phi^*(x)dx = \int -a^b L[\phi](x)\phi^*(x)dx \quad (30.2.17)$$

implying that:

$$(\lambda - \lambda^*) \int_a^b \phi^*(x)\phi(x)\rho(x)dx = 0 \quad (30.2.18)$$

Since the inner product is positive semi-definite, and  $\phi$  is non-trivial, (30.2.18) is only possible if  $\lambda = \lambda^*$ , that is if  $\lambda$  is real.

- (ii) We repeat the calculation above, but with two different eigenfunctions  $\phi_i(x), \phi_j(x)$  with distinct eigenvalues  $\lambda_i, \lambda_j$

$$L[\phi_i](x) = \lambda_i\phi_i(x)\rho(x) \quad (30.2.19)$$

$$L[\phi_j](x) = \lambda_j\phi_j(x)\rho(x) \quad (30.2.20)$$

Then we get that:

$$\int_a^b \phi_j^*(x) L[\phi_i](x) dx = \lambda_j \int_a^b \phi_j^*(x) \phi_i(x) \rho(x) dx \quad (30.2.21)$$

$$\int_a^b \phi_i(x) (L[\phi_j](x))^* dx = \lambda_i \int_a^b \phi_j^*(x) \phi_i(x) \rho(x) dx \quad (30.2.22)$$

(30.2.23)

Since  $L$  is hermitian, we may write that:

$$\int_a^b \phi_i(x) (L[\phi_j](x))^* dx = \int_a^b \phi_j^*(x) L^\dagger[\phi_i](x) dx = \int_a^b \phi_j^*(x) L[\phi_i](x) dx \quad (30.2.24)$$

implying that:

$$(\lambda_j - \lambda_i) \int_a^b \phi_j^*(x) \phi_i(x) \rho(x) dx = 0 \implies \langle \phi_i, \phi_j \rangle_\rho = 0 \quad (30.2.25)$$

as desired. ■

### 30.3 Sturm-Liouville problems

We now examine a special type of differential operators that props up very frequently in mathematical physics.

**Definition (Sturm-Liouville operator)** A **Sturm-Liouville operator** is a differential operator of the form:

$$L[y] = \frac{d}{dx} \left( p(x) \frac{dy}{dx} \right) - q(x)y \quad (30.3.1)$$

**Theorem (Sturm-Liouville is self-adjoint)**

All Sturm-Liouville operators are self-adjoint, and can be made hermitian on an interval  $[a, b]$  by requiring:

$$\left[ p(x) \left( \frac{df^*}{dx} g(x) - \frac{dg}{dx} f(x) \right) \right]_a^b = 0 \quad (30.3.2)$$

*Proof.* We find that:

$$\int_a^b (L[f])^*(x) g(x) dx = \int_a^b \left[ \frac{d}{dx} \left( p(x) \frac{df}{dx} \right) - q(x) f(x) \right]^* g(x) dx \quad (30.3.3)$$

$$= \int_a^b \left( p(x) \frac{df^*}{dx} \frac{dg}{dx} - q(x) f^*(x) g(x) dx \right) \quad (30.3.4)$$

$$- \left[ p(x) \frac{df^*(x)}{dx} g(x) \right]_a^b \quad (30.3.5)$$

and similarly

$$\int_a^b f^*(x)(L[g])(x)dx = \int_a^b \left[ \frac{d}{dx} \left( p(x) \frac{dg}{dx} \right) - q(x)g(x) \right] f^*(x)dx \quad (30.3.6)$$

$$= \int_a^b \left( p(x) \frac{dg}{dx} \frac{df^*}{dx} - q(x)f^*(x)g(x)dx \right) \quad (30.3.7)$$

$$- \left[ p(x) \frac{dg}{dx} f(x) \right]_a^b \quad (30.3.8)$$

Consequently we see that:

$$\int_a^b f^*(x)(L[g])(x)dx = \int_a^b (L[f])^*(x)g(x)dx + \left[ p(x) \left( \frac{df^*}{dx} g(x) - \frac{dg}{dx} f(x) \right) \right]_a^b \quad (30.3.9)$$

so  $L$  is indeed self-adjoint. Furthermore, if we impose the boundary condition:

$$\left[ p(x) \left( \frac{df^*}{dx} g(x) - \frac{dg}{dx} f(x) \right) \right]_a^b = 0 \quad (30.3.10)$$

then we do indeed find that  $L$  is hermitian. ■

It may seem unnatural to look at this special form of second order ODEs. However, much like how we can find an integrating factor to write  $y' + p(x)y = q$  into  $(p(x)\lambda(x))' = q(x)$ , we can write any second order linear ODE into a Sturm-Liouville ODE. To see why, consider:

$$p(x)y''(x) + r(x)y' + q(x)y = 0 \quad (30.3.11)$$

We multiply by  $\eta$  so that:

$$\eta(x)p(x)y''(x) + \eta(x)r(x)y' + \eta(x)q(x)y = \frac{d}{dx}(\eta(x)p(x)y'(x)) - \eta(x)q(x)y \quad (30.3.12)$$

implying

$$\eta py'' + \eta ry' = \eta'py' + \eta[py'' + p'y'] \quad (30.3.13)$$

$$\implies \eta'py' + \eta(p'y' - ry') = 0 \quad (30.3.14)$$

$$\implies \frac{\eta'}{\eta} = \frac{r - p'}{p} \implies \eta(x) = \exp \left( \int^x \frac{r(s) - p'(s)}{p(s)} ds \right) \quad (30.3.15)$$

is the required integrating factor. We can simplify it to

$$\eta(x) = \frac{1}{p(x)} \exp \left( \int^x \frac{r(s)}{p(s)} ds \right) \quad (30.3.16)$$

Since  $\eta(x)L$  is a Sturm-Liouville operator, its eigenfunctions are orthogonal with respect to the unit weight function, implying that the eigenfunctions of  $L$  are orthogonal with respect to the weight function  $w(x) = \eta(x)$ .

**Proposition (Sturm-Liouville integrating factor)**

The second order ODE

$$p(x)y''(x) + r(x)y'(x) + q(x)y(x) = 0 \quad (30.3.17)$$

may be turned into a Sturm-Liouville problem by multiplying by the following integrating factor:

$$\eta(x) = \frac{1}{p(x)} \exp\left(\int_a^x \frac{r(s)}{p(s)} ds\right) \quad (30.3.18)$$

This is equivalent to setting the weight function to be equal to the integrating factor  $w(x) = \eta(x)$ .

Now suppose we have an inhomogeneous Sturm-Liouville problem:

$$L[y](x) = \rho(x)F(x) \equiv f(x) \quad (30.3.19)$$

Since  $L$  is hermitian, its eigenfunctions will form an orthonormal basis  $\{\phi_k(x)\}$  of the Hilbert space. Thus the solution  $y(x)$  to (30.3.19) and  $F(x)$  may be expanded as:

$$y(x) = \sum_n y_n \phi_n(x), \quad y_n = \int_a^b \phi_n^*(x) y(x) dx \quad (30.3.20)$$

and

$$F(x) = \sum_n F_n \phi_n(x), \quad F_n = \int_a^b \phi_n^*(x) F(x) dx \quad (30.3.21)$$

Substituting this into (30.3.19) we get that:

$$L[y](x) = \sum_n y_n L[\phi_n](x) = \sum_n \lambda_n y_n \phi_n(x) = \sum_n F_n \phi_n(x) \rho(x) \quad (30.3.22)$$

$$\Rightarrow y(x) = \sum_n \frac{F_n}{\lambda_n} \phi_n(x) \quad (30.3.23)$$

where in the last step we used the linear independence of  $\{\phi_n(x)\}$ . We can rewrite the above as:

$$y(x) = \sum_n \frac{\langle \phi_n, F \rangle_\rho}{\lambda_n} \phi_n(x) = \sum_n \left( \frac{\phi_n(x)}{\lambda_n} \int_a^b \phi_n^*(x') \underbrace{F(x') \rho(x')}_{f(x')} dx' \right) \quad (30.3.24)$$

This motivates the following definition:

### Definition (Green's function)

Given a Sturm-Liouville operator  $L$ , we define its Green's function to be:

$$G(x, x') = \sum_n \frac{1}{\lambda_n} \phi_n^*(x) \phi_n(x) \quad (30.3.25)$$

so that the solution to  $L[y](x) = f(x)$  becomes an integral:

$$y(x) = \int_a^b G(x, x') f(x') dx' \quad (30.3.26)$$

**Proposition (Parseval's identity)**

Let  $\{\phi_n\}$  be an orthonormal basis for a Hilbert space with weight function  $\rho$ , and let  $f$  be a function with  $\langle \phi_n, f \rangle_\rho = f_n$ . Then:

$$\langle f, f \rangle_\rho = \sum_n |f_n|^2 \quad (30.3.27)$$

*Proof.* ■

## 30.4 Green's functions for BVPs

### Green's functions for homogeneous BCs

There is another approach to define Green's functions, a distributional approach which makes use of the delta function.

**Definition (Green's function)**

Let  $L$  be a differential operator. Then we define a Green's function  $G(x, x')$  to satisfy:

$$L[G](x, x') = \delta(x - x') \quad (30.4.1)$$

We would like to verify that this definition is equivalent to the linear algebra approach. Indeed let

$$G(x, x') = \sum_n \frac{1}{\lambda_n} Y_n(x) Y_n^*(x') \quad (30.4.2)$$

then operating  $L$  (with variable  $x$ , not  $x'$ ) on both sides

$$L[G](x, x') = \sum_n Y_n(x) Y_n^*(x') \quad (30.4.3)$$

$$\implies \int_a^b L[G](x, x') Y_m(x') w(x) dx = \sum_n Y_n(x) \langle Y_n, Y_m \rangle_w \quad (30.4.4)$$

$$\implies Y_m(x) = \int_a^b L[G](x, x') Y_m(x') w(x') dx' \quad (30.4.5)$$

implying that  $\delta(x - x')$  as desired.

The utility of Green's functions is best expressed in the following theorem.

**Theorem (Green's functions)**

Let  $G(x, x')$  be a Green's function of the differential operator  $L$  satisfying homogeneous boundary conditions  $G(a, x') = G(b, x') = 0$ . Then the particular solution of  $L[y] = f$  satisfying  $y(a) = y(b) = 0$  is given by

$$y(x) = \int_a^b G(x, x') f(x') dx' \quad (30.4.6)$$

*Proof.* Indeed note that

$$L[y] = \int_a^b L[G](x, x') f(x') dx' = \int_a^b \delta(x - x') f(x') dx = f(x) \quad (30.4.7)$$

and  $y(a) = y(b) = 0$  since  $G(a, x') = G(b, x') = 0$ , as desired. ■

We now derive a general procedure to calculate Green's functions with homogeneous boundary conditions. We consider the case where  $L$  is a second order operator, and begin by noting that  $L[G](x, x') = 0$  for  $x > x'$  or  $x < x'$  so in these two regions,  $G(x, x')$  will simply be equal to the complementary function of  $L$ . Therefore suppose we have found a fundamental set of solutions  $\{y_1, y_2\}$  which also satisfy one of the BCs i.e.  $y_1(a) = 0$  and  $y_2(b) = 0$ .

It follows that the Green's function to the left of  $x'$  is proportional to  $y_1$ . and to the right of  $x'$  it is proportional to  $y_2$ , in order to satisfy the boundary conditions:

$$G(x, x') = \begin{cases} A(x') y_1(x), & x \in [a, x'] \\ B(x') y_2(x), & x \in (x', b] \end{cases} \quad (30.4.8)$$

but what should happen at  $x'$ ? Suppose  $G$  were discontinuous at  $x = x'$ , for simplicity we assume a step discontinuity. Then

$$\frac{\partial G}{\partial x} = \delta(x - x') \implies \frac{\partial^2 G}{\partial x^2} = \delta'(x - x') \quad (30.4.9)$$

which can't be the case since  $L[G](x, x')$  does not contain derivatives of  $\delta$ . Therefore the Green's function must be continuous on  $[a, b]$ , especially at  $x = x'$ :

$$A(x') y_1(x') = B(x') y_2(x') \quad (30.4.10)$$

To further investigate the behaviour at  $x'$  let's integrate  $L[G]$  on  $(x' - \epsilon, x' + \epsilon)$ :

$$\int_{x'-\epsilon}^{x'+\epsilon} \left( \alpha(x) \frac{\partial^2 G}{\partial x^2} + \beta(x) \frac{\partial G}{\partial x} + \gamma(x) G \right) dx = 1 \quad (30.4.11)$$

Since  $G$  is continuous only the second derivative can contribute to the integral as  $\epsilon \rightarrow 0$  so

$$\lim_{\epsilon \rightarrow 0} \int_{x'-\epsilon}^{x'+\epsilon} \alpha(x) \frac{\partial^2 G}{\partial x^2} dx = \alpha(x') \left( \frac{\partial G}{\partial x} \Big|_{x'+} - \frac{\partial G}{\partial x} \Big|_{x'-} \right) = 1 \quad (30.4.12)$$

implying that

$$A(x') y'_1(x') - B(x') y'_2(x') = \frac{1}{\alpha(x')} \quad (30.4.13)$$

To summarize, we found the gluing conditions for  $G(x < x')$  and  $G(x > x')$ :

which are equivalent to

$$\begin{cases} A(x') y_1(x') = B(x') y_2(x') \\ A(x') y'_1(x') - B(x') y'_2(x') = \frac{1}{\alpha(x')} \end{cases} \quad (30.4.15)$$

These can be solved to give

$$A(x') = \frac{y_2(x')}{\alpha(x')W(x')}, \quad B(x') = \frac{y_1(x')}{\alpha(x')W(x')} \quad (30.4.16)$$

where  $W(x')$  is the Wronskian. Finally we have that

$$G(x, x') = \begin{cases} \frac{y_1(x)y_2(x')}{\alpha(x')W(x')}, & x \in [a, x'] \\ \frac{y_1(x')y_2(x)}{\alpha(x')W(x')}, & x \in [x', b] \end{cases} \quad (30.4.17)$$

### Example: Tension in a string

Suppose we have a string of density  $\mu$  that is fixed at its ends  $x = 0, L$ . Letting  $y(x, t)$  represent the vertical displacement of the string at  $(x, t)$ , then one can derive (see Analytical mechanics lecture notes):

$$T \frac{\partial^2 y}{\partial x^2} - \mu g = \mu \frac{\partial^2 y}{\partial t^2}, \quad x \in [0, L], \quad y(0) = y(L) = 0 \quad (30.4.18)$$

where  $T$  is the tension in the string. In the case of a stationary string we have that its profile will satisfy the following ODE:

$$\frac{d^2 y}{dx^2} = \frac{\mu(x)g}{T} \quad (30.4.19)$$

This equation can be solved using normal methods of for well-behaved  $\mu(x)$ . However, an interesting discussion arises from letting

$$\mu = m\delta(x - x') \quad (30.4.20)$$

equivalent to having a massless string and adding a point mass  $m$  at  $x'$ . The result is that we get a homogeneous Sturm-Liouville problem

$$\frac{d^2 y}{dx^2} \frac{mg}{T} \delta(x - x'), \quad y(0) = y(L) = 0 \quad (30.4.21)$$

The procedure from the preceding section can thus be put to use, the fundamental set of solutions can be chosen to be  $\{x, 1 - \frac{x}{L}\}$ . Then the Green's function can be decomposed as

$$G(x, x') = \begin{cases} A(x')x, & 0 \leq x < x' \\ B(x')\left(1 - \frac{x}{L}\right), & x' < x \leq L \end{cases} \quad (30.4.22)$$

The condition of continuity at  $x'$  yields

$$A(x')x' = B(x')\left(1 - \frac{x'}{L}\right) \implies B(x') = \frac{Lx'}{L - x'}A(x') \quad (30.4.23)$$

while the jump condition for the first derivative yields

$$A(x') + 1 = -\frac{1}{L}B(x') \implies A(x') + 1 = -\frac{x'}{L - x'}A(x') \quad (30.4.24)$$

Therefore we find that

$$A(x') = \frac{x'}{L}, \quad B(x') = -x' \quad (30.4.25)$$

giving the following Green's function

$$G(x, x') = \begin{cases} -\left(1 - \frac{x'}{L}\right)x, & 0 \leq x \leq x' \\ -\left(1 - \frac{x}{L}\right)x', & x' \leq x \leq L \end{cases} \quad (30.4.26)$$

Physically, this means that the string will form two straight lines, a negative slope on one side and a positive slope to the right.

We can therefore interpret the Green's function as the system's response at  $x$  to a unit impulse at  $x'$ . Moreover, the beauty of Green's functions is that it gives us the particular integral  $y(x)$  for any  $\mu(x)$ . Indeed the solution to (30.4.19) is

$$y(x) = \frac{g}{T} \left[ \frac{x-L}{L} \int_0^x x' \mu(x') dx' + x \int_x^L \frac{x'-L}{L} \mu(x') dx' \right] \quad (30.4.27)$$

Physically, we are building up our solution by decomposing  $\mu(x)$  into several impulses, calculating the system's response, and summing all the contributions to give the displacement at  $x$ .

### Green's functions for inhomogeneous BCs

Let us now turn to inhomogeneous boundary conditions. Suppose we have found a particular solution  $y_P$  of  $L[y] = 0$ . If we let  $G$  be the Green's functions with homogeneous boundary conditions then

$$y(x) = y_P(x) + \int_a^b G(x, x') f(x') dx' \quad (30.4.28)$$

is the particular solution with inhomogeneous BCs.

## 30.5 Green's functions for IVPs

Green's functions can be used to solve IVPs, such as

$$L[y] = f(t), \quad y(t_0) = y'(t_0) = 0 \quad (30.5.1)$$

Then the Green's function  $G(t, t')$  will solve the following

$$L[y] = \delta(t - t'), \quad G(t_0, t') = G'(t_0, t') = 0 \quad (30.5.2)$$

Suppose we have found a fundamental set of solutions  $\{y_1, y_2\}$ . For  $t_0 \leq t < t'$ , we find that

$$G^-(t, t') = A(t') y_1(t) + B(t') y_2(t), \quad t_0 \leq t < t' \quad (30.5.3)$$

and apply the initial value conditions

$$\begin{pmatrix} y_1(t_0) & y_2(t_0) \\ y_1'(t_0) & y_2'(t_0) \end{pmatrix} \begin{pmatrix} A \\ B \end{pmatrix} = \mathbf{0} \quad (30.5.4)$$

but since the Wronskian is non-zero,  $A = B = 0$  so that the Green's function  $G(t, t') = 0$  for  $t < t'$ . Physically, this means that the system cannot respond in advance at  $t$  to an impulse occurring at a later time  $t'$ , it implies causality.

For  $t_0 \leq t' < t$  we construct

$$G^+(t, t') = C(t')y_1(t) + D(t')y_2(t), \quad t > t' \quad (30.5.5)$$

and apply the Green's function boundary conditions at  $t = t'$ :

$$\begin{cases} C(t')y_1(t') + D(t')y_2(t') = 0 \\ C(t')y'_1(t) + D(t')y'_2(t') = \frac{1}{\alpha(t')} \end{cases} \quad (30.5.6)$$

Note that we do not have to apply the initial conditions since  $t > t_0$ . Solving this system will give  $G(t, t')$

$$G(t, t') = \begin{cases} 0, & t_0 \leq t < t' \\ G^+(t, t'), & t_0 < t' \leq t \end{cases} \quad (30.5.7)$$

**Example.** Consider the equation for a driven harmonic oscillator

$$\frac{d^2y}{dt^2} + y = f(t), \quad y(0) = y'(0) = 0 \quad (30.5.8)$$

The fundamental set of solutions is  $\{\sin t, \cos t\}$ . We need only consider the case where  $t > t'$ , where we need to solve

$$\begin{cases} C(t') \sin(t') + D(t') \cos(t') = 0 \\ C(t') \cos(t) - D(t') \sin(t') = 1 \end{cases} \quad (30.5.9)$$

We find that

$$C = -D \cot(t') \implies D = -\sin(t'), \quad C = \cos(t') \quad (30.5.10)$$

giving the following Green's function

$$G(t, t') = \begin{cases} 0, & t < t' \\ \sin(t - t'), & t > t' \end{cases} \quad (30.5.11)$$

and the following particular integral

$$y(t) = \int_0^t \sin(t - t')f(t') dt' \quad (30.5.12)$$



# Linear systems of ODEs

In chapter 5, we learned how to solve the Cauchy problem for any system of linear ODEs, provided we knew the fundamental system of solutions. We are now going to see, for the special case of constant coefficients, how to find this fundamental set.

## 31.1 Non-degenerate Eigenvalues

Consider the system of ODEs:

$$\dot{\mathbf{y}} = \mathbf{A} \cdot \mathbf{y} + \mathbf{f}(t) \quad (31.1.1)$$

To find the eigenvalues of the matrix  $\mathbf{A}$ , it suffices to solve the *eigenvalue equation*:

$$\det(\mathbf{A} - \lambda I) = 0 \quad (31.1.2)$$

and then find the corresponding eigenvectors using the *eigenvector equation*:

$$(\mathbf{A} - \lambda I) \cdot \mathbf{v}_i = 0. \quad (31.1.3)$$

### Proposition 3.

If a matrix has non-degenerate eigenvalues  $\lambda_i$ , then the corresponding eigenvectors  $\mathbf{v}_i$  form a basis.

*Proof.* Consider (6.3), if  $\{\mathbf{v}_i\}$  are linearly dependent, then it must be possible to write one as:

$$\mathbf{v}_n = \sum_j \alpha_j \mathbf{v}_j \quad (31.1.4)$$

Then, on one hand we find:

$$\mathbf{A} \cdot \mathbf{v}_n = \lambda_n \mathbf{v}_n = \lambda_n \sum_j \alpha_j \mathbf{v}_j \quad (31.1.5)$$

and on the other hand:

$$\mathbf{A} \cdot \mathbf{v}_n = \sum_j \alpha_j \mathbf{A} \cdot \mathbf{v}_j = \sum_j \lambda_j \alpha_j \mathbf{v}_j. \quad (31.1.6)$$

Equating the two gives:

$$\sum_j \alpha_j \mathbf{v}_j (\lambda_n - \lambda_j) = 0 \implies \alpha_j = 0 \quad (31.1.7)$$

since by assumption, the eigenvectors are independent and the eigenvalues are non degenerate.

■

Following the algebraic approach of the Wronskians, we pack the eigenvectors and eigenvalues into matrices as:

$$\begin{aligned} R &\equiv (\mathbf{v}_1 \dots \mathbf{v}_n), L \equiv \text{diag}\{\lambda_i\} \\ \implies A \cdot R &= R \cdot L \\ \implies A &= R \cdot L \cdot R^{-1} \end{aligned}$$

We may substitute this into the associated homogeneous equation to (6.1) and find:

$$\begin{aligned} \mathbf{y} &= R \cdot L \cdot R^{-1} \cdot \mathbf{y} \\ \implies \frac{d}{dt}(R^{-1} \cdot \mathbf{y}) &= L \cdot R^{-1} \cdot \mathbf{y} \end{aligned}$$

and hence we get:

$$\dot{\zeta} = L \cdot \zeta \implies \zeta_i(t) = C_i e^{\lambda_i t} \quad (31.1.8)$$

We transform back into the original coordinates so that:

$$\mathbf{y} = \sum_i C_i \mathbf{v}_i e^{\lambda_i t} \quad (31.1.9)$$

For simplicity, define the following matrix:

$$E(t) \equiv \text{diag} e^{\lambda_i t} \quad (31.1.10)$$

which recasts (6.4) as:

$$\mathbf{y} = \underbrace{R \cdot E(t) \cdot C}_{Y(t)} \quad (31.1.11)$$

Finally, let us set the Cauchy condition  $\mathbf{y}(0) = \mathbf{y}_0$

$$\mathbf{y}_0 = R \cdot C \implies C = R^{-1} \cdot \mathbf{y}_0 \quad (31.1.12)$$

and hence:

$$\mathbf{y} = R \cdot E(t) \cdot R^{-1} \cdot \mathbf{y}_0 \quad (31.1.13)$$

## 31.2 Matrix exponentiation

Consider now the more general Cauchy problem (6.1), which can be solved, as always, via the method of variation of parameters:

$$\mathbf{y}(t) = Y(t) \left[ Y^{-1}(t_0) \cdot \mathbf{y}_0 + \int_{t_0}^t Y^{-1}(t') \cdot \mathbf{f}(t') dt' \right] \quad (31.2.1)$$

Here:

$$\begin{aligned} Y^{-1}(t) &= E^{-1}(t) \cdot R^{-1} \implies Y(t) \cdot Y^{-1}(t') = R \cdot E(t) \cdot E^{-1}(t) \cdot R^{-1} = R \cdot E(t - t') \cdot R^{-1} \\ Y^{-1}(0) &= R^{-1} \end{aligned}$$

so that:

$$\mathbf{y}(t) = Y(t_0) \left[ Y^{-1}(t) \cdot \mathbf{y}_0 + \int_{t_0}^t R \cdot E(t - t') \cdot R^{-1} \cdot \mathbf{f}(t') dt' \right] \quad (31.2.2)$$

We may now compare this with the one-dimensional solution using the integrating factor method, we realize that a new expression for matrix exponentiation can be found:

$$e^{\mathbf{A}t} = \mathbf{R} \cdot \mathbf{E}(t) \cdot \mathbf{R}^{-1} \quad (31.2.3)$$

### 31.3 Higher Order Linear Constant Coefficient Equations

Consider the  $n$ th order linear ODE with constant coefficients:

$$y^{(n)} + a_{n-1}y^{(n-1)} + \dots + a_1\dot{y} + a_0y = f(t) \quad (31.3.1)$$

which can be reduced into a first order linear system of ODEs:

$$\begin{cases} \dot{p}_{n-1} = -a_{n-1}p_{n-1} - \dots - a_1p_1 - a_0y + f(t) \\ \dot{p}_{n-2} = p_{n-1} \\ \dots \\ \dot{p}_1 = p_2\dot{y} = p_1 \end{cases} \quad (31.3.2)$$

or in matrix form:

$$\frac{d}{dt} \begin{pmatrix} p_{n-1} \\ \vdots \\ p_1 \\ y \end{pmatrix} = \begin{pmatrix} -a_{n-1} & -a_{n-2} & \dots & -a_1 & -a_0 \\ 1 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} p_{n-1} \\ \vdots \\ p_1 \\ y \end{pmatrix} + \begin{pmatrix} f \\ 0 \\ \vdots \\ 0 \end{pmatrix} \quad (31.3.3)$$

It can then be shown that the eigenvalue equation and auxiliary equation are equivalent:

$$\det(\mathbf{A} - \lambda\mathbf{I}) = 0 \iff \lambda^n + a_{n-1}\lambda^{n-1} + \dots + \lambda a_1 + a_0 = 0 \quad (31.3.4)$$

If the matrix  $\mathbf{A}$  turns out to be diagonalisable, then the solution is:

$$y(t) = \sum_{i=1}^n C_i e^{\lambda_i t} \quad (31.3.5)$$

However, if the matrix is not diagonalisable, and the eigenvalues are thus degenerate and don't have distinct eigenvalues, then this approach will not work.

### 31.4 Triangulation

**Theorem 4. (Schur's Triangulation Theorem)**

*For any matrix  $\mathbf{A}$ , there is a unitary transformation that converts it into triangular form:*

$$\mathbf{U}^\dagger \cdot \mathbf{A} \cdot \mathbf{U} = T = \begin{pmatrix} \lambda_1 & stuff \\ 0 & \lambda_n \end{pmatrix} \quad (31.4.1)$$

*Proof.* Consider one eigenvector of our matrix:

$$A \cdot \mathbf{v}_1 = \lambda_1 \mathbf{v}_1 \quad (31.4.2)$$

and consider an orthonormal basis  $\{\mathbf{w}_i\}$  with  $\mathbf{w}_1 = \mathbf{v}_1$  and  $\mathbf{w}_i^* \cdot \mathbf{w}_j = \delta_{ij}$ . As a consequence of this orthonormality, the matrix:

$$R = (\mathbf{v}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_n) \quad (31.4.3)$$

is unitary, that is,  $R^{-1} = R^\dagger$ . Therefore, it follows that:

$$R^\dagger \cdot A \cdot R = \begin{pmatrix} \mathbf{v}_1 \\ \mathbf{w}_2 \\ \vdots \\ \mathbf{w}_n \end{pmatrix} \cdot R = (\lambda_1 \mathbf{v}_1 \ A \cdot \mathbf{w}_2 \ \dots \ A \cdot \mathbf{w}_n) = \begin{pmatrix} \lambda_1 & \mathbf{v}_1^* \cdot A \cdot \mathbf{w}_2 & \dots & \mathbf{v}_1^* \cdot A \cdot \mathbf{w}_2 \\ 0 & & & \\ \vdots & & & \mathbf{A}_{n-1} \\ 0 & & & \end{pmatrix} \quad (31.4.4)$$

where  $\mathbf{A}_{n-1}$  has elements  $\mathbf{w}_i^* \cdot A \cdot \mathbf{w}_j$ . We can keep on repeating this to  $\mathbf{A}_{n-1}$  until we reach an upper triangular matrix as required. ■

Once we have found the unitary transformation made up of an eigenvector and an orthonormal basis, we can then find a solution to the system of ODEs:

$$\dot{\mathbf{y}} = A \cdot \mathbf{y} = U \cdot T \cdot U^\dagger \implies \frac{d}{dt}(U^\dagger \cdot \mathbf{y}) = T \underbrace{U^\dagger \cdot \mathbf{y}}_{\zeta} \implies \dot{\zeta} = T \cdot \zeta \quad (31.4.5)$$

where:

$$T = \begin{pmatrix} \lambda_1 & T_{12} & T_{13} & \dots & T_{1n} \\ 0 & \lambda_2 & T_{23} & \dots & T_{2n} \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & 0 & \dots & \lambda_n \end{pmatrix} \quad (31.4.6)$$

This system can be solved component by component, since it is made up of first order ODEs. Once  $\zeta$  has been found, simply revert back to  $\mathbf{y}$ :

$$\mathbf{y} = U \cdot \zeta \quad (31.4.7)$$

## 31.5 Jordan Form

Finally, let us consider a simpler method for solving ODEs with non-diagonalisable matrices.

**Definition 6.1** Consider a degenerate eigenvalue  $\lambda_1$ , called *generator*, that repeats  $m$  times, with at least one associated eigenvector  $\mathbf{v}_1$ . Then,  $\{\mathbf{v}_1 \dots \mathbf{v}_k\}$  is called a *Jordan Chain* if they satisfy:

$$\begin{cases} A \cdot \mathbf{v}_1 = \lambda_1 \mathbf{v}_1 \\ A \cdot \mathbf{v}_2 = \lambda_1 \mathbf{v}_2 + \mathbf{v}_1 \\ \dots \\ A \cdot \mathbf{v}_k = \lambda_1 \mathbf{v}_k + \mathbf{v}_{k-1} \end{cases} \quad (31.5.1)$$

and are therefore linearly independent. We may do so for a set of eigenvalues  $\{\lambda_i\}$  and find the Jordan form:

$$R^{-1} \cdot A \cdot R = J = \text{diag } J_i \quad (31.5.2)$$

where for each eigenvalue  $\lambda_1$  corresponds a Jordan block:

$$J_i = \begin{pmatrix} \lambda_i & 1 & 0 & \dots & 0 \\ 0 & \lambda_i & 0 & \dots & 0 \\ \vdots & & & & \vdots \\ 0 & 0 & 0 & \dots & \lambda_i \end{pmatrix} \quad (31.5.3)$$

**Theorem. (Jordan's Theorem)**

*For any matrix  $A$ , there is always a basis in  $\mathbb{C}^n$  consisting of Jordan chains.*

As usual, we have:

$$\dot{y} = A \cdot y = R \cdot J \cdot R^{-1} \implies \dot{\zeta} = J \cdot \zeta \quad (31.5.4)$$

where  $\zeta = R^{-1} \cdot y$ .

For example, consider the Jordan block  $J_1$ :

$$\begin{cases} \dot{\zeta}_1 = \lambda_1 \zeta_1 + \zeta_2 \\ \dots \\ \dot{\zeta}_{k-1} = \lambda_1 \zeta_{k-1} + \zeta_k \\ \dot{\zeta}_k = \lambda_1 \zeta_k \end{cases} \implies \begin{cases} \zeta_k = C_k e^{\lambda_1 t} \\ \zeta_{k-1} = (C_{k-1} + C_k t) e^{\lambda_1 t} \\ \dots \\ \zeta_1 = \left( C_1 + C_2 t + \dots + C_k \frac{t^{k-1}}{(k-1)!} \right) e^{\lambda_1 t} \end{cases} \quad (31.5.5)$$

Finally, reverting back to  $y$ :

$$y(t) = \zeta_1 v_1 + \dots + \zeta_k v_k \quad (31.5.6)$$

# Series solutions methods

## 32.1 Power Series

We summarise below a set of properties of power series we encountered in the Analysis part:

1. The power series  $\sum_{n=0}^{\infty} a_n(x - x_0)^n$  converges at a point  $x$  if the following limit exists:

$$\lim_{n \rightarrow \infty} \sum_{n=0}^{\infty} a_n(x - x_0)^n \quad (32.1.1)$$

and converges absolutely at a point  $x$  if the associated power series:

$$\sum_{n=0}^{\infty} |a_n(x - x_0)^n| \quad (32.1.2)$$

converges.

2. If a power series converges absolutely, it converges. The converse isn't always true.

3. If  $a_n \neq 0$  and for a fixed  $x$ :

$$L = \lim_{n \rightarrow \infty} |x - x_0| \left| \frac{a_{n+1}}{a_n} \right| \quad (32.1.3)$$

then if  $L < 1$ , the series converges, if  $L > 1$ , the series diverges, and  $L = 1$  is inconclusive.

4. For a power series, there exists  $0 \leq \rho \leq \infty$  called *radius of convergence* such that it will converge for  $|x - x_0| < \rho$  and diverge for  $|x - x_0| > \rho$ .
5. Suppose  $\sum_{n=0}^{\infty} a_n(x - x_0)^n$  converges to  $f(x)$ . Then  $f$  is continuous and is infinitely differentiable over the interval of convergence  $|x - x_0| < \rho$ .
6. The value of  $a_n$  is then given by:

$$a_n = \frac{f^{(n)}(x_0)}{n!} \quad (32.1.4)$$

and the series is then called the *Taylor series*:

$$f(x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(x_0)}{n!} (x - x_0)^n \quad (32.1.5)$$

A function with Taylor series with non-zero radius of convergence is said to be analytic.

7. If  $\sum_{n=0}^{\infty} a_n(x - x_0)^n = \sum_{n=0}^{\infty} b_n(x - x_0)^n$ , then  $a_n = b_n$  for all  $n$ .

## 32.2 Series Solutions near ordinary points

Let us consider the second order linear homogeneous ODE:

$$P(x)y'' + Q(x)y' + R(x)y = 0 \quad (32.2.1)$$

### Definition (*Ordinary point*)

Consider the ODE  $P(x)y'' + Q(x)y' + R(x)y = 0$ , then any point  $x_0$  such that  $P(x_0) = 0$  is an **ordinary point**.

Since  $P$  is continuous, we may always find an interval containing  $x_0$ , and divide the ODE by  $P(x)$  obtaining the more approachable equation

$$y'' + p(x)y' + q(x)y = 0 \quad (32.2.2)$$

where  $p(x) = \frac{Q(x)}{P(x)}$ ,  $q(x) = \frac{R(x)}{P(x)}$ . We now look for series solutions of the form:

$$y = \sum_{n=0}^{\infty} a_n(x - x_0)^n \quad (32.2.3)$$

with some radius of convergence  $\rho > 0$ . One can then substitute this expression into the ODE and use the aforementioned properties of power series to deduce the coefficients.

To illustrate this method, let us solve **Airy's equation**

$$y'' - xy = 0, \quad -\infty < x < \infty. \quad (32.2.4)$$

We note that  $P(x) = 1$ ,  $Q(x) = 0$ ,  $R(x) = -x$ , hence any point is ordinary. We assume that:

$$y = \sum_{n=0}^{\infty} a_n(x)^n \quad (32.2.5)$$

converges in some interval  $|x| < \rho$ . Then:

$$\sum_{n=0}^{\infty} (n+2)(n+1)a_{n+2}x^n - \sum_{n=0}^{\infty} a_n(x)^{n+1} = 0 \quad (32.2.6)$$

and using a shift of index we may rewrite:

$$2a_2 + \sum_{n=1}^{\infty} (n+2)(n+1)a_{n+2}x^n - \sum_{n=1}^{\infty} a_{n-1}(x)^n = 0. \quad (32.2.7)$$

This is only possible for all  $x$  if the coefficients of like powers of  $x$  cancel each other out:

$$(n+2)(n+1)a_{n+2} - a_{n-1} = 0, \quad \text{for } n = 1, 2, 3\dots \quad (32.2.8)$$

This is a second order recurrence relation, so the coefficients will be determined in steps of 3. Note that:

$$a_2 = 0 \implies a_{3n+2} = 0 \quad (32.2.9)$$

Furthermore:

$$a_3 = \frac{a_0}{2 \cdot 3}, \quad a_6 = \frac{a_3}{5 \cdot 6} = \frac{a_0}{2 \cdot 3 \cdot 5 \cdot 6}, \quad a_9 = \frac{a_6}{8 \cdot 9} = \frac{a_0}{2 \cdot 3 \cdot 5 \cdot 6 \cdot 8 \cdot 9} \quad (32.2.10)$$

which suggests (as can be proven by induction):

$$a_{3n} = \frac{a_0}{2 \cdot 3 \cdot 5 \cdot 6 \dots (3n-1)(3n)} \quad (32.2.11)$$

Finally, we also find:

$$a_4 = \frac{a_1}{3 \cdot 4}, \quad a_7 = \frac{a_4}{6 \cdot 7} = \frac{a_1}{3 \cdot 4 \cdot 6 \cdot 7}, \quad a_{10} = \frac{a_7}{9 \cdot 10} = \frac{a_0}{3 \cdot 4 \cdot 6 \cdot 7 \cdot 9 \cdot 10} \quad (32.2.12)$$

which suggests:

$$a_{3n+1} = \frac{a_1}{3 \cdot 4 \cdot 6 \cdot 7 \dots (3n)(3n+1)} \quad (32.2.13)$$

Finally, we arrive at the general solution to Airy's equation:

$$y(x) = a_0 \left( 1 + \sum_{n=1}^{\infty} \frac{x^{3n}}{2 \cdot 3 \cdot 5 \dots (3n)} \right) + a_1 \left( x + \sum_{n=1}^{\infty} \frac{x^{3n+1}}{3 \cdot 4 \cdot 6 \dots (3n+1)} \right) \quad (32.2.14)$$

Also, notice that the first sum satisfies  $y(0) = 0, y'(0) = 0$  and the second sum satisfies  $y(0) = 0, y'(0) = 1$ , which implies that  $W(0) = 1 \neq 0$ , and hence we have truly found the general solution.

It remains to be justified that given an ODE of the form:

$$P(x)y''(x) + Q(x)y'(x) + R(x)y(x) = 0 \quad (32.2.15)$$

we can find a series solution. Consider for example that we have found a solution  $y = \phi(x)$  which can be expanded as a Taylor series:

$$\phi(x) = \sum_{n=0}^{\infty} a_n (x - x_0)^n \quad (32.2.16)$$

which converges in the interval  $|x - x_0| < \rho, \rho > 0$ . We can differentiate (32.2.16)  $m$  times to find:

$$m! \cdot a_m = \phi^{(m)}(x_0) \quad (32.2.17)$$

so the coefficients in the series expansion are determined by evaluating derivatives of  $\phi$  at the ordinary point  $x_0$ . Can this always be done? Suppose we have W.L.O.G the initial conditions  $\phi(x_0) = y_0, \phi'(x_0) = y'_0$  implying that  $a_0 = y_0, a_1 = y'_0$ . Now note that there exists interval around  $x_0$  for which  $P(x) \neq 0$ . Consequently we can determine  $a_2$ :

$$a_2 = \frac{1}{2} \phi''(x_0) = -\frac{1}{2} (p(x_0)\phi'(x_0) + q(x_0)p(x_0)) \quad (32.2.18)$$

where  $p(x_0) = Q(x_0)/P(x_0), q(x_0) = R(x_0)/P(x_0)$ . Similarly we can find  $a_3$  by differentiating (32.2.15) and evaluating at  $x_0$ , and continue for all other  $a_n$ . Note that since  $P, Q, R$  are polynomials they can be differentiated infinitely many times and  $P(x_0) \neq 0$  so we can keep repeating this process indefinitely.

It seems like we only need to assume that the functions  $p(x) = \frac{Q(x)}{P(x)}$  and  $q(x) = \frac{R(x)}{P(x)}$  are smooth near  $x_0$ . This condition is unfortunately too weak, the better condition is that they be analytic at  $x_0$ . We should therefore re-define our definition of ordinary point to the following:

**Definition (Ordinary point redefined)**

The point  $x_0$  is an ordinary point of  $P(x)y'' + Q(x)y' + R(x)y = 0$  if  $p(x) = \frac{Q(x)}{P(x)}$  and  $q(x) = \frac{R(x)}{P(x)}$  are analytic at  $x_0$ .

Regarding the convergence of our series solution, it is expected that it will be bounded below by the convergence of the Taylor series for  $p$  and  $q$ , as the following theorem (which we shall not prove) states.

**Theorem (Convergence radius of series solutions)**

Let  $x_0$  be an ordinary point of:

$$P(x)y'' + Q(x)y' + R(x)y = 0, \quad (32.2.19)$$

Then the general solution may be expressed as a power series about  $x_0$ :

$$y = \sum_{n=0}^{\infty} a_n(x - x_0)^n = a_0y_1 + a_1y_2 \quad (32.2.20)$$

where  $y_1$  and  $y_2$  form a fundamental set of solutions, and the radius of convergence for each of them is greater than or equal to the minimum of the radii of convergence  $\rho_p, \rho_q$  of the Taylor series for  $p = \frac{Q(x)}{P(x)}$  and  $q = \frac{R(x)}{P(x)}$  respectively:

$$\rho_y \geq \min\{\rho_p, \rho_q\}. \quad (32.2.21)$$

Note that if  $P(x), Q(x), R(x)$  are polynomials with common factors cancelled out, then it follows from a theorem in complex analysis that the Taylor series for  $p(x)$  and  $q(x)$  have a radius of convergence equal to the distance between  $x_0$  and the closest zero of  $P(x)$ . This gives us a very powerful tool to determine the convergence properties of series solutions quickly and efficiently. In the case of Airy functions, for example, it is clear that  $P(x) = 1$  has no zeros so the radius of convergence is infinite.

### 32.3 Euler equations

Before developing a full theory for series solutions about singular points, we should consider an illustrative example. The Euler equation reads:

$$x^2y''(x) + \alpha xy'(x) + \beta y(x) = 0, \quad \alpha, \beta \in \mathbb{R} \quad (32.3.1)$$

We begin by assuming  $x > 0$  and try the ansatz  $y = x^r$  then:

$$r(r-1)x^r + \alpha rx^r + \beta x^r = 0 \implies r^2 + (\alpha - 1)r + \beta = 0 \quad (32.3.2)$$

which has solutions:

$$r_{1,2} = -\frac{1}{2} \left[ (\alpha - 1) \pm \sqrt{(\alpha - 1)^2 - 4\beta} \right] \quad (32.3.3)$$

As in the case of second order constant coefficient ODEs, we can distinguish between three different cases:

- (i) If  $(\alpha - 1)^2 - 4\beta > 0$  then we have a fundamental set of solutions and we find:

$$y(x) = c_1 x^{r_1} + c_2 x^{r_2} \quad (32.3.4)$$

- (ii) If  $(\alpha - 1)^2 - 4\beta > 0$  then  $r_1 = \lambda + i\mu, r_2 = \bar{r}_1^*$  for some  $\lambda, \mu \in \mathbb{R}$ . Consequently

$$x^{r_1} = x^\lambda x^{i\mu} = x^\lambda e^{i\mu \ln x} = x^\lambda [\cos(\mu \ln x) + i \sin(\mu \ln x)] \quad (32.3.5)$$

giving the fundamental set of solutions with general form:

$$y(x) = x^\lambda [c_1 \sin(\mu \ln x) + c_2 \cos(\mu \ln x)] \quad (32.3.6)$$

- (iii) If  $(\alpha - 1)^2 - 4\beta = 0$  then  $r_1 = r_2$  so we have only found one solution  $y = x^{r_1}$ . We could find another one by method of reduction of order, but here we present a different derivation that is much faster. Note that:

$$L[x^r] = (r^2 + (\alpha - 1) + \beta)x^r = \left(r^2 + 2\frac{\alpha - 1}{2}r + \frac{(\alpha - 1)}{4}\right)x^r \quad (32.3.7)$$

$$= \left(r - \frac{\alpha - 1}{2}\right)^2 x^r = (r - r_1)^2 x^r \quad (32.3.8)$$

so that

$$L\left[\frac{\partial x^r}{\partial r}\right] = L[x^r \ln x] = \frac{\partial}{\partial r}(L[x^r]) = 2(r - r_1)x^r + (r - r_1)^2 x^r \ln x \quad (32.3.9)$$

For the above to vanish we need  $r = r_1$  so we have found another solution  $y = x^{r_1} \ln x$ . Consequently the general solution is given by:

$$y(x) = (c_1 + c_2 \ln x)x^{r_1} \quad (32.3.10)$$

Finally, note that if we let  $u = -x$  then the Euler equation remains unchanged, so our solutions for  $x > 0$  will also hold for  $x < 0$  by simply changing the sign of  $x$ .

## 32.4 Frobenius' method

We can rewrite the Euler equations as

$$y''(x) + \frac{\alpha}{x}y'(x) + \frac{\beta}{x^2}q(x) = 0 \quad (32.4.1)$$

which suggests generalizing to ODEs of the form

$$y''(x) + p(x)y'(x) + q(x) = 0 \quad (32.4.2)$$

where we can write the following series expansion

$$p(x) = \frac{p_0}{x} + p_1 + p_2 x + p_3 x^2 \dots \quad (32.4.3)$$

$$q(x) = \frac{q_0}{x^2} + \frac{q_1}{x} + q_2 + q_3 x \dots \quad (32.4.4)$$

**Definition (Regular singular points)** The equation (32.4.2) has a **regular singular point** at  $x = x_0$  if  $(x - x_0)p(x)$  and  $(x - x_0)^2q(x)$  are analytic at  $x_0$ . A non-regular singular point is **irregular**.

Let us multiply (32.4.2) by  $x^2$  to get:

$$x^2 y''(x) + x(xp(x))y'(x) + x^2 q(x)y(x) = 0 \quad (32.4.5)$$

We saw that when constants  $\alpha, \beta$  got promoted to functions  $p(x), q(x)$ , it helped to assume a power series solution. Inspired by this we may be able to obtain a solution by multiplying our ansatz  $x^r$  by a power series:

$$y(x) = x^r \sum_{n=0}^{\infty} a_n x^n = \sum_{n=0}^{\infty} a_n x^{n+r} \quad (32.4.6)$$

We see that:

$$\sum_{n=0}^{\infty} a_n (n+r)(n+r-1)x^{n+r} + \left( \sum_{m=0}^{\infty} p_m x^m \right) \sum_{n=0}^{\infty} a_n (n+r)x^{n+r} \quad (32.4.7)$$

$$+ \left( \sum_{m=0}^{\infty} q_m x^m \right) \sum_{n=0}^{\infty} a_n x^{n+r} = 0 \quad (32.4.8)$$

We can divide by  $x^r$  and collect like-terms. For example the  $x^0$  term reads:

$$a_0[r(r-1) + p_0 r + q_0] = 0 \quad (32.4.9)$$

while more generally the  $x^k$  ( $k > 0$ ) coefficient reads

$$a_k(k+r)(k+r-1) + \sum_{n=0}^{k-1} [p_{k-n}a_n(n+r) + q_{k-n}a_n] = 0 \quad (32.4.10)$$

or alternatively:

$$[a_k(k+r)(k+r-1) + p_0 a_k(k+r) + q_0 a_k] + \sum_{n=0}^{k-1} [p_{k-n}a_n(n+r) + q_{k-n}a_n] = 0 \quad (32.4.11)$$

Letting  $F(r) = r(r-1) + p_0 r + q_0$  then we get the following recurrence relation

$$F(k+r)a_k + \sum_{n=0}^{k-1} [p_{k-n}(n+r) + q_{k-n}]a_n = 0$$

(32.4.12)

Note that for  $a_0 \neq 0$  we find that  $a_0 F(r)x^r = 0$  so we get the **indicial equation**

$$F(r) = r(r - 1) + p_0r + q_0 = 0 \quad (32.4.13)$$

which, as before with the Euler equation has two solutions which may or may not be distinct. However, in our case we have a further problem when  $r_1$  and  $r_2$  differ by some integer. We provide the solution without proof.

**Theorem (Series about regular singular points)**

Suppose we have an ODE of the form given in (32.4.2) with a regular singular point at  $x = 0$  so that  $xp(x)$  and  $x^2q(x)$  are analytic at  $x = 0$ :

$$xp(x) = \sum_{n=0}^{\infty} p_n x^n, \quad x^2q(x) = \sum_{n=0}^{\infty} q_n x^n \quad (32.4.14)$$

with radius of convergence  $\rho$ . Let  $r_1$  and  $r_2$  be roots of:

$$r(r - 1) + p_0r + q_0 = 0 \quad (32.4.15)$$

with  $r_1 \geq r_2$ . Then we have two linearly independent solutions on  $0 < x < \rho$ :

(i) if  $r_1 - r_2$  is not an integer then:

$$y(x) = c_1 x^{r_1} \sum_{n=0}^{\infty} a_n x^n + c_2 x^{r_2} \sum_{n=0}^{\infty} b_n x^n \quad (32.4.16)$$

(ii) if  $r_1 = r_2$  then:

$$y(x) = c_1(1 + \ln x) x^{r_1} \sum_{n=0}^{\infty} a_n x^n + c_2 x^{r_1} \sum_{n=0}^{\infty} b_n x^n \quad (32.4.17)$$

(iii) if  $r_1 - r_2 = N$  is an integer then:

$$y(x) = c_1(1 + a \ln x) x^{r_1} \sum_{n=0}^{\infty} a_n x^n + c_2 x^{r_2} \sum_{n=0}^{\infty} b_n x^n \quad (32.4.18)$$

*Proof.* The special case of equal roots to the indicial equation can be solved using the same method used for the Euler equation. We know that one solution can be found:

$$y(x) = x^{r_1} \sum_{n=0}^{\infty} a_n x^n \quad (32.4.19)$$

Now recall that:

$$L[y] = a_0 F(r) x^r + \sum_{n=1}^{\infty} \sum_{k=1}^{\infty} \left\{ [a_k F(k+r) \right. \quad (32.4.20)$$

$$\left. + \sum_{n=0}^{k-1} [p_{k-n} a_n (n+r) + q_{k-n} a_n] \right\} x^{k+r} \quad (32.4.21)$$

so treating the coefficients  $a_m(r)$  as functions of  $r$  then we can set:

$$a_k(r) = -\frac{\sum_{n=0}^{k-1} [p_{k-n}a_n \cdot (n+r) + q_{k-n}a_n]}{F(k+r)} \implies L[y] = a_0 F(r)x^r \quad (32.4.22)$$

giving us the solution to  $L[y] = 0$

$$y(x) = x^{r_1} \sum_{n=0}^{\infty} a_n x^n \quad (32.4.23)$$

where  $r_1$  is the root of  $F(r)$ . Now once again we consider the derivative of  $L[y]$  with respect to  $r$ , and recall that  $F(r) = (r - r_1)^2$  for repeated roots:

$$\frac{\partial}{\partial r} L[y] = L\left[\frac{\partial y}{\partial r}\right] = 2a_0(r - r_1)x^r + a_0(r - r_1)^2x^r \ln x \quad (32.4.24)$$

which again vanishes when  $r = r_1$ . Hence we find that:

$$y_2(x) = \left. \frac{\partial y_1}{\partial r} \right|_{r=r_1} = y_1(x) \ln x + \sum_{n=0}^{\infty} b_n x u u^{n+r_1} \quad (32.4.25)$$

where  $b_n = a'_n(r)|_{r=r_1}$ . We have thus derived the required result. ■

---

# Special functions

- 33.1 Laguerre polynomials
- 33.2 Legendre polynomials
- 33.3 Spherical harmonics
- 33.4 Hermite polynomials
- 33.5 Chebyshev polynomials
- 33.6 Bessel functions

# Distributions

## 34.1 Introducing the Dirac delta

Suppose we have a random variable  $x$  with a gaussian probability distribution:

$$\rho(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-x^2/2\sigma^2} \quad (34.1.1)$$

Note that this being a probability density function implies that it is properly normalized, and thus satisfies:

$$\int_{\mathbb{R}} \rho(x) dx = 1 \quad (34.1.2)$$

We ask what would happen if we let  $\sigma \rightarrow 0$ . Clearly we must have that for  $x \neq 0$ ,  $\lim_{\sigma \rightarrow 0} \rho(x) \neq 0$  due to the exponential suppression overpowering the  $1/\sigma$  factor in the front. However, we must still have that:

$$\lim_{\sigma \rightarrow 0} \int_{\mathbb{R}} \rho(x) dx = 1 \quad (34.1.3)$$

so clearly the value of this function at 0 must be very peculiar.

### Definition (Dirac delta function)

We define the Dirac delta function to be the limit of a normalized gaussian in the zero standard deviation limit:

$$\delta(x) \equiv \lim_{\sigma \rightarrow 0} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-x^2/2\sigma^2} \quad (34.1.4)$$

which satisfies:

$$\begin{cases} \delta(x) = 0, & x \neq 0 \\ \int_{\mathbb{R}} \delta(x) dx = 1 \end{cases} \quad (34.1.5)$$

Statistically, the Dirac delta represents the probability distribution of a random variable with zero variance.

Let's now consider:

$$\mathbb{P}(X \leq x) = \int_{-\infty}^x \rho(x') dx' \quad (34.1.6)$$

which is the probability of the variable  $X$  having a value smaller than  $x$ . As long as  $\sigma > 0$  we find that:

$$\frac{d\mathbb{P}(X \leq x)}{dx} = \rho(x') \quad (34.1.7)$$

Should we expect this to hold for  $\sigma \rightarrow 0$ ? In this limit we have that:

$$\mathbb{P}(X \leq x) \rightarrow \Theta(x) = \begin{cases} 0, & x < 0 \\ 1, & x \geq 0 \end{cases} \quad (34.1.8)$$

since the probability of  $X$  being smaller than one should be zero. Extending (34.1.7) in the  $\sigma \rightarrow 0$  case gives:

$$\frac{d\Theta(x)}{dx} = \delta(x) \quad (34.1.9)$$

We would like to make this reasoning more rigorous, and will require us to introduce the notion of distributions.

To restate our point, consider for example a wire strung between two walls at  $x = \pm L$ , taught with tension  $T$  under a weight  $W$  hung at  $x = 0$ . We should expect there to be no curvature in the wire away from the mass:

$$\frac{d^2y}{dx^2} = 0, \quad x \neq 0 \quad (34.1.10)$$

Also,  $y(\pm L) = 0$  and at  $x = a$  the wire's profile must be continuous (or else it would break). Newton's second law gives us:

$$2T \sin \theta \approx 2T \tan \theta = T \left[ \frac{dy}{dx} \right]_{x=-L}^{x=L} = W \quad (34.1.11)$$

which can be solved:

$$y(x) = \begin{cases} -W(L+x)/2T, & -L < x < 0 \\ -W(L-x)/2T, & 0 < x < L \end{cases} \quad (34.1.12)$$

Note that this solution satisfies:

$$T \frac{d^2y}{dx^2} = W \delta(x) \quad (34.1.13)$$

### Theorem (Delta sifting property)

For an object  $\delta(x)$  satisfying (34.1.5), letting  $f(x) \in L^1(\mathbb{R})$  then:

$$\int_{\mathbb{R}} f(x)\delta(x)dx = f(0) \quad (34.1.14)$$

*Proof.* Let  $\delta > 0$  be small. Then we have that:

$$\int_{\mathbb{R}} f(x)\delta(x)dx = \int_{-\delta}^{\delta} f(x)\delta(x)dx \quad (34.1.15)$$

since  $\delta(x) = 0$  for  $x \neq 0$ . Next:

$$\int_{\mathbb{R}} f(x)\delta(x)dx = \int_{-\delta}^{\delta} (f(x) - f(0) + f(0))\delta(x)dx \quad (34.1.16)$$

$$= \int_{-\delta}^{\delta} (f(x) - f(0))\delta(x)dx + f(0) \quad (34.1.17)$$

Since  $f(x)$  is continuous we must have that for any  $\epsilon > 0$  there is an  $\delta > 0$  such that:

$$|f(x) - f(0)| < \epsilon, \forall x \text{ s.t. } |x| < \delta \quad (34.1.18)$$

so we can make the integrand  $(f(x) - f(0))$  arbitrarily small:

$$\int_{-\delta}^{\delta} (f(x) - f(0))\delta(x)dx = 0 \quad (34.1.19)$$

and thus

$$\int_{\mathbb{R}} f(x)\delta(x)dx = f(0) \quad (34.1.20)$$

as desired.

It is clear that to rigorously treat the delta function we should not treat it as a normal function as the name would suggest, but rather as a generalized function defined through its action on other, well-behaved functions. ■

## 34.2 Rigorous treatment-distributions

We begin by defining exactly how these “well-behaved functions”, known as test functions, must behave like. We will not be particularly interested in the exact forms of these test functions, as long as the following conditions are satisfied.

### Definition (*Test function*)

A function  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  is a test function if:

- (i)  $\phi(x) \in C^\infty$  (smoothness)
- (ii)  $\phi(x) = 0$  has compact support, it vanishes outside some interval.

This definition ensures that  $\phi(x) \rightarrow 0$  at  $\pm\infty$  and also ensures that  $\phi^{(n)}(x)$  is also a test function. These functions are called test functions because they are used to test the action of distributions (which we have yet to define) on them. Although the exact form of test functions is not particularly important, we construct an example for sake of clarity. Indeed one famous example of a test function can be constructed from:

$$\Phi(x) = \begin{cases} 0, & x \leq 0 \\ e^{-1/x}, & x > 0 \end{cases} \quad (34.2.1)$$

Note that this function is infinitely differentiable (with all derivatives vanishing at  $x = 0$  since its  $n$ th derivative will be to leading order  $o(x^n/e^x)$ ). Unfortunately  $\Phi(x)$  is not yet a test function since it does not vanish at  $+\infty$  but we can easily fix this by defining:

$$\phi(x) = \Phi(x)\Phi(1-x) \quad (34.2.2)$$

which does indeed have a compact support.

### Definition (*Action on test functions*)

Let  $f(x) \in L^1(\mathbb{R})$  be a measurable function. Then we define its action on a test function  $\phi(x)$

by:

$$\langle f, \phi \rangle = \int_{\mathbb{R}} f(x) \phi(x) dx \quad (34.2.3)$$

Note that this is not the same as a normal inner product since  $f$  and  $\phi$  need not lie in the same vector space. We should ask ourselves if this definition of action is well-defined: can two functions have the same action on a given test function?

**Theorem (Unique definition of distribution)**

Let  $f_{1,2} : \mathbb{R} \rightarrow \mathbb{R}$  be continuous functions. If  $\langle f_1, \phi \rangle = \langle f_2, \phi \rangle$  for all test functions  $\phi$  then  $f_1(x) = f_2(x)$ .

*Proof.* We begin by proving that if:

$$\langle f, \phi \rangle = \int_{\mathbb{R}} f(x) \phi(x) dx = 0 \quad (34.2.4)$$

for all test functions then  $f(x) = 0$ . Suppose that  $f(a) > 0$  at some point  $x = a$ . Due to the continuity of  $f$  there exists  $\delta > 0$  such that  $f(x) > 0$  for all  $x \in (a - \delta, a + \delta)$ . We can find a test function  $\phi(x)$  which vanishes outside  $(a - \delta, a + \delta)$  and is non-zero inside this interval, implying that:

$$\langle f, \phi \rangle = \int_{a-\delta}^{a+\delta} f(x) \phi(x) dx > 0 \quad (34.2.5)$$

This however is a contradiction, so we cannot have that  $f(a) > 0$  at some  $x = a$ . Applying this result to  $f(x) = f_1(x) - f_2(x)$  immediately reproduces the theorem. ■

Much like in functional analysis, we will also need to define a property of convergence for sequences of test functions  $\{\phi_n(x)\}$ .

**Definition (Convergence of test function sequences)** The sequence  $\{\phi_n\}$  of test functions converges to zero if:

- (i) for all  $x$  outside some interval  $I$ ,  $\phi_n(x) = 0$  for all  $n$ , so there is a shared compact support for all test functions in the sequence,
- (ii) for all  $k$ ,  $\phi_n^{(k)}$  converges uniformly to 0 as  $n \rightarrow \infty$ .

We are now ready to define distributions.

**Definition (Distribution)**

A distribution (or generalised function)  $\mathcal{F}$  is a continuous functional mapping from the set of test functions  $\mathbb{D}$  to  $\mathbb{R}$  defined by the action:

$$\phi \rightarrow \langle \mathcal{F}, \phi \rangle \in \mathbb{R} \quad (34.2.6)$$

It must satisfy:

- (i) Continuous: if  $\phi_n \rightarrow 0$  then  $\langle \mathcal{F}, \phi_n \rangle \rightarrow 0$ .
- (ii) Linear:  $\langle \mathcal{F}, \alpha\phi + \beta\psi \rangle = \alpha \langle \mathcal{F}, \phi \rangle + \beta \langle \mathcal{F}, \psi \rangle$  for real constants  $\alpha, \beta$ .

An example of a distribution is the Heaviside distribution  $\mathcal{T}$  generated from the Heaviside function  $\Theta(x)$ :

$$\langle \mathcal{T}, \phi \rangle = \int_0^\infty \phi(x) dx \quad (34.2.7)$$

Clearly this distribution is linear. It is continuous since if  $\phi_n(x) \rightarrow 0$  then  $\langle \mathcal{T}, \phi_n \rangle \rightarrow 0$  too.

### Proposition (Delta function as a distribution)

The delta function  $\delta(x)$  defined by the action  $\langle \delta, \phi \rangle = \phi(0)$  is a distribution.

*Proof.* Linearity is trivial. If  $\phi_n \rightarrow 0$  then by the uniform convergence of test functions  $\phi_n(0) \rightarrow 0$  so  $\langle \delta, \phi_n \rangle \rightarrow 0$ . ■

### Proposition (Delta function on $\mathcal{C}^0(\mathbb{R})$ )

If  $f(x)$  is continuous at  $x = 0$  then

$$\langle \delta_n, f \rangle \rightarrow f(0) \text{ as } n \rightarrow \infty \quad (34.2.8)$$

where

$$\delta_n = \begin{cases} \frac{n}{2}, & |x| < \frac{1}{n} \\ 0, & |x| \geq \frac{1}{n} \end{cases} \quad (34.2.9)$$

so  $\delta(x)$  can be applied on continuous functions too:

$$\langle \delta, f \rangle = f(0), \forall f \in \mathcal{C}^0(\mathbb{R}) \quad (34.2.10)$$

*Proof.* For any continuous function  $f(x) \in \mathcal{C}^0(\mathbb{R})$ :

$$\langle \delta_n, f \rangle = \frac{n}{2} \int_{-1/n}^{1/n} f(x) dx = \frac{n}{2} f(\eta_n) \int_{-1/n}^{1/n} dx = f(\eta_n) \rightarrow f(0) \quad (34.2.11)$$

■

This is a very important result because it means that the action of  $\delta(x)$  is not restricted to the space of test functions as is the default case for distributions. It makes sense to apply  $\delta(x)$  to normal, continuous functions too.

### Definitions (Operations with distributions)

Let  $F(x)$  be a distribution. Then we define:

$$\langle F(x-a), \phi(x) \rangle = \langle F(x), \phi(x+a) \rangle \quad (34.2.12)$$

$$\langle F(ax), \phi(x) \rangle = \frac{1}{|a|} \langle F(x), \phi(x/a) \rangle, a \neq 0 \quad (34.2.13)$$

Applying this to the delta function we see that:

$$\langle \delta(x-a), f(x) \rangle = \int_{\mathbb{R}} \delta(x-a) f(x) dx = f(a)$$

(34.2.14)

and

$$\langle \delta(ax), f(x) \rangle = \int_{\mathbb{R}} \delta(ax) f(x) dx = \frac{1}{|a|} f(0) \quad (34.2.15)$$

known as sifting properties.

To see where these definitions come from, let us turn to the integral approach and consider these distributions as normal functions. Then

$$\langle \mathcal{F}(x-a), \phi(x) \rangle = \int_{\mathbb{R}} \mathcal{F}(x-a) \phi(x) dx = \int_{\mathbb{R}} \mathcal{F}(s) \phi(s+a) ds = \langle \mathcal{F}(x), \phi(x+a) \rangle \quad (34.2.16)$$

and similarly:

$$\langle \mathcal{F}(ax), \phi(x) \rangle = \int_{\mathbb{R}} \mathcal{F}(ax) \phi(x) dx = \frac{1}{|a|} \int_{\mathbb{R}} \mathcal{F}(s) \phi(s/a) ds = \frac{1}{|a|} \langle \mathcal{F}(x), \phi(x/a) \rangle \quad (34.2.17)$$

where the modulus originates from the fact that if  $a < 0$  then the integral bounds are flipped, giving a negative sign.

Similarly, we may use this approach to motivate a definition for the derivative of a distribution. Then we see that:

$$\langle \mathcal{F}', \phi \rangle = \int_{\mathbb{R}} \mathcal{F}'(x) \phi(x) dx = [\mathcal{F}(x)\phi(x)]_{-\infty}^{\infty} - \int_{\mathbb{R}} \mathcal{F}(x) \phi'(x) dx = -\langle \mathcal{F}, \phi \rangle \quad (34.2.18)$$

prompting us to define:

### Definition (*Differentiating distributions*)

Let  $\mathcal{F}$  be a distribution. Then its derivative  $\mathcal{F}'$  is defined so as to satisfy  $\langle \mathcal{F}', \phi \rangle = -\langle \mathcal{F}, \phi \rangle$ .

We know that the derivative of a test function is a distribution. Are the derivatives of distributions also differentiable then?

**Theorem (*Derivatives of distributions are distributions*)** Let  $\mathcal{F}$  be a distribution, then so is its derivative.

*Proof.* The action of  $\mathcal{F}'$  is linear due to the linearity of  $\mathcal{F}$ . Also, if  $\phi_n \rightarrow 0$  then  $\phi' \rightarrow 0$  by uniform convergence so that:

$$\langle \mathcal{F}', \phi_n \rangle = -\langle \mathcal{F}, \phi'_n \rangle \rightarrow 0 \quad (34.2.19)$$

■

The above result also tells us that since test functions can be infinitely differentiated, any distribution can also be differentiated infinitely many times. This allows us to prove a wonderful result. Consider the derivative of the Heaviside function:

$$\langle \mathcal{T}', \phi \rangle = -\langle \mathcal{T}, \phi' \rangle = -\int_0^{\infty} \phi'(x) dx = -\phi(0) = -\langle \delta, \phi \rangle \implies \mathcal{T}' = -\delta \quad (34.2.20)$$

In other words we get the result we previously wanted to prove rigorously, that the derivative of the Heaviside function is the Dirac delta function:

$$\boxed{\frac{d\Theta(x)}{dx} = \delta(x)} \quad (34.2.21)$$

Similarly we see that the derivative of the delta function can be defined via the action:

$$\langle \delta', f \rangle = -\langle \delta, f' \rangle = -f'(0) \quad (34.2.22)$$

It seems like these distributions satisfy typical rules of differential calculus. For example, we can recover the product rule of differentiation using the definition and properties of distributions:

**Proposition (Leibniz rule)**

If  $\mathcal{F}$  is a distribution and  $f \in C^\infty(\mathbb{R})$  then  $(f\mathcal{F})' = f\mathcal{F}' + f'\mathcal{F}$ .

*Proof.* We have that:

$$\langle (f\mathcal{F})', \phi \rangle = -\langle f\mathcal{F}, \phi' \rangle = -\langle \mathcal{F}, \phi' f \rangle \quad (34.2.23)$$

$$= -\langle \mathcal{F}, (\phi f)' \rangle + \langle \mathcal{F}, f' \phi \rangle = \langle f\mathcal{F}', \phi \rangle + \langle f'\mathcal{F}, \phi \rangle \quad (34.2.24)$$

as desired. ■

# Laplace transform methods

## 35.1 Basic definition and properties of the Laplace transform

The Laplace transforms is a very important type of integral transform which is a precursor to its cousin, the Fourier transform. It, like its cousin, can be used to solve both ordinary and partial differential equations with more immediacy.

### Definition (Laplace transform)

The Laplace transform maps a function  $f(t) \in L^p(\mathbb{R}^+)$  to another function  $F(s) = \mathcal{L}\{f(t)\}$  defined by:

$$F(s) = \mathcal{L}\{f(t)\} \equiv \int_0^\infty e^{-st} f(t) dt \quad (35.1.1)$$

It is important to note that the Laplace transform is linear due to the linearity of integration so that:

$$\mathcal{L}\{\alpha f(t) + \beta g(t)\} = \alpha \mathcal{L}\{f(t)\} + \beta \mathcal{L}\{g(t)\} \quad (35.1.2)$$

**Example.** Compute the Laplace transform of the Heaviside function:

$$\Theta(t) = \begin{cases} 1, & t \geq 0 \\ 0, & t \leq 0 \end{cases} \quad (35.1.3)$$

We find that:

$$\mathcal{L}\{\Theta(t)\} = \int_0^\infty e^{-st} dt = -\frac{1}{s} \lim_{\tau \rightarrow \infty} (e^{-\tau s} - 1) = \frac{1}{s}, \quad s > 0 \quad (35.1.4)$$

where we must assume that  $s > 0$  for the Laplace transform to be well defined. By the linearity of  $\mathcal{L}$  we may write:

$$\mathcal{L}\{c\Theta(t)\} = \frac{c}{s} \quad (35.1.5)$$

Note also that for any function  $f(t)$  integrable on  $[0, \infty)$ :

$$\mathcal{L}\{f(t)\Theta(t)\} = \mathcal{L}\{f(t)\} \quad (35.1.6)$$

so from (35.1.5) we have found that:

$$\mathcal{L}\{c\} = \frac{c}{s} \quad (35.1.7)$$

Note that when performing the Laplace transform of some function we lose all information about its behaviour for  $t < 0$ , we say that the Laplace transform is not unitary. Consequently, one should not expect to be able to invert the Laplace transform and get back the initial function in its original domain since we don't know how it behaves for  $t < 0$ . There is an ambiguity in choosing the function's behaviour if we want to define it over  $\mathbb{R}$ . We can solve this problem by using the Heaviside function. Namely:

$$\mathcal{L}\{f(t)\Theta(t)\} = F(s) \implies \mathcal{L}\{F(s)\}^{-1} = f(t)\Theta(t) \quad (35.1.8)$$

We have no ambiguity here since we simply assume that  $f(t) = 0$ ,  $t < 0$ .

**Example.** Compute the Laplace transform of  $f(t) = t$ .

We need to evaluate the following integral:

$$\mathcal{L}\{t\} = \int_0^\infty te^{-st} dt \quad (35.1.9)$$

which can be solved using Feynman's trick. Let:

$$I(s) = \int_0^\infty te^{-st} dt \quad (35.1.10)$$

Then:

$$I(s) = -\frac{d}{ds} \left( \int_0^\infty e^{-st} dt \right) = -\frac{d}{ds} \left( \frac{1}{s} \right) = \frac{1}{s^2}, \quad s > 0 \quad (35.1.11)$$

so we have that:

$$\mathcal{L}\{t\} = \mathcal{L}\{t\Theta(t)\} = \frac{1}{s^2}, \quad s > 0 \quad (35.1.12)$$

◀

**Example.** Compute the Laplace transform of  $f(t) = \sin(kt)$  and  $g(t) = \cos(kt)$ . It is useful to first evaluate  $\mathcal{L}\{e^{kt}\}$ . This is a trivial integral:

$$\mathcal{L}\{e^{kt}\} = \mathcal{L}\{e^{kt}\Theta(t)\} = \int_0^\infty e^{-(s-k)t} dt = \frac{1}{s-k}, \quad s > k \quad (35.1.13)$$

It follows that:

$$\mathcal{L}\{e^{ikt}\} = \mathcal{L}\{\cos kt\} + i\mathcal{L}\{\sin kt\} = \frac{s}{s^2 + k^2} + i\frac{k}{s^2 + k^2}, \quad s > 0 \quad (35.1.14)$$

implying that:

$$\mathcal{L}\{\cos kt\} = \frac{s}{s^2 + k^2}, \quad \mathcal{L}\{\sin kt\} = \frac{k}{s^2 + k^2}, \quad s > 0 \quad (35.1.15)$$

◀

**Proposition (Properties of the Laplace transform)** Given a function  $f(t) \in L^p(\mathbb{R}^+)$  with Laplace transform:

$$F(s) = \mathcal{L}\{f(t)\} \quad (35.1.16)$$

we have that:

$$\mathcal{L}\{e^{kt}f(t)\} = F(s - k) \quad (35.1.17)$$

and

$$\mathcal{L}\{t^n f(t)\} = (-1)^n \frac{d^n F(s)}{ds^n} \quad (35.1.18)$$

*Proof.* Firstly:

$$\mathcal{L}\{e^{kt} f(t)\} = \int_0^\infty e^{-(s-k)t} f(t) dt = F(s - k) \quad (35.1.19)$$

Secondly:

$$\mathcal{L}\{t^n f(t)\} = \int_0^\infty t^n f(t) e^{-st} dt = (-1)^n \frac{d^n}{ds^n} \left( \int_0^\infty f(t) e^{-st} dt \right) = (-1)^n \frac{d^n F(s)}{ds^n} \quad (35.1.20)$$

■

### Theorem (Laplace transform of a derivative)

Let  $f(t)$  be a  $n$ -differentiable function at 0. Then if  $\mathcal{L}\{f(t)\} = F(s)$  then:

$$\mathcal{L}\{f^{(n)}(t)\} = s^n F(s) - \sum_i^{n-1} s^{n-i} f^{(i)}(0) \quad (35.1.21)$$

*Proof.* We proceed by induction. We have that:

$$\mathcal{L}\{\dot{f}(t)\} = \int_0^\infty \dot{f}(t) e^{-st} dt = \int_0^\infty \frac{d}{dt}(f(t) e^{-st}) dt + \int_0^\infty s f(t) e^{-st} dt \quad (35.1.22)$$

$$= sF(s) + [f(t)e^{-st}]_0^\infty \quad (35.1.23)$$

$$= sF(s) - f(0) \quad (35.1.24)$$

assuming that  $f(t)$  is dominated by the exponentially decaying  $e^{-st}$  as  $t \rightarrow \infty$ . Let us now suppose that (35.1.21) is true up to  $n$ . Then:

$$\mathcal{L}\{f^{(n+1)}(t)\} = \mathcal{L}\{\dot{f^{(n)}}(t)\} = s\mathcal{L}\{f^{(n)}(t)\} - f^{(n)}(0) \quad (35.1.25)$$

$$= s^{n+1} F(s) + \sum_i^n s^{n-i} f^{(i)}(0) \quad (35.1.26)$$

■

### Theorem (Laplace transform of an integral)

Let  $f(t')$  be integrable over  $[0, t]$ . Then:

$$\mathcal{L}\left\{\int_0^t f(t') dt'\right\} = \frac{1}{s} F(s) \quad (35.1.27)$$

*Proof.* Define:

$$g(t) = \int_0^t f(t') dt' \implies \dot{g}(t) = f(t) \quad (35.1.28)$$

Then:

$$\mathcal{L}\{\dot{g}(t)\} = F(s) = s\mathcal{L}\{g(t)\} \implies \mathcal{L}\{g(t)\} = \frac{1}{s}F(s) \quad (35.1.29)$$

as desired. ■

**Proposition (Second shift theorem)** We have that:

$$\mathcal{L}\{f(t-a)\Theta(t-a)\} = e^{-sa}F(s) \quad (35.1.30)$$

We have that:

$$\mathcal{L}\{f(t-a)\Theta(t-a)\} = \int_0^\infty \Theta(t-a)e^{-st}f(t-a)dt \quad (35.1.31)$$

$$= \int_a^\infty e^{-st}f(t-a)dt = e^{-sa} \int_0^\infty e^{-st'}f(t')dt' = e^{-sa}F(s) \quad (35.1.32)$$

## 35.2 Solving ODEs with Laplace transforms

Consider the general linear inhomogeneous second order ODE with constant coefficients:

$$a\ddot{x} + b\dot{x} + cx = f(t), \quad x(0) = x_0, \dot{x}(0) = 0 \quad (35.2.1)$$

where  $f(t)$  has Laplace transform  $F(s)$ . We can take the Laplace transform of this equation:

$$a(s^2\mathcal{L}\{x(t)\} - sx_0 - 0) + b(s\mathcal{L}\{x(t)\} - x_0) + c\mathcal{L}\{x(t)\} = F(s) \quad (35.2.2)$$

$$\implies (as^2 + bs + c)\mathcal{L}\{x(t)\} - asx_0 - bsx_0 = F(s) \quad (35.2.3)$$

$$\implies \mathcal{L}\{x(t)\} = \frac{F(s) + asx_0 + bx_0}{as^2 + bs + c} \quad (35.2.4)$$

It is interesting to note that the characteristic polynomial of the ODE popped up in the denominator of this Laplace transform! For the homogeneous case where  $f(t) = 0 \implies F(s) = 0$ , assuming that the characteristic polynomial has roots at  $\lambda_{1,2}$  then we find that:

$$\mathcal{L}\{x(t)\} = \frac{asx_0 + bx_0}{(s - \lambda_1)(s - \lambda_2)} = \frac{A}{(s - \lambda_1)} + \frac{B}{(s - \lambda_2)} \quad (35.2.5)$$

where

$$A = \frac{a\lambda_1 - b}{\lambda_1 - \lambda_2}x_0, \quad B = \frac{b - a\lambda_2}{\lambda_1 - \lambda_2}x_0 \quad (35.2.6)$$

We can invert the Laplace transform and find that:

$$x(t) = Ae^{\lambda_1 t} + Be^{\lambda_2 t} \quad (35.2.7)$$

as expected. The importance of Laplace transforms is now clear: it is a very useful tool in solving differential equations.

### 35.3 Convolutions

**Definition (Convolution integral)**

Given two functions  $f(t)$  and  $g(t)$ , their convolution  $(f * g)(t)$  is defined as:

$$(f * g)(t) = \int_{-\infty}^{\infty} f(s)g(t-s)ds \quad (35.3.1)$$

Intuitively, the convolution integral gives us the overlap between  $f$  and  $g$  as we shift one relative to the other by  $t$  along  $s$ .

**Proposition (Causal functions)**

A function  $f(t)$  is causal if  $f(t) = 0, \forall t < 0$ . The convolution of two causal functions  $f, g$  is given by:

$$(f * g)(t) = \int_0^t f(s)g(t-s)ds \quad (35.3.2)$$

*Proof.* We have that

$$\int_{-\infty}^{\infty} f(s)g(t-s)ds = \int_0^{\infty} f(s)g(t-s)ds = \int_0^t f(s)g(t-s)ds \quad (35.3.3)$$

since  $g(t-s) = 0$  for  $s > t$ . ■

**Theorem (Convolution theorem for Laplace transforms)**

The Laplace transform of the convolution of two causal functions is the product of their Laplace transforms:

$$\mathcal{L}\{(f * g)(t)\} = \mathcal{L}\{f(t)\}\mathcal{L}\{g(t)\} \quad (35.3.4)$$

*Proof.* We have that:

$$\mathcal{L}\{(f * g)(t)\} = \int_0^{\infty} e^{-st} \int_0^t f(x)g(t-x)dxdt \quad (35.3.5)$$

$$= \int_0^{\infty} \int_x^{\infty} e^{-st} g(t-x)dt f(x)dx \quad (35.3.6)$$

$$= \int_0^{\infty} \int_0^{\infty} e^{-s(v+x)} g(v+x-x)dv f(x)dx \quad (35.3.7)$$

$$= \int_0^{\infty} g(v)e^{-sv}dv \int_0^{\infty} f(x)e^{-sx}dx = \mathcal{L}\{g(t)\}\mathcal{L}\{f(t)\} \quad (35.3.8)$$

as desired. ■

**Example.** Let's find the inverse Laplace transform of:

$$F(s) = \frac{a}{s^4 + a^2 s^2} \quad (35.3.9)$$

We have that:

$$\mathcal{L}^{-1}(F(s)) = \mathcal{L}^{-1}\left(\frac{1}{s^2} \frac{a}{s^2 + a^2}\right) \quad (35.3.10)$$

and since  $\mathcal{L}^{-1}(1/s^2) = t$  and  $\mathcal{L}^{-1}(a/(s^2 + a^2)) = \sin(at)$  one finds:

$$\mathcal{L}^{-1}(F(s)) = \mathcal{L}^{-1}(\mathcal{L}\{t * \sin(at)\}) = \int_0^t (t-s) \sin(as) ds = \frac{at - \sin at}{a^2} \quad (35.3.11)$$

◀

---

# Phase plane analysis

36

# First order PDEs

## 37.1 Introduction

**Definition (Linear first order PDE)** A linear first order partial differential equation is of the form

$$\sum_{i=1}^n a_i(x^1, \dots, x^n) \frac{\partial u}{\partial x^i} + c(x, y)u = d(x, y) \quad (37.1.1)$$

where  $a_i, b, d$  are continuous functions. In two dimensions this reads

$$a(x, y) \frac{\partial u}{\partial x} + b(x, y) \frac{\partial u}{\partial y} + c(x, y)u = d(x, y) \quad (37.1.2)$$

where  $a, b, c, d$  are continuous functions.

As with first order ODEs, there are some special cases of linear first order PDEs that can be considered.

## 37.2 From ODEs to PDEs

**Definition (Separable linear first order PDE)** A first order PDE is said to be separable if it is of the form

$$\frac{\partial u}{\partial x^i} = f(x^1, \dots, x^{i-1}, x^{i+1}, \dots)g(x^i) \quad (37.2.1)$$

Separable PDEs can be solved by integrating directly, this time remembering that all but one variable are not integrating over and should thus be treated as constants. In other words

$$\frac{\partial u}{\partial x^i} = f(x^1, \dots, x^{i-1}, x^{i+1}, \dots)g(x^i) \implies \int \frac{\partial u}{\partial x^i} dx^i = \int f(x^1, \dots, x^{i-1}, x^{i+1}, \dots)g(x^i)dx^i \quad (37.2.2)$$

Consider for example the following

$$\frac{\partial u}{\partial y} = xy \implies u(x, y) - f(x) = \frac{1}{2}xy^2 \implies u(x, y) = \frac{1}{2}xy^2 + f(x) \quad (37.2.3)$$

Note that the general solution to this first order PDE in  $y$  contains an arbitrary function of  $x$ , just like the solution to a first order ODE contains an arbitrary constant.

Another case which can be solved in a manner analogous to ODEs are those of the type

$$\frac{\partial u}{\partial x} + a(x, y)u = b(x, y) \quad (37.2.4)$$

These can be solved using the method of integrating factor. Suppose we have found the integrating factor  $\Lambda$  using the usual method, then the above can be written as

$$\frac{\partial}{\partial x}(\Lambda(x, y)u(x, y)) = \Lambda(x, y)b(x, y) \quad (37.2.5)$$

which can be integrated to give

$$\Lambda(x, y)u(x, y) = B(x, y) + f(y) \implies y(x, y) = \frac{B(x, y)}{\Lambda(x, y)} + \frac{f(y)}{\Lambda(x, y)} \quad (37.2.6)$$

**Example.** For example, consider the following first order PDE

$$\frac{\partial u}{\partial y} + \frac{3}{y}u = y^2, \quad y \neq 0 \quad (37.2.7)$$

The integrating factor is

$$\Lambda(x, y) = \exp\left(\int \frac{3}{y} dy\right) = e^{3 \ln y} = y^3 \quad (37.2.8)$$

so we find

$$\frac{\partial}{\partial y}(y^3 u(x, y)) = y^5 \implies u(x, y) = \frac{1}{6}y^3 + \frac{f(x)}{y^3} \quad (37.2.9)$$

As another example, consider

$$\frac{\partial u}{\partial y} - xyu = y, \quad x \neq 0 \quad (37.2.10)$$

The integrating factor is

$$\Lambda(x, y) = e^{-xy^2/2} \implies \frac{\partial}{\partial y}(e^{-xy^2/2}u(x, y)) = ye^{-xy^2/2} \quad (37.2.11)$$

which can be integrated to give

$$e^{-xy^2/2}u(x, y) = -\frac{e^{-xy^2/2}}{x} + f(x) \implies u(x, y) = -\frac{1}{x} + f(x)e^{xy^2/2} \quad (37.2.12)$$

◀

We can solve certain second order PDEs by transforming it into a first order PDE, as the next example shows.

**Example.** Consider  $\frac{\partial^2 y}{\partial x \partial t} + \frac{1}{t} \frac{\partial u}{\partial x} = x$ ,  $t \neq 0$ . We can transform this into a first order PDE by letting  $f(x, y) = \frac{\partial u}{\partial x}$ :

$$\frac{\partial f}{\partial t} + \frac{1}{t}f = x \quad (37.2.13)$$

This can be solved by the method of integrating factors, it is easy to see that  $\Lambda(t) = t$  so that

$$\frac{\partial}{\partial t}(tf) = tx \implies f(x, t) = \frac{1}{2}tx + \frac{1}{t}h(x) \quad (37.2.14)$$

To find  $u(x, t)$  we substitute back  $f = \frac{\partial u}{\partial x}$  which yields

$$\frac{\partial u}{\partial x} = \frac{1}{2}tx + \frac{1}{t}h(x) \quad (37.2.15)$$

which can be integrated directly to give the general solution:

$$u(x, t) = \frac{1}{4}tx^2 + \frac{1}{t}h(x) + g(t) \quad (37.2.16)$$

◀

### 37.3 Change of variable and the chain rule

Substitutions of variables are often useful in simplifying PDEs, this requires the use of the chain rule.

For example, consider the first order PDE

$$\frac{\partial u}{\partial x} - \frac{\partial u}{\partial y} + u = 2 \quad (37.3.1)$$

We may define  $x = \eta$  and  $y = \phi - \eta$ . Then we see that

$$\frac{\partial u}{\partial \eta} = \frac{\partial u}{\partial x} - \frac{\partial u}{\partial y} \quad (37.3.2)$$

so that the PDE transforms to

$$\frac{\partial u}{\partial \eta} + u = 2 \quad (37.3.3)$$

We can solve this equation for  $u$  using the integrating factor method

$$\Lambda(x, y) = e^\eta \implies u(\eta) = 2 + e^{-\eta}f(\phi) \quad (37.3.4)$$

which in the old coordinates reads

$$f(x, y) = 2 + e^{-x}f(x + y) \quad (37.3.5)$$

### 37.4 The method of characteristics

Consider the PDE

$$\frac{\partial u}{\partial x} + \frac{\partial u}{\partial y} = 0 \quad (37.4.1)$$

We know that letting  $\chi = x$  and  $\eta = x - y$  then it will transform to

$$\frac{\partial u}{\partial \chi} = 0 \implies u(x, y) = f(x - y) \quad (37.4.2)$$

Note that by performing this special change of coordinates we got a PDE depending on only one variable. The corresponding solutions are therefore found by integrating over level curves of the missing variable, in the previous example by keeping  $x - y = \text{cnst}$ . Such curves are known as characteristic curves of the PDE and are very useful tools.

Consider the general first order PDE

$$a(x, y) \frac{\partial u}{\partial x} + b(x, y) \frac{\partial u}{\partial y} + c(x, y)u = d(x, y) \quad (37.4.3)$$

where  $a(x, y) \neq 0$ . Thus we may equivalently write

$$\frac{\partial u}{\partial x} + g(x, y) \frac{\partial u}{\partial y} + h(x, y)u = k(x, y) \quad (37.4.4)$$

We define two new variables  $\chi = \chi(x, y)$  and  $\eta = \eta(x, y)$  so that

$$\frac{\partial u}{\partial \chi} + h(\chi, \eta)u = k(\chi, \eta) \quad (37.4.5)$$

This can be done if we let

$$\frac{\partial u}{\partial \chi} = \frac{\partial u}{\partial x} \frac{\partial x}{\partial \chi} + \frac{\partial u}{\partial y} \frac{\partial y}{\partial \chi} = \frac{\partial u}{\partial x} + g(x, y) \frac{\partial u}{\partial y} \quad (37.4.6)$$

so by solving the following system of PDEs

$$\begin{cases} \frac{\partial x}{\partial \chi} = 1 \\ \frac{\partial y}{\partial \chi} = g(x, y) \end{cases} \quad (37.4.7)$$

Let  $\chi = x$ , then the first equation is satisfied, and along the characteristic curves  $\eta = \text{csnt}$  we find that

$$\frac{\partial y}{\partial \chi} = \frac{dy}{dx} = g(x, y) \quad (37.4.8)$$

Thus solving the above ODE will yield the general solutions  $\alpha(x, y) = c$  which are the required characteristic curves if we define  $\eta = \alpha(x, y)$ .

**Example.** Consider for example

$$x^2 \frac{\partial f}{\partial x} - xy \frac{\partial f}{\partial y} + y^2 = 0, \quad x \neq 0 \quad (37.4.9)$$

which we can write equivalently as

$$\frac{\partial f}{\partial x} - \frac{y}{x} \frac{\partial f}{\partial y} + \frac{y^2}{x^2} = 0 \quad (37.4.10)$$

To find the characteristic curves we need to solve the ODE

$$\frac{dy}{dx} = -\frac{y}{x} \implies xy = c \quad (37.4.11)$$

so we see that the characteristics are hyperbolae. Hence let  $\chi = x$  and  $\eta = xy$ . Then we have

that

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial \chi} \frac{\partial \chi}{\partial y} + \frac{\partial f}{\partial \eta} \frac{\partial \eta}{\partial y} = x \frac{\partial f}{\partial \eta} \quad (37.4.12)$$

and similarly

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial \chi} \frac{\partial \chi}{\partial x} + \frac{\partial f}{\partial \eta} \frac{\partial \eta}{\partial x} = \frac{\partial f}{\partial \chi} + y \frac{\partial f}{\partial \eta} \quad (37.4.13)$$

The PDE now transforms to

$$\frac{\partial f}{\partial \chi} + \frac{\eta^2}{\chi^4} = 0 \quad (37.4.14)$$

which we can integrate directly to find

$$f(\chi, \eta) = \frac{\eta^2}{3\chi^3} + f(\eta) \implies f(x, y) = \frac{y^2}{3x} + f(xy) \quad (37.4.15)$$

◀

We can apply this reasoning to some second order PDEs by transforming them into first order PDEs. This is best illustrated in the following example.

**Example.** Consider

$$\frac{\partial^2 u}{\partial x \partial y} + 3x^2 \frac{\partial^2 u}{\partial y^2} - 2 \frac{\partial u}{\partial y} = 0 \quad (37.4.16)$$

Let  $f = \frac{\partial u}{\partial y}$  this becomes a first order PDE

$$\frac{\partial f}{\partial x} + 3x^2 \frac{\partial f}{\partial y} - 2f = 0 \quad (37.4.17)$$

To find the characteristic curves we solve  $y' = 3x^2 \implies y - x^3 = c$ . Thus we define the new coordinates  $\chi = x$  and  $\eta = y - x^3$ . Then by the chain rule

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial \chi} - 3x^2 \frac{\partial f}{\partial \eta}, \quad \frac{\partial f}{\partial y} = 3x^2 \frac{\partial f}{\partial \eta} \quad (37.4.18)$$

We therefore get a PDE in  $\chi$  only

$$\frac{\partial f}{\partial \chi} - 2f = 0 \quad (37.4.19)$$

which can be integrated directly to give

$$f(\chi, \eta) = g(\eta)e^{2\chi} \implies f(x, y) = g(y - x^3)e^{2x} \quad (37.4.20)$$

This means that

$$\frac{\partial u}{\partial y} = g(y - x^3)e^{2x} \implies u(x, y) = h(y - x^3)e^{2x} + k(x) \quad (37.4.21)$$

is the general solution. ◀

---

# **Second order PDEs: an overview**

38

---

# Elliptic PDEs: Electrostatics

## 39.1 Existence and uniqueness theorem

# Hyperbolic PDEs: Waves

## 40.1 The wave equation

We are interested in the solutions of the following equation

$$\boxed{\frac{\partial^2 \phi}{\partial t^2} = c^2 \nabla^2 \phi} \quad (40.1.1)$$

but before we look at the solutions lets try to derive these from a physical perspective.

### Newtonian derivation

Consider an elastic, circular rubber membrane  $\Omega$  of uniform density  $\rho$  that is fixed along some boundary  $\partial\Omega$  due to a uniform tension force per unit length  $T$ . We will also deal with small oscillations only, so  $\phi(x, y, t)$  is small. Let's look at the forces acting on an infinitesimal surface element  $dA = dx dy$ . We see that there are four forces acting on four edges of the element, two have magnitude  $Tdy$  and the other two have magnitude  $Tdx$ . We can model these forces to act only on the midpoint of the edges, as shown below.

The force along  $x$  acting on the frontal edge of the membrane is given by

$$F_x^1 = T \cos \theta_2 dy - T \cos \theta_1 dy \approx 0 \quad (40.1.2)$$

$$F_y^1 = T \sin \theta_2 dy - T \sin \theta_1 dy \approx T(\tan \theta_2 - \tan \theta_1) dy \quad (40.1.3)$$

where we used the small angle approximation. Now note that

$$\tan \theta_2 = \left. \frac{\partial u}{\partial x} \right|_{x+dx, y_1}, \text{ and } \tan \theta_1 = \left. \frac{\partial u}{\partial x} \right|_{x, y_2} \quad (40.1.4)$$

where  $y_1$  and  $y_2$  are midpoints of the left and right edges respectively. Therefore we see that

$$F_y^1 = T \left( \left. \frac{\partial u}{\partial x} \right|_{x+dx, y_1} - \left. \frac{\partial u}{\partial x} \right|_{x, y_2} \right) dy = T \left. \frac{\partial^2 u}{\partial x^2} \right|_{x, y_1} dx dy \quad (40.1.5)$$

Similarly for the other end we find that

$$F_y^2 = T \left( \left. \frac{\partial u}{\partial x} \right|_{x_1, y+dy} - \left. \frac{\partial u}{\partial x} \right|_{x_2, y} \right) dx = T \left. \frac{\partial^2 u}{\partial y^2} \right|_{x_1, y} dx dy \quad (40.1.6)$$

Consequently using Newton's second law we find that

$$\frac{\partial^2 u}{\partial t^2} \rho dx dy = T \nabla^2 u dx dy \quad (40.1.7)$$

which gives the 2D wave equation

$$\frac{\partial^2 u}{\partial t^2} = c^2 \nabla^2 u, \quad c = \sqrt{\frac{T}{\rho}} \quad (40.1.8)$$

We see that the phase velocity of waves propagating on the membrane is given by  $\sqrt{\frac{T}{\rho}}$ . This derivation however is tricky to generalize if the tension is non-uniform, a simple way to account for this is by taking the lagrangian approach instead.

### Variational/Lagrangian derivation

The kinetic energy of the membrane is given by

$$K = \frac{1}{2} \int_{\Omega} \mu \left( \frac{\partial \phi}{\partial t} \right)^2 dx dy \quad (40.1.9)$$

and since we are only working to first order we can approximate

$$dA = \sqrt{1 + \frac{\partial \phi}{\partial x}^2} \sqrt{1 + \frac{\partial \phi}{\partial y}^2} \approx 1 \quad (40.1.10)$$

To work out the potential energy, note that if the displaced membrane has area  $dS$  then the associated potential energy is  $T(dS - dx dy)$  and therefore

$$V = \int_{\Omega} T \sqrt{1 + \frac{\partial \phi}{\partial x}^2} \sqrt{1 + \frac{\partial \phi}{\partial y}^2} dx dy - \int_{\Omega} T dx dy \quad (40.1.11)$$

$$\approx \frac{1}{2} \int_{\Omega} T \left[ \left( \frac{\partial \phi}{\partial x} \right)^2 + \left( \frac{\partial \phi}{\partial y} \right)^2 \right] dx dy \quad (40.1.12)$$

Combining these the Lagrangian reads

$$L = \frac{1}{2} \int_{\Omega} \left\{ \mu \left( \frac{\partial \phi}{\partial t} \right)^2 - T \left[ \left( \frac{\partial \phi}{\partial x} \right)^2 + \left( \frac{\partial \phi}{\partial y} \right)^2 \right] \right\} dx dy \quad (40.1.13)$$

and thus the Lagrangian density is

$$\mathcal{L} = \frac{1}{2} \mu \left( \frac{\partial \phi}{\partial t} \right)^2 - \frac{1}{2} T \left[ \left( \frac{\partial \phi}{\partial x} \right)^2 + \left( \frac{\partial \phi}{\partial y} \right)^2 \right] \quad (40.1.14)$$

The Euler-Lagrange equations are

$$\frac{\partial}{\partial y} \left( \frac{\partial \mathcal{L}}{\partial (\partial_y \phi)} \right) + \frac{\partial}{\partial x} \left( \frac{\partial \mathcal{L}}{\partial (\partial_x \phi)} \right) + \frac{\partial}{\partial t} \left( \frac{\partial \mathcal{L}}{\partial \dot{\phi}} \right) = \frac{\partial \mathcal{L}}{\partial \phi} \quad (40.1.15)$$

and consequently we get

$$\frac{\partial^2 \phi}{\partial t^2} = c^2 \nabla^2 \phi \quad (40.1.16)$$

as desired.

## 40.2 d'Alembert's solution

Let  $\chi = x - ct$  and  $\eta = x + ct$  so that

$$\frac{\partial \phi}{\partial x} = \frac{\partial \phi}{\partial \chi} \frac{\partial \chi}{\partial x} + \frac{\partial \phi}{\partial \eta} \frac{\partial \eta}{\partial x} \quad (40.2.1)$$

$$= \frac{\partial \phi}{\partial \chi} + \frac{\partial \phi}{\partial \eta} \quad (40.2.2)$$

and

$$\frac{\partial \phi}{\partial t} = \frac{\partial \phi}{\partial \chi} \frac{\partial \chi}{\partial t} + \frac{\partial \phi}{\partial \eta} \frac{\partial \eta}{\partial t} \quad (40.2.3)$$

$$= c \frac{\partial \phi}{\partial \chi} - c \frac{\partial \phi}{\partial \eta} \quad (40.2.4)$$

Consequently

$$\frac{\partial^2 \phi}{\partial x^2} = \frac{\partial^2 \phi}{\partial \chi^2} + \frac{\partial^2 \phi}{\partial \eta^2} + 2 \frac{\partial^2 \phi}{\partial \chi \partial \eta} \quad (40.2.5)$$

and

$$\frac{\partial^2 \phi}{\partial t^2} = c^2 \frac{\partial^2 \phi}{\partial \chi^2} + c^2 \frac{\partial^2 \phi}{\partial \eta^2} - 2c^2 \frac{\partial^2 \phi}{\partial \chi \partial \eta} \quad (40.2.6)$$

from which we get the wave equation in the new coordinates

$$\frac{\partial^2 \phi}{\partial \chi \partial \eta} = 0 \implies \phi(x, t) = f(x - ct) + g(x + ct) \quad (40.2.7)$$

## 40.3 The 1D wave equation: strings

## 40.4 The 2D wave equation: membranes

## 40.5 Existence and uniqueness theorem

---

# Parabolic PDEs: Heat and Diffusion

## 41.1 Existence and uniqueness theorem

---

# **Green's functions for PDEs**

**42**

# **Part V**

# **Linear Algebra**

# Vector spaces

## 43.1 Definitions

We begin by defining a fundamental mathematical concept used in several areas of physics (most notably Quantum Mechanics), the **Vector Space**.

Classically speaking, vectors are defined as objects with both a magnitude and direction. However, as we will see soon this definition is very limited, and breaking beyond the barrier of arrows with lengths and directions will enable us to create a broader mathematical structure.

### **Definition 34.1 (Vector space axioms)**

A linear space  $V$  over a field  $\mathbb{K}$  is a collection of **vectors**  $\mathbf{v}$  over which two binary operations  $+, \cdot$  are defined, such that  $\forall \mathbf{u}, \mathbf{v}, \mathbf{z} \in V$  and  $\forall \alpha_1, \alpha_2 \in \mathbb{K}$  the following are satisfied:

- (VS1) Closure under addition:  $\mathbf{u} + \mathbf{v} \in V$
- (VS2) Closure under scalar multiplication:  $\alpha_1 \mathbf{u} \in V$
- (VS3) Commutativity of addition:  $\mathbf{u} + \mathbf{v} = \mathbf{v} + \mathbf{u}$
- (VS4) Associativity of addition:  $\mathbf{u} + \mathbf{v}$
- (VS5) Associativity of addition:  $\mathbf{u} + (\mathbf{v} + \mathbf{z}) = (\mathbf{u} + \mathbf{v}) + \mathbf{z}$
- (VS6) Associativity of scalar multiplication:  $\alpha_1(\alpha_2 \mathbf{u}) = \alpha_1 \alpha_2 \mathbf{u}$
- (VS7) Right-distributivity:  $(\alpha_1 + \alpha_2) \mathbf{x}_1 = \alpha_1 \mathbf{x}_1 + \alpha_2 \mathbf{x}_1$
- (VS8) Left-distributivity:  $\alpha_1(\mathbf{u} + \mathbf{v}) = \alpha_1 \mathbf{u} + \alpha_1 \mathbf{v}$
- (VS9) Existence of **zero vector**:  $\exists \mathbf{0} \in V$  such that  $\mathbf{u} + \mathbf{0} = \mathbf{u}$
- (VS10) Existence of inverse under addition:  $\exists (-\mathbf{u}) \in V$  such that  $\mathbf{u} + (-\mathbf{u}) = \mathbf{0}$

### **Definition 34.2 (Vector subspace)**

A vector space  $W$  is said to be a vector subspace of a vector space  $V$  if  $W \subseteq V$ , and is a proper subspace if  $W$  is neither the zero subspace  $\{\mathbf{0}\}$  nor  $V$ .

So, if a vector space satisfying VS1-VS10 is a subset of some other vector space, then it is a vector subspace. Luckily, given a subset of  $V$ , one does not necessarily have to prove that all the vector space axioms hold, since some of them hold for all subsets of  $V$ . Indeed, it turns out that only VS1 and VS2 do not necessarily hold for a subset of a vector space. All the others do.

### **Proposition 34.3 (Criteria of vector subspaces)**

A subset  $W \subseteq V$  is said to be a **vector subspace** of  $V$  over  $\mathbb{K}$  iff:

- (S1) Closure under addition:  $\forall \mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_1 + \mathbf{w}_2 \in W$

- (S2) Closure under multiplication:  $\forall \alpha \in \mathbb{K}, \forall \mathbf{w} \in W, \alpha\mathbf{w} \in \mathbb{W}$
- (S3) Identity inclusion:  $\mathbf{0}_V \in W$ , where  $\mathbf{0}_V$  is the zero of  $V$ , such that  $\mathbf{0}_V + \mathbf{w} = \mathbf{w}, \forall \mathbf{w} \in W$ .

*Proof.* We proceed by showing that all the vector space axioms hold for  $W$ :

(VS1) is equivalent to S1

(VS2) is equivalent to S2

(VS3-VS8)  $\mathbf{w}_1, \mathbf{w}_2 \in W \implies \mathbf{w}_1, \mathbf{w}_2 \in V$ . Hence, since VS3-VS8 hold for all vectors of  $V$ , they must necessarily hold for all vectors of  $W$ .

(VS9) is equivalent to S3

(VS10) implication of S2 with  $\alpha$  being the negative identity of  $\mathbb{K}$ .

■

#### Definition 34.4 (*Span*)

The span of a set of vectors  $\{\mathbf{v}_1 \dots \mathbf{v}_k\}$  is defined as the set of all their linear combinations:

$$\text{Span}(\mathbf{v}_1 \dots \mathbf{v}_k) \equiv \left\{ \sum_{i=1}^k \alpha_i \mathbf{v}_i : \forall \alpha_i \in \mathbb{K} \right\} \quad (43.1.1)$$

#### Definition 34.5 (*Linear independence*)

Let  $V$  be a vector space over  $\mathbb{K}$  and let  $\alpha_1 \dots \alpha_k \in \mathbb{K}$ . Then we say that the set of vectors  $\{\mathbf{v}_1 \dots \mathbf{v}_k\}$  are linearly independent iff:

$$\sum_{i=1}^k \alpha_i \mathbf{v}_i = \mathbf{0} \implies \alpha_i = 0, \forall 1 \leq i \leq k \quad (43.1.2)$$

Otherwise, they are said to be linearly dependent.

Firstly note that by this definition (that uses *otherwise*) a set of vectors is either linearly dependent or linearly independent, it cannot be both or neither.

Also note that it suffices for only one coefficient of a set of vectors to not be zero for linear dependence to be satisfied. Linear independence occurs only when all coefficients  $\alpha_i$  must be zero.

An immediate result is the following:

#### Proposition 34.5 (*Linear dependence of sets containing $\mathbf{0}$* )

Any set of vectors  $\{\mathbf{0}, \mathbf{v}_1, \dots, \mathbf{v}_k\} \subseteq V$  is linearly dependent.

*Proof.* The set  $\{\mathbf{0}, \mathbf{v}_1, \dots, \mathbf{v}_k\} \subseteq V$  cannot be linearly independent, since for  $\alpha \neq 0$  we may write:

$$\alpha \cdot \mathbf{0} + \sum_{i=1}^k 0 \cdot \mathbf{v}_i = \mathbf{0} \quad (43.1.3)$$

Therefore, the vectors must be linearly dependent. ■

## 43.2 Basis and dimensions

### Definition 34.6 (Basis)

A set of vectors  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\} \subseteq V$  is said to be a basis of  $V$  iff:

- (B1) they are linearly independent
- (B2) they generate  $V$ :  $\text{Span}(\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k) = V$ .

Then,  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$  are said to be **basis vectors**.

So, given a basis  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$ , it is always possible to write any vector in  $V$  as a linear combination of these basis vectors.

### Proposition 34.7 (Uniqueness of linear combination)

Let  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$  be a basis of a vector space  $V$ . Then, any vector  $\mathbf{v} \in V$  can be expressed as:

$$\mathbf{v} = \sum_{i=1}^k \alpha_i \mathbf{v}_i \quad (43.2.1)$$

where  $\alpha_i$  are uniquely determined.

*Proof.* Suppose that  $\mathbf{v}$  can be expressed as two different linear combinations:

$$\mathbf{v} = \sum_{i=1}^k \alpha_i \mathbf{v}_i = \sum_{i=1}^k \alpha'_i \mathbf{v}_i \quad (43.2.2)$$

Then:

$$\sum_{i=1}^k (\alpha_i - \alpha'_i) \mathbf{v}_i = \mathbf{0} \quad (43.2.3)$$

However, since  $\mathbf{v}_i$  are linearly independent by (B1), this implies that  $\alpha_i = \alpha'_i$ , which is a contradiction. ■

### Theorem 34.8 (Steinitz Exchange theorem)

Let  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$  be a basis of  $V$ , and let  $\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_l\} \subsetneq V$ . If  $l > k$ , then  $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_l$  are linearly dependent.

*Proof.* If  $\mathbf{w}_1 = \mathbf{0}$ , then by Proposition 34.5  $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_l$  are linearly dependent.

Suppose  $\mathbf{w}_1 \neq \mathbf{0}$ . Then, we may write:

$$\mathbf{w}_1 = \sum_{i=1}^k \alpha_i \mathbf{v}_i \iff \mathbf{v}_1 = \frac{1}{\alpha_1} (\mathbf{w}_1 - \sum_{i=2}^k \alpha_i \mathbf{v}_i) \quad (43.2.4)$$

where we assume without loss of generality that  $\exists \alpha_1 \neq 0$ , since otherwise the sets would be linearly dependent as desired. Hence:

$$\text{Span}\left(\left(\mathbf{w}_1 - \sum_{i=2}^k \alpha_i \mathbf{v}_i\right), \mathbf{v}_2, \dots, \mathbf{v}_k\right) = V \quad (43.2.5)$$

where we omit  $\frac{1}{\alpha_1}$  since it is only a constant and will be lost when writing out the linear combination. Note however that  $\sum_{i=2}^k \alpha_i \mathbf{v}_i$  has already been included in the other vectors in the span, and can therefore be ignored:

$$\text{Span}(\mathbf{w}_1, \mathbf{v}_2, \dots, \mathbf{v}_k) = V \quad (43.2.6)$$

We can repeat this process by replacing  $\mathbf{v}_j$  with:

$$\frac{1}{\alpha'_j} \left( \mathbf{w}_j - \sum_{i=1}^{j-1} \alpha'_i \mathbf{w}_i - \sum_{i=j+1}^k \alpha'_i \mathbf{v}_i \right), \forall 1 < j \leq l \quad (43.2.7)$$

so that, by similar logic to before:

$$\text{Span}\left(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_{j-1}, \mathbf{w}_j - \sum_{i=1}^{j-1} \alpha'_i \mathbf{w}_i - \sum_{i=j+1}^k \alpha'_i \mathbf{v}_i, \mathbf{v}_{j+1}, \dots, \mathbf{v}_k\right) = V \quad (43.2.8)$$

Again, all the vectors in  $\sum_{i=1}^{j-1} \alpha'_i \mathbf{w}_i$  and  $\sum_{i=j+1}^k \alpha'_i \mathbf{v}_i$  have already been included in the Span, and can be neglected. Hence, we get:

$$\text{Span}(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_j, \mathbf{v}_{j+1}, \dots, \mathbf{v}_l) \quad (43.2.9)$$

Now since  $l \geq k$ , the end result of reiterating this algorithm will be:

$$\text{Span}(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k) = V \quad (43.2.10)$$

with  $l - k$  remaining  $\mathbf{w}_i$  vectors. This means that the vectors which were left out can be expressed as a linear combination of  $\mathbf{w}_1 \dots \mathbf{w}_k$ . For example:

$$\mathbf{w}_{l-k} = \sum_i^k \beta_i \mathbf{w}_i \implies \mathbf{w}_{l-k} - \sum_i^k \beta_i \mathbf{w}_i = \mathbf{0} \quad (43.2.11)$$

implying linear dependence. ■

The contrapositive of the Exchange lemma also provides an interesting result which we shall use soon.

### Proposition 34.9 (Contrapositive of the exchange theorem)

Let  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$  be a basis of  $V$ , and let  $\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_l\} \subsetneq V$ . If  $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_l$  are linearly independent, then  $l \leq k$ .

**Definition 34.10** (*Dimension of a finite dimensional vector space*)

The dimension of a finite dimensional vector space is the number of vectors in its basis. So if  $V$  has a basis of cardinality  $n$ , then the dimension over a field  $\mathbb{K}$  is :

$$\dim_{\mathbb{K}}(V) \equiv n \quad (43.2.12)$$

It may not be immediately clear that the dimension of a vector space is well-defined. How can we know that all the bases of a vector space contain the same number of vectors?

**Theorem 34.11** (*Well-definedness of vector space dimension*)

Any two bases of a finite dimensional vector space must contain the same number of basis vectors.

*Proof.* Suppose we have two bases,  $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m\}$  and  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$ , each of cardinality  $m$  and  $n$  respectively. Consequently, by B1, they must be both linearly independent. Using Proposition 34.9 then, we have that  $m \geq k$  and  $k \geq m$ , implying that  $k = m$  as desired. ■

It is interesting to note that the dimension of a vector space can depend on the field we define it on. For example, the vector space of all  $3 \times 3$  matrices over  $\mathbb{R}$  has dimension 9, whereas over  $\mathbb{C}$  it has dimension 18.

In general, we will omit inserting the field when it is clear from the context.

**Proposition 34.11** (*Properties of finite dimensional spaces*)

The following properties are satisfied by any finite dimensional vector space  $V$ :

- (D1)  $V$  has a basis
- (D2) every linearly independent subset of  $V$  can be expanded to form a basis
- (D3) If  $n = \dim(V)$ , then any linearly independent subset  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\} \subseteq V$  is a basis of  $V$
- (D4) if  $\dim(V) = \dim(W)$  and  $V \subseteq W$ , then  $V = W$ .

*Proof.*

- (D1)  $V$  is spanned by a finitely many vectors, by definition. Hence, we can always find  $k$  vectors that span  $V$ :

$$\text{Span}(\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k) = V \quad (43.2.13)$$

If  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$  are linearly independent, then we have found a basis.

Otherwise, one of the vectors can be expressed as a linear combination of the others, and can be dropped from the span. Repeat this process until the remaining set of vectors is linearly independent.

- (D2) Suppose  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k$  do not span  $V$  (since otherwise we would already have a basis). Then,  $\exists \mathbf{v}_{k+1} \neq \text{Span}(\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k)$ , which can be added to our set of linearly independent vectors. Continue until  $V$  has been generated.
- (D3) If  $\dim(V) = n$ , then any basis of  $V$  must necessarily contain  $n$  vectors. Suppose a linearly independent set of vectors  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$  does not form a basis. Then,  $\exists \mathbf{v}_{n+1} \neq \text{Span}(\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n)$

which can be added to the linearly independent set. However, by Proposition 34.9, any linearly independent set of  $m$  vectors must satisfy  $m \leq n$ , where  $n$  is the dimension of  $V$ . This would imply  $n + 1 \leq n$ , a contradiction.

- (D4) If  $\dim V = \dim W = n$ , then there exists a basis  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\} \subseteq V$ . We can deduce using (D3) that  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\} \subseteq V$  must be a basis of  $W$  too, and thus:

$$\text{Span}(\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n) = V = W \quad (43.2.14)$$

as desired. ■

### 43.3 Operations on subspaces

In this section we will more closely inspect the properties of vector subspaces, and the operations we can apply on them, namely sums, direct sums and direct products, as well as intersections and unions.

We begin by providing an alternative, often faster way to prove that some subsets are subspaces.

#### Proposition 34.12 (Subspace criterion)

For a subset  $W \subseteq V$  to be a vector subspace of  $V$  over  $\mathbb{K}$ , we need:

$$c\mathbf{w}_1 + \mathbf{w}_2 \in W, \forall \mathbf{w}_1, \mathbf{w}_2 \in W, \forall c \in \mathbb{K} \quad (43.3.1)$$

*Proof.* We wish to prove that (S1)-(S3) are equivalent to 34.3.1. Firstly, note that:

$$(\mathbf{w}'_1 + \mathbf{w}'_2 \in W) \wedge (c\mathbf{w}'_1 \in W, \forall \mathbf{w}'_1, \mathbf{w}'_2 \in W, \forall c \in \mathbb{K}) \quad (43.3.2)$$

$$\implies c\mathbf{w}_1 + \mathbf{w}_2 \in W, \forall \mathbf{w}_1, \mathbf{w}_2 \in W, \forall c \in \mathbb{K} \quad (43.3.3)$$

if we take  $\mathbf{w}'_1 = c\mathbf{w}_1$ . Similarly:

$$c\mathbf{w}_1 + \mathbf{w}_2 \in W, \forall \mathbf{w}_1, \mathbf{w}_2 \in W, \forall c \in \mathbb{K} \quad (43.3.4)$$

$$\implies (\mathbf{w}'_1 + \mathbf{w}'_2 \in W) \wedge (c\mathbf{w}'_1 \in W, \forall \mathbf{w}'_1, \mathbf{w}'_2 \in W, \forall c \in \mathbb{K}) \quad (43.3.5)$$

if we take  $c = 1$  and  $\mathbf{w}_2 = \mathbf{w}_1$  to prove the left and right statements of 34.3.5 respectively. Also:

$$c\mathbf{w}_1 + \mathbf{w}_2 \in W, \forall \mathbf{w}_1, \mathbf{w}_2 \in W, \forall c \in \mathbb{K} \implies \mathbf{0} \in W \quad (43.3.6)$$

if we take  $c = -1$  and  $\mathbf{w}_1 = \mathbf{w}_2$  ■

#### Theorem 34.13 (Spanning subspace)

Let  $S \subseteq V$ , then  $\text{Span}(S)$  is a vector subspace of  $V$ .

*Proof.* Let  $S = \{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k\} \subseteq V$ . Then,  $\mathbf{0} = \sum_{i=1}^k 0 \cdot \mathbf{u}_i \implies \mathbf{0} \in \text{Span}(S)$ .

Now let  $\mathbf{v}_1 = \sum_{i=1}^k \alpha_i \mathbf{u}_i \in \text{Span}(W)$  and  $\mathbf{v}_2 = \sum_{i=1}^k \beta_i \mathbf{u}_i \in \text{Span}(W)$  with  $\alpha_i, \beta_i \in \mathbb{K}$ , then:

$$c\mathbf{v}_1 + \mathbf{v}_2 = c \sum_{i=1}^k \alpha_i \mathbf{u}_i + \sum_{i=1}^k \beta_i \mathbf{u}_i = \sum_{i=1}^k (c\alpha_i + \beta_i) \mathbf{u}_i = \sum_{i=1}^k \gamma_i \mathbf{u}_i \in \text{Span}(W) \quad (43.3.7)$$

where  $\gamma_i = c\alpha_i + \beta_i \in \mathbb{K}$  due to the closure of fields. Hence, by Proposition 34.12,  $\text{Span}(W)$  is a subspace of  $V$ .  $\blacksquare$

### Proposition 34.14 (Dimension of subspace)

The dimension of a vector subspace of  $V$  is always less than or equal to the dimension of  $V$ .

*Proof.* Let  $V$  be a vector space of dimension  $\dim(V) = n$ , and let  $S \subseteq V$  be a subspace of  $V$  of dimension  $\dim(W) = m$ . Let  $\mathcal{B}_V = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$  be a basis of  $V$ , and  $\mathcal{B}_W = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_m\}$  be a basis of  $W$ , then it follows that  $\mathcal{B}_W$  is linearly independent. Hence, using Proposition 34.9,  $\dim(W) = m < n = \dim(V)$ .  $\blacksquare$

### Proposition 34.15 (Union and intersection of subspaces)

Let  $V$  be a finite dimensional vector space, with subspaces  $\{W_1, W_2, \dots, W_n\}$ . Then  $\forall 1 \leq k \leq n$ :

$$W = \bigcap_{i=1}^k W_i \text{ is a subspace of } V \quad \forall 1 \leq k \leq n \quad (43.3.8)$$

and for any two subspaces  $W_1, W_2$ :

$$W_1 \cup W_2 \text{ is a subspace of } V \text{ iff } W_1 \subseteq W_2 \text{ or } W_2 \subseteq W_1 \quad (43.3.9)$$

*Proof.* We begin by proving that  $\bigcap_{i=1}^k W_i$  is a subspace of  $V$ ,  $\forall 1 \leq k \leq n$ .  $\mathbf{0} \in W_i$  for all  $i$  by the subgroup axioms.

Moreover, if  $\mathbf{u}, \mathbf{v} \in W$ , then  $\mathbf{u}, \mathbf{v} \in W_i$  for all  $i$ . By proposition 34.12 then:

$$\alpha\mathbf{u} + \mathbf{v} \in W_i, \forall i \implies \alpha\mathbf{u} + \mathbf{v} \in W \quad (43.3.10)$$

as required. Hence, the subgroup criteria are met, and  $W = \bigcap_{i=1}^k W_i$  is a subspace of  $V$ .

Next we prove that  $W_1 \cup W_2$  is a subspace of  $V$  iff  $W_1 \subseteq W_2$  or  $W_2 \subseteq W_1$ .

( $\implies$ ) We proceed by contradiction. Suppose  $W_1 \cup W_2$  is a subspace of  $V$  and suppose  $W_1 \not\subseteq V$  and  $W_2 \not\subseteq V$ . Then  $\exists \mathbf{w}_1 \in W_1 \setminus W_2$  and  $\exists \mathbf{w}_2 \in W_2 \setminus W_1$ . Therefore, by the closure axiom of groups:

$$\mathbf{w}_1 + \mathbf{w}_2 \in V \quad (43.3.11)$$

Now suppose that  $\mathbf{w}_1 + \mathbf{w}_2 \in W_1$ . Then:

$$(-\mathbf{w}_1) + (\mathbf{w}_1 + \mathbf{w}_2) = \mathbf{w}_2 \in W_1 \quad (43.3.12)$$

which is a contradiction. So  $\mathbf{w}_1 + \mathbf{w}_2 \notin W_1$ . Similarly, suppose that  $\mathbf{w}_1 + \mathbf{w}_2 \in W_2$ . Then this would imply that:

$$(-\mathbf{w}_2) + (\mathbf{w}_1 + \mathbf{w}_2) = \mathbf{w}_1 \in W_1 \quad (43.3.13)$$

which is a contradiction. Thus, we conclude that  $\mathbf{w}_1 + \mathbf{w}_2 \notin W_1 \cup W_2$ . However, this violates the subspace criteria in Proposition 34.3. Hence, we must require  $W_1 \subseteq W_2$  or  $W_2 \subseteq W_1$ .

( $\Leftarrow$ ) Suppose that  $W_1 \subseteq W_2$  or  $W_2 \subseteq W_1$ . Then  $W_1 \cup W_2 = W_2$  or  $W_1 \cup W_2 = W_1$  respectively, and since both  $W_1, W_2$  are both subspaces of  $W$  then it follows that  $W_1 \cup W_2$  is in either cases a subspace of  $W$ . ■

### Definition 34.16 (Cosets, quotient spaces, sums of spaces)

Let  $V$  be a vector space and let  $W$  be a vector subspace of  $V$ . Then, a coset of  $V$  is:

$$\mathbf{v} + W \equiv \{\mathbf{v} + \mathbf{w} : \forall \mathbf{v} \in V\} \quad (43.3.14)$$

where  $\mathbf{w} \in W$ . The set of all cosets of  $V$  in  $W$  is called the quotient space of  $W$  modulo  $V$ :

$$V/W \equiv \{\mathbf{v} + W : \forall \mathbf{v} \in V\} \quad (43.3.15)$$

Finally, the sum of two vector subspaces  $U, W$  of  $V$  is defined as:

$$U + W \equiv \{\mathbf{u} + \mathbf{w} : \forall \mathbf{u} \in U, \mathbf{w} \in W\} \quad (43.3.16)$$

### Theorem 34.17 (Dimension of sum of spaces)

Let  $U, W$  be subspaces of a finite dimensional space  $V$ . Then:

$$\dim(U + W) = \dim(U) + \dim(W) - \dim(U \cap W). \quad (43.3.17)$$

*Proof.* We firstly prove that  $\dim(U + W) = \dim(U) + \dim(W) - \dim(U \cap W)$ .

Let  $\mathcal{B}_{\cap} = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r\}$  be a basis for  $U \cap W$  so that  $\dim(U \cap W) = r$ . From Proposition 34.15, we must have that  $U \cup W$  is a subspace of both  $U$  and  $W$ , and consequently (D2) of Proposition 34.11 implies that  $\mathcal{B}_{\cap}$  can be extended to form a basis of  $U$  and  $W$ .

Hence suppose that  $\mathcal{B}_U = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r, \mathbf{u}_1, \dots, \mathbf{u}_m\}$  and  $\mathcal{B}_W = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r, \mathbf{w}_1, \dots, \mathbf{w}_n\}$  are bases for  $U$  and  $W$  respectively.

Now let  $\mathbf{v}' \in U + W$ , then it may be expressed as:

$$\mathbf{v}' = \mathbf{u}' + \mathbf{w}' \quad (43.3.18)$$

$$= \underbrace{\left( \sum_{i=1}^r \alpha_i \mathbf{v}_i + \sum_{i=1}^m \alpha'_i \mathbf{u}_i \right)}_{\mathbf{u}'} + \underbrace{\left( \sum_{i=1}^r \beta_i \mathbf{v}_i + \sum_{i=1}^n \beta'_i \mathbf{w}_i \right)}_{\mathbf{w}'} \quad (43.3.19)$$

for some  $\mathbf{u}' \in U, \mathbf{w}' \in W$ . We can rearrange the above equation:

$$\mathbf{v}' = \left( \sum_{i=1}^r \alpha_i \mathbf{v}_i + \sum_{i=1}^m \alpha'_i \mathbf{u}_i \right) + \left( \sum_{i=1}^r \beta_i \mathbf{v}_i + \sum_{i=1}^n \beta'_i \mathbf{w}_i \right) \quad (43.3.20)$$

$$= \sum_{i=1}^r \gamma_i \mathbf{v}_i + \sum_{i=1}^m \alpha'_i \mathbf{u}_i + \sum_{i=1}^n \beta'_i \mathbf{w}_i \quad (43.3.21)$$

where  $\gamma_i = \alpha_i + \beta_i$ . Therefore:

$$\text{Span}(\mathbf{v}_1, \dots, \mathbf{v}_r, \mathbf{u}_1, \dots, \mathbf{u}_m, \mathbf{w}_1, \dots, \mathbf{w}_n) = U \quad (43.3.22)$$

Also, note that  $\{\mathbf{v}_1, \dots, \mathbf{v}_r, \mathbf{u}_1, \dots, \mathbf{u}_m, \mathbf{w}_1, \dots, \mathbf{w}_n\}$  is linearly independent. Indeed, suppose:

$$\sum_{i=1}^r \alpha_i \mathbf{v}_i + \sum_{i=1}^m \beta_i \mathbf{u}_i + \sum_{i=1}^n \gamma_i \mathbf{w}_i = \mathbf{0} \implies \sum_{i=1}^r \alpha_i \mathbf{v}_i + \sum_{i=1}^m \beta_i \mathbf{u}_i = - \sum_{i=1}^n \gamma_i \mathbf{w}_i \quad (43.3.23)$$

for some  $\alpha_i, \beta_i, \gamma_i \in \mathbb{K}$ . The above belongs to  $U$ , since  $\{\mathbf{v}_1, \dots, \mathbf{v}_r, \mathbf{u}_1, \dots, \mathbf{u}_m\}$  is a basis of  $U$ . Similarly, it also belongs to  $W$ . Therefore:

$$\sum_{i=1}^r \alpha_i \mathbf{v}_i + \sum_{i=1}^m \beta_i \mathbf{u}_i = - \sum_{i=1}^n \gamma_i \mathbf{w}_i \in U \cap W \quad (43.3.24)$$

and can therefore be written as:

$$\sum_{i=1}^r \alpha_i \mathbf{v}_i + \sum_{i=1}^m \beta_i \mathbf{u}_i = - \sum_{i=1}^n \gamma_i \mathbf{w}_i = \sum_{i=1}^r c_i \mathbf{v}_i \quad (43.3.25)$$

$$\implies \sum_{i=1}^m \beta_i \mathbf{u}_i = \sum_{i=1}^r d_i \mathbf{v}_i \text{ and } - \sum_{i=1}^n \gamma_i \mathbf{w}_i = \sum_{i=1}^r c_i \mathbf{v}_i \quad (43.3.26)$$

where  $d_i = c_i - \alpha_i$ . Recall that  $\mathbf{u}_i$  and  $\mathbf{v}_i$  must be linearly independent, since they form a basis of  $U$ . Thus we obtain  $\beta_i = 0$  and  $d_i = 0 \implies c_i = \alpha_i$ .

Similarly, we require that  $\mathbf{w}_i$  and  $\mathbf{v}_i$  be linearly independent since they form a basis for  $W$ . Consequently  $\gamma_i = c_i = 0 \implies \alpha_i = 0$ . Linear dependence is thus satisfied.

So we may claim that  $\{\mathbf{v}_1, \dots, \mathbf{v}_r, \mathbf{u}_1, \dots, \mathbf{u}_m, \mathbf{w}_1, \dots, \mathbf{w}_n\}$  is a basis of  $U + W$ . It follows that:

$$\dim(U + W) = (r + m) + (r + n) - r = \dim(U) + \dim(W) - \dim(U \cap W) \quad (43.3.27)$$

as was desired. ■

### Theorem 34.18 (Dimension of quotient spaces)

Let  $U, W$  be subspaces of a finite dimensional space  $V$ . Then:

$$\dim(V/W) = \dim(V) - \dim(W) \quad (43.3.28)$$

*Proof.* Let  $\mathcal{B}_W = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_n\}$  and  $\mathcal{B}_V = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_n, \mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m\}$  be bases for  $W$  and  $V$

respectively. Let  $\mathbf{v} + W \in V/W$ , where  $\mathbf{v} \in V$ , then it may be expressed as:

$$\mathbf{v} + W = \sum_{i=1}^n \alpha_i \mathbf{w}_i + \sum_{i=1}^m \beta_i \mathbf{v}_i + V \quad (43.3.29)$$

$$= \sum_{i=1}^n \alpha_i (\mathbf{w}_i + W) + \sum_{i=1}^m \beta_i (\mathbf{v}_i + W) \quad (43.3.30)$$

$$= \sum_{i=1}^m \beta_i (\mathbf{v}_i + W) \quad (43.3.31)$$

$$\in \text{Span}(\mathbf{v}_1 + W, \mathbf{v}_2 + W, \dots, \mathbf{v}_m + W) \quad (43.3.32)$$

where  $\sum_{i=1}^n \alpha_i (\mathbf{w}_i + W)$  disappears since it is equal to  $W$ , and can be reabsorbed into the second sum.

Also, we note that  $\{\mathbf{v}_1 + W, \mathbf{v}_2 + W, \dots, \mathbf{v}_m + W\}$  is linearly independent. Indeed:

$$\sum_{i=1}^m \alpha_i (\mathbf{v}_i + W) = \mathbf{0} + W \equiv W \quad (43.3.33)$$

implies:

$$\sum_{i=1}^m \alpha_i \mathbf{v}_i \in W \implies \sum_{i=1}^m \alpha_i \mathbf{v}_i = \sum_{i=1}^n \beta_i \mathbf{w}_i \quad (43.3.34)$$

But  $\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_n, \mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m\}$  is linearly independent since it forms a basis for  $V$ . Consequently,  $\alpha_i = 0$ , thus proving linear independence of  $\{\mathbf{v}_1 + W, \mathbf{v}_2 + W, \dots, \mathbf{v}_m + W\}$ .

So we may claim that  $\{\mathbf{v}_1 + W, \mathbf{v}_2 + W, \dots, \mathbf{v}_m + W\}$  is a basis for  $V/W$ . It follows immediately that:

$$\dim(V/W) = (m+n) - n = \dim(V) - \dim(W) \quad (43.3.35)$$

as required. ■

### Definition 34.19 (Direct sums)

Let  $U, W$  be subspaces of a finite dimensional space  $V$ . Then, we say that  $V$  is the (internal) direct sum of  $U$  and  $W$  if:

$$(DS1) \quad U + W = V$$

$$(DS2) \quad U \cap W = \{\mathbf{0}\}$$

We then say that  $U$  and  $W$  are complementary spaces, and denote the direct sum as:

$$U \oplus W = V \quad (43.3.36)$$

Instead, given two arbitrary vector spaces  $V_1, V_2$  then their (external) direct sum is defined as:

$$V_1 \oplus V_2 \equiv \{(\mathbf{v}_1, \mathbf{v}_2) : \mathbf{v}_1 \in V_1, \mathbf{v}_2 \in V_2\} \quad (43.3.37)$$

An immediate consequence of this definition is that:

$$\dim(U \oplus W) = \dim(U) + \dim(W) \quad (43.3.38)$$

or more generally that:

$$\dim\left(\bigoplus_{i=1}^n W_i\right) = \sum_{i=1}^n \dim(W)_i \quad (43.3.39)$$

---

# Euclidean geometry in $\mathbb{R}^3$

44

---

# **Matrix algebra**

**45**

# Linear transformations

## 46.1 What is a map?

We begin by restating some common results on maps that you should be familiar with.

### Definition 46.1 (Map, Domain, Image and Kernel)

A map between two sets  $X$  and  $Y$  assigns to each  $x \in X$  some  $y = f(x) \in Y$  referred to as the **image of  $x$  under  $f$** :

$$f : X \rightarrow Y \quad (46.1.1)$$

$$x \mapsto f(x) \quad (46.1.2)$$

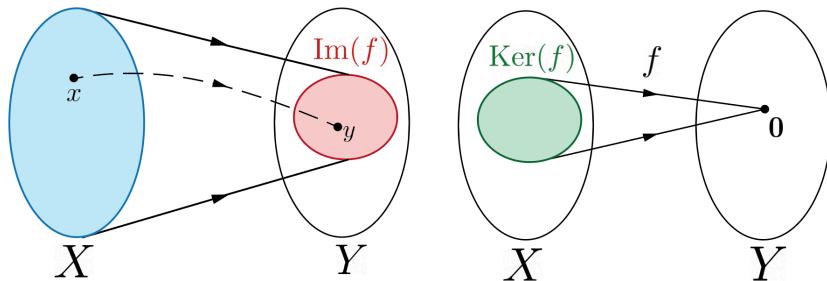
Here  $X$  is called the **domain** of  $f$ , denoted  $\text{dom}(f)$ . Instead, the set:

$$\text{Im}(f) = \{f(x) : x \in X\} \subseteq Y \quad (46.1.3)$$

is called the **image** of  $f$ . If  $f$  is a homomorphism, then the set:

$$\text{Ker}(f) = \{x \in X : f(x) = 0\} \quad (46.1.4)$$

is called the **kernel** of  $f$ .



**Figure 46.1.** Map from  $X$  to  $Y$ .

**Definition 46.2 (Injectivity, Surjectivity, Bijectivity)** Recall that a function  $f : X \rightarrow Y$  is said to be:

- (i) Injective: if every element of  $Y$  is the image of at most one element of  $X$  i.e.  $f(x) = f(x') \implies x = x', \forall x, x' \in X$ .

- (ii) Surjective: if every element of  $Y$  is the image of at least one element of  $X$  i.e.  $\text{Im}(f) = Y$ .
- (iii) Bijective: if it is both surjective and injective.

Recall that another way to state surjectivity is that if  $\forall y \in Y, \exists x \in X \text{ s.t. } f(x) = y$ . This latter definition is equivalent to  $y \in Y \implies y \in \text{Im}(f)$  so that  $Y \subseteq \text{Im}(f)$ . But  $\text{Im}(f) \subseteq Y$  so  $Y = \text{Im}(f)$ .

As we saw in Group theory, it is possible to compose different elements of a dihedral group. Similarly, one can also compose maps.

### Definition 46.3 (Map composition)

Given two maps  $f : X \rightarrow Y$  and  $g : Y \rightarrow Z$  then their composite map is defined as:

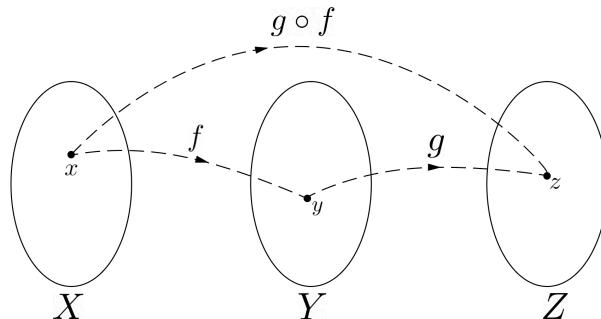
$$g \circ f : X \rightarrow Z \quad (46.1.5)$$

$$x \mapsto g(f(x)) \quad (46.1.6)$$

shown in the form of a commutative diagram below:

$$\begin{array}{ccc} A & \xrightarrow{f} & B \\ & \searrow g \circ f & \downarrow g \\ & C & \end{array} \quad (46.1.7)$$

It is very important that  $\text{Im}(f) \subseteq \text{dom } g$  since otherwise it would not be possible to evaluate  $g(f(x))$ . We can interpret the composite of two maps as another map which "jumps over" and bypasses  $Y$  as shown below:



**Figure 46.2.** Composite map  $g \circ f$

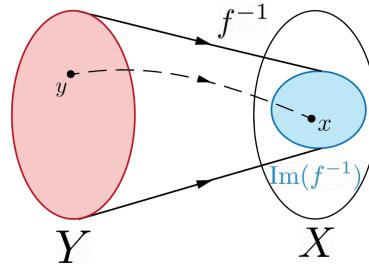
Suppose we wish to find two maps such that when composed together, they map back every element to itself. Such two maps would then be inverses of each other, so that when composed they give the identity transformation. We define them more rigorously below.

### Definition 46.4 (Inverse and identity maps)

The identity map  $\text{id}_X : X \rightarrow X$  maps all elements of  $X$  to themselves.

Given a map  $f : X \rightarrow Y$ , then  $g : Y \rightarrow X$  is its inverse is:

$$(g \circ f) = \text{id}_X \text{ and } (f \circ g) = \text{id}_Y \quad (46.1.8)$$



**Figure 46.3.** Inverse map  $f^{-1}$

It is not always a given that a map  $f$  has an inverse. We can use the following theorem to determine which map  $f$ .

**Theorem 46.4 (Bijectivity and invertibility)**

The map  $f : X \rightarrow Y$  has an inverse iff  $f$  is bijective, and this inverse is unique.

*Proof.*

( $\implies$ ) Suppose  $f$  has an inverse  $g : Y \rightarrow X$ . Then:

$$g \circ f = \text{id}_X \text{ and } f \circ g = \text{id}_Y \quad (46.1.9)$$

If  $f(x) = f(x')$  then  $g(f(x)) = \text{id}_X(x) = x$  and  $g(f(x')) = \text{id}_X(x') = x'$  so that  $x = x'$ , thus implying injectivity.

Let  $y \in Y$ . So  $g(y) = x$  for some  $x$ , hence  $f(g(y)) = \text{id}_Y(y) = y = f(x)$ . Hence for any  $y$  there exists some  $x$  so that  $y = f(x)$ , giving surjectivity.

( $\impliedby$ ) Suppose  $f$  is bijective, so  $\forall y \in Y$ , there exists  $x$  s.t.  $f(x) = y$ . If we define  $g(y) = x$  then for all  $y \in Y$ :

$$(f \circ g)(y) = f(x) = y \text{ and } (g \circ f)(x) = g(y) = x \quad (46.1.10)$$

as required,  $g$  is its inverse.

Suppose  $f$  is invertible with two inverses  $g, h$ . Then:

$$(f \circ g)(x) = (f \circ h)(x) \forall x \in X \quad (46.1.11)$$

then composing with  $g$   $g(x) = h(x) \implies g = h$  since this equality holds for any  $x$ . ■

One can see this more intuitively. Indeed, if  $f$  is surjective, then there may be some elements in  $Y$  that are not mapped. Hence we cannot find an  $x$  so that  $f^{-1}(y) = x$ , which is clearly a problem since all  $y$  must get mapped by  $f^{-1}$ . If  $f$  is injective, then there may be several elements  $x$  mapping to the same  $y$ , so that  $f^{-1}(y)$  is no longer well-defined. If however it is bijective, then every element

of  $x$  gets mapped exactly once to some  $y$ , and all  $y$  are a map of some  $x$ , so invertibility is easily satisfied.

### Proposition 46.6 (Important inverses)

Suppose  $f, g$  are bijective maps. Then  $f^{-1}, g^{-1}$  and  $f \circ g$  are bijective, with inverses:

$$(f \circ g)^{-1} = g^{-1} \circ f^{-1} \text{ and } (f^{-1})^{-1} = f \quad (46.1.12)$$

*Proof.* The first follows from:

$$(f \circ g)^{-1} \circ (f \circ g) = \text{id} \quad (46.1.13)$$

and:

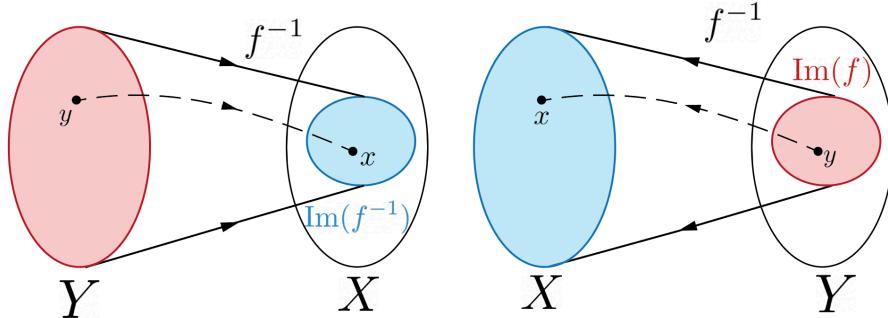
$$(g^{-1} \circ f^{-1}) \circ (f \circ g) = g^{-1} \circ \text{id} \circ g = \text{id} \quad (46.1.14)$$

so  $(f \circ g)^{-1}$  and  $g^{-1} \circ f^{-1}$  are inverses of  $f \circ g$ . But inverses are unique, hence the two must be equal.

Similarly:

$$(f^{-1})^{-1} \circ f^{-1} = \text{id} \text{ and } f \circ f^{-1} = \text{id} \quad (46.1.15)$$

hence by the same logic as before  $f = (f^{-1})^{-1}$  as desired. ■



**Figure 46.4.** Left: diagram of  $f^{-1}$  treating it following definition 46.1 (so all elements of  $Y$  get mapped). Right: diagram of  $f^{-1}$  treating it as the map that undoes  $f$

## 46.2 What is a linear map?

### Definition 46.7 (Linear map)

A **linear map** is a map  $f : V \rightarrow W$  between two vector spaces  $V$  and  $W$  over a field  $\mathbb{F}$  such that it preserves addition and scalar multiplication:

$$(L1) \quad f(\mathbf{v}_1 + \mathbf{v}_2) = f(\mathbf{v}_1) + f(\mathbf{v}_2), \forall \mathbf{v}_1, \mathbf{v}_2 \in V$$

$$(L2) \quad f(\alpha \mathbf{v}_1) = \alpha f(\mathbf{v}_1), \forall \mathbf{v}_1 \in V, \forall \alpha \in \mathbb{F}$$

The set of all linear maps from  $V$  to  $W$  is denoted  $\text{Hom}(V, W)$

### Proposition 46.8 (Properties of linear maps)

For any linear maps  $f : V \rightarrow W$ :

- (i)  $f(\mathbf{0}) = \mathbf{0}$ , that is,  $f$  fixes the zero vector.
- (ii)  $\text{Ker}(f)$  is a subspace of  $V$  and  $\text{Im}(f)$  is a subspace of  $W$ .
- (iii)  $f$  is surjective  $\iff \text{Im}(f) = W \iff \dim \text{Im}(f) = \dim W$ .
- (iv)  $f$  is injective  $\iff \text{Ker}(f) = \{\mathbf{0}\} \iff \dim \text{Ker}(f) = 0$ .
- (v)  $\alpha f$  is linear for  $\alpha \in \mathbb{K}$ .
- (vi) if  $g$  is linear then  $f + g$  is linear.
- (vii) if  $g$  is linear then  $f \circ g$  is linear.

*Proof.*

- (i) We have already shown that  $\mathbf{0} \in \text{Ker}(f)$  in point (i).  $f(\mathbf{0}) = f(0 \cdot \mathbf{0}) = 0f(\mathbf{0}) = \mathbf{0}$
- (ii) Let  $\mathbf{v}_1, \mathbf{v}_2 \in \text{Ker}(f)$ . Then  $f(\alpha\mathbf{v}_1 + \mathbf{v}_2) = \alpha f(\mathbf{v}_1) + f(\mathbf{v}_2) = \mathbf{0} \implies \alpha\mathbf{v}_1 + \mathbf{v}_2 \in \text{Ker}(f)$ . Proposition 34.12 then ensures that  $\text{Ker}(f)$  is a subspace of  $V$ .

We have already shown that  $\mathbf{0} \in \text{Im}(f)$  in point (i). Let  $\mathbf{w}_1, \mathbf{w}_2 \in \text{Im}(f)$ , so that  $\exists \mathbf{v}_1, \mathbf{v}_2 \in V$  such that  $\mathbf{w}_i = f(\mathbf{v}_i)$ . Then  $\alpha\mathbf{w}_1 + \mathbf{w}_2 = \alpha f(\mathbf{v}_1) + f(\mathbf{v}_2) = f(\alpha\mathbf{v}_1 + \mathbf{v}_2) \in \text{Im}(f)$ . Proposition 34.12 then ensures that  $\text{Im}(f)$  is a subspace of  $W$ .

- (iii) If  $f$  is surjective, then  $\forall \mathbf{w} \in W, \exists \mathbf{v} \in V$  such that  $\mathbf{w} = f(\mathbf{v}) \in \text{Im}(f)$ . So  $W \subseteq \text{Im}(f)$ , and  $\text{Im}(f) \subseteq W$  implying  $\text{Im}(f) \subseteq W$ .

If instead  $\text{Im}(f) = W$ , then  $W \subseteq \text{Im}(f)$ , and  $\text{Im}(f) \subseteq W$ . Hence  $\forall \mathbf{w} \in W, \exists \mathbf{v} \in V$  such that  $\mathbf{w} = f(\mathbf{v}) \in \text{Im}(f)$ , implying that  $f$  is surjective. It follows then that the two spaces have the same dimension.

- (iv) Suppose  $f$  is injective, and let  $\mathbf{v}_1, \mathbf{0} \in \text{Ker}(f)$ . Then  $f(\mathbf{v}_1 - \mathbf{0}) = \mathbf{0}$  since  $\text{Ker}(f)$  is a subspace. Consequently  $f(\mathbf{v}_1) = f(\mathbf{0})$  and so  $\mathbf{v}_1 = \mathbf{0}$ . So  $\text{Ker}(f) = \{\mathbf{0}\}$ .

Suppose  $\text{Ker}(f) = \{\mathbf{0}\}$ , and let  $\mathbf{v}_1, \mathbf{v}_2 \in \text{Ker}(f)$ . Then,  $f(\mathbf{v}_1) = f(\mathbf{v}_2) \implies f(\mathbf{v}_1 - \mathbf{v}_2) = \mathbf{0}$ . Hence  $\mathbf{v}_1 - \mathbf{v}_2 \in \text{Ker}(f)$  and so  $\mathbf{v}_1 = \mathbf{v}_2$ .

It follows that  $\dim \text{Ker}(f) = \dim \{\mathbf{0}\} = 0$ .

■

### Definition 46.9 (Rank and nullity)

The dimension of the image of a linear map  $f$  is called its **rank**:

$$\text{rk}(f) = \dim \text{Im}(f) \quad (46.2.1)$$

and the dimension of the kernel of a linear map  $f$  is called its **nullity**:

$$\text{null}(f) = \dim \text{Ker}(f) \quad (46.2.2)$$

### Theorem 46.10 (Rank-nullity theorem)

For any linear map  $f : V \rightarrow W$ :

$$\text{null}(f) + \text{rk}(f) = \dim(V) \quad (46.2.3)$$

*Proof.* Let  $n = \dim V$  and  $k = \text{null}(f)$ , and let  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$  be a basis for  $\text{Ker}(f)$ , which we can complete to  $\{\mathbf{v}_1, \dots, \mathbf{v}_k, \mathbf{v}_{k+1}, \dots, \mathbf{v}_n\}$  to form a basis for  $V$ . We wish to show that  $f(\mathbf{v}_{k+1} \dots \mathbf{v}_n)$  is a basis of  $\text{Im}(f)$ .

Let us firstly prove that  $\text{Im}(f) = \text{Span}(f(\mathbf{v}_{k+1} \dots \mathbf{v}_n))$ . Indeed consider  $\mathbf{w} \in \text{Im}(f)$ . Then there exists some  $\mathbf{v}$  such that:

$$\mathbf{w} = f(\mathbf{v}) = f\left(\sum_{i=1}^n \alpha_i \mathbf{v}_i\right) = \sum_{i=1}^n \alpha_i f(\mathbf{v}_i) \quad (46.2.4)$$

However, for  $i \leq k$  we have that  $f(\mathbf{v}_i) = \mathbf{0}$  and so  $\text{Im}(f) = \text{Span}(f(\mathbf{v}_{k+1} \dots \mathbf{v}_n))$ .

Now let us prove that  $\{\mathbf{v}_{k+1} \dots \mathbf{v}_n\}$  are linearly independent. Consider:

$$\sum_{i=k+1}^n \alpha_i f(\mathbf{v}_i) = \mathbf{0} \implies f\left(\sum_{i=k+1}^n \alpha_i \mathbf{v}_i\right) = \mathbf{0} \quad (46.2.5)$$

and consequently  $\sum_{i=k+1}^n \alpha_i \mathbf{v}_i \in \text{Ker}(f)$ . Using the fact that  $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$  is a basis of  $\text{Ker}(f)$  we find that:

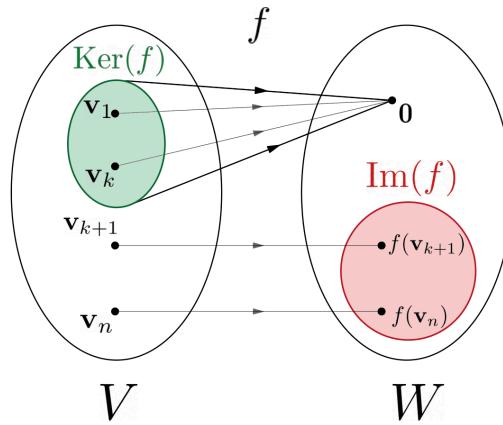
$$\sum_{i=k+1}^n \alpha_i \mathbf{v}_i = \sum_{i=1}^k \beta_i \mathbf{v}_i \implies \sum_{i=1}^n \gamma_i \mathbf{v}_i = \mathbf{0} \quad (46.2.6)$$

where  $\gamma_i = \alpha_i$  for  $k < i \leq n$ , and  $\gamma_i = \beta_i$  for  $1 \leq i \leq k$ . However, we know that  $\{\mathbf{v}_i, \dots, \mathbf{v}_n\}$  forms a basis, and must therefore be linearly independent. Hence  $\gamma_i = 0$  and so  $\alpha_i = 0$  as well.

Since  $\{\mathbf{v}_{k+1} \dots \mathbf{v}_n\}$  both generate  $\text{Im}(f)$  and are also linearly independent, they must form a basis of  $\text{Im}(f)$ . So  $\text{rk}(f) = n - k$  and consequently:

$$\dim V = n = k + (n - k) = \text{null}(f) + \text{rk}(f) \quad (46.2.7)$$

■



**Figure 46.5.** Visual depiction of the rank-nullity theorem

Throughout this proof we have also demonstrated how to construct a basis for  $\text{Im}(f)$ , which is to simply take the image of the basis of the domain.

It is then easy to see that for  $f$  to be invertible/bijective, then we need  $\text{Ker}(f)$  to only contain  $\mathbf{0}$ , and no other vectors. In other words, we need  $\text{null}(f) = 0$  and hence  $\text{rk}(f) = \dim W$ .

Let's justify this more rigorously.

**Proposition 46.11 (Consequence of rank-nullity)**

Let  $f : V \rightarrow W$  be a linear transformation with  $n = \dim V$ . Then:

- (i)  $f$  is bijective  $\implies \dim V = \dim W$ .
- (ii)  $f$  is bijective  $\iff \text{null}(f) = 0 \iff \text{rk}(f) = n$ .
- (iii) If it exists,  $f^{-1} : W \rightarrow V$  is linear.

*Proof.*

- (i) if  $f$  is bijective, then we showed in Proposition 46.8 that  $\text{null}(f) = 0$  and  $\text{rk}(f) = \dim W$ .

Using the rank nullity theorem:

$$\dim V = \text{rk}(f) = \dim W \quad (46.2.8)$$

as required.

- (ii) Suppose  $\dim V = \dim W = n$  (which follows from bijectivity), then from the rank-nullity theorem:

$$n = \dim V = \text{rk}(f) + \text{null}(f) = \dim W \quad (46.2.9)$$

and since  $f$  is bijective then:

$$\text{null}(f) = 0 \implies \text{rk}(f) = n \quad (46.2.10)$$

- (iii) Set  $\mathbf{w}_1 = f(\mathbf{v}_1)$  and  $\mathbf{w}_2 = f(\mathbf{v}_2)$ . Then:

$$f^{-1}(\alpha\mathbf{w}_1 + \mathbf{w}_2) = f^{-1}(\alpha f(\mathbf{v}_1) + f(\mathbf{v}_2)) \quad (46.2.11)$$

$$= f^{-1}(f(\alpha\mathbf{v}_1 + \mathbf{v}_2)) \quad (46.2.12)$$

$$= \alpha\mathbf{v}_1 + \mathbf{v}_2 \quad (46.2.13)$$

$$= \alpha f^{-1}(\mathbf{v}_1) + f^{-1}(\mathbf{v}_2) \quad (46.2.14)$$

■

**Proposition 46.12 (Linear map given basis of domain)**

Let  $V$  and  $W$  be vector spaces, with  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  as a basis of  $V$ . For  $\mathbf{w}_1, \dots, \mathbf{w}_n \in W$  there exists one linear map  $f : V \rightarrow W$  such that  $f(\mathbf{v}_i) = \mathbf{w}_i$ .

*Proof.* One such linear map is:

$$f(\mathbf{v}) = \sum \alpha_i \mathbf{w}_i \quad (46.2.15)$$

where  $\alpha_i$  are the coefficients of  $\mathbf{v}_i$ :

$$\mathbf{v} = \sum \alpha_i \mathbf{v}_i \quad (46.2.16)$$

It is linear since:

$$f(\beta\mathbf{v} + \mathbf{v}') = f\left(\sum \beta\alpha_i \mathbf{v}_i + \sum \alpha'_i \mathbf{v}_i\right) \quad (46.2.17)$$

$$= f\left(\sum (\beta\alpha_i + \alpha'_i) \mathbf{v}_i\right) \quad (46.2.18)$$

$$= \sum (\beta\alpha_i + \alpha'_i) \mathbf{w}_i \quad (46.2.19)$$

$$= \beta \sum \alpha_i \mathbf{v}_i + \sum \alpha'_i \mathbf{v}_i \quad (46.2.20)$$

Obviously, it maps  $\mathbf{v}_i \mapsto \mathbf{w}_i$ . Finally, it is unique, since if  $g$  was also a linear map with such properties then for all  $\mathbf{v} \in V$ :

$$g(\mathbf{v}) = \sum \alpha_i \mathbf{w}_i = f(\mathbf{v}) \implies f = g \quad (46.2.21)$$

as required. ■

## 46.3 Isomorphisms

### Definition 46.13 (Isomorphism)

Two vector spaces  $V$  and  $W$  are said to be **isomorphic** whenever we can find an invertible linear map  $f : V \rightarrow W$ , called an **isomorphism** of  $V$  onto  $W$ .

It follows immediately from (i) of Proposition 46.11 that if two spaces are isomorphic, then they must have the same dimension. It turns out that the converse is also true.

### Theorem 46.14 (Equivalent statement of isomorphicity)

Two finite dimensional vector spaces  $V$  and  $W$  are isomorphic iff  $\dim V = \dim W$ .

*Proof.* We have already proven  $\implies$ .

Now suppose that  $\dim V = \dim W$  and let  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  and  $\{\mathbf{w}_1, \dots, \mathbf{w}_n\}$  be bases for  $V$  and  $W$  respectively. Then we know from Proposition 46.12 that there exists a linear map  $f$  mapping the basis of  $V$  to the basis of  $W$ . Let us show that it is invertible, that is, bijective. Firstly, it is injective, since:

$$f(\mathbf{v}) = f(\mathbf{v}') \implies \sum \alpha_i \mathbf{v}_i = \sum \beta_i \mathbf{v}_i \implies \alpha_i = \beta_i \quad (46.3.1)$$

and so  $\mathbf{v} = \mathbf{v}'$ . It is also surjective since  $\text{Im}(f) = \{f(\mathbf{v}_i) : i \in \mathbb{N}\} = \{\mathbf{w}_i : i \in \mathbb{N}\}$ . Hence  $f$  is invertible, provided  $\dim V = \dim W$ . ■

## 46.4 Linear maps and matrices

### Definition 46.15 (Coordinate map)

For a vector space  $V$  over  $\mathbb{K}$ , endowed with a basis  $\beta = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ , then we define the

**coordinate map**  $\varphi : V \mapsto \mathbb{K}^n$  by:

$$\varphi_\beta(\mathbf{v}) = \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_n \end{pmatrix} = [\mathbf{v}]_\beta, \forall \mathbf{v} \in V \quad (46.4.1)$$

where  $\mathbf{v} = \sum_{i=1}^n \alpha_i \mathbf{v}_i$ . Here,  $[\mathbf{v}]_\beta$  is known as the **coordinate vector of  $\mathbf{v}$  relative to  $\beta$** .

The inverse of  $\varphi_\beta$  is  $\phi_\beta$  which given a basis  $\beta$  associates to a list of scalars  $\alpha_1, \alpha_2, \dots, \alpha_n$  a vector  $\mathbf{v}$ .

Note that  $\varphi$  and  $\phi$  are clearly linear. Moreover,  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  forms a basis of  $V$  then  $\dim V = n = \dim \mathbb{K}^n$  implying that  $f$  is bijective and invertible. Consequently  $\varphi$  is an isomorphism, giving us the next insightful result:

### Proposition 46.16 (Isomorphisms of $V$ onto $\mathbb{K}^n$ )

Every  $n$  dimensional vector space  $V$  is isomorphic to  $\mathbb{F}^n$  through a coordinate map  $\phi$ .

Therefore, given a vector space  $V$  and a basis  $\alpha$  then the maps  $\phi_\alpha$  and  $\varphi_\alpha$ :

$$\mathbb{K}^n \xrightarrow{\phi_\alpha} V \xrightarrow{\varphi_\alpha} \mathbb{K}^n$$

### Definition 46.17 (Matrix representation)

Let us now consider a basis  $\beta = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  for  $V$  and a basis  $\gamma = \{\mathbf{w}_1, \dots, \mathbf{w}_m\}$  for  $W$ . Let  $f$  be a linear transformation. Then, there exist unique scalars  $a_{ij} \in \mathbb{K}$  for  $1 \leq i \leq m$  such that for each  $1 \leq j \leq n$ :

$$f(\mathbf{v}_j) = \sum_{i=1}^m a_{ij} \mathbf{w}_i, \quad (46.4.2)$$

The scalars  $a_{ij}$  form a matrix  $A$  of size  $m \times n$ , called the **matrix representation** of  $f$  in the bases  $\beta$  and  $\gamma$ , denoted as  $A = [f]_\beta^\gamma$ .

We can draw a **commutative diagram** to demonstrate how matrix representations work:

$$\begin{array}{ccc} V & \xrightarrow{f} & W \\ \phi_\beta \uparrow & & \uparrow \phi_\gamma \\ \mathbb{K}^n & \xrightarrow{[f]_\beta^\gamma} & \mathbb{K}^m \end{array}$$

In this diagram, we designate the fields over which we define  $V$  and  $W$  in the bottom row. These contain the coordinate vectors of any  $\mathbf{v} \in V, \mathbf{w} \in W$ . To map from  $\mathbb{K}^n$  to  $V$  given a basis  $\beta$ , we need  $\phi_\beta$  as defined in Definition 46.15. Similarly to map from  $\mathbb{K}^m$  to  $W$  we need  $\phi_\gamma$ . Finally, if  $f : V \rightarrow W$  then we can map from  $\mathbb{K}^n$  to  $\mathbb{K}^m$  (map from one coordinate vector to another) by multiplying by the matrix representation.

**Proposition 46.18 (Columns of matrix representation)**

Given a matrix representation  $A$  of  $f$  in the bases  $\beta$  and  $\gamma$  then:

$$A = ([f(\mathbf{v}_1)]_\gamma \ [f(\mathbf{v}_2)]_\gamma \ \dots \ [f(\mathbf{v}_n)]_\gamma) \quad (46.4.3)$$

so that  $f(\mathbf{v}) = A\mathbf{v}$  for all  $\mathbf{v} \in V$  (or alternatively  $[f(\mathbf{v})]_\gamma = A[\mathbf{v}]_\beta$ ).

*Proof.* Note that for  $\mathbf{v} \in V$ :

$$f(\mathbf{v}) = f\left(\sum_{j=1}^n v_j \mathbf{v}_j\right) = \sum_{j=1}^n v_j f(\mathbf{v}_j) \quad (46.4.4)$$

so that if we let  $[f(\mathbf{v}_j)]_\gamma = (a_{1j} \ a_{2j} \ \dots \ a_{mj})^T$  then:

$$[f(\mathbf{v})]_\gamma = \begin{pmatrix} a_{11}v_1 + a_{12}v_2 + \dots + a_{1n}v_n \\ a_{21}v_1 + a_{22}v_2 + \dots + a_{2n}v_n \\ \vdots \\ a_{m1}v_1 + a_{m2}v_2 + \dots + a_{mn}v_n \end{pmatrix} \quad (46.4.5)$$

We can express this more compactly as:

$$[f(\mathbf{v})]_\gamma = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix} \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{pmatrix} \quad (46.4.6)$$

$$= ([f(\mathbf{v}_1)]_\gamma \ [f(\mathbf{v}_2)]_\gamma \ \dots \ [f(\mathbf{v}_n)]_\gamma) [\mathbf{v}]_\beta \quad (46.4.7)$$

$$\implies [f(\mathbf{v})]_\gamma = A[\mathbf{v}]_\beta \quad (46.4.8)$$

More abstractly, one may write:

$$f(\mathbf{v}) = f\left(\sum_{j=1}^n v_j \mathbf{v}_j\right) = \sum_{j=1}^n v_j \sum_{i=1}^m a_{ij} \mathbf{w}_i = \sum_{i=1}^m \left( \sum_{j=1}^n v_j a_{ij} \right) \mathbf{w}_i \quad (46.4.9)$$

so that we end up with:

$$(f(\mathbf{v}))_i = \sum_{j=1}^n v_j a_{ij} = (A\mathbf{v})_i \quad (46.4.10)$$

This however is the expression of matrix multiplication  $A\mathbf{v}$ , as desired. ■

**Theorem 46.19 (Uniqueness of matrix representation)**

The matrix representation of a linear representation with respect to the bases  $\beta$  and  $\gamma$  is unique.

*Proof.* Suppose we had two different matrix representations A and B as shown below:

$$A = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{pmatrix} \text{ and } B = \begin{pmatrix} b_{11} & b_{12} & \dots & b_{1n} \\ b_{21} & b_{22} & \dots & b_{2n} \\ \vdots & \vdots & & \vdots \\ b_{m1} & b_{m2} & \dots & b_{mn} \end{pmatrix} \quad (46.4.11)$$

Then the image of any basis vector  $\mathbf{v}_j \in \beta$  is from Proposition 46.18:

$$[f(\mathbf{v}_j)]_\gamma = (a_{1j} \ a_{2j} \ \dots \ a_{mj})^T = (b_{1j} \ b_{2j} \ \dots \ b_{mj})^T \quad (46.4.12)$$

so  $j$ th column of A and B coincide. Since this is true for any  $1 \leq j \leq n$  we have that  $A = B$ . ■

### Proposition 46.20 (Properties of matrix representations)

Let  $V$  and  $W$  be finite dimensional vector spaces with bases  $\beta$  and  $\gamma$  respectively, and let  $f, g$  be linear maps of  $V$  onto  $W$ . Then:

- (i)  $[f + g]_\beta^\gamma = [f]_\beta^\gamma + [g]_\beta^\gamma$
- (ii)  $[\alpha f]_\beta^\gamma = \alpha[f]_\beta^\gamma$

*Proof.* Let  $\beta = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  and  $\gamma = \{\mathbf{w}_1, \dots, \mathbf{w}_m\}$  then there exists unique  $a_{ij}$  and  $b_{ij}$  such that:

$$f(\mathbf{v}_j) = \sum_{i=1}^m a_{ij} \mathbf{w}_i, \text{ and } g(\mathbf{v}_j) = \sum_{i=1}^m b_{ij} \mathbf{w}_i \quad (46.4.13)$$

Then:

$$(\alpha f + g)(\mathbf{v}_j) = \sum_{i=1}^n (\alpha a_{ij} + b_{ij}) \mathbf{w}_i \quad (46.4.14)$$

and so:

$$([f + g]_\beta^\gamma)_{ij} = \alpha a_{ij} + b_{ij} = \alpha[f]_\beta^\gamma + [g]_\beta^\gamma \quad (46.4.15)$$

■

### Theorem 46.21 (Matrix representation of composition)

Let  $V, W, U$  be finite dimensional vector spaces with bases  $\alpha, \beta$  and  $\gamma$  respectively. Then, if  $g : V \rightarrow W$  and  $f : W \rightarrow U$  are linear maps:

$$[f \circ g]_\alpha^\gamma = [f]_\beta^\gamma [g]_\alpha^\beta \quad (46.4.16)$$

*Proof.* Let  $C = [f \circ g]_\alpha^\gamma$ ,  $A = [f]_\beta^\gamma$ ,  $B = [g]_\alpha^\beta$ . Also, let  $\alpha = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ ,  $\beta = \{\mathbf{w}_1, \dots, \mathbf{w}_m\}$ ,  $\gamma = \{\mathbf{u}_1, \dots, \mathbf{u}_l\}$ . We can draw the following commutative diagram:

$$\begin{array}{ccccc} V & \xrightarrow{g} & W & \xrightarrow{f} & U \\ \phi_\alpha \uparrow & & \phi_\beta \uparrow & & \phi_\gamma \uparrow \\ \mathbb{K}^n & \xrightarrow{[g]_\alpha^\beta} & \mathbb{K}^m & \xrightarrow{[f]_\beta^\gamma} & \mathbb{K}^l \end{array}$$

Then for  $\mathbf{v}_i \in \alpha$ :

$$(f \circ g)(\mathbf{v}_i) = f\left(\sum_{k=1}^m B_{ki} \mathbf{w}_k\right) = \sum_{k=1}^m B_{ki} f(\mathbf{w}_k) \quad (46.4.17)$$

$$= \sum_{k=1}^m B_{ki} \left( \sum_{j=1}^l A_{jk} \mathbf{u}_k \right) \quad (46.4.18)$$

$$= \sum_{k=1}^m \sum_{j=1}^l A_{jk} B_{ki} \mathbf{u}_k = \sum_{k=1}^l C_{ji} \mathbf{u}_k \quad (46.4.19)$$

where:

$$C_{ji} = \sum_{j=1}^l A_{jk} B_{ki} = \mathbf{A}_j \cdot \mathbf{B}^i \quad (46.4.20)$$

However, this is exactly the definition of matrix multiplication we encountered in the previous chapter, hence:

$$[f \circ g]_\alpha^\gamma = \mathbf{C} = \mathbf{AB} = [f]_\beta^\gamma [g]_\alpha^\beta \quad (46.4.21)$$

■

### Proposition 46.22 (*Invertibility of linear maps*)

Let  $f : V \rightarrow W$  be a linear map, where  $\alpha$  and  $\beta$  are bases of  $V$  and  $W$  respectively. Then  $f$  is invertible iff  $[f]_\beta^\gamma$  is invertible, then  $[f^{-1}]_\beta^\gamma = ([f]_\beta^\gamma)^{-1}$ .

*Proof.* Suppose that  $f$  has an inverse, so that  $\dim V = \dim W = n$ . Then  $[f]_\beta^\gamma$  is a square matrix of size  $n$ , and:

$$\mathbb{I}_n = [\text{id}_V]_\beta = [f^{-1} \circ f]_\beta = [f^{-1}]_\gamma^\beta [f]_\beta^\gamma \quad (46.4.22)$$

Similarly  $[f]_\beta^\gamma [f^{-1}]_\gamma^\beta$ , so  $[f]_\beta^\gamma$  is invertible, and  $[f^{-1}]_\beta^\gamma = ([f]_\beta^\gamma)^{-1}$ .

Now suppose  $\mathbf{A} = [f]_\beta^\gamma$  is invertible, then there exists  $\mathbf{B}$  so that  $\mathbf{AB} = \mathbf{BA} = \mathbb{I}_n$ . Then there exists a map  $g \in \text{Hom}(W, V)$  such that:

$$g(\mathbf{w}_j) = \sum_{i=1}^n B_{ij} \mathbf{v}_i \quad (46.4.23)$$

where  $\gamma = \{\mathbf{w}_1, \dots, \mathbf{w}_n\}$  and  $\beta = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  are bases of  $W$  and  $V$  respectively. Then:

$$[g \circ f]_\beta = [g]_\gamma^\beta [f]_\beta^\gamma = \mathbf{BA} = \mathbb{I}_n = [\text{id}_V]_\beta \quad (46.4.24)$$

so we conclude that  $g \circ f = \text{id}_V$  (and similarly  $f \circ g = \text{id}_V$ ) by the uniqueness of matrix representations.

■

### Theorem 46.23 (*All linear maps have a representation*)

The vector spaces  $\text{Mat}_{n,m}(\mathbb{K})$  and  $\text{Hom}(V, W)$  are isomorphic.

*Proof.* Let  $\beta = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  and  $\gamma = \{\mathbf{w}_1, \dots, \mathbf{w}_m\}$  be bases for  $V$  and  $W$  respectively. Then:

$$\Phi : \text{Hom}(V, W) \rightarrow \text{Mat}_{n,m}(\mathbb{K}) \quad (46.4.25)$$

$$f \mapsto [f]_{\beta}^{\gamma} \quad (46.4.26)$$

is an isomorphism. Indeed, it is linear, as shown in proposition 46.20. Furthermore, it is injective, since matrix representations are unique as was shown in Theorem 46.19. Finally, it is surjective, since given any matrix  $A$  then we can always find a linear map such that:

$$f(\mathbf{v}_i) = \sum_{i=1}^m A_{ij} \mathbf{w}_i \quad (46.4.27)$$

as warranted by Proposition 46.12. But this implies that  $\Phi(f) = [f]_{\beta}^{\gamma} = A$  as required. ■

An immediate consequence is that given two vector spaces  $V, W$  of dimensions  $n$  and  $m$  then  $\text{Hom}(V, W)$  has dimension  $n \cdot m$ .

This is a fundamental result since it proves that given any matrix, we can associate it to some linear map. Similarly, to every linear map we can associate some matrix representations. **Generally, a map is linear iff it has a matrix representation.**

Since every matrix defines a linear map, we can define its nullity and rank.

#### Definition 46.24 (Matrix rank)

Given a matrix  $A : \mathbb{F}^n \rightarrow \mathbb{F}^m$ , its associated linear map is:

$$A\mathbf{v} = \sum_{i=1}^m v_i f(\mathbf{v}_i) = \sum_{i=1}^m v_i A^i \quad (46.4.28)$$

then the image of  $A$  is  $\text{Im}(A) = \text{Span}(A^1, \dots, A^m)$ . Its column rank,  $\text{rk}(A)$ , is the number of linearly independent columns of  $A$ . Its row rank is the number of linearly independent rows of  $A$ , so  $\text{rk}(A^T)$ .

#### Proposition 46.25 (Column rank and row rank)

Column rank and row rank are the same.

*Proof.* Suppose that  $A_1$  can be written as a linear combination:

$$A_1 = \sum_{j=1}^n \alpha_j A_j \quad (46.4.29)$$

and let  $\alpha = (\alpha_2 \ \alpha_3 \ \dots \ \alpha_n)^T$ . Let:

$$A^i = \begin{pmatrix} a_i \\ \mathbf{b}_i \end{pmatrix}, \mathbf{b}_i = \begin{pmatrix} b_{2i} \\ b_{3i} \\ \vdots \\ b_{ni} \end{pmatrix} \implies A = \begin{pmatrix} a_1 & a_2 & \dots & a_m \\ b_{21} & b_{22} & \dots & b_{2m} \\ b_{31} & b_{32} & \dots & b_{3m} \\ \vdots & \vdots & & \vdots \\ b_{n1} & b_{n2} & \dots & b_{nm} \end{pmatrix} \quad (46.4.30)$$

Consequently

$$a_i = (\mathbf{A}^i)_1 = (\mathbf{A}_1)_i = \sum_{j=2}^n \alpha_j A_{ji} = \sum_{j=2}^n \alpha_j (\mathbf{A}^i)_j \quad (46.4.31)$$

$$= \alpha_2 b_{2i} + \alpha_3 b_{3i} + \dots + \alpha_n b_{ni} \quad (46.4.32)$$

$$= \boldsymbol{\alpha} \cdot \mathbf{b}_i \quad (46.4.33)$$

so that

$$\mathbf{A}^i = \begin{pmatrix} \boldsymbol{\alpha} \cdot \mathbf{b}_i \\ \mathbf{b}_i \end{pmatrix} \quad (46.4.34)$$

Therefore, if one row is a linear combination of the others, then we can drop it, leaving the row rank unchanged obviously. However, the column rank also remains unchanged. Indeed dropping the first row of  $\mathbf{A}$  we get:

$$\mathbf{A} = \begin{pmatrix} \boldsymbol{\alpha} \cdot \mathbf{b}_1 & \boldsymbol{\alpha} \cdot \mathbf{b}_2 & \dots & \boldsymbol{\alpha} \cdot \mathbf{b}_m \\ b_{21} & b_{22} & \dots & b_{2m} \\ b_{31} & b_{32} & \dots & b_{3m} \\ \vdots & \vdots & & \vdots \\ b_{n1} & b_{n2} & \dots & b_{nm} \end{pmatrix} \longrightarrow \mathbf{A}' = \begin{pmatrix} b_{21} & b_{22} & \dots & b_{2m} \\ b_{31} & b_{32} & \dots & b_{3m} \\ \vdots & \vdots & & \vdots \\ b_{n1} & b_{n2} & \dots & b_{nm} \end{pmatrix} \quad (46.4.35)$$

the column rank remains the same. Indeed, the first element of any column already contains a linear combination of all elements below it. Consequently, the column spans of  $\mathbf{A}$  and  $\mathbf{A}'$  are generated by the same number of elements. In other words, the column rank is unchanged by removing a linearly dependent row.

More rigorously, we need to prove that removing the row containing  $\boldsymbol{\alpha} \cdot \mathbf{b}_i$  will not alter the linear dependence of the columns. Suppose that  $\mathbf{A}^1, \mathbf{A}^2, \dots, \mathbf{A}^l$  is the maximal linearly independent set of columns of  $\mathbf{A}$ . Then note that

$$c_1 \begin{pmatrix} \boldsymbol{\alpha} \cdot \mathbf{b}_1 \\ \mathbf{b}_1 \end{pmatrix} + c_2 \begin{pmatrix} \boldsymbol{\alpha} \cdot \mathbf{b}_2 \\ \mathbf{b}_2 \end{pmatrix} + \dots + c_l \begin{pmatrix} \boldsymbol{\alpha} \cdot \mathbf{b}_l \\ \mathbf{b}_l \end{pmatrix} = \begin{pmatrix} \boldsymbol{\alpha} \cdot (c_1 \mathbf{b}_1 + \dots + c_l \mathbf{b}_l) \\ c_1 \mathbf{b}_1 + \dots + c_l \mathbf{b}_l \end{pmatrix} \quad (46.4.36)$$

So if  $\mathbf{A}^1, \mathbf{A}^2, \dots, \mathbf{A}^l$  are linearly independent, then  $\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_l$  are also linearly independent. Moreover, adding any other  $\mathbf{b}_i$  will result in linear dependence. Indeed, if this were not the case, then  $c'_1 \mathbf{b}_1 + \dots + c'_l \mathbf{b}_l + c'_{l+1} \mathbf{b}_{l+1} \implies c'_i = 0$  and so:

$$\begin{pmatrix} \boldsymbol{\alpha} \cdot (c'_1 \mathbf{b}_1 + \dots + c'_{l+1} \mathbf{b}_{l+1}) \\ c'_1 \mathbf{b}_1 + \dots + c'_{l+1} \mathbf{b}_{l+1} \end{pmatrix} = c'_1 \begin{pmatrix} \boldsymbol{\alpha} \cdot \mathbf{b}_1 \\ \mathbf{b}_1 \end{pmatrix} + \dots + c'_{l+1} \begin{pmatrix} \boldsymbol{\alpha} \cdot \mathbf{b}_{l+1} \\ \mathbf{b}_{l+1} \end{pmatrix} = \mathbf{0} \implies c'_i = 0 \quad (46.4.37)$$

so the columns  $\mathbf{A}^1, \mathbf{A}^2, \dots, \mathbf{A}^{l+1}$  would be linearly independent, a contradiction.

A similar reasoning can be used to show that the row rank is unchanged by removing a linearly dependent column.

So, if we continue this process, removing linearly dependent rows/columns, eventually we will end up with a final matrix  $\mathbf{A}$ , whose row and column ranks will not have been altered, and whose rows and columns will be linearly independent. This final matrix must be forcibly square. If it were  $n \times m$ , assume WLOG  $n < m$  then the  $m$  vectors in  $\mathbb{F}^n$  must be linearly dependent, a contradiction.

Hence, the number of linearly independent rows and columns are the same, as desired. ■

## 46.5 Change of basis and equivalence

### Definition 46.26 (Standard basis)

The standard basis of a coordinate vector space  $\mathbb{K}^n$  is:

$$\mathbf{e}_i = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} \leftarrow i\text{ith component} \quad (46.5.1)$$

and its coordinate map is denoted  $\varphi_{\text{id}}$ .

Suppose we are given the coordinate vector of some vector  $\mathbf{v} \in V$  in the **standard coordinates**:

$$\mathbf{v} = \sum_{i=1}^n v_i \mathbf{e}_i \quad (46.5.2)$$

How do we find its coordinates in some other basis  $\beta = \mathbf{v}_1, \dots, \mathbf{v}_n$ ?

Consider the commutative diagram below:

$$\begin{array}{ccc} V & \xleftarrow{\text{id}_V} & V \\ \uparrow \phi_{\text{id}} & & \uparrow \phi_{\beta} \\ \mathbb{K}_n & \xleftarrow{P_{\beta}} & \mathbb{K}_n \end{array}$$

then:

$$P_{\beta} = [\text{id}_V]_{\beta}^{\text{id}} = ([\mathbf{v}_1]_{\text{id}} \ [\mathbf{v}_2]_{\text{id}} \ \dots \ [\mathbf{v}_n]_{\text{id}}) \quad (46.5.3)$$

whereas:

$$M_{\beta} = P_{\beta}^{-1} = [\text{id}_V]_{\beta}^{\text{id}} = ([\mathbf{e}_1]_{\beta} \ [\mathbf{e}_2]_{\beta} \ \dots \ [\mathbf{e}_n]_{\beta}) \quad (46.5.4)$$

Indeed, notice that:

$$M_{\beta}[\mathbf{v}]_{\text{id}} = ([\mathbf{e}_1]_{\beta} \ [\mathbf{e}_2]_{\beta} \ \dots \ [\mathbf{e}_n]_{\beta})[\mathbf{v}]_{\text{id}} \quad (46.5.5)$$

$$= v_1[\mathbf{e}_1]_{\beta} + v_2[\mathbf{e}_2]_{\beta} + \dots + v_n[\mathbf{e}_n]_{\beta} \quad (46.5.6)$$

$$= [v_1\mathbf{e}_1 + v_2\mathbf{e}_2 + \dots + v_n\mathbf{e}_n]_{\beta} \quad (46.5.7)$$

$$= [\mathbf{v}]_{\beta} = [\text{id}_V(\mathbf{v})]_{\beta} \quad (46.5.8)$$

where  $v_1, v_2, \dots, v_n$  are the components of  $[\mathbf{v}]_{\text{id}}$ . Similarly:

$$P_\beta[\mathbf{v}]_\beta = ([\mathbf{v}_1]_{\text{id}} \ [\mathbf{v}_2]_{\text{id}} \ \dots \ [\mathbf{v}_n]_{\text{id}})[\mathbf{v}]_\beta \quad (46.5.9)$$

$$= v'_1[\mathbf{v}_1]_{\text{id}} + v'_2[\mathbf{v}_2]_{\text{id}} + \dots + v'_n[\mathbf{v}_n]_{\text{id}} \quad (46.5.10)$$

$$= [v'_1\mathbf{v}_1 + v'_2\mathbf{v}_2 + \dots + v'_n\mathbf{v}_n]_{\text{id}} \quad (46.5.11)$$

$$= [\mathbf{v}]_{\text{id}} = [\text{id}_V(\mathbf{v})]_{\text{id}} \quad (46.5.12)$$

where  $v'_1, v'_2, \dots, v'_n$  are the components of  $[\mathbf{v}]_\beta$ .

### Definition 46.27 (*Transition matrix*)

If  $\beta = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$  is a basis of a vector space  $V$  then:

$$M_\beta = ([\mathbf{e}_1]_\beta \ [\mathbf{e}_2]_\beta \ \dots \ [\mathbf{e}_n]_\beta) \quad (46.5.13)$$

called the **transition matrix** maps  $\varphi_{\text{id}}(\mathbf{v}) \mapsto \varphi_\beta(\mathbf{v})$  so that:

$$[\mathbf{v}]_\beta = M_\beta[\mathbf{v}]_{\text{id}} \quad (46.5.14)$$

We can view the transition matrix as the **matrix representation** of  $\varphi_\beta$ . Similarly  $P_\beta$  is the matrix representation of  $\phi_\beta$ .

Let us generalize this result for any two bases:

### Proposition 46.28 (*Change of coordinate matrix*)

Let  $\beta = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$  and  $\gamma$  be two bases of a vector space  $V$ . Then:

$$M_{\beta \rightarrow \gamma} = ([\mathbf{v}_1]_\gamma \ [\mathbf{v}_2]_\gamma \ \dots \ [\mathbf{v}_n]_\gamma) \quad (46.5.15)$$

is the **change of coordinate matrix** mapping  $\varphi_\beta(\mathbf{v}) \mapsto \varphi_\gamma(\mathbf{v})$  so that:

$$[\mathbf{v}]_\gamma = M_{\beta \rightarrow \gamma}[\mathbf{v}]_\beta \quad (46.5.16)$$

*Proof.* Consider the following commutative diagram:

$$\begin{array}{ccccc} & & M_{\beta \rightarrow \gamma} & & \\ & \nearrow P_\beta & & \searrow M_\gamma & \\ \mathbb{K}_n & \xrightarrow{\quad} & \mathbb{K}_n & \xrightarrow{\quad} & \mathbb{K}_n \\ \downarrow \phi_\beta & & \downarrow \phi_{\text{id}} & & \downarrow \phi_\gamma \\ V & \xrightarrow{\text{id}_V} & V & \xrightarrow{\text{id}_V} & V \end{array}$$

so that it is clear that:

$$M_{\beta \rightarrow \gamma} = [\varphi_\gamma \circ \phi_\beta]_{\text{id}} = [\text{id}]_\beta^\gamma \quad (46.5.17)$$

Then

$$M_{\beta \rightarrow \gamma} = M_\gamma P_\beta = M_\gamma([\mathbf{v}_1]_{\text{id}} \ [\mathbf{v}_2]_{\text{id}} \ \dots \ [\mathbf{v}_n]_{\text{id}}) \quad (46.5.18)$$

$$= (\phi_\gamma([\mathbf{v}_1]_{\text{id}}) \ \phi_\gamma([\mathbf{v}_2]_{\text{id}}) \ \dots \ \phi_\gamma([\mathbf{v}_n]_{\text{id}})) \quad (46.5.19)$$

$$= ([\mathbf{v}_1]_\gamma \ [\mathbf{v}_2]_\gamma \ \dots \ [\mathbf{v}_n]_\gamma) \quad (46.5.20)$$

as desired. Alternatively, using proposition 46.18 with  $f = \text{id}_V$  then:

$$M_{\beta \rightarrow \gamma} = [\text{id}]_{\beta}^{\gamma} = ([v_1]_{\gamma} \ [v_2]_{\gamma} \ \dots \ [v_n]_{\gamma}) \quad (46.5.21)$$

as found previously. ■

**Proposition 46.29 (Invertibility of change of coordinate matrix)**

The inverse of  $M_{\beta \rightarrow \gamma}$  is  $M_{\beta \rightarrow \gamma}^{-1} = M_{\gamma \rightarrow \beta}$ .

*Proof.* Let us firstly prove that  $M_{\gamma \rightarrow \beta} M_{\beta \rightarrow \gamma} = I_n$ . Indeed, let  $\beta = \{v_1, \dots, v_n\}$  so that:

$$M_{\gamma \rightarrow \beta} M_{\beta \rightarrow \gamma} = [\text{id}_V]_{\gamma}^{\beta} [\text{id}_V]_{\beta}^{\gamma} = [\text{id}_V]_{\beta}^{\beta} = ([v_1]_{\beta} \ [v_2]_{\beta} \ \dots \ [v_n]_{\beta}) = I_n \quad (46.5.22)$$

where we used Theorem 46.21. Similarly, one can also find that  $M_{\beta \rightarrow \gamma} M_{\gamma \rightarrow \beta} = I_n$ . ■

**Theorem 46.30 (Change of basis of linear transformation)**

Let  $\beta, \beta'$  be two bases for  $V$  and let  $\gamma, \gamma'$  be two bases for  $W$ . If  $f : V \rightarrow W$  is a linear map, then:

$$A' = M_{\gamma \rightarrow \gamma'} A M_{\beta' \rightarrow \beta} \quad (46.5.23)$$

*Proof.* Consider the commutative diagram below:

$$\begin{array}{ccccccc} V & \xrightarrow{\text{id}_V} & V & \xrightarrow{f} & W & \xleftarrow{\text{id}_V} & W \\ \phi_{\beta'} \uparrow & & \phi_{\beta} \uparrow & & \phi_{\gamma} \uparrow & & \phi_{\gamma'} \uparrow \\ \mathbb{K}^m & \xrightarrow{M_{\beta' \rightarrow \beta}} & \mathbb{K}^m & \xrightarrow{A} & \mathbb{K}^n & \xleftarrow{M_{\gamma' \rightarrow \gamma}} & \mathbb{K}^n \\ & & & & \searrow A' & & \end{array}$$

Then note that:

$$A = \varphi_{\gamma} \circ f \circ \phi_{\beta} \quad (46.5.24)$$

$$A' = \varphi_{\gamma'} \circ f \circ \phi_{\beta'} \quad (46.5.25)$$

hence:

$$A' = \varphi_{\gamma'} \circ f \circ \phi_{\beta'} \quad (46.5.26)$$

$$= \varphi_{\gamma'} \circ (\phi_{\gamma} \circ \varphi_{\gamma}) \circ f \circ (\phi_{\beta} \circ \varphi_{\beta}) \circ \phi_{\beta'} \quad (46.5.27)$$

$$= (\varphi_{\gamma'} \circ \phi_{\gamma}) \circ A \circ (\varphi_{\beta} \circ \phi_{\beta'}) \quad (46.5.28)$$

$$= M_{\gamma' \rightarrow \gamma}^{-1} A M_{\beta' \rightarrow \beta} \quad (46.5.29)$$

but we have proven that  $M_{\gamma' \rightarrow \gamma}^{-1} = M_{\gamma \rightarrow \gamma'}$  so that:

$$A' = M_{\gamma \rightarrow \gamma'} A M_{\beta' \rightarrow \beta} \quad (46.5.30)$$

as desired. ■

We can interpret this result with some intuition by considering the action of each matrix in 46.5.29 on  $[\mathbf{v}]_{\beta'}$ . Indeed, the matrix  $M_{\beta' \rightarrow \beta}$  converts it to  $[\mathbf{v}]_{\beta}$ . Then  $A$  maps it to  $A[\mathbf{v}]_{\beta} = [f(\mathbf{v})]_{\gamma}$  by Proposition 46.18. Finally  $M_{\gamma \rightarrow \gamma'}$  maps it to  $[f(\mathbf{v})]_{\gamma'}$ .

### Definition 46.31 (Equivalent matrices)

Let  $A, B \in \text{Mat}_{n,m}(\mathbb{K})$  are **equivalent matrices** if there exists invertible matrices  $P, Q \in \text{Mat}_m(\mathbb{K})$  such that:

$$B = Q^{-1}AP \quad (46.5.31)$$

We see immediately that if two matrices represent the same linear map with respect to different bases, then they are similar.

### Proposition 46.32 (Similarity to special matrix)

Any matrix  $A \in \text{Mat}_{m,n}(\mathbb{K})$  is equivalent to:

$$\begin{pmatrix} \mathbb{I}_r & 0 \\ 0 & 0 \end{pmatrix} \quad (46.5.32)$$

where  $r = \text{rk}(A)$ .

*Proof.* We begin by proving that any linear map  $f : V \rightarrow W$  has some set of bases  $\beta$  and  $\gamma$  of  $V$  and  $W$  respectively such that:

$$[f]_{\beta}^{\gamma} = \begin{pmatrix} \mathbb{I}_r & 0 \\ 0 & 0 \end{pmatrix} \quad (46.5.33)$$

Set  $r$  so that  $\text{null}(f) = n-r$ , so that  $\ker f$  has basis  $\mathbf{v}_{r+1}, \dots, \mathbf{v}_n$  which we extend to  $\mathbf{v}_1, \dots, \mathbf{v}_r, \mathbf{v}_{r+1}, \dots, \mathbf{v}_n$  a basis for  $V$ . We know that  $f(\mathbf{v}_1), f(\mathbf{v}_2), \dots, f(\mathbf{v}_n)$  is a basis for  $\text{Im}(f)$ , and can be extended to a basis  $\gamma$  for  $W$ . Consequently:

$$[f]_{\beta}^{\gamma} = ([f(\mathbf{v}_1)]_{\gamma} [f(\mathbf{v}_2)]_{\gamma} \dots [f(\mathbf{v}_r)]_{\gamma} [f(\mathbf{v}_{r+1})]_{\gamma} \dots f(\mathbf{v}_n)]_{\gamma}) \quad (46.5.34)$$

but  $f(\mathbf{v}_i) = \mathbf{0}$  for  $r < i \leq n$  so that:

$$[f]_{\beta}^{\gamma} = ([f(\mathbf{v}_1)]_{\gamma} [f(\mathbf{v}_2)]_{\gamma} \dots [f(\mathbf{v}_r)]_{\gamma} \mathbf{0} \dots \mathbf{0}) = \begin{pmatrix} \mathbb{I}_r & 0 \\ 0 & 0 \end{pmatrix} \quad (46.5.35)$$

Now, if we choose some matrix  $A \in \text{Mat}_{m,n}(\mathbb{K})$ , then by Theorem 46.23 it must be the matrix representation of some linear map  $f$ . Hence, it must be equivalent to:

$$\begin{pmatrix} \mathbb{I}_r & 0 \\ 0 & 0 \end{pmatrix} \quad (46.5.36)$$

which is also another representation of  $f$  as desired. ■

We can use this theorem to prove the equivalence of row and column rank. Indeed, if we let  $A \in$

$\text{Mat}_{m,n}(\mathbb{K})$  then there exist  $Q, P \in \text{Mat}_m(\mathbb{K})$  such that:

$$Q^{-1}AP = \begin{pmatrix} I_r & 0 \\ 0 & 0 \end{pmatrix} \quad (46.5.37)$$

implying that:

$$(Q^{-1}AP)^T = P^T A^T (Q^{-1})^T = \begin{pmatrix} I_r & 0 \\ 0 & 0 \end{pmatrix} \quad (46.5.38)$$

so that  $A^T$  is also equivalent to  $\begin{pmatrix} I_r & 0 \\ 0 & 0 \end{pmatrix}$ , and hence  $A^T$  and  $A$  both represent the same map and must therefore have the same rank. So column rank and row rank are the same.

# Solving linear equations

## 47.1 Structure of solutions

Let's consider a linear map  $f : V \rightarrow W$ . Suppose we wish to find the solutions to

$$f(\mathbf{x}) = \mathbf{b}, \mathbf{x} \in V, \mathbf{b} \in W \quad (47.1.1)$$

For  $\mathbf{b} \neq 0$ , this is known as the inhomogeneous linear equation, whose associated homogeneous equation is:

$$f(\mathbf{x}) = \mathbf{0} \quad (47.1.2)$$

Clearly, the solution of the latter is  $\text{Ker}(f)$ .

### **Proposition (Structure of solutions)**

Suppose  $\mathbf{x}_0 \in V$  is a solution to the inhomogeneous equation (47.1.1). Then the general solution is given by:

$$\mathbf{x} = \mathbf{x}_0 + \text{Ker}(f) \quad (47.1.3)$$

*Proof.* Suppose  $\mathbf{x} \in V$  is a general solution to (47.1.1) so that  $f(\mathbf{x}) = \mathbf{b}$ . Since  $\mathbf{x}_0$  is a solution, we have that  $f(\mathbf{x}_0) = \mathbf{b}$ , so we may write that  $f(\mathbf{x} - \mathbf{x}_0) = \mathbf{0}$ . This implies that  $\mathbf{x} - \mathbf{x}_0 \in \text{Ker}(f)$  or alternatively that  $\mathbf{x} = \mathbf{x}_0 + \text{Ker}(f)$ . ■

Suppose we have a linear transformation  $f$  represented by a matrix  $A : \mathbb{K}^n \rightarrow \mathbb{K}^m$ . Then, for  $\mathbf{x} \in \mathbb{K}^n$  and  $\mathbf{b} \in \mathbb{K}^m$  we consider the linear system of equations  $A\mathbf{x} = \mathbf{b}$ :

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n = b_1 \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n = b_2 \\ \vdots \\ a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n = b_m \end{cases} \quad (47.1.4)$$

We must now consider the following three scenarios:

- (i) if  $\text{rk}(A) = m$ , then a solution exists for any choice of  $\mathbf{b}$ . Indeed, since  $\text{Im}(A) = \mathbb{K}^m$  (since  $\text{Im}(A) \subseteq \mathbb{K}^m$  and they have the same dimensionality), it follows that any vector  $\mathbf{b} \in \mathbb{K}^m$  may be expressed as an image of  $f$ . The number of free parameters is given by  $\dim \text{Ker}(A) = n - m$ . Geometrically, if we let  $m = n = 3$  and  $\mathbb{K} = \mathbb{R}$ , then we see that if  $\text{rk}(A) = 3$  then the linear map  $f$  maps the typical Euclidean space to itself.

- (ii) if  $\text{rk}(A) < m$  and  $\mathbf{b} \in \text{Im}(A)$  then solution exists. Indeed, in this case a generic case of  $\mathbf{b}$  will no longer work, we must be careful and make sure that it belongs to the image of  $A$ . The number of free parameters will be  $\dim \text{Ker}(A) = n - \text{rk}(A)$ . Geometrically, this corresponds to the linear map  $f$  mapping the Euclidean space to a subspace of itself, such as a plane. We may add any vector that maps to the origin to a solution. For example in the case where  $\text{Im}(A)$  is a plane the kernel will be a line, any vector on this line may be added, giving a free parameter.
- (iii) if  $\text{rk}(A) < m$  but  $\mathbf{b} \notin \text{Im}(A)$  then solution doesn't exist. Indeed if the vector  $\mathbf{b}$  doesn't lie in the space spanned by  $A$  then clearly a solution will not exist, since no vector  $\mathbf{x}$  will get mapped to  $\mathbf{b}$ .

**Definition (Augmented matrix)**

Consider the linear system of equations  $A\mathbf{x} = \mathbf{b}$  where  $A : \mathbb{K}^n \rightarrow \mathbb{K}^m$  is a matrix and  $\mathbf{x} \in \mathbb{K}^n$  and  $\mathbf{b} \in \mathbb{K}^m$ . We define its augmented matrix to be:

$$(A|\mathbf{b}) = \left( \begin{array}{cccc|c} a_{11} & a_{12} & \dots & a_{1n} & b_1 \\ a_{21} & a_{22} & \dots & a_{2n} & b_2 \\ \vdots & \vdots & & \vdots & \dots \\ a_{m1} & a_{m2} & \dots & a_{mn} & b_m \end{array} \right) \quad (47.1.5)$$

**Theorem (Rank of augmented matrix)**

For a matrix  $A : \mathbb{K}^n \rightarrow \mathbb{K}^m$  and  $\mathbf{b} \in \mathbb{K}^n$ :

$$\mathbf{b} \in \text{Im}(A) \iff \text{rk}(A) = \text{rk}((A|\mathbf{b})) \quad (47.1.6)$$

*Proof.*  $\implies$  Suppose that  $\mathbf{b} \in \text{Im}(A)$ . Then adding this vector to  $A$  will not alter the space it spans, so that  $\text{rk}(A) = \text{rk}((A|\mathbf{b}))$ .

$\impliedby$  Suppose that  $\text{rk}(A) = \text{rk}((A|\mathbf{b}))$ . Then, this means that:

$$\text{Span}(\mathbf{A}^1, \mathbf{A}^2, \dots, \mathbf{A}^n, \mathbf{b}) = \text{Span}(\mathbf{A}^1, \mathbf{A}^2, \dots, \mathbf{A}^n) \quad (47.1.7)$$

showing that  $\mathbf{b} \in \text{Im}(A)$ . ■

Therefore, if we can find the rank of the augmented matrix and show that it is equal to the rank of the coefficient matrix  $A$ , then we have shown that a solution must indeed exist. We can find the rank of matrices using elementary matrix operations.

## 47.2 Elementary matrix operations

**Definition (Elementary row operations)**

The following are **elementary row operations**:

- (R1) exchange two rows
  - (R2) scale a row by a non-zero scalar
  - (R3) add a non-zero multiple of a row to another row
- These operations do not alter the rank of a matrix.

Because these operations do not alter the rank of a matrix, we may use them to transform a given matrix into a simpler one where it is easier to determine the span of its column/row vectors.

### Definition ((Reduced) row echelon form)

A matrix is said to be in **row echelon form** if:

- (i) a leading entry in a non-zero row is strictly to the right of the leading entry in the row above
- (ii) zero rows are at the bottom

so it has general form:

$$\begin{pmatrix} \dots & a_{ij_1} & \dots & \dots & \dots & * \\ \vdots & & a_{2j_2} & & & \vdots \\ \vdots & & & & & \vdots \\ \vdots & & & & a_{rj_r} & \dots \\ 0 & 0 & \dots & \dots & \dots & 0 \end{pmatrix} \quad (47.2.1)$$

Instead, a matrix is said to be in **reduced row echelon form** if:

- (i) it is in row echelon form
- (ii) each leading entry is a 1
- (iii) each leading 1 is the only non-zero entry in its column

For example, the following matrix:

$$\begin{pmatrix} 0 & 1 & 0 & 2 & 0 & 7 \\ 0 & 0 & 1 & -3 & 0 & 2 \\ 0 & 0 & 0 & 0 & 1 & 7 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \quad (47.2.2)$$

is not in reduced row echelon form since the leading 1 in the fourth row is not the only non-zero element in its column. However it is in row echelon form.

Instead, the following matrix:

$$\begin{pmatrix} 1 & 0 & 0 & 3 \\ 0 & 0 & 1 & \frac{1}{2} \\ 0 & 0 & 0 & 0 \end{pmatrix} \quad (47.2.3)$$

is indeed in reduced row echelon form.

This gives us a strategy, known as Gauss-Jordan elimination to solve systems of linear equations.

### Strategy (Gauss-Jordan elimination)

- (i) Apply row operations to  $(A|b)$  until the matrix  $A$  within it is in row-echelon form.
- (ii) Let  $r = \text{rk}(A)$ . If  $b_i \neq 0$  for some  $i > r$ , then this means that  $\text{rk}(A) \neq \text{rk}((A|b))$ , and hence  $b \notin \text{Im}(A)$ . The system therefore has no solutions.
- (iii) Otherwise, convert to reduced row echelon form, and solve the resulting system of equations.

**Example.** Consider:

$$\begin{cases} 3x_1 - 11x_2 - 3x_3 = 3 \\ 2x_1 - 6x_2 - 2x_3 = 1 \\ 5x_1 - 17x_2 - 6x_3 = 2 \\ 4x_1 - 8x_2 = 7 \end{cases} \quad (47.2.4)$$

We construct the augmented matrix:

$$(A|\mathbf{b}) = \left( \begin{array}{ccc|c} 3 & -11 & -3 & 3 \\ 2 & -6 & -2 & 1 \\ 5 & -17 & -6 & 2 \\ 4 & -8 & 0 & 7 \end{array} \right) \quad (47.2.5)$$

which we reduce to row echelon form:

$$(A|\mathbf{b}) = \left( \begin{array}{ccc|c} 3 & -11 & -3 & 3 \\ 2 & -6 & -2 & 1 \\ 5 & -17 & -6 & 2 \\ 4 & -8 & 0 & 7 \end{array} \right) \rightarrow \left( \begin{array}{ccc|c} 3 & -11 & -3 & 3 \\ 0 & \frac{4}{3} & 0 & -1 \\ 0 & \frac{4}{3} & -1 & -3 \\ 0 & \frac{20}{3} & 4 & 3 \end{array} \right) \quad (47.2.6)$$

$$\rightarrow \left( \begin{array}{ccc|c} 3 & -11 & -3 & 3 \\ 0 & \frac{4}{3} & 0 & -1 \\ 0 & 0 & -1 & -2 \\ 0 & 0 & 4 & 8 \end{array} \right) \rightarrow \left( \begin{array}{ccc|c} 3 & 0 & -3 & -\frac{21}{4} \\ 0 & \frac{4}{3} & 0 & -1 \\ 0 & 0 & -1 & -2 \\ 0 & 0 & 0 & 0 \end{array} \right) \quad (47.2.7)$$

(47.2.8)

So we see that  $\text{rk}(A) = \text{rk}((A|\mathbf{b})) = 3$  implying that  $\mathbf{b} \in \text{Im}(A)$ , and that a solution to the system exists.

We therefore convert the augmented matrix into reduced row echelon form:

$$(A|\mathbf{b}) \rightarrow \left( \begin{array}{ccc|c} 3 & 0 & -3 & -\frac{21}{4} \\ 0 & \frac{4}{3} & 0 & -1 \\ 0 & 0 & -1 & -2 \\ 0 & 0 & 0 & 0 \end{array} \right) \rightarrow \left( \begin{array}{ccc|c} 3 & 0 & 0 & \frac{3}{4} \\ 0 & \frac{4}{3} & 0 & -1 \\ 0 & 0 & -1 & -2 \\ 0 & 0 & 0 & 0 \end{array} \right) \quad (47.2.9)$$

$$\rightarrow \left( \begin{array}{ccc|c} 1 & 0 & 0 & \frac{1}{4} \\ 0 & 1 & 0 & -\frac{3}{4} \\ 0 & 0 & 1 & 2 \\ 0 & 0 & 0 & 0 \end{array} \right) \quad (47.2.10)$$

which gives the solution:

$$x_1 = \frac{1}{4}, x_2 = -\frac{3}{4}, x_3 = 2 \quad (47.2.11)$$

◀

A useful tool to check whether or not an arithmetic mistake has been made while performing row reduction is to write the sum of the entries in each row as an extra column, so for example. For

example

$$\left( \begin{array}{ccc|cc} 3 & -11 & -3 & 3 & -8 \\ 0 & \frac{4}{3} & 0 & -1 & \frac{1}{3} \\ 0 & \frac{4}{3} & -1 & -3 & -\frac{8}{3} \\ 0 & \frac{20}{3} & 4 & 3 & \frac{41}{3} \end{array} \right) \quad (47.2.12)$$

When we perform a row operation, we perform it on this extra column as well. If the numbers in this final column still correspond to the sum of the elements in the corresponding row, then no mistakes have been made. In our example, we get:

$$\left( \begin{array}{ccc|cc} 3 & -11 & -3 & 3 & -8 \\ 0 & \frac{4}{3} & 0 & -1 & \frac{1}{3} \\ 0 & 0 & -1 & -2 & -3 \\ 0 & 0 & 4 & 8 & 12 \end{array} \right) \quad (47.2.13)$$

so we do indeed find that the sum of all the rows are given in the transformed fifth column.

If say we had gotten, say, 13 in the last row, then an arithmetic mistake must have been made.

### 47.3 Inverting matrices

#### Definition (*Elementary matrix*)

The matrices obtained by performing row operations on  $\mathbb{1}$  are known as elementary matrices.

Note that elementary matrices are important because they represent row operations. Suppose we have some row operation, which when acted on  $\mathbb{1}_n$  gives the elementary matrix E. Then applying the same row operation on another  $n \times n$  matrix A we will get  $A' = EA$ . For example, if E represents the exchange of rows  $i$  and  $j$ , then:

$$\mathbb{1} = (\mathbf{e}_1 \dots \mathbf{e}_j \dots \mathbf{e}_i \dots \mathbf{e}_n)^T \implies E = (\mathbf{e}_1 \dots \mathbf{e}_j \dots \mathbf{e}_i \dots \mathbf{e}_n)^T \quad (47.3.1)$$

so that:

$$E_1 A = (\mathbf{e}_1 \dots \mathbf{e}_j \dots \mathbf{e}_i \dots \mathbf{e}_n)^T (A_1 \dots A_i \dots A_j \dots A_n) \quad (47.3.2)$$

$$= \begin{pmatrix} \mathbf{e}_1 \cdot A_1 & \mathbf{e}_1 \cdot A_2 & \dots & \mathbf{e}_1 \cdot A_n \\ \vdots & \vdots & & \vdots \\ \mathbf{e}_j \cdot A_1 & \mathbf{e}_j \cdot A_2 & \dots & \mathbf{e}_j \cdot A_n \\ \vdots & \vdots & & \vdots \\ \mathbf{e}_i \cdot A_1 & \mathbf{e}_i \cdot A_2 & \dots & \mathbf{e}_i \cdot A_n \\ \vdots & \vdots & & \vdots \\ \mathbf{e}_n \cdot A_1 & \mathbf{e}_n \cdot A_2 & \dots & \mathbf{e}_n \cdot A_n \end{pmatrix} \quad (47.3.3)$$

which is indeed the version of  $\mathbf{A}$  with the  $i$ th and  $j$ th rows exchanged:

$$\mathbf{A} = \mathbb{1}\mathbf{A} = (\mathbf{e}_1 \dots \mathbf{e}_j \dots \mathbf{e}_i \dots \mathbf{e}_n)^T (\mathbf{A}_1 \dots \mathbf{A}_i \dots \mathbf{A}_j \dots \mathbf{A}_n) \quad (47.3.4)$$

$$= \begin{pmatrix} \mathbf{e}_1 \cdot \mathbf{A}_1 & \mathbf{e}_1 \cdot \mathbf{A}_2 & \dots & \mathbf{e}_1 \cdot \mathbf{A}_n \\ \vdots & \vdots & & \vdots \\ \mathbf{e}_i \cdot \mathbf{A}_1 & \mathbf{e}_i \cdot \mathbf{A}_2 & \dots & \mathbf{e}_i \cdot \mathbf{A}_n \\ \vdots & \vdots & & \vdots \\ \mathbf{e}_j \cdot \mathbf{A}_1 & \mathbf{e}_j \cdot \mathbf{A}_2 & \dots & \mathbf{e}_j \cdot \mathbf{A}_n \\ \vdots & \vdots & & \vdots \\ \mathbf{e}_n \cdot \mathbf{A}_1 & \mathbf{e}_n \cdot \mathbf{A}_2 & \dots & \mathbf{e}_n \cdot \mathbf{A}_n \end{pmatrix} \quad (47.3.5)$$

Similarly, if  $E_2$  represents multiplication of the  $i$ th row by a scalar  $\lambda$  then clearly:

$$E_2\mathbf{A} = (\mathbf{e}_1 \dots \lambda\mathbf{e}_i \dots \dots \mathbf{e}_n)^T (\mathbf{A}_1 \dots \mathbf{A}_i \dots \dots \mathbf{A}_n) \quad (47.3.6)$$

$$= \begin{pmatrix} \mathbf{e}_1 \cdot \mathbf{A}_1 & \mathbf{e}_1 \cdot \mathbf{A}_2 & \dots & \mathbf{e}_1 \cdot \mathbf{A}_n \\ \vdots & \vdots & & \vdots \\ \lambda\mathbf{e}_i \cdot \mathbf{A}_1 & \lambda\mathbf{e}_i \cdot \mathbf{A}_2 & \dots & \lambda\mathbf{e}_i \cdot \mathbf{A}_n \\ \vdots & \vdots & & \vdots \\ \mathbf{e}_n \cdot \mathbf{A}_1 & \mathbf{e}_n \cdot \mathbf{A}_2 & \dots & \mathbf{e}_n \cdot \mathbf{A}_n \end{pmatrix} \quad (47.3.7)$$

which is the version of  $\mathbf{A}$  with the  $i$ th row multiplied by  $\lambda$ .

Finally, suppose that  $E_3$  represents adding a  $\lambda$ -multiple of the  $j$ th row to the  $i$ th row. Then we find that:

$$E_1\mathbf{A} = (\mathbf{e}_1 \dots \mathbf{e}_i \dots \mathbf{e}_j + \lambda\mathbf{e}_i \dots \mathbf{e}_n)^T (\mathbf{A}_1 \dots \mathbf{A}_i \dots \mathbf{A}_j \dots \mathbf{A}_n) \quad (47.3.8)$$

$$= \begin{pmatrix} \mathbf{e}_1 \cdot \mathbf{A}_1 & \mathbf{e}_1 \cdot \mathbf{A}_2 & \dots & \mathbf{e}_1 \cdot \mathbf{A}_n \\ \vdots & \vdots & & \vdots \\ \mathbf{e}_i \cdot \mathbf{A}_1 & \mathbf{e}_i \cdot \mathbf{A}_2 & \dots & \mathbf{e}_i \cdot \mathbf{A}_n \\ \vdots & \vdots & & \vdots \\ \mathbf{e}_j \cdot \mathbf{A}_1 + \lambda\mathbf{e}_i \cdot \mathbf{A}_1 & \mathbf{e}_j \cdot \mathbf{A}_2 + \lambda\mathbf{e}_i \cdot \mathbf{A}_2 & \dots & \mathbf{e}_j \cdot \mathbf{A}_n + \lambda\mathbf{e}_i \cdot \mathbf{A}_n \\ \vdots & \vdots & & \vdots \\ \mathbf{e}_n \cdot \mathbf{A}_1 & \mathbf{e}_n \cdot \mathbf{A}_2 & \dots & \mathbf{e}_n \cdot \mathbf{A}_n \end{pmatrix} \quad (47.3.9)$$

which is indeed the version of  $\mathbf{A}$  with the  $\lambda$ -multiple of the  $i$ th row added to the  $j$ th row.

Note also that elementary row matrices are all invertible, because the row operations they represent are all invertible.

### Theorem (*Invertibility theorem*)

- (a) An  $n \times n$  square matrix  $\mathbf{A}$  is invertible iff its reduced row echelon form is  $\mathbb{1}$ , so if  $\text{rk}(\mathbf{A}) = n$ .
- (b) Any sequence of row operations that transform  $\mathbf{A}$  to  $\mathbb{1}$  also transform  $\mathbb{1}$  to  $\mathbf{A}^{-1}$ .

*Proof.* Let  $\mathbf{A}$  be an  $n \times n$  matrix whose reduced row echelon form is:

$$\mathbf{U} = E_k E_{k-1} \dots E_1 \mathbf{A} = \mathbf{B} \mathbf{A} \quad (47.3.10)$$

where  $E_k, E_{k-1}, \dots, E_1$  are elementary matrices. Since they are all invertible, we have that  $\mathbf{B}^{-1}$  exists.

$\Rightarrow$  Suppose  $A$  is invertible. Then  $U$  is the product of invertible matrices, and is therefore invertible itself. Hence it cannot have any zero rows, since such matrices are not invertible. It follows from the conditions of the reduced row echelon form of matrices that the only possible choice of  $U$  is  $\mathbb{1}$ . Indeed, it is upper triangular (so in row echelon) with each leading entry as 1. Because it has  $n$  leading 1s and  $n$  rows, these leading entries must be on the diagonal. Finally, since there are no other entries on each column with a leading 1, and all columns have a leading 1, we get the identity matrix.

$\Leftarrow$  Suppose  $U = \mathbb{1}$ , then:

$$BA = \mathbb{1} \Rightarrow A = B^{-1}\mathbb{1} \Rightarrow AB = \mathbb{1} \quad (47.3.11)$$

Note however that  $\mathbb{1}$  and  $B^{-1}$  are invertible, so  $A$  will also be invertible, with  $A^{-1} = B$ .

Therefore, we find that:

$$A^{-1} = B = E_k E_{k-1} \dots E_1 \mathbb{1} \Rightarrow (A|\mathbb{1}) \rightarrow (\mathbb{1}|A^{-1}) \quad (47.3.12)$$

so we find  $A^{-1}$  by applying the same row operations that row reduce  $A$  to  $\mathbb{1}$ . ■

We can use the invertibility theorem to find the inverse of matrices.

**Example.** Let's find the inverse of the following matrix:

$$A = \begin{pmatrix} 1 & 4 & 1 \\ 1 & 6 & 3 \\ 2 & 3 & 0 \end{pmatrix} \quad (47.3.13)$$

We find that:

$$\left( \begin{array}{ccc|ccc} 1 & 4 & 1 & 1 & 0 & 0 \\ 1 & 6 & 3 & 0 & 1 & 0 \\ 2 & 3 & 0 & 0 & 0 & 1 \end{array} \right) \rightarrow \left( \begin{array}{ccc|ccc} 1 & 4 & 1 & 1 & 0 & 0 \\ 0 & 2 & 2 & -1 & 1 & 0 \\ 0 & -5 & -2 & -2 & 0 & 1 \end{array} \right) \rightarrow \left( \begin{array}{ccc|ccc} 1 & 4 & 1 & 1 & 0 & 0 \\ 0 & 2 & 2 & -1 & 1 & 0 \\ 0 & 0 & 3 & -\frac{9}{2} & \frac{5}{2} & 1 \end{array} \right) \quad (47.3.14)$$

$$\rightarrow \left( \begin{array}{ccc|ccc} 1 & 4 & 1 & 1 & 0 & 0 \\ 0 & 2 & 0 & 2 & -\frac{2}{3} & -\frac{2}{3} \\ 0 & 0 & 1 & -\frac{3}{2} & \frac{5}{6} & \frac{1}{3} \end{array} \right) \rightarrow \left( \begin{array}{ccc|ccc} 1 & 4 & 0 & \frac{5}{2} & -\frac{5}{6} & -\frac{1}{3} \\ 0 & 2 & 0 & 2 & -\frac{2}{3} & -\frac{2}{3} \\ 0 & 0 & 1 & -\frac{3}{2} & \frac{5}{6} & \frac{1}{3} \end{array} \right) \quad (47.3.15)$$

$$\rightarrow \left( \begin{array}{ccc|ccc} 1 & 0 & 0 & -\frac{3}{2} & \frac{1}{2} & 1 \\ 0 & 2 & 0 & 2 & -\frac{2}{3} & -\frac{2}{3} \\ 0 & 0 & 1 & -\frac{3}{2} & \frac{5}{6} & \frac{1}{3} \end{array} \right) \rightarrow \left( \begin{array}{ccc|ccc} 1 & 0 & 0 & -\frac{3}{2} & \frac{1}{2} & 1 \\ 0 & 1 & 0 & 1 & -\frac{1}{3} & -\frac{1}{3} \\ 0 & 0 & 1 & -\frac{3}{2} & \frac{5}{6} & \frac{1}{3} \end{array} \right) \quad (47.3.16)$$

so we see that:

$$A^{-1} = \frac{1}{6} \begin{pmatrix} -9 & 3 & 6 \\ 6 & -2 & -2 \\ -9 & 5 & 2 \end{pmatrix} \quad (47.3.17)$$

To check:

$$\frac{1}{6} \begin{pmatrix} -9 & 3 & 6 \\ 6 & -2 & -2 \\ -9 & 5 & 2 \end{pmatrix} \begin{pmatrix} 1 & 4 & 1 \\ 1 & 6 & 3 \\ 2 & 3 & 0 \end{pmatrix} = \frac{1}{6} \begin{pmatrix} 6 & 0 & 0 \\ 0 & 6 & 0 \\ 0 & 0 & 6 \end{pmatrix} = \mathbb{1} \quad (47.3.18)$$

as expected. ◀

The invertibility of matrices is especially important when solving linear systems of equations, since they can be used to find solutions whenever the associated system has only trivial solutions (that is, when the kernel of the matrix is null, and hence the rank is maximal).

**Proposition (Linear systems and invertibility)**

For an  $n \times n$  matrix  $A$ , the following statements are equivalent:

- (a)  $A$  is invertible
- (b) The system  $Ax = b$  has a unique solution for any  $b$
- (c) The system  $Ax = 0$  only has a trivial solution.

*Proof.*

- (a)  $\Rightarrow$  (b) Let  $A$  be invertible. Suppose  $Ax = b$ , then multiplying by  $A^{-1}$  then  $A^{-1}Ax = A^{-1}b \Rightarrow x = A^{-1}b$ . Instead, if  $x = A^{-1}b$  then multiplying by  $A$  we get  $Ax = b$ .
- (b)  $\Rightarrow$  (c) Suppose  $Ax = b$  has a unique solution for any  $b$ . Then  $Ax = 0$  also has a unique solution, which can only be the trivial solution
- (c)  $\Rightarrow$  (a) Suppose that  $Ax = 0$  only has a trivial solution. Then this means that when we reduce the augmented matrix:

$$\left( \begin{array}{ccc|c} a_{11} & a_{12} & a_{13} & 0 \\ a_{21} & a_{22} & a_{23} & 0 \\ a_{31} & a_{32} & a_{33} & 0 \end{array} \right) \quad (47.3.19)$$

then we must get that:

$$\left( \begin{array}{ccc|c} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{array} \right) \quad (47.3.20)$$

since there can only be trivial solutions. So  $A$  has a reduced row echelon form of  $\mathbb{I}$ , proving by the Invertibility theorem that it is invertible.

Alternatively, we could note that if  $\dim \text{Ker}(A) = 0$  then  $\text{rk}(A) = n$ . Therefore  $A$  is invertible. ■

**Example.** Let's consider the system:

$$\left\{ \begin{array}{l} x + 4y + z = 4 \\ x + 6y + 3z = 6 \\ 2x + 3y = 9 \end{array} \right. \quad (47.3.21)$$

which may be written in matrix form as:

$$Ax = b, A = \begin{pmatrix} 1 & 4 & 1 \\ 1 & 6 & 3 \\ 2 & 3 & 0 \end{pmatrix}, b = \begin{pmatrix} 4 \\ 6 \\ 9 \end{pmatrix} \quad (47.3.22)$$

We have already found  $A^{-1}$ , so we find that:

$$\mathbf{x} = A^{-1}\mathbf{b} = \frac{1}{6} \begin{pmatrix} -9 & 3 & 6 \\ 6 & -2 & -2 \\ -9 & 5 & 2 \end{pmatrix} \begin{pmatrix} 4 \\ 6 \\ 9 \end{pmatrix} = \frac{1}{6} \begin{pmatrix} 36 \\ -6 \\ 12 \end{pmatrix} = \begin{pmatrix} 6 \\ -1 \\ 2 \end{pmatrix} \quad (47.3.23)$$

so the solution to the system of equations is  $x = 6, y = -1, z = 2$ . ◀

### **Proposition (*Inverse matrix properties*)**

Let  $A, B$  be invertible matrices. Then:

- (a)  $(A^T)^{-1} = (A^{-1})^T$
- (b)  $(A^{-1})^{-1} = A$
- (c)  $AB$  is invertible and  $(AB)^{-1} = B^{-1}A^{-1}$

*Proof.* (a) We find that:

$$(A^T)(A^{-1})^T = (A^{-1}A)^T = \mathbb{1}^T = \mathbb{1} \quad (47.3.24)$$

and similarly:

$$(A^{-1})^T(A^T) = (AA^{-1})^T = \mathbb{1}^T = \mathbb{1} \quad (47.3.25)$$

(b) We find that:

$$(A^{-1})(A) = \mathbb{1} = (A)(A^{-1}) \quad (47.3.26)$$

(c) We find that:

$$ABB^{-1}A^{-1} = AA^{-1} = \mathbb{1} \quad (47.3.27)$$

and

$$B^{-1}A^{-1}AB = B^{-1}B = \mathbb{1} \quad (47.3.28)$$

as desired. ■

# Determinants

## 48.1 The determinant of a matrix

### Definition (Determinant)

The determinant  $\det : \mathbb{K}^n \rightarrow \mathbb{K}$  maps  $n$  vectors  $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_N \in \mathbb{K}^n$  to a scalar  $\det(\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_N) \in \mathbb{K}$ . It satisfies the following properties:

- (a)  $\det(\dots, \alpha\mathbf{a} + \beta\mathbf{b}, \dots) = \alpha \det(\dots, \mathbf{a}, \dots) + \beta \det(\dots, \mathbf{b}, \dots)$ .
- (b)  $\det(\dots, \mathbf{a}, \dots, \mathbf{b}, \dots) = -\det(\mathbf{b}, \dots, \mathbf{a}, \dots)$
- (c)  $\det(\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n) = 1$

The determinant of a matrix  $\mathbf{A}$  with column vectors  $\mathbf{A}^i$ ,  $1 \leq i \leq n$  is the determinant of these column vectors:

$$\det \mathbf{A} \equiv \det(\mathbf{A}^1, \mathbf{A}^2, \dots, \mathbf{A}^n) \quad (48.1.1)$$

Note that:

$$\det(\dots, \mathbf{0}, \dots) = \det(\dots, \mathbf{v} - \mathbf{v}, \dots) = \det(\dots, \mathbf{v}, \dots) - \det(\dots, \mathbf{v}, \dots) = 0 \quad (48.1.2)$$

and:

$$\det(\dots, \mathbf{a}, \dots, \mathbf{a}, \dots) = -\det(\dots, \mathbf{a}, \dots, \mathbf{a}, \dots) \implies \det(\dots, \mathbf{a}, \dots, \mathbf{a}, \dots) = 0 \quad (48.1.3)$$

So the determinant of a matrix with a zero column is null, and so is the determinant of a matrix with a repeated column vector.

Also, we have that:

$$\det(\dots, \mathbf{a}, \dots, \alpha\mathbf{a}, \dots) = \alpha \det(\dots, \mathbf{a}, \dots, \mathbf{a}, \dots) = 0 \quad (48.1.4)$$

so the determinant of a matrix with two columns that are proportional to each other will also be zero. We can combine these results to state that the determinant of a matrix where one row is a linear combination of some of the others is also zero.

We summarize these results in the next theorem:

### Proposition (Zero determinant matrices)

The following matrices have a zero determinant:

- (a) an entire row (or column) of zeros
- (b) a row (or column) that is a linear combination of other rows (or columns)

**Proposition (Special determinants)** The matrix of a diagonal matrix  $A$  is given by the product of its diagonal elements.

The matrix of an upper or lower triangular matrix  $A$  is also given by the product of its diagonal elements.

*Proof.* Consider a matrix  $A = \text{diag}(a_{11}, a_{22}, \dots, a_{nn})$ . Then:

$$\det A = \det(a_{11}\mathbf{e}_1, a_{22}\mathbf{e}_2, \dots, a_{nn}\mathbf{e}_n) = a_{11}a_{22}\dots a_{nn} \quad (48.1.5)$$

Instead, for an upper triangular matrix:

$$\det A = \det\left(a_{11}\mathbf{e}_1, a_{12}\mathbf{e}_1 + a_{22}\mathbf{e}_2, \dots, \sum_i a_{in}\mathbf{e}_i\right) \quad (48.1.6)$$

$$= \det\left(a_{11}\mathbf{e}_1, a_{12}\mathbf{e}_1, \dots, \sum_i^0 a_{in}\mathbf{e}_i\right) + \det\left(a_{11}\mathbf{e}_1, a_{22}\mathbf{e}_2, \dots, \sum_i a_{in}\mathbf{e}_i\right) \quad (48.1.7)$$

$$= \det(a_{11}\mathbf{e}_1, a_{22}\mathbf{e}_2, \dots, a_{nn}\mathbf{e}_n) \quad (48.1.8)$$

$$= a_{11}a_{22}\dots a_{nn} \quad (48.1.9)$$

as desired. ■

### Theorem (Determinant of matrix)

For a given  $n \times n$  matrix  $A$ , we have that:

$$\det A = \sum_{\sigma \in S_n} \text{sgn}(\sigma) \prod_{k=1}^n a_{\sigma(i_k)k} \quad (48.1.10)$$

*Proof.* We start by writing:

$$\det A = \det(\mathbf{A}^1, \mathbf{A}^2, \dots, \mathbf{A}^n) \quad (48.1.11)$$

$$= \det\left(\sum_{i_1} a_{i_1 1}\mathbf{e}_{i_1}, \sum_{i_2} a_{i_2 2}\mathbf{e}_{i_2}, \dots, \sum_k a_{i_n n}\mathbf{e}_{i_n}\right) \quad (48.1.12)$$

$$= \sum_{i_1} a_{i_1 1} \det\left(\mathbf{e}_{i_1}, \sum_{i_2} a_{i_2 2}\mathbf{e}_{i_2}, \dots, \sum_{i_n} a_{i_n n}\mathbf{e}_{i_n}\right) \quad (48.1.13)$$

$$= \sum_{i_1 i_2 \dots i_n} a_{i_1 1} a_{i_2 2} \dots a_{i_n n} \det(\mathbf{e}_{i_1}, \mathbf{e}_{i_2}, \dots, \mathbf{e}_{i_n}) \quad (48.1.14)$$

Now note that the only terms that survive out of this sum are  $i_1 \neq i_2 \neq i_3 \neq \dots \neq i_n$ . In other

words, the only terms surviving are all the permutations of  $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n$ , so we write:

$$\det \mathbf{A} = \sum_{\sigma \in S_n} \text{sgn}(\sigma) a_{\sigma(i_1)1} a_{\sigma(i_2)2} \dots a_{\sigma(i_n)n} \det(\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n) \quad (48.1.15)$$

$$= \sum_{\sigma \in S_n} \text{sgn}(\sigma) \prod_{k=1}^n a_{\sigma(i_k)k} \quad (48.1.16)$$

as desired. ■

### Proposition (Determinant properties)

Let's consider two  $n \times n$  matrices  $\mathbf{A}, \mathbf{B}$ . Then:

- (a)  $\det(\mathbf{A}^T) = \det \mathbf{A}$
- (b)  $\det(\mathbf{AB}) = \det \mathbf{A} \cdot \det \mathbf{B}$
- (c)  $\mathbf{A}$  is bijective  $\iff \det \mathbf{A} \neq 0$  and  $\det \mathbf{A}^{-1} = \frac{1}{\det \mathbf{A}}$

*Proof.* (a) Let  $a'_{ij}$  be the matrix elements of  $\mathbf{A}^T$  so that  $a'_{ij} = a_{ji}$ . Then:

$$\det \mathbf{A}^T = \sum_{\sigma \in S_n} \text{sgn}(\sigma) \prod_{k=1}^n a'_{\sigma(i_k)k} = \sum_{\sigma \in S_n} \text{sgn}(\sigma) \prod_{k=1}^n a_{k\sigma(i_k)} \quad (48.1.17)$$

$$= \sum_{\sigma \in S_n} \text{sgn}(\sigma^{-1}) \prod_{k=1}^n a_{\sigma^{-1}(i_k)k} = \sum_{\rho \in S_n} \text{sgn}(\rho^{-1}) \prod_{k=1}^n a_{\rho^{-1}(i_k)k} = \det \mathbf{A} \quad (48.1.18)$$

(b) Recall that  $(\mathbf{AB})_{ik} = \sum_j a_{ij} b_{jk}$  implying that:

$$(\mathbf{AB})^k = \sum_{ij} a_{ij} b_{jk} \mathbf{e}_i = \sum_j b_{jk} \mathbf{A}^j \quad (48.1.19)$$

Hence, we get that:

$$\det \mathbf{AB} = \det \left( \sum_{j_1} b_{j_1 1} \mathbf{A}^{j_1}, \sum_{j_2} b_{j_2 2} \mathbf{A}^{j_2}, \dots, \sum_{j_n} b_{j_n n} \mathbf{A}^{j_n} \right) \quad (48.1.20)$$

$$= \sum_{j_1, \dots, j_n} b_{j_1 1} b_{j_2 2} \dots b_{j_n n} \det(\mathbf{A}^{j_1}, \mathbf{A}^{j_2}, \dots, \mathbf{A}^{j_n}) \quad (48.1.21)$$

Again, we see that the only terms that survive are those where  $j_1 \neq j_2 \neq \dots \neq j_n$ , so we will get:

$$\det \mathbf{AB} = \sum_{\sigma \in S_n} \text{sgn}(\sigma) b_{\sigma(j_1)1} b_{\sigma(j_2)2} \dots b_{\sigma(j_n)n} \det(\mathbf{A}^1, \mathbf{A}^2, \dots, \mathbf{A}^n) \quad (48.1.22)$$

$$= \det \mathbf{A} \cdot \det \mathbf{B} \quad (48.1.23)$$

as desired.

(c) ( $\implies$ ) Suppose that  $\mathbf{A}$  is bijective, and thus invertible. Then:

$$\det(\mathbf{AA}^{-1}) = \det \mathbf{A} \cdot \det \mathbf{A}^{-1} = \det \mathbb{1} = 1 \quad (48.1.24)$$

implying that  $\det \mathbf{A} \neq 0$ .

( $\Leftarrow$ ) Suppose that  $\mathbf{A}$  is not bijective, so that  $\text{rk}(\mathbf{A}) < n$ . Therefore, there is at least one column vector, say  $\mathbf{A}^i$ , which may be written as a linear combination of some of the others:

$$\mathbf{A}^i = \sum_j \alpha_j \mathbf{A}^j \quad (48.1.25)$$

It then follows that

$$\det \mathbf{A} = \det(\mathbf{A}^1, \dots, \mathbf{A}^i, \dots, \mathbf{A}^n) = \det(\mathbf{A}^1, \dots, \mathbf{A}^i, \dots, \mathbf{A}^n) \quad (48.1.26)$$

$$= \det\left(\mathbf{A}^1, \dots, \sum_j \alpha_j \mathbf{A}^j, \dots, \mathbf{A}^n\right) \quad (48.1.27)$$

$$= \sum_j \alpha_j \det(\mathbf{A}^1, \dots, \mathbf{A}^j, \dots, \mathbf{A}^j, \dots, \mathbf{A}^n) \quad (48.1.28)$$

$$= 0 \quad (48.1.29)$$

as desired.

Using (48.1.24):

$$\det \mathbf{A}^{-1} = \frac{1}{\det \mathbf{A}} \quad (48.1.30)$$

as desired. ■

## 48.2 Laplace expansion

### Definition (Cofactor matrix)

Consider an  $n \times n$  matrix  $\mathbf{A}$ , then we define the associated  $(i, j)$  matrix  $\mathbf{A}(i, j)$  as the matrix  $\mathbf{A}$  with the  $i$ th row and  $j$ th column substituted with  $\mathbf{e}_i$  and  $\mathbf{e}_j^T$  respectively:

$$\mathbf{A}(i, j) = \begin{pmatrix} & & & \text{jth col} \\ & & 0 & \\ \mathbf{A} & \vdots & & \mathbf{A} \\ 0 \dots 0 & 1 & 0 \dots 0 \\ & \vdots & & \mathbf{A} \\ \mathbf{A} & \vdots & & 0 \end{pmatrix} \leftarrow i\text{th row} \quad (48.2.1)$$

The  $(i, j)$  cofactor coefficient is then defined as the determinant of  $\mathbf{A}(i, j)$

$$C_{ij} = \det(\mathbf{A}(i, j)) \quad (48.2.2)$$

The cofactor matrix is the matrix  $\mathbf{C}$  whose elements are  $C_{ij}$ . The cofactor expansion in the  $i$ th row is defined as:

$$\text{cof}_i \mathbf{A} = \sum_k a_{ik} C_{ik} = a_{i1} C_{i1} + a_{i2} C_{i2} + \dots + a_{in} C_{in} \quad (48.2.3)$$

It turns out that the cofactor matrix is especially important in evaluating the inverse of matrices. It is therefore important to be able to calculate the determinant of  $A_{ij}$  more easily.

**Proposition (Cofactor matrix calculation)**

Let  $A$  be a  $n \times n$  matrix and let  $\tilde{A}_{ij}$  be the matrix  $A$  with the  $i$ th row and  $j$ th column removed.

Then:

$$C_{ij} = (-1)^{i+j} \det(\tilde{A}(i, j)) \quad (48.2.4)$$

*Proof.* Note that the determinant only acquires a sign change when moving columns, and the same goes for rows since  $\det(A^T) = \det A$ . Hence, we may move the  $i$ th row and  $j$ th column to the first row and column respectively. To do so we must perform  $i - 1$  row exchanges followed by  $j - 1$  exchanges.<sup>1</sup> If we define  $B(i, j)$  to be the matrix with the  $i$ th and  $j$ th rows removed, then if:

$$B(i, j) \equiv \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & & & \\ \vdots & \tilde{A}(i, j) & & \\ 0 & & & \end{pmatrix} \quad (48.2.5)$$

$$\implies \det(B(i, j)) = (-1)^{i+j+2} \det(A(i, j)) = \det(\tilde{A}(i, j)) \quad (48.2.6)$$

so that:

$$\det(A(i, j)) = (-1)^{i+j} \det(\tilde{A}(i, j)) \quad (48.2.7)$$

■

This is a much easier formula to use when evaluating the cofactor matrix, since instead of evaluating the determinant of an  $n \times n$  matrix, we're evaluating the determinant of an  $(n - 1) \times (n - 1)$  matrix.

**Theorem (Laplace expansion)**

For a given  $n \times n$  matrix  $A$  with cofactor matrix  $C$ :

$$(\det A)\mathbb{1} = C^T A \quad (48.2.8)$$

so that:

$$\det A = \sum_k (-1)^{k+i} A_{kj} \det(\tilde{A}(k, i)) \quad (48.2.9)$$

*Proof.* We find that:

$$(C^T A)_{ij} = \sum_k C_{ki} A_{kj} \quad (48.2.10)$$

$$= \sum_k \det(A(k, i)) A_{kj} \quad (48.2.11)$$

$$= \sum_k A_{kj} \det(A^1 - A_{k1}\mathbf{e}_k, A^2 - A_{k2}\mathbf{e}_k, \dots, \mathbf{e}_k, \dots, A^n - A_{kn}\mathbf{e}_n) \quad (48.2.12)$$

<sup>1</sup>we can't just exchange the  $i$ th row with the first row, since this would alter the order of the rows, and would not give the matrix  $A$  with the  $i$ th row removed.

Now note that

$$\det(\mathbf{A}^1 - A_{k1}\mathbf{e}_k, \dots, \mathbf{e}_k, \dots) = \det(\mathbf{A}^1, \dots, \mathbf{e}_k, \dots) - A_{k1} \det(\mathbf{e}_k, \dots, \mathbf{e}_k, \dots) \quad (48.2.13)$$

$$= \det(\mathbf{A}^1, \dots, \mathbf{e}_k, \dots) \quad (48.2.14)$$

Repeating this process we find that:

$$(\mathbf{C}^T \mathbf{A})_{ij} = \sum_k A_{kj} \det(\mathbf{A}^1, \mathbf{A}^2, \dots, \mathbf{e}_k, \dots, \mathbf{A}^n) \quad (48.2.15)$$

$$= \sum_k \det(\mathbf{A}^1, \mathbf{A}^2, \dots, \mathbf{A}^j, \dots, \mathbf{A}^n) \quad (48.2.16)$$

$$= \det(\mathbf{A}^1, \mathbf{A}^2, \dots, \mathbf{A}^i, \dots, \mathbf{A}^n) \delta_{ij} \quad (48.2.17)$$

$$= \det \mathbf{A} \delta_{ij} \quad (48.2.18)$$

implying that:

$$\mathbf{C}^T \mathbf{A} = \det \mathbf{A} \mathbf{1} \quad (48.2.19)$$

as desired. ■

**Example.** Consider a  $3 \times 3$  matrix:

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix} \quad (48.2.20)$$

We calculate the cofactor expansion in the first column: Then we find that:

$$C_{11} = a_{22}a_{33} - a_{23}a_{32}, \quad C_{21} = a_{13}a_{32} - a_{12}a_{23}, \quad C_{31} = a_{12}a_{23}a_{22}a_{13} \quad (48.2.21)$$

so:

$$\det \mathbf{A} = C_{11}a_{11} + C_{21}a_{21} + C_{31}a_{31} \quad (48.2.22)$$

$$= a_{11}(a_{22}a_{33} - a_{23}a_{32}) + a_{21}(a_{13}a_{32} - a_{12}a_{23}) + a_{31}(a_{12}a_{23}a_{22}a_{13}) \quad (48.2.23)$$

Note that had we used the cofactor expansion in any other column (or also row), the d ◀

### Proposition (Cofactor orthogonality with rows)

We have that the  $j$ th row  $\mathbf{A}^j$  of a matrix  $\mathbf{A}$  is orthogonal to the  $i$ th row of its cofactor matrix  $\mathbf{C}$ , for  $i \neq j$ . So:

$$\mathbf{A}^j \cdot \mathbf{C}_i = \sum_k a_{jk} C_{ik} = 0 \quad (48.2.24)$$

*Proof.* For  $i \neq j$ , we have that

$$\mathbf{A}^j \cdot \mathbf{C}_i = \sum_k a_{jk} C_{ik} = a_{j1}C_{i1} + a_{j2}C_{i2} + \dots + a_{jn}C_{in} \quad (48.2.25)$$

Our goal is to find some other matrix  $\mathbf{B}$  whose cofactor expansion is like this. To find the form of

this matrix, we note that:

$$a_{j1}C_{i1} + a_{j2}C_{i2} + \dots + a_{jn}C_{in} = b_{i1}C_{i1} + b_{i2}C_{i2} + \dots + b_{in}C_{in} \quad (48.2.26)$$

Firstly, since  $C_{ij}$  takes all the entries of A without the  $i$ th row and  $j$ th column, it follows that for these cofactors to coincide with those of B, all the elements of B except for the  $i$ th row are identical to those of A.

The only change is that for  $a_{jk} = b_{ik}$ ,  $k = 1, 2, \dots, n$ , so the  $i$ th row of B is the  $j$ th row of A. However, note that the  $j$ th row of B must also be the  $j$ th row of A as we argued in the previous paragraph, so B has two repeated rows. Therefore, its determinant/cofactor expansion must vanish:

$$a_{j1}C_{i1} + a_{j2}C_{i2} + \dots + a_{jn}C_{in} = b_{i1}C_{i1} + b_{i2}C_{i2} + \dots + b_{in}C_{in} = 0 \quad (48.2.27)$$

■

### Definition (Adjoint matrix)

The adjoint  $\text{adj } A$  of a matrix A is the transpose of its cofactor matrix C:

$$\text{adj } A = C^T \implies (\det A)\mathbb{1} = (\text{adj } A)A \quad (48.2.28)$$

### Theorem (Inverse of matrix)

Let A be an invertible  $n \times n$  matrix, then:

$$A^{-1} = \frac{1}{\det A} \text{adj } A \quad (48.2.29)$$

*Proof.* We find that:

$$(\det A)\mathbb{1} = C^T A \implies (\det A)A^{-1} = \text{adj } A \implies A^{-1} = \frac{1}{\det A} \text{adj } A \quad (48.2.30)$$

as desired. ■

**Example.** Consider the  $2 \times 2$  matrix:

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}, ad - bc \neq 0 \quad (48.2.31)$$

We find that:

$$\det A = ad - bc \quad (48.2.32)$$

Furthermore:

$$C_{11} = d, C_{12} = -c, C_{21} = -b, C_{22} = d \quad (48.2.33)$$

giving the cofactor matrix:

$$C = \begin{pmatrix} d & -c \\ -b & d \end{pmatrix} \quad (48.2.34)$$

Consequently

$$\text{adj } \mathbf{A} = \begin{pmatrix} d & -b \\ -c & a \end{pmatrix} \Rightarrow \mathbf{A}^{-1} = \frac{1}{ad - bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix} \quad (48.2.35)$$

◀

## 48.3 Cramer's rule

With our newfound knowledge of determinants, we are now ready to formulate yet another method to solve linear systems. Previously we have discussed Gauss-Jordan elimination, as well as inverting matrices as methods of solutions.

### Theorem (Cramer's rule)

Consider a linear system of equations  $\mathbf{Ax} = \mathbf{b}$ . Let  $\mathbf{B}_i$  be the matrix  $\mathbf{A}$  with the  $i$ th column replaced with  $\mathbf{b}$ . Then the solution to the system is given by:

$$x_i \det \mathbf{A} = \det \mathbf{B}_i \quad (48.3.1)$$

*Proof.* We find that:

$$\det \mathbf{B}_i = \det(\mathbf{A}^1, \dots, \mathbf{b}, \dots, \mathbf{A}^n) \quad (48.3.2)$$

$$= \sum_i b_i \det(\mathbf{A}^1, \dots, \mathbf{e}_i, \dots, \mathbf{A}^n) \quad (48.3.3)$$

$$= \sum_i A_{ij} x_j \det(\mathbf{A}^1, \dots, \mathbf{e}_i, \dots, \mathbf{A}^n) \quad (48.3.4)$$

$$= \sum_i x_j \det(\mathbf{A}^1, \dots, \mathbf{A}^j, \dots, \mathbf{A}^n) \quad (48.3.5)$$

$$= x_i \det(\mathbf{A}^1, \dots, \mathbf{A}^i, \dots, \mathbf{A}^n) = x_i \det \mathbf{A} \quad (48.3.6)$$

as desired. Hence, if  $\mathbf{A}$  is invertible then the solutions are given by:

$$x_i = \frac{\det \mathbf{B}_i}{\det \mathbf{A}} \quad (48.3.7)$$

■

**Example.** Consider the following system:

$$\begin{cases} x + 2y + 3z = 0 \\ 2x + 3y + 4z = 1 \\ 3x + 4y + 6z = 2 \end{cases} \quad (48.3.8)$$

Then we see that:

$$\mathbf{A} = \begin{pmatrix} 1 & 2 & 3 \\ 2 & 3 & 4 \\ 3 & 4 & 6 \end{pmatrix}, \mathbf{b} = \begin{pmatrix} 0 \\ 1 \\ 2 \end{pmatrix} \quad (48.3.9)$$

so that:

$$\mathbf{B}_1 = \begin{pmatrix} 0 & 2 & 3 \\ 1 & 3 & 4 \\ 2 & 4 & 6 \end{pmatrix}, \quad \mathbf{B}_2 = \begin{pmatrix} 1 & 0 & 3 \\ 2 & 1 & 4 \\ 3 & 2 & 6 \end{pmatrix}, \quad \mathbf{B}_3 = \begin{pmatrix} 1 & 2 & 0 \\ 2 & 3 & 1 \\ 3 & 4 & 2 \end{pmatrix} \quad (48.3.10)$$

Now we evaluate the determinants using the Laplace expansion:

$$\det \mathbf{A} = (18 - 16) - 2(12 - 12) + 3(8 - 9) = -1 \quad (48.3.11)$$

$$\det \mathbf{B}_1 = -2(6 - 8) + 3(4 - 6) = -2 \quad (48.3.12)$$

$$\det \mathbf{B}_2 = (6 - 8) + 3(4 - 3) = 1 \quad (48.3.13)$$

$$\det \mathbf{B}_3 = (6 - 4) - 2(4 - 3) = 0 \quad (48.3.14)$$

so that:

$$x = 2, \quad y = -1, \quad z = 0 \quad (48.3.15)$$



# Inner product spaces

## 49.1 Inner products

### Definition (Inner products)

An inner product on a vector space  $V$  defined over  $\mathbb{K} = \mathbb{R}$  (or  $\mathbb{C}$ ) is a map:

$$\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{K} \quad (49.1.1)$$

satisfying:

- S1.  $\langle \mathbf{v}, \mathbf{w} \rangle = \langle \mathbf{w}, \mathbf{v} \rangle$  for a symmetric scalar product  
(or  $\langle \mathbf{v}, \mathbf{w} \rangle = \langle \mathbf{w}, \mathbf{v} \rangle^*$  for hermitian inner product)
  - S2.  $\langle \mathbf{v}, \alpha \mathbf{u} + \beta \mathbf{w} \rangle = \alpha \langle \mathbf{v}, \mathbf{u} \rangle + \beta \langle \mathbf{v}, \mathbf{w} \rangle$
  - S3.  $\langle \mathbf{v}, \mathbf{v} \rangle \geq 0$  with equality holding only for  $\mathbf{v} = 0$
- for all  $\mathbf{v}, \mathbf{w}, \mathbf{u} \in V, \alpha, \beta \in \mathbb{K}$ .

### Example.

- i. Minkowski product: for  $\mathbf{v} = v^\mu \mathbf{e}_\mu \in \mathbb{R}^4$  and  $\mathbf{w} = w^\mu \mathbf{e}_\mu \in \mathbb{R}^4$ , the Minkowski product is defined as:

$$\langle \mathbf{v}, \mathbf{w} \rangle = \mathbf{v}^T \boldsymbol{\eta} \mathbf{w} \quad (49.1.2)$$

where  $\boldsymbol{\eta} = \text{diag}(1, -1, -1, -1)$ .

- ii. Functional products: for  $f, g \in C^0([a, b])$ , we define their inner product as:

$$\langle f, g \rangle = \int_a^b f(x)^* g(x) dx \quad (49.1.3)$$



### Proposition (Inner product is well-defined)

Let  $f, g \in \text{End}(V)$  be two endomorphisms on  $V$ , then:

$$\langle \mathbf{v}, f(\mathbf{w}) \rangle = \langle \mathbf{v}, g(\mathbf{w}) \rangle, \forall \mathbf{v}, \mathbf{w} \in V \implies f = g \quad (49.1.4)$$

*Proof.* Let  $\mathcal{B} = \{\mathbf{e}_i\}$  be an orthonormal basis for  $V$ . It follows that for any  $\mathbf{v}, \mathbf{w} \in V$  (we assume a

hermitian inner product, the proof for a symmetric scalar product is similar):

$$\langle \mathbf{v}, f(\mathbf{w}) \rangle = \langle \mathbf{v}, g(\mathbf{w}) \rangle \iff \langle \mathbf{v}, f(\mathbf{w}) - g(\mathbf{w}) \rangle = 0 \quad (49.1.5)$$

$$\iff \sum_{ij} v_i^* w_j \langle \mathbf{e}_i, f(\mathbf{e}_j) - g(\mathbf{e}_j) \rangle = 0 \quad (49.1.6)$$

Since this applies for any  $\mathbf{v}, \mathbf{w}$ , we must therefore have that:

$$\langle \mathbf{e}_i, h(\mathbf{e}_j) \rangle \equiv \langle \mathbf{e}_i, f(\mathbf{e}_j) - g(\mathbf{e}_j) \rangle = 0, \forall i, j = 1, 2, \dots \quad (49.1.7)$$

where we defined  $h = f - g \in \text{End}(V)$ . Consequently,  $h(\mathbf{e}_j)$  cannot have any component along any of the basis vectors for  $V$ , and must therefore be zero. Hence  $f = g$  as desired. ■

**Definition (Orthogonal complement)** For a subspace  $W$  of  $V$ , we define the orthogonal complement  $W^\perp$  as:

$$W^\perp = \{\mathbf{v} \in V : \langle \mathbf{v}, \mathbf{w} \rangle = 0, \forall \mathbf{w} \in W\} \quad (49.1.8)$$

## 49.2 Projectors

**Proposition (Properties of orthogonal complements)** For a subspace  $W$  of  $V$  with orthogonal complement  $W^\perp$ , we have that:

- (i)  $W^\perp \subset V$
- (ii)  $W \cap W^\perp = \{\mathbf{0}\}$
- (iii)  $\dim W + \dim W^\perp = \dim V$

*Proof.* (i) Trivial.

(ii),(iii) Firstly note that  $W \oplus W^\perp = V$ . Indeed, given a vector  $\mathbf{v} \in V$ , we can decompose it as:

$$\mathbf{v} = \underbrace{\langle \mathbf{w}, \mathbf{v} \rangle}_{\in W} \mathbf{w} + \underbrace{(\mathbf{v} - \langle \mathbf{w}, \mathbf{v} \rangle \mathbf{w})}_{\in W^\perp} \quad (49.2.1)$$

Therefore, we must have that  $\dim W + \dim W^\perp = \dim V$ , as well as  $W \cap W^\perp = \{\mathbf{0}\}$  from the properties of direct sums. ■

**Definition ((Orthogonal) Projection operators)** Let  $V = U \oplus W$ . We define a projection from  $V$  to  $W$  as a map  $\Pi$  satisfying:

$$\Pi : V \rightarrow W \quad (49.2.2)$$

$$\mathbf{u} + \mathbf{w} \rightarrow \mathbf{w} \quad (49.2.3)$$

with  $\mathbf{u} + \mathbf{w}$ . Clearly,  $\Pi^2 = \Pi$  so it is idempotent.

If  $U = V^\perp$  then we say that  $\Pi$  is an orthogonal projection operator.

**Theorem (All idempotent maps are projective)**

All idempotent maps are projections.

*Proof.* Let  $\Pi$  be an idempotent map so that  $\Pi^2 = \Pi$ . Let  $\mathbf{v}$ , then clearly we have that:

$$\mathbf{v} = \Pi(\mathbf{v}) + (\mathbf{v} - \Pi(\mathbf{v})) = \Pi(\mathbf{v}) + (\mathbb{1} - \Pi)(\mathbf{v}) \quad (49.2.4)$$

Now let us define  $W = \{\Pi(\mathbf{v}) \mid \forall \mathbf{v} \in V\} = \text{Im}(\Pi)$  and  $U = \{\mathbf{v} - \Pi(\mathbf{v}) : \forall \mathbf{v} \in V\} = \text{Im}(\mathbb{1} - \Pi)$ . Since these are both images of linear transformations, we have that  $U, W$  are subspaces of  $V$ . Note also that  $\Pi(\mathbf{w}) = 0$  for all  $\mathbf{w} \in W$  and  $\Pi(\mathbf{u}) = \mathbf{u}$ . Consequently  $V = U \otimes W$ , with:

$$\Pi(\mathbf{v}) = \Pi(\mathbf{u} + \mathbf{w}) = \Pi(\mathbf{u}) + \Pi(\mathbf{w}) = \mathbf{u} \quad (49.2.5)$$

proving that  $\Pi$  is indeed a projector. ■

### 49.3 Inner products and matrices

**Definition (Adjoint, Hermitian and Unitary linear map)**

For a linear map  $f \in \text{End}(V)$  on  $V$  defined over  $\mathbb{K}$ , its adjoint map  $f^\dagger \in \text{End}(V)$  is defined so that it satisfies:

$$\langle \mathbf{v}, f(\mathbf{w}) \rangle = \langle f^\dagger(\mathbf{v}), \mathbf{w} \rangle, \quad \forall \mathbf{v}, \mathbf{w} \in V \quad (49.3.1)$$

If  $\mathbb{K} = \mathbb{R}$  (or  $\mathbb{C}$ ), then a symmetric (or hermitian) linear map  $f$  satisfies  $f = f^\dagger$ , while an orthogonal (or unitary) linear map  $U$  satisfies  $U^{-1} = U^\dagger$ .

**Proposition (Adjoint map properties)**

For a linear map  $f \in \text{End}(V)$ , the following must hold:

- (i)  $f^\dagger$  is unique
- (ii)  $(f^\dagger)^\dagger = f$
- (iii)  $(f \circ g)^\dagger = g^\dagger \circ f^\dagger$
- (iv)  $(f^{-1})^\dagger = (f^\dagger)^{-1}$

- (i) Let  $g, h$  be adjoint maps of  $f$ . Then  $\langle \mathbf{v}, f(\mathbf{w}) \rangle = \langle g(\mathbf{v}), \mathbf{w} \rangle = \langle \mathbf{v}, f(\mathbf{w}) \rangle = \langle h(\mathbf{v}), \mathbf{w} \rangle$  which by previous proposition implies that  $g = h$ .
- (ii) For all  $\mathbf{v}, \mathbf{w} \in V$ , we have that  $\langle \mathbf{v}, f(\mathbf{w}) \rangle = \langle f^\dagger(\mathbf{v}), \mathbf{w} \rangle = \langle \mathbf{v}, (f^\dagger)^\dagger(\mathbf{w}) \rangle$  which by the same proposition as before implies that  $(f^\dagger)^\dagger = f$ .
- (iii)  $\langle \mathbf{v}, (f \circ g)^\dagger(\mathbf{w}) \rangle = \langle f(g(\mathbf{v})), \mathbf{w} \rangle = \langle g(\mathbf{v}), f^\dagger(\mathbf{w}) \rangle = \langle \mathbf{v}, (g^\dagger \circ f^\dagger)(\mathbf{w}) \rangle$ .
- (iv)  $f \circ f^{-1} = \text{id}_V \implies (f^{-1})^\dagger \circ f^\dagger = \text{id}_V \implies (f^{-1})^\dagger = (f^\dagger)^{-1}$  where we used (iii) to take the adjoint of both sides in the first implication.

Note that if we have an orthonormal basis  $\mathcal{B} = \{\mathbf{e}_i\}$  for  $V$ , then given any  $f \in \text{End}(V)$ , we have that:

$$f(\mathbf{e}_i) = \sum_j A_{ij} \mathbf{e}_j \implies A_{ij} = \langle \mathbf{e}_j, f(\mathbf{e}_i) \rangle \quad (49.3.2)$$

so we can use inner products to find the matrix elements of linear maps. It is easy to see that if  $f \in \text{End}(V)$  defined over  $\mathbb{C}$  has matrix elements  $A_{ij}$  in a given basis, then its adjoint  $f^\dagger$  will have matrix elements  $B_{ij}$  in the same basis given by:

$$B_{ij} = \langle \mathbf{e}_i, f^\dagger(\mathbf{e}_j) \rangle = \langle f(\mathbf{e}_i), \mathbf{e}_j \rangle = \langle \mathbf{e}_j, f(\mathbf{e}_i) \rangle^* = A_{ji}^* \quad (49.3.3)$$

Hence the matrices  $A, A^\dagger$  representing  $f, f^\dagger$  respectively satisfy  $A^\dagger = (A^*)^T$ .

It follows that Hermitian maps have matrix representations  $A = A^T = (A^*)^T$  and Unitary maps have matrix representations  $A^{-1} = A^T = (A^*)^T$ .

**Proposition (Alternative definition of Unitarity)** Let  $U \in \text{End}(V)$  is a unitary map, then  $\langle U(\mathbf{v}), U(\mathbf{w}) \rangle = \langle \mathbf{v}, \mathbf{w} \rangle$  for all  $\mathbf{v}, \mathbf{w} \in V$  is an equivalent definition.

*Proof.* We find  $\langle U(\mathbf{v}), U(\mathbf{w}) \rangle = \langle \mathbf{v}, (U^\dagger \circ U)(\mathbf{w}) \rangle = \langle \mathbf{v}, \mathbf{w} \rangle$ . Hence by the well-definedness of inner products,  $U^\dagger \circ U = \text{id}_V \iff U^{-1} = U^\dagger$ . ■

## 49.4 Bilinear and Sesquilinear forms

In the previous section we looked at properties of inner product spaces over real or complex vector spaces. It turns out that when we remove the condition for the inner product to be positive semi-definite we get some interesting new forms.

**Definition (Bilinear/Sesquilinear form)** A Bilinear form on a vector space  $V$  defined over  $\mathbb{R}$  is a map  $T : V \times V \rightarrow \mathbb{R}$  linear in both of its terms:

$$T(\alpha\mathbf{v} + \beta\mathbf{w}, \mathbf{u}) = \alpha T(\mathbf{v}, \mathbf{u}) + \beta T(\mathbf{w}, \mathbf{u}) \quad (49.4.1)$$

$$T(\mathbf{u}, \alpha\mathbf{v} + \beta\mathbf{w}) = \alpha T(\mathbf{u}, \mathbf{v}) + \beta T(\mathbf{u}, \mathbf{w}) \quad (49.4.2)$$

A bilinear form is symmetric if  $T(\mathbf{v}, \mathbf{w}) = T(\mathbf{w}, \mathbf{v})$ .

A Sesquilinear form on a vector space  $V$  defined over  $\mathbb{C}$  is a map  $T : V \times V \rightarrow \mathbb{C}$  linear in both of its terms:

$$T(\alpha\mathbf{v} + \beta\mathbf{w}, \mathbf{u}) = \alpha^* T(\mathbf{v}, \mathbf{u}) + \beta^* T(\mathbf{w}, \mathbf{u}) \quad (49.4.3)$$

$$T(\mathbf{u}, \alpha\mathbf{v} + \beta\mathbf{w}) = \alpha T(\mathbf{u}, \mathbf{v}) + \beta T(\mathbf{u}, \mathbf{w}) \quad (49.4.4)$$

A sesquilinear form is hermitian if  $T(\mathbf{v}, \mathbf{w}) = T(\mathbf{w}, \mathbf{v})$ .

**Example.** An important example of a Bilinear form often used in Special relativity is:

$$T(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{A} \mathbf{y} \quad (49.4.5)$$

where  $\mathbf{x}, \mathbf{y} \in V$  and  $\mathbf{A} \in \text{Mat}_n(V)$ . Indeed, we prove linearity in the first argument as follows:

$$T(\alpha\mathbf{x} + \beta\mathbf{y}, \mathbf{z}) = (\alpha\mathbf{x} + \beta\mathbf{y})^T \mathbf{A}\mathbf{z} \quad (49.4.6)$$

$$= (\alpha\mathbf{x}^T + \beta\mathbf{y}^T) \mathbf{A}\mathbf{z} \quad (49.4.7)$$

$$= \alpha\mathbf{x}^T \mathbf{A}\mathbf{z} + \beta\mathbf{y}^T \mathbf{A}\mathbf{z} \quad (49.4.8)$$

$$= \alpha T(\mathbf{x}, \mathbf{z}) + \beta T(\mathbf{y}, \mathbf{z}) \quad (49.4.9)$$

The proof is similar for the linearity in second argument. We can extend this example to sesquilinear forms by defining:

$$T(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\dagger \mathbf{A} \mathbf{y} \quad (49.4.10)$$

where  $\mathbf{x}, \mathbf{y} \in V$  and  $\mathbf{A} \in \text{Mat}_n(V)$ . Here  $\dagger$  denotes conjugate transposition. We again prove linearity in the first argument as follows:

$$T(\alpha\mathbf{x} + \beta\mathbf{y}, \mathbf{z}) = (\alpha\mathbf{x} + \beta\mathbf{y})^\dagger \mathbf{A}\mathbf{z} \quad (49.4.11)$$

$$= (\alpha\mathbf{x}^\dagger + \beta\mathbf{y}^\dagger) \mathbf{A}\mathbf{z} \quad (49.4.12)$$

$$= \alpha\mathbf{x}^\dagger \mathbf{A}\mathbf{z} + \beta\mathbf{y}^\dagger \mathbf{A}\mathbf{z} \quad (49.4.13)$$

$$= \alpha T(\mathbf{x}, \mathbf{z}) + \beta T(\mathbf{y}, \mathbf{z}) \quad (49.4.14)$$

Note also that if  $\mathbf{A}$  is a symmetric matrix then:

$$T(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \mathbf{A}) \mathbf{y} = (\mathbf{A}^T \mathbf{x})^T \mathbf{y} = \mathbf{y}^T \mathbf{A} \mathbf{x} \quad (49.4.15)$$

so  $T$  is a symmetric bilinear form. We can extend this result to sesquilinear forms quite easily by letting  $\mathbf{A}$  be hermitian. Then

$$T(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^\dagger \mathbf{A}) \mathbf{y} = (\mathbf{A}^\dagger \mathbf{x})^\dagger \mathbf{y} = \mathbf{y}^\dagger \mathbf{A} \mathbf{x} \quad (49.4.16)$$

◀

It turns out that all bilinear/sesquilinear forms may be expressed in the form of the previous example. Indeed, let  $T$  be a (bilinear) sesquilinear form on a (real) complex vector space  $V$ . Let  $\{\mathbf{e}_i\}$  be an ordered basis of  $V$ , then:

$$T(\mathbf{x}, \mathbf{y}) = \sum_{ij} x_i^* y_j T(\mathbf{e}_i, \mathbf{e}_j) \quad (49.4.17)$$

If we let  $A_{ij} = T(\mathbf{e}_i, \mathbf{e}_j)$  then clearly:

$$\mathbf{x}^\dagger \mathbf{A} \mathbf{y} = \sum_{ij} x_i^* A_{ij} y_j = \sum_{ij} x_i^* y_j T(\mathbf{e}_i, \mathbf{e}_j) = T(\mathbf{x}, \mathbf{y}) \quad (49.4.18)$$

as desired.

**Theorem (Matrix representation of forms)** A bilinear/sesquilinear form  $T$  over a real/complex vector space  $V$  has an associated matrix representation in a given basis  $\{\mathbf{e}_i\}$

of  $V$ :

$$A = \begin{pmatrix} T(\mathbf{e}_1, \mathbf{e}_1) & T(\mathbf{e}_1, \mathbf{e}_2) & \dots \\ T(\mathbf{e}_2, \mathbf{e}_1) & T(\mathbf{e}_2, \mathbf{e}_2) & \dots \\ \vdots & \ddots & \vdots \end{pmatrix} \quad (49.4.19)$$

Analogously to linear maps, one can also perform changes of basis for sesquilinear/bilinear forms. Suppose that in the basis  $\{\mathbf{e}_i\}$  the form  $T$  is represented by  $A$  so that

$$A_{ij} = T(\mathbf{e}_i, \mathbf{e}_j) \quad (49.4.20)$$

Let us introduce a new basis  $\{\mathbf{e}'_i\}$  such that:

$$\mathbf{e}'_i = \sum_m c_{mi} \mathbf{e}_m \quad (49.4.21)$$

then we find

$$A'_{ij} = T(\mathbf{e}'_i, \mathbf{e}'_j) = \sum_{mn} c_{mi}^* c_{nj} T(\mathbf{e}_m, \mathbf{e}_n) = \sum_{mn} c_{mi}^* A_{mn} c_{nj} \quad (49.4.22)$$

Consequently, if we define a change of basis matrix  $P$  with components  $P_{mn} = c_{mn}$  then we get:

$$A' = P^\dagger A P \quad (49.4.23)$$

We interpret this result as usual.  $P$  converts our vector from the original unprimed basis to the new primed basis:

$$\mathbf{x} = \sum_i x'_i \mathbf{e}'_i = \sum_{ij} x'_i c_{ji} \mathbf{e}_j = \sum_j x_j \mathbf{e}_j \quad (49.4.24)$$

where

$$x'_j = \sum_i x'_i c_{ji} \iff [\mathbf{x}]' = P[\mathbf{x}] \quad (49.4.25)$$

Therefore if we want to calculate the form  $\mathbf{x}^\dagger A \mathbf{y}$  then we need a  $P^\dagger$  to the left of  $A$  to convert the components of  $\mathbf{x}^\dagger$  to the primed basis, and a  $P$  to the right of  $A$  to convert the components of  $\mathbf{y}$ .

# Eigen-everything

## 50.1 Finding eigenvalues and eigenvectors

### **Definition (Eigenvalue and eigenvector)**

Let  $f : V \rightarrow V$  be a linear map on  $V$  over  $\mathbb{F}$ . We say that  $\lambda \in \mathbb{F}$  is an eigenvalue of  $f$  if  $\exists \mathbf{v} \in V$ , s.t.  $\mathbf{v} \neq \mathbf{0}$ , known as an eigenvector, such that:

$$f(\mathbf{v}) = \lambda \mathbf{v} \quad (50.1.1)$$

The eigenspace of  $\lambda$  is defined as:

$$\text{Eig}_f(\lambda) \equiv \text{Ker}(f - \lambda \text{id}_V) \subseteq V \quad (50.1.2)$$

For a map  $f$  to have non-trivial eigenvalues, we require that:

$$\dim \text{Eig}_f(\lambda) = \text{Ker}(f - \lambda \text{id}_V) > 0 \implies \text{rk}(f - \lambda \text{id}_V) < \dim V \quad (50.1.3)$$

This is equivalent, by the proposition on linear systems and invertibility, to setting

$$\det(f - \lambda \text{id}_V) = 0 \quad (50.1.4)$$

We could have also seen this by noting that if  $f - \lambda \text{id}_V$  were invertible, then there would only be one  $\mathbf{v}$  in  $\text{Ker}(f - \lambda \text{id}_V)$ , which must be  $\mathbf{0}$ . Therefore  $f - \lambda \text{id}_V$  cannot be invertible, yielding (50.1.4).

### **Definition (Characteristic polynomial)**

The characteristic polynomial of a map  $f : V \rightarrow V$  is defined as:

$$\chi_f(\lambda) = \det(f - \lambda \text{id}_V) \quad (50.1.5)$$

To find the eigenvalues of a matrix, it suffices to:

- (i) Compute its characteristic polynomial.
- (ii) Find the roots  $\lambda$  of  $\chi_f(\lambda)$ .
- (iii) For each solution  $\lambda$ , find the corresponding eigenspace by solving:

$$(f - \lambda \text{id}_V) \mathbf{v} = 0 \quad (50.1.6)$$

using one of the methods introduced for solving linear systems.

**Example.** Let us find the eigenvalues and eigenspaces of

$$A = \begin{pmatrix} 4 & 0 & 4 \\ 0 & 4 & 4 \\ 4 & 4 & 8 \end{pmatrix} \quad (50.1.7)$$

Its characteristic polynomial is:

$$\chi_A(\lambda) = \begin{vmatrix} 4 - \lambda & 0 & 4 \\ 0 & 4 - \lambda & 4 \\ 4 & 4 & 8 - \lambda \end{vmatrix} \quad (50.1.8)$$

$$= (4 - \lambda)((4 - \lambda)(8 - \lambda) - 16) - 16(4 - \lambda) \quad (50.1.9)$$

$$= (4 - \lambda)(\lambda^2 - 12\lambda) \quad (50.1.10)$$

$$= \lambda(4 - \lambda)(\lambda - 12) \quad (50.1.11)$$

Clearly, the solutions to  $\chi_A(\lambda) = 0$  are  $\lambda = 0, 4, 12$ .

For  $\lambda_1 = 0$ , we need:

$$\begin{pmatrix} 4 & 0 & 4 \\ 0 & 4 & 4 \\ 4 & 4 & 8 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \implies x = -z, y = -z \quad (50.1.12)$$

giving an eigenspace:

$$\text{Eig}_A(\lambda_1) = \left\{ k \begin{pmatrix} -1 \\ -1 \\ 1 \end{pmatrix}, \forall k \in \mathbb{R}^* \right\} \quad (50.1.13)$$

Similarly, for  $\lambda_1 = 4$ , we need:

$$\begin{pmatrix} 0 & 0 & 4 \\ 0 & 0 & 4 \\ 4 & 4 & 4 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \implies x = -y, z = 0 \quad (50.1.14)$$

giving an eigenspace:

$$\text{Eig}_A(\lambda_1) = \left\{ k \begin{pmatrix} 1 \\ -1 \\ 0 \end{pmatrix}, \forall k \in \mathbb{R}^* \right\} \quad (50.1.15)$$

Finally, for  $\lambda_1 = 12$ , we need:

$$\begin{pmatrix} -8 & 0 & 4 \\ 0 & -8 & 4 \\ 4 & 4 & -4 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \implies z = 2x, y = x \quad (50.1.16)$$

giving an eigenspace:

$$\text{Eig}_A(\lambda_1) = \left\{ k \begin{pmatrix} 1 \\ 1 \\ 2 \end{pmatrix}, \forall k \in \mathbb{R}^* \right\} \quad (50.1.17)$$


**Proposition (Characteristic polynomial)**

For a matrix  $A \in \text{Mat}_n(\mathbb{K})$  with characteristic polynomial  $\chi_A(\lambda) = \sum_{i=0}^n c_i \lambda^i$ , the following hold:

- (i)  $\chi_{PAP^{-1}} = \chi_A$  for all  $P \in \text{Mat}_n(\mathbb{K})$
- (ii)  $c_n = (-1)^n, c_{n-1} = (-1)^{n-1} \text{tr}A, c_0 = \det A$

*Proof.* (i) We have that:

$$\det(PAP^{-1} - \lambda \mathbb{1}) = \det(P(A - \lambda \mathbb{1})P^{-1}) = \det(A - \lambda \mathbb{1}) \quad (50.1.18)$$

as desired.

(ii) We have that:

$$c_0 = \chi_A(0) = \det A \quad (50.1.19)$$

Furthermore

$$\chi_A(\lambda) = \prod_{i=1}^n (A_{ii} - \lambda) + o(\lambda^{n-2}) \quad (50.1.20)$$

$$= (-1)^n \lambda^n + (-1)^{n-1} \sum_{i=1}^n A_{ii} + o(\lambda^{n-1}) \quad (50.1.21)$$

$$= (-1)^n \lambda^n + (-1)^{n-1} \text{tr}A + o(\lambda^{n-1}) \quad (50.1.22)$$

implying that  $c_n = (-1)^n$  and  $c_{n-1} = (-1)^{n-1} \text{tr}A$ .



## 50.2 Matrix diagonalization

**Definition (Diagonalized map)**

A linear map  $f : V \rightarrow V$  can be diagonalized iff there exists a basis of  $V$  which makes the matrix representation of  $f$  diagonal, that is, if it is similar to a diagonal matrix.

**Theorem (Diagonalizability)**

A linear map  $f : V \rightarrow V$  can be diagonalised iff there exists a basis of  $V$  consisting of eigenvectors of  $f$ . In this basis, the matrix representation of  $f$  is  $\text{diag}(\lambda_i)$  where  $\lambda_i$  are the eigenvalues of  $f$ .

*Proof.* ( $\implies$ ) Suppose that  $f$  can be diagonalized into the form  $\text{diag}(c_i)$  in some basis  $\{\mathbf{v}_i\}$ . This implies that:

$$f(\mathbf{v}_i) = \sum_j A_{ji} \mathbf{v}_j = \sum_j \delta_{ji} c_i \mathbf{v}_i \quad (50.2.1)$$

which gives  $f(\mathbf{v}_i) = c_i \mathbf{v}_i$ , as desired.

( $\Leftarrow$ ) Let  $P = (\mathbf{v}_1 \ \mathbf{v}_2 \ \dots \ \mathbf{v}_n)$  be the transition matrix from some basis  $\mathcal{B}$  to the set  $\{\mathbf{v}_i\}$  of eigenvectors (with respective eigenvalues  $\lambda_i$ ). Of course, to perform the required change of basis we need the set of eigenvectors  $\{\mathbf{v}_i\}$  to form a basis of  $V$ . Then, we find that the matrix representation  $A$  of  $f$  in  $\mathcal{B}$ :

$$A' = P^{-1}AP = P^{-1}A(\mathbf{v}_1 \ \mathbf{v}_2 \ \dots \ \mathbf{v}_n) \quad (50.2.2)$$

$$= P^{-1}(A\mathbf{v}_1 \ A\mathbf{v}_2 \ \dots \ A\mathbf{v}_n) \quad (50.2.3)$$

$$= P^{-1}(\lambda_1\mathbf{v}_1 \ \lambda_2\mathbf{v}_2 \ \dots \ \lambda_n\mathbf{v}_n) \quad (50.2.4)$$

$$= P^{-1}P\text{diag}(\lambda_i) \quad (50.2.5)$$

$$= \text{diag}(\lambda_i) \quad (50.2.6)$$

Thus, the new matrix in the basis  $P$  is indeed diagonal, with entries equal to the eigenvalues. ■

Recall that the traces and determinants of a matrix are independent of the chosen basis, so if  $A$  is diagonalizable then:

$$\text{tr}A = \sum_i \lambda_i, \ \det A = \prod_i \lambda_i \quad (50.2.7)$$

This gives us a nice way to check if any arithmetic mistakes have been made in evaluating the eigenvalues.

We also gain some geometrical insight behind what eigenvectors really are. Indeed, if we consider any diagonalizable linear map  $f$  acting on vectors in  $V$ , it follows that its action will be to stretch  $V$  by a factor of  $\lambda_i$  along  $\mathbf{v}_i$ . This can be readily verified by looking at the diagonalized form of  $f$ , and noting that the  $i$ th column of a matrix representation gives the vector that the corresponding  $i$ th basis vector gets mapped to.

In other words, the eigenvectors are vectors which, when acted upon by a linear map  $f$ , only change by a phase, but still point in the same "direction".

**Example.** Let's consider the one of the Pauli matrices

$$\sigma = \begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix} \quad (50.2.8)$$

Its characteristic equation is:

$$\chi_\sigma(\lambda) = \lambda^2 - 1 = 0 \implies \lambda = \pm 1 \quad (50.2.9)$$

For  $\lambda_1 = 1$  we find that:

$$\begin{pmatrix} -1 & -i \\ i & -1 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \implies x = -iy \quad (50.2.10)$$

Its eigenspace is:

$$\text{Eig}_\sigma(\lambda_1) = \left\{ k \begin{pmatrix} -i \\ 1 \end{pmatrix} : k \in \mathbb{R}^* \right\} \quad (50.2.11)$$

Similarly, for  $\lambda_2 = -1$  we find that:

$$\begin{pmatrix} 1 & -i \\ i & 1 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \implies x = iy \quad (50.2.12)$$

Its eigenspace is:

$$\text{Eig}_\sigma(\lambda_1) = \left\{ k \begin{pmatrix} i \\ 1 \end{pmatrix} : k \in \mathbb{R}^* \right\} \quad (50.2.13)$$

We therefore choose the eigenbasis with  $k = 1$ , whose transition matrix is:

$$P = \begin{pmatrix} -i & i \\ 1 & 1 \end{pmatrix} \quad (50.2.14)$$

whose inverse is (since  $\det P = -i - i = -2i$ ):

$$P^{-1} = -\frac{1}{2i} \begin{pmatrix} 1 & -i \\ -1 & -i \end{pmatrix} = \frac{i}{2} \begin{pmatrix} 1 & -i \\ -1 & -i \end{pmatrix} = \frac{1}{2} \begin{pmatrix} i & 1 \\ -i & 1 \end{pmatrix} \quad (50.2.15)$$

We find that:

$$\sigma' = \frac{1}{2} \begin{pmatrix} i & 1 \\ -i & 1 \end{pmatrix} \begin{pmatrix} 0 & -i \\ i & 0 \end{pmatrix} \begin{pmatrix} -i & i \\ 1 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix} \quad (50.2.16)$$

◀

### Theorem (Simultaneous diagonalization)

Two linear maps  $f, g$  are simultaneously diagonalizable, that is, they are diagonalized by the same matrix  $P$ , iff they commute:

$$[f, g] \equiv f \circ g - g \circ f = 0 \quad (50.2.17)$$

*Proof.* ( $\implies$ ) Suppose that  $f, g$  are simultaneously diagonalizable, so that they share a set of eigenvectors  $\mathbf{v}_i$  with eigenvalues  $\lambda_i$  and  $\lambda'_i$  respectively. Then:

$$[f, g](\mathbf{v}_i) = \lambda_i g(\mathbf{v}_i) - \lambda'_i f(\mathbf{v}_i) = \lambda_i \lambda'_i - \lambda_i \lambda'_i = 0 \quad (50.2.18)$$

Any vector can be expanded as a linear combination of  $\mathbf{v}_i$ , so  $[f, g](\mathbf{v}) = 0$  holds for all  $\mathbf{v} \in V$ . It follows that  $[f, g] = 0$ , the two maps commute.

( $\impliedby$ ) Suppose that  $f, g$  commute, and suppose that  $f$  has eigenvectors  $\mathbf{v}_i$  with eigenvalues  $\lambda_i$ . Then:

$$(f \circ g)(\mathbf{v}_i) = f(g(\mathbf{v}_i)) = (g \circ f)(\mathbf{v}_i) = \lambda_i g(\mathbf{v}_i) \quad (50.2.19)$$

implying that  $g(\mathbf{v}_i) \in \text{Eig}_f(\lambda_i)$ . For non-degenerate eigenvalues, this means that  $g(\mathbf{v}_i) = \alpha \mathbf{v}_i$  for some non-zero  $\alpha$ . Hence  $\mathbf{v}_i$  is an eigenvector of both  $f$  and  $g$ , the two maps are simultaneously diagonalizable. ■

## 50.3 Orthogonal diagonalization

**Theorem (Spectral properties for hermitian matrices)**

Let  $A$  be a hermitian matrix (so that  $A^\dagger = A$ ). Then all its eigenvalues are real, and the eigenvectors corresponding to distinct eigenvalues are orthogonal.

*Proof.* Let  $\mathbf{v}$  be an eigenvector of  $A$  with eigenvalue  $\lambda$ . Then, it follows that

$$\langle \mathbf{v}, A\mathbf{v} \rangle = \lambda = \langle A\mathbf{v}, \mathbf{v} \rangle = \lambda^* \implies \lambda \in \mathbb{R} \quad (50.3.1)$$

With this established, an immediate consequence is that if  $A\mathbf{v}_i = \lambda_i \mathbf{v}_i$  and  $A\mathbf{v}_j = \lambda_j \mathbf{v}_j$  with  $\lambda_i \neq \lambda_j$  then:

$$\langle \mathbf{v}_i, A\mathbf{v}_j \rangle = \lambda_j \langle \mathbf{v}_i, \mathbf{v}_j \rangle = \lambda_i \langle \mathbf{v}_i, \mathbf{v}_j \rangle \quad (50.3.2)$$

giving:

$$(\lambda_j - \lambda_i) \langle \mathbf{v}_i, \mathbf{v}_j \rangle = 0 \quad (50.3.3)$$

and since by assumption  $\lambda_i \neq \lambda_j$  then this can only occur  $\langle \mathbf{v}_i, \mathbf{v}_j \rangle$ . ■

Interestingly, Hermitian matrices can always be diagonalized.

**Proposition (Hermitian diagonalizability)**

Let  $V$  be an  $n$  dimensional vector space over  $\mathbb{C}$ , if  $f$  is a Hermitian map defined on  $V$  then it has an orthonormal eigenbasis  $\mathcal{E} = \{\mathbf{v}_i\}$ .

*Proof.* We consider non-degenerate maps  $f$ , and proceed by induction.

If  $\dim V = 1$ , then the result is trivial.

Suppose for all  $k < n$  we have shown that all hermitian maps have an orthonormal eigenbasis. Then, let's consider a map  $f$  on an  $n$ -dimensional  $V$ . We have that its characteristic equation must have at least one root  $\lambda$  over  $\mathbb{C}$ . Its eigenspace  $W \equiv \text{Eig}_f(\lambda)$  is such that  $\dim W > 0$ , and we consider its orthogonal complement  $W^\perp = \{\mathbf{u} : \in V : \langle \mathbf{u}, \mathbf{v}_i \rangle = 0 \forall \mathbf{v}_i \in W\}$ . We have that for  $\mathbf{u} \in W^\perp$ :

$$\langle \mathbf{w}, f(\mathbf{u}) \rangle = \langle f(\mathbf{w}), \mathbf{u} \rangle = \lambda \langle \mathbf{w}, \mathbf{u} \rangle = 0 \quad (50.3.4)$$

so  $f(\mathbf{u}) \in W^\perp$ . Consequently, we may restrict  $f$  to  $W^\perp$ , and define its restriction as  $g \equiv f|_{W^\perp}$ . Since  $\dim W^\perp = k < n$ , we can use the induction hypothesis to deduce that it has an orthonormal eigenbasis  $\{\mathbf{v}_i\}$ . Furthermore, since  $\dim W^\perp + \dim W = \dim V = n$ , we have that  $\{\mathbf{v}_i\} \cup W$  will give a set of  $n$  linearly independent eigenvectors of  $f$ , as desired. ■

Hermitian matrices play an important role, since their diagonalization is often easier to perform.

**Theorem (Diagonalizing hermitian matrices)** Let  $f$  be a hermitian linear map on  $V$ .

Then, its diagonalized form is found through the similarity transformation:

$$A' = P^T A P \quad (50.3.5)$$

where  $\mathbf{A}$  is the matrix representation of  $f$  in some basis  $\mathcal{B}$ , and

$$\mathbf{P} = ([\mathbf{v}_1]_{\mathcal{B}} \ [\mathbf{v}_2]_{\mathcal{B}} \dots [\mathbf{v}_n]_{\mathcal{B}}) \quad (50.3.6)$$

*Proof.* We have that:

$$\mathbf{P}^T \mathbf{P} = \begin{pmatrix} [\mathbf{v}_1]_{\mathcal{B}}^T \\ [\mathbf{v}_2]_{\mathcal{B}}^T \\ \vdots \\ [\mathbf{v}_n]_{\mathcal{B}}^T \end{pmatrix} ([\mathbf{v}_1]_{\mathcal{B}} \ [\mathbf{v}_2]_{\mathcal{B}} \dots [\mathbf{v}_n]_{\mathcal{B}}) \quad (50.3.7)$$

$$= \begin{pmatrix} [\mathbf{v}_1]_{\mathcal{B}}^T [\mathbf{v}_1]_{\mathcal{B}} & [\mathbf{v}_1]_{\mathcal{B}}^T [\mathbf{v}_2]_{\mathcal{B}} & \dots & [\mathbf{v}_1]_{\mathcal{B}}^T [\mathbf{v}_n]_{\mathcal{B}} \\ [\mathbf{v}_2]_{\mathcal{B}}^T [\mathbf{v}_1]_{\mathcal{B}} & [\mathbf{v}_2]_{\mathcal{B}}^T [\mathbf{v}_2]_{\mathcal{B}} & \dots & [\mathbf{v}_2]_{\mathcal{B}}^T [\mathbf{v}_n]_{\mathcal{B}} \\ \vdots & \vdots & \ddots & \vdots \\ [\mathbf{v}_n]_{\mathcal{B}}^T [\mathbf{v}_1]_{\mathcal{B}} & [\mathbf{v}_n]_{\mathcal{B}}^T [\mathbf{v}_2]_{\mathcal{B}} & \dots & [\mathbf{v}_n]_{\mathcal{B}}^T [\mathbf{v}_n]_{\mathcal{B}} \end{pmatrix} \quad (50.3.8)$$

$$= \mathbb{1} \quad (50.3.9)$$

Consequently  $\mathbf{P}$  is unitary, and hence when diagonalizing  $\mathbf{A}'$  according to the general procedure:

$$\mathbf{A}' = \mathbf{P}^{-1} \mathbf{A} \mathbf{P} = \mathbf{P}^T \mathbf{A} \mathbf{P} \quad (50.3.10)$$

as desired. ■

**Example.** Let's diagonalize the following hermitian matrix:

$$\mathbf{A} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 2 & 1 \\ 0 & 1 & 2 \end{pmatrix} \quad (50.3.11)$$

Its characteristic equation is:

$$\begin{vmatrix} 1 - \lambda & 0 & 0 \\ 0 & 2 - \lambda & 1 \\ 0 & 1 & 2 - \lambda \end{vmatrix} = (1 - \lambda)((2 - \lambda)^2 - 1) \quad (50.3.12)$$

$$= (1 - \lambda)(1 - \lambda)(3 - \lambda) = 0 \quad (50.3.13)$$

giving  $\lambda = 1, 3$ . It may seem like this matrix is not diagonalisable, since we only have two eigenvalues. However, we know that this can't be the case, Hermitian matrices are always diagonalisable. Indeed, although we only have two eigenvalues, it turns out that the first  $\lambda_1 = 1$  will have a two dimensional eigenspace, so we will be able to find two orthonormal eigenvectors associated to this eigenvalue.

For  $\lambda_1 = 1$  we get that

$$\begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 1 & 1 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \implies y = -z \quad (50.3.14)$$

This gives a two-dimensional eigenspace:

$$\text{Eig}_A(\lambda_1) = \left\{ t \begin{pmatrix} k \\ 1 \\ -1 \end{pmatrix}, \forall k, t \in \mathbb{R} \right\} \quad (50.3.15)$$

so we can choose as our orthonormal eigenvectors:

$$\mathbf{v}_1 = \frac{1}{\sqrt{2}} \begin{pmatrix} 0 \\ 1 \\ -1 \end{pmatrix}, \mathbf{v}_2 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \quad (50.3.16)$$

Finally, for  $\lambda_2 = 3$  then we get that:

$$\begin{pmatrix} -2 & 0 & 0 \\ 0 & -1 & 1 \\ 0 & 1 & -1 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \implies x = 0, y = z \quad (50.3.17)$$

giving the following eigenspace:

$$\text{Eig}_A(\lambda_2) = \left\{ t \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix}, \forall t \in \mathbb{R}^* \right\} \quad (50.3.18)$$

We choose the following eigenvector

$$\mathbf{v}_3 = \frac{1}{\sqrt{2}} \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix} \quad (50.3.19)$$

Hence, the orthonormal eigenbasis has a transition matrix:

$$P = \frac{1}{\sqrt{2}} \begin{pmatrix} 0 & \sqrt{2} & 0 \\ 1 & 0 & 1 \\ -1 & 0 & 1 \end{pmatrix} \implies P^T = \frac{1}{\sqrt{2}} \begin{pmatrix} 0 & 1 & -1 \\ \sqrt{2} & 0 & 0 \\ 0 & 1 & 1 \end{pmatrix} \quad (50.3.20)$$

The diagonalized form of  $A$  is then:

$$\frac{1}{2} \begin{pmatrix} 0 & 1 & -1 \\ \sqrt{2} & 0 & 0 \\ 0 & 1 & 1 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 2 & 1 \\ 0 & 1 & 2 \end{pmatrix} \begin{pmatrix} 0 & \sqrt{2} & 0 \\ 1 & 0 & 1 \\ -1 & 0 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 3 \end{pmatrix} \quad (50.3.21)$$

as desired. This concludes our process of orthogonal diagonalization. ◀

### Definition (*Normal linear map*)

Let  $f : V \rightarrow V$  be a linear map. Then,  $f$  is normal iff  $[f, f^\dagger] = 0$ .

Note that unitary and hermitian maps are special cases of normal linear maps.

**Proposition (Hermitian adjoint of normal map)**

Let  $f : V \rightarrow V$  be a normal linear map over  $V$ , and let  $\mathbf{v} \in \text{Eig}_f(\lambda) \subseteq V$  be an eigenvector of  $f$  with eigenvalue  $\lambda$ . Then, we have that  $f^\dagger(\mathbf{v}) = \lambda^*\mathbf{v}$ .

*Proof.* Let  $g = f - \lambda \text{id}$ . Then we find that:

$$g \circ g^\dagger = (f - \lambda \text{id}) \circ (f^\dagger - \lambda^* \text{id}) = f \circ f^\dagger - \lambda^* f - \lambda f^\dagger + |\lambda|^2 \text{id} \quad (50.3.22)$$

$$= f^\dagger \circ f - \lambda^* f - \lambda f^\dagger + |\lambda|^2 \text{id} \quad (50.3.23)$$

$$= g^\dagger \circ g \quad (50.3.24)$$

implying that  $g$  is a normal linear map. Consequently:

$$0 = \langle g(\mathbf{v}), g(\mathbf{v}) \rangle = \langle \mathbf{v}, g^\dagger \circ g(\mathbf{v}) \rangle = \langle \mathbf{v}, g \circ g^\dagger(\mathbf{v}) \rangle = \langle g^\dagger(\mathbf{v}), g^\dagger(\mathbf{v}) \rangle \quad (50.3.25)$$

implying that  $g^\dagger(\mathbf{v}) = f^\dagger(\mathbf{v}) - \lambda^*\mathbf{v} = 0$  as desired.  $\blacksquare$

**Theorem (Spectral theorem for normal maps)**

Let  $f : V \rightarrow V$  be a linear map. Then  $f$  is normal iff it has an orthonormal eigenvector basis, that is, it is diagonalizable.

*Proof.* ( $\implies$ ) Suppose  $f$  is a normal linear map on  $V$ , we proceed by induction.

If  $\dim V = n = 1$ , then the result is trivially verified.

Suppose that for  $\dim V = k < n$  we have shown that all normal linear maps have an orthonormal eigenvector basis. Then, let's consider a normal linear map  $f$  on an  $n$ -dimensional  $V$ . We have that its characteristic equation must have at least one root  $\lambda$  over  $\mathbb{C}$ . Its eigenspace  $W \equiv \text{Eig}_f(\lambda)$  is such that  $\dim W = 1$ , and we consider its orthogonal complement  $W^\perp = \{\mathbf{u} : \langle \mathbf{u} \in V : \mathbf{u}, \mathbf{v}_i \rangle = 0 \forall \mathbf{v}_i \in W\}$  with  $\dim W^\perp = n - 1$ . We have that for  $\mathbf{u} \in W^\perp$ :

$$\langle f(\mathbf{u}), \mathbf{v} \rangle = \langle \mathbf{u}, f^\dagger(\mathbf{v}) \rangle = \lambda^* \langle \mathbf{w}, \mathbf{v} \rangle = 0 \quad (50.3.26)$$

$$\langle f^\dagger(\mathbf{u}), \mathbf{v} \rangle = \langle \mathbf{u}, f(\mathbf{v}) \rangle = \lambda \langle \mathbf{w}, \mathbf{v} \rangle = 0 \quad (50.3.27)$$

These two results imply that  $f(W^\perp), f^\dagger(W^\perp) \subseteq W$ , and we may therefore consider the restriction  $f|_{W^\perp}$ . By the induction assumption, this normal map has  $n - 1$  orthonormal eigenvectors. Adding  $\frac{\mathbf{v}}{|\mathbf{v}|}$  gives the desired list of  $n$  linearly independent orthonormal eigenvectors.

( $\impliedby$ ) Suppose that  $f$  has an orthonormal eigenvector basis  $\{\mathbf{n}_i\}$ . Then:

$$f \circ f^\dagger(\mathbf{n}_i) = \lambda_i^* f(\mathbf{n}_i) = |\lambda_i|^2 \mathbf{n}_i \quad (50.3.28)$$

and

$$f^\dagger \circ f(\mathbf{n}_i) = \lambda_i f(\mathbf{n}_i) = |\lambda_i|^2 \mathbf{n}_i \quad (50.3.29)$$

implying that

$$[f, f^\dagger](\mathbf{n}_i) = 0 \quad (50.3.30)$$

Given any vector  $\mathbf{v} \in V$ , it may be expanded in the eigenbasis as  $\mathbf{v} = \sum_i \alpha_i \mathbf{n}_i$  so that:

$$[f, f^\dagger](\mathbf{v}) = \sum_i \alpha_i [f, f^\dagger](\mathbf{n}_i) = 0, \forall \mathbf{v} \in V \implies [f, f^\dagger] = 0 \quad (50.3.31)$$

as desired. ■

Notice the resemblance between this proof and the proof that all hermitian maps are diagonalizable. We proceeded by showing that there must be some eigenvector, and that its orthogonal complement is invariant under the map we are interested in. Diagonalizing the restriction of the map to the orthogonal complement gives an extra set of eigenvectors which we can use to complete the proof.

## 50.4 Classifying conics

Suppose we have a conic with general equation:

$$Ax^2 + Bxy + Cy^2 + Fx + Gy + H = 0 \quad (50.4.1)$$

Our goal will be to classify this conic as either a parabola, hyperbola or ellipse, and determine some of its fundamental features.

### Aligning the axes

We can write (50.4.1) as a product of matrices:

$$(x \ y) \begin{pmatrix} A & \frac{B}{2} \\ \frac{B}{2} & C \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} + (F \ G) \begin{pmatrix} x \\ y \end{pmatrix} + H = 0 \quad (50.4.2)$$

Let us define:

$$\mathbf{A} = \begin{pmatrix} A & \frac{B}{2} \\ \frac{B}{2} & C \end{pmatrix}, \mathbf{J} = (F \ G), \mathbf{x} = \begin{pmatrix} x \\ y \end{pmatrix} \quad (50.4.3)$$

then (50.4.2) turns into

$$\mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{J}^T \mathbf{x} + H = 0 \quad (50.4.4)$$

It is important to note that  $\mathbf{A}$  is a symmetric matrix, and can therefore be orthogonally diagonalized. Suppose that  $D = P^T \mathbf{A} P$ , then we get that:

$$\mathbf{x}^T P D P^T \mathbf{x} + \mathbf{J}^T \mathbf{x} + H = 0 \quad (50.4.5)$$

Let us define the coordinate vectors in the eigenbasis as:  $\mathbf{x}' = P^T \mathbf{x}$  (recall that  $P$  represents a change from the eigenbasis, so its inverse/transpose will represent a change to the eigenbasis). Then we find that:

$$(\mathbf{x}')^T D \mathbf{x}' + \mathbf{J}^T P \mathbf{x}' + H = 0 \quad (50.4.6)$$

The process we have gone through can be viewed geometrically as rotating  $\mathbb{R}^2$  to align the axes with the eigenvector basis of  $\mathbf{A}$ . Indeed, since  $\mathbf{A}$  is symmetric, its transition matrix will be unitary, it will represent a rotation/reflection. By performing a change of basis  $\mathbf{x} \rightarrow \mathbf{x}'$  we were really just rotating  $\mathbb{R}^2$ . Suppose  $D = \text{diag}(\lambda_1, \lambda_2)$  and  $\mathbf{J}^T P = (f' \ g')$ , then we find:

$$\lambda_1 x'^2 + \lambda_2 y'^2 + f' x' + g' y' + H = 0 \quad (50.4.7)$$

### Translating the origin

The final step is translating the origin to get a standard conic. For ellipses and hyperbolas we do so by completing the square:

$$\lambda_1 x'^2 + \lambda_2 y'^2 + fx' + gy' + H = 0 \quad (50.4.8)$$

$$\implies \lambda_1 \left( x' + \frac{f}{2\lambda_1} \right)^2 + \lambda_2 \left( y' + \frac{g}{2\lambda_2} \right)^2 + H - \frac{f^2}{4\lambda_1} - \frac{g^2}{4\lambda_2} = 0 \quad (50.4.9)$$

Letting the translated axes be defined by:

$$\mathbf{x}'' = \mathbf{x}' + \frac{1}{2} \begin{pmatrix} \frac{f}{\lambda_1} \\ \frac{g}{\lambda_2} \end{pmatrix} \quad (50.4.10)$$

then we find that:

$$\lambda_1 x''^2 + \lambda_2 y''^2 + h = 0, \quad h = H - \frac{f^2}{4\lambda_1} - \frac{g^2}{4\lambda_2} \quad (50.4.11)$$

which is a conic. It can be rearranged into the more useful form :

$$\frac{x''^2}{a^2} + \frac{y''^2}{b^2} = 1 \quad (50.4.12)$$

where  $a^2 = -\frac{h}{\lambda_1}$  and  $b^2 = -\frac{h}{\lambda_2}$ . Depending on the values of  $a, b$  this will be either a hyperbola or ellipse.

Conic	Standard form
Hyperbola	$\frac{x^2}{a^2} - \frac{y^2}{b^2} = 1$
Parabola	$y^2 - 4ax = 0$
Ellipse	$\frac{x^2}{a^2} - \frac{y^2}{b^2} = 1$

If instead we are dealing with a parabola, then we will find that one of the eigenvalues is 0. Suppose WLOG that  $\lambda_1 = 0$ , then we find:

$$\lambda_2 y'^2 + fx' + gy' + H = 0 \quad (50.4.13)$$

$$\implies \lambda_2 \left( y' + \frac{g}{2\lambda_2} \right)^2 + fx' + H - \frac{g^2}{4\lambda_2} = 0 \quad (50.4.14)$$

Letting the translated axes be defined by:

$$\mathbf{x}'' = \mathbf{x}' + \left( \frac{\frac{H}{f} - \frac{g^2}{4f\lambda_2}}{\frac{g}{2\lambda_2}} \right) \quad (50.4.15)$$

then we find that:

$$y''^2 + \frac{f}{\lambda_2} x'' = 0 \quad (50.4.16)$$

which is a parabola. Had  $\lambda_2 = 0$  then we would have found in complete analogy to before:

$$x''^2 + \frac{g}{\lambda_1} y'' = 0 \quad (50.4.17)$$

**Example.** Let's classify the conic described by:

$$x^2 - 4xy + 4y^2 - 6x - 8y + 5 = 0 \quad (50.4.18)$$

which may be re-expressed in matrix form as:

$$\mathbf{x}^T \begin{pmatrix} 1 & -2 \\ -2 & 4 \end{pmatrix} \mathbf{x} + (-6 \ -8) \mathbf{x} + 5 = 0 \quad (50.4.19)$$

The eigenvalues of A are easily found to obey

$$(1 - \lambda)(4 - \lambda) - 4 = 0 \implies \lambda(\lambda - 5) = 0 \implies \lambda_1 = 0, \lambda_2 = 5 \quad (50.4.20)$$

For  $\lambda_1 = 0$  we get the eigenspace

$$\text{Eig}_A(\lambda_1) = \left\{ k \begin{pmatrix} 2 \\ 1 \end{pmatrix} : k \in \mathbb{R}^* \right\} \quad (50.4.21)$$

Similarly, for  $\lambda_2 = 5$  we get the eigenspace:

$$\text{Eig}_A(\lambda_2) = \left\{ k \begin{pmatrix} 1 \\ -2 \end{pmatrix} : k \in \mathbb{R}^* \right\} \quad (50.4.22)$$

so we may choose the orthonormal basis  $\left\{ \frac{1}{\sqrt{5}} \begin{pmatrix} 2 \\ 1 \end{pmatrix}, \frac{1}{\sqrt{5}} \begin{pmatrix} 1 \\ -2 \end{pmatrix} \right\}$  with transition matrix:

$$P = \frac{1}{\sqrt{5}} \begin{pmatrix} 2 & 1 \\ 1 & -2 \end{pmatrix} \implies J^T P = \frac{1}{\sqrt{5}} (-20 \ 10) \quad (50.4.23)$$

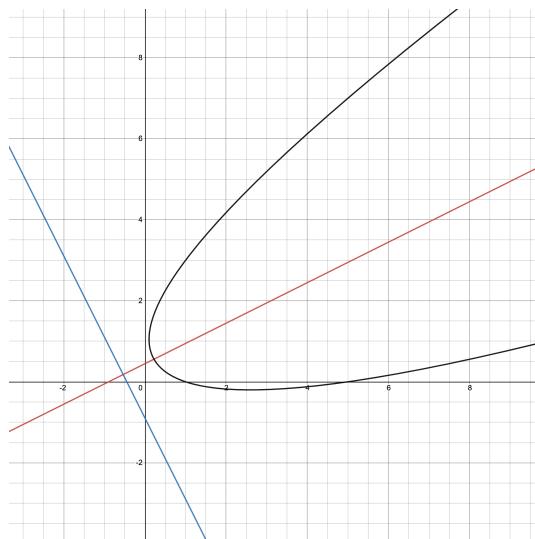
Consequently we find:

$$5y'^2 - 4\sqrt{5}x' + 2\sqrt{5}y' + 5 = 0 \implies y'^2 - \frac{4\sqrt{5}}{5}x' + \frac{2\sqrt{5}}{5}y' + 1 = 0 \quad (50.4.24)$$

We now complete the square:

$$(y' + \frac{\sqrt{5}}{5})^2 - \frac{4\sqrt{5}}{5}x' + \frac{4}{5} = 0 \implies y''^2 = \frac{4\sqrt{5}}{5}x'' \quad (50.4.25)$$

where  $y'' = y' + \frac{\sqrt{5}}{5}$  and  $x'' = x' - \frac{\sqrt{5}}{5}$ . This is therefore a parabola.



## 50.5 Matrix exponentials and Lie algebras

One final important application of diagonalization is in determining matrix exponents.

### Proposition (Matrix exponents)

Suppose that  $A$  is a diagonalizable matrix with  $A' = P^{-1}AP$ . Then:

$$A^n = PA'^n P^{-1} \quad (50.5.1)$$

*Proof.* This follows immediately from:

$$A^n = \underbrace{(PAP^{-1})(PAP^{-1})\dots(PAP^{-1})}_{n \text{ times}} = PA'^n P^{-1} \quad (50.5.2)$$

### Definition (Matrix exponential)

Suppose that  $A$  matrix, then its exponential is defined as:

$$e^A = \sum_{i=0}^{\infty} \frac{1}{n!} A^n = \mathbb{1} + A + \frac{1}{2} A^2 + \dots \quad (50.5.3)$$

## 50.6 Schur's triangulation theorem

### Theorem (Schur's triangulation theorem)

Let  $A \in \text{Mat}_n(\mathbb{C})$  with eigenvalues  $\lambda_1, \lambda_2, \dots$  which may be degenerate.. Then  $A$  is unitarily equivalent to an upper triangular matrix:

$$A = UTU^\dagger \quad (50.6.1)$$

where:

$$T = \begin{pmatrix} \lambda_1 & & \cdots \\ 0 & \lambda_2 & & \\ \vdots & & & \ddots \end{pmatrix} \quad (50.6.2)$$

*Proof.* We proceed by induction. For  $n = 1$  the result is trivial. Suppose we have shown that any  $m \times m$  matrix where  $m \leq n - 1$  is unitarily equivalent to an upper triangular matrix. Let  $A \in \text{Mat}_n(\mathbb{C})$  have eigenvalues  $\lambda_1, \lambda_2, \dots$  and eigenvectors  $\mathbf{v}_1, \mathbf{v}_2, \dots$ , which may be degenerate, and are assumed to have unit norm.

Now  $\mathbf{v}_1$  can be used to form an orthonormal basis  $\{\mathbf{v}_1, \mathbf{u}_2, \dots, \mathbf{u}_n\}$ . Then, the resulting matrix when we change to this basis will be unitarily equivalent to  $A$ :

$$A = V \left( \begin{array}{c|cccc} \lambda_1 & a_{12} & \dots & a_{1n} \\ \hline 0 & & & \\ \vdots & & \tilde{A} & \\ 0 & & & \end{array} \right) V^\dagger \quad (50.6.3)$$

Clearly, we must have that  $\chi_A(\lambda) = (\lambda_1 - \lambda)\chi_{\tilde{A}}(\lambda)$ , implying that  $\tilde{A}$  has eigenvalues  $\lambda_2, \lambda_3, \dots$  identical to  $A$ , and which could be degenerate. We can now use the induction hypothesis, since  $\tilde{A} \in \text{Mat}_{n-1}(\mathbb{C})$  we have that it is unitarily equivalent to some upper triangular matrix:

$$\tilde{A} = \tilde{W} \left( \begin{array}{ccccc} \lambda_2 & \tilde{a}_{12} & \dots & \tilde{a}_{1n} \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \tilde{a}_{n-1n} \\ 0 & \dots & 0 & \lambda_n \end{array} \right) \tilde{W}^\dagger \quad (50.6.4)$$

Note also that:

$$\left( \begin{array}{c|ccccc} 1 & 0 & \dots & 0 \\ 0 & & & & & \\ \vdots & & \tilde{W} & & & \\ 0 & & & & & \end{array} \right)^\dagger \left( \begin{array}{c|ccccc} \lambda_1 & a_{12} & \dots & a_{1n} \\ 0 & & & & & \\ \vdots & & \tilde{A} & & & \\ 0 & & & & & \end{array} \right) \left( \begin{array}{c|ccccc} 1 & 0 & \dots & 0 \\ 0 & & & & & \\ \vdots & & \tilde{W} & & & \\ 0 & & & & & \end{array} \right) = \left( \begin{array}{c|ccccc} 1 & 0 & \dots & 0 \\ 0 & & & & & \\ \vdots & & \tilde{W}^\dagger & & & \\ 0 & & & & & \end{array} \right) \left( \begin{array}{c|ccccc} \lambda_1 & b_{12} & \dots & b_{1n} \\ 0 & & & & & \\ \vdots & & \tilde{A}\tilde{W} & & & \\ 0 & & & & & \end{array} \right) \quad (50.6.5)$$

$$= \left( \begin{array}{c|ccccc} \lambda_1 & c_{12} & \dots & c_{1n} \\ 0 & & & & & \\ \vdots & & \tilde{W}^\dagger\tilde{A}\tilde{W} & & & \\ 0 & & & & & \end{array} \right) = \left( \begin{array}{ccccc} \lambda_1 & \tilde{b}_{12} & \dots & \dots & \tilde{b}_{1n} \\ 0 & \lambda_2 & \ddots & \ddots & \vdots \\ 0 & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \tilde{a}_{n-1n} \\ 0 & \dots & \dots & 0 & \lambda_n \end{array} \right) \quad (50.6.6)$$

Defining:

$$W = \left( \begin{array}{c|ccccc} 1 & 0 & \dots & 0 \\ 0 & & & & & \\ \vdots & & \tilde{W} & & & \\ 0 & & & & & \end{array} \right) \quad (50.6.7)$$

we see that  $W$  is unitary and that

$$A = VW \left( \begin{array}{c|ccccc} \lambda_1 & a_{12} & \dots & a_{1n} \\ 0 & & & & & \\ \vdots & & \tilde{A} & & & \\ 0 & & & & & \end{array} \right) W^\dagger V^\dagger \quad (50.6.8)$$

Thus  $A$  is unitarily triangularizable via  $U = VW$ . ■

### Theorem (The Cayley-Hamilton theorem)

Let  $A \in \text{Mat}_n(\mathbb{C})$  have characteristic polynomial  $\chi_A(\lambda)$ . Then, we have that  $\chi_A(A) = 0$  where 0 is the zero element of  $\text{Mat}_n(\mathbb{C})$ ,

*Proof.* We can factorize the characteristic polynomial into the following form due to the Fundamental theorem of algebra:

$$\chi_A(\lambda) = (\lambda_1 - \lambda)(\lambda_2 - \lambda)\dots(\lambda_n - \lambda) \quad (50.6.9)$$

implying that

$$\chi_A(A) = (\lambda_1 \mathbb{1} - A)(\lambda_2 \mathbb{1} - A)\dots(\lambda_n \mathbb{1} - A) \quad (50.6.10)$$

Since  $A \in \text{Mat}_n(\mathbb{C})$ , Schur's theorem tells us that it can be triangularized unitarily  $A = UTU^\dagger$  where  $T$  is upper triangular. Therefore:

$$\chi_A(A) = U(\lambda_1 \mathbb{1} - T)U^\dagger U(\lambda_2 \mathbb{1} - T)U^\dagger \dots U(\lambda_n \mathbb{1} - T)U^\dagger \quad (50.6.11)$$

$$= U(\lambda_1 \mathbb{1} - T)(\lambda_2 \mathbb{1} - T)\dots(\lambda_n \mathbb{1} - T)U^\dagger \quad (50.6.12)$$

Each of the factors  $\lambda_i \mathbb{1} - T$  will be upper triangular with the  $i$ th diagonal element equal to zero. It is easy to verify that a product of such matrices must be null. Let  $A = \lambda_1 \mathbb{1} - T$  and  $B = \lambda_2 \mathbb{1} - T$ ,

and define  $C = AB$ . We have that  $C_{11} = C_{22} = 0$  since in general for triangular matrices:

$$C_{ii} = \sum_j A_{ij}B_{ji} = A_{ii}B_{ii} \quad (50.6.13)$$

Instead,  $C_{12} = \sum_j A_{1j}B_{j2} = 0$  since  $A_{11}$  and  $B_{22}$  are both zero. The first two rows of  $C$  are thus equal to zero.

Suppose we have repeated this process up to the factor  $(\lambda_m \mathbb{1} - T)$  so that the first  $m - 1$  rows are all zero. Then letting  $C = \prod_{i=1}^m (\lambda_i \mathbb{1} - T) = A(\lambda_m \mathbb{1} - T)$  we get:

$$C_{lm} = \sum_{l \leq j \leq m} A_{lj}B_{jm} \quad (50.6.14)$$

Now since  $A_{lj} = 0$  for all  $j < m$ , the only term that will survive will be that with  $j = m$ . Consequently:

$$C_{lm} = A_{lm}B_{mm} = 0 \quad (50.6.15)$$

so the  $m$ th column will also be zero. It follows by induction that  $(\lambda_1 \mathbb{1} - A) \dots (\lambda_n \mathbb{1} - A) = 0$  and thus  $\chi_A(A) = 0$  as desired. ■

We can use Schur's triangulation theorem to prove the Spectral theorem more generally.

**Proposition (Normal triangular matrices)** A triangular matrix is normal iff it is diagonal.

*Proof.* The  $\iff$  is trivial. Suppose  $A \in \text{Mat}_n(\mathbb{C})$  is a normal triangular matrix. The case  $n = 1$  is obvious. Suppose the proposition is true for all  $m \times m$  matrices where  $m \leq n - 1$ . Then, writing  $A$  as:

$$A = \begin{pmatrix} a_{11} & \mathbf{a} \\ 0 & \tilde{A} \end{pmatrix} \implies A^\dagger = \begin{pmatrix} a_{11}^* & 0 \\ \mathbf{a}^* & \tilde{A}^\dagger \end{pmatrix} \quad (50.6.16)$$

we see that

$$A^\dagger A = \begin{pmatrix} |a_{11}|^2 & \dots \\ \dots & \tilde{A}^\dagger \tilde{A} + \|\mathbf{a}\|^2 \end{pmatrix} \quad (50.6.17)$$

$$AA^\dagger = \begin{pmatrix} |a_{11}|^2 + \|\mathbf{a}\|^2 & \dots \\ \dots & \tilde{A} \tilde{A}^\dagger \end{pmatrix} \quad (50.6.18)$$

For these to be equal, we need  $a_{11} = 0$ ,  $\mathbf{a} = 0$  and  $\tilde{A}^\dagger \tilde{A} = \tilde{A} \tilde{A}^\dagger$ . Also, since  $A \in \text{Mat}_{n-1}(\mathbb{C})$  is upper triangular and normal, it must be diagonal. Consequently  $A$  is diagonal, as required. ■

**Theorem (The Spectral Theorem)** Let  $A \in \text{Mat}_n(\mathbb{C})$  with eigenvalues  $\lambda_1, \dots, \lambda_n$ , the following are equivalent:

- (i)  $A$  is normal
- (ii)  $A$  is unitarily diagonalizable
- (iii)  $\sum_{ij} |A_{ij}|^2 = \sum_i |\lambda_i|^2$

*Proof.*

(i)  $\implies$  (ii) If  $A$  is unitarily diagonalizable then  $A = UDU^\dagger$  and thus

$$A^\dagger A = (UD^\dagger U^\dagger)(UDU^\dagger) = U D^\dagger D U^\dagger \quad (50.6.19)$$

$$AA^\dagger = (UDU^\dagger)(UD^\dagger U^\dagger) = U D D^\dagger U^\dagger \quad (50.6.20)$$

but  $D^\dagger D = DD^\dagger$  so  $A$  is normal.

(ii)  $\implies$  (i) Suppose  $A$  is normal. By Schur's theorem it is unitarily equivalent to an upper triangular matrix.

**Lemma.** Normality is preserved under unitary transformations.

Indeed if  $A^\dagger A = AA^\dagger$  then

$$(UAU^\dagger)^\dagger(UAU^\dagger) = UA^\dagger AU^\dagger = UAA^\dagger U^\dagger = (UAU^\dagger)(UAU^\dagger)^\dagger \quad (50.6.21)$$

Consequently, the upper triangular decomposition of  $A$  must be normal, and thus diagonal, as desired.

(ii)  $\implies$  (iii) Since  $A$  is unitarily diagonalizable, we have that  $A = UDU^\dagger$  where  $D = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$  and thus:

$$\sum_{ij} |a_{ij}|^2 = \text{tr}(A^\dagger A) = \text{tr}(UD^\dagger DU^\dagger) = \text{tr}(U^\dagger U) \text{tr}(D) = \sum_i |\lambda_i|^2 \quad (50.6.22)$$

(iii)  $\implies$  (ii) By Schur's theorem we have that  $A$  is triangularizable into  $T$  with  $\sum_i |T_{ii}|^2 = \sum_i |\lambda_i|^2$ . Note also that (iii) implies:

$$\sum_{ij} |a_{ij}|^2 = \text{tr}(A^\dagger A) = \text{tr}(T^\dagger T) = \sum_{ij} |T_{ij}|^2 = \sum_i |\lambda_i|^2 \quad (50.6.23)$$

We therefore find that  $\sum_i |T_{ii}|^2 = \sum_{ij} |T_{ij}|^2$  which is only possible if  $T$  is diagonal.

■

## 50.7 Jordan canonical form (make sure to write by October)

# Multilinear algebra

## 51.1 Review of linear algebra

Recall the normal definition of vector spaces.

### Definition (*Vector space*)

A vector space  $(V, +, \cdot)$  over a field  $\mathbb{F}$  is a set  $V$  and two maps:

$$+ : V \times V \rightarrow v \tag{51.1.1}$$

$$\cdot : \mathbb{F} \times V \rightarrow V \tag{51.1.2}$$

known as addition and scalar multiplication such that for all  $u, v, w \in V, \alpha, \beta \in \mathbb{F}$ :

- (i)  $v + w = w + v$
- (ii)  $(u + v) + w = u + (v + w)$
- (iii)  $\exists 0 \in V$  such that  $v + 0 = v$
- (iv)  $\exists (-v) \in V$  such that  $v + (-v) = 0$
- (v)  $\alpha \cdot (\beta \cdot v) = (\alpha\beta) \cdot v$
- (vi)  $\alpha \cdot (v + w) = \alpha \cdot v + \alpha \cdot w$
- (vii)  $(\alpha + \beta) \cdot v = \alpha \cdot v + \beta \cdot v$
- (viii) given the identity element 1 of  $\mathbb{F}$ , then  $1 \cdot v = v$

The element of a vector space is informally referred to as a **vector**. It is easy to see that the set  $\mathcal{P}_n$  of all polynomial functions  $p$  up to order  $n \in \mathbb{N}$  is indeed a vector space:

$$\mathcal{P}_n = \{p(x) = \sum_{m=0}^n p_m x^m : p_m \in \mathbb{R}\} \tag{51.1.3}$$

### Definition (*Linear map*)

Let  $(V, +_V, \cdot_V)$  and  $(W, +_W, \cdot_W)$  be two vector spaces. Then the map  $\phi : V \rightarrow W$  is a **linear map** if it is structure preserving:

- (i)  $\phi(u +_V v) = \phi(u) +_W \phi(v)$
- (ii)  $\phi(\alpha \cdot_V v) = \alpha \cdot_W \phi(v)$

If such a map between  $V$  and  $W$  exists then we write that  $V \cong W$ .

**Definition (Set of linear maps)**

We define the set of all linear maps between two vector spaces  $(V, +_V, \bullet_V)$  and  $(W, +_W, \bullet_W)$  (the latter defined over  $\mathbb{F}_W$ ) as  $\text{Hom}(V, W)$ . Together with the operations  $\oplus$  and  $\odot$  defined below this set becomes a vector space:

$$\oplus : \text{Hom}(V, W) \times \text{Hom}(V, W) \rightarrow \text{Hom}(V, W) \quad (51.1.4)$$

$$(\phi, \varphi) \mapsto \phi \oplus \varphi \quad (51.1.5)$$

and

$$\odot : \mathbb{F}_W \times \text{Hom}(V, W) \rightarrow \text{Hom}(V, W) \quad (51.1.6)$$

$$(\alpha, \varphi) \mapsto \alpha \odot \varphi \quad (51.1.7)$$

where

$$(\phi \oplus \varphi)(v) = \phi(v) +_W \varphi(v), \forall v \in V \quad (51.1.8)$$

$$(\alpha \odot \varphi)(v) = \alpha \bullet_W \phi(v), \forall v \in V, \forall \alpha \in \mathbb{F}_W \quad (51.1.9)$$

**Definition (Dual vector space)**

Given a vector space  $V$  defined over a field  $\mathbb{F}$ , then the set  $V^*$  is defined as:

$$V^* \equiv \{\phi : V \rightarrow \mathbb{F} : \phi \text{ is linear}\} = \text{Hom}(V, \mathbb{F}) \quad (51.1.10)$$

Equipped with  $\oplus$  and  $\otimes$  defined previously then  $V^*$  is defined as the **dual vector space**.

Informally, the element  $\phi \in V^*$  is referred to as a **covector**.

**Definition (Dual basis)**

Given a basis  $\mathcal{B} = \{e_1, e_2, \dots, e_n\} \subset V$  for  $V$  then the **dual basis** of the dual space  $V^*$  is defined as  $\mathcal{B}^* = \{\epsilon^1, \epsilon^2, \dots, \epsilon^n\} \subset V^*$  such that:

$$\epsilon^i(\epsilon_j) = \delta_i^j, \forall i, j \quad (51.1.11)$$

For the case of  $\mathcal{P}_n$ , we can choose the basis  $\mathcal{B} = \{e_i(x) = x^i : i = 0, 1, \dots, n\}$ . The dual basis  $\mathcal{B}^* = \{\epsilon_i : i = 0, 1, \dots, n\}$  can be easily identified as:

$$\epsilon^i = \frac{1}{i!} \left. \frac{d^i}{dx^i} \right|_{x=0} \quad (51.1.12)$$

Indeed one finds that

$$\epsilon^i(e_j) = \left. \frac{1}{i!} \frac{d^i}{dx^i} (x^j) \right|_{x=0} = \begin{cases} 0, & i > j \\ 1, & i = j \\ 0, & i < j \end{cases} = \delta_{ij} \quad (51.1.13)$$

## 51.2 Multilinear maps

### Definition (Tensor)

Let  $(V, +, \cdot)$  be a vector space defined over a field  $\mathbb{F}$ . Then an  $(r, s)$ -tensor  $T$  over  $V$  is a map:

$$T : (V^*)^r \times V^s \mapsto \mathbb{F} \quad (51.2.1)$$

that is linear in all its entries, thus a **multi-linear map**.

Consider for example a  $(1, 1)$ -tensor  $T$ . Then we must have that:

$$T(\phi + \varphi, v) = T(\phi, v) + T(\varphi, v) \quad (51.2.2)$$

$$T(\alpha\phi, v) = \alpha T(\phi, v) \quad (51.2.3)$$

$$T(\phi, v + w) = T(\phi, v) + T(\phi, w) \quad (51.2.4)$$

$$T(\phi, \alpha v) = \alpha T(\phi, v) \quad (51.2.5)$$

$$(51.2.6)$$

We can connect this to the typical view of a  $(1, 1)$ -tensor as a map that takes a vector and maps it to another vector. Indeed let us construct:

$$\phi_T : V \rightarrow V \quad (51.2.7)$$

$$v \mapsto T(\cdot, v) \quad (51.2.8)$$

which takes in  $v$  and maps it to another map  $T(\cdot, v)$  which takes covectors to  $\mathbb{F}$ . But this map is of the type  $V^* \rightarrow \mathbb{F}$  and thus  $T(\cdot, v) \in (V^*)^* = V$  provided  $\dim V < \infty$ . It follows that  $T(\cdot, v)$  is a vector, so  $\phi_T$  does indeed map a vector to vectors of  $V$ .

### Definition (Vectors vs Covectors)

Let  $\phi \in V^*$  be a covector. Then by definition it must be a  $(0, 1)$ -tensor since  $\phi : V \rightarrow \mathbb{F}$  is linear.

Similarly, let  $v \in V = (V^*)^*$  be a vector. Then by definition it must be a  $(1, 0)$ -tensor since  $v : V^* \rightarrow \mathbb{F}$ .

### Definition (Tensor components)

Let  $T$  be an  $(r, s)$ -tensor over a  $n$ -dimensional vector space  $V$  with basis  $\mathcal{B} = \{e_i : i = 0, 1, \dots, n\}$  and dual space  $V^*$  with dual basis  $\mathcal{B}^* = \{\epsilon_i : i = 0, 1, \dots, n\}$ . Then we define the components of  $T$  as:

$$T_{j_1 \dots j_s}^{i_1 \dots i_r} \equiv T(\epsilon^{i_1}, \dots, \epsilon^{i_r}, e_{j_1}, \dots, e_{j_s}) \quad (51.2.9)$$

so that:

$$T(\phi, v) = \sum_{ij} \phi_i v^j T_j^i \quad (51.2.10)$$

where  $\phi = \sum_i \phi_i \epsilon^i \in V^*$  and  $v = \sum_j v^j e_j \in V$ .

For example, given a  $(1, 1)$ -tensor  $T$  then  $T_j^i = T(\epsilon^i, e_j)$ . Also, due to the linearity of the tensor map

the tensor components are very useful because they allow us to expand:

$$T(\phi, v) = T\left(\sum_i \phi_i \epsilon^i, \sum_j v^j e_j\right) \quad (51.2.11)$$

$$= \sum_{ij} \phi_i v^j T(\epsilon^i, e_j) \quad (51.2.12)$$

$$= \sum_{ij} \phi_i v^j T_j^i \quad (51.2.13)$$

Note the careful placements of subscripts and superscripts, where basis vectors are given subscripts and basis covectors are given superscripts. This ensures that multiply labelled indices always appear in an up-down arrangement. Moreover, if we assume that an index appearing both up and down is summed over, then we retrieve the **Einstein summation convention**. For example, (54.4.23) becomes

$$T(\phi, v) = \phi_i v^j T_j^i \quad (51.2.14)$$

The study of tensors in specific bases will be the subject of the next two chapters, and will culminate in the mathematical framework of differential geometry.

# Tensor algebra

We have studied several types of mathematical objects until now (especially in linear algebra), and have seen that many tend to transform under changes of basis in different ways. Scalars are invariant under basis changes, whereas vectors transform under matrix multiplication. Linear operators instead transform using similarity transformations.

Tensor algebra studies the way we may categorize the ways objects transform under changes of basis, and the properties that follow from such classifications. In the next chapter on Tensor calculus, we study how we can differentiate these objects.

## 52.1 Einstein summation convention

Consider a vector space  $V$  with a basis  $\mathcal{B} = \{v_1, v_2, \dots, v_n\}$ , which we modify to  $\mathcal{B}' = \{v'_1, v'_2, \dots, v'_n\}$ . From our study of linear algebra we know that we can express this change of basis as:

$$\begin{pmatrix} v'_1 \\ \vdots \\ v'_n \end{pmatrix} = \begin{pmatrix} A_{11} & \dots & A_{1n} \\ \dots & \dots & \dots \\ A_{n1} & \dots & A_{nn} \end{pmatrix} \begin{pmatrix} v_1 \\ \vdots \\ v_n \end{pmatrix} \quad (52.1.1)$$

Alternatively:

$$v'_i = \sum_{j=1}^n A_{ij} v_j \quad (52.1.2)$$

for  $1 \leq i \leq n$ . Here, we call the index  $i$  as the **running index**, whereas the index  $j$  is called the **dummy index**.

Often times calculations in fields such as General relativity or QFT require the use of several dummy indices, and could lead to a clutter of summation symbols  $\Sigma$ . Indeed, Einstein himself faced this problem when trying to work out the differential geometry of his theory of space-time, and to fix this issue, he devised a notation known as **Einstein notation**.

In this notation, if an index appears more than once in a summation, then the corresponding  $\Sigma$  symbol may be omitted, since summation is implied. In other words, we would find that:

$$v'_i = \sum_{j=1}^n A_{ij} v_j \longrightarrow A_{ij} v_j \quad (52.1.3)$$

Suppose we wish to calculate the product of three matrices:

$$(ABC)_{il} = \sum_{j=1}^n \sum_{k=1}^n A_{ij} B_{jk} C_{kl} \quad (52.1.4)$$

In einstein notation, this becomes:

$$(ABC)_{il} = A_{ij} B_{jk} C_{kl} \quad (52.1.5)$$

which is considerably shorter.

In general, the greek alphabet is reserved for indices of space time components only. These indices therefore only take values of 0, 1, 2, 3..., where 0 is the temporal component. Instead, the normal alphabet is reserved for indices of spatial components only, so these indices run from 1, 2, 3.... according to the dimension of the space we are working in.

## 52.2 Cartesian tensors

A cartesian coordinate system in an n-dimensional space associates a coordinate  $(x_1, x_2, \dots, x_n)$  to each point in this space, with reference to a set of  $n$  basis vectors.

Suppose we have some vector  $\mathbf{r}$  in this space, whose components are  $(x_1, x_2, x_3)$  in the cartesian system  $\mathcal{C}$  and  $(x'_1, x'_2, x'_3)$  in the primed cartesian system  $\mathcal{C}'$ . The former has a basis  $\{\mathbf{e}_i\}$  and the latter has a basis  $\{\mathbf{e}'_i\}$ .

### Definition (*Rigid rotation*)

A **rigid rotation** of the cartesian axes is the transformation of components of one cartesian system to another:

$$x'_j = R_{jk} x_k \quad (52.2.1a)$$

$$x_k = R_{jk} x'_j \quad (52.2.1b)$$

where  $R_{jk} = \mathbf{e}'_j \cdot \mathbf{e}_k$ .

Indeed, a general vector  $\mathbf{r}$  may be expressed as:

$$\mathbf{r} = x_k \mathbf{e}_k = x'_j \mathbf{e}'_j \quad (52.2.2)$$

and since  $\mathbf{e}_k = (\mathbf{e}'_j \cdot \mathbf{e}_k) \mathbf{e}'_j = R_{jk} \mathbf{e}'_j$  we find that:

$$x_k R_{jk} \mathbf{e}'_j = x'_j \mathbf{e}'_j \implies x_k R_{jk} = x'_j \quad (52.2.3)$$

since  $\{\mathbf{e}'_j\}$  is linearly independent. Similarly, we may write  $\mathbf{e}'_j = (\mathbf{e}'_j \cdot \mathbf{e}_k) \mathbf{e}_k = R_{jk} \mathbf{e}_k$  so that:

$$x'_j \mathbf{e}'_j = x'_j R_{jk} \mathbf{e}_k = x_k \mathbf{e}_k \implies x_k R_{jk} = x'_j \quad (52.2.4)$$

as desired.

**Example.** Suppose we want to rotate the cartesian system with basis  $\{\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3\}$  by some angle  $\phi$  about the  $\mathbf{e}_3$  (anti-clockwise when viewed from the tip of  $\mathbf{e}_3$ ), to get another cartesian

system with basis  $\{\mathbf{e}'_1, \mathbf{e}'_2, \mathbf{e}'_3\}$  then:

$$x'_1 = (\mathbf{e}_1 \cdot \mathbf{e}'_1)x_1 + (\mathbf{e}_2 \cdot \mathbf{e}'_1)x_2 + (\mathbf{e}_3 \cdot \mathbf{e}'_1)x_3 = \cos \phi x_1 + \sin \theta x_2 \quad (52.2.5)$$

$$x'_2 = (\mathbf{e}_1 \cdot \mathbf{e}'_2)x_1 + (\mathbf{e}_2 \cdot \mathbf{e}'_2)x_2 + (\mathbf{e}_3 \cdot \mathbf{e}'_2)x_3 = -\sin \phi x_1 + \cos \theta x_2 \quad (52.2.6)$$

$$x'_3 = (\mathbf{e}_1 \cdot \mathbf{e}'_3)x_1 + (\mathbf{e}_2 \cdot \mathbf{e}'_3)x_2 + (\mathbf{e}_3 \cdot \mathbf{e}'_3)x_3 = x_3 \quad (52.2.7)$$

Therefore, we define the **rotation vector**:

$$\mathbf{R} = \begin{pmatrix} \cos \phi & \sin \phi & 0 \\ -\sin \phi & \cos \phi & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad (52.2.8)$$

so that a vector  $\mathbf{x}$ :

$$\mathbf{x}' = \mathbf{R}\mathbf{x} \quad (52.2.9)$$



### Theorem (Rotation operator $\mathbf{R}$ )

The rotation operator in three dimensions, defined as:

$$\mathbf{R} = \begin{pmatrix} \cos \phi & \sin \phi & 0 \\ -\sin \phi & \cos \phi & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad (52.2.10)$$

is unitary/orthogonal, satisfying  $\mathbf{R}^T \mathbf{R} = \mathbf{R} \mathbf{R}^T = \mathbb{I}$  and  $\det \mathbf{R} = \pm 1$ .

*Proof.* We can write that:

$$\mathbf{x}' = \mathbf{R}\mathbf{x} \iff \mathbf{x} = \mathbf{R}^{-1}\mathbf{x}' \quad (52.2.11)$$

or in component form:

$$x'_i = R_{ij}x_j \iff x_j = (\mathbf{R}^{-1})_{ji}x_i \quad (52.2.12)$$

implying that  $(\mathbf{R}^{-1})_{ji} = R_{ij}$ , or in other words:

$$\boxed{\mathbf{R}^T = \mathbf{R}^{-1}} \quad (52.2.13)$$

In other words, the rotation matrix  $\mathbf{R}$  is unitary, and in hindsight it is quite obvious why it should be unitary. A rigid cartesian rotation, as the name suggests, is rigid, and hence will preserve angles between two vectors (we will prove this soon).

Immediately, we also find that:

$$1 = \det \mathbb{I}_n = \det(\mathbf{R}\mathbf{R}^T) = \det \mathbf{R} \det \mathbf{R}^T = (\det \mathbf{R})^2 \quad (52.2.14)$$

so that:

$$\boxed{\det \mathbf{R} = \pm 1} \quad (52.2.15)$$

Another important property of rigid rotation matrices is that their components are unitary:

$$\boxed{R_{ik}R_{jk} = \delta_{ij}, R_{ik}R_{il} = \delta_{kl}} \quad (52.2.16)$$

Indeed:

$$\mathbf{e}'_i \cdot \mathbf{e}'_j = (\mathbf{e}'_i \cdot \mathbf{e}_k)(\mathbf{e}'_j \cdot \mathbf{e}_k) \mathbf{e}_k \cdot \mathbf{e}_k = R_{ik}R_{jk} = \delta_{ij} \quad (52.2.17)$$

Similarly:

$$\mathbf{e}_k \cdot \mathbf{e}_l = (\mathbf{e}'_i \cdot \mathbf{e}_k)(\mathbf{e}'_i \cdot \mathbf{e}_l) \mathbf{e}'_i \cdot \mathbf{e}'_i = R_{ik}R_{il} = \delta_{kl} \quad (52.2.18)$$

as desired.

We could also see this by using the unitarity of  $\mathbf{R}$  and write:

$$\mathbb{I} = \mathbf{R}\mathbf{R}^{-1} = \mathbf{R}\mathbf{R}^T \implies \delta_{ij} = R_{ik}R_{jk} \quad (52.2.19)$$

■

Geometrically, we see that this result makes sense. Indeed  $R_{ik}$  and  $R_{jk}$  are the angles between  $\mathbf{e}'_i, \mathbf{e}_k$  and  $\mathbf{e}'_j, \mathbf{e}_k$ . Since  $\mathbf{e}'_i$  and  $\mathbf{e}'_j$  are orthogonal, we cannot have that both terms in the product  $R_{ik}R_{jk}$  be non-zero unless  $i = j$ .

### Definition (First order Cartesian tensor)

A **first order Cartesian tensor** (or vector) is defined as a geometric object  $\mathbf{v}$  represented by the components  $v_i$  in the cartesian system  $\mathcal{C}$  and represented by the components  $v'_i$  in the cartesian system  $\mathcal{C}'$ , such that they transform under rigid cartesian rotations as :

$$v'_i = R_{ij}v_j \quad (52.2.20a)$$

$$v_i = R_{ki}v'_k \quad (52.2.20b)$$

where  $\mathbf{R}$  is the rotation matrix as defined (52.2.8).

**Example.** Consider the quantity  $\mathbf{v} = (x_1^2, x_2^2)$ , which transforms under rotations as:

$$v'_1 = (x'_1)^2 = (x_1 \cos \theta + x_2 \sin \theta)^2 \quad (52.2.21)$$

$$v'_2 = (x'_2)^2 = (-x_1 \sin \theta + x_2 \cos \theta)^2 \quad (52.2.22)$$

If this were a first order cartesian tensor, we would need:

$$v'_1 = v_1 \cos \theta + v_2 \sin \theta = x_1^2 \cos \theta + x_2^2 \sin \theta \quad (52.2.23)$$

$$v'_2 = -v_1 \sin \theta + v_2 \cos \theta = -x_1^2 \sin \theta + x_2^2 \cos \theta \quad (52.2.24)$$

which clearly isn't the case. Consequently this is not a first order cartesian tensor.

Consider instead the quantity  $\mathbf{u} = (x_2, -x_1)$ , which transforms under rotations as:

$$u'_1 = x'_2 = -x_1 \sin \theta + x_2 \cos \theta \quad (52.2.25)$$

$$u'_2 = -x'_1 = x_1 \cos \theta + x_2 \sin \theta \quad (52.2.26)$$

If this were a first order cartesian tensor, we would need:

$$u'_1 = u_1 \cos \theta + u_2 \sin \theta = x_2 \cos \theta - x_1 \sin \theta \quad (52.2.27)$$

$$u'_2 = -u_1 \sin \theta + u_2 \cos \theta = x_2 \sin \theta - x_1 \cos \theta \quad (52.2.28)$$

which is true for all  $\theta$ . Consequently  $\mathbf{u}$  is a first order cartesian tensor.  $\blacktriangleleft$

### Proposition (Scalar product invariance)

The scalar product of two vectors,  $\mathbf{u} \cdot \mathbf{v}$ , is invariant under rotations.

*Proof.* We consider:

$$u'_i v'_i = R_{ij} u_j R_{ik} u_k = R_{ij} R_{ik} u_j v_k = \delta_{jk} u_j v_k = u_j v_j \quad (52.2.29)$$

as desired. Hence the scalar product of first order cartesian tensors is a zeroth order cartesian tensor.  $\blacksquare$

The definition of a second order cartesian tensor is quite similar to that of a first order cartesian tensor, only that rotations must be repeated twice due to the presence of two indices.

### Definition (Second order Cartesian tensor)

A **second order Cartesian tensor** is defined as a geometric object  $T$  represented by the components  $T_{ij}$  in the cartesian system  $\mathcal{C}$  and represented by the components  $T'_{ij}$  in the cartesian system  $\mathcal{C}'$ , such that they transform under rigid cartesian rotations:

$$T'_{ij} = R_{ik} R_{jl} T_{kl} \quad (52.2.30a)$$

$$T_{kl} = R_{mk} R_{nl} T_{mn} \quad (52.2.30b)$$

or alternatively  $T' = RTR^T$ .

**Example.** The gradient of a vector  $\mathbf{v}^a$ , denoted  $\nabla \mathbf{v}$  is a second order cartesian tensor. Indeed, the components of  $\nabla \mathbf{v}$  are:

$$(\nabla \mathbf{v})_{ij} = \frac{\partial v_i}{\partial x_j} \quad (52.2.31)$$

We may regard  $\left\{ \frac{\partial}{\partial x_j} \right\}$  as a basis, known as the **holonomic basis** or **coordinate basis**. Indeed:

$$\frac{\partial}{\partial x'_i} = \frac{\partial x_j}{\partial x'_i} \frac{\partial}{\partial x_j} = R_{ij} \frac{\partial}{\partial x_j} \implies R_{ij} = \frac{\partial x_j}{\partial x'_i} \quad (52.2.32)$$

from which it follows that  $\nabla$  is a first rank cartesian tensor<sup>b</sup>.

Hence the components of  $\nabla \mathbf{v}$  transform as:

$$(\nabla \mathbf{v})'_{ij} = \frac{\partial v'_i}{\partial x'_j} \quad (52.2.34)$$

$$= \frac{\partial v'_i}{\partial x_k} \frac{\partial x_k}{\partial x'_j} \quad (52.2.35)$$

$$= \frac{\partial}{\partial x_k} (R_{il} v_l) \frac{\partial x_k}{\partial x'_j} \quad (52.2.36)$$

$$= R_{il} \frac{\partial v_l}{\partial x_k} \frac{\partial x_k}{\partial x'_j} \quad (52.2.37)$$

$$= R_{il} R_{jk} (\nabla \mathbf{v})_{lk} \quad (52.2.38)$$

as would be expected from a second order tensor.  $\blacktriangleleft$

<sup>a</sup>this is a first order cartesian tensor

<sup>b</sup>Similarly:

$$R_{ji} = \frac{\partial x'_i}{\partial x_j} \quad (52.2.33)$$

### Definition (Outer product)

The outer product of two vectors  $\mathbf{v}, \mathbf{u}$  is defined as:

$$(\mathbf{u} \otimes \mathbf{v})_{ij} = u_i v_j \quad (52.2.39)$$

and is a second order tensor.

It is easy to see that:

$$(\mathbf{u} \otimes \mathbf{v})'_{ij} = u'_i v'_j = R_{ik} R_{jl} u_k v_l = R_{ik} R_{jl} (\mathbf{u} \otimes \mathbf{v})_{kl} \quad (52.2.40)$$

as desired. Moreover, since  $\mathbf{u} = u_i \mathbf{e}_i$  and  $\mathbf{v} = v_i \mathbf{e}_i$  then:

$$\mathbf{u} \otimes \mathbf{v} = u_i v_j \mathbf{e}_i \otimes \mathbf{e}_j \quad (52.2.41)$$

where  $\mathbf{e}_i \otimes \mathbf{e}_j$  is a sparse matrix with the only non-zero element  $(\mathbf{e}_i \otimes \mathbf{e}_j)_{ij} = 1$ .

**Example.** Consider the matrix:

$$\mathbf{T} = \begin{pmatrix} x_2^2 & -x_1 x_2 \\ -x_1 x_2 & x_1^2 \end{pmatrix} \quad (52.2.42)$$

Using  $s \equiv \sin \theta$  and  $c \equiv \cos \theta$  for shorthand we get that the components of  $\mathbf{T}$  transform as:

$$T'_{11} = (x'_2)^2 = (-x_1 s + x_2 c)^2 = x_1^2 s^2 + x_2^2 c^2 - 2x_1 x_2 c s \quad (52.2.43)$$

$$T'_{12} = -x'_1 x'_2 = -(x_1 c + x_2 s)(-x_1 s + x_2 c) = x_1^2 s c - x_2^2 s c + x_1 x_2 (s^2 - c^2) \quad (52.2.44)$$

$$T'_{21} = x_1^2 s c - x_2^2 s c + x_1 x_2 (s^2 - c^2) \quad (52.2.45)$$

$$T'_{22} = (x'_1)^2 = (x_1 c + x_2 s)^2 = x_1^2 c^2 + x_2^2 s^2 + 2x_1 x_2 s c \quad (52.2.46)$$

so that:

$$T' = \begin{pmatrix} x_1^2 s^2 + x_2^2 c^2 - 2x_1 x_2 c s & x_1^2 s c - x_2^2 s c + x_1 x_2 (s^2 - c^2) \\ x_1^2 s c - x_2^2 s c + x_1 x_2 (s^2 - c^2) & x_1^2 c^2 + x_2^2 s^2 + 2x_1 x_2 s c \end{pmatrix} \quad (52.2.47)$$

If  $T$  were a tensor then we would find

$$T' = RTR^T = \begin{pmatrix} c & s \\ -s & c \end{pmatrix} \begin{pmatrix} x_2^2 & -x_1 x_2 \\ -x_1 x_2 & x_1^2 \end{pmatrix} \begin{pmatrix} c & -s \\ s & c \end{pmatrix} \quad (52.2.48)$$

$$= \begin{pmatrix} c & s \\ -s & c \end{pmatrix} \begin{pmatrix} x_2^2 c - x_1 x_2 s & -x_2^2 s - x_1 x_2 c \\ -x_1 x_2 c + s^2 s & x_1 x_2 s + x_1^2 c \end{pmatrix} \quad (52.2.49)$$

$$= \begin{pmatrix} x_1^2 s^2 + x_2^2 c^2 - 2x_1 x_2 c s & x_1^2 s c - x_2^2 s c + x_1 x_2 (s^2 - c^2) \\ x_1^2 s c - x_2^2 s c + x_1 x_2 (s^2 - c^2) & x_1^2 c^2 + x_2^2 s^2 + 2x_1 x_2 s c \end{pmatrix} \quad (52.2.50)$$

so  $T$  is indeed a second order cartesian tensor.

More simply, we could have noticed that  $T = (x_2, -x_1) \otimes (x_2, -x_1)$  and since  $(x_2, -x_1)$  was proven to be a first order cartesian tensor,  $T$  will be a second order cartesian tensor as desired.  $\blacktriangleleft$

**Definition (Cartesian Tensor)** In general, a cartesian tensor is defined as a geometric object  $T$  represented by the components  $T_{ij..k}$  in the cartesian system  $\mathcal{C}$  and represented by the components  $T'_{ij...k}$  in the cartesian system  $\mathcal{C}'$ , such that they transform under rigid cartesian rotations:

$$T'_{ij...k} = R_{ip} R_{jq} \dots R_{kr} T_{pq...r} \quad (52.2.51a)$$

$$T_{ij...k} = R_{pi} R_{qj} \dots R_{rk} T'_{pq...r} \quad (52.2.51b)$$

### Theorem (Quotient law)

Suppose that  $\mathbf{B}$  and  $\mathbf{C}$  are tensors such that its components in any rotated basis satisfy:

$$A_{pq...k...m} B_{ij...k...n} = C_{pq...mij...n} \quad (52.2.52)$$

then  $\mathbf{A}$  must also be a tensor.

*Proof.* We consider the case for second order tensors (the general case follows exactly the same logic, just with more indices). We are given that:

$$A_{pk} B_{ik} = C_{pi}, \quad A'_{pk} B'_{ik} = C'_{pi} \quad (52.2.53)$$

and that  $B_{jl} = R_{mj} R_{nl} B'_{mn}$ ,  $C'_{pi} = R_{pq} R_{ij} C_{qj}$  then we get that:

$$A'_{pk} B'_{ik} = C'_{pi} \quad (52.2.54)$$

$$= R_{pq} R_{ij} C_{qj} \quad (52.2.55)$$

$$= R_{pq} R_{ij} A_{ql} B_{jl} \quad (52.2.56)$$

$$= R_{pq} R_{ij} A_{ql} R_{mj} R_{nl} B'_{mn} \quad (52.2.57)$$

$$= R_{pq} R_{nl} A_{ql} B'_{in} \quad (52.2.58)$$

so that:

$$(A'_{pk} - R_{pq}R_{nl}A_{ql})B'_{ik} = 0 \quad (52.2.59)$$

Since this must hold for any  $B'_{ik}$  we must have that:

$$A'_{pk} = R_{pq}R_{nl}A_{ql} \quad (52.2.60)$$

as desired. ■

The quotient law is a much faster way to prove that a quantity is a tensor, since it suffices to contract this quantity with a known tensor and ensure that the resulting quantity is also a tensor.

**Example.** Let's prove that:

$$\mathbf{T} = \begin{pmatrix} x_2^2 & -x_1x_2 \\ -x_1x_2 & x_1^2 \end{pmatrix} \quad (52.2.61)$$

is a tensor. We have already done this in two ways, but a third way is by using the quotient law:

$$T_{11}x_1^2 = x_1^2x_2^2 \quad (52.2.62)$$

$$T_{12}x_1x_2 = -x_1^2x_2^2 \quad (52.2.63)$$

$$T_{21}x_2x_1 = -x_1^2x_2^2 \quad (52.2.64)$$

$$T_{22}x_2^2 = x_1^2x_2^2 \quad (52.2.65)$$

so that  $T_{ij}x_i x_j = 0$  which is a tensor. Since  $x_i x_j$  is an outer product and thus a tensor, it follows that  $\mathbf{T}$  is also a tensor. ◀

## 52.3 The $\delta_{ij}$ and $\epsilon_{ijk}$ tensors

## 52.4 Physical examples of cartesian tensors

Consider the angular momentum  $\mathbf{L} = \mathbf{r} \times \mathbf{p} = m\mathbf{r} \times \dot{\mathbf{r}}$ , this is a first order cartesian tensor.

Indeed if we write the components of  $\mathbf{L}$  as:

$$L_i = m\epsilon_{ijk}r_j \dot{r}_k \quad (52.4.1)$$

they they will transform under rotations as:

$$L'_i = m\epsilon_{ijk}r'_j \dot{r}'_k \quad (52.4.2)$$

and since  $r'_j = R_{jm}r_m$  and  $\dot{r}'_k = \frac{d}{dt}(R_{kn}r_n) = R_{kn}\dot{r}_n$  we will find that.

$$L'_i = m\epsilon_{ijk}R_{jm}R_{kn}r_m \dot{r}_n \quad (52.4.3)$$

Consider the  $i, j$  component of the **inertia tensor**:

$$I_{ij} = \int (\delta_{ij}r^2 - x_i x_j) dm \quad (52.4.4)$$

In dirac notation, it is easy to see that  $I_{ij} = \langle i|I|j\rangle$ ,  $\delta_{ij} = \langle i|j\rangle$ ,  $r^2 = \langle r|r\rangle$ ,  $x_i = \langle r|i\rangle = \langle i|r\rangle$ . Consequently:

$$\langle i|I|j\rangle = \langle i| \int \langle r|r\rangle \mathbb{I} - |r\rangle \langle r| dm |j\rangle \quad (52.4.5)$$

from which it follows that:

$$I = \int r^2 \mathbb{I} - dm \quad (52.4.6)$$

This is a second order cartesian tensor.

## 52.5 Non-cartesian tensors

### Definition (*Contravariant and covariant bases*)

We saw that for general curvilinear coordinates  $(u_1, u_2, u_3)$ , so that a position vector may be expressed as  $\mathbf{r}(u_1, u_2, u_3)$  then we have two important sets of basis vectors:

$$\mathbf{e}_i = \frac{\partial \mathbf{r}}{\partial u^i}, \quad \mathbf{e}^i = \nabla u_i \quad (52.5.1)$$

We have slightly modified our summation convention to include superscripts. From now on, we will assume that any lower case index that appears exactly once as a superscript and once as a subscript must be summed over.

### Proposition (*Reciprocity relation*)

The contravariant and covariant bases are orthonormal to each other:

$$\mathbf{e}_i \cdot \mathbf{e}^j = \delta_i^j \quad (52.5.2)$$

*Proof.* We find that:

$$\mathbf{e}_i \cdot \mathbf{e}^j = \frac{\partial \mathbf{r}}{\partial u^i} \nabla u_j \quad (52.5.3)$$

$$= \frac{\partial u_1}{\partial u^i} \frac{\partial u_j}{\partial u_1} + \frac{\partial u_2}{\partial u^i} \frac{\partial u_j}{\partial u_2} + \frac{\partial u_3}{\partial u^i} \frac{\partial u_j}{\partial u_3} \quad (52.5.4)$$

$$= \frac{\partial u_j}{\partial u^i} = \delta_i^j \quad (52.5.5)$$

as desired. ■

Consequently, given some vector  $\mathbf{a}$  it may be expanded in both the contravariant and covariant bases. We find that if

$$\mathbf{a} = a^i \mathbf{e}_i = a_i \mathbf{e}^i \quad (52.5.6)$$

where  $a^i$  are the contravariant components while  $a_i$  are the covariant components, then:

$$\mathbf{a} \cdot \mathbf{e}^j = a^i \delta_i^j = a^j \quad (52.5.7)$$

$$\mathbf{a} \cdot \mathbf{e}_j = a_i \delta_j^i = a_j \quad (52.5.8)$$

as expected.

Let's now consider an infinitesimal vector displacement  $d\mathbf{r} = du^i \mathbf{e}_i$ . Then, the infinitesimal arc length is:

$$(ds)^2 = du^i du^j \mathbf{e}_i \cdot \mathbf{e}_j = g_{ij} du^i du^j \quad (52.5.9)$$

where  $g_{ij} \equiv \mathbf{e}_i \mathbf{e}_j$  is defined as the metric tensor.

### Definition (Metric tensor)

For a given set of curvilinear coordinates  $(u_i)$  with contravariant and covariant bases  $\{\mathbf{e}_i = \frac{\partial \mathbf{r}}{\partial u_i}\}$  and  $\{\mathbf{e}^i = \nabla u_i\}$  respectively, the metric tensor is defined to be:

$$g_{ij} \equiv \mathbf{e}_i \mathbf{e}_j \quad (52.5.10)$$

Furthermore, the volume element may be expressed as:

$$dV = \sqrt{g} du^i \quad (52.5.11)$$

where  $g = \det[g_{ij}]$  is the determinant of the metric tensor.

**Example.** Let's evaluate the metric tensor in spherical polar coordinates. We have that the position vector of some general point  $(u_1, u_2, u_3) = (r, \theta, \phi)$  is given by:

$$\mathbf{r} = r \sin \theta \cos \phi \mathbf{e}_x + r \sin \theta \sin \phi \mathbf{e}_y + r \cos \theta \mathbf{e}_z \quad (52.5.12)$$

The covariant basis is easily found to be:

$$\mathbf{e}_1 = \frac{\partial \mathbf{r}}{\partial r} = \sin \theta \cos \phi \mathbf{e}_x + \sin \theta \sin \phi \mathbf{e}_y + \cos \theta \mathbf{e}_z \quad (52.5.13)$$

$$\mathbf{e}_2 = \frac{\partial \mathbf{r}}{\partial \theta} = r \cos \theta \cos \phi \mathbf{e}_x + r \cos \theta \sin \phi \mathbf{e}_y - r \sin \theta \mathbf{e}_z \quad (52.5.14)$$

$$\mathbf{e}_3 = \frac{\partial \mathbf{r}}{\partial \phi} = -r \sin \theta \sin \phi \mathbf{e}_x + r \sin \theta \cos \phi \mathbf{e}_y \quad (52.5.15)$$

Consequently:

$$[g_{ij}] = \begin{pmatrix} 1 & 0 & 0 \\ 0 & r^2 & 0 \\ 0 & 0 & r^2 \sin^2 \theta \end{pmatrix} \quad (52.5.16)$$

As expected from an orthogonal coordinate system, the metric is diagonal. Moreover, we find that the infinitesimal arc length may be expressed as

$$(ds)^2 = dr^2 + r^2 d\theta^2 + r^2 \sin^2 \theta d\phi^2 \quad (52.5.17)$$

while the volume element is found to be:

$$dV = r^2 \sin \theta dr d\theta d\phi \quad (52.5.18)$$

Just like vectors, tensors may also be expressed in both bases:

$$\mathbf{T} = T^{ij} \mathbf{e}_i \otimes \mathbf{e}_j = T_{ij} \mathbf{e}^i \otimes \mathbf{e}^j = T_i^j \mathbf{e}^i \otimes \mathbf{e}_j \quad (52.5.19)$$

Here  $T^{ij}, T_{ij}, T_i^j$  are known as the contravariant, covariant and mixed tensor components respectively.

We can use the metric tensor to express the scalar product of two vectors. Indeed, using the covariant and contravariant expansions:

$$\mathbf{a} \cdot \mathbf{b} = a^i \mathbf{e}_i b^j \mathbf{e}_j = g_{ij} a^i b^j = a_i \mathbf{e}^i b_j \mathbf{e}^j = g^{ij} a_i b_j \quad (52.5.20)$$

Finally, we also find that:

$$\mathbf{a} \cdot \mathbf{b} = a^i \mathbf{e}_i b_j \mathbf{e}^j = a^i b_j \delta_i^j = a^i b_i = a_i b^i \quad (52.5.21)$$

Consequently:

$$g_{ij} a^i b^j = a^i b_i, \quad g_{ij} a_i b_j = a_i b^i \quad (52.5.22)$$

This holds for any arbitrary  $\mathbf{a}$ , so we get the following very useful result:

### Theorem (Raising and lowering indices)

For a given vector  $\mathbf{b}$  expressed in a set of curvilinear coordinates with metric tensor  $g_{ij}$ , its covariant and contravariant components are related by:

$$g_{ij} b^j = b_i, \quad g^{ij} b_j = b^i \quad (52.5.23)$$

The special case when  $\mathbf{b} = \mathbf{e}_i$  then:

$$\mathbf{e}^i = g^{ij} \mathbf{e}_j, \quad \mathbf{e}_i = g_{ij} \mathbf{e}^j \quad (52.5.24)$$

### Proposition (Contravariant and covariant components are inverses)

The contravariant  $g^{ij}$  and covariant  $g_{ij}$  components of a metric tensor obey:

$$g^{ij} g_{jk} = \delta_k^i \quad (52.5.25)$$

*Proof.* Consider:

$$a^i = \delta_k^i a^k = g^{ij} a_j = g^{ij} g_{jk} a^k \implies \delta_k^i = g^{ij} g_{jk} \quad (52.5.26)$$

as desired. ■

## 52.6 Covariance and contravariance

Let's now see what happens how covariant and contravariant components transform when we perform a change of curvilinear coordinates, since this is how we initially defined a tensor.

Suppose we have a coordinate system  $\{u^i\}$  which we transform to another system  $\{u'^i\}$ . The new sets of basis vectors are:

$$\mathbf{e}'_i = \frac{\partial \mathbf{r}}{\partial u'^i}, \quad \mathbf{e}'^i = \nabla u'^i \quad (52.6.1)$$

The new covariant basis can be related to the old one by:

$$\mathbf{e}'_i = \frac{\partial \mathbf{r}}{\partial u'^i} = \frac{\partial u^j}{\partial u'^i} \frac{\partial \mathbf{r}}{\partial u^j} = \frac{\partial u^j}{\partial u'^i} \mathbf{e}_j \quad (52.6.2)$$

or alternatively:

$$\mathbf{e}_j = \frac{\partial u'^i}{\partial u^j} \mathbf{e}'_i \quad (52.6.3)$$

Therefore, expanding an arbitrary vector in both covariant bases:

$$\mathbf{a} = a^i \mathbf{e}_i = a^i \frac{\partial u'^j}{\partial u^i} \mathbf{e}'_j = a'^j \mathbf{e}'_j \quad (52.6.4)$$

implying that:

$$a'^j = \frac{\partial u'^j}{\partial u^i} a^i \quad (52.6.5)$$

Similarly, the new contravariant basis can be related to the old one by:

$$\mathbf{e}'^i = \nabla u'^i = \frac{\partial u'^i}{\partial u^j} \nabla u^j = \frac{\partial u'^i}{\partial u^j} \mathbf{e}^j \quad (52.6.6)$$

or alternatively:

$$\mathbf{e}^i = \frac{\partial u^i}{\partial u'^j} \mathbf{e}'^j \quad (52.6.7)$$

Therefore, expanding an arbitrary vector in both contravariant bases we find that:

$$\mathbf{a} = a_i \mathbf{e}^i = a_i \frac{\partial u^i}{\partial u'^j} \mathbf{e}'^j = a'_j \mathbf{e}'^j \quad (52.6.8)$$

implying that covariant components transform as:

$$a'_j = \frac{\partial u^i}{\partial u'^j} a_i \quad (52.6.9)$$

**Theorem (Transformation of co(ntra)variant components)** For a given vector  $\mathbf{a}$ , its covariant and contravariant components transform under a change of basis as:

$$a'^j = \frac{\partial u'^j}{\partial u^i} a^i \quad (52.6.10)$$

$$a'_j = \frac{\partial u^i}{\partial u'^j} a_i \quad (52.6.11)$$

In a completely analogous way, we can show how the co(ntra)variant and mixed components of a second rank tensor transform:

$$T'^{ij} = \frac{\partial u'^i}{\partial u^k} \frac{\partial u'^j}{\partial u^l} T^{kl} \quad (52.6.12)$$

$$T'^i_j = \frac{\partial u'^i}{\partial u^k} \frac{\partial u^l}{\partial u'^j} T^k_l \quad (52.6.13)$$

$$T'_{ij} = \frac{\partial u^k}{\partial u'^i} \frac{\partial u^l}{\partial u'^j} T_{kl} \quad (52.6.14)$$

*Proof.* We have already proven the result for vectors. For second rank tensors, we find:

$$T'^{ij} \mathbf{e}'_i \otimes \mathbf{e}'_j = T^{ij} \mathbf{e}_i \otimes \mathbf{e}_j \quad (52.6.15)$$

$$T'_i{}^j' \mathbf{e}'^i \otimes \mathbf{e}'_j = T_i{}^j \mathbf{e}^i \otimes \mathbf{e}_j \quad (52.6.16)$$

$$T'_{ij} \mathbf{e}'^i \otimes \mathbf{e}'^j = T_{ij} \mathbf{e}^i \otimes \mathbf{e}^j \quad (52.6.17)$$

For example, we would find that:

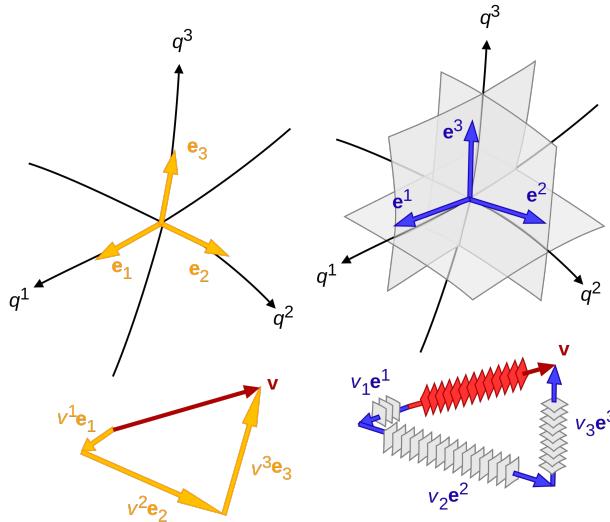
$$T'^{ij} \mathbf{e}'_i \otimes \mathbf{e}'_j = T^{kl} \mathbf{e}_k \otimes \mathbf{e}_l = T^{kl} \frac{\partial u'^i}{\partial u^k} \frac{\partial u'^j}{\partial u^l} \mathbf{e}'_i \otimes \mathbf{e}'_j \quad (52.6.18)$$

implying

$$T'^{ij} = \frac{\partial u'^i}{\partial u^k} \frac{\partial u'^j}{\partial u^l} T^{kl} \quad (52.6.19)$$

as desired. ■

Visually, we can explain covariant and contravariant components as follows: “contravariant components transform as position vector components, while covariant components transform as gradient vector components”. This makes physically sense, as the contravariant basis is parallel everywhere to its coordinate curves, while the covariant basis is orthonormal everywhere to its coordinate surfaces. This is in alignment with the intuition that coordinate curves transform as position vectors, while coordinate surfaces transform as gradient vectors.



**Figure 52.1.** By Maschen - Own work, CC0, <https://commons.wikimedia.org/w/index.php?curid=21043814>

Clearly, if we increase the lengths of the contravariant basis vectors, the contravariant components decrease (less arrows along each basis vector are needed), and vice versa. The contravariant components “contra-vary” with the change of basis.

On the other hand, if I increase the lengths of the covariant basis vectors, the covariant components increase (“more” planes can be packed along each basis vector), and vice versa. The covariant components “co-vary” with the change of basis.

**Theorem (Metric tensor is a second rank tensor)**

The metric tensor  $g_{\mu\nu}$  is a second-rank tensor.

*Proof.* We have that :

$$g_{\mu\nu}dx^\mu dx^\nu = g'_{\alpha\beta}dx'^\alpha dx^\beta \quad (52.6.20)$$

$$= g'_{\alpha\beta} \frac{\partial x'^\alpha}{\partial x^\mu} \frac{\partial x'^\beta}{\partial x^\nu} dx^\mu dx^\nu \quad (52.6.21)$$

$$\implies \left( g_{\mu\nu} - g'_{\alpha\beta} \frac{\partial x'^\alpha}{\partial x^\mu} \frac{\partial x'^\beta}{\partial x^\nu} \right) dx^\mu dx^\nu = 0 \quad (52.6.22)$$

Seeing as this must hold for all  $dx^\mu, dx^\nu$ , it follows that:

$$g_{\mu\nu} = \frac{\partial x'^\alpha}{\partial x^\mu} \frac{\partial x'^\beta}{\partial x^\nu} g'_{\alpha\beta} \quad (52.6.23)$$

proving the required transformation rule. ■

## 52.7 Application to special relativity: four-vectors

In special relativity we deal with four vectors whose components transform like:

$$v'^\mu = \frac{\partial x'^\mu}{\partial x^\nu} v^\nu \quad (52.7.1)$$

where  $\frac{\partial x'^\mu}{\partial x^\nu}$  are components of the Lorentz transformation. An important four vector is the displacement vector:

$$dx^\mu = (cdt, dx, dy, dz) = (dx^0, dx^1, dx^2, dx^3) = (cdt, d\mathbf{r}) \quad (52.7.2)$$

We also know that the following quantity, known as the proper time, is an invariant under Lorentz transformations:

$$(cd\tau)^2 = c^2 dt^2 - d\mathbf{r} \cdot d\mathbf{r} \quad (52.7.3)$$

This suggests introducing the following metric tensor, known as the Minkowski metric:

$$[\eta_{\mu\nu}] = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \\ 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & -1 \end{pmatrix} \implies (cd\tau)^2 = g_{\mu\nu} dx^\mu dx^\nu \quad (52.7.4)$$

This allows us to find the covariant components:

$$dx_\mu = (cdt, -d\mathbf{r}) \quad (52.7.5)$$

The importance of the proper time arises when we try to define a velocity four vector. Suppose we naively define it to be:

$$u^\mu = \frac{\partial x^\mu}{\partial t} = \left( c, \frac{d\mathbf{r}}{dt} \right) \quad (52.7.6)$$

This result is quite worrisome, since the first component  $c$  of this supposed four vector is constant, and therefore does not transform at all under a change of basis.

Suppose we instead define

$$u^\mu = \frac{\partial x^\mu}{\partial \tau} = \frac{\partial x^\mu}{\partial t} \frac{\partial t}{\partial \tau} = \left( c, \frac{d\mathbf{r}}{dt} \right) \frac{dt}{d\tau} \quad (52.7.7)$$

Then, since

$$c^2 d\tau^2 = c^2 dt^2 - d\mathbf{r} \cdot d\mathbf{r} = dt^2 \left( c^2 - \frac{d\mathbf{r}}{dt} \cdot \frac{d\mathbf{r}}{dt} \right) \quad (52.7.8)$$

$$\implies \frac{dt}{d\tau} = \frac{1}{\sqrt{1 - v^2/c^2}} \equiv \gamma \quad (52.7.9)$$

where we defined  $\mathbf{v} \equiv \frac{d\mathbf{r}}{dt}$ , we find that the velocity four vector reads:

$$u^\mu = \gamma(c, \mathbf{v}) \quad (52.7.10)$$

and thus:

$$u_\mu = \gamma(c, -\mathbf{v}) \quad (52.7.11)$$

# Tensor calculus

## 53.1 Christoffel symbols

Let's consider how we can take derivatives of basis vectors. For example, consider  $\frac{\partial \mathbf{e}_i}{\partial u^j}$ . Since this is itself a vector, it can be expanded in the covariant basis as::

$$\boxed{\frac{\partial \mathbf{e}_i}{\partial u^j} = \Gamma_{ij}^k \mathbf{e}_k} \quad (53.1.1)$$

It is easy to see that:

$$\Gamma_{ij}^k = \mathbf{e}^k \cdot \frac{\partial \mathbf{e}_i}{\partial u^j} \quad (53.1.2)$$

Also, we can differentiate the relation  $\mathbf{e}^i \cdot \mathbf{e}_j$  to find:

$$\frac{\partial}{\partial u^k} (\mathbf{e}^i \cdot \mathbf{e}_j) = \frac{\partial \mathbf{e}^i}{\partial u^k} \cdot \mathbf{e}_j + \mathbf{e}^i \cdot \frac{\partial \mathbf{e}_j}{\partial u^k} = 0 \quad (53.1.3)$$

$$\implies \frac{\partial \mathbf{e}^i}{\partial u^k} \cdot \mathbf{e}_j = -\Gamma_{jk}^i \quad (53.1.4)$$

$$\implies \boxed{\frac{\partial \mathbf{e}^i}{\partial u^j} = -\Gamma_{kj}^i \mathbf{e}^k} \quad (53.1.5)$$

The  $\Gamma_{ij}^k$  are known as Christoffel symbols. Although it looks like a third rank tensor, it actually does not follow the required transformation laws. Indeed:

$$\Gamma'_{ij}^k = \mathbf{e}'^k \cdot \frac{\partial \mathbf{e}'_i}{\partial u'^j} \quad (53.1.6)$$

$$= \frac{\partial u'^k}{\partial u^n} \mathbf{e}^n \cdot \frac{\partial}{\partial u'^k} \left( \frac{\partial u^l}{\partial u'^i} \mathbf{e}_l \right) \quad (53.1.7)$$

$$= \frac{\partial u'^k}{\partial u^n} \mathbf{e}^n \cdot \left( \frac{\partial^2 u^l}{\partial u'^j \partial u'^i} \mathbf{e}_l + \frac{\partial u^l}{\partial u'^i} \frac{\partial \mathbf{e}_l}{\partial u^m} \frac{\partial u^m}{\partial u'^j} \right) \quad (53.1.8)$$

$$= \frac{\partial u'^k}{\partial u^n} \frac{\partial^2 u^l}{\partial u'^j \partial u'^i} \delta_l^n + \frac{\partial u'^k}{\partial u^n} \frac{\partial u^l}{\partial u'^i} \frac{\partial u^m}{\partial u'^j} \mathbf{e}_n \cdot \frac{\partial \mathbf{e}_l}{\partial u^m} \quad (53.1.9)$$

$$= \frac{\partial u'^k}{\partial u^l} \frac{\partial^2 u^l}{\partial u'^j \partial u'^i} + \Gamma_{lm}^n \frac{\partial u'^k}{\partial u^n} \frac{\partial u^l}{\partial u'^i} \frac{\partial u^m}{\partial u'^j} \quad (53.1.10)$$

We can find another expression for the Christoffel symbols that allow for faster computation. Consider the derivative of the metric tensor:

$$\frac{dg_{ij}}{du^k} = \frac{\partial \mathbf{e}_i}{\partial u^k} \cdot \mathbf{e}_j + \mathbf{e}_i \cdot \frac{\partial \mathbf{e}_j}{\partial u^k} \quad (53.1.11)$$

$$= \Gamma_{ik}^l \mathbf{e}_l \cdot \mathbf{e}_j + \mathbf{e}_i \cdot \mathbf{e}_l \Gamma_{jk}^l \quad (53.1.12)$$

$$\implies \frac{dg_{ij}}{du^k} = \Gamma_{ik}^l g_{lj} + \Gamma_{jk}^l g_{il} \quad (53.1.13)$$

Now note that  $\Gamma_{ik}^l$  is symmetric, since:

$$\frac{\partial \mathbf{e}_i}{\partial u^j} = \frac{\partial^2 \mathbf{r}}{\partial u^i \partial u^j} = \frac{\partial \mathbf{e}_j}{\partial u^i} \quad (53.1.14)$$

Consequently, we can simply permute the indices in (53.1.13) and find the following:

$$\frac{dg_{ik}}{du^j} = \Gamma_{ij}^l g_{lk} + \Gamma_{jk}^l g_{il} \quad (53.1.15)$$

$$\frac{dg_{kj}}{du^i} = \Gamma_{ik}^l g_{lj} + \Gamma_{ji}^l g_{kl} \quad (53.1.16)$$

implying that:

$$\frac{dg_{ij}}{du^k} + \frac{dg_{ik}}{du^j} - \frac{dg_{kj}}{du^i} = 2\Gamma_{jk}^l g_{il} \quad (53.1.17)$$

$$\iff \left( \frac{dg_{ij}}{du^k} + \frac{dg_{ik}}{du^j} - \frac{dg_{kj}}{du^i} \right) g^{im} = 2\Gamma_{jk}^l g_{il} g^{im} \quad (53.1.18)$$

$$\boxed{\iff \Gamma_{jk}^m = \frac{1}{2} g^{im} \left( \frac{dg_{ij}}{du^k} + \frac{dg_{ik}}{du^j} - \frac{dg_{kj}}{du^i} \right)} \quad (53.1.19)$$

## 53.2 Differentiating tensors

Let's consider a contravariant vector:

$$x'^i = \frac{\partial x'^i}{\partial x^j} x^j \quad (53.2.1)$$

Its time derivative is transforms following:

$$\frac{\partial x'^i}{\partial u'^j} = \frac{\partial u^k}{\partial u'^j} \frac{\partial x'^i}{\partial u^k} \quad (53.2.2)$$

$$= \frac{\partial u^k}{\partial u'^j} \frac{\partial}{\partial u^k} \left( \frac{\partial x'^i}{\partial u^l} x^l \right) \quad (53.2.3)$$

$$= \frac{\partial u^k}{\partial u'^j} \left( \frac{\partial^2 u'^i}{\partial u^l \partial u^k} x^l + \frac{\partial u'^i}{\partial u^l} \frac{\partial x^l}{\partial u^k} \right) \quad (53.2.4)$$

$$= \frac{\partial u^k}{\partial u'^j} \frac{\partial^2 u'^i}{\partial u^j \partial u^k} x^j + \frac{\partial u^k}{\partial u'^j} \frac{\partial u'^i}{\partial u^l} \frac{\partial x^l}{\partial u^k} \quad (53.2.5)$$

The  $\frac{\partial u^k}{\partial u'^j} \frac{\partial^2 u'^i}{\partial u^j \partial u^k} x^j$  term ruins the transformation rule, and if it is non-zero then  $\frac{\partial x'^i}{\partial x^j}$  cannot be considered a tensor. This is quite problematic, as from our experience things such as the velocity vector should be a tensor.

Luckily, we can use the Christoffel symbols to deal with this issue. Consider:

$$\frac{\partial \mathbf{x}}{\partial u^j} = \frac{\partial x^i}{\partial u^j} \mathbf{e}_i + x^i \Gamma_{ij}^k \mathbf{e}_k \quad (53.2.6)$$

$$= \left( \frac{\partial x^i}{\partial u^j} + x^k \Gamma_{kj}^i \right) \mathbf{e}_i \quad (53.2.7)$$

leading us to defining the covariant derivative as follows:

### Definition (Covariant derivative)

Given a contravariant representation of a vector  $\mathbf{v} = v^i \mathbf{e}_i$ , its **covariant derivative** is defined as:

$$v_{;j}^i = \frac{\partial x^i}{\partial u^j} + x^k \Gamma_{kj}^i \quad (53.2.8)$$

and similarly for the covariant representation:

$$v_{i;j} = \frac{\partial x_i}{\partial u^j} - x_k \Gamma_{ij}^k \quad (53.2.9)$$

We see that the covariant derivative is different from the normal partial derivative in that the basis vectors that don't change over space for cartesian coordinates are variable in general coordinates. Due to the product derivative rule, this creates an extra term in the derivative which we identify using a Christoffel symbol.

**Example.** Let's work out the covariant derivative of the contravariant components of a second order tensor  $\mathbb{T}$ .

The contravariant components  $T^{ij}$  satisfy:

$$\mathbb{T} = T^{ij} (\mathbf{e}_i \otimes \mathbf{e}_j) \quad (53.2.10)$$

Consequently:

$$\frac{\partial \mathbb{T}}{\partial u^k} = \frac{\partial T^{ij}}{\partial u^k} (\mathbf{e}_i \otimes \mathbf{e}_j) + T^{ij} \frac{\partial}{\partial u^k} (\mathbf{e}_i \otimes \mathbf{e}_j) \quad (53.2.11)$$

We can simplify the second term on the RHS:

$$\frac{\partial}{\partial u^k} (\mathbf{e}_i \otimes \mathbf{e}_j) = \frac{\partial \mathbf{e}_i}{\partial u^k} \otimes \mathbf{e}_j + \mathbf{e}_i \otimes \frac{\partial \mathbf{e}_j}{\partial u^k} \quad (53.2.12)$$

$$= \Gamma_{ik}^l \mathbf{e}_l \otimes \mathbf{e}_j + \mathbf{e}_i \otimes \mathbf{e}_l \Gamma_{jk}^l \quad (53.2.13)$$

giving

$$\frac{\partial \mathbb{T}}{\partial u^k} = \frac{\partial T^{ij}}{\partial u^k} (\mathbf{e}_i \otimes \mathbf{e}_j) + T^{ij} (\Gamma_{ik}^l \mathbf{e}_l \otimes \mathbf{e}_j + \mathbf{e}_i \otimes \mathbf{e}_l \Gamma_{jk}^l) \quad (53.2.14)$$

$$= \left( \frac{\partial T^{ij}}{\partial u^k} + T^{lj} \Gamma_{lk}^i + T^{il} \Gamma_{lk}^j \right) (\mathbf{e}_i \otimes \mathbf{e}_j) \quad (53.2.15)$$

$$= T_{;k}^{ij} (\mathbf{e}_i \otimes \mathbf{e}_j) \quad (53.2.16)$$

where we defined a covariant derivative:

$$T_{:k}^{ij} \equiv \frac{\partial T^{ij}}{\partial u^k} + T^{lj}\Gamma_{lk}^i + T^{il}\Gamma_{lk}^j \quad (53.2.17)$$

◀

### 53.3 Application to geometry: curvilinear coordinates

Let's apply our knowledge of tensor calculus to study general curvilinear coordinates.

#### Gradient

The gradient of a scalar field  $\phi$  is just:

$$\nabla\phi \equiv \phi_{:i}\mathbf{e}^i = \frac{\partial\phi}{\partial u^i}\mathbf{e}^i \quad (53.3.1)$$

#### Divergence

The divergence of a vector field  $\mathbf{v}$  is:

$$\nabla \cdot \mathbf{x} \equiv v_{:i}^i = \frac{\partial v^i}{\partial u^i} + x^k\Gamma_{ki}^i \quad (53.3.2)$$

Note that the Christoffel symbol simplifies significantly:

$$\Gamma_{ki}^i = \frac{1}{2}g^{il}\left(\frac{\partial g_{il}}{\partial u^k} + \frac{\partial g_{kl}}{\partial u^i} - \frac{\partial g_{ki}}{\partial u^l}\right) \quad (53.3.3)$$

$$= \frac{1}{2}g^{il}\left(\frac{\partial g_{il}}{\partial u^k} + \frac{\partial g_{kl}}{\partial u^i} - \frac{\partial g_{kl}}{\partial u^i}\right) \quad (53.3.4)$$

$$= \frac{1}{2}g^{il}\frac{\partial g_{il}}{\partial u^k} \quad (53.3.5)$$

giving:

$$\nabla \cdot \mathbf{x} = \frac{\partial v^i}{\partial u^i} + \frac{1}{2}x^k g^{il} \frac{\partial g_{il}}{\partial u^k} \quad (53.3.6)$$

#### Proposition (Important determinant identity)

Let  $A = (a_{ij})$ ,  $B = (b^{ij})$  and  $A = B^{-1}$ . Then:

$$\frac{\partial|A|}{\partial u^k} = |A|b^{ji}\frac{\partial a_{ij}}{\partial u^k} \quad (53.3.7)$$

*Proof.* We have that if the cofactor of the element  $a_{ij}$  is  $c^{ij}$  then:

$$b^{ij} = \frac{1}{|A|}c^{ji} \quad (53.3.8)$$

Consequently, since  $|\mathbf{A}| = a_{ij}c^{ij}$  with fixed  $i$ :

$$\frac{\partial |\mathbf{A}|}{\partial a_{ij}} = c^{ij} = |\mathbf{A}|b^{ij} \quad (53.3.9)$$

implying:

$$\frac{\partial |\mathbf{A}|}{\partial u^k} = \frac{\partial |\mathbf{A}|}{\partial a_{ij}} \frac{\partial a_{ij}}{\partial u^k} = |\mathbf{A}|b^{ij} \frac{\partial a_{ij}}{\partial u^k} \quad (53.3.10)$$

as desired. ■

We can apply this proposition to the determinant of the metric tensor:

$$\frac{\partial |g|}{\partial u^k} = |g|g^{ij} \frac{\partial g_{ij}}{\partial u^k} \implies g^{ij} \frac{\partial g_{ij}}{\partial u^k} = \frac{1}{|g|} \frac{\partial |g|}{\partial u^k} = \frac{1}{\sqrt{|g|}} \frac{\partial \sqrt{|g|}}{\partial u^k} \quad (53.3.11)$$

and hence when substituted into (53.3.6) we get:

$$\nabla \cdot \mathbf{x} = \frac{\partial v^i}{\partial u^i} + \frac{1}{\sqrt{|g|}} \frac{\partial \sqrt{|g|}}{\partial u^i} x^i \quad (53.3.12)$$

or alternatively:

$\nabla \cdot \mathbf{v} = \frac{1}{\sqrt{|g|}} \frac{\partial (\sqrt{|g|} x^i)}{\partial u^i}$

(53.3.13)

### Laplacian

### Curl

## 53.4 Geodesics

## **Part VI**

# **Differential Geometry**

# Topology

## 54.1 Why differential geometry?

Broadly speaking, a manifold is a space that looks locally like  $\mathbb{R}^n$  even though it may not be so globally. For example, the surface of the Earth may be viewed as a locally flat surface embedded in  $\mathbb{R}^3$ . Thus most tools from multivariable calculus may be employed. In general however, we will discover that space-time does not seem to be naturally embedded in some higher-dimensional Euclidean space. Traditional calculus is unequipped to deal with such objects, we must develop a new formalism which is known nowadays as modern differential geometry.

## 54.2 Topology

We begin by introducing some important basic concepts from topology. Indeed, at its most fundamental (i.e. with least structure) level space-time may be viewed as a simple collection of points. This alone is not enough to talk about things such as derivatives etc..., extra conditions must be satisfied. However, we do not want to endow space-time with too much structure, just enough so that good old calculus may be used. For example, the way we may classify these points into sub-collections will be useful in giving structure to space-time. The first step is thus to define a topological space.

### **Definition (Topological space)**

A **topological space**  $(X, \mathcal{O})$  is a set  $X$  equipped with a collection  $\mathcal{O} \subseteq \mathcal{P}(X)$  of its subsets, which we refer to as **open sets** forming the **topology**  $\mathcal{O}$  of  $X$ , such that the following conditions are satisfied:

- (i) the empty set  $\emptyset$  and  $X$  are open sets
- (ii) if  $A, B \in \mathcal{O}$  then  $A \cap B$  is also open
- (iii) if  $A_i \in \mathcal{O}$  for some  $\{i\}$  then  $\bigcup_i A_i$  is open.

### **Definition (Closed set)**

Let  $(X, \mathcal{O})$  be a topological space. Then  $A \subseteq X$  is closed (relative to the topology  $\mathcal{O}$ ) if  $X \setminus A \in \mathcal{O}$ .

There are trivial examples of topologies. Let  $X$  be any set, then  $\mathcal{O} = \{\emptyset, X\}$  is a topology known as the **chaotic topology**. Also,  $\mathcal{O} = \mathcal{P}(X)$  is a topology known as the **discrete topology**.

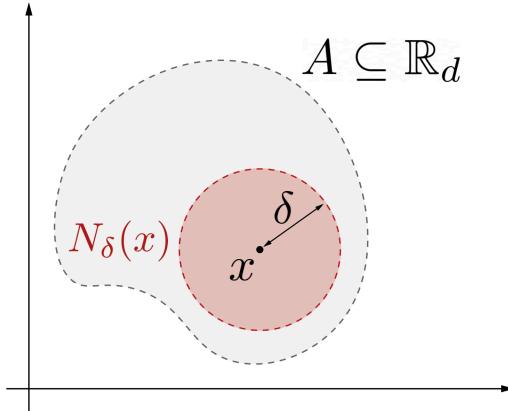
A very relevant example is the standard topology of the  $d$ -dimensional Euclidean space  $\mathbb{R}^d$ . Before defining it let us discuss what it means to be in the neighborhood of a point  $x$ .

**Definition (Neighborhood)**

Let us define for all  $x \in \mathbb{R}^d$  and  $d \in \mathbb{R}^+$ :

$$N_d(x) = \{y \in \mathbb{R}^d : \|y - x\| < d\} \quad (54.2.1)$$

which is an open (since the inequality is strict) ball centered at  $x$  with radius  $d$ . This is the **neighborhood** of  $x$  of radius  $d$ .

**Theorem (Standard topology)**

The following is a topology on  $\mathbb{R}^d$  known as the **standard topology**  $\mathcal{O}_{\mathbb{R}^d}$ :

$$A \in \mathcal{O}_{\mathbb{R}^d} \iff \forall x \in A, \exists \delta \in \mathbb{R}^+ \text{ s.t. } N_\delta(x) \subseteq A \quad (54.2.2)$$

that is, an open set is a member of this topology if we can always form a neighborhood of  $x$  contained within  $A$ .

*Proof.* We check that the three topological axioms are satisfied:

- (i) Firstly  $\emptyset \in \mathcal{O}_{\mathbb{R}^d}$  since the neighborhood of the empty set is empty for any  $r \in \mathbb{R}^+$  and  $\emptyset \subseteq \emptyset$ . Similarly  $\mathbb{R}^d \in \mathcal{O}_{\mathbb{R}^d}$  since the neighborhood of any arbitrary point for any  $r \in \mathbb{R}^+$  must by definition lie in  $\mathbb{R}^d$
- (iii) Let  $A, B \in \mathcal{O}_{\mathbb{R}^d}$  be open. This implies that  $\exists r_A, r_B \in \mathbb{R}^+$  such that  $N_{r_A}(x) \subseteq A$  and  $N_{r_B}(x) \subseteq B$ . Letting  $r = \min\{r_A, r_B\}$  then  $N_r(x) \subseteq A \cap B$  implying that  $N_r(x) \subseteq U \cup V$  as desired.
- (iv) Let  $A_i \in \mathcal{O}$  for given  $\{i\}$ . Then for any  $j$  there exists  $r_j \in \mathbb{R}^+$  such that  $N_{r_j}(x) \subseteq A_j$  implying that  $N_{r_j}(x) \subseteq \bigcup_i A_i$  as desired.

■

**Theorem (Induced topology)**

Let  $(X, \mathcal{O})$  be a topological space, and let  $B \subseteq X$ . Then:

$$\mathcal{O}|_B = \{A \cap B | A \in \mathcal{O}\} \quad (54.2.3)$$

is a topology on  $b$  known as the **induced topology**.

*Proof.* We verify the topology axioms:

- (i) since  $\emptyset \in \mathcal{O}$  we have that  $\emptyset \cap B = \emptyset$  is also an open set. Similarly since  $X \in \mathcal{O}$  we have that  $X \cap B = B$  is open as desired.
- (ii) let  $U, V \in \mathcal{O}|_N$ . Then  $\exists S, T \in \mathcal{O}$  such that  $U = S \cap B$  and  $V = T \cap B$ , implying that:

$$U \cap V = (S \cap B) \cap (T \cap B) = \underbrace{(S \cap T)}_{\in \mathcal{O}} \cap B \quad (54.2.4)$$

- (ii) let  $U_i \in \mathcal{O}|_N$  for some  $\{i\}$ . Then  $\exists S_i \in \mathcal{O}$  such that  $U_i = S_i \cap B$  so that

$$\bigcup_i U_i = \bigcup_i (S_i \cap B) = \overbrace{\left( \bigcup_i S_i \right)}^{\in \mathcal{O}} \cup B \quad (54.2.5)$$

■

An interesting example is the set  $\mathcal{C} = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| = 1\}$ . We can establish a topology on  $\mathcal{C}$  by inducing the standard topology on  $\mathcal{S}$ :

$$\mathcal{O} = \mathcal{O}_{\mathbb{R}^d}|_{\mathcal{S}} \quad (54.2.6)$$

This is because the empty set is the intersection of the  $\emptyset \in \mathcal{O}_{\mathbb{R}^d}$  and  $\mathcal{S}$ , similarly  $\mathcal{S}$  is the intersection of  $\mathbb{R}^d \in \mathcal{O}_{\mathbb{R}^d}$ . Also, given any subset of  $\mathcal{S}$ , it can be written as the intersection of another circle and  $\mathcal{S}$  as shown below.

### Definition (*Product topology*)

Let  $(A, \mathcal{O}_A)$  and  $(B, \mathcal{O}_B)$ . Then  $\mathcal{O}_{A \times B}$  is a topology, the **product topology** on  $A \times B$  defined by:

$$U \in \mathcal{O}_{A \times B} \iff \forall x = (a, b) \in U, \exists S \in \mathcal{O}_A, T \in \mathcal{O}_B \text{ s.t. } S \times T \subseteq U \text{ and } a \in S, b \in T \quad (54.2.7)$$

*Proof.* We verify the topology axioms:

- (i) of course  $\emptyset \in \mathcal{O}_{A \times B}$  since  $\nexists x \in \emptyset$ , so any statement following this assumption is trivially satisfied. Similarly  $U = A \times B \in \mathcal{O}_{A \times B}$  since given  $x = (a, b) \in U$  then  $a \in A \in \mathcal{O}_A, b \in B \in \mathcal{O}_B$  with  $A \times B \subseteq U$  as desired.
- (ii) Let  $U, V \in \mathcal{O}_{A \times B}$ , and let  $S_{U,V} \in \mathcal{O}_A$  and  $T_{U,V} \in \mathcal{O}_B$  be the open sets satisfying the topological conditions:

$$\forall x = (a, b) \in U, \exists S_U \in \mathcal{O}_A, T_U \in \mathcal{O}_B \text{ s.t. } S_U \times T_U \subseteq U \text{ and } a \in S_U, b \in T_U \quad (54.2.8)$$

$$\forall x = (a, b) \in V, \exists S_V \in \mathcal{O}_A, T_V \in \mathcal{O}_B \text{ s.t. } S_V \times T_V \subseteq V \text{ and } a \in S_V, b \in T_V \quad (54.2.9)$$

Let  $x = (a, b) \in U \cap V$ , and consider  $S_U \cap S_V \in \mathcal{O}_A$  and  $T_U \cap T_V \in \mathcal{O}_B$ . Then clearly  $a \in S_U \cap S_V$  and  $b \in T_U \cap T_V$ , and moreover:

$$(S_U \cap S_V) \times (T_U \cap T_V) = (S_U \times T_U) \cap (S_V \times T_V) \subseteq U \times V \quad (54.2.10)$$

as desired.

(iii) Let  $U_i \in \mathcal{O}_{A \times B}$  for some set  $\{i\}$ . Then, we must have that:

$$\forall i : \forall x = (a, b) \in U_i, \exists S_i \in \mathcal{O}_A, T_i \in \mathcal{O}_B \text{ s.t. } S_i \times T_i \subseteq U_i \text{ and } a \in S_i, b \in T_i \quad (54.2.11)$$

Let  $x = (a, b) \in \bigcup_i U_i$ , and let us consider the unions  $\bigcup_i S_i$  and  $\bigcup_i T_i$ . Clearly  $a \in \bigcup_i S_i$  and  $b \in \bigcup_i T_i$ . Moreover:

$$\left( \bigcup_i S_i \right) \times \left( \bigcup_i T_i \right) = \bigcup_i (S_i \times T_i) \subseteq \bigcup_i U_i \quad (54.2.12)$$

as desired. ■

This definition can be easily extended to any number of Cartesian products. Also, you can always fit a circle in a square and vice-versa,  $S \times T$  and  $N_d(x)$  can be fit into each other. Therefore,  $\mathcal{O}_{\mathbb{R}^d} = \mathcal{O}_{\mathbb{R} \times \dots \times \mathbb{R}}$ .

### Definition (Convergence)

A sequence  $q : \mathbb{N} \rightarrow X$  on a topological space  $(X, \mathcal{O})$  converges to a point  $a \in X$  if:

$$\forall U \in \mathcal{O} \text{ with } a \in U, \exists N \in \mathbb{N} \text{ such that } \forall n > N, q(n) \in U \quad (54.2.13)$$

### Theorem (Sequence convergence)

The sequence  $q : \mathbb{N} \rightarrow \mathbb{R}$  converges to  $a \in \mathbb{R}^d$  if:

$$\forall \epsilon > 0, \exists N \in \mathbb{N} \text{ s.t. } \|q(n) - a\| < \epsilon, \forall n > N \quad (54.2.14)$$

*Proof.* Suppose the sequence  $q : \mathbb{N} \rightarrow \mathbb{R}^d$  with standard topology converges to  $a$ . Let  $\epsilon > 0$ , and define:

$$U_\epsilon = \{x \in \mathbb{R}^d : \|x - a\| < \epsilon\} \quad (54.2.15)$$

This, as we shall prove soon, is an open set. Also, we have that  $a \in U_\epsilon$ , so the definition of sequence convergence implies that:

$$\exists N \in \mathbb{N} \text{ such that } \|q(n) - a\| < \epsilon, \forall n > N \quad (54.2.16)$$

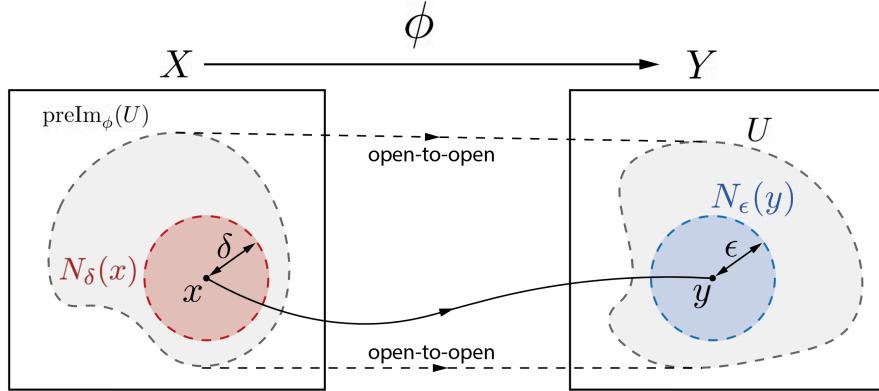
as desired. ■

### Definition (Continuity)

Let  $(X, \mathcal{O}_X)$  and  $(Y, \mathcal{O}_Y)$  be topological spaces. Then the map  $\phi : X \rightarrow Y$  is **continuous** if for all open sets  $U \in \mathcal{O}_Y$  the set  $\text{preIm}_\phi U \in \mathcal{O}_X$ .

In other words, a continuous function maps open sets to open sets without disrupting the topologies of the domain and image. Consider for example  $\phi : X \rightarrow Y$  where  $X$  is equipped with the topology  $\mathcal{P}(X)$  while  $Y$  is equipped with another topology  $\mathcal{O}_Y$ . Then  $\phi$  must be continuous since the pre-image of any  $V \in \mathcal{O}_Y$  is necessarily a subset of the domain  $X$ , and must therefore be included in the power set  $\mathcal{P}(X)$ . Similarly  $\phi : X \rightarrow Y$  where  $X$  is equipped with any topology and  $Y$  is equipped with the chaotic topology  $\{\emptyset, Y\}$  is also continuous.

We can visualize the definition of a continuous function as a map that preserved the “open-ness” of subsets. It is good sanity check to show that this topological definition of continuity yields the



epsilon-delta definition of continuity we are familiar with from Real analysis.

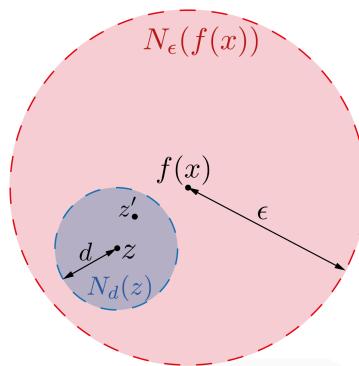
### Theorem (Continuity of $\mathbb{R}^n$ maps)

Let  $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$  be a map, where  $\mathbb{R}^m$  and  $\mathbb{R}^n$  are endowed with their standard topologies. The map  $f$  is continuous according to the topological definition iff it is continuous according to the epsilon-delta definition.

*Proof.*  $\implies$  Suppose that the map  $f$  is topologically continuous. First, we claim that, given  $\epsilon > 0$ , the neighbourhood  $N_\epsilon(f(x))$  is open

$$N_\epsilon(f(x)) = \{z \in \mathbb{R}^n : \|f(x) - z\| < \epsilon\} \quad (54.2.17)$$

Indeed given  $z \in N_\epsilon(f(x))$  then we see that choosing  $0 < d + \|z - f(x)\| < \epsilon$



and letting  $z' \in N_d(z)$  so that  $\|z' - z\| < d$  then:

$$\|z' - f(x)\| < \|z' - z\| + \|z - f(x)\| < d + \|z - f(x)\| < \epsilon \quad (54.2.18)$$

implying that

$$N_d(z) = \{z' \in \mathbb{R}^n : \|z' - z\| < d\} \subseteq N_\epsilon(f(x)) \quad (54.2.19)$$

as desired. Now since  $N_\epsilon(f(x))$  is an open set, its preimage  $\text{preIm}_f(N_\epsilon(f(x)))$  must also be open, and since  $x$  is in this preimage (because  $f(x) \in N_\epsilon(f(x))$ ) we have that  $\exists \delta > 0$  such that

$$N_\delta(x) = \{y \in \mathbb{R}^m : \|x - y\| < \delta\} \subseteq \text{preIm}_f(N_\epsilon(f(x))) \implies f(N_\delta(x)) \subseteq N_\epsilon(f(x)) \quad (54.2.20)$$

that is, if  $\|x - y\| < \delta$  then  $\|f(y) - f(x)\| < \epsilon$  as desired.

$\Leftarrow$  Suppose that the map  $f$  is continuous according to the epsilon-delta definition, so that if  $x \in \mathbb{R}^m$ , then for all  $\epsilon > 0$  there exists  $\delta > 0$  such that  $\|x - y\| < \delta \implies \|f(x) - f(y)\| < \epsilon$ . Consequently:

$$N_\delta(x) = \{y \in \mathbb{R}^m : \|x - y\| < \delta\} \subseteq \{y \in \mathbb{R}^m : \|f(x) - f(y)\| < \epsilon\} \quad (54.2.21)$$

Let  $A \in \mathbb{R}^n$  be an open set in the standard topology, and let  $z \in A$ . Then there exists  $\epsilon > 0$  such that  $N_\epsilon(z) \subseteq A$ . Consider  $\text{preIm}_f(A)$ , we claim this is open. Indeed let  $x \in \text{preIm}_f(A)$ . Then by the definition of  $\epsilon - \delta$  continuity:

$$f(N_\delta(x)) = \{f(y) : \|x - y\| < \delta\} \subseteq \{y : \|f(x) - f(y)\| < \epsilon\} \subseteq \{z : \|f(x) - z\| < \epsilon\} = N_\epsilon(z) \subseteq A \quad (54.2.22)$$

implying that  $N_\delta(x) \subseteq \text{preIm}_f(A)$ , as desired. So if  $A$  is open then  $\text{preIm}_f(A)$  is also open, as desired.  $\blacksquare$

### Definition (*Homeomorphism*)

Let  $\phi : X \rightarrow Y$  be a bijective map, with  $(X, \mathcal{O}_X)$  and  $(Y, \mathcal{O}_Y)$  being topological spaces. Then  $\phi$  is a **homeomorphism** if

- (a)  $\phi : X \rightarrow Y$  is continuous
- (b)  $\phi^{-1} : Y \rightarrow X$  is continuous If such a  $\phi$  exists then  $(X, \mathcal{O}_X) \cong (Y, \mathcal{O}_Y)$ , they are homeomorphic.

## 54.3 Topological manifolds

Now that we have defined a topological space we can add additional structure, namely that of a manifold which, as hinted earlier, is a topological space that locally looks Euclidean.

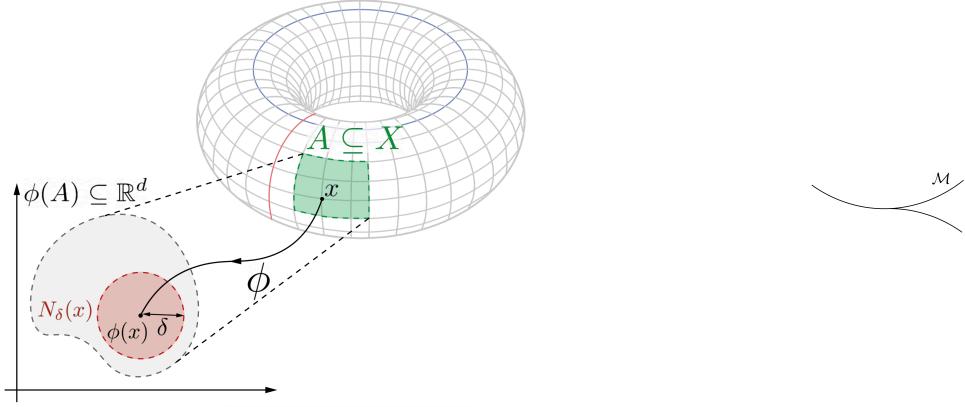
### Definition (*Topological manifold*)

We define a  $d$ -dimensional **topological manifold** if it is a topological space  $(X, \mathcal{O})$  if for all  $x \in X$ , there exists an open set  $A \subseteq X$  with  $x \in A$  such that there exists a map  $\exists \phi : A \rightarrow \phi(A) \in \mathcal{O}_{\mathbb{R}^d}$  satisfying:

- (i)  $\phi$  is invertible
- (ii)  $\phi$  is continuous
- (iii)  $\phi^{-1}$  is continuous

where we choose the topology  $\mathcal{O}$  on  $A$  and  $\mathcal{O}_{\mathbb{R}^d}$  on  $\mathbb{R}^d$  (i.e.  $A$  is homeomorphic to  $\mathbb{R}^d$ ). The choice  $(A, \phi)$  is defined as a **chart**.

What this means is that all points  $x$  in the topological manifold  $(X, \mathcal{O})$  lie in some open set  $A$  that is homeomorphic to  $\mathbb{R}^d$ . Consider as an example a wire  $M \subset \mathbb{R}^2$  that bifurcates into two branches. We



can embed the topology  $\mathcal{O}_{\mathbb{R}^2}|_M$  on  $M$  creating a topological space. This however is not a topological manifold, since for the bifurcation point we cannot find an open set it lies in that can be mapped continuously bijectively to  $\mathbb{R}^2$ .

### Definition (Atlas)

The set of charts  $\mathcal{A} = \{(A_i, \phi_i : i \in S\}$  is an **atlas** of the topological space  $(X, \mathcal{O})$  provided that:

$$X = \bigcup_{i \in S} A_i \quad (54.3.1)$$

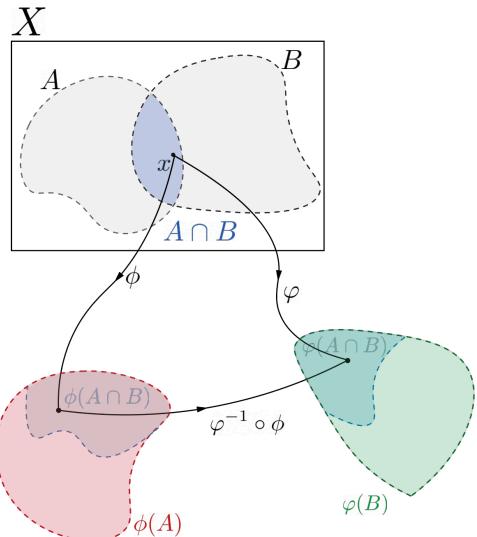
### Definition (Chart components)

For a given chart  $(A, \phi)$  of an  $n$ -manifold we define its **chart components**  $\phi^\mu$  as the maps:

$$\phi^\mu : A \rightarrow \mathbb{R}^n \quad (54.3.2)$$

$$p \mapsto [\phi(p)]^\mu \quad (54.3.3)$$

which map  $p \in A$  to the  $\mu$ th component of  $\phi(p) \in \mathbb{R}^n$ .



Consider two charts  $(A, \phi)$  and  $(B, \varphi)$  with  $A \cap B \neq \emptyset$ . Then this implies that given some  $x \in A \cap B$ , it can be mapped to  $\mathbb{R}^d$  using two charts,  $\phi$  or  $\varphi$ . Can we find a way to map between  $\phi(x)$  and  $\varphi(x)$ ? Well it is easy to see from the commutative diagram that  $\psi \equiv \varphi \circ \phi^{-1}$  does precisely this. We call  $\psi$  the chart transition map, it is akin to a change of basis map.

The reason topological manifolds are so useful in physics is that it allows us to describe objects (for ex-

ample a curve  $\gamma : \mathbb{R} \rightarrow X$ ) in the real world by looking at how they may be charted (e.g. consider the charts  $(\phi, X)$  and  $(\varphi, X)$ ). The real physical object is  $\gamma$  but physically we describe them using  $\phi \circ \gamma$  or  $\varphi \circ \gamma$ . We must however make sure that whatever properties we are looking, it does not matter what charts we are looking at. The advantage of this approach is that we have all the tools of multivariable analysis in  $\mathbb{R}^d$  at our disposal. The disadvantage is that we are blurring the lines between the real world and the charts we are studying them in by replacing  $\gamma$ , the real world object, with its mathematical description  $\phi \circ \gamma$  within the chart  $(X, \phi)$ .

For example suppose  $\phi \circ \gamma$  is continuous, how do we know that  $\varphi \circ \gamma$  is also continuous? From the commutative diagram we see that  $\varphi \circ \gamma = (\varphi \circ \phi^{-1}) \circ (\phi \circ \gamma)$  so  $\varphi \circ \gamma$  is indeed continuous.

$$\begin{array}{ccc}
 & \varphi(X) \subseteq \mathbb{R}^d & \\
 & \nearrow \varphi \circ \gamma \quad \uparrow \varphi & \\
 \mathbb{R} & \xrightarrow{\gamma} & X \\
 & \searrow \phi \circ \gamma \quad \downarrow \phi & \\
 & & \phi(X) \subseteq \mathbb{R}^d
 \end{array}$$

Can we say the same about differentiability? If we have that  $\phi \circ \gamma$  is differentiable, we still do not know if this differentiability is shared in all charts since  $\varphi \circ \phi^{-1}$  is not necessarily differentiable (they are continuous but can have kinks). We need to find some way to examine differentiability using charts.

## 54.4 Multilinear algebra

Recall the normal definition of vector spaces.

### Definition (Vector space)

A vector space  $(V, +, \cdot)$  over a field  $\mathbb{F}$  is a set  $V$  and two maps:

$$+ : V \times V \rightarrow V \tag{54.4.1}$$

$$\cdot : \mathbb{F} \times V \rightarrow V \tag{54.4.2}$$

known as addition and scalar multiplication such that for all  $\mathbf{u}, \mathbf{v}, \mathbf{w} \in V, \alpha, \beta \in \mathbb{F}$ :

- (i)  $\mathbf{v} + \mathbf{w} = \mathbf{w} + \mathbf{v}$
- (ii)  $(\mathbf{u} + \mathbf{v}) + \mathbf{w} = \mathbf{u} + (\mathbf{v} + \mathbf{w})$
- (iii)  $\exists 0 \in V$  such that  $\mathbf{v} + 0 = \mathbf{v}$
- (iv)  $\exists (-\mathbf{v}) \in V$  such that  $\mathbf{v} + (-\mathbf{v}) = 0$
- (v)  $\alpha \cdot (\beta \cdot \mathbf{v}) = (\alpha\beta) \cdot \mathbf{v}$
- (vi)  $\alpha \cdot (\mathbf{v} + \mathbf{w}) = \alpha \cdot \mathbf{v} + \alpha \cdot \mathbf{w}$
- (vii)  $(\alpha + \beta) \cdot \mathbf{v} = \alpha \cdot \mathbf{v} + \beta \cdot \mathbf{v}$
- (viii) given the identity element  $1$  of  $\mathbb{F}$ , then  $1 \cdot \mathbf{v} = \mathbf{v}$

The element of a vector space is informally referred to as a **vector**. It is easy to see that the set  $\mathcal{P}_n$  of all polynomial functions  $p$  up to order  $n \in \mathbb{N}$  is indeed a vector space:

$$\mathcal{P}_n = \{p(x) = \sum_{m=0}^n p_m x^m : p_m \in \mathbb{R}\} \quad (54.4.3)$$

### Definition (Linear map)

Let  $(V, +_V, \cdot_V)$  and  $(W, +_W, \cdot_W)$  be two vector spaces. Then the map  $\phi : V \rightarrow W$  is a **linear map** if it is structure preserving:

- (i)  $\phi(\mathbf{u} +_V \mathbf{v}) = \phi(\mathbf{u}) +_W \phi(\mathbf{v})$
- (ii)  $\phi(\alpha \cdot_V \mathbf{v}) = \alpha \cdot_W \phi(\mathbf{v})$

If such a map between  $V$  and  $W$  exists then we write that  $V \cong W$ .

### Definition (Set of linear maps)

We define the set of all linear maps between two vector spaces  $(V, +_V, \cdot_V)$  and  $(W, +_W, \cdot_W)$  (the latter defined over  $\mathbb{F}_W$ ) as  $\text{Hom}(V, W)$ . Together with the operations  $\oplus$  and  $\odot$  defined below this set becomes a vector space:

$$\oplus : \text{Hom}(V, W) \times \text{Hom}(V, W) \rightarrow \text{Hom}(V, W) \quad (54.4.4)$$

$$(\phi, \varphi) \mapsto \phi \oplus \varphi \quad (54.4.5)$$

and

$$\odot : \mathbb{F}_W \times \text{Hom}(V, W) \rightarrow \text{Hom}(V, W) \quad (54.4.6)$$

$$(\alpha, \varphi) \mapsto \alpha \odot \varphi \quad (54.4.7)$$

where

$$(\phi \oplus \varphi)(\mathbf{v}) = \phi(\mathbf{v}) +_W \varphi(\mathbf{v}), \forall \mathbf{v} \in V \quad (54.4.8)$$

$$(\alpha \odot \varphi)(\mathbf{v}) = \alpha \cdot_W \phi(\mathbf{v}), \forall \mathbf{v} \in V, \forall \alpha \in \mathbb{F}_W \quad (54.4.9)$$

### Definition (Dual vector space)

Given a vector space  $V$  defined over a field  $\mathbb{F}$ , then the set  $V^*$  is defined as:

$$V^* \equiv \{\phi : V \rightarrow \mathbb{F} : \phi \text{ is linear}\} = \text{Hom}(V, \mathbb{F}) \quad (54.4.10)$$

Equipped with  $\oplus$  and  $\odot$  defined previously then  $V^*$  is defined as the **dual vector space**.

Informally, the element  $\phi \in V^*$  is referred to as a **covector**, for reasons that shall be clearer soon.

### Definition (Dual basis)

Given a basis  $\mathcal{B} = \{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\} \subset V$  for  $V$  then the **dual basis** of the dual space  $V^*$  is defined as  $\mathcal{B}^* = \{\epsilon^1, \epsilon^2, \dots, \epsilon^n\} \subset V^*$  such that:

$$\epsilon^i(\mathbf{e}_j) = \delta_i^j, \forall i, j \quad (54.4.11)$$

For the case of  $\mathcal{P}_n$ , we can choose the basis  $\mathcal{B} = \{e_i(x) = x^i : i = 0, 1, \dots, n\}$ . The dual basis  $\mathcal{B}^* = \{\epsilon_i : i = 0, 1, \dots, n\}$  can be easily identified as:

$$\epsilon^i = \frac{1}{i!} \left. \frac{d^i}{dx^i} \right|_{x=0} \quad (54.4.12)$$

Indeed one finds that

$$\epsilon^i(e_j) = \left. \frac{1}{i!} \frac{d^i}{dx^i}(x^j) \right|_{x=0} = \begin{cases} 0, & i > j \\ 1, & i = j \\ 0, & i < j \end{cases} = \delta_{ij} \quad (54.4.13)$$

so the derivative operator may be regarded as a covector in the dual space of  $P_n$ . This will turn out to be a more general property of differential geometry.

### Definition (Tensor)

Let  $(V, +, \bullet)$  be a vector space defined over a field  $\mathbb{F}$ . Then an  $(r, s)$ -tensor  $T$  over  $V$  is a map:

$$T : (V^*)^r \times V^s \mapsto \mathbb{F} \quad (54.4.14)$$

that is linear in all its entries, thus a **multi-linear map**.

Consider for example a  $(1, 1)$ -tensor  $T$ . Then we must have that:

$$T(\phi + \varphi, \mathbf{v}) = T(\phi, \mathbf{v}) + T(\varphi, \mathbf{v}) \quad (54.4.15)$$

$$T(\alpha\phi, \mathbf{v}) = \alpha T(\phi, \mathbf{v}) \quad (54.4.16)$$

$$T(\phi, \mathbf{v} + \mathbf{w}) = T(\phi, \mathbf{v}) + T(\phi, \mathbf{w}) \quad (54.4.17)$$

$$T(\phi, \alpha\mathbf{v}) = \alpha T(\phi, \mathbf{v}) \quad (54.4.18)$$

$$(54.4.19)$$

We can connect this to the typical view of a  $(1, 1)$ -tensor as a map that takes a vector and maps it to another vector. Indeed let us construct:

$$\phi_T : V \rightarrow V \quad (54.4.20)$$

$$\mathbf{v} \mapsto T(\cdot, \mathbf{v}) \quad (54.4.21)$$

which takes in  $\mathbf{v}$  and maps it to another map  $T(\cdot, \mathbf{v})$  which takes covectors to  $\mathbb{F}$ . But this map is of the type  $V^* \rightarrow \mathbb{F}$  and thus  $T(\cdot, \mathbf{v}) \in (V^*)^* = V$  provided  $\dim V < \infty$ . It follows that  $T(\cdot, \mathbf{v})$  is a vector, so  $\phi_T$  does indeed map a vector to vectors of  $V$ .

### Theorem (Vectors vs Covectors)

Let  $\phi \in V^*$  be a covector. Then by definition it must be a  $(0, 1)$ -tensor since  $\phi : V \rightarrow \mathbb{F}$  is linear.

Similarly, let  $v \in V = (V^*)^*$  be a vector. Then by definition it must be a  $(1, 0)$ -tensor since  $\mathbf{v} : V^* \rightarrow \mathbb{F}$ .

### Definition (Tensor components)

Let  $T$  be an  $(r, s)$ -tensor over a  $n$ -dimensional vector space  $V$  with basis  $\mathcal{B} = \{\mathbf{e}_i : i =$

$0, 1, \dots, n\}$  and dual space  $V^*$  with dual basis  $\mathcal{B}^* = \{\epsilon_i : i = 0, 1, \dots, n\}$ . Then we define the components of  $T$  as:

$$T_{j_1 \dots j_s}^{i_1 \dots i_r} \equiv T(\epsilon^{i_1}, \dots, \epsilon^{i_r}, \mathbf{e}_{j_1}, \dots, \mathbf{e}_{j_s}) \quad (54.4.22)$$

so that:

$$T(\phi, \mathbf{v}) = \sum_{ij} \phi_i v^j T_j^i \quad (54.4.23)$$

where  $\phi = \sum_i \phi_i \epsilon^i \in V^*$  and  $\mathbf{v} = \sum_j v^j \mathbf{e}_j \in V$ .

For example, given a  $(1, 1)$ -tensor  $T$  then  $T_j^i = T(\epsilon^i, \mathbf{e}_j)$ . Also, due to the linearity of the tensor map the tensor components are very useful because they allow us to expand:

$$T(\phi, \mathbf{v}) = T\left(\sum_i \phi_i \epsilon^i, \sum_j v^j \mathbf{e}_j\right) \quad (54.4.24)$$

$$= \sum_{ij} \phi_i v^j T(\epsilon^i, \mathbf{e}_j) \quad (54.4.25)$$

$$= \sum_{ij} \phi_i v^j T_j^i \quad (54.4.26)$$

Note the careful placements of subscripts and superscripts, where basis vectors are given subscripts and basis covectors are given superscripts. This ensures that multiply labelled indices always appear in an up-down arrangement. Moreover, if we assume that an index appearing both up and down is summed over, then we retrieve the **Einstein summation convention**. For example, (54.4.23) becomes

$$T(\phi, \mathbf{v}) = \phi_i v^j T_j^i \quad (54.4.27)$$

# Differentiable manifolds

## 55.1 Differenti

We are now ready to tackle the problem mentioned in section 54.3 about the consistency in studying the differentiability of topological spaces. The key lies in choosing the charts on the topological space from a restricted set (which is still an atlas) so that all transition functions are differentiable. Manifolds on which this can be done are known as **differentiable manifolds**.

### **Definition (Compatibility)**

Let  $\square$  be some “property” (e.g. differentiability) of a topological space. Two charts  $(A, \phi_A)$  and  $(B, \phi_B)$  are  $\square$ -**compatible** if:

- (i)  $A \cap B = \emptyset$  or
- (ii)  $\phi_B \circ \phi_A^{-1}$  and  $\phi_A \circ \phi_B^{-1}$  restricted to the domain  $A \cap B$  are  $\square$ .

Moreover, an atlas  $\mathcal{A}_\square$  is an  $\square$ -compatible atlas if given any two charts in  $\mathcal{A}$ , they are  $\square$ -compatible.

If this happens then we can define a corresponding type of  $\square$ -manifold  $(X, \mathcal{O}, \mathcal{A}_\square)$ .

It can be shown that any  $C^k$ -atlas  $\mathcal{A}_{C^k}$  contains a  $C^\infty$ -atlas. This is a substantial result because for physical applications we can always consider  $C^\infty$ -manifolds (smooth manifolds) without loss of generality.

### **Definition (Topological isomorphism)**

Two topological spaces  $(A, \mathcal{O}_A)$  and  $(B, \mathcal{O}_B)$  are **topologically isomorphic** if there exists a bijective homomorphism between them  $\phi : A \rightarrow B$ .

### **Definition (Smooth map)**

A map  $\varphi : M \rightarrow N$  for two smooth manifolds  $(M, \mathcal{O}_M)$  and  $(N, \mathcal{O}_N)$  is said to be **smooth** if for any two charts  $(X, \phi_X)$  of  $(M, \mathcal{O}_M)$  and  $(Y, \phi_Y)$  of  $(N, \mathcal{O}_N)$ , the map  $\phi_Y \circ \varphi \circ \phi_X^{-1}$  is smooth.

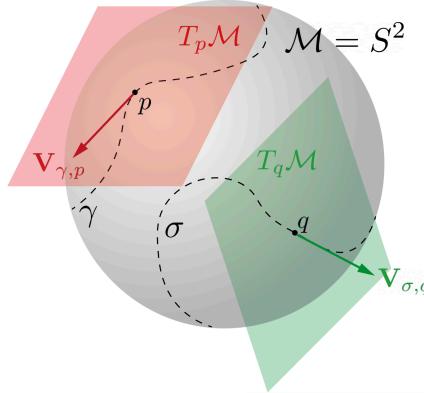
### **Definition (Diffeomorphism)**

Two smooth manifolds  $(M, \mathcal{O}_M, \mathcal{A}_M)$  and  $(N, \mathcal{O}_N, \mathcal{A}_N)$  are **diffeomorphic** if there exists a bijection  $\phi : M \rightarrow N$  such that both  $\phi$  and  $\phi^{-1}$  are smooth.

**Definition (Smooth structure of topological manifolds)**

In  $d = 4$  dimensions, the number of smooth manifolds that can be made out of a given topological manifold (up to a diffeomorphism) is uncountably infinite.

We can now start talking about tangent vectors and spaces. Consider for example a sphere  $X = S^2$  embedded in  $\mathbb{R}^3$ . We know that given a point  $p$  on this sphere, we can define a plane tangent to it at  $p$ . Any vector on this plane would then be a tangent vector. This visualization however relies on the fact that the sphere can be embedded in a larger Euclidean space, something that cannot be generally done. Thus we must develop a more abstract concept of tangent vector that does not live outside of the manifold it is tangent to.



It turns out that the correct definition of a tangent vector is as a directional derivative.

**Definition (Tangent vector)**

Let  $(X, \mathcal{O}, \mathcal{A})$  be a smooth manifold, and let  $\gamma : \mathbb{R} \rightarrow X$  be a smooth parametrized curve, with  $\gamma(t_p) = p \in X$ . Then the **tangent vector**  $\mathbf{V}_{\gamma,p}$  of  $\gamma$  at the point  $p$  is the linear map:

$$\mathbf{V}_{\gamma,p} : C^\infty(X) \rightarrow \mathbb{R} \quad (55.1.1)$$

$$f \mapsto (f \circ \gamma)'(t_p) \quad (55.1.2)$$

where  $C^\infty(X) = \{f : X \rightarrow \mathbb{R} | f \text{ is smooth}\}$  and ' represents differentiation with respect to the parametrization parameter  $t$ .

In other words, suppose  $f$  is some scalar field, for example temperature.  $f \circ \gamma$  gives the temperature as we go along the path  $\gamma$ , and thus the velocity is defined as the rate of change of this temperature field as we run along  $\gamma$ . This is shown well in the following commutative diagram:

$$\begin{array}{ccccc} \mathbb{R} & \xrightarrow{\gamma} & X & \xrightarrow{\phi} & \mathbb{R}^d \\ & \searrow f \circ \gamma & \downarrow f & \swarrow f \circ \phi^{-1} & \\ & & \mathbb{R} & & \end{array}$$

Differentiating  $f$  itself doesn't make a lot of sense since it is not a typical type of map that is encountered in analysis. We resolved this issue by looking at the parametrization of  $f$ , which is a "normal" map (i.e. from  $\mathbb{R}$  to  $\mathbb{R}$ ) and can thus be differentiated in the usual way.

We can write this result using the less rigorous, but more familiar differential notation

$$\mathbf{V}_{\gamma,p}(f) = \frac{d}{dt} f(\gamma(t)) \Big|_{t=t_p} \quad (55.1.3)$$

### Definition (*Tangent space*)

For each point  $p \in X$  we define the **tangent space to  $X$  at  $p$**  as the set of tangent vectors at  $p$  for all smooth curves  $\gamma$  on  $X$ .

$$T_p X = \{\mathbf{V}_{\gamma,p} : \gamma \text{ is a smooth curve}\} \quad (55.1.4)$$

Note that the view that the tangent space is a plane works fine for classical geometry, but it requires the concept of an ambient space surrounding the topological space  $X$ . If for example we're considering the universe, then it doesn't make sense to refer to a tangent plane outside of it. The strength of this new definition of tangent space is that it makes no reference to things outside of  $X$ .

It is interesting to look at the structure of this tangent space, is it a vector space for example? We can equip the tangent space  $T_p X$  with the following operators:

$$\oplus : T_p X \times T_p X \rightarrow \text{Hom}(C^\infty(X), \mathbb{R}) \quad (55.1.5)$$

$$(\mathbf{V}_{\gamma_1,p} \oplus \mathbf{V}_{\gamma_2,p})(f) \mapsto \mathbf{V}_{\gamma_1,p}(f) +_{\mathbb{R}} \mathbf{V}_{\gamma_2,p}(f) \quad (55.1.6)$$

and

$$\odot : \mathbb{R} \times T_p X \rightarrow \text{Hom}(C^\infty(X), \mathbb{R}) \quad (55.1.7)$$

$$(\alpha \odot \mathbf{V}_{\gamma,p})(f) \mapsto \alpha \cdot_{\mathbb{R}} \mathbf{V}_{\gamma,p}(f), \forall \alpha \in \mathbb{R} \quad (55.1.8)$$

with which the vector space axioms are satisfied. It is yet not understood however if the closure relation is satisfied, that is if  $\alpha \cdot_{\mathbb{R}} \mathbf{V}_{\gamma,p}(f)$  and  $\mathbf{V}_{\gamma_1,p}(f) +_{\mathbb{R}} \mathbf{V}_{\gamma_2,p}(f)$  belong to the tangent space themselves, and are thus tangent vectors to some smooth curve on  $X$ .

Thus, suppose there is some curve  $\sigma$  so that  $\mathbf{V}_{\sigma,p} = \alpha \cdot_{\mathbb{R}} \mathbf{V}_{\gamma,p}(f)$ . Intuitively we should expect (from now on any  $+$ ,  $\cdot$  operation is implicitly defined on the field  $\mathbb{R}$ ):

$$\sigma : \mathbb{R} \rightarrow X \quad (55.1.9)$$

$$t \mapsto \gamma(at + t_p) \quad (55.1.10)$$

to do the job. Indeed  $\sigma(0) = \gamma(t_p) = p$  so this curve passes through the required point. Also, defining  $\mu_\alpha : \lambda \rightarrow \alpha\lambda + \lambda_0$  then (we do the calculation two ways, first using the rigorous map notation and then using the familiar calculus notation):

$$\mathbf{V}_{\sigma,p} = (f \circ \sigma)'(0) = (f \circ \gamma \circ \mu_\alpha)'(0) = \alpha \cdot (f \circ \gamma)'(\lambda_0) = \alpha \cdot v_{\gamma,p} \quad (55.1.11)$$

$$\mathbf{V}_{\sigma,p} = \frac{d}{dt} f(\sigma(t)) \Big|_{t=0} = \frac{d}{dt} f(\gamma(at + t_p)) \Big|_{t=0} = \alpha \cdot \frac{d}{dt} f(\gamma(t)) \Big|_{t=t_p} = \alpha \cdot \mathbf{V}_{\gamma,p} \quad (55.1.12)$$

The proof for  $\oplus$  is slightly more involved, but the calculation will be used heavily later. We choose a chart  $(A, \phi)$  where  $p \in A$ , keeping in mind that if we make reference to a specific property of this chart that is not compatible with the chosen smooth atlas then the proof will be faulty. We define

the curve:

$$\sigma : \mathbb{R} \rightarrow X \quad (55.1.13)$$

$$t \mapsto \phi^{-1}((\phi \circ \gamma_1)(t + t_p) + (\phi \circ \gamma_2)(t + t'_p) - (\phi \circ \gamma)(t_p)) \quad (55.1.14)$$

where  $\gamma_1(t_p) = \gamma_2(t'_p) = p$ . We summarize the situation in the situation in the commutative diagrams below :

$$\begin{array}{ccc} \mathbb{R} & \xrightarrow{\sigma} & A & \xrightarrow{f} & \mathbb{R} \\ & \searrow \phi \circ \sigma & \downarrow \phi & \nearrow f \circ \phi^{-1} & \\ & & \mathbb{R}^d & & \end{array} \qquad \begin{array}{ccccc} t & \xrightarrow{\sigma} & \sigma(t) & \xrightarrow{f} & f(\sigma(t)) \\ \searrow \phi \circ \sigma & & \downarrow \phi & & \nearrow f \circ \phi^{-1} \\ & & \mathbf{x}(t) & & \end{array}$$

where we defined  $\mathbf{x}(t) = (\phi \circ \gamma)(t) \in \mathbb{R}^d$ . Firstly, note that  $\sigma(0) = \phi^{-1}(\phi(p) + \phi(p) - \phi(p)) = p$  so this curve does indeed contain  $p$ . Moreover:

$$\mathbf{V}_{\sigma,p}(f) = (f \circ \sigma)'(0) = ((f \circ \phi^{-1}) \circ (\phi \circ \sigma))'(0) \quad (55.1.15)$$

$$\mathbf{V}_{\sigma,p}(f) = \left. \frac{d}{dt} f(\sigma(t)) \right|_{t=0} = \left. \frac{d}{dt} (f \circ \phi^{-1})(\phi(\sigma(t))) \right|_{t=0} \quad (55.1.16)$$

We recognize  $\mathbf{x}(t) = (\phi \circ \sigma)(t)$ . Note that  $(f \circ \phi^{-1})$  is a map from  $\mathbb{R}^d$  to  $\mathbb{R}$  so we must use the multi-dimensional chain rule:

$$\mathbf{V}_{\sigma,p}(f) = (\phi \circ \sigma)^{\mu'}(0) \cdot (\partial_\mu(f \circ \phi^{-1}))(\mathbf{x}_p) \quad (55.1.17)$$

$$\mathbf{V}_{\sigma,p}(f) = \left. \frac{d}{dt} \phi^\mu(\sigma(t)) \right|_{t=0} \cdot \left. \frac{\partial}{\partial x^\mu} (f \circ \phi^{-1})(\mathbf{x}) \right|_{\mathbf{x}=\mathbf{x}_p} \quad (55.1.18)$$

We substitute (55.1.13) into the first term in the above:

$$(\phi \circ \sigma)^{\mu'}(0) = (\phi \circ \gamma_1)^{\mu'}(t_p) + (\phi \circ \gamma_2)^{\mu'}(t'_p) \quad (55.1.19)$$

$$\left. \frac{d}{dt} \phi^\mu(\sigma(t)) \right|_{t=0} = \left. \frac{d}{dt} \phi^\mu(\gamma_1(t)) \right|_{t=t_p} + \left. \frac{d}{dt} \phi^\mu(\gamma_2(t)) \right|_{t=t'_p} \quad (55.1.20)$$

so that:

$$\mathbf{V}_{\sigma,p}(f) = (\phi \circ \gamma_1)^{\mu'}(t_p) \cdot (\partial_\mu(f \circ \phi^{-1}))(\mathbf{x}_p) + (\phi \circ \gamma_2)^{\mu'}(t'_p) \cdot (\partial_\mu(f \circ \phi^{-1}))(\mathbf{x}_p) \quad (55.1.21)$$

$$\mathbf{V}_{\sigma,p}(f) = \left. \frac{d}{dt} \phi^\mu(\gamma_1(t)) \right|_{t=t_p} \cdot \left. \frac{\partial}{\partial x^\mu} (f \circ \phi^{-1})(\mathbf{x}) \right|_{\mathbf{x}=\mathbf{x}_p} + \left. \frac{d}{dt} \phi^\mu(\gamma_2(t)) \right|_{t=t'_p} \cdot \left. \frac{\partial}{\partial x^\mu} (f \circ \phi^{-1})(\mathbf{x}) \right|_{\mathbf{x}=\mathbf{x}_p} \quad (55.1.22)$$

Using the inverse chain rule this becomes:

$$\mathbf{V}_{\sigma,p}(f) = (f \circ \gamma)'(\lambda_0) + (f \circ \gamma')'(\lambda_1) = v_{\gamma,p}(f) + v_{\gamma',p}(f) \quad (55.1.23)$$

as desired. Note that we made no restriction on what chart to use, any chart would have worked as well as long as it was in a smooth atlas. To conclude, we have proven that:

**Theorem (The tangent space is a vector space)** The tangent space  $T_p X$  is a vector space when equipped with the operations:

$$\oplus : T_p X \times T_p X \rightarrow T_p X \quad (55.1.24)$$

$$(v_{\gamma,p} \oplus v_{\gamma',p})(f) \mapsto v_{\gamma,p}(f) +_{\mathbb{R}} v_{\gamma',p}(f) \quad (55.1.25)$$

and

$$\odot : \mathbb{R} \times T_p X \rightarrow T_p X \quad (55.1.26)$$

$$(\alpha \odot v_{\gamma,p})(f) \mapsto \alpha \cdot_{\mathbb{R}} v_{\gamma,p}(f), \forall \alpha \in \mathbb{R} \quad (55.1.27)$$

Let  $(A, \phi) \in \mathcal{A}_{C^\infty}$  be a chart in a smooth atlas. Let us define the curve:

$$\gamma : \mathbb{R} \rightarrow A \quad (55.1.28)$$

$$t \mapsto \gamma(t) \quad (55.1.29)$$

such that  $\gamma(t_p) = p \in A$  and let  $\mathbf{x}(t_p) = \phi(\gamma(t_p)) = \mathbf{x}_p$ . Then we have that by the same calculation as before:

$$\mathbf{V}_{\gamma,p}(f) = (f \circ \gamma)'(t_p) = ((f \circ \phi^{-1}) \circ (\phi \circ \gamma))'(t_p) \quad (55.1.30)$$

$$= (\phi \circ \gamma)^{\mu'}(t_p) \cdot (\partial_\mu(f \circ \phi^{-1}))(\mathbf{x}_p) \quad (55.1.31)$$

$$= \frac{d}{dt} \phi^\mu(\gamma(t)) \Big|_{t=t_p} \cdot \frac{\partial}{\partial x^\mu}(f \circ \phi^{-1})(\mathbf{x}) \Big|_{\mathbf{x}=\mathbf{x}_p} \quad (55.1.32)$$

We now define the partial derivative for  $f$  as the usual partial derivative of the chart representation of  $f$ :

$$\frac{\partial}{\partial x^\mu} \Big|_{t=t_p} f = \frac{\partial}{\partial x^\mu} (f \circ \gamma)(t) \Big|_{t=t_p} = \frac{\partial}{\partial x^\mu} (f \circ \phi^{-1})(\mathbf{x}) \Big|_{\mathbf{x}=\mathbf{x}_p} \quad (55.1.33)$$

giving:

$$\mathbf{V}_{\gamma,p}(f) = \frac{dx^\mu(t)}{dt} \Big|_{t=t_p} \cdot \frac{\partial}{\partial x^\mu} \Big|_{t=t_p} f \quad (55.1.34)$$

We may thus define the components of the velocity vector:

### Definition (Velocity components)

Let  $(A, x) \in \mathcal{A}$  be a chart a smooth atlas. Let us define the curve:

$$\gamma : \mathbb{R} \rightarrow A \quad (55.1.35)$$

$$t \mapsto \gamma(t) \quad (55.1.36)$$

Then the velocity vector of  $\gamma$  at  $p$  is:

$$\mathbf{V}_{\gamma,p} = \frac{dx^\mu(t)}{dt} \Big|_{t=t_p} \cdot \frac{\partial}{\partial x^\mu} \Big|_{t=t_p} \quad (55.1.37)$$

where  $\frac{dx^\mu}{dt} \Big|_{t=t_p}$  are the **velocity components** and  $\frac{\partial}{\partial x^\mu} \Big|_{t=t_p}$  are the **basis elements of  $T_p X$**  induced by the chart  $(A, x)$ .

**Definition (Chart induced basis for tangent space)**

Given a chart  $(A, \phi)$  with  $p \in A$  in a smooth atlas then  $\left\{ \frac{\partial f}{\partial x^\mu} \Big|_p \right\} \subset T_p X$  is a basis of  $T_p A$ .

*Proof.* We have seen that all velocity vectors  $\mathbf{V}_{\gamma, p}$  may be written in this basis so  $\text{Span}\left\{ \frac{\partial f}{\partial x^\mu} \Big|_p \right\} = T_p A$ . To check linear independence we let  $f = \phi^\nu$ :

$$V^\mu \frac{\partial}{\partial x^\mu} \Big|_{t=t_p} \phi^\nu = 0 \implies V^\mu \frac{\partial}{\partial x^\mu} (\phi^\nu \circ \gamma)(t) \Big|_{t=t_p} = 0 \implies V^\mu \frac{\partial x^\nu}{\partial x^\mu} = 0 \implies V^\mu \delta_\mu^\nu = V^\nu = 0 \quad (55.1.38)$$

as desired. ■

## 55.2 Vectors and 1-forms

It is interesting to see how the tangent vector components transform under a change of chart. Consider two overlapping charts  $(A, \phi)$  and  $(\tilde{A}, \varphi)$  with  $p \in U \cap \tilde{U}$ . Let  $\mathbf{V} \in T_p M$  be a tangent vector at  $p$  along some smooth curve  $\gamma$ , and which can therefore be expanded in the two chart induced bases  $\left\{ \frac{\partial}{\partial x^\mu} \right\}$  and  $\left\{ \frac{\partial}{\partial \tilde{x}^\mu} \right\}$  <sup>1</sup>:

$$\mathbf{V} = V^\mu \frac{\partial}{\partial x^\mu} = \tilde{V}^\nu \frac{\partial}{\partial \tilde{x}^\nu} \quad (55.2.1)$$

We see that for  $f \in C^\infty(X)$ :

$$\frac{\partial}{\partial x^\mu} f = \partial_\mu (f \circ \phi^{-1})(\mathbf{x}) \quad (55.2.2)$$

$$= \partial_\mu ((f \circ \tilde{\phi}^{-1}) \circ (\tilde{\phi} \circ \phi^{-1}))(\mathbf{x}) \quad (55.2.3)$$

$$= (\partial_\mu (\tilde{\phi} \circ \phi^{-1}))^\nu(\mathbf{x}) \cdot \partial_\nu (f \circ \tilde{\phi}^{-1})(\tilde{\mathbf{x}}) \quad (55.2.4)$$

$$= (\partial_\mu (\tilde{\phi}^\nu \circ \phi^{-1}))(\mathbf{x}) \cdot \partial_\nu (f \circ \tilde{\phi}^{-1})(\tilde{\mathbf{x}}) \quad (55.2.5)$$

$$= \frac{\partial \tilde{x}^\nu}{\partial x^\mu} \frac{\partial}{\partial \tilde{x}^\nu} f \quad (55.2.6)$$

Consequently:

$$\tilde{V}^\nu = \frac{\partial \tilde{x}^\nu}{\partial x^\mu} V^\mu \quad (55.2.7)$$

so the tangent vector components  $V^\nu$  transforms contravariantly, they are **contravariant vector components**.

Since the tangent space  $T_p X$  is a vector space, we can consider its dual vector space.

**Definition (Cotangent space)**

Given the tangent space  $T_p X$  for  $p \in x$  then its **cotangent space** is defined as  $T_p^* X \equiv \text{Hom}(T_p X, \mathbb{R})$ .

For example, let  $f \in C^\infty(X)$  be a smooth function. We introduce the following map in the cotangent

<sup>1</sup>where  $\tilde{x}^\mu = (\tilde{\phi} \circ \gamma)(t)$  and  $x^\mu = (\phi \circ \gamma)(t)$

space:

$$df|_p : T_p X \rightarrow \mathbb{R} \quad (55.2.8)$$

$$\mathbf{V} \rightarrow \mathbf{V}(f) \quad (55.2.9)$$

known as the **gradient** of  $f$  at  $p$ . This is a  $(0, 1)$ -tensor over  $T_p X$ , a covector as defined earlier, also known as a **1-form**. Its components in the chart induced basis are:

$$(df|_p)_\mu = df|_p(\mathbf{e}_\mu) = df|_p\left(\frac{\partial}{\partial x^\mu}\Big|_p\right) = \frac{\partial}{\partial x^\mu}\Big|_p f \quad (55.2.10)$$

giving the standard definition of gradient from vector calculus.

**Theorem (Chart induced basis for cotangent space)**

Let  $(A, \varphi)$  be a chart in the smooth atlas. Then  $\{d\varphi^1|_p, d\varphi^2|_p, \dots, d\varphi^d|_p\}$  is the dual basis of  $T_p^* X$ .

*Proof.* We have that:

$$d\varphi^\mu|_p(\mathbf{e}_\nu) = d\varphi^\mu|_p\left(\frac{\partial}{\partial x^\nu}\Big|_p\right) = \frac{\partial}{\partial x^\nu}\Big|_p \varphi^\mu = \frac{\partial x^\mu}{\partial x^\nu}\Big|_p = \delta_\nu^\mu \quad (55.2.11)$$

so  $\{d\varphi^1|_p, d\varphi^2|_p, \dots, d\varphi^d|_p\}$  is indeed the dual basis. ■

Again, much like with tangent vectors we are interested in how the components of a 1-form transform. We have that given  $\omega \in T_p^* X$  it can be written as:

$$\omega = \omega_\mu dx^\mu|_p = \tilde{\omega}_\nu d\tilde{x}^\nu|_p \quad (55.2.12)$$

Therefore:

$$d\varphi^\mu|_p(\mathbf{V}) = \mathbf{V}(\varphi^\mu) = (\varphi^\mu \circ \gamma)'(t_p) \quad (55.2.13)$$

$$= ((\varphi^\mu \circ (\tilde{\varphi}^{-1})^\nu) \circ (\tilde{\varphi}^\nu \circ \gamma))'(t_p) \quad (55.2.14)$$

$$= (\tilde{\varphi}^\nu \circ \gamma)'(t_p) \cdot (\varphi^\mu \circ (\tilde{\varphi}^{-1})^\nu)'(\tilde{x}_p^\nu) \quad (55.2.15)$$

$$= \mathbf{V}(\tilde{\varphi}^\nu) \cdot \tilde{\partial}_\nu(\varphi^\mu \circ (\tilde{\varphi}^{-1})^\nu)(\tilde{x}_p^\nu) \quad (55.2.16)$$

and since:

$$\tilde{\partial}_\nu(\varphi^\mu \circ (\tilde{\varphi}^{-1})^\nu)(x_p^\nu) = \frac{\partial}{\partial \tilde{x}^\nu} \varphi^\mu \circ (\tilde{\varphi}^{-1})^\nu \Big|_{\tilde{x}_p^\nu} = \frac{\partial \varphi^\mu(p)}{\partial \tilde{x}^\nu} = \frac{\partial x^\mu}{\partial \tilde{x}^\nu} \quad (55.2.17)$$

we get that:

$$d\varphi^\mu|_p(\mathbf{V}) = \frac{\partial x^\mu}{\partial \tilde{x}^\nu} d\tilde{x}^\nu|_p(\mathbf{V}) \quad (55.2.18)$$

and hence the covector components of  $\omega$  transform as:

$$\tilde{\omega}_\nu = \frac{\partial x^\mu}{\partial \tilde{x}^\nu} \omega_\mu \quad (55.2.19)$$

This transformation law is in sharp contrast to (55.2.7), the transforming matrices are inverses of

each other:

$$\tilde{\mathbf{V}} = \mathbf{J}\mathbf{V}, \tilde{\omega} = \mathbf{J}^{-1}\omega \quad (55.2.20)$$

where  $\mathbf{J}$  is known as the **Jacobian matrix**.

### Definition (*Push-forward*)

Let  $\phi : X \rightarrow Y$  be a smooth map between two smooth manifolds. Then we define the **push-forward**  $\phi_*$  of  $\phi$  at the point  $p \in X$  is the linear map:

$$\phi_* : T_p X \rightarrow T_{\phi(p)} Y \quad (55.2.21)$$

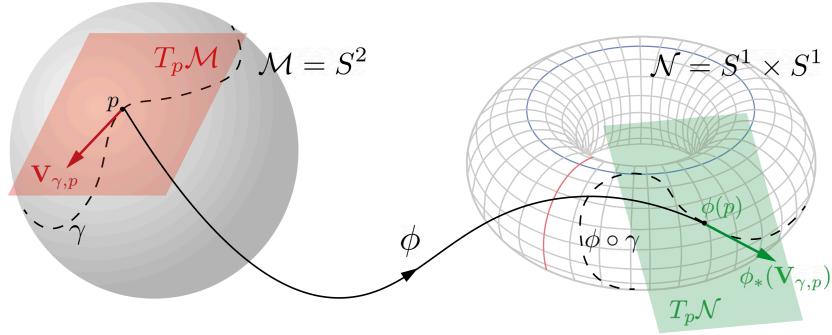
$$\mathbf{V} \mapsto \phi_*(\mathbf{V}) \quad (55.2.22)$$

where we define for any  $f \in C^\infty(X)$

$$\phi_*(\mathbf{V})(f) = \mathbf{V}(f \circ \phi) \quad (55.2.23)$$

We call  $\phi_*$  the **derivative** of  $\phi$ .

Note that  $f \circ \phi \in C^\infty(X)$ . Also note that the definition of  $\phi^*(\mathbf{V})f$  is forced upon us since there are no other linear maps from the two tangent spaces. Let's apply this definition to the tangent vector



$\mathbf{V}_{\gamma,p} \in T_p X$  of some curve  $\gamma$  at  $p \in X$ . Then:

$$\phi_*(\mathbf{V}_{\gamma,p})(f) = \mathbf{V}_{\gamma,p}(f \circ \phi) = ((f \circ \phi) \circ \gamma)'(t_p) = (f \circ (\phi \circ \gamma))(t_p) \quad (55.2.24)$$

but the tangent vector to  $\phi \circ \gamma$  at  $\phi(p)$  is:

$$\mathbf{V}_{\phi \circ \gamma, \phi(p)}(f) = (f \circ (\phi \circ \gamma))(t_p) \quad (55.2.25)$$

since  $(\phi \circ \gamma)(t_p) = \phi(p)$ .

$$\begin{array}{ccc} T_p X & \xrightarrow{\phi_*} & T_{\phi(p)} Y \\ \downarrow \pi_{T_p X} & & \downarrow \pi_{T_{\phi(p)} Y} \\ X & \xrightarrow{\phi} & Y \xrightarrow{f} \mathbb{R} \\ & \searrow f \circ \phi & \end{array}$$

Therefore, the push forward map maps a tangent vector at a point along some curve to the tangent vector at the image of that point along the image of that curve:

$$\phi_*(\mathbf{V}_{\gamma,p}) = \mathbf{V}_{\phi \circ \gamma, \phi(p)} \quad (55.2.26)$$

### Definition (*Pull-back*)

Let  $\phi : X \rightarrow Y$  be a smooth map between two smooth manifolds. Then we define the **pull-back**  $\phi^*$  of  $\phi$  at the point  $\phi(p) \in Y$  is the linear map:

$$\phi^* : T_p^* X \leftarrow T_{\phi(p)}^* Y \quad (55.2.27)$$

$$\phi^*(\omega)\omega \quad (55.2.28)$$

where we define for any  $\omega \in T_{\phi(p)}^* X, \mathbf{V} \in T_p X$ :

$$\phi^*(\omega)(\mathbf{V}) = \omega(\phi_*(\mathbf{V})) \quad (55.2.29)$$

Working in a local basis  $y^\mu$  for  $Y$  and  $\partial_\nu$  for  $X$  we find that

$$(\phi^*)^\mu{}_\nu = \phi^*(dy^\mu)(\partial_\nu) = dy^\mu(\phi_*(\partial_\nu)) = (\phi_*)^\mu{}_\nu \quad (55.2.30)$$

so the components of the push-forwards and pull-backs are the same! These can be expressed more simply as:

$$dy^\mu(\phi_*(\partial_\nu)) = \phi_*(\partial_\nu)y^\mu = \frac{\partial}{\partial x^\nu}(y^\mu \circ \phi) = \frac{\partial X^\mu}{\partial x^\nu} \quad (55.2.31)$$

where we defined  $X^\mu = y^\mu \circ \phi$ .

## 55.3 Interlude: Embeddings and immersions

The concept of an embedding/immersion is quite intuitive: we want some smooth manifold to “sit” within some  $\mathbb{R}^n$ . As we shall learn, there are two known ways a manifold can sit, one way gives an immersion while the other gives an embedding.

### Definition (*Immersion*)

Let  $\phi : M \rightarrow \mathbb{R}^n$  be a smooth map. Then  $\phi$  is an **immersion** of  $M$  into  $\mathbb{R}^n$  if  $\phi_*$  is injective at any  $p \in M$ .

Consider the map  $\phi : S^1 \rightarrow \mathbb{R}^2$  as shown below.

Even though  $\phi$  is not injective, it is easy to see that  $\phi_*$  is injective so this is indeed an immersion.

### Definition (*Embedding*)

Let  $\phi : M \rightarrow \mathbb{R}^n$  be an immersion. Then it is an **embedding** if  $\phi(M)$  is homeomorphic to  $N$ .

We see that  $\phi$  defined previously is not an embedding since  $\phi(S^1)$  is not a manifold and cannot thus be homeomorphic.

**Theorem (Whitney embedding theorems)**

Any smooth manifold  $M$  can be:

- (i) embedded in  $\mathbb{R}^{2 \dim M}$ .
- (ii) immersed in  $\mathbb{R}^{2 \dim M - 1}$

There is a weaker form of Whitney's theorem that is quite amusing:

**Theorem (Immersion theorem)**

Any smooth manifold  $M$  can be immersed in  $\mathbb{R}^{2 \dim M - a(\dim M)}$  where  $a(\dim M)$  is the number of 1s appearing in the binary form of  $\dim M$ .

## 55.4 The tangent bundle

We begin by introducing the concept of bundles, fibres and sections that will be fundamental in defining tensor fields.

**Definition (Bundle)**

A **bundle** is a triple  $(X, Y, \pi)$  where  $X, Y$  are topological manifolds and  $\pi$  is a continuous surjective map from  $X$  to  $Y$ .

**Definition (Fibre bundle)**

Let  $p \in Y$  and let  $(X, Y, \pi)$  be a bundle. Then  $\text{preIm}_\pi(p)$  is the **fibre** at point  $p$ . Moreover, let  $Z$  be another manifold. If for all  $p \in Y$ ,  $\text{preIm}_\pi(p)$  is homeomorphic to some manifold  $M$ , then  $(X, Y, \pi)$  is a **fibre bundle**.

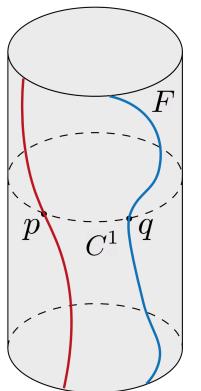
**Definition (Section)**

Let  $(X, Y, \pi)$  be a bundle. Then  $\sigma : Y \rightarrow X$  is a **section** of  $(X, Y, \pi)$  if  $\pi \circ \sigma = \text{id}_M$ .

We see that the section  $\sigma$  maps  $p \in Y$  to some  $\sigma(p) \in \text{preIm}_\pi(p)$  so that  $\pi$  maps it back to itself.

We can visualize bundles, fibres and sections as in Figure ???. Here we consider the surjective map  $\pi : F \rightarrow \mathcal{C}^1$  mapping the cylindrical manifold  $F$  to the ring manifold  $\mathcal{C}^1$ . The green line crossing  $p$  can be viewed as the fibre at  $p$  since it is  $\pi$ 's preimage at that point. For the same reason the blue line is the fibre at  $q$ . A section  $\sigma$  would then map  $p$  to a point on the green line.

A special type of bundle is the tangent bundle, defined as the disjoint union of the tangent spaces at all points on a manifold.

**Definition (Tangent bundle)**

Let  $X$  be a smooth manifold. Then its tangent bundle  $TX$  is defined as:

$$TX = \coprod_{p \in X} T_p X \quad (55.4.1)$$

Similarly we may define the projective map of the tangent bundle as follows.

### Definition (Bundle projection)

Let  $TM$  be the tangent bundle of some smooth manifold  $X$ . Then the bundle projection of the tangent bundle is

$$\pi : TX \rightarrow X \quad (55.4.2)$$

$$\mathbf{V} \mapsto p \quad (55.4.3)$$

where  $p$  is the point in  $X$  such that  $\mathbf{V} \in T_p X$ .

Note that since the tangent spaces  $T_p X$  are disjoint, there exists only one such  $p$ . Note however that to claim that this indeed a bundle would require  $TX$  to be itself a manifold. We now prove this.

We can make  $TX$  a smooth manifold by using the smooth atlas  $\mathcal{A}_X$  on  $X$ . Take some chart  $(U, \phi) \in \mathcal{A}_X$  and consider  $(\text{preIm}_\pi(U), \eta)$  where  $\eta$  is some chart map to  $\mathbb{R}^{2 \dim X}$ .

Let  $\mathbf{V} \in TX$ . Then we can define:

$$\eta(\mathbf{V}) = (\phi^1 \pi(\mathbf{V}), \dots, \phi^{\dim X} \pi(\mathbf{V}), \dots) \quad (55.4.4)$$

where we are missing  $\dim X$  coordinates that specify how “far”  $\mathbf{V}$  is from  $U$ . The basis of  $T_{\pi(\phi)} X$  (this tangent space contains  $\mathbf{V}$ ) is

$$\left\{ \frac{\partial}{\partial x^\mu} \Big|_{\pi(x)} \right\} \quad (55.4.5)$$

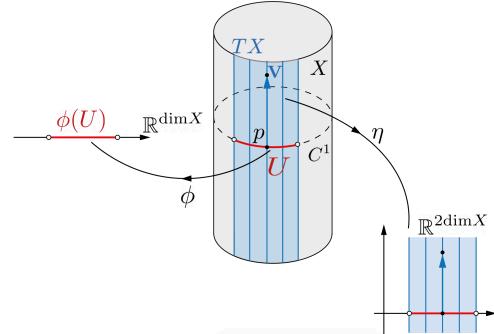
which can be used to expand  $\mathbf{V}$  as:

$$\mathbf{V} = V^\mu \frac{\partial}{\partial x^\mu} \Big|_{\pi(x)} \quad (55.4.6)$$

We define the missing components as  $V^\mu$ :

$$\eta(\mathbf{V}) = (\phi^1 \pi(\mathbf{V}), \dots, \phi^{\dim X} \pi(\mathbf{V}), V^1, \dots, V^{\dim X}) \quad (55.4.7)$$

It is easy to verify that any two such charts are  $C^\infty$ -compatible making the tangent bundle a smooth manifold. Now that we know that the tangent bundle is a smooth manifold, we can provide a rigorous definition for tensor fields (that is not simply “something that assigns a tensor at every point on a manifold”)..



## 55.5 Tensor fields

### Definition (Vector field)

Let  $X$  be a smooth manifold and let  $TX$  be its tangent bundle with bundle projection  $\pi$ . Then a **vector field** is a smooth section of  $TX$ . We further define the set of all vector fields  $\Gamma(TX)$ .

Intuitively, this means that a vector field smoothly maps to each point  $p$  in the manifold  $X$  a tangent vector in its fibre  $\mathbf{V} \in T_p X$ . Note that had  $TX$  not been a smooth manifold we would not have been able to define the vector field as a smooth section.

We can equip  $\Gamma(TX)$  with two operations:

$$\oplus : \Gamma(TX) \times \Gamma(TX) \rightarrow \Gamma(TX) \quad (55.5.1)$$

$$(\sigma, \tau) \mapsto \sigma \oplus \tau \quad (55.5.2)$$

where

$$(\sigma \oplus \tau)(p) = \sigma(p) + \tau(p) \quad (55.5.3)$$

and

$$\odot : C^\infty(X) \times \Gamma(TX) \rightarrow \Gamma(TX) \quad (55.5.4)$$

$$(G, \sigma) \mapsto f \odot \sigma \quad (55.5.5)$$

where

$$(f \odot \sigma)(p) = f(p)\dot{\sigma}(p) \quad (55.5.6)$$

It is easy to check that these two operations satisfy the vector space axioms.

Note that  $(C^\infty(X), +, \cdot)$  where  $\cdot$  is the familiar scalar multiplication is a vector space over  $\mathbb{R}$ . In contrast,  $(C^\infty(X), +, \cdot)$  is not a vector space but has the structure of a ring since if a function in  $C^\infty(X)$  has a zero somewhere then it will not have a multiplicative inverse function. Therefore,  $\Gamma(TX)$  satisfies the typical vector space axioms when equipped with  $\oplus$  and  $\odot$  only that it is defined over a ring rather than a field. Such spaces are known as ring modules, as we shall now see.

### Definition (Ring)

A **ring** is a triplet  $(R, +, \cdot)$  satisfying  $C^+ A^+ N^+ I^+ (C^\bullet) A^\bullet (N^\bullet) (I^\bullet) D$  where  $\cdot$ -commutativity only applies to commutative rings,  $\cdot$ -neutrality only applies to unital element, and  $\cdot$ -invertibility only applies to division rings. Commutative, unital division rings are fields.

Consider once again  $(C^\infty(X), +, \cdot)$ . It is easy to see that it is a commutative, unital ring.

### Definition (Ring module)

Let  $(M, \oplus, \odot)$  defined over a ring  $R$  equipped with two operations:

$$\oplus : M \times M \rightarrow M \quad (55.5.7)$$

$$\odot : R \times M \rightarrow M \quad (55.5.8)$$

satisfying the vector space axioms is an  $R$ -module (ring-module).

Therefore,  $(\Gamma(TX), \oplus, \odot)$  is a  $C^\infty(X)$ -module.

Note that this definition implies that modules don't necessarily have a basis (but division modules do).

Consider for example the sphere  $S^2$  and let  $\mathbf{v} \in \Gamma(TS^2)$  be a vector field on it. There is an important result in algebraic topology, known as the **hairy ball theorem**, that states that one cannot have a

vector field on a sphere that is smooth and non-zero everywhere. This is a problem because it means that one cannot produce a set of linearly independent vector fields on  $S^2$ .

Generally the following result holds:

**Theorem (Division ring basis)**

Let  $D$  be a division ring and let  $V$  be a  $D$ -module. Then  $V$  has a basis.

**Definition ( $R$ -module terminology)**

A **free module** is a module over a ring that possesses a basis. A **projective module** is a module  $\Gamma$  over a ring  $R$  such that there exists another  $R$ -module  $Q$  whose direct sum gives a free module:

$$\exists Q \text{ s.t. } \Gamma \oplus Q = F \quad (55.5.9)$$

**Theorem (Serre, Swan et al)**

The set of all smooth functions of a vector fibre bundle over a smooth manifold  $X$  is finitely generated projective  $C^\infty(X)$ -module  $\Gamma(E)$ .

**Theorem (Homomorphism space is finitely generated)**

Let  $P, Q$  be finitely generated (projective) modules over a commutative ring  $R$ . Then:

$$\text{Hom}_R(P, Q) = \{\varphi : P \rightarrow Q, \varphi \text{ is linear}\} \quad (55.5.10)$$

equipped with  $\oplus$  and  $\odot$  is a finitely generated (projective) module.

This result in particular shows that:

$$\text{Hom}_{C^\infty(X)}(\Gamma(TX), C^\infty(X)) = \Gamma(TX)^* = \Gamma(T^*X) \quad (55.5.11)$$

is a module.

**Definition (Tensor field)**

An  $(r, s)$ -tensor field  $T$  on a smooth manifold  $X$  is a  $C^\infty(X)$ -multilinear map:

$$T : \underbrace{\Gamma(T^*X) \times \dots \times \Gamma(T^*X)}_r \times \underbrace{\Gamma(TX) \times \dots \times \Gamma(TX)}_s \rightarrow C^\infty(X) \quad (55.5.12)$$

and the space of all  $(r, s)$ -tensor fields is denoted by  $\mathcal{T}_s^r$ .

We can define the differential operator the same way for tensor fields:

$$df : \mathcal{T}_s^r \rightarrow C^\infty(M) \quad (55.5.13)$$

$$T \mapsto T(\omega_1, \dots, \mathbf{V}_1, \dots) \quad (55.5.14)$$

# Differentiable forms

## 56.1 Differentiable forms

### Definition (Differentiable form)

Let  $X$  be a smooth manifold, then a **differentiable  $n$ -form** is a  $(0, n)$ -tensor field  $\omega$ :

$$\omega : \underbrace{\Gamma(TX) \times \dots \times \Gamma(TX)}_n \rightarrow C^\infty(X) \quad (56.1.1)$$

that is totally anti-symmetric:

$$\omega(\mathbf{V}_1, \dots, \mathbf{V}_n) = \text{sgn}(\pi)\omega(\mathbf{V}_{\pi(1)}, \dots, \mathbf{V}_{\pi(n)}), \quad \forall \pi \in S_n, \mathbf{V}_i \in \Gamma(TM) \quad (56.1.2)$$

where  $S_n$  is the symmetric group of order  $n$ . We denote the set of all  $n$ -forms on  $X$  by  $\Omega^n(X)$ .

It is important to note that  $\Omega^0(X) = C^\infty(X)$  and  $\Omega^1(X) = \Gamma(T^*X)$ .

### Definition (Wedge product)

The **wedge product** is defined as:

$$\wedge : \Omega^n(X) \times \Omega^m(X) \rightarrow \Omega^{n+m}(X) \quad (56.1.3)$$

$$(\omega_1, \omega_2) \mapsto \omega_1 \wedge \omega_2 \quad (56.1.4)$$

where:

$$(\omega_1 \wedge \omega_2)(\mathbf{V}_1, \dots, \mathbf{V}_{n+m}) = \frac{1}{n!} \frac{1}{m!} \sum_{\pi \in S_{n+m}} \text{sgn}(\pi) (\omega_1 \otimes \omega_2)(\mathbf{V}_{\pi(1)}, \dots, \mathbf{V}_{\pi(n+m)}) \quad (56.1.5)$$

For example, if  $\omega_1, \omega_2 \in \Omega^1(X)$  are one forms then:

$$\omega_1 \wedge \omega_2 = \omega_1 \otimes \omega_2 - \omega_2 \otimes \omega_1 \quad (56.1.6)$$

### Theorem (Basis for $\Gamma^n(X)$ )

Let  $X$  be a smooth  $n$ -manifold with a chart  $(U, \varphi)$ . Taking the wedge product of the chart

induced covariant basis  $\{d\varphi^\mu\}$ :

$$\{d\varphi^{\mu_1} \wedge d\varphi^{\mu_2} \wedge \dots \wedge d\varphi^{\mu_k} : 1 \leq \mu_1 < \dots < \mu_k \leq n\} \quad (56.1.7)$$

gives an ordered basis for  $\Omega^k(X)$ .

*Proof.* We begin by proving that this basis spans  $\Omega^n(X)$ . Let  $\omega \in \Omega^k(X)$  and define:

$$\omega_{\mu_1 \dots \mu_n} = \omega(\mathbf{e}_{\mu_1}, \dots, \mathbf{e}_{\mu_n}) \quad (56.1.8)$$

where  $\{\mathbf{e}_\mu\}$  is the chart-induced contravariant basis. We claim that

$$\omega = \omega_{\mu_1 \dots \mu_k} d\varphi^{\mu_1} \wedge \dots \wedge d\varphi^{\mu_k} \quad (56.1.9)$$

To see why this holds, note that since  $\omega$  is multi-linear we just need to look at the action of the candidate basis on  $(\mathbf{e}_{\nu_1}, \dots, \mathbf{e}_{\nu_k})$ :

$$\omega_{\mu_1 \dots \mu_k} (d\varphi^{\mu_1} \wedge \dots \wedge d\varphi^{\mu_k})(\mathbf{e}_{\nu_1}, \dots, \mathbf{e}_{\nu_k}) \quad (56.1.10)$$

Due to the duality of  $\mathbf{e}_{\nu_j}$  and  $d\varphi^{\mu_i}$ , the contributing terms will be delta functions  $\delta_{\nu_j}^{\mu_i}$ . It is helpful to consider the  $k = 2$  case:

$$\omega_{\alpha\beta} (d\varphi^\alpha \wedge d\varphi^\beta)(\mathbf{e}_\mu, \mathbf{e}_\nu) = \omega_{\alpha\beta} (\delta_\mu^\alpha \delta_\nu^\beta - \delta_\nu^\alpha \delta_\mu^\beta) = \frac{1}{2} (\omega_{\mu\nu} - \omega_{\nu\mu}) = \omega_{\mu\nu} \quad (56.1.11)$$

where the  $\frac{1}{2}$  comes because the sum is restricted to  $\alpha < \beta$ . Note that these Kronecker-deltas act independently of each other, that is given a  $(\mu, \nu)$  then only one of these products of  $\delta$  will be non-zero, yielding:

$$\omega_{12} (d\varphi^1 \wedge d\varphi^2)(\mathbf{e}_1, \mathbf{e}_2) = \omega_{12} \quad (56.1.12)$$

$$\omega_{12} (d\varphi^1 \wedge d\varphi^2)(\mathbf{e}_2, \mathbf{e}_1) = -\omega_{12} = \omega_{21} \quad (56.1.13)$$

as desired. Hence more generally we have that:

$$(d\varphi^{\mu_1} \wedge \dots \wedge d\varphi^{\mu_k})(\mathbf{e}_{\nu_1}, \dots, \mathbf{e}_{\nu_k}) = \sum_{\pi \in S_k} \text{sgn}(\pi) (d\varphi^{\mu_1} \otimes \dots \otimes d\varphi^{\mu_k})(\mathbf{e}_{\pi(\nu_1)}, \dots, \mathbf{e}_{\pi(\nu_k)}) \quad (56.1.14)$$

$$= \sum_{\pi \in S_k} \text{sgn}(\pi) \delta_{\pi(\nu_1)}^{\mu_1} \dots \delta_{\pi(\nu_k)}^{\mu_k} \quad (56.1.15)$$

When we contract this with  $\omega_{\mu_1 \dots \mu_k}$ , the sum over  $\mu_1 \dots \mu_k$  is restricted to  $1 \leq \mu_1 < \mu_2 < \dots < \mu_{k-1} < \mu_k \leq n$ , giving  $\frac{1}{k!}$  many terms as the unrestricted sum<sup>1</sup>. Consequently

$$\omega_{\mu_1 \dots \mu_k} (d\varphi^{\mu_1} \wedge \dots \wedge d\varphi^{\mu_k})(\mathbf{e}_{\nu_1}, \dots, \mathbf{e}_{\nu_k}) = \frac{1}{k!} \sum_{\pi \in S_k} \text{sgn}(\pi) \omega_{\mu_1 \dots \mu_k} \delta_{\pi(\nu_1)}^{\mu_1} \dots \delta_{\pi(\nu_k)}^{\mu_k} \quad (56.1.16)$$

$$= \frac{1}{k!} \sum_{\pi \in S_k} \text{sgn}(\pi) \omega_{\pi(\nu_1) \dots \pi(\nu_k)} = \omega(\mathbf{e}_{\nu_1}, \dots, \mathbf{e}_{\nu_k}) \quad (56.1.17)$$

<sup>1</sup>indeed if  $\omega_{\mu_1 \dots \mu_k} (d\varphi^{\mu_1} \wedge \dots \wedge d\varphi^{\mu_k})$  appears in the restricted sum, then the unrestricted sum will contain  $\omega_{\mu_1 \dots \mu_k} (d\varphi^{\mu_1} \wedge \dots \wedge d\varphi^{\mu_k})$  plus all  $k! - 1$  permutations of the indices. Since  $\omega$  is anti-symmetric and so is the wedge product, permuting indices gives no sign change so the unrestricted sum will indeed be  $k!$  times larger.

where by anti-symmetry  $\text{sgn}(\pi)\omega_{\pi(\nu_1)\dots\pi(\nu_k)} = \omega(\mathbf{e}_{\nu_1}, \dots, \mathbf{e}_{\nu_k})$ . Thus every form  $\omega \in \Omega^k(X)$  can be expanded locally in the basis  $\{d\varphi^{\mu_1} \wedge \dots \wedge d\varphi^{\mu_k}\}$ .

Linear independence immediately follows from the fact that if:

$$\omega_{\mu_1\dots\mu_k}(d\varphi^{\mu_1} \wedge \dots \wedge d\varphi^{\mu_k}) = 0 \quad (56.1.18)$$

then  $\omega_{\mu_1\dots\mu_k}(d\varphi^{\mu_1} \wedge \dots \wedge d\varphi^{\mu_k})(\mathbf{e}_{\nu_1}, \dots, \mathbf{e}_{\nu_k}) = \omega_{\nu_1\dots\nu_k} = 0$  as desired.  $\blacksquare$

Since we can define bases for differential forms we should be interested in their transformation properties.

### Theorem (Transformation of $n$ -form)

Let  $M$  be a smooth  $n$ -manifold and let  $\omega, \tau$  be two differentiable 1-forms on  $M$  related by:

$$\omega^\mu = A^\mu{}_\nu \tau^\nu \quad (56.1.19)$$

Then:

$$\omega^{\mu_1} \wedge \dots \wedge \omega^{\mu_n} = (\det A) \tau^{\nu_1} \wedge \dots \wedge \tau^{\nu_n} \quad (56.1.20)$$

*Proof.* We have that:

$$\omega^{\mu_1} \wedge \dots \wedge \omega^{\mu_n} = (A_{\nu_1}^{\mu_1} \tau^{\nu_1}) \wedge \dots \wedge (A_{\nu_n}^{\mu_n} \tau^{\nu_n}) \quad (56.1.21)$$

$$= \sum_{\sigma \in S_n} A_{\sigma(\nu_1)}^{\mu_1} \dots A_{\sigma(\nu_n)}^{\mu_n} \tau^{\sigma(\nu_1)} \wedge \dots \wedge \tau^{\sigma(\nu_n)} \quad (56.1.22)$$

We can now rearrange the  $\tau^{\sigma(\nu_1)}$  so that the run from  $\tau^{\nu_1}$  to  $\tau^{\nu_n}$  which of course comes with a  $\text{sgn}(\sigma)$  factor:

$$\omega^{\mu_1} \wedge \dots \wedge \omega^{\mu_n} = \sum_{\sigma \in S_n} \text{sgn}(\sigma) A_{\sigma(\nu_1)}^{\mu_1} \dots A_{\sigma(\nu_n)}^{\mu_n} \tau^{\nu_1} \wedge \dots \wedge \tau^{\nu_n} \quad (56.1.23)$$

$$= (\det A) \tau^{\nu_1} \wedge \dots \wedge \tau^{\nu_n} \quad (56.1.24)$$

as desired.  $\blacksquare$

An immediate consequence of this theorem is that:

$$d\varphi^{\mu_1} \wedge \dots \wedge d\varphi^{\mu_n} = \det \left( \frac{\partial x^\mu}{\partial \tilde{x}^\nu} \right) d\tilde{\varphi}^{\nu_1} \wedge \dots \wedge d\tilde{\varphi}^{\nu_n} \quad (56.1.25)$$

which will be a fundamental result in generalizing our notion of integration on manifolds. Of course, this also implies that when a  $n$ -form  $\omega$  is expanded in these two bases as <sup>2</sup>:

$$\omega = ad\varphi^1 \wedge \dots \wedge d\varphi^n = \tilde{a}d\tilde{\varphi}^1 \wedge \dots \wedge d\tilde{\varphi}^n \quad (56.1.26)$$

<sup>2</sup>note that there is only one way to order the differentials so that the superscripts are increasing

the components transform according to

$$a' = a \det \left( \frac{\partial x^\mu}{\partial \tilde{x}^\nu} \right) \quad (56.1.27)$$

It will also be useful to extend our notion of push-forward and pull-back to  $n$ -forms.

### Definition (Pull-back)

Let  $\omega \in \Omega^n(Y)$  and let  $\phi : X \rightarrow Y$  be smooth. Then we define the pull-back of  $\omega$  as:

$$\Phi^*(\omega)(\mathbf{V}_1, \dots, \mathbf{V}_n) = \omega(\Phi_*(\mathbf{V}_1), \dots, \Phi_*(\mathbf{V}_n)) \quad (56.1.28)$$

### Theorem (Pull-back of wedge product)

The pull-back distributes over the wedge product:

$$\Phi^*(\omega_1 \wedge \omega_2) = \Phi^*(\omega_1) \wedge \Phi^*(\omega_2) \quad (56.1.29)$$

*Proof.* It is easy to verify that if  $\omega_1$  is a  $p$ -form and  $\omega_2$  is a  $q$ -form and  $n = p + q$  then:

$$\Phi^*(\omega_1 \wedge \omega_2)(\mathbf{V}_1, \dots, \mathbf{V}_n) = (\omega_1 \wedge \omega_2)(\Phi_*(\mathbf{V}_1), \dots, \Phi_*(\mathbf{V}_n)) \quad (56.1.30)$$

$$= \frac{1}{n!} \frac{1}{m!} \sum_{\pi \in S_{n+m}} \text{sgn}(\pi) (\omega_1 \otimes \omega_2)(\Phi_*(\mathbf{V}_{\pi(1)}), \dots, \Phi_*(\mathbf{V}_{\pi(n)})) \quad (56.1.31)$$

$$= \frac{1}{n!} \frac{1}{m!} \sum_{\pi \in S_{n+m}} \text{sgn}(\pi) (\Phi^*(\omega_1) \otimes \Phi^*(\omega_2))(\mathbf{V}_{\pi(1)}, \dots, \mathbf{V}_{\pi(n)}) \quad (56.1.32)$$

$$= \Phi^*(\omega_1) \wedge \Phi^*(\omega_2) \quad (56.1.33)$$

as desired. ■

It would be nice to have a space that is closed under the wedge product.

Define the following space:

$$\Omega(X) = \bigoplus_{i=0}^{\dim X} \Omega^i(X) \quad (56.1.34)$$

Then  $(\Omega(X), +, \cdot, \wedge)$  is known as the **Grassmann algebra** on  $X$ , and:

$$\wedge : \Omega(X) \times \Omega(X) \rightarrow \Omega(X) \quad (56.1.35)$$

For example, let  $\sigma = \omega_1 + \omega_2$  where  $\omega_1 \in \Omega^1(X)$  and  $\omega_2 \in \Omega^2(X)$ , and let  $\tau \in \Omega^n(X)$ . Then:

$$\tau \wedge \sigma = \tau \wedge (\omega_1 + \omega_2) := \tau \wedge \omega_1 + \tau \wedge \omega_2 \quad (56.1.36)$$

where since  $\tau \wedge \omega_1 \in \Omega^{n+1}(X)$  and  $\tau \wedge \omega_2 \in \Omega^{n+3}(X)$ , the addition  $+$  of these forms must be in the Grassmann algebra  $\Omega(X)$ .

Let's look at the commutativity of the wedge product:

**Proposition (Wedge commutativity)**

Let  $\omega_1 \in \Omega^n(X)$  and  $\omega_2 \in \Omega^m(X)$  be two forms. Then:

$$\omega_1 \wedge \omega_2 = (-1)^{nm} \omega_2 \wedge \omega_1 \quad (56.1.37)$$

*Proof.* We have that:

$$(\omega_1 \wedge \omega_2)(\mathbf{V}_1, \dots, \mathbf{V}_{n+m}) = \frac{1}{n!m!} \sum_{\pi \in S_{n+m}} \text{sgn}(\pi) (\omega_1 \otimes \omega_2)(\mathbf{V}_{\pi(1)}, \dots, \mathbf{V}_{\pi(n+m)}) \quad (56.1.38)$$

$$= \frac{(-1)^n}{n!m!} \sum_{\pi \in S_{n+m}} \text{sgn}(\pi) (\omega_1 \otimes \omega_2)(\mathbf{V}_{\pi(n+1)}, \dots, \mathbf{V}_{\pi(n)}, \mathbf{V}_{\pi(1)}, \dots, \mathbf{V}_{\pi(n+m)}) \quad (56.1.39)$$

⋮

$$= \frac{(-1)^{nm}}{n!m!} \sum_{\pi \in S_{n+m}} \text{sgn}(\pi) (\omega_1 \otimes \omega_2)(\mathbf{V}_{\pi(n+1)}, \dots, \mathbf{V}_{\pi(n+m)}, \mathbf{V}_{\pi(1)}, \dots, \mathbf{V}_{\pi(n)}) \quad (56.1.40)$$

$$= \frac{(-1)^{nm}}{n!} \frac{1}{m!} \sum_{\pi \in S_{n+m}} \text{sgn}(\pi) (\omega_2 \otimes \omega_1)(\mathbf{V}_{\pi(1)}, \dots, \mathbf{V}_{\pi(n+m)}) \quad (56.1.41)$$

$$= (-1)^{nm} \omega_2 \wedge \omega_1 \quad (56.1.42)$$

as desired. ■

## 56.2 The exterior derivative

**Definition (Exterior derivative)**

We define the **exterior derivative**:

$$d : \Omega^n(X) \rightarrow \Omega^{n+1}(X) \quad (56.2.1)$$

$$\omega \mapsto d\omega \quad (56.2.2)$$

where  $\forall \mathbf{V}_\mu \in \Gamma(TX)$ :

$$(d\omega)(\mathbf{V}_1, \dots, \mathbf{V}_{n+1}) = \sum_i (-1)^{i+1} \mathbf{V}_i (\omega(\mathbf{V}_1, \dots, \cancel{\mathbf{V}_i}, \dots, \mathbf{V}_{n+1})) \quad (56.2.3)$$

$$+ \sum_{i \leq j} \omega([\mathbf{V}_i, \mathbf{V}_j], \mathbf{V}_1, \dots, \cancel{\mathbf{V}_i}, \dots, \cancel{\mathbf{V}_j}, \dots, \mathbf{V}_{n+1}) \quad (56.2.4)$$

where  $\cancel{\mathbf{V}_i}$  means that the vector field  $\mathbf{V}_i$  is omitted.

For example, we have that if  $n = 1$  then the exterior derivative reads

$$d\omega(\mathbf{V}, \mathbf{W}) = \mathbf{V}(\omega(\mathbf{W})) - \mathbf{W}(\omega(\mathbf{V})) - \omega([\mathbf{V}, \mathbf{W}]) \quad (56.2.5)$$

Then we see that:

$$d\omega(\mathbf{V}, \mathbf{W}) = -d\omega(\mathbf{W}, \mathbf{V}) \quad (56.2.6)$$

by standard commutator rules. Moreover:

$$d\omega(\mathbf{V}, f\mathbf{W}) = \mathbf{V}(\omega(f\mathbf{W})) - f\mathbf{W}(\omega(\mathbf{V})) - \omega(\mathbf{V}(f\mathbf{W}) + f\mathbf{W}\mathbf{V}) \quad (56.2.7)$$

$$= f\mathbf{V}(\omega(\mathbf{W})) + \mathbf{V}(f)\omega(\mathbf{W}) - f\mathbf{W}(\omega(\mathbf{V})) - \omega(\mathbf{V}(f)\mathbf{W} + f\mathbf{V}\mathbf{W} - f\mathbf{W}\mathbf{V}) \quad (56.2.8)$$

The last term can be written as:

$$\mathbf{V}(f)\omega(\mathbf{W}) + f\omega([\mathbf{V}, \mathbf{W}]) \quad (56.2.9)$$

so that:

$$d\omega(\mathbf{V}, f\mathbf{W}) = f\mathbf{V}\omega(\mathbf{W}) - f\mathbf{W}(\omega(\mathbf{V})) - f\omega([\mathbf{V}, \mathbf{W}]) \quad (56.2.10)$$

$$= fd\omega(\mathbf{V}, \mathbf{W}) \quad (56.2.11)$$

as desired.

Now let's look at how  $d$  acts on the wedge product.

### Proposition (Exterior derivative of wedge product)

Let  $\omega_1 \in \Omega^n(X)$  and  $\omega_2 \in \Omega^m(X)$ . Then:

$$d(\omega_1 \wedge \omega_2) = d\omega \wedge \omega_2 + (-1)^n \omega_1 \wedge d\omega_2 \quad (56.2.12)$$

*Proof.*

### Theorem (Exterior derivative commutes with pull-back)

Let  $\phi : X \rightarrow Y$  be a smooth map between two manifolds. Then exterior derivative  $d$  commutes with the pull-back  $\Phi^*$ :

$$\Phi^*(d\omega) = d(\Phi^*(\omega)), \forall \omega \in T^*Y \quad (56.2.13)$$

where  $\omega$  is an  $n$ -form on  $Y$ .

*Proof.*

### Definition (Anti-symmetrisation bracket)

Let  $A_{\mu_1 \dots \mu_n}$  be some object with  $n$  indices. Then we define:

$$A_{[\mu_1 \dots \mu_n]} = \frac{1}{n!} \sum_{\pi \in S_n} \text{sgn}(\pi) A_{\pi(\mu_1) \dots \pi(\mu_n)} \quad (56.2.14)$$

The same goes for superscripts.

It follows from this definition is:

$$A_{\mu\nu}B^{[\mu\nu]} = \frac{1}{n!} \sum_{\pi \in S_n} \text{sgn}(\pi) A_{\mu\nu} B^{\pi(\mu_1)\pi(\mu_2)} \quad (56.2.15)$$

$$= \frac{1}{n!} \sum_{\pi \in S_n} \text{sgn}(\pi) A_{\pi^{-1}(\mu)\pi^{-1}(\nu)} B^{\mu\nu} \quad (56.2.16)$$

$$= \frac{1}{n!} \sum_{\pi \in S_n} \text{sgn}(\pi) A_{\pi(\mu)\pi(\nu)} B^{\mu\nu} \quad (56.2.17)$$

$$= A_{[\mu\nu]} B^{\mu\nu} \quad (56.2.18)$$

It also follows from this definition that:

### Definition (*Exact and closed forms*)

Let  $\omega \in \Omega^n(X)$ . Then  $\omega$  is **exact** if  $\omega \in \text{Im}(d_n)$  and closed if  $\omega \in \ker(d_{n+1})$ . We let  $B^n$  denote the set of exact  $n$ -forms on  $X$  and  $Z^n$  denote the set of closed  $n$ -forms on  $X$ .

## 56.3 de Rham cohomology and Electromagnetism

### Theorem ( $d^2 = 0$ )

Given any  $n$ -form  $\omega \in \Omega^n(M)$ , then the  $n + 2$ -form  $d^2\omega \equiv (d \circ d)(\omega) = 0$  is closed.

Let  $\omega \in \Omega^n(X)$ , then given a chart  $(U, \phi)$  with induced local basis  $\{d\phi^\mu\}$ :

$$d\omega = (\partial_\nu \omega_{\mu_1, \dots, \mu_n}) d\phi^\nu \wedge d\phi^{\mu_1} \wedge \dots \wedge d\phi^{\mu_n} \quad (56.3.1)$$

$$\implies d^2\omega = (\partial_\alpha \partial_\nu \omega_{\mu_1, \dots, \mu_n}) d\phi^\alpha d\phi^\nu \wedge d\phi^{\mu_1} \wedge \dots \wedge d\phi^{\mu_n} \quad (56.3.2)$$

Note that the  $\alpha, \nu$  indices in the wedge product are anti-symmetric by definition. However, as long as  $\omega_{\mu_1, \dots, \mu_n} \in C^2$  then  $\alpha, \nu$  are symmetric since partial derivatives can be commuted. The contraction of symmetric and anti-symmetric indices gives zero so we find that  $d^2\omega = 0$  as desired.

Let's now look at the following sequence of maps:

$$\Omega^{n-1}(X) \xrightarrow{d_n} \Omega^n(X) \xrightarrow{d_{n+1}} \Omega^{n+1}(X) \quad (56.3.3)$$

Since  $d^2 = 0$  it follows that  $\text{Im}d_n \subseteq \ker d_{n+1}$ . In other words **all exact forms are closed**:

$$B^n \subseteq Z^n \quad (56.3.4)$$

However, it is not generally true that  $B^n$  is equivalent to  $Z^n$ . One important case is when we are working in euclidean topologies:

### Theorem (*Poincare lemma*)

If  $X = \mathbb{R}^m$  then  $B^n = Z^n$  for  $n > 0$ .

Let  $F$  be the electromagnetic field strength. Then we know that

$$F = dA \implies dF = 0 \quad (56.3.5)$$

for some  $A \in \Omega^1(X)$ . This reproduces the homogeneous Maxwell equations!

**Definition (*de Rham cohomology groups*)**

The  $n$ -th **deRham cohomology group** is the quotient vector space:

$$H^n(X) = Z^n / B^n \quad (56.3.6)$$

# Connections and parallel transport

## 57.1 Covariant derivatives

Recall that a vector field  $\mathbf{V} \in \Gamma(TM)$  can be used to produce a directional derivative  $\mathbf{V}(f)$  of a smooth map  $f \in C^\infty(X)$ . For notational simplicity we will define:

$$\nabla_{\mathbf{V}} f \equiv \mathbf{V}(f) \quad (57.1.1)$$

One might wonder why we should have this cumbersome notation, especially since  $df$  was already defined so that  $df(\mathbf{V}) = \mathbf{V}(f)$ . However, note that

$$X : C^\infty(X) \rightarrow C^\infty(X) \quad (57.1.2)$$

$$df : \Gamma(TX) \rightarrow C^\infty(X) \quad (57.1.3)$$

so  $\mathbf{V}$  and  $df$  are different things. Moreover,  $\nabla_{\mathbf{V}}$ , which as of now only works on smooth functions ((0, 0)-tensors), will soon be extended so as to act on  $C^\infty(X)$ -tensor fields, thus requiring a different symbol.

A directional derivative should have a list of properties that we would like to be satisfied.

### Definition (*Connection*)

A **connection**  $\nabla$  on a smooth manifold  $X$  is a map that maps a pair of a vector (field)  $\mathbf{V}$  and a  $(p, q)$  tensor field  $T$  to  $(p, q)$ -tensor (field)  $\nabla_{\mathbf{V}} T$ , such that:

- (i) Extension rule:  $\nabla_{\mathbf{V}} f = \mathbf{V}(f)$ ,  $\forall f \in C^\infty(M)$
- (ii) Additivity 1 rule:  $\nabla_{\mathbf{V}}(T + S) = \nabla_{\mathbf{V}} T + \nabla_{\mathbf{V}} S$ ,  $\forall T, S \in \mathcal{T}_q^p(M)$
- (iii) Additivity 2 rule:  $\nabla_{f\mathbf{V}+\mathbf{U}} T = f\nabla_{\mathbf{V}} T + \nabla_{\mathbf{U}} T$
- (iv) Leibniz rule:  $\nabla_{\mathbf{V}} \underbrace{T(\omega, \mathbf{W})}_{\in C^\infty(X)} = (\nabla_{\mathbf{V}} T)(\omega, \mathbf{W}) + T(\nabla_{\mathbf{V}} \omega, \mathbf{W}) + T(\omega, \nabla_{\mathbf{V}} \mathbf{W})$  and analogously for any  $(p, q)$  tensor field  $T$ .

How many such connections are there, is the definition above enough to fix just one connection? Let  $\mathbf{V}, \mathbf{W}$  be vector fields, then working in a local basis:

$$\nabla_{\mathbf{V}} \mathbf{W} = \nabla_{V^\mu \partial_\mu} W^\nu \partial_\nu \quad (57.1.4)$$

$$= V^\mu \left[ \left( \nabla_{\partial_\mu} W^\nu \right) \partial_\nu + W^\nu \nabla_{\partial_\mu} \partial_\nu \right] \quad (57.1.5)$$

$$= V^\mu \frac{\partial W^\nu}{\partial x^\mu} \frac{\partial}{\partial x^\nu} + V^\mu W^\nu \nabla_{\partial_\mu} \left( \frac{\partial}{\partial x^\nu} \right) \quad (57.1.6)$$

Now the second term is not fixed since we do not know what the action of the connection is on a vector. However, we note that the result of  $\nabla_{\partial_\mu} \frac{\partial}{\partial x^\nu}$  will itself be a vector so it may be expanded as:

$$\nabla_{\partial_\mu} \left( \frac{\partial}{\partial x^\nu} \right) = \Gamma_{\nu\mu}^\alpha \frac{\partial}{\partial x^\alpha} \iff \Gamma_{\nu\mu}^\alpha = d\varphi^\alpha \left( \nabla_{\partial_\mu} \frac{\partial}{\partial x^\nu} \right) \quad (57.1.7)$$

where  $\Gamma_{\nu\mu}^\alpha$  are the **connection coefficient functions**, and can be recognized as **Christoffel symbols**. Therefore:

$$(\nabla_{\mathbf{V}} \mathbf{W})^\alpha = \mathbf{V}(W^\alpha) + \Gamma_{\nu\mu}^\alpha V^\mu W^\nu \quad (57.1.8)$$

We are not done yet, we still have to see if given the connection coefficient functions, the action of  $\nabla_{\mathbf{V}}$  on a 1-form is also fixed. By similar reasoning as before:

$$\nabla_{\mathbf{V}} \omega = \nabla_{V^\mu \partial_\mu} \omega_\nu d\varphi^\nu \quad (57.1.9)$$

$$= V^\mu \left[ \left( \nabla_{\partial_\mu} \omega_\nu \right) d\varphi^\nu + \omega_\nu \nabla_{\partial_\mu} d\varphi^\nu \right] \quad (57.1.10)$$

$$= V^\mu \frac{\partial \omega_\nu}{\partial x^\mu} d\varphi^\nu + V^\mu \omega_\nu \nabla_{\partial_\mu} (d\varphi^\nu) \quad (57.1.11)$$

Note however that:

$$\nabla_{\partial_\mu} (d\varphi^\nu (\mathbf{e}_\alpha)) = \nabla_{\partial_\mu} (\delta_\alpha^\nu) = 0 \quad (57.1.12)$$

and using the Leibnitz rule:

$$\nabla_{\partial_\mu} (d\varphi^\nu (\mathbf{e}_\alpha)) = \nabla_{\partial_\mu} (d\varphi^\nu) (\mathbf{e}_\alpha) + d\varphi^\nu \left( \nabla_{\partial_\mu} (\mathbf{e}_\alpha) \right) = 0 \quad (57.1.13)$$

$$\implies \nabla_{\partial_\mu} (d\varphi^\nu) (\mathbf{e}_\alpha) = -d\varphi^\nu \left( \Gamma_{\alpha\mu}^\beta \frac{\partial}{\partial x^\beta} \right) \quad (57.1.14)$$

$$\implies \nabla_{\partial_\mu} (d\varphi^\nu) (\mathbf{e}_\alpha) = -\Gamma_{\alpha\mu}^\nu \iff \nabla_{\partial_\mu} (d\varphi^\nu) = -\Gamma_{\alpha\mu}^\nu d\varphi^\alpha \quad (57.1.15)$$

Consequently:

$$(\nabla_{\mathbf{V}} \omega)_\alpha = \mathbf{V}(\omega_\alpha) - \Gamma_{\alpha\mu}^\nu V^\mu \omega_\nu \quad (57.1.16)$$

Similarly we get that for a  $(1, 1)$ -tensor  $T$ , the covariant derivative acts as:

$$(\nabla_{\mathbf{V}} T)_\beta^\alpha = \mathbf{V}(T_\beta^\alpha) + \Gamma_{\nu\mu}^\alpha V^\mu T_\beta^\nu - \Gamma_{\beta\mu}^\nu V^\mu T_\nu^\alpha \quad (57.1.17)$$

This can be seen by applying the Leibnitz rule to  $T = \mathbf{T} \otimes \tau$ , where  $\mathbf{T} \in \Gamma(TX)$  and  $\tau \in \Gamma(T^*X)$ .

### Definition (Divergence)

Let  $\mathbf{V}$  be a vector field on a manifold. The divergence of  $\mathbf{V}$  is given by:

$$\text{div } \mathbf{V} = (\nabla_{\partial_\mu} \mathbf{V})^\mu \quad (57.1.18)$$

We have not checked the change of the components of  $\Gamma$  when we perform a chart transition. Let  $(U, \varphi)$  and  $(\tilde{U}, \tilde{\varphi})$  be two smooth charts with  $U \cap \tilde{U} \neq \emptyset$ . Then, starting from  $\Gamma'$  and transforming

all  $\tilde{x}$  to  $x$ :

$$\tilde{\Gamma}_{\nu\mu}^\alpha = d\tilde{\varphi}^\alpha \left( \nabla_{\tilde{\partial}_\mu} \frac{\partial}{\partial \tilde{x}^\nu} \right) = \frac{\partial \tilde{x}^\alpha}{\partial x^\beta} d\varphi^\beta \left( \nabla_{\tilde{\partial}_\mu} x^\delta \partial_\delta \frac{\partial x^\eta}{\partial \tilde{x}^\nu} \frac{\partial}{\partial x^\eta} \right) \quad (57.1.19)$$

$$= \frac{\partial \tilde{x}^\alpha}{\partial x^\beta} d\varphi^\beta \left\{ \frac{\partial x^\delta}{\partial \tilde{x}^\mu} \left[ \nabla_{\tilde{\partial}_\delta} \left( \frac{\partial x^\eta}{\partial \tilde{x}^\nu} \right) \frac{\partial}{\partial x^\eta} + \frac{\partial x^\eta}{\partial \tilde{x}^\nu} \nabla_{\tilde{\partial}_\delta} \left( \frac{\partial}{\partial x^\eta} \right) \right] \right\} \quad (57.1.20)$$

$$= \frac{\partial \tilde{x}^\alpha}{\partial x^\beta} \frac{\partial x^\delta}{\partial \tilde{x}^\mu} \left[ \frac{\partial}{\partial x^\delta} \left( \frac{\partial x^\eta}{\partial \tilde{x}^\nu} \right) \delta_\eta^\beta + \frac{\partial x^\eta}{\partial \tilde{x}^\nu} \Gamma_{\eta\delta}^\beta \right] \quad (57.1.21)$$

$$\Rightarrow \boxed{\tilde{\Gamma}_{\nu\mu}^\alpha = \frac{\partial \tilde{x}^\alpha}{\partial x^\beta} \frac{\partial x^\delta}{\partial \tilde{x}^\mu} \frac{\partial x^\eta}{\partial \tilde{x}^\nu} \Gamma_{\eta\delta}^\beta + \frac{\partial \tilde{x}^\alpha}{\partial x^\beta} \frac{\partial}{\partial \tilde{x}^\mu} \left( \frac{\partial x^\beta}{\partial \tilde{x}^\nu} \right)} \quad (57.1.22)$$

So we see that in general the connection coefficient functions do not transform as tensor components. Note that even if  $\Gamma$  is zero in some chart, it may not be zero in another chart.

### Theorem (Symmetric part of $\Gamma$ vanish)

Let  $p \in X$  be a point in a smooth connection. Then one can always construct a chart  $(U, \varphi)$  containing  $p$  such that the symmetric part of the connection coefficient functions vanish at that point  $\Gamma_{(\nu\mu)}^\alpha = 0$ . The corresponding coordinates are known as **normal coordinates**.

*Proof.* Let  $(U, \varphi)$  be some smooth chart with  $p \in V$ , so that  $\Gamma$  does not necessarily vanish. We construct a new chart  $(\tilde{U}, \tilde{\varphi})$  with  $p \in \tilde{U}$  with chart transition map:

$$(\tilde{\varphi} \circ \varphi^{-1})^\alpha(x^1, \dots, x^d) = x^\mu - \Gamma_{(\nu\mu)}^\alpha(p)x^\mu x^\nu \quad (57.1.23)$$

where  $(x^1, \dots, x^d) \in \mathbb{R}^d$  and  $\Gamma_{(\nu\mu)}^\alpha$  is evaluated at the point  $p$  and is thus a constant. Then we see that:

$$\frac{\partial \tilde{x}^\beta}{\partial x^\nu} = \frac{\partial}{\partial x^\nu} (\tilde{\varphi}^\beta \circ \phi) = \delta_\nu^\mu - \Gamma_{(\nu\mu)}^\beta(p)x^\mu \quad (57.1.24)$$

$$\Rightarrow \frac{\partial}{\partial x^\mu} \frac{\partial \tilde{x}^\beta}{\partial x^\nu} = -\Gamma_{(\nu\mu)}^\beta(p) \quad (57.1.25)$$

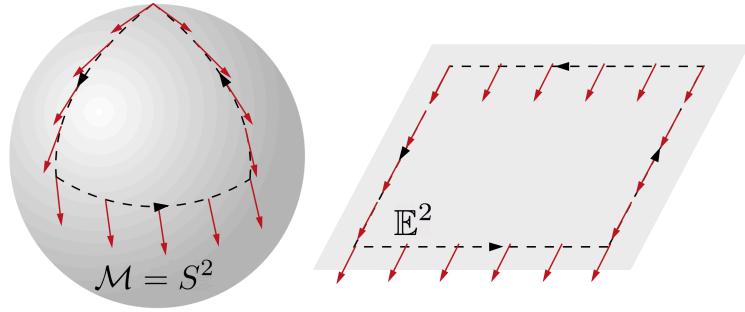
Consequently the new connection coefficient function at  $p$  takes the form:

$$\tilde{\Gamma}_{(\nu\mu)}^\alpha(p) = 0 \quad (57.1.26)$$

as desired. ■

## 57.2 Parallel transport

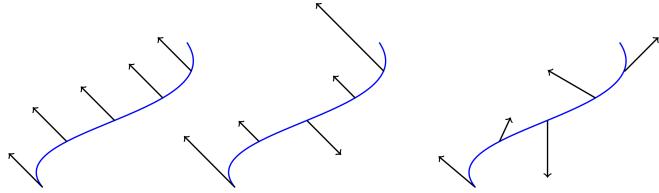
Suppose you are on an expedition in the North pole, with your nose pointing out of this page, and start walking down until you reach the equator. You then move east along one fourth of the equator, and up back to the initial point, always keeping your nose pointing in the same direction relative to your body. It is clear that your nose at the end will be pointing in a different direction to when it started! This is an example of parallel transport, the movement of a vector field  $\mathbf{W}$  along a path  $\gamma$  on a smooth manifold so that it stays parallel to the connection  $\nabla$  we impose on it.



**Definition (Parallel transport)** Let  $X$  be a smooth manifold with a connection  $\nabla$  and let  $\mathbf{W}$  be a vector field on  $X$ . Then we say that  $\mathbf{W}$  is **parallelly transported** along a smooth curve  $\gamma$  on  $X$  if:

$$V^\mu \nabla_\mu \mathbf{W} = \nabla_{\mathbf{V}} \mathbf{W} = 0 \quad (57.2.1)$$

Consider the three curves and vector fields below:



In the first case the vector field is parallelly transported along  $\gamma$ . In the second case the vector field is just parallel to  $\gamma$ , but since its magnitude changes along it there is no parallel transports. Finally, in the third case the curve is neither parallel nor parallelly transported.

**Definition (Autoparallel transport)**

Let  $\gamma : \mathbb{R} \rightarrow X$  be a smooth curve on a smooth manifold  $X$  with tangent vector  $\mathbf{V}$ . Then  $\gamma$  is **autoparallelly transported** if:

$$\nabla_{\mathbf{V}} \mathbf{V} = 0 \quad (57.2.2)$$

Such curves  $\gamma$  are known as **geodesics**.

Suppose we choose a chart  $(U, \phi) \in \mathcal{A}_{C^\infty}$  and consider the portion of  $\gamma$  in  $U$ . Working in the local chart-induced basis, we can write  $\nabla_\mu \equiv \nabla_{\partial_\mu}$  as long as we are working just in one chart:

$$\nabla_{\mathbf{V}} \mathbf{V} = V^\mu \nabla_\mu (V^\nu \partial_\nu) \quad (57.2.3)$$

$$= V^\mu \left( (\nabla_\mu V^\nu) \frac{\partial}{\partial x^\nu} + V^\nu \Gamma_{\nu\mu}^\sigma \frac{\partial}{\partial x^\sigma} \right) \quad (57.2.4)$$

$$= \left( V^\mu \frac{\partial V^\sigma}{\partial x^\mu} + V^\mu V^\nu \Gamma_{\nu\mu}^\sigma \right) \frac{\partial}{\partial x^\sigma} \quad (57.2.5)$$

Recalling that if  $\gamma$  is parametrised by  $t$  then  $V^\mu = \frac{dx^\mu}{dt}$ , we see that  $V^\mu \frac{\partial V^\sigma}{\partial x^\mu} = \frac{d^2 x^\sigma}{dt^2}$ , giving the

following, very important equation:

$$\boxed{\frac{d^2x^\sigma}{dt^2} + \Gamma_{\mu\nu}^\sigma \frac{dx^\mu}{dt} \frac{dx^\nu}{dt} = 0} \quad (57.2.6)$$

known as the **Geodesic equation**. This is the condition for a curve  $\gamma$  with coordinates  $x^\mu$  in a given chart to be autoparallelly transported.

For example, in Euclidean space  $\mathbb{R}^2$  with standard topology, and where  $\Gamma_{jl}^i$  by definition, we would have that straight lines:

$$\frac{d^2x^\sigma}{dt^2} = 0 \implies x^\sigma(t) = A^\sigma t + B^\sigma \quad (57.2.7)$$

are the geodesics.

Consider a universe with at least two particles interacting gravitationally. Here Newton's first law is completely useless since there is no particle such that the force acting upon it is zero. To salvage the first law one could envisage gravity not as a force but a curvature of space-time. Then we could have a particle with no force acting upon it (gravity is not a force anymore), and the path that it takes would be the geodesics described by the geodesic equation which is mathematically equivalent to the path given by the second law. This is qualitatively the same description given by General relativity, as we shall see in the subsection on Newtonian space-time.

### Definition (Torsion)

The **torsion** of a connection  $\nabla$  is:

$$T(\omega, \mathbf{V}, \mathbf{W}) = \omega(\nabla_{\mathbf{V}}\mathbf{W} - \nabla_{\mathbf{W}}\mathbf{V} - [\mathbf{V}, \mathbf{W}]) \quad (57.2.8)$$

where

$$[\mathbf{V}, \mathbf{W}](f) = \mathbf{V}(\mathbf{W}(f)) - \mathbf{W}(\mathbf{V}(f)) \quad (57.2.9)$$

*Proof.* We must check that this is a tensor field. Linearity in the first argument is trivial, but for the other two some care is needed. Indeed, given  $f \in C^\infty$ :

$$T(\omega, f\mathbf{V}, \mathbf{W}) = \omega(\nabla_{f\mathbf{V}}\mathbf{W} - \nabla_{\mathbf{W}}(f\mathbf{V}) - [f\mathbf{V}, \mathbf{W}]) \quad (57.2.10)$$

$$= \omega(f\nabla_{\mathbf{V}}\mathbf{W} - (\mathbf{W}(f))\mathbf{V} - f\nabla_{\mathbf{W}}\mathbf{V} - f[\mathbf{V}, \mathbf{W}] + (\mathbf{W}(f))\mathbf{V}) \quad (57.2.11)$$

$$= f\omega(\nabla_{\mathbf{V}}\mathbf{W} - \nabla_{\mathbf{W}}\mathbf{V} - [\mathbf{V}, \mathbf{W}]) = fT(\omega, \mathbf{V}, \mathbf{W}) \quad (57.2.12)$$

where to go from the second to third step we used the property:

$$[f\mathbf{V}, \mathbf{W}](g) = f\mathbf{V}(\mathbf{W}(g)) - \mathbf{W}(f\mathbf{V}(g)) \quad (57.2.13)$$

$$= f\mathbf{V}(\mathbf{W}(g)) - \mathbf{W}(f)\mathbf{V}(g) - f\mathbf{W}(\mathbf{V}(g)) \quad (57.2.14)$$

$$= f[\mathbf{V}, \mathbf{W}](g) - \mathbf{W}(f)\mathbf{V}(g) \quad (57.2.15)$$

Similarly we find that for additivity:

$$T(\omega, \mathbf{V}_1 + \mathbf{V}_2, \mathbf{W}) = \omega(\nabla_{\mathbf{V}_1+\mathbf{V}_2}\mathbf{W} - \nabla_{\mathbf{W}}(\mathbf{V}_1 + \mathbf{V}_2) - [\mathbf{V}_1 + \mathbf{V}_2, \mathbf{W}]) \quad (57.2.16)$$

$$= \omega(\nabla_{\mathbf{V}_1}\mathbf{W} + \nabla_{\mathbf{V}_2}\mathbf{W} - \nabla_{\mathbf{W}}\mathbf{V}_1 - \nabla_{\mathbf{W}}\mathbf{V}_2 - [\mathbf{V}_1, \mathbf{W}] - [\mathbf{V}_2, \mathbf{W}]) \quad (57.2.17)$$

$$= T(\omega, \mathbf{V}_1, \mathbf{W}) + T(\omega, \mathbf{V}_2, \mathbf{W}) \quad (57.2.18)$$

as desired. Linearity in the second argument follows from the anti-linearity of the tensor field in  $\mathbf{V}$  and  $\mathbf{W}$ .  $\blacksquare$

### Definition (*Torsion-free manifold*)

A manifold with connection  $\nabla$  is torsion free if  $T = 0$ .

Working in a chart, the condition for a manifold with connection to be torsion free is that:

$$T^\alpha_{\mu\nu} = T(dx^\alpha, \partial_\mu, \partial_\nu) = dx^\mu (\nabla_\mu \partial_\nu - \nabla_\nu \partial_\mu) = \Gamma_{\nu\mu}^\alpha - \Gamma_{\mu\nu}^\alpha = 0 \implies \Gamma_{[\mu\nu]}^\alpha = 0 \quad (57.2.19)$$

so the Christoffel symbols must have vanishing anti-symmetric components. Note that for a torsion-free manifold we can always locally set the Christoffel symbols to zero since we can make both its symmetric and anti-symmetric components equal to zero.

We can visualize torsion as the failure of parallelograms to close. Indeed let us take two vectors fields  $\mathbf{X}$  and  $\mathbf{Y}$  in  $T_x M$ . It is clear that the points  $r$  and  $s$  will have coordinates

$$r : x^\mu + \epsilon X^\mu, s : x^\mu + \epsilon Y^\mu \quad (57.2.20)$$

where  $\epsilon$  is infinitesimal. Suppose we parallelly transport  $\mathbf{X}$  along  $\mathbf{Y}$  giving a new vector

$$X'^\mu = X^\mu - \epsilon \Gamma_{\mu\nu}^\rho X^\mu Y^\nu \quad (57.2.21)$$

The corresponding point  $q$  then has coordinates

$$q : x^\mu + \epsilon(X^\mu + Y^\mu) - \epsilon^2 \Gamma_{\mu\nu}^\rho X^\mu Y^\nu \quad (57.2.22)$$

On the other hand the point  $t$  has coordinates

$$t : x^\mu + \epsilon(X^\mu + Y^\mu) - \epsilon^2 \Gamma_{\mu\nu}^\rho Y^\mu X^\nu \quad (57.2.23)$$

In general, unless torsion vanishes we will find that  $q$  and  $t$  do not coincide, in other words the parallelogram formed by parallelly transporting two vectors does not necessarily close. It is clear how for most physical applications we can safely take torsion to vanish as we do not want any openings in our space-time manifold.

## 57.3 Riemannian curvature

### Definition (*Riemann curvature*)

The **Riemannian curvature** of a connection  $\nabla$  is the tensor field:

$$R(\omega, \mathbf{U}, \mathbf{V}, \mathbf{W}) = \omega(\nabla_{\mathbf{V}} \nabla_{\mathbf{W}} \mathbf{U} - \nabla_{\mathbf{W}} \nabla_{\mathbf{V}} \mathbf{U} - \nabla_{[\mathbf{V}, \mathbf{W}]} \mathbf{U}) \quad (57.3.1)$$

Working in one chart:

$$R^\alpha_{\sigma\mu\nu} = dx^\alpha (\nabla_\mu \nabla_\nu \partial_\sigma - \nabla_\nu \nabla_\mu \partial_\sigma - \nabla_{\partial_\mu \partial_\nu} \partial_\sigma + \nabla_{\partial_\nu \partial_\mu} \partial_\sigma) \quad (57.3.2)$$

$$= dx^\alpha [\nabla_\mu (\Gamma_{\sigma\nu}^\alpha \partial_\alpha) - \nabla_\nu (\Gamma_{\sigma\mu}^\alpha \partial_\alpha) - \nabla_{\partial_\mu \partial_\nu} \partial_\sigma + \nabla_{\partial_\nu \partial_\mu} \partial_\sigma] \quad (57.3.3)$$

$$= dx^\alpha [\Gamma_{\sigma\nu}^\tau \Gamma_{\tau\mu}^\gamma \partial_\gamma - \Gamma_{\sigma\mu}^\tau \Gamma_{\tau\nu}^\gamma + (\partial_\mu \Gamma_{\sigma\nu}^\tau) \partial_\tau - (\partial_\nu \Gamma_{\sigma\mu}^\tau) \partial_\mu] \quad (57.3.4)$$

$$\implies R^\alpha_{\sigma\mu\nu} = \Gamma_{\sigma\nu}^\tau \Gamma_{\tau\mu}^\alpha - \Gamma_{\sigma\mu}^\tau \Gamma_{\tau\nu}^\alpha + \partial_\mu \Gamma_{\sigma\nu}^\alpha - \partial_\nu \Gamma_{\sigma\mu}^\alpha \quad (57.3.5)$$

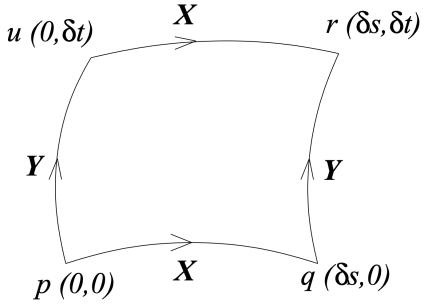
Geometrically, the Riemann curvature tensor captures the failure of vector field  $\mathbf{U}$  to be parallelly transported to the same vector along  $\mathbf{V}$  or  $\mathbf{W}$  on a smooth manifold  $X$ . The more curved our manifold is, the greater the discrepancies between the different parallel transports. Indeed, consider the figure beside. We take two vector fields  $\mathbf{X}$  and  $\mathbf{Y}$  along which we transport another vector field  $\mathbf{Z}$  taking two different paths,  $p \rightarrow q \rightarrow r$  and  $p \rightarrow u \rightarrow r$  on a torsion-free Riemannian manifold. We can thus take  $x^\mu$  to be normal coordinates in which the Christoffel symbols evaluated at  $p$  vanish. We perform all approximations up to second order.

To go from  $p$  to  $q$  it is clear that:

$$(\nabla_X Z)^\rho = X^\nu \frac{dZ^\mu}{dx^\nu} + \Gamma_{\mu\nu}^\rho Z^\mu X^\nu = 0 \quad (57.3.6)$$

Letting  $X^\mu \frac{\partial}{\partial x^\mu} = \frac{\partial}{\partial s}$  then

$$\frac{d^2 Z^\rho}{ds^2} = -X^\sigma \partial_\sigma (\Gamma_{\mu\nu}^\rho Z^\mu X^\nu) \quad (57.3.7)$$



We can now perform a Taylor expansion about  $p$ :

$$Z^\rho(q) = Z^\rho(p) + \left. \frac{dZ^\mu}{ds} \right|_p \delta s + \frac{1}{2} \left. \frac{d^2 Z^\mu}{ds^2} \right|_p (\delta s)^2 \quad (57.3.8)$$

but it is clear that  $\left. \frac{dZ^\mu}{ds} \right|_p$  must vanish so we find that:

$$Z^\rho(q) - Z^\rho(p) = -\frac{1}{2} (X^\sigma Z^\mu X^\nu \partial_\sigma \Gamma_{\mu\nu}^\rho) \Big|_p (\delta s)^2 \quad (57.3.9)$$

Now we go from  $q$  to  $r$ :

$$Z^\rho(r) - Z^\rho(q) = \left. \frac{dZ^\mu}{ds} \right|_q \delta t + \frac{1}{2} \left. \frac{d^2 Z^\mu}{ds^2} \right|_q (\delta t)^2 \quad (57.3.10)$$

$$= -(\Gamma_{\mu\nu}^\rho Z^\mu Y^\nu) \Big|_q \delta t - \frac{1}{2} [X^\sigma \partial_\sigma (\Gamma_{\mu\nu}^\rho Z^\mu Y^\nu)] \Big|_q (\delta t)^2 \quad (57.3.11)$$

but recall that expanding  $(\Gamma_{\mu\nu}^\rho Z^\mu Y^\nu) \Big|_q$  to first order (since we already have a  $\delta t$ )

$$(\Gamma_{\mu\nu}^\rho Z^\mu Y^\nu) \Big|_q = \underbrace{(\Gamma_{\mu\nu}^\rho Z^\mu Y^\nu) \Big|_p}_0 + X^\sigma \partial_\sigma (\Gamma_{\mu\nu}^\rho Z^\mu Y^\nu) \Big|_p \delta s \quad (57.3.12)$$

$$= X^\sigma Z^\mu Y^\nu \partial_\sigma \Gamma_{\mu\nu}^\rho \Big|_p \delta s \quad (57.3.13)$$

hence

$$Z^\rho(r) - Z^\rho(q) = -X^\sigma Z^\mu Y^\nu \partial_\sigma \Gamma_{\mu\nu}^\rho|_p \delta s \delta t - \frac{1}{2} [Y^\sigma \partial_\sigma (\Gamma_{\mu\nu}^\rho Z^\mu Y^\nu)]|_q (\delta t)^2 \quad (57.3.14)$$

Similarly we may expand  $[X^\sigma \partial_\sigma (\Gamma_{\mu\nu}^\rho Z^\mu Y^\nu)]|_q$ , this time to just zeroth order since it comes with  $(\delta t)^2$ . We can thus evaluate this term at  $p$  and simplify:

$$Z^\rho(r) - Z^\rho(q) = -Y^\sigma Z^\mu Y^\nu \partial_\sigma \Gamma_{\mu\nu}^\rho|_p \delta s \delta t - \frac{1}{2} (Y^\sigma Z^\mu Y^\nu \partial_\sigma \Gamma_{\mu\nu}^\rho)|_p (\delta t)^2 \quad (57.3.15)$$

Thus, we find that:

$$(Z^\rho(r) - Z^\rho(p))_{pqr} = -\frac{1}{2} \partial_\sigma \Gamma_{\mu\nu}^\rho (X^\sigma X^\nu (\delta s)^2 + Y^\sigma Y^\nu (\delta t)^2 + 2X^\sigma Y^\nu \delta s \delta t) Z^\mu|_p \quad (57.3.16)$$

By simply replacing  $\mathbf{X} \leftrightarrow \mathbf{Y}$  and  $\delta s \leftrightarrow \delta t$  we find that

$$(Z^\rho(r) - Z^\rho(p))_{pur} = -\frac{1}{2} \partial_\sigma \Gamma_{\mu\nu}^\rho (Y^\sigma Y^\nu (\delta t)^2 + X^\sigma X^\nu (\delta s)^2 + 2Y^\sigma X^\nu \delta s \delta t) Z^\mu|_p \quad (57.3.17)$$

We can now compute the difference between  $Z^\rho(r)$  transported along  $\gamma' : pqr$  and along  $\gamma : pur$

$$\frac{Z_{(\gamma')}^\rho - Z_{(\gamma)}^\rho}{\delta s \delta t} = \partial_\sigma \Gamma_{\mu\nu}^\rho (Y^\sigma X^\nu - X^\sigma Y^\nu) Z^\mu|_p \quad (57.3.18)$$

$$= Z^\mu X^\nu Y^\sigma (\partial_\sigma \Gamma_{\mu\nu}^\rho - \partial_\nu \Gamma_{\mu\sigma}^\rho)|_p \quad (57.3.19)$$

$$\implies \boxed{\frac{Z_{(\gamma')}^\rho - Z_{(\gamma)}^\rho}{\delta s \delta t} = (Z^\mu Y^\sigma X^\nu R_{\mu\sigma\nu}^\rho)|_p} \quad (57.3.20)$$

where we used the fact that at  $p$ , the Riemann curvature tensor reads  $R_{\mu\sigma\nu}^\rho = \partial_\sigma \Gamma_{\mu\nu}^\rho - \partial_\nu \Gamma_{\mu\sigma}^\rho$  in normal coordinates. Note also that this result is not an artifact of working within a special chart. Indeed our end result is a relation between tensor so there is no coordinate dependence.

So when we parallelly transport a vector field along two different paths their final differences is encoded in the Riemann curvature tensor.

# Metrics and Curvature

## 58.1 The metric tensor

We are missing a notion of length and angles in our formulation of differential geometry, which we can establish by imposing a metric on our manifold.

### Definition (Metric tensor)

Let  $X$  be a smooth manifold, then a **metric**  $g$  on  $X$  is a  $(0, 2)$ -tensor field such that:

- (i) symmetry:  $g(\mathbf{V}, \mathbf{W}) = g(\mathbf{W}, \mathbf{V})$
- (ii) non-degeneracy: define the map  $\sharp : \Gamma(TX) \rightarrow \Gamma(T^*X)$  such that  $(\sharp(\mathbf{V}))(\mathbf{W}) = g(\mathbf{V}, \mathbf{W})$ .  
Then  $\sharp$  is a  $C^\infty$ -isomorphism.

We not discuss a bit of convention. we can denote the components of the  $\sharp$ -map acting on  $\mathbf{V}$  as simply:

$$(\sharp(\mathbf{V}))_\mu \equiv V_\mu = g_{\mu\nu} V^\nu \quad (58.1.1)$$

where the last equality was established by definition.

### Definition (Inverse metric)

The symmetric  $(2, 0)$ -tensor field  $g^{-1}$  related to the metric tensor field  $g$  is defined by:

$$g^{-1} : \Gamma(T^*X) \times \Gamma(T^*X) \rightarrow C^\infty(X) \quad (58.1.2)$$

$$(\omega, \tau) \mapsto \omega(\sharp^{-1}(\tau)) \quad (58.1.3)$$

This is not really the inverse of  $g$  since the domains and images don't match up, but can be naively viewed as such when looking at it purely as a matrix. Indeed we have that

$$\sharp(\partial_\nu) = g_{\mu\nu} dx^\mu \implies \sharp^{-1}(g_{\mu\nu} dx^\mu) = \partial_\nu \quad (58.1.4)$$

and hence

$$(g^{-1})^{\mu\nu} g_{\nu\sigma} = dx^\mu (\sharp^{-1}(dx^\nu)) g_{\nu\sigma} = dx^\mu (\sharp^{-1}(g_{\nu\sigma} dx^\nu)) = dx^\mu (\partial_\sigma) = \delta_\sigma^\mu \quad (58.1.5)$$

Hence, since  $\sharp$  is an isomorphism it is invertible we may write:

$$(\sharp^{-1}(\omega)) \equiv \omega^\mu = \omega_\nu \sharp^{-1}((g^{-1})^{\nu\mu} g_{\mu\sigma} dx^\sigma) = (g^{-1})^{\mu\nu} \omega_\nu \partial_\mu \quad (58.1.6)$$

$$\implies \omega^\mu = (g^{-1})^{\mu\nu} \omega_\nu \quad (58.1.7)$$

For simplicity we will often write  $g^{-1})^{\mu\nu} \equiv g^{\mu\nu}$  since the indices being up already indicate that we are taking components of  $g^{-1}$  and not  $g$ . Hence we have found that we can raise and lower the operators of covectors and vectors using the metric:

$$V_\mu = g_{\mu\nu} V^\nu, \omega_\mu = g^{\mu\nu} \omega_\nu \quad (58.1.8)$$

We claim that a metric tensor can always be reduced to a diagonal form where the entries are either  $+1, -1, 0$ . If this diagonal form of the metric contains  $p$  1's and  $q$  -1's then we say that the metric has  $(p, q)$  signature.

### Definition (Metric types)

A metric is said to be:

- (i) **Riemannian** if the signature is  $(+ + \dots +)$ .
- (ii) **pseudo-Riemannian** otherwise, with the special case of a **Lorentzian** metric if the signature is  $(+ - \dots -)$

A manifold equipped with a Riemannian metric is a Riemannian manifold.

### Definition (Curve length)

Let  $X$  be a Riemannian manifold with metric  $g$ . Then the speed of a smooth curve  $\gamma(0, 1) \mapsto X$  parametrised by  $t$  at  $p = (\phi \circ \gamma)(t_p) \in X$  is given by

$$s(t_p) = \sqrt{g(\mathbf{V}_{\gamma,p}, \mathbf{V}_{\gamma,p})} \quad (58.1.9)$$

The length of  $\gamma$  is then given by:

$$L[\gamma] = \int_0^1 dt s(t) = \int_0^1 dt \sqrt{g(\mathbf{V}_\gamma, \mathbf{V}_\gamma)} \quad (58.1.10)$$

Of course, the length of a curve should not depend on the way it is parametrised. This is indeed the case:

### Theorem (Curve reparametrisation)

Let  $\gamma : (0, 1) \rightarrow X$  be a curve and let  $\sigma : (0, 1) \rightarrow (0, 1)$  be a smooth, increasing bijection. Then:

$$L[\gamma] = L[\gamma \circ \sigma] \quad (58.1.11)$$

*Proof.* We have that:

$$V_{\gamma \circ \sigma, t_p} = \frac{d}{dt} f(\gamma(\sigma(t))) \Big|_{t_p} = \frac{d\sigma(t)}{dt} \frac{d}{d\sigma(t)} f(\gamma(\sigma(t))) \Big|_{t_p} = \frac{d\sigma(t)}{dt} \frac{d}{dt} f(\gamma(t)) \Big|_{\sigma(t_p)} \quad (58.1.12)$$

$$= \frac{d\sigma}{dt} V_{\gamma, \sigma(t_p)} \quad (58.1.13)$$

This then gives

$$L[\gamma \circ \sigma] = \int_0^1 dt \sqrt{g(\mathbf{V}_{\gamma \circ \sigma, t}, \mathbf{V}_{\gamma \circ \sigma, t})} \quad (58.1.14)$$

$$= \int_0^1 dt \sqrt{g(\mathbf{V}_{\gamma, \sigma(t)}, \mathbf{V}_{\gamma, \sigma(t)})} \quad (58.1.15)$$

$$= \int_0^1 dt \frac{d\sigma(t)}{dt} \sqrt{g(\mathbf{V}_{\gamma, \sigma(t)}, \mathbf{V}_{\gamma, \sigma(t)})} \quad (58.1.16)$$

We now use change variables  $s = \sigma(t) \implies ds = \frac{d\sigma}{dt} dt$ . The bounds of integration remain the same since  $\sigma$  is bijective and increasing, implying that  $\sigma(0) = 0$ ,  $\sigma(1) = 1$ . Thus:

$$L[\gamma \circ \sigma] = \int_0^1 ds \sqrt{g(\mathbf{V}_{\gamma, s}, \mathbf{V}_{\gamma, s})} = L[\gamma] \quad (58.1.17)$$

as desired. ■

In the context of metric manifolds, we should define geodesics as curves that have the shortest possible length defined by the metric.

### Definition (Geodesic)

A curve  $\gamma$  is a **geodesic** on a Riemannian manifold if it is a stationary curve of  $L$ .

We would of course like this geodesic to be the same geodesic defined by autoparallel transports, thus requiring the connection on a manifold to be determined by the metric we impose on it. Of course we can express the minimum length principle of geodesics using the Euler-Lagrange equation:

$$\frac{\partial \mathcal{L}}{\partial x^\mu(t)} = \frac{d}{dt} \frac{\partial \mathcal{L}}{\partial \dot{x}^\mu(t)} \quad (58.1.18)$$

using the Lagrangian:

$$\mathcal{L} : TX \rightarrow \mathbb{R} \quad (58.1.19)$$

$$\mathbf{V} \mapsto \sqrt{g(\mathbf{V}, \mathbf{V})} \quad (58.1.20)$$

so that  $\mathcal{L}(x^\mu, \dot{x}^\mu) = \sqrt{g_{\mu\nu} \dot{x}^\mu \dot{x}^\nu}$ . Letting  $\sqrt{g(\mathbf{V}, \mathbf{V})} = g$  we find that

$$\frac{\partial \mathcal{L}}{\partial x^\alpha} = \frac{1}{2g} \partial_\alpha g_{\mu\nu}(x) \dot{x}^\mu \dot{x}^\nu \quad (58.1.21)$$

Also:

$$\frac{\partial \mathcal{L}}{\partial \dot{x}^\alpha} = \frac{1}{2g} 2g_{\alpha\mu}(x) \dot{x}^\mu = \frac{1}{g} g_{\alpha\mu}(x) \dot{x}^\mu \quad (58.1.22)$$

so that

$$\frac{d}{dt} \frac{\partial \mathcal{L}}{\partial \dot{x}^\alpha} = \frac{d}{dt} \left( \frac{1}{g} \right) g_{\alpha\mu}(x) \dot{x}^\mu + \frac{1}{g} \left( g_{\alpha\mu}(x) \ddot{x}^\mu + \dot{x}^\nu \frac{\partial}{\partial x^\nu} (g_{\alpha\mu}(x)) \dot{x}^\mu \right) \quad (58.1.23)$$

Now recall that we can re-parametrise our curve without affecting the length functional. Thus we can choose a parametrisation such that  $g(\mathbf{V}, \mathbf{V}) = \text{const}$ , that is such that the speed of the curve is constant (note that the metric can never give zero so this can always be done). The first term then

vanishes, and hence:

$$g_{\alpha\mu}\ddot{x}^\mu + \partial_\nu(g_{\alpha\mu})\dot{x}^\mu\dot{x}^\nu - \frac{1}{2}\partial_\alpha g_{\mu\nu}\dot{x}^\mu\dot{x}^\nu = 0 \quad (58.1.24)$$

We can now use the inverse matrix to raise the  $\alpha$  index, and use  $g^{\alpha\beta}g_{\beta\rho} = \delta_\rho^\alpha$ :

$$\ddot{x}^\sigma + g^{\alpha\sigma}\left(\partial_\nu g_{\alpha\mu} - \frac{1}{2}\partial_\alpha g_{\mu\nu}\right)\dot{x}^\mu\dot{x}^\nu = 0 \quad (58.1.25)$$

Since the  $i, j$  indices are symmetric, we may double the term in parenthesis by writing its copy with  $i, j$  exchanged, and hence write that:

$$\boxed{\frac{d^2x^\sigma}{dt^2} + \frac{1}{2}g^{\alpha\sigma}\left(\partial_\nu g_{\alpha\mu} + \partial_\mu g_{\alpha\nu} - \partial_\alpha g_{\mu\nu}\right)\frac{dx^\mu}{dt}\frac{dx^\nu}{dt} = 0} \quad (58.1.26)$$

Now suppose we wish to establish a connection on our Riemannian manifold, such that the autoparallel geodesics are exactly the geodesics described above. Comparing (57.2.6) with (58.1.26) it follows that the Christoffel symbols must be defined so that:

$$\boxed{\Gamma_\mu^\sigma = \frac{1}{2}g^{\alpha\sigma}\left(\partial_\nu g_{\alpha\mu} + \partial_\mu g_{\alpha\nu} - \partial_\alpha g_{\mu\nu}\right)} \quad (58.1.27)$$

The connection with these Christoffel symbols is known as the **Levi-Civita** connection. There are several useful curvatures that one may define with the Levi-Civita connection.

### Definition (*Important curvatures*)

On a Riemannian manifold with the Levi-Civita connection, the

- (i) **Riemann Christoffel curvature** is defined as

$$R_{\rho\sigma\mu\nu} = g_{\rho\alpha}R^\alpha{}_{\sigma\mu\nu} \quad (58.1.28)$$

- (ii) **Ricci curvature** is defined as:

$$R_{\mu\nu} = R^\sigma_{\mu\sigma\nu} \quad (58.1.29)$$

- (iii) **Scalar curvature**:

$$R = g^{\mu\nu}R_{\mu\nu} \quad (58.1.30)$$

- (iv) **Einstein curvature**:

$$G_{\mu\nu} = R_{\mu\nu} - \frac{1}{2}g_{\mu\nu}R \quad (58.1.31)$$

## 58.2 Lie derivatives and symmetry

We saw how we could push forward vectors between tangent spaces, and pull back 1-forms between dual tangent spaces, but having seen metrics it is reasonable to wonder whether the same can be done with a metric tensor.

If we start with a metric manifold  $N$  and embed a smaller manifold  $M$  in it, then we can define a metric on  $M$  by pushing forward vectors on  $M$  using the embedding and use the metric for  $N$  on them. This defines the so-called induced metric.

**Definition (Induced metric)**

Suppose we have a metric  $g$  on a manifold  $N$ , and suppose  $\phi$  is an embedding (injective map) of a manifold  $M$  into  $N$ . We define the **induced metric**  $g_M$  on  $M$  as:

$$g_M(\mathbf{V}, \mathbf{W}) = g(\phi_*(\mathbf{V}), \phi_*(\mathbf{W})) \quad (58.2.1)$$

or alternatively

$$(g_M)_{\mu\nu} = g_{\alpha\beta} \frac{\partial X^\alpha}{\partial x^\mu} \frac{\partial X^\beta}{\partial x^\nu} \quad (58.2.2)$$

where recall that  $X^\alpha = y^\alpha \circ \phi$ .

**Definition (Integral curve)**

Let  $M$  be a smooth manifold and let  $\mathbf{W}$  be a vector field on  $M$ . A curve  $\gamma$  is an **integral curve** of  $\mathbf{W}$  if the value of the vector field on the curve coincides everywhere with the tangent vectors  $\mathbf{V}$  along the curve:

$$\mathbf{V}_{\gamma, \gamma(t)} = \mathbf{W}_{\gamma(t)} \quad (58.2.3)$$

Furthermore,  $\mathbf{W}$  is complete if all its integral curves have a domain that can be extended to  $\mathbb{R}$ .

It would also be nice if given some complete vector field, we could create a family of integral curves passing through all the points on the manifold. This would generate a flow of the vector field, with every point flowing along it by some distance  $t$ .

**Definition (Vector field flow)**

The **flow of a complete vector field**  $\mathbf{W}$  on  $M$  is:

$$h^{\mathbf{W}} : \mathbb{R} \times M \rightarrow M \quad (58.2.4)$$

$$(t, p) \mapsto \gamma_p(t) \quad (58.2.5)$$

where  $\gamma_p : \mathbb{R} \rightarrow M$  is an integral curve of  $\mathbf{W}$  traversing  $p$  at  $t = 0$ . This flow is also a **one parameter group of diffeomorphisms**.

We can fix  $t$  in the flow map to get  $h_t^{\mathbf{W}}$ . Its action is to take every point  $p$  on the manifold and map it to the point  $\gamma_p(t)$  on the integral curve passing through it at  $t = 0$ , thus shifting it a parameter distance  $t$ .

The differential equation that the coordinates  $y^\mu(t)$  of an integral curve must satisfy is thus

$$\frac{dy^\mu}{dt} = X^\mu, \quad y^\mu(0) = x^\mu \quad (58.2.6)$$

where  $x_\mu^0$  are the coordinates of  $p$ . For infinitesimal flows then we see that to leading order in  $t$

$$y^\mu(h_t^{\mathbf{W}}) = x^\mu + tW^\mu \quad (58.2.7)$$

**Theorem (Vector fields form Lie algebra)**

Let  $\Gamma(TM)$  denote the set of all vector fields on a manifold  $M$ , and let  $[\cdot, \cdot]$  be the commutator:

$$[\mathbf{V}, \mathbf{W}]f = \mathbf{V}(\mathbf{W}(f)) - \mathbf{W}(\mathbf{V}(f)) \quad (58.2.8)$$

which satisfies the following properties:

- (i)  $[\mathbf{V}, \mathbf{W}] = -[\mathbf{W}, \mathbf{V}]$
- (ii)  $[\lambda\mathbf{V} + \mathbf{U}, \mathbf{W}] = \lambda[\mathbf{V}, \mathbf{W}] + [\mathbf{U}, \mathbf{W}]$
- (iii)  $[\mathbf{U}, [\mathbf{V}, \mathbf{W}]] + [\mathbf{W}, [\mathbf{U}, \mathbf{V}]] + [\mathbf{V}, [\mathbf{W}, \mathbf{U}]] = 0$

Then  $(\Gamma(TM), [\cdot, \cdot])$  is a **Lie algebra**.

*Proof.* We prove Jacobi's identity (iii) only as the first two are trivial. ■

### Definition (Lie subalgebra)

Let  $\{\mathbf{V}_i\}_{i=0,1,\dots}$  be vector fields on  $X$  such that:

$$[\mathbf{V}_i, \mathbf{V}_j] = C_{ij}^k \mathbf{V}_k \quad (58.2.9)$$

where  $C_{ij}^k$  are known as **structure constants** of the **Lie subalgebra** ( $\text{span}_{\mathbb{R}}\{\mathbf{V}_i\}, [\cdot, \cdot]$ ).

### Definition (Symmetry)

A Lie subalgebra  $L$  on a smooth metric manifold is a **symmetry** of a metric  $g$  if for any vector fields  $\mathbf{V}$  in  $L$ <sup>a</sup> then:

$$((h_t^{\mathbf{V}})^* g)(\mathbf{X}, \mathbf{Y}) \equiv g((h_t^{\mathbf{V}})_*(\mathbf{X}), (h_t^{\mathbf{V}})_*(\mathbf{Y})) = g(\mathbf{X}, \mathbf{Y}), \forall \mathbf{X}, \mathbf{Y} \in \Gamma(TM), \forall t \quad (58.2.10)$$

<sup>a</sup>which are guaranteed to be complete by a very nice theorem assuming they are compactly supported

This definition makes intuitively sense, it says that if  $L$  is a symmetry, then moving the metric (backwards) along the flows generated by any of its vector fields has no effect at all. We can take the metric of two vector fields at some point, move them along the flows generated by a vector field in the symmetry, and get the same result back.

Despite the intuitiveness behind this definition it is often very difficult to check for symmetries by looking at the pull-back of the metric, it is a very tedious task. A much faster approach involves defining the Lie derivative.

### Definition (Lie derivative)

We define the **Lie derivative** of a tensor  $T$  along some vector field  $\mathbf{V}$  in a Lie sub-algebra  $L$  as

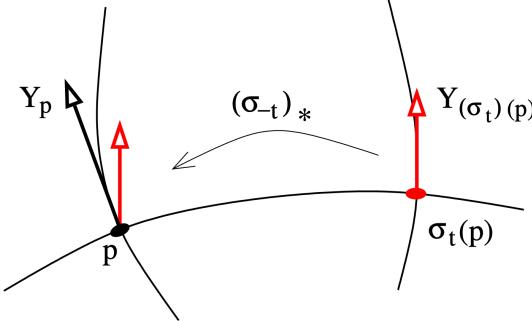
$$\mathcal{L}_{\mathbf{V}} T \equiv \lim_{t \rightarrow 0} \frac{((h_{-t}^{\mathbf{V}})^* T)_p - T_p}{t} \quad (58.2.11)$$

so that if the Lie subalgebra  $L$  is a symmetry then  $\mathcal{L}_{\mathbf{V}} g = 0, \forall \mathbf{V} \in L$ .

This definition is quite close to how we define normal derivatives. Indeed a naive guess would be to simply write

$$\mathcal{L}_{\mathbf{V}} T = \lim_{t \rightarrow \infty} \frac{T_{p+t} - T_p}{t} \quad (58.2.12)$$

but there are two problems with this definition. Firstly, the notion of adding points on a manifold makes no sense,  $p + t$  is not a thing. We need an alternative notion that is physically equivalent, that of flowing the point  $p$  along a  $\mathbf{V}$  by an amount  $t$ , which is precisely what the flow map does.



Secondly, the tensors  $T_{h_t(p)}$  and  $T_p$  do not belong to the same tensor space so it doesn't make sense to simply take their difference (this is like taking the difference of tangent vectors at different points). We can resolve this problem by pushing forward  $T_{h_t(p)}$  to the tangent space at  $p$ , so that the difference may be evaluated. The end result is the definition we provided previously

We can define special coordinates to write the Lie derivative in a nice way. Suppose that we have generated the integral curves of some vector field  $\mathbf{V}$  in a  $n$ -manifold. Let  $p \in M$  and let us create a  $n - 1$  dimensional hypersurface  $\Sigma$  of  $M$  that is not tangent to the integral curve through  $p$ . Then we can let the parametrisation variable  $t$  be the  $n$ th coordinate together with the  $x^i$  coordinates of  $\Sigma$ .

It follows that the push-forward of the flow  $h_t^*$  maps a point  $p$  with coordinates  $(t_p, x_p^i)$  to the point  $q$  with coordinates  $(t_p + t, x_p^i)$ . Therefore

$$((h_t)_*)^\mu_\nu = \frac{\partial y^\mu}{\partial x^n u} = \delta_\nu^\mu \quad (58.2.13)$$

implying that

$$((h_{-t})_* T^{\mu_1 \dots \mu_n}{}_{\nu_1 \dots \nu_n})_p = \frac{\partial y^{\mu_1}}{\partial x^{\sigma_1}} \dots \frac{\partial y^{\mu_n}}{\partial x^{\sigma_n}} \frac{\partial x^{\rho_1}}{\partial y^{\nu_1}} \dots \frac{\partial x^{\rho_n}}{\partial y^{\nu_n}} (T^{\sigma_1 \dots \sigma_n}{}_{\rho_1 \dots \rho_n})_{\phi(p)} \quad (58.2.14)$$

$$= (T^{\mu_1 \dots \mu_n}{}_{\nu_1 \dots \nu_n})(t_p + t, x^i) \quad (58.2.15)$$

Hence we get that

$$(\mathcal{L}_V T)^{\mu_1 \dots \mu_n}{}_{\nu_1 \dots \nu_n} = \frac{\partial}{\partial t} T^{\mu_1 \dots \mu_n}{}_{\nu_1 \dots \nu_n} \Big|_{(t_p, x^i)} \quad (58.2.16)$$

so the Lie derivative in this particular coordinate frame is just the partial derivative of the tensor components.

It can be shown that this limit definition of the Lie derivative is equivalent to the following algebraic definition:

### Definition (Lie derivative)

The Lie derivative  $\mathcal{L}$  on a smooth manifold maps a vector field and a  $(p, q)$  tensor field to

another  $(p, q)$  tensor field such that:

- (i)  $\mathcal{L}_V f = \nabla f$
- (ii)  $\mathcal{L}_V W = [V, W]$
- (iii)  $\mathcal{L}_V(T + S) = \mathcal{L}_V T + \mathcal{L}_V S$
- (iv)  $\mathcal{L}_V T(\omega, W) = (\mathcal{L}_V T)(\omega, W) + T(\mathcal{L}_V \omega, W) + T(\omega, \mathcal{L}_V W)$  and similarly for any other tensor  $T$
- (v)  $\mathcal{L}_{V+W} T = \mathcal{L}_V T + \mathcal{L}_W T$

It may seem like the Lie derivative is pretty much like a covariant derivative, but with even less structure. Indeed the Lie derivative has no  $C^\infty$  linearity which can be verified by looking at (ii) and evaluating  $[fV, W] \neq f[V, W]$ . We could not have used this Lie derivative in place of the covariant derivative for exactly this reason, we need the linearity to differentiate tensors and talk about parallel transport, geodesics etc...

### Theorem (Lie derivative components)

In a local basis, we have that:

$$(\mathcal{L}_V W)^\mu = V^\nu \frac{\partial W^\mu}{\partial x^\nu} - \frac{\partial V^\mu}{\partial x^\nu} W^\nu \quad (58.2.17)$$

*Proof.* The proof is immediate from the definition of the commutator

$$(\mathcal{L}_V W)f = V(W(f)) - W(V(f)) = V^\nu \frac{\partial}{\partial x_\nu} \left( W^\mu \frac{\partial}{\partial x^\mu} f \right) - W^\mu \frac{\partial}{\partial x^\mu} \left( V^\nu \frac{\partial}{\partial x_\nu} f \right) \quad (58.2.18)$$

$$= \left( V^\nu \frac{\partial W^\mu}{\partial x^\nu} - \frac{\partial V^\mu}{\partial x^\nu} W^\nu \right) \frac{\partial}{\partial x^\mu} f \quad (58.2.19)$$

as desired. ■

Compare this to the definition of the covariant derivative:

$$(\nabla_V W)^\mu = V^\mu \frac{\partial W^\nu}{\partial x^\mu} - \Gamma_{\rho\nu}^\mu V^\nu W^\rho \quad (58.2.20)$$

then we see that while  $\nabla$  does not require information about the vector fields outside of the point it is evaluated at,  $\mathcal{L}$  does due to the partial derivatives.

We can use the Leibniz rule to write:

$$(\mathcal{L}_V T)_\nu^\mu = V^\alpha \frac{\partial}{\partial x^\alpha} T_\nu^\mu + \frac{\partial V^\alpha}{\partial x^\nu} T_\alpha^\mu - \frac{\partial V^\mu}{\partial x^\alpha} T_\nu^\alpha \quad (58.2.21)$$

There is a second, more intuitive definition of the Lie derivative which is equivalent to the first, and can be used in treating symmetries.

# Integration on manifolds

Let us take a smooth manifold  $M$  and consider a function  $f : M \rightarrow \mathbb{R}$ . Working in a chart  $(U, \phi)$  we integrate  $f \circ \phi$  in  $\phi(U)$ , we want to transition to another chart  $(U, \varphi)$  and integrate  $f \circ \phi$  in  $\varphi(U)$ . This situation is summarized in the following commutative diagram.

$$\begin{array}{ccccc}
 & & \varphi(U) & & \\
 & \nearrow \varphi & & \searrow f_\varphi & \\
 \varphi \circ \phi^{-1} & \leftarrow U \xrightarrow{f} & \mathbb{R} & & \\
 & \downarrow \phi & & \nearrow f_\phi & \\
 & \phi(U) & & &
 \end{array}$$

where we defined  $f_\phi = f \circ \phi^{-1}$  and  $f_\varphi = f \circ \varphi^{-1}$ . We can integrate in  $\mathbb{R}^d$  and use the typical rule of change of variables to find:

$$\int_{\varphi(U)} d^n \tilde{x} f_\varphi(\tilde{x}) = \int_{\phi(U)} d^n x \left| \det \left( \partial_a (\varphi \circ \phi^{-1})^b \right)(x) \right| (f_\varphi \circ (\varphi \circ \phi^{-1}))(x) \quad (59.0.1)$$

$$= \int_{\phi(U)} d^n x \left| \det \left( \frac{\partial y^b}{\partial x^a} \right)(\phi^{-1}(x)) \right| f_\phi(x) \quad (59.0.2)$$

$$\neq \int_{\phi(U)} d^n x f_\phi(x) \quad (59.0.3)$$

so if we want to define integration so that it is not chart-dependent then we must insert a factor that takes care of the Jacobian. One object which we know transforms inversely to the Jacobian are differential forms. Let for example  $\alpha$  be a differential form which we expand in some local basis as:

$$\alpha = a(x^1, \dots, x^n) d\varphi^1 \wedge \dots \wedge d\varphi^n \quad (59.0.4)$$

Then if we define

$$\int_U \alpha = \int_{\varphi(U)} a(x^1, \dots, x^n) d\varphi^1 \wedge \dots \wedge d\varphi^n \quad (59.0.5)$$

we find that working in some other chart  $(U, \tilde{\varphi})$ :

$$\int_U \alpha = \int_{\tilde{\varphi}(U)} \tilde{a}(\tilde{x}^1, \dots, \tilde{x}^n) d\tilde{\varphi}^1 \dots d\tilde{\varphi}^n \quad (59.0.6)$$

$$= \int_{\varphi(U)} \det\left(\frac{\partial x^\mu}{\partial \tilde{x}^\nu}\right) a(x^1, \dots, x^n) \det\left(\frac{\partial \tilde{x}^\mu}{\partial x^\nu}\right) d\varphi^1 \dots d\varphi^n \quad (59.0.7)$$

$$= \int_{\varphi(U)} a(x^1, \dots, x^n) d\varphi^1 \dots d\varphi^n \quad (59.0.8)$$

as desired. It is therefore clear that we can use special differential forms, known as volume forms, to clear up the mess from the Jacobians.

### **Definition (Volume form)**

Consider a smooth,  $n$ -dimensional manifold  $M$ . Then a  $(0, n)$ -tensor field  $\Omega$  such that:

- (a)  $\Omega$  does not vanish anywhere
- (b)  $\Omega$  is totally antisymmetric

is known as a **volume form**.

### **Definition (Metric volume form)**

On a metric manifold we can always construct a **metric volume form** in a chart  $(U, \phi)$  from the metric  $g$ :

$$\Omega^\phi = \sqrt{\det(g_{\mu\nu}^\phi)} d\phi^1 \wedge \dots \wedge d\phi^n \quad (59.0.9)$$

*Proof.* Our definition is seemingly chart-dependent so we must make sure that transformation under chart transitions are well-defined. Suppose we have two charts  $(U, \phi)$  and  $(U, \varphi)$  with local bases  $x$  and  $\tilde{x}$  respectively. We then have that:

$$\Omega^\varphi = \sqrt{\det(g_{\mu\nu}^\varphi)} d\varphi^1 \wedge \dots \wedge d\varphi^n \quad (59.0.10)$$

$$= \sqrt{\det(g_{\alpha\beta}^\phi)} \det\left(\frac{\partial x^\alpha}{\partial \tilde{x}^\mu} \frac{\partial x^\beta}{\partial \tilde{x}^\nu}\right) \det\left(\frac{\partial \tilde{x}}{\partial x}\right) d\phi^1 \wedge \dots \wedge d\phi^n \quad (59.0.11)$$

$$= \sqrt{\det(g_{\alpha\beta}^\phi)} \left| \det\left(\frac{\partial x}{\partial \tilde{x}}\right) \right| \det\left(\frac{\partial \tilde{x}}{\partial x}\right) d\phi^1 \wedge \dots \wedge d\phi^n \quad (59.0.12)$$

$$= \sqrt{\det(g_{\alpha\beta}^\phi)} \operatorname{sgn}\left(\det\left(\frac{\partial x}{\partial \tilde{x}}\right)\right) d\phi^1 \wedge \dots \wedge d\phi^n \quad (59.0.13)$$

so this is only equal to  $\Omega^\phi$  if the Jacobian is positive, that is we must impose

$$\operatorname{sgn}\left(\det\left(\frac{\partial x}{\partial \tilde{x}}\right)\right) = 1 \quad (59.0.14)$$

In this case then we do indeed find that  $\Omega^\phi = \Omega^\varphi$  so we may simply write the volume form as  $\Omega$ . ■

### **Definition (Integration on a metric manifold)**

---

Let  $f$  be a function on a chart domain  $U$ . Then:

$$\int_U f \equiv \int_{\phi(U)} d^n x \sqrt{\det(g_{\mu\nu})} f_\phi(x) \quad (59.0.15)$$

## **Part VII**

# **Complex analysis**

# Complex numbers

## 60.1 What are complex numbers

### Definition (*Complex numbers and functions*)

Let  $x, y \in \mathbb{R}$  and let us denote by  $i$  the quantity such that  $i^2 = -1$ . Then we say that  $z = x + iy$  is a **complex number** expressed in **Cartesian form**, with  $\operatorname{Re} z = x$  representing the **real part** of  $z$ , and  $\operatorname{Im} z = y$  representing the **imaginary part** of  $z$ . The set of all complex numbers is  $\mathbb{C}$  and has the structure of a field with the operations:

$$+ : \mathbb{C} \times \mathbb{C} \rightarrow \mathbb{C} \quad (60.1.1)$$

$$(z_1, z_2) \mapsto z_1 + z_2 \quad (60.1.2)$$

and

$$\cdot : \mathbb{C} \times \mathbb{C} \rightarrow \mathbb{C} \quad (60.1.3)$$

$$(z_1, z_2) \mapsto z_1 \cdot z_2 \quad (60.1.4)$$

where if  $z_1 = x_1 + iy_1$  and  $z_2 = x_2 + iy_2$  then:

$$z_1 + z_2 = (x_1 + x_2) + i(y_1 + y_2) \quad (60.1.5)$$

$$z_1 \cdot z_2 = x_1 x_2 - y_1 y_2 + i(x_1 y_2 + x_2 y_1) \quad (60.1.6)$$

### Definition (*Complex conjugate*)

Let  $z = x + iy \in \mathbb{C}$ . Then its **complex conjugate** is defined as  $\bar{z} = x - iy$ .

There are several useful properties that come with the complex conjugate:

### Proposition (*Complex conjugate properties*)

Let  $z_1, z_2 \in \mathbb{C}$ , then:

- (i)  $z_1 + \bar{z}_1 = 2 \operatorname{Re} z_1$
- (ii)  $z_1 - \bar{z}_1 = 2i \operatorname{Im} z_1$
- (iii)  $\overline{(z_1)} = z_1$
- (iv)  $\overline{z_1 \pm z_2} = \bar{z}_1 \pm \bar{z}_2$
- (v)  $\overline{z_1 \cdot z_2} = \bar{z}_1 \cdot \bar{z}_2$
- (vi)  $\overline{\left(\frac{z_1}{z_2}\right)} = \frac{\bar{z}_1}{\bar{z}_2}$

Since the arithmetic of real numbers applies equally well to complex numbers, it is not so surprising that the Binomial theorem holds in  $\mathbb{C}$ .

**Theorem (Complex binomial theorem)**

Let  $z_1, z_2 \in \mathbb{C}$  and  $n \in \mathbb{N}$ , then:

$$(z_1 + z_2)^n = \sum_{k=0}^n \binom{n}{k} (z_1)^{n-k} z_2^k \quad (60.1.7)$$

*Proof.* Identical to the  $\mathbb{R}$  counterpart. ■

Similarly, the following result from real analysis still holds:

**Theorem (Geometric series)**

Let  $z_1, z_2 \in \mathbb{C}$  and  $n \in \mathbb{N}$ , then:

$$z_1^n - z_2^n = (z_1 - z_2)(z_1^{n-1} + z_1^{n-2} z_2 + \dots + z_2^{n-1}) \quad (60.1.8)$$

Due to the isomorphism between  $\mathbb{C}$  and  $\mathbb{R}^2$ , we can represent complex numbers in a two-dimensional plane, known as the **complex plane**, with the  $x$ -axis replaced by the real axis, and the  $y$ -axis replaced by the imaginary axis. Recall also that in real analysis, the modulus of a real number gave the distance between the origin and this point. Similarly, we may define the complex counterpart as follows

**Definition (Complex modulus)**

Let  $z = x + iy \in \mathbb{C}$  be a complex number. Then the **complex modulus** of  $z$  is defined as:

$$|z| = \sqrt{x^2 + y^2} \quad (60.1.9)$$

We list some important properties of the modulus.

**Proposition (Complex modulus properties)**

Let  $z \in \mathbb{C}$ , then:

- (i)  $|z| \geq 0$ , with equality only iff  $z = 0$
- (ii)  $|\bar{z}| = |z|$  and  $|-z| = |z|$
- (iii)  $|z|^2 = z\bar{z}$
- (iv)  $|z_1 z_2| = |z_1||z_2|$
- (v)  $\left| \frac{z_1}{z_2} \right| = \frac{|z_1|}{|z_2|}$

Note that the modulus of a complex number is not enough to fully specify it (unless it is zero). Indeed the set of all numbers in  $\mathbb{C}$  with the same modulus  $r$  form a circle of radius  $r$ . It follows that the second piece of information needed to specify a complex number is the angle it makes with the real axis, known as the argument.

**Definition (Argument)**

Let  $z = x + iy \in \mathbb{C}$ . Then we define its **argument**  $\theta \in \mathbb{R}$  so that:

$$\sin \theta = \frac{y}{r}, \cos \theta = \frac{x}{r} \quad (60.1.10)$$

where  $r = |z|$ .

Recall that by convention, positive angles are measured anti-clockwise. Also, note that a given number there has an infinite number of arguments. Indeed if  $\theta$  is an argument then so is  $\theta + 2k\pi$  for any  $k \in \mathbb{Z}$ .

### Definition (Principal argument)

The argument  $\theta$  of  $z \in \mathbb{C}$  satisfying  $-\pi < z \leq \pi$  is known as the **principal argument**  $\arg z$ .

Since the modulus and argument of  $z$  are enough to fully specify it, we may define a new form other than the Cartesian form in which complex numbers may be expressed.

### Definition (Polar form)

Let  $z$  be a complex number with  $|z| = r, \arg z = \theta$ . Then  $z$  may be expressed in the form

$$z = r(\cos \theta + i \sin \theta) = re^{i\theta} \quad (60.1.11)$$

known as its **polar form**.

This polar form allows us to express the product/quotient of two complex numbers with the same ease we express the sum/difference of two complex numbers in cartesian form. Let  $z_1 = r_1 e^{i\theta_1}, z_2 = r_2 e^{i\theta_2}$ , then

$$z_1 z_2 = r_1 r_2 e^{i(\theta_1 + \theta_2)} = r_1 r_2 (\cos(\theta_1 + \theta_2) + i \sin(\theta_1 + \theta_2)) \quad (60.1.12)$$

$$\frac{z_1}{z_2} = \frac{r_1}{r_2} e^{i(\theta_1 - \theta_2)} = \frac{r_1}{r_2} (\cos(\theta_1 - \theta_2) + i \sin(\theta_1 - \theta_2)) \quad (60.1.13)$$

This provides an alternative derivation of the identities  $|z_1 z_2| = |z_1| |z_2|$ ,  $\left| \frac{z_1}{z_2} \right| = \frac{|z_1|}{|z_2|}$ . Also, note that while the sum of an argument of  $z_1$  and  $z_2$  is an argument of  $z_1 z_2$ , the same does not hold for the principal argument. Indeed the sum of two principal arguments could be larger than  $\pi$ , and thus not qualify. Instead we have the weaker result that

$$\arg(z_1 z_2) = \arg z_1 + \arg z_2 + 2n\pi, n = -1, 0, 1 \quad (60.1.14)$$

In the special case where  $z_1 = 1 = e^{i\cdot 0}$  then letting  $z_2 = z = re^{i\theta}$  we find that

$$\frac{1}{z} = \frac{1}{r} e^{-i\theta} = \overline{\left( \frac{1}{r} e^{i\theta} \right)} = \frac{\bar{z}}{|z|^2} \quad (60.1.15)$$

from which it follows that:

$$\arg \bar{z} = \arg z^{-1} = -\arg z \quad (60.1.16)$$

**Theorem (De Moivre's theorem)**

It holds that if  $\theta \in \mathbb{R}$  and  $n \in \mathbb{N}$  then:

$$(\cos \theta + i \sin \theta)^n = \cos(n\theta) + i \sin(n\theta) \quad (60.1.17)$$

*Proof.* We notice that  $\cos \theta + i \sin \theta = e^{i\theta}$  so:

$$(\cos \theta + i \sin \theta)^n = (e^{i\theta})^n = e^{i(n\theta)} = \cos(n\theta) + i \sin(n\theta) \quad (60.1.18)$$

as desired. ■

**Theorem (Complex roots)**

Let  $w = \rho(\cos \phi + i \sin \phi)$  be a non-zero complex number, and define the solutions to  $z^n = w$  to be the  **$n$ th roots** of  $w$ . These roots are given by:

$$z_m = \rho^{1/n} \left[ \cos \left( \frac{\phi}{n} + \frac{2m\pi}{n} \right) + i \sin \left( \frac{\phi}{n} + \frac{2m\pi}{n} \right) \right], \quad m \in \mathbb{Z}_n \quad (60.1.19)$$

*Proof.* We need to solve  $z^n = w$ , so let  $z = r(\cos \theta + i \sin \theta)$ . Using DeMoivre's theorem yields

$$z^n = r^n(\cos(n\theta) + i \sin(n\theta)) = \rho(\cos \phi + i \sin \phi) \quad (60.1.20)$$

It is then clear that  $r^n = \rho$  by equating moduli. Moreover, arguments can only differ by a multiple of  $2\pi$  so  $n\theta = \phi + 2m\pi$  where  $m \in \mathbb{Z}$ . Consequently solutions are of the form:

$$z_m = \rho^{1/n} \left[ \cos \left( \frac{\phi}{n} + \frac{2m\pi}{n} \right) + i \sin \left( \frac{\phi}{n} + \frac{2m\pi}{n} \right) \right], \quad m \in \mathbb{Z}_n \quad (60.1.21)$$

Note that there are not infinitely many solutions, one for each  $m$ . However, if  $m_1$  and  $m_2$  differ by an integer multiple  $kn$  then:

$$\frac{\phi}{n} + \frac{2m_2\pi}{n} = \frac{\phi}{n} + \frac{2m_1\pi}{n} + 2k\pi \quad (60.1.22)$$

so  $z_{m_1} = z_{m_2}$ . Consequently the only distinct solutions can be formed from  $k = 0, 1, 2, \dots, n - 1$ . Note that the angle between  $z_k$  and  $z_{k+1}$  is the same for all  $k$ , and their moduli are all the same. Consequently these roots lie on a circle and form the vertices of a regular  $n$ -polygon. ■

Complex inequalities work the same way as real inequalities. We repeat the standard rules for rearranging inequalities for completeness:

- (i)  $a < b \iff 0 < b - a$
- (ii)  $a < b \iff a + c < b + c$
- (iii)  $a < b \iff ac < bc$  for  $c > 0$  and  $a < b \iff ac > bc$  for  $c < 0$
- (iv)  $a > b \iff \frac{1}{a} < \frac{1}{b}$  for  $a, b > 0$
- (v)  $a < b \iff a^p < b^p$  if  $a, b \geq 0$  and  $p > 0$
- (vi)  $|a| < b \iff -b < a < b$

- (vii)  $a < b$  and  $b < c \implies a < c$
- (viii)  $a < b$  and  $c < d \implies a + c < b + d$
- (ix)  $a < b$  and  $c < d \implies ac < bd$  if  $a, c \geq 0$

The triangle inequalities are also quite useful:

**Proposition (Triangle inequalities)**

If  $z_1, z_2 \in \mathbb{C}$  then:

$$|z_1 + z_2| \leq |z_1| + |z_2| \quad (60.1.23)$$

$$||z_1| - |z_2|| \leq |z_1 - z_2| \quad (60.1.24)$$

*Proof.*

$$|z_1 + z_2|^2 = (z_1 + z_2)\overline{(z_1 + z_2)} \quad (60.1.25)$$

$$= (z_1 + z_2)(\overline{z_1} + \overline{z_2}) \quad (60.1.26)$$

$$= |z_1|^2 + z_1\overline{z_2} + \overline{z_1}z_2 + |z_2|^2 \quad (60.1.27)$$

$$= |z_1|^2 + 2\operatorname{Re}(z_1\overline{z_2}) + |z_2|^2 \quad (60.1.28)$$

$$\leq |z_1|^2 + 2|z_2\overline{z_2}| + |z_2|^2 \quad (60.1.29)$$

$$= (|z_1| + |z_2|)^2 \quad (60.1.30)$$

Applying one of the rules for rearranging inequalities gives (60.1.23).

Also, note that:

$$|z_1| \leq |z_1 - z_2| + |z_2| \implies |z_1| - |z_2| \leq |z_1 - z_2| \quad (60.1.31)$$

and similarly:

$$|z_2| \leq |z_2 - z_1| + |z_1| \implies |z_2| - |z_1| \leq |z_2 - z_1| \quad (60.1.32)$$

Consequently:

$$-|z_1 - z_2| \leq |z_2| - |z_1| \leq |z_1 - z_2| \implies ||z_2| - |z_1|| \leq |z_1 - z_2| \quad (60.1.33)$$

■

The geometrical interpretation of the first triangle inequality is that the diagonal of a parallelogram cannot be longer than the sum of the lengths of its two adjacent sides. The geometrical interpretation of the second triangle inequality is that the distance between two points on two concentric circles cannot be larger than the difference in the circles' radii.

As a consequence of the fact that  $\mathbb{C}$  is not a complex field, we cannot write an inequality between two complex numbers  $z_1 \leq z_2$ , it simply does not make sense. Nevertheless one can use the modulus and argument of complex numbers to construct meaningful inequalities with geometric interpretations.

**Definition (Subsets of  $\mathbb{C}$ )**

A **line** is a set of the type  $\{z : a \operatorname{Re} z + b \operatorname{Im} z > c\}$  where  $a, b, c \in \mathbb{R}$  with the first two not both equal to zero. An **open half-plane** is a set of the type  $\{z : a \operatorname{Re} z + b \operatorname{Im} z > c\}$  and a **closed half-plane** is a set of the type  $\{z : a \operatorname{Re} z + b \operatorname{Im} z \geq c\}$ .

The **circle** centered at  $\alpha \in \mathbb{C}$  with radius  $r > 0$  is the set  $\{z : |z - \alpha| = r\}$ . Analogously, an **open disc** is a set of the type  $\{z : |z - \alpha| < r\}$  and a **closed disc** is a set of the type  $\{z : |z - \alpha| \leq r\}$ .

A **half-line** is a set of the form  $\{z : \arg(z - \alpha) = \theta\}$  where  $\alpha \in \mathbb{C}$  and  $-\pi < \theta \leq \pi$ . Hence an **open sector** is a set of the form  $\{z : a < \arg(z - \alpha) < b\}$  or  $\{z : \arg(z - \alpha) < a \vee \arg(z - \alpha) > b\}$ .

## 60.2 Complex functions

**Definition (Complex functions)**

Let  $A, B \subseteq \mathbb{C}$ . Then we define the map:

$$f : A \rightarrow B \quad (60.2.1)$$

$$z \mapsto f(z) \quad (60.2.2)$$

to be a **complex function**. The **domain** of  $f$  is  $A$ , the **co-domain** of  $f$  is  $B$ , and the image of  $f$  is:

$$f(A) \equiv \{f(z) : z \in A\} \quad (60.2.3)$$

If  $f(A) = B$  then  $f$  is **surjective**, while if  $f(z_1) = f(z_2) \implies z_1 = z_2$  then  $f$  is **injective**.

**Example.** Consider for example  $f(z) = \frac{3z+1}{z+i}$ . It is implied that the domain of  $f$  is defined to be the subset of  $\mathbb{C}$  where its rule is applicable. In this case, the domain of  $f$  is thus  $A = \mathbb{C} \setminus \{-i\}$  and the codomain is  $\mathbb{C}$ . Consequently:

$$f(A) = \left\{ w = \frac{3z+1}{z+i} : z \in \mathbb{C} \setminus \{-i\} \right\} \quad (60.2.4)$$

Now note that:

$$w = \frac{3z+1}{z+i} \implies w(z^2 + 1) = 3z^2 + (1 - 3i)z - i \quad (60.2.5)$$

$$\implies (w - 3)z^2 - (1 - 3i)z + (w + i) = 0 \quad (60.2.6)$$

which has a well defined solution:

$$z = \frac{(1 - 3i) \pm \sqrt{(1 - 3i)^2 - 4(w - 3)(w + i)}}{2(w - 3)} = \quad (60.2.7)$$

as long as  $w \neq 3$ . Thus  $f(A) = \mathbb{C} \setminus \{3\}$ , and due to this  $f$  is not surjective. However it is

injective, since:

$$f(z_1) = f(z_2) \Rightarrow \frac{3z_1 + 1}{z_1 + i} = \frac{3z_2 + 1}{z_2 + i} \Rightarrow 3z_1 z_2 + i + 3iz_1 + z_2 = 3z_2 z_1 + i + 3iz_2 + z_1 \quad (60.2.8)$$

$$\Rightarrow 3i(z_1 - z_2) = z_1 - z_2 \Rightarrow z_1 = z_2 \quad (60.2.9)$$

as desired.  $\blacktriangleleft$

The standard operations of sums, products/quotients and compositions can be performed on complex functions just like with real functions.

**Example.** Let  $f(z) = \frac{1}{z}$  with domain  $z \in \mathbb{C} \setminus \{0\}$  and  $g(z) = \frac{z+3i}{z^2-z}$  with domain  $z \in \mathbb{C} \setminus \{0, 1\}$ . We see that the complex functions:

$$(f+g)(z) = \frac{2z-1+3i}{z^2-z}, \quad (fg)(z) = \frac{z+3i}{z(z^2-z)} \quad (60.2.10)$$

have domain  $\mathbb{C} \setminus \{0, 1\}$ . Similarly, we have that:

$$\frac{f}{g} = \frac{z^2-z}{z^2+3iz} = \frac{z-1}{z+3i} \quad (60.2.11)$$

has domain  $\mathbb{C} \setminus \{0, 1, -3i\}$ . Also, we can compose these functions to obtain:

$$(g \circ f)(z) = \frac{\frac{1}{z} + 3i}{\frac{1}{z^2} - \frac{1}{z}} = \frac{3iz^2 + z}{1 - z} \quad (60.2.12)$$

Its domain is:

$$\{z \in \mathbb{C} \setminus \{0\} : f(z) \in \mathbb{C} \setminus \{0, 1\}\} = \mathbb{C} \setminus \{0, 1\} \quad (60.2.13)$$

Similarly, we obtain that

$$(f \circ g)(z) = \frac{z^2-z}{z+3i} \quad (60.2.14)$$

whose domain is

$$\{z \in \mathbb{C} \setminus \{0, 1\} : g(z) \in \mathbb{C} \setminus \{0\}\} = \mathbb{C} \setminus \{0, 1, -3i\} \quad (60.2.15)$$

$\blacktriangleleft$

### Definition (*Inverse of complex function*)

Let  $f : A \rightarrow B$  be a complex injective function. Then the inverse function of  $f$  is defined as:

$$f^{-1} : B \rightarrow B \quad (60.2.16)$$

$$z \mapsto f^{-1}(z) \quad (60.2.17)$$

such that:

$$f^{-1}(w) = z \iff f(w) = z \quad (60.2.18)$$

It follows that the existence of an inverse function boils down to whether or not one can find for every  $w \in B$  a  $z \in A$  such that  $f(z) = w$ . In many cases this cannot be done for a given  $B$ , but can be done if we somehow restrict  $B$  to some new codomain  $B'$ , redefining our function to  $f|_B$ . Thus, a function is invertible if for all  $w \in f(A)$  there exists a unique  $z \in A$  such that  $f(z) = w$ .

**Example.** Let  $A = \{0\} \cup \{z : -\pi/3 < \arg z \leq \pi/3\}$  be the domain of  $f(z) = z^3$ . This function has an inverse. Indeed, note that the images set of  $A$  is:

$$f(A) = \{z^3 : z \in \{0\} \cup \{z : -\pi/3 < \arg z \leq \pi/3\}\} \quad (60.2.19)$$

$$= \{z^3 : z = re^{i\theta}, -\pi/3 < \theta \leq \pi/3, r \geq 0\} \quad (60.2.20)$$

$$= \{r^3 e^{3i\theta} : -\pi/3 < \theta \leq \pi/3, r \geq 0\} \quad (60.2.21)$$

$$= \mathbb{C} \quad (60.2.22)$$

Hence, given  $w \in \mathbb{C}$  we need to find a unique  $z \in A$  such that  $z^3 = w$ . If  $w = 0$  then clearly  $f(0) = 0 \implies f^{-1}(0) = 0$ . If instead  $w \neq 0$  then:

$$z^3 = w = re^{i\theta}, \quad r > 0, -\pi < \theta \leq \pi \quad (60.2.23)$$

but using De Moivre's theorem we have that:

$$z = r^{1/3} e^{i(\theta+2k\pi)/3}, \quad \forall k = 0, 1, 2 \quad (60.2.24)$$

Since we need  $z \in A$ , the only possible choice is  $k = 0$ . Consequently, the inverse of  $f$  is:

$$f^{-1}(0) = 0, \quad f^{-1}(w) = r^{1/3} e^{i\theta/3}, \quad w \neq 0 \quad (60.2.25)$$

with domain  $\mathbb{C}$  and co-domain  $A$ . ◀

### Definition (Real-valued functions)

Given a complex function  $f$ , we say that it is **real valued** if  $f(A) \subseteq \mathbb{R}$  and a **real function** if  $A \subseteq \mathbb{R}$  and  $B \subseteq \mathbb{R}$ .

Note that any complex function can be decomposed into real functions. For example consider  $f(z) = \frac{1}{z}$ . Then:

$$(\operatorname{Re} f)(z) = \operatorname{Re}(f(z)) = \operatorname{Re}\left(\frac{1}{z}\right) = \frac{\operatorname{Re} z}{|z|^2}, \quad (z \in \mathbb{C} \setminus \{0\}) \quad (60.2.26)$$

and similarly:

$$(\operatorname{Im} f)(z) = \operatorname{Im}(f(z)) = \operatorname{Im}\left(\frac{1}{z}\right) = -\frac{\operatorname{Im} z}{|z|^2}, \quad (z \in \mathbb{C} \setminus \{0\}) \quad (60.2.27)$$

so that:

$$f(z) = (\operatorname{Re} f)(z) + i(\operatorname{Im} f)(z) = \frac{\operatorname{Re} z - i \operatorname{Im} z}{|z|^2} = \frac{z^*}{|z|^2} \quad (z \in \mathbb{C} \setminus \{0\}) \quad (60.2.28)$$

as expected.

### Definition (Complex function parametrisation)

A **path**  $\Gamma \subseteq \mathbb{C}$  is the image set of a continuous function  $\gamma : I \rightarrow \mathbb{C}$  where  $I \subseteq \mathbb{R}$ , known as

the parametrisation of  $\Gamma$ .

Consider for example the unit circle  $\Gamma$  in the complex plane. It can be parametrised by:

$$\gamma(t) = \cos t + i \sin t, t \in [0, 2\pi] \quad (60.2.29)$$

which traverses the circle anti-clockwise. Consider the effect of letting  $z$  vary on  $\Gamma$  following the parametrisation  $\gamma$ , and seeing the effects on  $f(z) = \frac{1}{z}$ . We saw that letting  $z = x + iy$  and  $w = f(z) = u + iv$  then

$$u = \frac{x}{x^2 + y^2}, v = -\frac{y}{x^2 + y^2} \quad (60.2.30)$$

Thus, letting  $x = \cos t$  and  $y = \sin t$  then we find that:

$$u = \cos t, v = -\sin t, t \in [0, 2\pi] \implies z = u + iv = \cos(-t) + i \sin(-t), t \in [0, 2\pi] \quad (60.2.31)$$

Consequently, the unit circle in the  $w$  plane will be traversed clockwise as  $z$  traverses the unit circle anti-clockwise.

## 60.3 Mappings under complex functions

### 60.4 Special complex functions

We now extend notions of exponential, logarithmic, trigonometric and hyperbolic functions to the complex domain.

**Definition (Complex exponential)**

Let  $z \in \mathbb{C}$ . Then we have define:

$$e^z = e^{\operatorname{Re} z}(\cos(\operatorname{Im} z) + i \sin(\operatorname{Im} z)) \quad (60.4.1)$$

**Proposition (Exponential identities)**

Let  $z_1, z_2 \in \mathbb{C}$ , then:

- (i)  $e^{z_1+z_2} = e^{z_1}e^{z_2}$
- (ii)  $|e^z| = e^{\operatorname{Re} z}$
- (iii)  $e^{-z} = \frac{1}{e^z}$
- (iv)  $e^{z+2\pi i} = e^z$

*Proof.* (i) We find that

■

**Example.** Note that for  $z \in \mathbb{C}$  then

$$|e^z| = |e^{\operatorname{Re} z}(\cos(\operatorname{Im} z) + i \sin(\operatorname{Im} z))| = e^{\operatorname{Re} z} \quad (60.4.2)$$

and since  $\operatorname{Re} z \leq |z|$  we have that  $|e^z| \leq e^{|z|}$  by the monotonicity of the real exponential function.

◀

Interestingly, we have that  $e^{z+2in\pi} = e^z$  for all  $n \in \mathbb{Z}$ , so going up the complex plane in steps of  $2i\pi$  brings you back to your starting point under complex exponentiation. This suggests that we take the line  $\Gamma : \gamma(t) = a + it$ , and look at how it is mapped by  $w = f(z) = e^z$ . We have that if  $z = x + iy$  then:

$$w = u + iv = e^x(\cos y + i \sin y) \implies u = e^x \cos y, v = e^x \sin y \quad (60.4.3)$$

thus:

$$u = e^a \cos t, v = e^a \sin t \quad (60.4.4)$$

so that  $u^2 + v^2 = e^{2a}$ . As expected, the straight line gets mapped to a circle. Similarly, taking the path  $\Gamma : \gamma(t) = t + ib$  then:

$$w = u + iv = e^t(\cos b + i \sin b) \implies u = e^t \cos b, v = e^t \sin b \quad (60.4.5)$$

thus:

$$u = e^t \cos b, v = e^t \sin b \quad (60.4.6)$$

which is a half-line with  $\arg z = b$ . Hence if we take a grid of  $x = \text{cnst}$  and  $y = \text{cnst}$  then it will get mapped to outward radial lines with concentric circles. A rectangular region gets mapped to a quarter annulus. Points in the left half-plane with  $\operatorname{Re} z < 0$  get mapped to points inside the unit circle while points in the right halfplane get mapped outside the unit circle as a result of  $|e^z| = e^{\operatorname{Re} z}$ . Finally, note that the image of the strip  $\{x + iy : -\pi < y \leq \pi\}$  under  $f(z) = e^z$  is  $\mathbb{C} \setminus \{0\}$ .

We have seen that  $f(z) = e^z$  is not injective on  $\mathbb{C}$ , but we can make it injective by restricting its domain, for example, to a strip of width  $2\pi$  (with one edge removed), such as:

$$A = \{x + iy : -\pi < y \leq \pi\} \quad (60.4.7)$$

Then we can show that  $f(z) = e^z$  has an inverse on this domain.

**Example.** We firstly look at the image set of  $f(z)$ :

$$f(A) = \{w = e^{x+iy} : -\pi < y \leq \pi\} = \{w = e^x e^{iy} : -\pi < y \leq \pi\} = \{w = \rho e^{i\theta} : \rho > 0, -\pi < \theta \leq \pi\} \quad (60.4.8)$$

$$= \mathbb{C} \setminus \{0\} \quad (60.4.9)$$

Now, given any  $w \in \mathbb{C} \setminus \{0\}$  then  $w = \rho e^{i\theta}$  with  $\rho > 0$  and  $-\pi < \theta \leq \pi$ . We claim that  $z = \log(\rho) + i\theta \in A$  is the required inverse (well defined since  $\rho \neq 0$ ):

$$f(z) = e^z = e^{\log(\rho)} e^{i\theta} = \rho e^{i\theta} = w \quad (60.4.10)$$

as desired. We could have also chosen  $z = \log(\rho) + i\theta + 2in\pi$  but this would not be in  $A$ . Consequently,  $f$  is an injective function with image set  $\mathbb{C} \setminus \{0\}$  and inverse:

$$f^{-1}(z) = \log|z| + i \arg z \quad (60.4.11)$$

with domain  $\mathbb{C} \setminus \{0\}$ . ◀

Much like how a complex number can have infinitely many arguments, its logarithm is also multiply valued. However we can define the principal logarithm by restricting the argument of  $w$  to its principal argument.

**Definition (Principal logarithm)**

Let  $z \in \mathbb{C} \setminus \{0\}$ . The **principal logarithm** of  $z$  is defined as:

$$\log(z) = \log|z| + i \arg z \quad (60.4.12)$$

We collect some useful properties of the principal logarithm below:

**Proposition (Logarithm identities)**

For  $\arg z_1, \arg z_2 \in (-\pi/2, \pi/2]$  then:

$$\log(z_1 z_2) = \log z_1 + \log z_2 \quad (60.4.13)$$

and for  $\arg z \in (-\pi, \pi)$  then:

$$\log\left(\frac{1}{z}\right) = -\log z \quad (60.4.14)$$

*Proof.* Suppose  $\arg z_1, \arg z_2 \in (-\pi/2, \pi/2]$  so that  $\arg(z_1 z_2) = \arg z_1 + \arg z_2 \in (-\pi, \pi]$ . Then:

$$\log(z_1 z_2) = \log|z_1 z_2| + i \arg(z_1 z_2) = \log|z_1| + \log|z_2| + i \arg z_1 + i \arg z_2 = \log(z_1) + \log(z_2) \quad (60.4.15)$$

Similarly, if  $\arg z \in (-\pi, \pi)$  then  $\arg \frac{1}{z} = -\arg z$  so that:

$$\log\left(\frac{1}{z}\right) = \log\left|\frac{1}{z}\right| + i \arg \frac{1}{z} = -\log|z| - i \arg z = -\log z \quad (60.4.16)$$

as desired. ■

We can view the geometric effect of the complex logarithm in the same (or inverse) way for the complex exponential. Going back to the complex exponential, we can invert Euler's identity to express  $\cos z$  and  $\sin z$  as follows:

$$\begin{cases} e^{i\theta} = \cos \theta + i \sin \theta \\ e^{-i\theta} = \cos \theta - i \sin \theta \end{cases} \implies \begin{cases} \cos \theta = \frac{1}{2}(e^{i\theta} + e^{-i\theta}) \\ \sin \theta = \frac{1}{2i}(e^{i\theta} - e^{-i\theta}) \end{cases} \quad (60.4.17)$$

This motivates us to make the following definitions:

**Definition (Complex trigonometric functions I)**

We define for all  $z \in \mathbb{C}$ :

$$\cos z = \frac{1}{2}(e^{iz} + e^{-iz}), \sin z = \frac{1}{2i}(e^{iz} - e^{-iz}) \quad (60.4.18)$$

Interestingly, the zeros of the complex trigonometric functions are the same as their real counterparts. Indeed:

$$\sin z = 0 \implies e^{iz} = e^{-iz} \implies e^{2iz} = 1 = e^0 \implies 2z = 2n\pi \implies z = n\pi \forall n \in \mathbb{Z} \quad (60.4.19)$$

and similarly:

$$\cos z = 0 \implies e^{iz} = -e^{-iz} \implies e^{2iz} = -1 = e^{i\pi} \quad (60.4.20)$$

$$\implies 2z = \pi + 2n\pi \implies z = \left(n + \frac{1}{2}\right)\pi, \forall n \in \mathbb{Z} \quad (60.4.21)$$

Thus, we may define the tangent, secant, cotangent and cosecant as follows:

### Definition (Complex trigonometric functions II)

We define for all  $z \in \mathbb{C} \setminus \{z = (n + 1/2)\pi : n \in \mathbb{Z}\}$ :

$$\tan z = \frac{\sin z}{\cos z}, \sec z = \frac{1}{\cos z} \quad (60.4.22)$$

and for all  $z \in \mathbb{C} \setminus \{z = n\pi : n \in \mathbb{Z}\}$ :

$$\cot z = \frac{\cos z}{\sin z}, \csc z = \frac{1}{\sin z} \quad (60.4.23)$$

### Proposition (Trigonometric identities)

(i) Addition:

$$\sin(z_1 + z_2) = \sin z_1 \cos z_2 + \cos z_1 \sin z_2 \quad (60.4.24)$$

$$\cos(z_1 + z_2) = \cos z_1 \cos z_2 - \sin z_1 \sin z_2 \quad (60.4.25)$$

$$\tan(z_1 + z_2) = \frac{\tanh z_1 + \tanh z_2}{1 - \tanh z_1 \tanh z_2} \quad (60.4.26)$$

(ii) Squares:

$$\cos^2 z + \sin^2 z = 1 \quad (60.4.27)$$

$$\sec^2 z = 1 + \tan^2 z \quad (60.4.28)$$

$$\csc^2 z = 1 + \cot^2 z \quad (60.4.29)$$

*Proof.* The proof is identical to the real case. ■

Suppose we remove the  $i$  in the exponentials defining  $\sin z$  and  $\cos z$ . We know from real analysis that the result are hyperbolic functions:

$$\sinh x = \frac{1}{2}(e^x - e^{-x}), \cosh x = \frac{1}{2}(e^x + e^{-x}) \quad (60.4.30)$$

Thus, extending these definitions to  $\mathbb{C}$ , we get:

### Definition (Complex hyperbolic functions)

We define for all  $z \in \mathbb{C}$ :

$$\cosh z = \frac{1}{2}(e^z + e^{-z}), \sinh z = \frac{1}{2}(e^z - e^{-z}) \quad (60.4.31)$$

We define for all  $z \in \mathbb{C} \setminus \{z = i(n + 1/2)\pi : n \in \mathbb{Z}\}$ :

$$\tanh z = \frac{\sinh z}{\cosh z}, \quad \operatorname{sech} z = \frac{1}{\cosh z} \quad (60.4.32)$$

and for all  $z \in \mathbb{C} \setminus \{z = in\pi : n \in \mathbb{Z}\}$ :

$$\coth z = \frac{\cosh z}{\sinh z}, \quad \operatorname{csch} z = \frac{1}{\sinh z} \quad (60.4.33)$$

### Proposition (*Trigonometric functions are hyperbolic*)

For all  $z \in \mathbb{C}$ :

$$\cosh(iz) = \cos z, \quad \sinh(iz) = i \sin z \quad (60.4.34)$$

*Proof.* We find:

$$\cosh(iz) = \frac{1}{2}(e^{iz} + e^{-iz}) = \cos z \quad (60.4.35)$$

$$\sinh(iz) = i \frac{1}{2i}(e^{iz} - e^{-iz}) = i \sin z \quad (60.4.36)$$

as desired. ■

### Proposition (*Hyperbolic identities*)

(i) Addition:

$$\sinh(z_1 + z_2) = \sinh z_1 \cosh z_2 + \cosh z_1 \sinh z_2 \quad (60.4.37)$$

$$\cosh(z_1 + z_2) = \cosh z_1 \cosh z_2 + \sinh z_1 \sinh z_2 \quad (60.4.38)$$

$$\tanh(z_1 + z_2) = \frac{\tanh z_1 + \tanh z_2}{1 + \tanh z_1 \tanh z_2} \quad (60.4.39)$$

(ii) Squares:

$$\cosh^2 z - \sinh^2 z = 1 \quad (60.4.40)$$

$$\operatorname{sech}^2 z = 1 - \tanh^2 z \quad (60.4.41)$$

$$\operatorname{csch}^2 z = \coth^2 z - 1 \quad (60.4.42)$$

*Proof.* Can be found from the trigonometric identities using the previous proposition's correspondence between  $\sin z, \sinh z$  and  $\cos z, \cosh z$ . ■

# Continuity of complex functions

## 61.1 Complex sequences

We begin our study of continuity by extending some results from the rReal analysis of sequences to  $\mathbb{C}$ .

### **Definition (Convergence in $\mathbb{C}$ )**

The complex sequence  $(z_n)$  **converges** to  $\alpha$  if  $\forall \epsilon > 0$  there exists an integer  $N$  such that

$$|z_n - \alpha| < \epsilon, \forall n > N \quad (61.1.1)$$

If the sequence converges to zero then it is a **null sequence**. Equivalently,  $(z_n)$  converges to  $\alpha$  if  $(z_n - \alpha)$  is null. If a sequence is not convergent then it is **divergent**.

The geometrical intuition behind the epsilon-delta definition is clear. If  $(z_n)$  converges to  $\alpha$  then given a circle centered at  $\alpha$  with finite radius, all terms after a given number will fall within this circle. This implies that every convergent sequence is bounded.

### **Proposition (Bounded $\iff$ convergent)**

Every convergent sequence is bounded.

*Proof.* Let  $(z_n)$  converge to  $\alpha$ . Applying the triangle inequality to the epsilon-delta definition with  $\epsilon = 1$  we get

$$|z_n - \alpha| \leq |z_n| - |\alpha| < 1 \implies |z_n| < 1 + |\alpha|, \forall n > N \quad (61.1.2)$$

Choosing  $M = \max\{z_1, z_2, \dots, z_N, 1 + |\alpha|\}$  it is clear that  $|z_n| < M$  for all  $n$ . ■

**Example.** We want to prove that the sequence  $(z_n)$  with  $z_n = \frac{1+i}{n}$  is convergent. We claim that it converges to 0. Indeed let  $\epsilon > 0$ , we need to find an integer  $N$  such that

$$\left| \frac{1+i}{n} \right| < \epsilon \implies \frac{\sqrt{2}}{n} < \epsilon, \forall n > N \quad (61.1.3)$$

which is clearly satisfied if  $N > \frac{\sqrt{2}}{\epsilon}$ . Thus  $(z_n)$  is a null sequence. ◀

As usual most sequence converges theorems from real analysis still apply. One important example is the Squeeze rule.

**Theorem (Squeeze rule)**

If  $(a_n)$  is a null non-negative sequence and if  $(z_n)$  is a sequence such that

$$|z_n| \leq a_n, \quad n = 1, 2, \dots \quad (61.1.4)$$

then  $(z_n)$  is a null sequence.

**Example.** Let us prove that  $(z_n)$  with  $z_n = \left(\frac{i}{2}\right)^n$  is a null sequence. We have that

$$|z_n| = \frac{1}{2^n} \leq \frac{1}{n} \quad (61.1.5)$$

but  $\left(\frac{1}{n}\right)$  is a null sequence. Hence by the Squeeze rule  $(z_n)$  is null.  $\blacktriangleleft$

**Theorem (Standard null sequences)**

The following sequences are null

- (a)  $(\frac{1}{n^p})$  for  $p > 0$
- (b)  $\alpha^n$  for  $|\alpha| < 1$

**Example.** Let us find the limit of  $(z_n)$  with  $z_N = \frac{(3+i)^n + (2+2i)^n}{(1+2i)^n + 2(3+i)^n}$ . We can guess that the dominant term will be  $(3+i)^n$  so we divide everything by it:

$$z_n = \frac{1 + (2+2i)^n/(3+i)^n}{2 + (1+2i)^n/(3+i)^n} \quad (61.1.6)$$

We have that  $\frac{(2+2i)^n}{(3+i)^n}$  is a null sequence since

$$\left| \frac{(2+2i)^n}{(3+i)^n} \right| = \left( \sqrt{\frac{4}{5}} \right)^n \quad (61.1.7)$$

Similarly  $\frac{(1+2i)^n}{(3+i)^n}$  is a null sequence since

$$\left| \frac{(1+2i)^n}{(3+i)^n} \right| = \frac{1}{\sqrt{2}^n} \quad (61.1.8)$$

Consequently we find that

$$\lim_{n \rightarrow \infty} z_n = \frac{1}{2} \quad (61.1.9)$$

so  $z_n$  converges to  $\frac{1}{2}$ .  $\blacktriangleleft$

**Proposition (Limit properties)**

If  $\lim_{n \rightarrow \infty} z_n = \alpha$  then

- (a)  $\lim_{n \rightarrow \infty} |z_n| = |\alpha|$
- (b)  $\lim_{n \rightarrow \infty} \overline{z_n} = \overline{\alpha}$

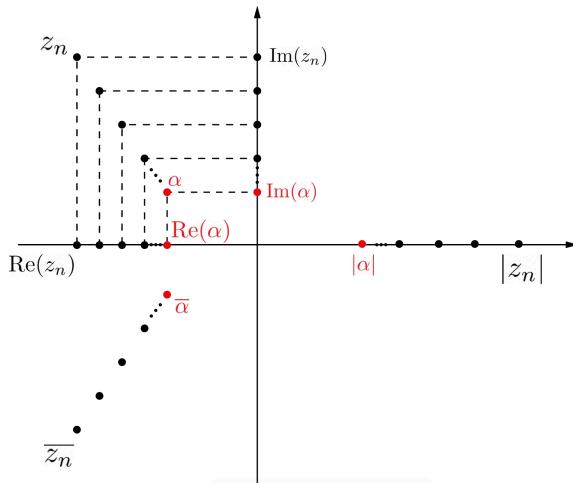
- (c)  $\lim_{n \rightarrow \infty} \operatorname{Re}(z_n) = \operatorname{Re}(\alpha)$
- (d)  $\lim_{n \rightarrow \infty} \operatorname{Im}(z_n) = \operatorname{Im}(\alpha)$

*Proof.* (a) We have that  $\|z_n - \alpha\| \leq |z_n - \alpha| \leq \epsilon$  by definition for all  $n > N$ . Hence for a given  $\epsilon > 0$ , a value of  $N$  which proves the convergence of  $z_n$  will also prove the convergence of  $|z_n|$ . In other words, the result follows from the Squeeze rule.

- (b) Immediate from  $|\bar{z_n} - \bar{\alpha}| = |z_n - \alpha|$
- (c) Immediate from  $|\operatorname{Re}(z_n) - \operatorname{Re}(\alpha)| = |\operatorname{Re}(z_n - \alpha)| \leq |z_n - \alpha|$  and the application of the Squeeze rule.
- (d) Immediate from  $|\operatorname{Im}(z_n) - \operatorname{Im}(\alpha)| = |\operatorname{Im}(z_n - \alpha)| \leq |z_n - \alpha|$  and the application of the Squeeze rule.

■

We can visualize the above proposition in the following figure



Finally, it is also useful to be able to combine different limits using the standard operations of addition, multiplication and division.

### Proposition (Limit combinations)

Let  $\lim_{n \rightarrow \infty} z_n = \alpha$  and  $\lim_{n \rightarrow \infty} w_n = \beta$ . Then

- (i)  $\lim_{n \rightarrow \infty} (z_n + w_n) = \alpha + \beta$
- (ii)  $\lim_{n \rightarrow \infty} (\lambda z_n) = \lambda \alpha$  for  $\lambda \in \mathbb{C}$
- (iii)  $\lim_{n \rightarrow \infty} (z_n w_n) = \alpha \beta$
- (iv)  $\lim_{n \rightarrow \infty} \frac{z_n}{w_n} = \frac{\alpha}{\beta}$  if  $\beta \neq 0$

*Proof.* Since  $z_n$  converges to  $\alpha$  and  $w_n$  converges to  $\beta$  we have that for all  $\epsilon_z, \epsilon_w > 0$  there exists  $N_z$  and  $N_w$  such that

$$|z_n - \alpha| < \epsilon_z, \forall n > N_z \quad (61.1.10)$$

$$|w_n - \beta| < \epsilon_w, \forall n > N_w \quad (61.1.11)$$

(i) Let  $\epsilon > 0$  and set  $\epsilon_z = \epsilon_w = \frac{1}{2}\epsilon$ . Then, letting  $N = \max(N_z, N_w)$  we have that

$$|(z_n + w_n) - (\alpha + \beta)| \leq |z_n - \alpha| + |w_n - \beta| < \epsilon, \forall n > N \quad (61.1.12)$$

as desired.

(ii) Let  $\epsilon > 0$  and set  $|\epsilon_z| = \frac{\epsilon}{|\lambda|}$ . Then, letting  $N = N_z$  we have that

$$|\lambda z_n - \lambda \alpha| = |\lambda||z_n - \alpha| < |\lambda|\epsilon_z = \epsilon, \forall n > N \quad (61.1.13)$$

as desired.

(iii) Let  $\epsilon > 0$ . We have that

$$\lim_{n \rightarrow \infty} (z_n w_n - \alpha \beta) = \lim_{n \rightarrow \infty} [z_n(w_n - \beta) + \beta(z_n - \alpha)] \quad (61.1.14)$$

The second term vanishes by the convergence of  $z_n$  and the scalar multiple property we proved in (ii). Furthermore, since  $z_n$  is convergent it is also bounded by some real number  $M$ . Hence

$$z_n w_n - \alpha \beta \leq M(w_n - \beta) \implies \lim_{n \rightarrow \infty} (z_n w_n - \alpha \beta) = 0 \quad (61.1.15)$$

as desired.

(iv) Let  $\epsilon > 0$ . We have that

$$\left| \frac{z_n}{w_n} - \frac{\alpha}{\beta} \right| = \frac{|\beta||z_n - \alpha| - |\alpha||w_n - \beta|}{|\beta||w_n|} \quad (61.1.16)$$

Our main worry is that  $|w_n| = 0$  for some  $n$  making this quantity ill-defined. However,  $|w_n|$  must eventually be positive, since

$$|\beta| - |w_n| < |w_n - \beta| < \frac{1}{2}|\beta| \implies |w_n| > \frac{1}{2}|\beta| \forall n > N \quad (61.1.17)$$

Consequently

$$\left| \frac{z_n}{w_n} - \frac{\alpha}{\beta} \right| < \frac{2}{|\beta|^2} (|\beta||z_n - \alpha| - |\alpha||w_n - \beta|) \quad (61.1.18)$$

and since the RHS defines a null sequence, the LHS must also be null by the Squeeze rule. ■

We have given a brief rundown of convergent complex sequences. Now let us look at the opposite case, that of divergent sequences. One very common family of divergent sequences are those that tend to infinity which we define below.

### Definition (*Infinite sequence*)

A sequence  $(z_n)$  **tends to infinity** if  $\forall M > 0$  there exists an integer  $N$  such that  $|z_n| >$

$$M, \forall n > N.$$

In other words a sequence converges to infinity if its points in the complex plane cannot be enclosed by a circle centered at the origin of finite radius. This is opposite to the case of a null sequence, where all terms after a certain point fall into any given circle of finite radius. Recalling that taking the reciprocal of a complex number is equivalent to scaling it by its norm and reflecting it in the real axis, it seems like the reciprocal of a null sequence tends to infinity and vice-versa.

**Theorem (Reciprocal rule)**

Let  $(z_n)$  be a sequence. Then  $(z_n)$  tends to infinity iff  $(\frac{1}{z_n})$  is null.

*Proof.* Suppose that  $(z_n)$  is null and let  $\epsilon > 0$  so that  $|z_n| \leq \epsilon, \forall n > N$  for some integer  $N$ . Then

$$\left| \frac{1}{z_n} \right| = \frac{1}{|z_n|} > \epsilon, \forall n > N \quad (61.1.19)$$

So given  $\epsilon > 0$  then by choosing the  $N$  which satisfies the epsilon-delta convergence of  $(z_n)$  we also find the  $N$  required to prove the divergence of  $(\frac{1}{z_n})$ . The converse proof is identical. ■

**Example.** We have that  $z_n = n^3 - in^2 + (1+i)n$  tends to infinity since

$$\frac{1}{z_n} = \frac{1}{n^3 - in^2 + (1+i)n} = \frac{1}{n^3} \frac{1}{1 - i/n + (1+i)/n^2} \rightarrow 0 \cdot 1 = 0 \quad (61.1.20)$$

◀

**Theorem (Subsequence rules)**

**First subsequence rule:** the sequence  $(z_n)$  diverges if it has two convergent subsequences with different limits.

**Second subsequence rule:** the sequence  $(z_n)$  diverges if  $(z_n)$  has a subsequence that tends to infinity.

*Proof.* (i) Trivial by application of epsilon-delta definition.

(ii) Trivial by boundedness of convergent sequences.

■

**Example.** Consider the sequence  $(z_n)$  with  $z_n = n^2 \sin(n\pi/3)$ . The subsequence  $z_{3k+1} = \frac{\sqrt{3}}{2}(3k+1)^2$  clearly diverges. Indeed its reciprocal

$$\frac{1}{z_{3k+1}} = \frac{2}{\sqrt{3}} \frac{1}{(3k+1)^2} \leq \frac{2}{\sqrt{3}} \frac{1}{k^2} \quad (61.1.21)$$

is a null sequence. Hence by the second subsequence rule  $(z_n)$  also diverges. ■

**Proposition (Subsequence decomposition)**

If a sequence  $(z_n)$  can be decomposed into two subsequences converging to the same value then  $(z_n)$  also converges to the same value.

*Proof.* ■

## 61.2 Continuity of complex functions

**Definition (Sequential continuity)**

Let  $f : A \rightarrow \mathbb{C}$  and  $\alpha \in A$ , then  $f$  is continuous at  $\alpha$  if for each sequence  $(z_n)$  in  $A$  such that  $z_n \rightarrow \alpha$ :

$$f(z_n) \rightarrow f(\alpha) \quad (61.2.1)$$

If  $f$  is continuous at all  $\alpha \in A$  then we say that it is continuous on  $A$ . If  $f$  is not continuous at  $\alpha$  then it is discontinuous.

**Example.** We consider some examples

- (i) Let  $f(z) = \bar{z}$  whose trivial domain is  $\mathbb{C}$ . We have that if  $(z_n)$  is a sequence on  $\mathbb{C}$  converging to some  $\alpha \in \mathbb{C}$  then  $\lim_{n \rightarrow \infty} f(z_n) = \lim_{n \rightarrow \infty} \bar{z_n} = \bar{\alpha} = f(\alpha)$  so  $f$  is indeed continuous.
- (ii) Let  $f(z) = \operatorname{Im}(z)$  whose trivial domain is  $\mathbb{C}$ . We have that if  $(z_n)$  is a sequence on  $\mathbb{C}$  converging to some  $\alpha \in \mathbb{C}$  then  $\lim_{n \rightarrow \infty} f(z_n) = \lim_{n \rightarrow \infty} \operatorname{Im}(z_n) = \operatorname{Im}(\alpha) = f(\alpha)$  so  $f$  is indeed continuous.
- (iii) Let  $f(z) = |z|$  whose trivial domain is  $\mathbb{C}$ . We have that if  $(z_n)$  is a sequence on  $\mathbb{C}$  converging to some  $\alpha \in \mathbb{C}$  then  $\lim_{n \rightarrow \infty} f(z_n) = \lim_{n \rightarrow \infty} |z_n| = |\alpha| = f(\alpha)$  so  $f$  is indeed continuous.
- (iv) Consider  $f(z) = \arg(z)$ , defined on  $\mathbb{C}$ , and let  $z_n = e^{i(\pi+1/n)}$  for  $n = 1, 2, \dots$ . Then we have that  $f(\lim_{n \rightarrow \infty} z_n) = \arg(e^{i\pi}) = \pi$  while  $\lim_{n \rightarrow \infty} f(z_n) = \lim_{n \rightarrow \infty} (\frac{1}{n} - \pi) = -\pi$ . Hence  $f(z)$  is discontinuous at  $z = -1$ . Consequently by the multiple rule for complex limits,  $f(z)$  is discontinuous on the negative real line.

In general we can combine continuous functions to get another continuous function back.

**Theorem (Combination rules)**

Let  $f, g$  be continuous functions at  $\alpha$  and let  $\lambda \in \mathbb{C}$ , then

- (i)  $f + g$  is continuous at  $\alpha$
- (ii)  $\lambda f$  is continuous at  $\alpha$
- (iii)  $fg$  is continuous at  $\alpha$
- (iv)  $f/g$  is continuous at  $\alpha$  provided  $g(\alpha) \neq 0$

*Proof.* These are immediate from the combination of limit theorems. ■

We can also compose continuous functions to get other continuous functions.

**Theorem (Composition rule)**

Let  $f$  be continuous at  $\alpha$  and let  $g$  be continuous at  $f(\alpha)$ . Then  $g \circ f$  is continuous at  $\alpha$ .

*Proof.* We have that for any sequence  $(z_n)$  on  $\mathbb{C}$ , if  $z_n \rightarrow \alpha$  then  $f(z_n) \rightarrow f(\alpha)$ . Furthermore,  $g$  is continuous at  $f(\alpha)$  so given any sequence  $w_n \rightarrow \beta$ , we must have that  $g(w_n) \rightarrow g(\beta)$ . In our case we consider the sequence  $f(z_n)$ , then  $g(f(z_n)) \rightarrow g(f(\alpha))$ . It thus follows that for any sequence  $(z_n) \rightarrow \alpha$ ,  $(g \circ f)(z_n) \rightarrow (g \circ f)(\alpha)$  as desired. ■

**Proposition (Restriction rule)**

Let  $f, g$  be functions defined on  $A, B$  respectively such that  $A \subseteq B$ . If

- (i)  $f(z) = g(z)$  for all  $z \in A$
  - (ii)  $g$  is continuous at  $\alpha \in A$
- then  $f$  is continuous at  $\alpha$ .

*Proof.* Let  $z_n$  be a sequence on  $A$ , and thus also in  $B$  such that  $z_n \rightarrow \alpha$ . Then we know that  $g(z_n) = f(z_n) \rightarrow g(\alpha) = f(\alpha)$ . Hence  $f$  is continuous at  $\alpha$ . ■

**Example.** We consider some examples

- (i) Consider  $f(x) = \frac{x^2+i}{x^2-i}$  defined for  $x \in \mathbb{R}$ . This is a restriction of  $g(z) = \frac{z^2+i}{z^2-i}$  over  $\mathbb{C} \setminus \{e^{i\pi/4}\}$ . We have that  $z^2 - i \neq 0$  for  $z \in \mathbb{R}$  so by the quotient combination rule we have that  $g(z)$  is continuous for all  $z \in \mathbb{R}$ , and thus so is  $f(x)$ .
- (ii) Consider  $f(z) = \log|z|$  defined over  $\mathbb{C}$ . We have that  $g(z) = |z|$  is a continuous function (this was proven earlier), and we know from real analysis that  $h(x) = \log x$  is continuous for  $x > 0$ . Since  $|z| > 0$  for all  $z \in \mathbb{C}$ , it follows from the composition rule that  $f(z)$  is continuous.

We now provide an equivalent definition of continuous functions.

**Definition ( $\epsilon - \delta$  continuity)**

Let  $f : A \rightarrow \mathbb{C}$  and  $\alpha \in A$ . Then  $f$  is continuous at  $\alpha$  if for all  $\epsilon > 0$  there exists a  $\delta > 0$  such that

$$|f(z) - f(\alpha)| < \epsilon, \forall z \in A \text{ s.t. } |z - \alpha| < \delta \quad (61.2.2)$$

Geometrically, this means that given any open disc of radius  $\epsilon$  around  $f(\alpha)$ , there exists an open disc of some radius  $\delta$  around  $\alpha$  such that its image under  $f$  lies within the disc of radius  $\epsilon$ .

**Theorem (Sequential and  $\epsilon - \delta$  continuity are equivalent)**

A function  $f(z)$  is continuous according to the sequential definition iff it is continuous according to the  $\epsilon - \delta$  definition.

*Proof.* Let  $\epsilon > 0$ . Suppose that  $|f(z) - f(\alpha)| < \epsilon$  for all  $z$  such that  $|z - \alpha| < \delta$ , where  $\delta > 0$ . Let  $(z_n)$  be a sequence on  $A$  with  $z_n \rightarrow \alpha$ , so there exists  $N$  such that  $\forall n > N$ ,  $|z_n - \alpha| < \delta$ . This however implies that  $|f(z_n) - f(\alpha)| < \epsilon$ , proving that  $f$  is continuous according to the  $\epsilon - \delta$  definition.

Now suppose that there is some  $\epsilon > 0$  such that for all  $\delta > 0$ ,  $|f(z) - f(\alpha)| \geq \epsilon$  for  $z$  with  $|z - \alpha| < \delta$ . So setting  $\delta = \frac{1}{n}$  then  $|f(z_n) - f(\alpha)| \geq \epsilon$  for  $z_n \in A$  with  $|z_n - \alpha| < \frac{1}{n}$ . Clearly  $z_n \rightarrow \alpha$  but  $f(z_n) \not\rightarrow f(\alpha)$ .  $\blacksquare$

**Example.** Let us show that  $f(z) = \arg(z)$  is continuous on  $A = \mathbb{C} \setminus \{x \in \mathbb{R} : x \leq 0\}$ .

Let  $\alpha \in A$ , and let  $\epsilon > 0$ , then we need to find  $\delta > 0$  such that

$$|z - \alpha| < \delta \implies |f(z) - f(\alpha)| < \epsilon \quad (61.2.3)$$

We choose  $\delta$  so that the disc  $\{z \in \mathbb{C} : |z - \alpha| < \delta\}$  does not cross the negative real axis. This is important to ensure that the disc lies entirely in  $A$ . We also choose it so that the angular size of the circle from the origin is less than  $2\epsilon$ , which is equivalent to saying that we choose  $\delta > 0$  such that  $\sin^{-1}(\delta/|\alpha|) < \epsilon$ .

Then we find that if  $z$  is in this disc then

$$|f(z) - f(\alpha)| = |\arg(z) - \arg(\alpha)| < \epsilon \quad (61.2.4)$$

as desired.  $\blacktriangleleft$

## 61.3 Limits of complex functions

### Definition (Accumulation point)

Let  $A \subseteq \mathbb{C}$ . Suppose there exists a sequence  $(z_n)$  in  $A \setminus \{\alpha\}$  such that  $z_n \rightarrow \alpha$ . Then  $\alpha$  is an accumulation point of  $A$ .

### Definition (Limit of a function)

Let  $\alpha$  be an accumulation point of  $A$  and let  $f(z)$  be defined on  $A$ . Then we say that  $\lim_{z \rightarrow \alpha} f(z) \rightarrow \beta$  if for any sequence  $z_n$  in  $A \setminus \{\alpha\}$  such that  $z_n \rightarrow \alpha$ , we have that  $f(z_n) \rightarrow \beta$ .

### Example.

(i) Consider the following limit

$$\lim_{z \rightarrow i} \frac{z^3 + i}{z - i} \quad (61.3.1)$$

The function  $f(z) = (z^3 + i)/(z - i)$  is defined on  $\mathbb{C} \setminus \{i\}$ , and  $i$  is a limit point in this set. Indeed consider  $z_n = i - \frac{1}{n}$ , which lies entirely in the given region, it is clear that  $z_n \rightarrow i$  as  $n \rightarrow \infty$ .

Therefore, if  $z_n$  is a sequence in  $\mathbb{C} \setminus \{i\}$  such that  $z_n \rightarrow i$  then we have that

$$f(z_n) = \frac{z_n^3 + i}{z_n - i} = (z_n^2 + iz_n - 1) \rightarrow -3 \quad (61.3.2)$$

so

$$\lim_{z \rightarrow i} \frac{z^3 + i}{z - i} = -3 \quad (61.3.3)$$

(ii) Consider the limit

$$\lim_{z \rightarrow 0} \frac{z}{\operatorname{Re}(z)} \quad (61.3.4)$$

Note that  $f(z) = \frac{z}{\operatorname{Re}(z)}$  is defined on  $\mathbb{C} \setminus \{z : \operatorname{Re}(z) = 0\}$  and thus 0 is a limit point in this domain. Consider  $x \in \mathbb{R} \setminus \{0\}$  we have that  $\frac{x}{\operatorname{Re}(x)} = 1$ . Instead  $\frac{ix}{\operatorname{Re}(ix)}$  is ill-defined. Consequently let us consider the sequence  $z_n = \frac{1}{n} \rightarrow 0$ . We have that

$$f(z_n) = \frac{1/n}{1/n} \rightarrow 1 \quad (61.3.5)$$

On the other hand, consider the sequence  $w_n = \frac{1}{n} + \frac{i}{n}$ . Then we have that

$$f(w_n) = \frac{(1+i)/n}{1/n} \rightarrow 1+i \quad (61.3.6)$$

so there is no  $\beta$  such that  $f(z_n)$  converges to the same limit for all sequences  $z_n$ .

◀

### Theorem (Continuity with limits)

Let  $f$  be a function defined on  $A$  and let  $\alpha \in A$  be a limit point of  $A$ . Then  $f$  is continuous at  $\alpha$  iff  $\lim_{z \rightarrow \alpha} f(z) = f(\alpha)$ .

*Proof.* Suppose that  $f$  be continuous at  $\alpha$ . Then given any sequence  $z_n \rightarrow \alpha$  we have that  $f(z_n) \rightarrow f(\alpha)$ , which is the definition of  $\lim_{z \rightarrow \alpha} f(z) = f(\alpha)$ .

Now suppose that  $\lim_{z \rightarrow \alpha} f(z) = f(\alpha)$ . Then let  $(z_n)$  be a sequence on  $A \setminus \{\alpha\}$  such that  $z_n \rightarrow \alpha$ . Then we can separate  $(z_n)$  into a subsequence  $(z_{n_k})$  such that  $f(z_{n_k}) = f(\alpha)$  for all  $k$ , and a subsequence  $(z_{m_k})$  such that  $f(z_{m_k}) \neq f(\alpha)$ . Clearly  $f(z_{n_k}) \rightarrow f(\alpha)$  and  $f(z_{m_k}) \rightarrow f(\alpha)$  by assumption, so we have that  $f$  is continuous at  $\alpha$ . ■

Finally, to facilitate computations with limits we have the typical combination rules which can be proven using the sequence limit rules

### Theorem (Combination rules)

Let  $f, g$  be functions with domains  $A, B$  respectively, and let  $\alpha$  be a limit point of  $A \cup B$  so that

$$\lim_{z \rightarrow \alpha} f(z) = \beta, \quad \lim_{z \rightarrow \alpha} g(z) = \gamma \quad (61.3.7)$$

Then

- (i)  $\lim_{z \rightarrow \infty} (f(z) + g(z)) = \beta + \gamma$
- (ii)  $\lim_{z \rightarrow \infty} (\lambda f(z)) = \lambda \beta$  for  $\lambda \in \mathbb{C}$
- (iii)  $\lim_{z \rightarrow \infty} (f(z)g(z)) = \beta \gamma$

$$(iii) \lim_{z \rightarrow \infty} \frac{f(z)}{g(z)} = \frac{\beta}{\gamma} \text{ if } \gamma \neq 0$$

## 61.4 Topology on $\mathbb{C}$

### Definition (Open set)

A set  $A \subseteq \mathbb{C}$  is open if  $\forall \alpha \in A$ , there exists an open disc  $D_r(\alpha) = \{z \in \mathbb{C} : |z - \alpha| < r\}$  of some radius  $r > 0$  that lies entirely in  $A$ , that is  $D_r(\alpha) \subseteq A$ .

### Example.

- (i) The set  $A = \mathbb{C} \setminus \{0\}$  is open. Let  $\alpha \in A$ , then clearly  $D_{|\alpha|}(\alpha) = \{z \in \mathbb{C} : |z - \alpha| < |\alpha|\}$  lies entirely in  $A$ .
- (ii) The set  $A = \{z : -2 < \operatorname{Re}(z) < 2, -1 < \operatorname{Im}(z) < 1\}$  is open. Let  $\alpha \in A$ , then let  $r = \min\{\operatorname{Re}(\alpha) - 2, \operatorname{Re}(\alpha) + 2, \operatorname{Im}(\alpha) - 1, \operatorname{Im}(\alpha) + 1\}$  which is non-zero since  $\alpha = \pm 2, \pm i$ . Then  $D_r(\alpha)$  is entirely contained in  $A$ .
- (iii) The set  $A = \{z : 1 < |z| < 2\}$  is open. Let us define for all  $\alpha \in A$  the following quantity

$$r = \begin{cases} |\alpha| - 1, & |\alpha| < \frac{3}{2} \\ 2 - |\alpha|, & |\alpha| > \frac{3}{2} \end{cases} \quad (61.4.1)$$

Then  $D_r(\alpha)$  is entirely contained in  $A$ .

- (iv) Let  $A = \{z : \pi/3 < \arg z < 2\pi/3\}$ , and let  $\alpha \in A$ . Let  $\theta = \arg(\alpha)$ , and define

$$r = \begin{cases} |\alpha| \sin\left(\frac{2\pi}{3} - \theta\right), & \alpha \leq 0 \\ |\alpha| \sin\left(\theta - \frac{\pi}{3}\right), & \alpha > 0 \end{cases} \quad (61.4.2)$$

Then clearly  $D_r(\alpha)$  lies entirely in  $A$ .

◀

### Theorem (Combination of open sets)

If  $A_1$  and  $A_2$  are closed sets then

- (i)  $A_1 \cup A_2$
- (ii)  $A_1 \cap A_2$

are also open.

*Proof.*

- (i) Let  $\alpha \in A_1 \cup A_2$ , and suppose  $\alpha \in A_1$  wlog. Then, there exists some  $r > 0$  such that  $D_r(\alpha) \subseteq A_1 \subseteq A_1 \cup A_2$ , as desired.
- (ii) Let  $\alpha \in A_1 \cap A_2$ , so there exist  $r_1, r_2 > 0$  such that  $D_{r_1}(\alpha) \subseteq A_1$  and  $D_{r_2}(\alpha) \subseteq A_2$ . Choosing  $r = \min\{r_1, r_2\}$  then  $D_r(\alpha) \subseteq A_1 \cap A_2$ .

■

**Definition (Closed set)**

A set  $E \subseteq \mathbb{C}$  is closed if its complement  $\overline{E} \equiv \mathbb{C} \setminus E$  is open.

The combination rules for open sets apply equally to closed sets.

**Theorem (Combination of closed sets)**

If  $E_1$  and  $E_2$  are open sets then

- (i)  $E_1 \cup E_2$
- (ii)  $E_1 \cap E_2$

are also open.

*Proof.*

- (i) Since  $E_1$  and  $E_2$  are closed, their complements  $\mathbb{C} \setminus E_1$  and  $\mathbb{C} \setminus E_2$  are open. Consequently

$$\mathbb{C} \setminus (E_1 \cup E_2) = (\mathbb{C} \setminus E_1) \cap (\mathbb{C} \setminus E_2) \quad (61.4.3)$$

is open since it is the intersection of two open sets.

- (ii) Similarly

$$\mathbb{C} \setminus (E_1 \cap E_2) = (\mathbb{C} \setminus E_1) \cup (\mathbb{C} \setminus E_2) \quad (61.4.4)$$

is open since it is the union of two open sets.

■

**Definition (Interior, exterior and boundary)**

Let  $A \subseteq \mathbb{C}$  and let  $\alpha \in \mathbb{C}$ . Then  $\alpha$  is an **interior point** of  $A$  if there is an open disc centered at  $\alpha$  that lies in  $A$ . Similarly  $\alpha$  is an **exterior point** of  $A$  if there is an open disc centered at  $\alpha$  that lies outside of  $A$ . Finally,  $\alpha$  is a **boundary point** of  $A$  if any open disc centered at  $\alpha$  contains at least one point in  $A$  and one point in  $\overline{A}$ .

The set of interior points of  $A$  forms the interior  $\text{int } A$ , the set of exterior points forms the exterior  $\text{ext } A$  and the set of boundary points forms the boundary  $\partial A$ .

Note that any point in an open set is an interior point by definition.

**Definition (Connectedness)**

A set  $A \subseteq \mathbb{C}$  is (pathwise) connected if for any two distinct points  $\alpha, \beta \in A$ , there exists a path  $\Gamma \subseteq A$  joining the two.

**Example.** Let us prove that if  $A, B$  are both connected and  $A \cap B \neq \emptyset$  then  $A \cup B$  is connected.

Let  $\alpha \in A \cup B$ , and assume wlog that  $\alpha \in A$  and  $\beta \in B$  (if both points lie in  $A$  then the proof is trivial). Since  $A \cap B \neq \emptyset$  we may assume that there exists at least some  $\delta \in A \cap B$ . Since  $\delta \in A$  there exists a path  $\Gamma_A$  connecting  $\alpha$  with  $\delta$ . Since  $\delta \in B$  there exists a path  $\Gamma_B$  connecting  $\delta$  with  $\beta$ . Consequently  $\Gamma_1 \cup \Gamma_2$  connects  $\alpha$  with  $\beta$  as desired. ■

**Proposition** (*Continuous functions preserve connectedness*)

Let  $f$  be a continuous function on a connected domain  $A$ . Then  $f(A)$  is also connected.

*Proof.* Since  $A$  is connected, given  $\alpha, \beta \in A$  we can find a curve  $\Gamma$  parametrised continuously by a function

$$\gamma : [a, b] \rightarrow A, \gamma(a) = \alpha, \gamma(b) = \beta \quad (61.4.5)$$

Since  $f$  is continuous on  $A$ , it follows from the composition rule that  $f \circ \gamma$  is continuous on  $[a, b]$ . Moreover,  $(f \circ \gamma)(a) = f(\alpha)$  and  $(f \circ \gamma)(b) = \beta$ , so  $f \circ \gamma$  parametrises a curve  $f(\Gamma) = (f \circ \gamma)[a, b] \in f(A)$ . ■

**Definition** (*Region*)

A region  $\mathcal{R}$  is a non-empty, open, connected subset of  $\mathbb{C}$ .

**Theorem** (*Point-removal of regions*)

If  $\mathcal{R}$  is a region and  $\alpha_0 \in \mathcal{R}$  then  $\mathcal{R} \setminus \{\alpha_0\}$  is also a region.

*Proof.* Since  $\mathcal{R}$  is a region, and thus open, it cannot be a singleton and contains more than one element. Consequently  $\mathcal{R} \setminus \{\alpha_0\}$  is not empty. Now note that  $\mathcal{R} \setminus \{\alpha_0\} = \mathcal{R} \cap (\mathbb{C} \setminus \{\alpha_0\})$  is the intersection of two open sets, and so must also be open. Finally, let  $\alpha, \beta \in \mathcal{R} \setminus \{\alpha_0\}$ . Since  $\alpha, \beta \in \mathcal{R}$  they can be joined by a curve  $\Gamma$ . If  $\alpha_0 \notin \Gamma$  then this curve will lie in  $\mathcal{R} \setminus \{\alpha_0\}$ . If instead  $\alpha_0 \in \Gamma$ , we consider a disc  $D_r(\alpha_0) \subset \mathcal{R}$  and deform the curve so as to avoid  $\alpha_0$ . ■

Note that the above proof can be extended to the removal of a finite number of points, but breaks down for infinitely many points. Indeed the removal of infinitely many points should of course not necessarily lead to a non-empty set, so the point-removal theorem is in general false for infinitely many points.

**Example.** Let us prove that  $\mathcal{R} = \mathbb{C} \setminus \{(n + \frac{1}{2})\pi : n \in \mathbb{Z}\}$  is a region. It is obviously non-empty.

Unfortunately we cannot use the point-removal theorem since we are removing an infinite number of points, we go back to the original definition. Firstly note that  $\mathcal{R}$  is non-empty, since for example  $i \in \mathcal{R}$ . Now let  $\alpha \in \mathcal{R}$ , and let  $m$  be the integer such that  $(m + \frac{1}{2})\pi$  is the closest removed point to  $\alpha$ . Then letting  $r = |\alpha - (m + \frac{1}{2})\pi|$  we have that  $D_r(\alpha)$  lies entirely in  $\mathcal{R}$ . Finally, let  $\alpha, \beta \in \mathcal{R}$  and let  $\Gamma$  be the straight segment in  $\mathbb{C}$  joining them. If no removed point lies on  $\Gamma$  then this curve also lies in  $\mathcal{R}$ . If one of the removed point lies on  $\Gamma$  then we can take a neighborhood of the point and deform it so as to skip around it. In the worst case scenario only finitely many such deformations will be necessary, and the curve will lie entirely in  $\mathcal{R}$ . ◀

We make one final definition that will be useful in the next section.

**Definition** (*Compact sets*)

A set  $E \subset \mathbb{C}$  that is contained in some closed disc is bounded. Otherwise it is unbounded.

If  $E$  is both closed and bounded then it is compact.

## 61.5 Extreme value theorem

### Theorem (Extreme value theorem)

Let  $f$  be a function that is continuous on a compact set  $E$ . Then  $\exists \alpha, \beta \in E$  such that

$$|f(\beta)| \leq |f(z)| \leq |f(\alpha)|, \forall z \in E \quad (61.5.1)$$

*Proof.* We begin by proving the following lemma

**Lemma.** If  $E$  is a closed set and  $(z_n)$  is convergent in  $E$  to  $\alpha$ , then  $\alpha \in E$ .

Indeed suppose that  $\alpha \in \overline{E}$ , which is open and thus contains an open disc  $D_r(\alpha)$ . Then if  $z_n \rightarrow \alpha$  then  $D_r(\alpha) \subseteq \overline{E}$  contains all but a finite number of points of  $z_n$ . This contradicts the fact that  $z_n$  is in  $E$ .

We will also need the following result

### Theorem (Nested rectangles theorem)

Let  $R_0, R_1, R_2, \dots$  be a sequence of closed rectangular regions whose sides are parallel to the real and imaginary axes, and whose diagonals' lengths form a sequence  $s_0, s_1, s_2, \dots$  such that

$$R_0 \supseteq R_1 \supseteq R_2 \supseteq \dots, \lim_{n \rightarrow \infty} s_n = 0 \quad (61.5.2)$$

Then

$$\bigcap_i R_i = \{\alpha\} \quad (61.5.3)$$

for a unique complex number  $\alpha$ . Furthermore, given  $\epsilon > 0$  then there is an integer  $N$  such that

$$R_n \subseteq D_\epsilon(\alpha), \forall n > N \quad (61.5.4)$$

*Proof.* Let the  $n$ th rectangle be defined as

$$R_n = \{x + iy : a_n \leq x \leq c_n, b_n \leq y \leq d_n\} \quad (61.5.5)$$

so that since the rectangles are nested within each other

$$a_0 \leq a_1 \leq a_2 \leq \dots \leq c_2 \leq c_1 \leq c_0 \quad (61.5.6)$$

$$b_0 \leq b_1 \leq b_2 \leq \dots \leq d_2 \leq d_1 \leq d_0 \quad (61.5.7)$$

Since  $a_n$  is an increasing, bounded sequence ( $a_n \leq c_0$  for all  $n$ ) it follows from the monotone convergence theorem that  $a_n$  converges. Similarly  $b_n$  also converges. Let  $\lim_{n \rightarrow \infty} a_n = a$  and  $\lim_{n \rightarrow \infty} b_n = b$ .

We also have that

$$0 \leq c_n - a_n \leq s_n, 0 \leq d_n - b_n \leq s_n \quad (61.5.8)$$

implying that  $c_n \rightarrow a$  and  $d_n \rightarrow b$  by the Squeeze theorem. Let  $\alpha = a + ib$ , then since by the monotone convergence theorem  $a_n \leq a \leq c_n$  and  $b_n \leq b \leq d_n$  for all  $n$  it follows that  $\alpha$  is contained in all the rectangles:

$$\alpha \in \bigcap_i R_i \quad (61.5.9)$$

Now let  $z \in R_n$ , so that  $|z - \alpha| \leq s_n$  since any two points in a rectangle are at most a diagonal length apart. By assumption  $s_n \rightarrow 0$  so for a given  $\epsilon > 0$  there exists  $N$  such that  $s_n \leq \epsilon$  for all  $n > N$  and thus such that  $|z_n - \alpha| < \epsilon$ . It follows that

$$R_n \subseteq \{z : |z - \alpha| < \epsilon\}, \forall n > N \quad (61.5.10)$$

Consequently, suppose that there is another complex number  $\beta$  satisfying

$$\beta \in \bigcap_i R_i \quad (61.5.11)$$

but then there for any  $\epsilon$  we have that  $|\beta - \alpha| < \epsilon$ . If we let  $\epsilon = |\beta - \alpha| > 0$  then we obtain a contradiction, so  $\alpha$  is the only element in the intersection of nested rectangles. ■

We can now tackle the Extreme value theorem. We begin by proving that there exists an  $\alpha \in E$  such that  $|f(z)| \leq |f(\alpha)|$  for all  $z \in E$ .

Since  $E$  is compact it can be surrounded by a rectangle  $R_0$  of diagonal length  $s_0$ . This rectangle can be subdivided into four equivalent rectangles  $T_1, T_2, T_3, T_4$  of diagonal length  $\frac{s_0}{2}$ . Then at least one of the rectangles,  $T_j$ , is such that for all  $z \in E$  there exists  $w \in E \cap T_j$  such that  $|f(z)| \leq |f(w)|$ . In other words, there is at least one rectangle containing the value of  $z$  which bounds  $f$  from above. Indeed if this were false, then there is a complex  $z_k \in E$  with  $|f(z_k)| > |f(w)|$  for all  $w \in E \cap T_k$ . But then

$$\max\{|f(z_1)|, |f(z_2)|, |f(z_3)|, |f(z_4)|\} > |f(w)|, \forall w \in E \quad (61.5.12)$$

which is false since  $>$  should be replaced with a  $\geq$ .

Now let  $R_1 = T_j$  and repeat the process of subdividing and finding the quadrant containing the maximal value of  $f$ . Iterating this indefinitely we will get a sequence  $R_n$  of nested closed rectangular regions with the property that

$$R_{n+1} \subseteq R_n, s_n = \frac{1}{2^n} s_0 \rightarrow 0 \quad (61.5.13)$$

and for each  $z \in E \cap R_n$  there exists some  $w \in E \cap R_{n+1}$  such that  $|f(z)| < |f(w)|$ . Using the nested rectangle theorem we know that only one complex number  $\alpha$  will lie in all rectangles, and furthermore for a given  $\epsilon > 0$  there is a  $N$  such that  $R_n \subseteq D_\epsilon(\alpha)$  for all  $n > N$ . We now claim that the number  $\alpha$  must lie in  $E$  and that

$$f(z) \leq f(\alpha), \forall z \in E \quad (61.5.14)$$

To see why, let  $z_0 \in E$  and construct from it a sequence such that

$$z_n \in E \cap R_n, (z_0 \in E \cap R_0 = E \text{ is trivially satisfied}) \quad (61.5.15)$$

and

$$|f(z_n)| \leq |f(z_{n+1})| \quad (61.5.16)$$

Since  $z_n$  is a sequence on a closed set  $E$  converging to  $\alpha$  it follows that  $\alpha \in E$ , as desired. Moreover, by the continuity of  $f$  on  $E$  it follows that

$$\lim_{n \rightarrow \infty} |f(z_n)| = |f(\alpha)| \quad (61.5.17)$$

implying that

$$|f(z_0)| \leq |f(z_1)| \leq \dots \leq |f(\alpha)| \quad (61.5.18)$$

since  $|f(z_0)|$  is an increasing sequence. Notice that this argument applies to any  $z_0 \in E$  so it follows that  $|f(z)| \leq |f(\alpha)|$  for any  $z$ .

We have proven the first part of the theorem, now we prove that there is  $\beta$  such that  $|f(\beta)| \leq |f(z)|$  for all  $z \in E$ .

If  $f(w) = 0$  for some  $w$  then  $\beta = w$  will do the job. Otherwise  $f$  does not vanish on  $E$ , and thus  $g(z) = \frac{1}{f(z)}$  is a continuous function by the quotient rule. We can apply the argument used to prove the first part of the theorem to prove that such a complex number  $\beta$  must exist. ■

### Proposition (Map of compact set is compact)

Let  $f$  be a function that is continuous on a compact set  $E$ . Then  $f(E)$  is also compact.

*Proof.* We know that  $f(E)$  must be bounded by the Extreme value theorem.

To prove that  $f(E)$  is closed, we prove that  $A = \mathbb{C} \setminus f(E)$  is open. Let  $\alpha \in A$  and define

$$g(z) = f(z) - \alpha \quad (61.5.19)$$

By the extreme value theorem there is some  $\beta \in E$  such that

$$r \equiv |g(\beta)| \leq |g(z)|, \forall z \in E \quad (61.5.20)$$

implying that

$$|f(z) - \alpha| \geq r, \forall z \in E \quad (61.5.21)$$

Therefore if  $|f(z) - \alpha| < r$  then  $z \notin E$  so  $f(z) \notin f(E)$ . Consequently we have that  $D_r(\alpha) \subseteq A$ , as desired. ■

# Differentiating complex functions

## 62.1 Derivatives of complex functions

Having discussed the concept of a limit of a complex function and its relation to continuity, we are now ready to introduce the concept of differentiating complex functions.

### Definition (*Differentiability*)

Let  $f$  be a complex function on  $A$ . Then it is said to be **differentiable** at  $\alpha$  if the limit

$$\lim_{h \rightarrow 0} \frac{f(\alpha + h) - f(\alpha)}{h} \quad (62.1.1)$$

exists. This limit is known as the **derivative of  $f$  at  $\alpha$** . If  $f$  is differentiable on every point in a set  $B$  then it is differentiable on  $B$ . A differentiable function is differentiable on its domain. A function differentiable on  $\mathbb{C}$  is defined to be **entire**.

Note that for a function to be differentiable at a point we need that point to be a limit point of the function's domain.

**Example.** Let  $f(z) = \frac{1}{z}$  be defined on  $\mathbb{C} \setminus \{0\}$ . Letting  $\alpha \in \mathbb{C} \setminus \{0\}$  then we find that

$$\lim_{z \rightarrow \alpha} \frac{f(z) - f(\alpha)}{z - \alpha} = \lim_{z \rightarrow \alpha} \frac{(\alpha - z)/(z\alpha)}{z - \alpha} = - \lim_{z \rightarrow \alpha} \frac{1}{z\alpha} = - \frac{1}{\alpha^2} \quad (62.1.2)$$

so the derivative of  $f(z) = \frac{1}{z}$  is  $f'(z) = -\frac{1}{z^2}$ , and is defined on  $\mathbb{C} \setminus \{0\}$ . ◀

Note that while  $f(z)$  is not entire, it is still differentiable on a region. This property is useful and deserves its own name, such functions are known as **analytic** functions.

### Definition (*Analyticity*)

A function  $f(z)$  that is differentiable on a region  $\mathcal{R}$  is said to be **analytic on  $\mathcal{R}$** . If this region is the domain of  $f$  then it is simply **analytic**. Also  $f$  is **analytic at  $\alpha$**  if there exists a region containing  $\alpha$  on which it is differentiable.

Notice that if a function is differentiable at a point then it is not necessarily analytic at that point, although the converse is true.

As usual we can also combine differentiable functions

**Theorem (Derivative properties)** Let  $f, g$  have domains  $A$  and  $B$ , and be differentiable at  $\alpha$ . Then

- (i)  $(f + g)'(\alpha) = f'(\alpha) + g'(\alpha)$ .
- (ii)  $(\lambda f)'(\alpha) = \lambda f'(\alpha)$  for  $\lambda \in \mathbb{C}$
- (iii)  $(fg)'(\alpha) = f'(\alpha)g(\alpha) + f(\alpha)g'(\alpha)$
- (iv)  $\left(\frac{f}{g}\right)'(\alpha) = \frac{f'(\alpha)g(\alpha) - f(\alpha)g'(\alpha)}{(g(\alpha))^2}$  provided  $g(\alpha) \neq 0$ .

**Theorem (Differentiability  $\implies$  continuity)**

Let  $f$  be differentiable at  $\alpha$ , then it must be continuous at  $\alpha$ .

*Proof.* Suppose that  $z \rightarrow \alpha$ , then

$$\lim_{z \rightarrow \alpha} (f(z) - f(\alpha)) = \lim_{z \rightarrow \alpha} \frac{f(z) - f(\alpha)}{z - \alpha} \cdot \lim_{z \rightarrow \alpha} (f(z) - f(\alpha)) = f'(\alpha) \cdot 0 = 0 \quad (62.1.3)$$

so  $f(z) \rightarrow f(\alpha)$  as desired. ■

Consider for example the principal logarithm function  $\log z = \log |z| + i\arg(z)$ . We know that at no point on the negative real axis is this function continuous, so it is also not differentiable there.

A more striking example is that of the modulus function  $f(z) = |z|$ . In real analysis this function is differentiable at all non-zero points, but in complex analysis we shall see that this no longer holds.

The proof that  $f(z) = |z|$  is discontinuous at  $z = 0$  is the same as in real analysis, so let  $\alpha \neq 0$  and consider the circle centered at the origin passing through  $\alpha$ . Assume Let us define  $(z_n)$  to be the sequence

$$z_n = |\alpha| \exp\left(i\arg(\alpha) + \frac{i}{n}\right) \quad (62.1.4)$$

and similarly let us define  $(w_n)$  so that

$$w_n = \left(|\alpha| + \frac{1}{n}\right) e^{i\arg(\alpha)*} \quad (62.1.5)$$

Clearly  $z_n \rightarrow \alpha$  and  $w_n \rightarrow \alpha$ . However, note that

$$\lim_{n \rightarrow \infty} \frac{f(z_n) - f(\alpha)}{z_n - \alpha} = \lim_{n \rightarrow \infty} \frac{|\alpha| - |\alpha|}{z_n - \alpha} = 0 \quad (62.1.6)$$

while

$$\lim_{n \rightarrow \infty} \frac{f(w_n) - f(\alpha)}{w_n - \alpha} = \lim_{n \rightarrow \infty} \frac{1/n}{e^{i\arg(\alpha)}(|\alpha| + 1/n - |\alpha|)} = e^{-i\arg(\alpha)} \quad (62.1.7)$$

Consequently the limit

$$\lim_{z \rightarrow \alpha} \frac{f(z) - f(\alpha)}{z - \alpha} \quad (62.1.8)$$

does not exist, and  $f(z) = |z|$  is not differentiable anywhere on  $\mathbb{C}$ .

**Example.** Consider the function  $f(z) = \bar{z}$  defined over  $\mathbb{C}$ . There are no points where  $f(z)$  is differentiable. Indeed let  $\alpha \in \mathbb{C}$  and consider the sequence  $z_n = \alpha - \frac{1}{n}$ . Then

$$\lim_{n \rightarrow \infty} \frac{f(z_n) - f(\alpha)}{z_n - \alpha} = \lim_{n \rightarrow \infty} \frac{-1/n}{-1/n} = 1 \quad (62.1.9)$$

Now consider the sequence  $w_n = \alpha - \frac{i}{n}$ . Then

$$\lim_{n \rightarrow \infty} \frac{f(z_n) - f(\alpha)}{z_n - \alpha} = \lim_{n \rightarrow \infty} \frac{i/n}{-i/n} = -1 \quad (62.1.10)$$

Since these two limits are not equal the function  $f$  is not differentiable at  $\alpha$  which we assumed to be any complex number.  $\blacktriangleleft$

## 62.2 Cauchy-Riemann equations

### Theorem (Cauchy-Riemann theorem)

Let  $f(x+iy) = u(x, y) + iv(x, y)$  be a complex function defined on a region  $\mathcal{R}$  which contains  $a + ib$ . If  $f$  is differentiable at  $a + ib$  then the following **Cauchy-Riemann equations** are satisfied:

$$\frac{\partial u}{\partial x} \Big|_{(a,b)} = \frac{\partial v}{\partial y} \Big|_{(a,b)}, \quad \frac{\partial v}{\partial x} \Big|_{(a,b)} = -\frac{\partial u}{\partial y} \Big|_{(a,b)} \quad (62.2.1)$$

*Proof.* Let that  $z_n$  be any sequence on  $\mathcal{C} \setminus \{\alpha\}$  where  $\alpha = a + ib$  and let  $z_n = x_n + iy_n$  for all  $n$ . Firstly we observe that

$$\frac{f(z_n) - f(\alpha)}{z_n - \alpha} = \frac{u(x_n, y_n) - u(a, b)}{(x_n - a) + i(y_n - b)} + i \frac{v(x_n, y_n) - v(a, b)}{(x_n - a) + i(y_n - b)} \quad (62.2.2)$$

Consequently, let us choose  $x_n$  to be a sequence on  $\mathcal{R} \setminus \{a\}$  and define  $z_n = x_n + ib$  which converges to  $\alpha$ . Then we see that

$$\frac{f(z_n) - f(\alpha)}{z_n - \alpha} = \frac{u(x_n, b) - u(a, b)}{(x_n - a)} + i \frac{v(x_n, b) - v(a, b)}{(x_n - a)} \quad (62.2.3)$$

so after taking the limit we find

$$\frac{df}{dz} \Big|_{\alpha} = \frac{\partial u}{\partial x} \Big|_{(a,b)} + i \frac{\partial u}{\partial x} \Big|_{(a,b)} \quad (62.2.4)$$

Now let us choose  $y_n$  to be a sequence on  $\mathcal{R} \setminus \{b\}$  and define  $z_n = a + iy_n$  which converges to  $\alpha$ . Then we see that

$$\frac{f(z_n) - f(\alpha)}{z_n - \alpha} = \frac{u(a, y_n) - u(a, b)}{i(y_n - b)} + i \frac{v(a, y_n) - v(a, b)}{i(y_n - b)} \quad (62.2.5)$$

so after taking the limit we find

$$\frac{df}{dz} \Big|_{\alpha} = -i \frac{\partial u}{\partial y} \Big|_{(a,b)} + \frac{\partial v}{\partial y} \Big|_{(a,b)} \quad (62.2.6)$$

Since  $f$  is differentiable at  $\alpha$  we require (??) and (??) to be equal to each other. This gives the desired Cauchy-Riemann equations. ■

This theorem is extremely useful as it also tells us that if the Cauchy-Riemann equations fail then the function is not differentiable. Let us apply it to  $f(x + iy) = x - iy$  which we have previously determined to not be differentiably anywhere on  $\mathbb{C}$ . Now we see that

$$\frac{\partial u}{\partial x} \Big|_{(a,b)} = 1 \neq \frac{\partial v}{\partial y} \Big|_{(a,b)} = -1 \quad (62.2.7)$$

as desired.

Note that the converse of the Cauchy-Riemann theorem is not true, if a function satisfies (62.2.2) then it is not necessarily differentiable. If we add some extra conditions however a converse-like theorem can be proven.

**Theorem (Cauchy-Riemann converse theorem)**

Let  $f(x + iy) = u(x, y) + iv(x, y)$  be defined on a region  $\mathcal{R}$  and let  $\alpha = a + ib \in \mathcal{R}$ . If the derivatives  $\partial_x u, \partial_y u, \partial_x v, \partial_y v$

- (i) exist at  $(x, y)$  for all  $x + iy \in \mathcal{R}$
- (ii) are continuous at  $(a, b)$
- (iii) satisfy the Cauchy-Riemann equations

then  $f$  is differentiable at  $a + ib$  and its derivative is given by

$$f'(a + ib) = \frac{\partial u}{\partial x} \Big|_{(a,b)} + i \frac{\partial v}{\partial x} \Big|_{(a,b)} \quad (62.2.8)$$

**Example.** Let us show that  $f(z) = \sin z$  is entire. Indeed we have that

$$\sin(x + iy) = \sin x \cosh y + i \sinh y \cos x \quad (62.2.9)$$

so that  $u(x, y) = \sin x \cosh y$  and  $v(x, y) = \sinh y \cos x$ . Consequently

$$\frac{\partial u}{\partial x} = \cos x \cosh y \quad \frac{\partial v}{\partial y} = \cosh y \cos x \quad (62.2.10)$$

$$\frac{\partial u}{\partial y} = \sin x \sinh y \quad \frac{\partial v}{\partial x} = -\sin x \sinh y \quad (62.2.11)$$

These derivatives exist and are continuous everywhere on  $\mathbb{C}$ . Also they satisfy the Cauchy-Riemann equations:

$$\frac{\partial u}{\partial x} \Big|_{(a,b)} = \cos a \cosh b \frac{\partial v}{\partial y} \Big|_{(a,b)} \quad \frac{\partial v}{\partial x} \Big|_{(a,b)} = \sin a \sinh b = -\frac{\partial u}{\partial y} \Big|_{(a,b)} \quad (62.2.12)$$

Since all criteria of the converse Cauchy-Riemann theorem are satisfied, we find that

$$f'(a + ib) = \cos a \cosh b - i \sin a \sinh b = \cos(a + ib) \quad (62.2.13)$$

## 62.3 Derivative rules

We now prove the chain rule and the inverse function rule. They take the same form as in real analysis. These properties are particularly useful in the computation of the derivatives of more complex functions.

### Theorem (Chain rule)

Let  $f, g$  be complex functions, let  $D$  be the domain of  $g \circ f$  and let  $\alpha$  be a limit point in  $D$ . If  $f$  is differentiable at  $\alpha$  and  $g$  at  $f(\alpha)$  then  $g \circ f$  is differentiable at  $\alpha$  and

$$(g \circ f)'(\alpha) = g'(f(\alpha))f'(\alpha) \quad (62.3.1)$$

*Proof.* Let us define the following function with the same domain as  $g$

$$h(w) = \begin{cases} \frac{g(w)-g(\beta)}{w-\beta}, & w \neq \beta \\ g'(\beta), & w = \beta \end{cases} \quad (62.3.2)$$

where  $w = f(z)$  and  $\beta = f(\alpha)$ . Notice that since  $g$  is differentiable at  $\beta$ ,  $g'(\beta) = \lim_{w \rightarrow \beta} \frac{g(w)-g(\beta)}{w-\beta}$  so  $h(w)$  is continuous at  $\beta$ . Then we see that

$$\frac{g(f(z))-g(f(\alpha))}{z-\alpha} = h(f(z))\left(\frac{f(z)-f(\alpha)}{z-\alpha}\right) \quad (62.3.3)$$

for all  $z \neq \alpha$  in the domain of  $f$ . Indeed if  $w = \beta$  then both sides vanish, while if  $w \neq \beta$  then

$$h(f(z))\left(\frac{f(z)-f(\alpha)}{z-\alpha}\right) = \frac{g(f(z))-g(\beta)}{f(z)-f(\alpha)} \frac{f(z)-f(\alpha)}{z-\alpha} \quad (62.3.4)$$

Taking the limit as  $z \rightarrow \alpha$  we find

$$\lim_{z \rightarrow \alpha} \frac{g(f(z))-g(f(\alpha))}{z-\alpha} = \lim_{z \rightarrow \alpha} h(f(z)) \lim_{z \rightarrow \alpha} \frac{f(z)-f(\alpha)}{z-\alpha} = g'(\beta)f'(\alpha) \quad (62.3.5)$$

since  $h$  is continuous at  $f(\alpha)$ ,  $g$  differentiable at  $\beta$  and  $f$  is differentiable at  $\alpha$ . ■

### Theorem (Inverse function derivative)

Let  $f : A \rightarrow B$  be invertible, and suppose  $f^{-1}$  is continuous at  $\beta \in B$ . If  $f'$  is non-zero at  $f^{-1}(\beta) \in A$  then  $f^{-1}$  is differentiable at  $\beta$  with

$$(f^{-1})'(\beta) = \frac{1}{f'(f^{-1}(\beta))} \quad (62.3.6)$$

*Proof.* We know that  $\alpha = f^{-1}(\beta)$  is a limit point of  $A$  so let  $z_n$  be a sequence in  $A - \{\alpha\}$  converging to  $\alpha$ . Then  $f(z_n)$  will be a sequence in  $B - \{\beta\}$  since  $f(z_n) \neq \alpha$  (by injectivity). Since  $f$  is continuous at  $\alpha$ ,  $f(z_n)$  converges to  $\beta$ , thus proving that  $\beta$  is a limit point of  $B$ .

Thus let  $w_n$  be a sequence in  $B - \{\beta\}$  converging to  $\beta$ . Also,  $f^{-1}$  is continuous at  $\beta$  so if we define  $z_n = f^{-1}(w_n)$  then  $z_n \rightarrow f^{-1}(\beta) = \alpha$ . Consequently

$$\lim_{n \rightarrow \infty} \frac{f^{-1}(w_n) - f^{-1}(\beta)}{w_n - \beta} = \lim_{n \rightarrow \infty} \frac{z_n - \alpha}{f(z_n) - f(\alpha)} \quad (62.3.7)$$

and since  $f$  is injective,  $f(z_n) \neq f(\alpha)$  showing that

$$(f^{-1})'(\beta) = \frac{1}{f'(f^{-1}(\beta))} \quad (62.3.8)$$

as desired. ■

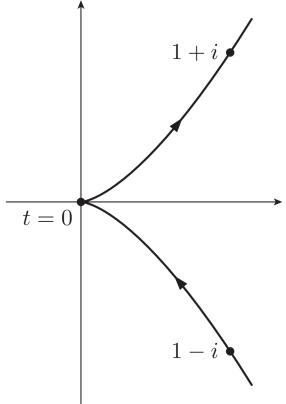
## 62.4 Smooth paths

### Definition (Path)

Let  $I \subseteq \mathbb{R}$  be an interval and let  $\gamma : I \rightarrow \Gamma \subseteq \mathbb{C}$  be a continuous function. The set  $\Gamma$  is known as the **path** parametrised by  $\gamma$ .

It is often helpful, given a parametrisation  $\gamma(t)$  of a path  $\Gamma$ , to express it in cartesian form

$$\gamma(t) = \phi(t) + i\psi(t), t \in I \quad (62.4.1)$$



where  $\phi$  and  $\psi$  are real functions on  $I$ . It is clear that the differentiability of both  $\phi$  and  $\psi$  is equivalent to the differentiability of  $\gamma$ .

Differentiability however is not enough to guarantee smoothness. For example, the following parametrisation

$$\gamma(t) = t^2 + it^3, t \in \mathbb{R} \quad (62.4.2)$$

is certainly differentiable on  $\mathbb{R}$ , but as we can see below  $\Gamma$  contains a kink at the origin where the slope is zero.

### Definition (Smooth parametrisation)

Let  $\gamma : I \rightarrow \mathbb{R}$  be a parametrisation such that

- (i)  $\gamma$  is differentiable
- (ii)  $\gamma'$  is continuous
- (iii)  $\gamma'$  is non-zero

on  $I$ . Then  $\gamma$  is a **smooth parametrisation** and the associated path  $\Gamma$  is **smooth**.

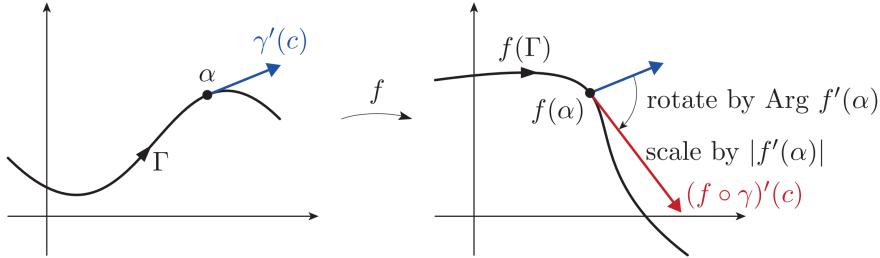
The notion of paths allows us to provide a geometrical interpretation of derivatives. Indeed let  $f$  be analytic on  $\mathcal{R}$  and let  $\Gamma$  be a smooth path in  $\mathcal{R}$  parametrised by  $\gamma : I \rightarrow \mathbb{C}$ . The image of  $\Gamma$  under  $f$ ,  $f(\Gamma)$  is then parametrised by  $f(\gamma(t))$ . We may view the derivatives of  $\gamma(t)$  and  $f(\gamma(t))$  at  $t = c$  as tangent vectors to  $\Gamma$  and  $f(\Gamma)$  at  $\gamma(c)$  and  $f(\gamma(c))$  respectively.

Consider  $\alpha = \gamma(c) \in \Gamma$  for some  $c \in I$ . Then  $f(\alpha) \in f(\Gamma)$ , and since  $\Gamma$  and  $f(\Gamma)$  are smooth it

follows that

$$(f \circ \gamma)'(c) = f'(\alpha)\gamma'(c) \quad (62.4.3)$$

So the derivatives to the two paths are related by a scaling factor of  $f'(\alpha)$ . In other words, one should rotate the tangent vector at  $\alpha$  by  $\text{Arg}(f'(\alpha))$  and scaled by  $|f'(\alpha)|$ .



Viewing the derivative of the parametrisation as tangent vectors allows us to define the angle between two curves at a point. Suppose  $\Gamma_1$  and  $\Gamma_2$  parametrised by  $\gamma_1$  ( $I_1$ ) and  $\gamma_2$  ( $I_2$ ) intersect at  $\alpha = \gamma_1(t_1) = \gamma_2(t_2)$ . Note that

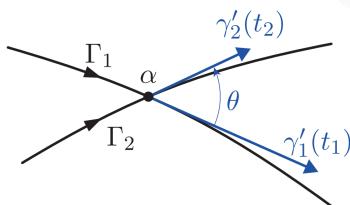
$$\gamma_2'(t_2) = \frac{\gamma_2'(t_2)}{\gamma_1'(t_1)}\gamma_1'(t_1) \quad (62.4.4)$$

so the tangent vector to one path is given by the tangent vector to the other path rotated by the argument of  $\frac{\gamma_2'(t_2)}{\gamma_1'(t_1)}$ .

### Definition (Angle between paths)

Let  $\Gamma_1$  and  $\Gamma_2$  be smooth paths with parametrisations  $\gamma_1$  and  $\gamma_2$  respectively, intersecting at  $\alpha = \gamma_1(t_1) = \gamma_2(t_2)$ . Then the angle from  $\Gamma_1$  to  $\Gamma_2$  at  $\alpha$  is

$$\theta = \text{Arg}\left(\frac{\gamma_2'(t_2)}{\gamma_1'(t_1)}\right) \quad (62.4.5)$$



There are special types of analytic functions which conserve the angle between any two paths in its domain. This angle-preserving property is known as conformality.

### Definition (Conformal functions)

Let  $f$  be analytic at  $\alpha$ . Then it is **conformal** at  $\alpha$  if the angle from any smooth path through  $\alpha$  to any other smooth path through  $\alpha$  is preserved under  $f$ .

Suppose  $f$  is analytic at  $\alpha$ . Then the tangent vectors to  $f(\Gamma)$  at  $f(\alpha)$  are given by the tangent vectors to  $\Gamma$  at  $\alpha$  through a rotation by  $\text{Arg}(f'(\alpha))$ , as long as  $f'(\alpha) \neq 0$ . Consequently the angle between two paths at a point is preserved by analytic functions whose derivative does not vanish at that point.

We will later prove the converse of this result, namely that conformal functions at a point have non-zero derivative at that point.

**Theorem (*Support of conformal function derivatives*)**

Let  $f$  be analytic at  $\alpha$ , then  $f$  is conformal at  $\alpha$  iff.  $f'(\alpha) \neq 0$ .

---

# Integrating complex functions

63

---

# **Taylor and Laurent series**

**64**

---

# **Residues**

**65**

---

# Zeros and extrema

66

---

# Conformal mappings

67

---

# **Applications to fluid flows**

**68**

---

# The Mandelbrot set and complex dynamics

69

**Part VIII**

**Calculus of Variations**

# **Part IX**

# **Fourier Analysis**

# Fourier series

## 70.1 Dirichlet conditions

The Fourier series technique may be employed to express functions that are not analytic (can't be expanded into a Taylor series) as power series nonetheless. However, a number of conditions, known as **Dirichlet conditions** must still be met.

### **Definition (Dirichlet conditions)**

The Dirichlet conditions are:

- (i) the function has finite fundamental period
- (ii) has at most a finite number of discontinuities
- (iii) finite number of stationary points per fundamental period
- (iv) the integral of  $|f(x)|$  converges over a fundamental period

where the **fundamental period** of a periodic function is the smallest  $T > 0$  such that  $f(t + T) = f(t)$ .

---

# Fourier transforms

71

---

# Convolutions

72

## **Part X**

# **Functional Analysis and Operator theory**

---

## Acknowledgments

This is the most common positions for acknowledgments. A macro is available to maintain the same layout and spelling of the heading.

**Note added.** This is also a good position for notes added after the paper has been written.

# Bibliography

- [1] Author, *Title, J. Abbrev.* **vol** (year) pg.
- [2] Author, *Title*, arxiv:1234.5678.
- [3] Author, *Title*, Publisher (year).