

Math and Programming at the Scale of Life



$$\nabla = \frac{4}{3} \pi r^3$$

$$\lim_{x \rightarrow x_0^-} f(x) := \lim_{\hbar \rightarrow 0^-} f(x_0 + \hbar)$$

$$\zeta(s) = \sum_{n=1}^{\infty} \frac{1}{n^s} \text{ for } \operatorname{Re}(s) > 1$$

$$F_b = -PgV$$

$$e := \lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n$$

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$



$$\sum_{n=1}^{\infty} \frac{1}{n^s} \text{ for } \operatorname{Re}(s) > 1$$

Author

Jessy Grizzle, Jerry W. and Carol L. Levin Professor of Engineering,
Elmer G. Gilbert Distinguished University Professor, Professor of Robotics

Cover design
T. Renay Hutchings

Inspiration

Chad Jenkins, Associate Chair of Undergraduate Studies, Professor of Robotics



Work together. Create smart machines. Serve society.

Cover design by T. Renay Hutchings, Student Administrative Assistant, Michigan Robotics

© 2025, 2024, 2023, Jessy Grizzle, Professor of Robotics, University of Michigan
JERRY W. AND CAROL L. LEVIN PROFESSOR OF ENGINEERING
ELMER G. GILBERT DISTINGUISHED UNIVERSITY PROFESSOR

Composed between April 2023 and December 2024 for use in ROB 201, Calculus for the Modern Engineer, a second-year course.

First release, August 26, 2024.

Contents

Preface	ix
Philosophy of the Course	xiii
Topic Coverage by the Traditional Paradigm	xvii
1 Pre-calculus: Notation, Functions, and Various Algebraic Facts	1
1.1 Introduction	2
1.1.1 Julia and LLM/GenAI Resources	2
1.1.2 Got Calculus Dread?	2
1.1.3 Resources for a Traditional Approach to Learning Calculus	3
1.1.4 Julia-based Sources Recognizing the Important Role of Programming in Bringing Math to Life	3
1.1.5 Efforts toward Interactive Calculus Texts	4
1.1.6 Additional Exercises	4
1.2 Notation or the Language of Mathematics	4
1.3 The Approximation Principle: The Essence of Calculus through the Lens of Irrational Numbers	8
1.3.1 Approximating π as Archimedes Did Around 200 BCE	10
1.3.2 Approximating $\sqrt{2}$ as the Greeks only Dreamed of Doing	11
1.3.3 Defining and Approximating Euler's Constant, e	14
1.3.4 Future Icons in Mathematics	15
1.4 Algebraic Manipulation and Inequalities	15
1.5 Functions, Domains, Ranges, Inverses, and Compositions	18
1.5.1 Intuitive Notion of a Function and Related Concepts	19
1.5.2 Formal Notions of Domain, Codomain, and Range of a Function	21
1.5.3 Composing Functions	23
1.5.4 Inverse Functions	24
1.5.5 Strictly Monotonic Functions and Relation to Existence of Inverse Functions	27
1.6 Trigonometric and Inverse Trigonometric Functions	30
1.6.1 Trig Functions	30
1.6.2 Trig Identities	30
1.6.3 A 3-Link Robot Kinematic Chain in the Plane	30
1.6.4 Notation for Inverse Trigonometric Functions	34
1.6.5 Understanding Inverse Trigonometric Functions	34
1.6.6 Composing Trig Functions and Inverse Trig Functions	39
1.7 Powers and Roots: Common, but Often Less Understood	39
1.7.1 n -th Roots, aka Integer Roots	40
1.7.2 Rational Exponents or Powers	43
1.7.3 Real Exponents	46
1.8 Exponentials and Logarithms	49
1.9 One Logarithm to Rule Them All	54
1.10 Set Notation and How to Make Sense of It	56
1.11 Euler's Formula	58
1.12 Hyperbolic Trig Functions and Relation to Euler's Formula	60
1.13 Summing Symbol and Change of Indices	61

1.14	Binomial Theorem	62
1.14.1	Binomial Coefficients	62
1.14.2	Expansion of $(x + y)^n$	63
1.15	Some Special Functions	65
1.16	Shifting and Scaling	66
1.17	A Peek at Feedback Control and a Partial Roadmap to this Textbook	67
1.18	(Optional Read:) Binomial Theorem meets Euler	72
1.18.1	Binomial Theorem meets Euler's Constant	72
1.18.2	Binomial Theorem and Powers of Sine and Cosine	73
1.19	(Optional Read:) Proofs Associated with the Chapter	76
2	Calculus Foundations: Proofs, Finite Sums, Limits at Infinity, and Geometric Sums	83
2.1	Setting Expectations	84
2.1.1	How to Use the Textbook	84
2.1.2	The Role of Lectures	84
2.1.3	Is this a Programming Course?	85
2.1.4	Why Have Written Homework at All?	85
2.1.5	What if You Don't Like Julia?	85
2.2	What is a Mathematical Proof? And Why are Mathematical Proofs Important?	85
2.3	Vocabulary and Helpful Notation	86
2.4	Proofs by Induction, Some Finite Sums, and Their Applications	88
2.4.1	Two Finite Sums via Induction	88
2.4.2	The Amazing Power of the Rectangle: Finding Area under a Curve	92
2.5	Limits at Infinity: The Hard Way	98
2.6	Limits at Infinity: The Easy Way	106
2.6.1	Rational Functions	106
2.6.2	Limit as an Integer Tends to Infinity	109
2.6.3	Products and Ratios of Exponential and Monomial Terms	113
2.6.4	Algebraic Operations with Infinity: Determinate and Indeterminate Forms	116
2.6.5	Limit Algebra	117
2.7	Geometric Sums	118
2.8	Countable Sets	121
2.8.1	The Rational Numbers are Countable	124
2.8.2	The Real Numbers are Uncountable	125
2.9	(Optional Read:) How to Really Define Euler's Number	126
2.10	(Optional Read:) Where to Learn Proofs at Michigan	128
2.11	(Optional Read:) Proofs Associated with the Chapter	128
3	Definite Integration as the Signed Area Under a Curve	135
3.1	The Riemann Integral (aka, Riemann-Darboux Integral)	136
3.1.1	Simple Version of Riemann Lower and Upper Sums	137
3.1.2	Riemann Integral of a Continuous Function over a Bounded Interval	139
3.1.3	Illustration of the Riemann Indefinite Integral of a Monomial	141
3.1.4	Are all Functions Riemann Integrable?	143
3.2	A Few Properties of the Riemann Integral	143
3.2.1	A Basic Additivity Property of Area Under a Curve	143
3.2.2	Dummy Variables of Integration	144
3.2.3	Integrating Linear Combinations	144
3.2.4	Making Sense of an Integral when its Lower Limit is Greater than its Upper Limit	144
3.2.5	Generalized First Additivity Property of Integrals	145
3.2.6	Shifting and Integration	145
3.2.7	Scaling and Integration	148
3.3	Numerical Methods for Approximating Riemann Integrals	150
3.3.1	Trapezoidal Rule	151

3.3.2	Simpson's Rule	157
3.3.3	Julia Packages	160
3.3.4	(Optional Read:) Even and Odd Functions as Great Test Cases	163
3.4	Applications of the Definite Integral	165
3.4.1	From Speed to Change in Position	166
3.4.2	Ballistic Motion	169
3.4.3	Area between Two Functions	172
3.4.4	Using the Power of the Rectangle to Find the Center of Mass of a (Planar) Robot Link	176
3.4.5	(Optional Read:) Using Again the Power of the Rectangle to Compute the Inertia of a (Planar) Robot Link	187
3.4.6	Solids of Revolution or More on the Power of the Rectangle	190
3.5	(Optional Read:) Proofs Associated with the Chapter	196
4	Properties of Functions: Left and Right Limits, Types of Continuity, Boundedness, and Generalizations of Max and Min	205
4.1	Intuition for Limits from the Left and the Right	206
4.2	Formal Definition of Limits from the Left and the Right	208
4.3	The Power of Continuity Done Right	220
4.3.1	Continuity at a Point and Continuous Everywhere on the Domain of Definition	220
4.3.2	Two-sided Limits and Continuous Functions	223
4.3.3	Key Properties of Continuous Functions or Why We Care about Continuity	226
4.3.4	Application to Integrating the Natural Exponential, e^x	231
4.4	The Squeeze Theorem	236
4.5	Piecewise Continuity	238
4.6	Bounded vs Unbounded Functions and Asymptotes	241
4.7	When does a Function have a Definite Riemann Integral?	245
4.8	Maximum and Minimum Values of Sets and Their Generalizations	246
4.9	Maximum and Minimum Values of Functions and Their Generalizations	249
4.10	(Optional Read:) Intermediate Value Theorem and the Mean Value Theorem	254
4.11	(Optional Read:) Proofs Associated with the Chapter	256
5	Differentiation	261
5.1	The Derivative as the Local Slope of a Function	262
5.2	The Derivative as a Local Linear Approximation of a Function	267
5.3	Software Tools for Computing Derivatives	269
5.3.1	Symbolic Derivatives	269
5.3.2	Numerical Differentiation	272
5.3.3	Automatic Differentiation	273
5.3.4	Guidelines on Choosing Differentiation Methods	275
5.4	Differentiation Rules that all Engineers are Expected to Understand	276
5.4.1	Differentiating by the Rules	276
5.4.2	The How and Why of the Rules	277
5.4.3	Bonus Examples	283
5.5	Use Cases of the Single-variable Derivative	285
5.5.1	From Position to Velocity	285
5.5.2	If the Derivative does not Change Sign, the Function is Monotonic	290
5.5.3	L'Hôpital's Rule for Evaluating Limits	291
5.5.4	Higher-order Derivatives Lead to Taylor's and Maclaurin's Expansions	296
5.5.5	(Optional Read:) Y2K or BlackPenRedPen Strikes Again	301
5.6	Partial Derivatives	302
5.6.1	Intuition	302
5.6.2	Partial Derivatives: Single-variable Derivatives meet Multivariable Functions	303
5.7	Packaging Partial Derivatives to form Jacobians, Gradients, and Hessians	305
5.7.1	The Jacobian	305
5.7.2	The Gradient	316

5.7.3	A Handy Product Rule for the Gradient and the Jacobian	317
5.7.4	(Optional Read:) The Hessian	318
5.7.5	(Optional Read:) More Matrix-Vector Calculus	319
5.8	The Total Derivative or the Chain Rule on Steroids	321
5.9	(Optional Read): Proofs Associated with the Chapter	328
6	Engineering Applications of the Derivative	331
6.1	Path Length or Arc Length	332
6.2	Root Finding and Unconstrained Minimization	338
6.2.1	Root Finding	338
6.2.2	Minimization without Constraints	343
6.2.3	Gradient Descent Algorithm	344
6.2.4	Second Derivative Tests for Local Min and Max	350
6.3	Constrained Optimization	356
6.3.1	Motivating Problems	357
6.3.2	Vocabulary of Constrained Optimization	357
6.3.3	Gradient Descent with a Single Equality Constraint	359
6.3.4	Gradient Descent with a Vector of Equality Constraints	362
6.3.5	Lagrange Multiplier for a Problem with a Single Equality Constraint	368
6.3.6	(Optional Read:) Lagrange Multipliers for a Problem with a Vector of Equality Constraints	374
6.3.7	3-Link Manipulator Meets Inequality Constraints	382
6.4	Dynamics à la Lagrange	385
6.4.1	The Dot Notation for Derivatives	386
6.4.2	Kinetic and Potential Energy of a Point Mass in the Plane	386
6.4.3	Kinetic and Potential Energy of the 3-Link Manipulator	388
6.4.4	Remarks on Simplifying Symbolic Equations	393
6.4.5	Lagrange's Equations	394
6.4.6	Lagrange's Equations of Motion for the 3-Link Manipulator	398
6.5	(Optional Read:) More on Lagrangian Dynamics and Mechanical Properties of Planar Bodies	401
6.5.1	A Crude Method to Numerically Integrate the Robot Equations	401
6.5.2	How the Robot Equations Arise from Lagrange's Equations of Motion	404
6.5.3	Including Actuators in Lagrange's Equations	405
6.5.4	Moment of Inertia	406
6.5.5	A Simple Planar Biped	408
6.6	(Optional Read:) Proof behind Lagrange Multipliers for Equality Constraints	412
6.7	(Optional Read:) Proof Underlying Gradient Descent with Equality Constraints	414
6.8	Words of Wisdom from one of the Planet's Greatest Living Mathematicians	417
7	Antiderivatives and the Fundamental Theorems of Calculus	419
7.1	Introduction	420
7.2	Uniting Integration and Differentiation through the Fundamental Theorems of Calculus	424
7.2.1	Fundamental Theorems	424
7.2.2	Using the Second Fundamental Theorem of Calculus for Definite Integration	424
7.2.3	Conflating Integration and Antiderivatives is a Pedagogical Pitfall in Calculus	427
7.3	The Art of the Antiderivative: Inverting Differentiation Rules to Find Antiderivatives	427
7.3.1	The Power Rule of Integration: Inverts the Power Rule of Differentiation	428
7.3.2	The Fundamental Rule: Integrating the Differential	428
7.3.3	Inverting the Chain Rule: Integration by Substitution, aka u-Substitution	428
7.3.4	Inverting the Product Rule: Integration by Parts	433
7.3.5	Antiderivatives of Rational Functions by Partial Fraction Expansion (PFE)	436
7.3.6	Trigonometric Substitutions for Radicals	442
7.4	Coming Up for Air and Taking Stock	450
7.5	Software Tools for Finding Antiderivatives	451
7.5.1	SymPy in Julia	451
7.5.2	More on Wolfram Alpha	453

7.5.3 UM-GPT: Show me how to Work Example 7.36	457
7.6 (Optional Read:) Proofs Associated with the Chapter	459
8 Improper Integrals	463
8.1 Introduction	464
8.2 Type-I Improper Integrals: Unbounded Limits of Integration	464
8.2.1 Comparision Test	468
8.2.2 Absolute Integrability	474
8.3 Type-II Improper Integrals: Vertical Asymptotes	476
8.4 Improper Integrals through the Lens of Probability	485
8.4.1 What is a Probability Density?	486
8.4.2 The Role of Improper Integrals	486
8.5 (Optional Read:) Proofs Associated with the Chapter	489
9 Ordinary Differential Equations	491
9.1 Introduction	492
9.2 Let's Start Simple: One equation, One Unknown, One Derivative	494
9.2.1 Analytical Solutions via Antiderivatives	495
9.2.2 Numerical Solutions	501
9.2.3 Finite Escape Time	505
9.2.4 One-Dimensional ODE with Multiple Solutions	506
9.2.5 (Optional Read:) Can an ODE Have No Solution?	507
9.2.6 (Optional Read:) The Independent Variable in an ODE can be any Strictly Increasing Quantity	508
9.3 Higher-order ODEs and a Direct Current (DC) Motor Model	508
9.4 Vector ODEs: Multidimensional Dynamics, Single Independent Variable	512
9.4.1 The Return of the Robot Equations	512
9.4.2 Solutions of First-order Vector ODEs	518
9.4.3 Existence and Uniqueness of Solutions	519
9.4.4 (Optional Read:) More General Notion of a Solution	525
9.5 Linear Systems of ODEs	525
9.5.1 Examples from Circuits and Mechanics	525
9.5.2 Expressing Higher-order Linear ODEs in First-order Form	528
9.5.3 Linearization of Nonlinear ODE Models	529
9.5.4 The Matrix Exponential	537
9.5.5 Software Tools	540
9.5.6 Solutions to Linear ODEs	542
9.5.7 Eigenvalues and Eigenvectors to the Rescue	546
9.5.8 Exponential Stability of Linear Systems of ODEs with Implications for Nonlinear ODEs	559
9.6 When in a Bind, Euler's Method is Always There for You	573
9.7 Resonance in ODEs	578
9.8 (Optional Read:) More on Julia Packages for Solving ODEs and PDEs	580
9.9 (Optional Read:) Proofs Associated with the Chapter	581
10 Laplace Transforms through the Lens of Feedback Control	587
10.1 Introduction	588
10.2 Single-Input Single-Output (SISO) Linear Systems	588
10.2.1 Input-Output and State-Variable Models	589
10.2.2 Input-Output Models from Circuits	589
10.2.3 A State-Variable Model of a (Planar) Segway Transporter	591
10.3 The Laplace Transform and an Algebraic Look at Linear ODEs	593
10.3.1 Common Laplace Transform Pairs	596
10.3.2 (Optional Read:) How to Compute Laplace Transform Pairs by Hand	596
10.3.3 Computing Inverse Laplace Transforms by Hand	599
10.3.4 Algebraic Approach to Solving Linear ODEs	603

10.3.5 Transfer Functions	609
10.3.6 The Segway Transporter à la Laplace	610
10.4 Poles, Zeros, and BIBO Stability	614
10.4.1 Poles and Zeros	614
10.4.2 BIBO Stability	615
10.5 Unity Feedback Systems	617
10.5.1 Closed-loop Transfer Functions	618
10.5.2 Closed-loop Poles and Zeros	621
10.6 Cascade Control Architecture	622
10.6.1 Two Common Compensators: Proportional (P) and Proportional-Derivative (PD)	623
10.6.2 Steady-State Error	625
10.7 Transient Response of First- and Second-order Systems	629
10.7.1 Common Performance Specifications	629
10.7.2 First-order System without a Zero	630
10.7.3 Second-order System without a Zero	631
10.7.4 Effects of Zeros	636
10.7.5 (Optional Read:) Why Zeros Alter the Transient Response	637
10.8 Design of Cascade Compensators and Pre-compensators	638
10.8.1 First-order Systems	638
10.8.2 Second-order Systems	641
10.8.3 Dealing with More Complicated Systems	644
10.9 Feedback Design for a Linearized Model of a Planar Segway Transporter	644
10.9.1 Controlling Body Lean Angle	644
10.9.2 Controlling Speed and Lean Angle	649
10.9.3 The Real Deal: Implementing the Controller on the Nonlinear Model	654
10.10(Optional Read:) The Impulse “Function”	662
10.11(Optional Read:) Robot Equations for a Planar Segway	662
10.12(Optional Read:) Proofs Associated with the Chapter	663

Preface

In the ever-evolving landscape of engineering education, calculus stands as a foundational pillar, yet its teaching methods have remained largely unchanged for decades. This textbook, a trailblazing effort in calculus education, is born out of a deep-seated need for reform. Since my tenure at the University of Michigan began in 1987, I have witnessed firsthand the challenges and limitations of the traditional calculus sequence in mathematics education, particularly within our College of Engineering (CoE).

Traditional calculus courses struggle to incorporate relevant, real-world examples. My colleagues in mathematics have been given an impossible job: present realistic examples in domains far from their fields of expertise. Students see right through this. The disconnect between mathematical theory and its real-world uses hinders the learning process and fails to fully engage students in the practical applications of calculus. A second core issue lies not in the examples used but in the selection of material and the overemphasis on lengthy, manual calculations that lack meaningful context.

This book represents a radical departure from conventional approaches, inspired by my experience developing and teaching ROB 101 *Computational Linear Algebra*. It aims to breathe life into mathematics through programming and realistic examples. In collaboration with my colleague Prof. Chad Jenkins, we have reimagined the calculus curriculum, placing integration at the beginning of the learning process. This approach, starting with sums and progressing to integration, is not only more intuitive but also lends itself well to programming applications, such as calculating the change in position of mobile robots from velocity curves.

Following integration, the book delves into differentiation (both single-variable and multivariable) and its myriad applications, culminating in a comprehensive chapter that instructors can tailor to their needs. Traditional methods of computing integrals, such as antiderivatives and substitution methods, are then introduced, and framed within the context of both manual calculations and modern symbolic tools.

The journey continues through improper integrals, ordinary differential equations (ODEs), Laplace Transforms, and Feedback Control, encompassing a breadth of topics rarely seen in a single calculus course. This ambitious curriculum is not just theoretical; it is being put to the test in a pilot program, reminiscent of the successful pilot of the year-one first-semester course, ROB 101, which included LU Factorization in its fifth week.

Integration with Michigan Robotics' Curriculum

This course is carefully designed to align with the needs of higher-level courses in Robotics. For instance, ROB 310 *Robot Sensors and Signals* uses Laplace transforms, ROB 311 *Build Robots and Make Them Move* needs dynamical models and feedback control, and both ROB 320 *Robot Operating Systems* and ROB 330 *Localization, Mapping, and Navigation* will benefit from the inclusion of inverse kinematics and constrained gradient descent, respectively. Furthermore, the foundational knowledge provided in this course will enhance a student's ability to excel in ROB 422 *Introduction to Algorithmic Robotics* and ROB 489-002/3 *Robot Control*.

Integration with Michigan Robotics' Values

This course, rooted deeply in the values of the Robotics Department, champions inclusivity and accessibility as its core principles. It is first of all important to recognize that Calculus often serves as a daunting barrier, particularly for women and underrepresented minorities, inadvertently hindering their progression within the engineering discipline. Furthermore, the disparity in high school mathematics preparation compounds the challenge of mastering calculus, affecting students' confidence and, ultimately, their election to pursue a STEM field. By integrating calculus with computational methods and real-world applications, this course is meticulously designed to dismantle these barriers. It sets reasonable expectations, elucidates concepts with clarity, and ensures that all students, regardless of their race, gender, economic, or educational background, find the subject equitable and approachable. Through this means, we are committed to making calculus not just a subject to be learned, but a gateway to empowerment and success in engineering for every student.

Prerequisites and Future Pathways

The only prerequisite for this course is ROB 101, ensuring that students are familiar with the Julia programming language and have experienced the practical applications of mathematics. To round out our Michigan students' mathematical education, they must earn 8 credits of mathematics in addition to *Calculus for the Modern Engineer*. If they already have credit for Calc I, Calc II, or Calc IV before taking Calculus for the Modern Engineer, they can still count them. To take any of these courses after completing *Calculus for the Modern Engineer* requires permission from an advisor. Students at Michigan are urged to further their mathematical education with additional courses from a [carefully curated list](#), including Math 215 *Multivariable and Vector Calculus*, Math 217 *Linear Algebra (with proofs)*, Math 312 *Applied Modern Algebra*, Math 351 *Principles of Analysis*, Math 371 (Engin 371) *Numerical Methods*, Math 412 *Abstract Algebra*, Math 416 *Theory of Algorithms*, and Math 451 *Advanced Calculus I*, a proof-based version of Calculus taught through the lens of Real Analysis in \mathbb{R}^n . In Robotics, mathematical knowledge is power, especially when you have the experience of rapidly turning it into code. This approach complements traditional calculus courses, adding depth and breadth to our students' mathematical toolkits.

Why Julia?

This textbook proudly uses Julia for its computational elements, embracing its open-source nature and cost-free availability as fundamental to our commitment to accessibility. Unlike the commonly used MATLAB, which comes with substantial licensing fees, Julia ensures that all students, regardless of their university's resources, can freely access state-of-the-art computational tools. This choice reflects our commitment to leveling the educational playing field, allowing every student to explore and excel in engineering without financial barriers.

A New Paradigm in Calculus Education

This textbook is not just a guide to calculus theory; it is a manual for applying calculus in the real world of engineering. It balances the need for theoretical understanding, computational skills, and practical application. While some manual calculations are included to bridge the gap with traditional calculus education, the emphasis is on programming-based solutions to complex, real-life engineering problems. This approach not only enhances understanding but also provides a clear, verifiable, and efficient method for problem-solving.

In summary, this textbook is more than just a collection of mathematical principles; it is a blueprint for a new way of teaching and learning calculus. It is designed not just to educate, but to inspire, opening doors to a world where calculus is not merely understood but applied with confidence and creativity.

Dependency on ROB 101 *Computational Linear Algebra*

By building on the foundation laid in [ROB 101 Computational Linear Algebra](#), students entering this course have already experienced the power of mathematics in solving real-world engineering problems, such as analyzing extensive LiDAR datasets and creating a control algorithm to balance a planar model of a Segway. This background supports the innovative approach taken in this textbook, where programming and practical application are central to understanding and utilizing calculus in engineering contexts. If ROB 101 is not assumed, then the instructor should add a Julia programming lab to the course, perhaps modifying the (open-source) [manual](#) created for ROB 101.

Acknowledgement

This textbook is dedicated to the group of students who were brave enough to take the pilot offering of ROB 201, *Calculus for the Modern Engineer: Math at the Scale of Life*, in Fall 2024. It is also in honor of the undergraduate Instructional Assistants, Kaylee Johnson (lead), Advaith (Adi) Balaji, Madeline (Maddy) Bezzina, Justin Boverhof, Elaina Mann, Reina Mezher, and Maxwell (Max) West, who collaborated with me in Winter and Summer 2024, diligently working to create homework sets, projects, and refine the initial draft of the textbook. Their invaluable insights and feedback have been instrumental in shaping this work. I am forever indebted to all of you.

ROB 201 was conceived by Prof. Chad Jenkins and the author as part of a complete undergraduate curriculum in robotics. Chad and I both thank past Associate Dean for Undergraduate Education, Prof. Joanna Mirecki Millunchick, for her encouragement in rethinking how mathematics is taught.

Jessy Grizzle

Ann Arbor, Michigan USA, August 2024

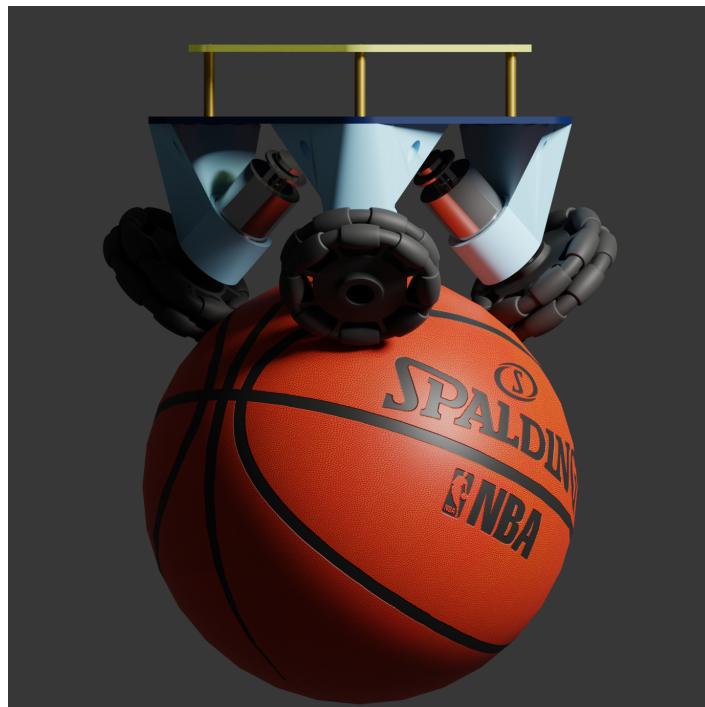
Philosophy of the Course

Talent is uniformly distributed. Opportunity is not.¹

Jessy Grizzle, University of Michigan Robotics, August 2020

The Robotics faculty seek to prepare students for the era of Information, AI, Data, and, of course, Robotics. To accomplish this, we must leave behind the Sputnik-inspired means of educating engineers, where we put our students through four semesters of Calculus before they are able to engage with any exciting engineering problems. We are replacing the standard approach with the idea that all aspects of the theory of mathematics should be brought alive through computation and compelling case studies. We are combining math education with modern computational tools and replacing most exams with accomplishment-based assessments through projects based on real engineering examples.

This book covers the essence of integral and differential Calculus of single variable functions, Jacobians and gradients of multivariable functions, and enough material on ordinary differential equations to understand the existence and uniqueness of solutions and how to compute solutions numerically. To bring the material to life, the course takes the design of a stabilizing feedback controller for the BallBot used in **ROB 311 Build Robots and Make them Move²** as an end goal. Along the way, we'll develop tools for understanding a dynamic model of the BallBot.



This book treats the following topics:

¹Trademark

²The course was designed by Prof. Elliott Rouse and piloted in Fall 2022 with the assistance of PhD student Yves Nazon and Engineer Senthur Raj Ayyappan. Image credit: Senthur Raj Ayyappan, Summer 2023.

- A pre-Calculus review with an emphasis on careful statements of the key results, with correct notation. Special attention is given to topics where students often struggle, such as domains, codomains, and ranges of functions, inverse functions, logarithms, and the Binomial Theorem. A multi-link robot manipulator is introduced, and its forward kinematics are computed using standard trigonometry. Euler's Formula is included early on due to its importance in understanding solutions of differential equations later in the textbook. The majority of this material is meant for self-study.
- As part of the pre-Calculus review, the dread that many students bring with them into a Calculus course is acknowledged, along with resources for confidence building.
- The Approximation Principle, or how to compute quantities with lower and upper bounds on their accuracy. This introduces the spirit of Calculus before the conceptual hurdles of integration and differentiation.
- The spirit of Calculus is furthered through a discussion of why mathematicians use proofs: simply put, what often seems obviously wrong or correct might be neither! Cantor's discovery of uncountable sets and his repudiation by the mathematical community of his day is used as a cautionary tale, treated at an appropriate level for Y2 students. The most used proof technique in the book, induction, is covered with application to finite sums that appear in integration. Importantly, these sums are rational functions of the number of terms in the sum. Next, the problem of approximating the area under curves is broached and solved for monomials before ever talking about limits or integration. A first introduction to limits is done in the context of rational functions, where the supporting arguments can be made transparent and where important motivating examples are already at hand from the discussion of the area under monomials. This will lessen the resistance to limits when they are done correctly in a later chapter in the context of continuous functions.
- Definite Integrals for functions of a single variable from fundamental and numerical perspectives, with application to:
 - determining a robot's path from knowledge of its speed or velocity,
 - ballistic motion,
 - area between two functions, with this topic applied to determining important parameters for building planar dynamical models of robots, such as total mass, center of mass, and moment of inertia.
- Limits are a workhorse of Calculus. And while students often dislike them, especially the epsilon-delta definitions, there is no getting around their centrality to really understanding how Calculus works. Consequently, limits are explored from both fundamental and computational perspectives. Limits are then applied to the fundamentals of Calculus itself, through the study of continuous functions showing the power of taking limits inside of continuous functions. Other key topics of a rather technical nature, such as the existence of minima and maxima find their place with this material.
- Derivatives of functions of a single variable and several variables, with emphasis on how to compute derivatives of complicated functions using modern software tools. For a practicing engineer, there is no drudgery involved with computing derivatives and students of the course understand why that is so.
- The applications of derivatives are so important and numerous that they are assembled into a separate chapter. Topics include,
 - root finding for scalar and vector functions, along with appropriate software tools,
 - gradient descent, along with appropriate software tools,
 - minimization with equality constraints, along with appropriate software tools, with application to reaching into a highly constrained space with a manipulator,
 - determining a robot's speed from knowledge of its trajectory,
 - two important quantities for building dynamical models of robots, namely, kinetic energy and potential energy, and
 - combining them to form the Lagrangian of a planar mechanical system comprised of rigid bodies,
 - Lagrange's method is introduced as a source of the robot equations!
- Antiderivatives and the two Fundamental Theorems of Calculus, with emphasis on modern software tools for doing the computations. In most Calc II courses, antiderivatives are taught as the primary way to do integration, which is a huge disservice to STEM students, who, outside of a College Campus, will never have a sufficiently simple integrand to which any of the standard methods for antiderivatives apply. Hence, antiderivatives are placed in a more realistic context of being able to understand simple intergrals very rapidly and to provide insight, with numerical tools being the real workhorse in practice. ChatGPT and the Wolfram plugin are highlighted as the modern equivalent of the “Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables,” which we all used back in the day.

- By the time Ordinary Differential Equations or ODEs roll around, examples of ODEs have occurred on three different occasions. Hence, the question is not so much what is an ODE, but how does one solve an ODE and how are ODEs used in real life case studies. One-dimensional ODEs from drag chutes and parachutes are introduced and then solved via separation of variables, until literally, a one-dimensional quadratic drag model is treated for which no closed-form solution is known. This motivates numerical solvers, which also treat high-dimensional problems for which all of the traditional methods taught in an ODE course fail. The existence and uniqueness of solutions are addressed and applied to small- and medium-sized models of systems. An intensive study of linear systems of ODEs is done, drawing on the students' background in linear algebra. The matrix exponential is analyzed, and conditions for exponential stability of an equilibrium are presented, drawing upon prior experience with eigenvalues and eigenvalues (after appropriate reminders). Jacobian linearizations of nonlinear ODEs provide a rich source of real-life examples.
- A very targeted introduction to the Laplace Transform is given, with the goal of understanding poles and zeros of transfer functions, and how PD controllers work. The main application treated in the chapter is the stabilization of a Segway mobile robot, which sets students up for their final project, designing a feedback controller for a simulation model of the BallBot.
- A broader introduction to linear feedback control is provided in an Appendix.

That's a lot for one semester; hence, we've had to be careful about what we cover versus omit so as to provide the learner with an expedited understanding of Calculus and how to apply it. To be clear, the course does not seek to replace the full four semesters of Calculus. To round out our Michigan students' mathematical education, we will require them to take two additional Math courses from a **carefully curated list**. Hopefully, many students will be inspired to take Math 215 Multivariable and Vector Calculus and proof-based courses such as, Math 217 Linear Algebra, where, already knowing quite a bit about linear algebra through ROB 101, they can focus on learning how proofs work, or Math 451 Advanced Calculus I, a proof-based version of Calculus taught through the lens of Real Analysis in \mathbb{R}^n .

Caveat: While we believe this book will prepare students uncommonly well for real engineering, you, the student, still must succeed in your Engineering courses, most of which will assume you learned Calculus in the standard manner. You'll have to compute derivatives and integrals on timed, written exams in those classes without access to modern software tools. In this textbook, we do our best to prepare you for this as well. Until our teaching philosophy is more widely adopted in Engineering, you must learn how to perform certain methods by hand in order to obtain your BS Degree in Engineering. We'll highlight some of that material in our lectures. Your ultimate source for the Calculus techniques you must know how to complete by hand will be the lectures and HW problems in your Engineering courses themselves. We are confident that the material in this book will provide an excellent reference throughout your undergraduate education.

Jessy Grizzle

Ann Arbor, Michigan USA, August 2024

Topic Coverage by the Traditional Paradigm

In the following, we re-organize the Table of Contents of Calculus for the Modern Engineer into the more traditional Calculus I through IV groupings³. The textbook does not cover infinite series and their associated convergence tests, multivariable integration and its associated topics of Green's functions, divergence, and curl, nor linear algebra, the topic of ROB 101 *Computational Linear Algebra*. It is hoped that this layout will allay fears that whole swaths of Calculus have been gutted (other than the aforementioned topics). We encourage students to take a traditional Multivariable Calculus course for the missing material.

About the Course

1. Philosophy of the Course
2. Introduction
3. Julia and LLM/GenAI Resources
4. Got Calculus Dread?
5. Resources for a Traditional Approach to Learning Calculus
6. Julia-based Sources Recognizing the Important Role of Programming in Bringing Math to Life

Pre-calculus (Mostly for Student Review, Treated in HW 1)

1. Notation or the Language of Mathematics
2. The Approximation Principle: The Essence of Calculus through the Lens of Irrational Numbers
3. Algebraic Manipulation and Inequalities
4. Functions, Domains, Ranges, Inverses, and Compositions
5. Strictly Monotonic Functions and Relation to Existence of Inverse Functions
6. Trigonometric and Inverse Trigonometric Functions
7. 3-Link Manipulator
8. Powers and Roots for Integer and Rational Exponents
9. Real Exponents
10. Exponentials and Logarithms
11. Euler's Formula
12. Hyperbolic Trig Functions and Relation to Euler's Formula
13. Binomial Theorem
14. (Optional Read:) Proofs Associated with the Chapter

³Calculus III, Multivariable Calculus, is only represented through partial derivatives and Jacobians.

Calculus I: Differentiation

1. Theoretical Notions
 - (a) Vocabulary and Helpful Notation
 - (b) What is a Mathematical Proof? And Why are Mathematical Proofs Important?
 - (c) Countable Sets
 - (d) Proofs by Induction
 - (e) Maximum and Minimum Values of Sets and Their Generalizations
 - (f) Maximum and Minimum Values of Functions and Their Generalizations
2. Specific Finite Sums, Geometric Sums, and Their Applications
3. Limits at Infinity
 - (a) Rational Functions
 - (b) Products and Ratios of Exponential and Monomial Terms
4. Finite Limits
 - (a) Intuition for Limits from the Left and the Right
 - (b) Formal Definition of Limits from the Left and the Right
 - (c) Two-sided Limits and Continuous Functions
 - (d) The Power of Continuity Done Right: Key Properties of Continuous Functions or Why We Care about Continuity
 - (e) Limit Algebra
 - (f) The Squeeze Theorem
 - (g) L'Hôpital's Rule for Evaluating Limits (taught after differentiation, of course)
5. Differentiation
 - (a) The Derivative as the Local Slope of a Function
 - (b) The Derivative as a Local Linear Approximation of a Function
 - (c) Software Tools
 - i. Symbolic Differentiation
 - ii. Numerical Differentiation
 - iii. Automatic Differentiation
 - iv. Guidelines on Choosing Differentiation Methods
 - (d) Differentiation Rules that all Engineers are Expected to Understand
 - i. Derivative of a sum
 - ii. Product Rule
 - iii. Ratio Rule
 - iv. Chain Rule
 - (e) Use Cases of the Single-variable Derivative
 - i. From Position to Velocity
 - ii. If the Derivative does not Change Sign, the Function is Monotonic
 - iii. Extreme values of functions
 - iv. Taylor's and Maclaurin's Polynomials and Series
 - (f) Typically covered in Calc III, but part of Differentiation due to ROB 101 being a Pre-requisite
 - i. Partial Derivatives
 - ii. Packaging Partial Derivatives to form Jacobians, Gradients, and Hessians
 - iii. The Total Derivative or the Chain Rule on Steroids

- (g) Engineering Applications of the Derivative
 - i. Root Finding (Vector-valued functions)
 - ii. Minimization without Constraints (multivariable, scalar-valued)
 - A. Gradient Descent Algorithm
 - B. Second Derivative Tests for Local Min and Max
 - iii. Lagrange Multipliers and Constrained Optimization
 - A. Motivating Problems and Vocabulary of Constrained Optimization
 - B. Lagrange Multiplier for a Problem with a Single Equality Constraint
 - C. (Optional Read:) Lagrange Multipliers for a Problem with a Vector of Equality Constraints
 - D. (Optional Read:) Proof behind Lagrange Multipliers for Equality Constraints
 - E. 3-Link Manipulator Meets Inequality Constraints
 - iv. Dynamics à la Lagrange
 - A. Kinetic and Potential Energy
 - B. Symbolic Computational Tools
 - C. Lagrange's Equations
- 6. (Optional Read:) Proofs Associated with the Chapter

Calculus II: Integration

- 1. The Riemann Integral (aka, Riemann-Darboux Integral)
 - (a) Simple Version of Riemann Lower and Upper Sums
 - (b) Riemann Integral of a Continuous Function over a Closed Bounded Interval
 - (c) Illustration of the Riemann Indefinite Integral of a Monomial
 - (d) Are all Functions Riemann Integrable?
- 2. Properties of the Riemann Integral
 - (a) A Basic Additivity Property of Area Under a Curve
 - (b) Integrating Linear Combinations
 - (c) Making Sense of an Integral when its Lower Limit is Greater than its Upper Limit
 - (d) Generalized First Additivity Property of Integrals
 - (e) Dummy Variables of Integration
 - (f) Shifting and Integration
 - (g) Scaling and Integration
- 3. Numerical Methods for Approximating Riemann Integrals
 - (a) Trapezoidal Rule
 - (b) Simpson's Rule
 - (c) Julia Packages
 - (d) (Optional Read:) Even and Odd Functions as Great Test Cases
- 4. Applications of the Definite Integral
 - (a) From Speed to Change in Position
 - (b) Ballistic Motion
 - (c) Area between Two Functions
 - (d) Modeling the Links in a Robot
 - (e) Solids of Revolution or the Surprising Power of the Rectangle

- (f) Path Length or Arc Length
- (g) (Optional Read:) Derivation of the Center of Mass Equations for a Continuous Object

5. Fundamental Theorems

- (a) Uniting Integration and Differentiation through the Fundamental Theorems of Calculus
- (b) Using the Second Fundamental Theorem of Calculus for Definite Integration
 - i. The Art of the Antiderivative: Inverting Differentiation Rules to Find Antiderivatives
 - ii. The Fundamental Rule: Integrating the Differential
 - iii. Inverting the Chain Rule: Integration by Substitution, aka u-Substitution
 - iv. Inverting the Product Rule: Integration by Parts
 - v. Trigonometric Substitutions for Radicals
 - vi. Antiderivatives of Rational Functions by Partial Fraction Expansion (PFE)

6. Why Conflating Integration and Antiderivatives is a Pedagogical Pitfall in Calculus

7. Unbounded Limits of Integration or Unbounded Functions

- (a) Type-I Improper Integrals: Unbounded Limits of Integration
- (b) Comparision Test
- (c) Absolute Integrability
- (d) Type-II Improper Integrals: Vertical Asymptotes

8. (Optional Read:) Proofs Associated with the Chapter

Calculus III: Multivariable Calculus

1. Jacobians, gradients, Hessians
2. Nothing more. This is a deliberate choice.

Calculus IV: ODEs

1. Introduction
 - (a) Let's Start Simple: One equation, One Unknown, One Derivative
 - i. Analytical Solutions via Antiderivatives (Separation of Variables and Integrating Factors)
 - ii. Numerical Solutions
 - iii. Finite Escape Time
 - iv. One-Dimensional ODE with Multiple Solutions
 - (b) (Optional Read:) Can an ODE Have No Solution?
 - (c) (Optional Read:) The Independent Variable in an ODE can be any Strictly Increasing Quantity
2. More Complex ODEs
 - (a) Higher-order ODEs and a Direct Current (DC) Motor Model
 - (b) Vector ODEs: Multidimensional Dynamics, Single Independent Variable
 - (c) The Return of the Robot Equations
3. What is a Solution to a First-order Vector ODE ($\dot{x} = f(x)$)
4. Existence and Uniqueness of Solutions
5. ODE Examples from a Plethora of Engineering Domains

6. Solving ODEs via Software
7. Linear Systems of ODEs
 - (a) Examples from Circuits and Mechanics
 - (b) Higher-order Linear ODEs
 - (c) Linearization of Nonlinear ODE Models
 - (d) The Matrix Exponential
 - (e) Software Tools for Linear ODEs
 - (f) Properties of Solutions to Linear ODEs
 - (g) Eigenvalues and Eigenvectors to the Rescue
 - (h) Exponential Stability of Linear Systems of ODEs with Implications for Nonlinear ODEs
8. Euler's Method
9. Resonance in ODEs
10. (Optional Read:) Proofs Associated with the Chapter

Calculus IV: Laplace Transforms

1. Setting the Stage for Developing Laplace Transforms in the Context of Feedback Control
 - (a) Single-Input Single-Output (SISO) Linear Systems
 - (b) Input-Output and State-Variable Models
 - (c) Input-Output Models from Circuits
 - (d) A State-Variable Model of a (Planar) Segway Transporter
2. The Laplace Transform and an Algebraic Approach to Linear ODEs
 - (a) Definition and a Key Property
 - (b) Common Laplace Transform Pairs
 - (c) (Optional Read:) How to Compute Laplace Transform Pairs by Hand
 - (d) Software Tools
 - (e) Transfer Functions
 - (f) The Segway Transporter à la Laplace
3. Poles, Zeros, and BIBO Stability
4. Unity Feedback Systems
 - (a) Closed-loop Transfer Functions
 - (b) Two Common Compensators: Proportional (P) and Proportional-Derivative (PD)
5. Performance Specifications
 - (a) Steady-State Error
 - (b) Transient Response of First- and Second-order Systems
6. Relating Transient Response to Poles and Zeros
 - (a) First-order System without a Zero
 - (b) Second-order System without a Zero
 - (c) Effects of Zeros

7. Design of Cascade Compensators and Pre-compensators
 - (a) First-order Systems
 - (b) Second-order Systems
 - (c) Remarks on Dealing with More Complicated Systems
8. Feedback Design for a Linearized Model of a Planar Segway Transporter
 - (a) Controlling Body Lean Angle
 - (b) Controlling Speed and Lean Angle
 - (c) The Real Deal: Implementing the Controller on the Nonlinear Model
 - (d) (Optional Read:) Robot Equations for a Planar Segway

Chapter 1

Pre-calculus: Notation, Functions, and Various Algebraic Facts

Learning Objectives

By the end of this chapter, the student should be able to:

- Recognize calculus as the science of approximations.
- Interpret mathematical notation and appreciate its efficacy in conveying complex mathematical ideas succinctly.
- Revisit and refresh knowledge on key mathematical concepts that are crucial for understanding calculus.
- Develop an understanding of the importance of algorithms in mathematical problem-solving and analysis.
- Cultivate a taste for careful and precise mathematical reasoning.
- Explore Euler's number, e , and learn how it arose from a simple everyday question.

Outcomes

Upon successful completion of this chapter, students will be able to:

- Recognize the utility and importance of mathematical notation in the precise expression of mathematical concepts.
- Observe the Approximation Principle at work through the study of numbers like π , $\sqrt{2}$, and e .
- Understand and apply the Bisection Algorithm as an example of the Approximation Principle in numerical methods.
- Review and properly apply rules for manipulating inequalities.
- Reaffirm understanding of fundamental concepts such as functions, domains, ranges, and inverses.
- Conduct a thorough examination of roots and powers and their properties.
- Review and consolidate knowledge of the key characteristics of exponential and logarithmic functions.
- Utilize Euler's Formula to simplify complex trigonometric expressions effectively.
- Revisit (or learn) the Binomial Theorem and its applications in algebraic expansions.
- Learn how to effectively apply shifting and scaling operations to functions for various analytical purposes.

1.1 Introduction

This Chapter “reviews” a selection of topics from Algebra I & II, and a tiny bit of Trigonometry, as typically taught in American High Schools. “Review” is in quotes because, for Algebra I & II, we make an attempt to state the results in more precise terms than a typical student would have seen them in High School. For Trigonometry, we provide an idea of how it relates to building mathematical models of robots, another topic not typically covered in High School.

1.1.1 Julia and LLM/GenAI Resources

Throughout the textbook, including this chapter, we’ll be taking advantage of ROB 101 *Computational Linear Algebra*, especially your knowledge of **Julia**. So that you have adequate resources at your fingertips, we put them upfront:

- **ROB 101 Specific Sites:**

- ROB 101 *Computational Linear Algebra* textbook what today’s undergraduate engineer needs to know about linear algebra to excel in courses, internships, and jobs.
- ROB 101 Lab Manual teaches Julia from scratch in the context of linear algebra.
- ROB 101 Website, ROB 101 General GitHub, and Public Textbook GitHub (in LaTeX format, for the ROB 101: *Computational Linear Algebra* course at the University of Michigan Robotics Department). You can fork it and personalize the content. **You cannot use the book to make money or trade for services, even after making your own modifications.**

- **Julia Sites**

- Julia Programming Language.
- JuliaPacakges.com (paraphrased) “Hello Julians, Big news from JuliaPackages.com. We’ve given the site a major refresh ... your perfect package is just a click away.
- JuliaHub: Ask AI. Access a Julia-speaking LLM.

- UofM Access point to ChatGPT and more. How does this help with Calculus?

- By pasting in a Julia code stub/block and error message, you can ask ChatGPT to fix it. Once ChatGPT has your code, you can just give it feedback via any new error messages.
- If you provide ChaptGPT a draft of your function and the names of key variables, you can ask it to complete the code for you. The results vary with your level of specificity. **You must ALWAYS have simple examples ready on which to test the code because ChatGPT can/does make mistakes.** With practice, you’ll become efficient at using ChatGPT as your coding buddy.
- If you give ChatGPT explicit enough instructions, you can get a good start on a draft piece of code. “Hello, ChatGPT! I’m interested in implementing the Gram-Schmidt algorithm, but with a twist. I’d like the algorithm to handle linearly dependent vectors as well. Additionally, I’d like the output to include both normalized and unnormalized orthogonal vectors. Could you provide a detailed explanation and a Julia implementation of this enhanced Gram-Schmidt process?”
- **You can also ask for background information on any topic, for example:** “Hello, ChatGPT! I’d like to explore the concepts of domains, codomains, and ranges of functions in more detail. Can you provide a comprehensive explanation of these concepts, highlighting their differences and significance? Additionally, I’d appreciate practical examples in the Julia programming language that demonstrate how to determine and work with these aspects of functions. This will help enhance my understanding of these foundational algebraic concepts” When this answer is too superficial for your needs, it is then up to you to request more information, such as “I’d like to delve deeper into the concepts of domains, codomains, and ranges of functions. Specifically, I’m interested in understanding the clear distinction between the range and the codomain of a function. Could you provide a detailed explanation, perhaps using a function where the range and codomain are different? Additionally, I’d appreciate a demonstration in Julia that visually represents this difference.” From here on, each of you will have a unique experience with ChatGPT.
- At the time when the first draft of this textbook was being composed, Google’s AI tool, (**first Bard, then Gemini**) did not provide adequate math or coding support to make it a reliable partner for your author. Keep trying. Tools are constantly evolving.

1.1.2 Got Calculus Dread?

In case you are dreading Calculus because “you feel that you are bad a math” ...

- Your author recommends reading this string of [reddit posts](#). They emphasize that math is like any other skill: it takes practice, and putting in the time to practice is easier in the context of examples that you find meaningful. In this book, we've taken a novel approach to revealing the secrets of Calculus, emphasizing its real-world applications in engineering, particularly in the realm of robotics. While traditional Calculus courses often delve deep into theory before allowing students to tackle "illustrative" problems, we believe in bringing mathematical concepts to life through computation and compelling case studies. This approach is evident in our use of the BallBot design as a central theme. However, it's worth noting that while this book offers a fresh perspective on Calculus, it doesn't necessarily cover every traditional topic in-depth, and some chapters may not be as rich in engineering examples as others. Nevertheless, we hope that the large number of solved problems for each topic will serve as a valuable review and resource for you. For those who might have skipped the preface, it's essential to understand that our goal is to equip you with practical skills for the modern era of Information, AI, Data Science, and Robotics. Yet, we also emphasize the importance of mastering traditional Calculus techniques, as they remain a foundational requirement in many engineering courses (that you must pass!)
- In [The Hard Truth About Intelligence and Learning](#), The Math Sorcerer discusses intelligence, learning, not being smart enough, how talent may have gotten you through High School, but in College, work habits win the game. The comments on this video are more informative than most. Here is one example of many: *@samuelakwensivie1734 As a person who aced almost everything in high school, the university really taught me this lesson. I had never failed in my studies as I did in the University. It was a really humbling experience. It has made me appreciate the process of failure more and has helped me improve my work ethic in my quest to become a robotics engineer.* The Math Sorcerer also provides a list of resources.
- And then there is [Is Success Luck or Hard Work?](#) by Derek Alexander Muller of Veritasium. Stick through the first minute of the video, and you'll be hooked until the end. His point: "In a competitive world, tiny advantages can make all the difference." Your author's hope: this course gives you a competitive advantage in achieving career success.
- [Online Anxiety Toolbox](#) This three-session video workshop focuses on understanding anxiety, learning strategies to manage anxiety, and developing a plan to apply the strategies on a daily basis. Speaking from experience, your author emphasizes that it is better to understand this now than at age 55! Ignoring it will only hold you back.

1.1.3 Resources for a Traditional Approach to Learning Calculus

- [Calculus Single and Multivariable, 8th Edition](#), textbook used at UofM in W-2022. It's over a thousand pages long. You will want to browse the appropriate math courses at your school to find a list of relevant topics before seeking them in this textbook. It is the most user-friendly traditional Calculus book known to your author.
- The content in [Paul's Online Math Notes](#) provides excellent tutorials and notes for different levels of mathematics: [Algebra](#), [Calculus I](#), [Calculus II](#), [Calculus III](#), [Differential Equations](#), and [Review/Extras](#). Additionally, there is a [Complete Calculus Cheat Sheet](#) which contains common facts, definitions, and properties of limits, derivatives, and integrals, mostly covering material taught in Calculus I and some from Calculus II.
- By Prof. Linda Green: [Pre-Calculus](#), [Algebra](#), [Calculus I](#), [Calculus II](#), [Calculus III](#). These videos give you access to a skilled teacher covering important topics in undergraduate mathematics, in a uniform manner. You can look through her playlists and find topics where you need more work.
- [Krista King](#) has a number of interesting playlists that you can explore at [Krista King Math](#). Her discussion style is down-to-earth and clear. She has a degree in Psychology and both a talent and passion for math. It's a great combination.

1.1.4 Julia-based Sources Recognizing the Important Role of Programming in Bringing Math to Life

- The site [Calculus With Julia](#) provides a "thorough exploration" into calculus topics using the Julia programming language. Specifically, the section on limits discusses fundamental concepts and computational approaches to understanding limits in calculus. The site employs Julia to create an interactive learning environment, enhancing the understanding of [Limits](#) and other calculus topics. If it has a drawback, it would be that the topics and examples are straight out of a traditional Calculus textbook. The plus would be its interactivity.
- The course [MTH 229: Calculus Computer Laboratory](#) was taught by [Prof. John Verzani](#) in Spring 2014 and used Julia at CUNY. Detailed course materials can be found on the official website: [MTH 229 Course Page](#). Additionally, related projects and resources are available on the [MTH 229 Projects page](#).

1.1.5 Efforts toward Interactive Calculus Texts

Interactive textbooks interleave text and programming cells, allowing learners to run examples as they read the material. The most notable efforts in this regard are¹:

- the [Active Calculus group](#), who are developing free and interactive Calculus textbooks;
- the Julia site we cited earlier, [Calculus With Julia](#);
- and [SimCalc](#).
- These build on much earlier efforts by [Callahan](#) and [Iverson](#).

One must applaud the efforts of these authors to bring Calculus alive through interactive textbooks. An important downside, however, is their faithful following of a traditional Calculus curriculum. What your textbook, *Calculus for the Modern Engineer*, does differently, besides mostly eschewing the interactive format, is:

- it changes the order of topics to fit a more seamless presentation of the core concepts of Calculus, while at the same time allowing the most challenging technical aspects of the theory to be developed at a more manageable pace;
- so that three courses can fit in one semester, selecting material that Robotics engineers will employ in coursework or practice (probably the same choices are relevant for ME and AERO);
- and, while doing enough hand computations that learners will still succeed in traditional engineering courses, it emphasizes AT SCALE computations using modern software tools, thereby equipping learners for real-world uses of Calculus.
- The book has a “tiny bit of interactivity”; in recognition and appreciation of the thousands of digestible YouTube videos on Calculus topics, it provides a curated list of linked videos that range from elementary to advanced treatments of key topics in the textbook, attempting to meet learners where they are at any given moment of their journey.

1.1.6 Additional Exercises

Learners often seek additional practice problems. Here are some sources:

- [Essential Calculus Skills Practice Workbook](#) by Chris McMullen; includes solutions.
- [Calculus Workbook](#). Lots of problems, mostly on the easier side. Solutions are not included, but your favorite LLM/ChatBot can solve them all.
- [Exercises and Problems in Calculus](#). Includes answers to odd-numbered exercises.
- Upload a portion of the textbook to your favorite LLM and ask it to generate exercises with solutions. Currently, if you try to upload the entire PDF of the textbook, you will overwhelm the system.

1.2 Notation or the Language of Mathematics

Good notation is very useful and hard to invent. What? Notation is an invention! Absolutely. In 1637, [René Descartes](#) “invented the convention of representing unknowns in equations by x, y, and z, and known quantities by a, b, and c”. Historically, cubic equations were studied in many cultures without using formulas. Can you imagine? Here’s a paraphrasing of a problem in the style of Renaissance Italian mathematicians such as Scipione del Ferro and Niccolò Fontana Tartaglia, who contributed to the understanding of cubic equations: “Suppose there is a number. You take this number and cube it, then you subtract six times the square of this number. If the result is equal to 20 times the number minus 36, what is this number?”

In modern notation, this would be written as $x^3 - 6x^2 = 20x - 36$, which can be rewritten as a standard cubic equation

$$x^3 - 6x^2 - 20x + 36 = 0. \quad (1.1)$$

Which do you find clearer? The text or the symbols? To further drive home the point, an ancient Babylonian mathematician might have described the process of solving a quadratic equation $ax^2 + bx + c = 0$ as follows: “Suppose you are asked to find a number

¹A special thanks to Prof. Mark Guzdial for sharing the references to the Active Calculus Group, SimCalc, and the pioneering efforts by Callahan and Iverson.

such that when it is multiplied by itself and then increased by the same number, sums to 10. Do the following: Take half the number that is being added (in this case, $1/2$), square it (giving $1/4$), and add that to both sides of the balance (making the left side a perfect square and the right side $41/4$). Now, find the square root of the number on the right side ($\sqrt{41}/2$, approximately 3.20256). Finally, subtract the number you originally squared (that is, $1/2$) from the square root. The result ($\sqrt{41}/2 - 1/2$, approximately 2.70) is the number you are looking for.” The process just described is essentially a verbal form of “completing the square,” a method still used for solving quadratic equations. Of course, the Babylonians

- did not have a symbol for **zero**,
- did not use negative numbers, so they did not find both answers, $-0.5 \pm \frac{\sqrt{41}}{2}$,
- they used a base 60 number system rather than base 10, so the actual numbers they wrote down would have looked very different, and
- they did not have the concept of a “general form of a problem”, meaning they could only consider a quadratic equation with specific numbers in it, such as $x^2 + x = 10$ and NOT something like $ax^2 + bx + c = 0$.

To appreciate what it was like to be a mathematical pioneer, here is the Babylonian process updated to include Descartes’ notation for unknowns, but still lacking the idea of solving a general problem: We are given a quadratic equation,

$$x^2 + x = 10. \quad (1.2)$$

The process then involves completing the square:

- take half the coefficient of the unknown (that is, $1/2$); square it (giving $1/4$); and then, add that to both sides of the equation, yielding

$$x^2 + x + \frac{1}{4} = 10 + \frac{1}{4}. \quad (1.3)$$

- This turns the left side of the equation into a perfect square, namely,

$$\left(x + \frac{1}{2}\right)^2 = \frac{41}{4}. \quad (1.4)$$

- Solve for the unknown by taking the square root of both sides (all the while ignoring that there is also a negative solution, because you are Babylonian)

$$x + \frac{1}{2} = \frac{\sqrt{41}}{2}. \quad (1.5)$$

- Finally, solve for the unknown by subtracting $1/2$ from both sides,

$$x = \frac{\sqrt{41}}{2} - \frac{1}{2}. \quad (1.6)$$

Was that any better? Your author prefers the quadratic formula,

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}. \quad (1.7)$$

Modern mathematical notation is wonderful!

Given the above as motivation, we define a set of symbols that we will frequently use in the course:

Definition 1.1. Symbols for Defining Other Symbols and Frequently Used Sets of Numbers:

(a) \coloneqq means that whatever is on the left side of “ \coloneqq ” is, **by definition**, equal to what is on the right. You may be more familiar with “ \triangleq ”, a triangle over the equals sign. As an example, $\Delta \coloneqq b^2 - 4ac$ means that Δ is by definition equal to $b^2 - 4ac$ and it is telling you that the symbols a , b , and c (which are on the side of the equals sign without the colon)

must already be known to you (in other words, they must have been defined somewhere previously in the text, or known to everyone who arrives at this level of mathematical training) and that Δ is a new symbol assuming the value $b^2 - 4ac$. The symbol $:=$ is a definition and thus is different from an equals sign, as in the equation $ax^2 + bx + c = 0$, or when $a \neq 0$, we write that the solutions to a quadratic equation as

$$x = \frac{-b \pm \sqrt{\Delta}}{2a}.$$

The symbol x is not being defined here as it was already used as the unknown in the quadratic equation $ax^2 + bx + c = 0$.

- (b) $=:$ means that whatever is on the right side of “ $=:$ ” is, by **definition**, equal to what is on the left. You cannot do this with the more familiar “ \equiv ”. We will not need this flexibility very often, but just in case!
- (c) \in means “**element of**” as in “ $x \in A$ ” stands for the entire phrase “the unknown x is an element of the set A .”
- (d) \notin means “**not an element of**” as in “ $x \notin A$ ” stands for the entire phrase “the unknown x is not an element of the set A .”
- (e) $\mathbb{N} := \{1, 2, 3, \dots\}$ the set of **natural numbers or counting numbers**. Here, you are assumed to know the numbers 1, 2, 3, and so on from your previous mathematical training. What is being defined here is the notation \mathbb{N} for the set of all counting numbers or the set of all natural numbers. Sometimes 0, the number zero, is included in the **natural numbers**, and sometimes not! Here, we are NOT including zero.
- (f) $\mathbb{Z} := \{\dots, -3, -2, -1, 0, 1, 2, 3, \dots\}$ the set of all **integers or whole numbers**. Once again, you are assumed to know already the counting numbers, their negatives, and the number zero. What is being defined here is the notation \mathbb{Z} for the set of all integers or whole numbers.

- (g) $\mathbb{Q} := \left\{ \frac{p}{q} \mid p \in \mathbb{Z}, q \in \mathbb{N}, \text{no common factors other than one (that is, reduce all fractions)} \right\}$ the set of **rational numbers**.

There is a lot going on here. First of all, because the denominator is a natural (counting) number, it cannot be zero. Secondly, in English, you would have to write “the set of rational numbers is defined to be the collection of all ratios of two integers p and q , where $q > 0$, and moreover, p and q do not have any common factors other than one.” You are assumed to know from your past mathematical training that $3/12 = 1/4$, $-1/(-4) = 1/4$, and so on. Why does one insist on removing common factors? It is so that rational numbers are unique. If we do not remove common factors, then every rational number can be written in an infinite number of ways, such as $\frac{1}{2} = \frac{k}{2k}$ for all $k \in \mathbb{N}$. In practice, we don’t worry about it too much. For example, when doing computations, a computer is just as happy with $3/12$ as it is with $1/4$.

Note: A fancier way to say “ p and q have no common factors other than one” is “the greatest common divisor of p and q is one”, which we abbreviate as $\gcd(p, q) = 1$. We recall that the **greatest common divisor** is the largest POSITIVE integer (i.e., largest counting number) that divides both p and q .

- (h) $\mathbb{R} :=$ is the set of **real numbers**. In this course, we will rely on the fact that you have used real numbers for what seems forever, and hence, they are very familiar to you. We’ll point out that you are likely familiar with the term “quantum mechanics”, but few of us really know what it means. Analogously, the first math course at Michigan that really develops the definition of real numbers is Math 451 Advanced Calculus I (it provides an introduction to Real Analysis^a). So, what are the real numbers? Very loosely speaking, they are the rational numbers unioned with the irrational numbers, where the irrational numbers are all the additional “numbers” that we can “approximate arbitrarily well by a sequence of rational numbers.” Examples include $\sqrt{2}$ and Euler’s number, e , and we have all seen that $\pi \approx 3.1415926535897932384626433832795028841971\dots$, and that any approximation of π by a rational number **can never be exact**, even if most of us are pretty good with $\pi \approx 3.14159$, or even, $\pi \approx 3.14$

The YouTube Channel **Numerphile** specializes in approachable videos about numbers and mathematics. You might start with their video on the classification of **numbers** if you’d like to learn more than the brief summary presented here.

- (i) ∞ denotes **infinity**, and stands for “**that which is boundless, endless, or larger than any natural [or real] number**”. **Infinity (∞) is a concept and not an actual number**. We use $-\infty$ to denote that which decreases without end. Sometimes we use $+\infty$ instead of ∞ to emphasize that we are concerned with increasing without bound. You may enjoy **Mathematician (Emily Riehl) Explains Infinity in 5 Levels of Difficulty** by Wired.

(j) For $a < b$ real numbers, the set $[a, b]$ is a **bounded closed interval**, (a, b) is a **bounded open interval**, while $[a, b)$ and $(a, b]$ are **bounded half-open intervals**. The sets $[a, \infty)$, $(-\infty, b]$, and (∞, ∞) are **unbounded closed intervals**, while (a, ∞) , $(-\infty, b)$, and $(-\infty, \infty)$ are **unbounded open intervals**^b. A single point and an empty set are not considered intervals. For definiteness, $[a, b] := \{x \in \mathbb{R} \mid a \leq x \leq b\}$, $[a, b) := \{x \in \mathbb{R} \mid a \leq x < b\}$, and $(-\infty, \infty) := \mathbb{R}$. The learner can fill in the other definitions. When a point a or b is included in an interval, it is called an **endpoint**. For definiteness, a and b are endpoints of $[a, b]$, while only b is an endpoint of $(a, b]$. Infinity can never be an endpoint.

(k) **Decimal expansions (also known as decimal representations) of real numbers.** Consider $[0, 1) = \{x \in \mathbb{R} \mid 0 \leq x < 1\}$ and note that the number 1.0 has been excluded from the set. You have seen that every $x \in [0, 1)$ has a decimal expansion of the form

$$x = 0.d_1d_2d_3\dots = d_110^{-1} + d_210^{-2} + d_310^{-3} + \dots =: \sum_{i=1}^{\infty} d_i10^{-i}, \quad d_i \in \{0, 1, 2, \dots, 9\},$$

where, for now, we are not formally defining the sum to infinity except to say that it stands for

$$\sum_{i=1}^N d_i10^{-i} := d_110^{-1} + d_210^{-2} + \dots + d_N10^{-N},$$

with N allowed to grow without bound.

You may or may not recall that for decimal representations to be unique, you must exclude those that terminate in an infinite string of nines. Why? Consider, for example, 0.9999... We agree that $1.0 - 0.9 = 0.1$, $1.0 - 0.99 = 0.01$, $1.0 - 0.999 = 0.001$, etc. Hence, $1.0 - 0.9999\dots = 0.0$, and therefore, $0.9999\dots = 1.0$. The same goes for $0.073999\dots = 0.074$.

(l) **i the imaginary number** whose square equals -1 , that is, $\mathbf{i}^2 := -1$. If you prefer to write this as $\mathbf{i} := \sqrt{-1}$, that is perfectly fine, as long as you already know how to compute the square root of negative one! What? You don't know how to do that? Well, that is why we define **i** to be a new symbol whose value satisfies $\mathbf{i} \cdot \mathbf{i} = -1$, just as for all $x \geq 0$, the symbol \sqrt{x} is defined to be a real number $y \in \mathbb{R}$ satisfying^c $y \cdot y = x$.

(m) $\mathbb{C} := \{\alpha + \mathbf{i}\beta \mid \alpha, \beta \in \mathbb{R}\}$ the set of **complex numbers**. You have used them to solve quadratic equations $ax^2 + bx + c = 0$, when the discriminant is negative, that is, $b^2 - 4ac < 0$. That is all we will say for now.

^aMath 451 has two complementary goals: (1) a rigorous development of the fundamental ideas of calculus, and (2) a further development of the student's ability to deal with abstract mathematics and mathematical proofs. The key words here are "rigor" and "proof;" almost all of the material of the course is geared toward understanding and constructing definitions, theorems (propositions, lemmas, etc.), and proofs. This is considered one of the more difficult among the undergraduate mathematics courses, and students should be prepared to make a strong commitment to the course. In particular, it is strongly recommended that some course which requires proofs (such as Math 412) be taken before Math 451.

^bYes, $(-\infty, \infty)$ is both open and closed. This becomes clear once you take a course on Real Analysis. For now, it is not that important.

^cFor this to make sense, you must have proven that all non-negative real numbers x there exists a real number y such that $y^2 = x$. Moreover, you eventually learn that if y satisfies this property, so does $-y$, leading us to employ \sqrt{x} and $-\sqrt{x}$ to denote their separate values.

Fact 1.2: Decimal Representation of Real Numbers

It is a fact that every real number can be written as $k + x$, for $k \in \mathbb{Z}$ and $x \in [0, 1)$. Alternatively, every non-negative real number can be written as $k + x$, for $x \in [0, 1)$ and $k \in \{0\} \cup \mathbb{N}$. The negative real numbers would then be obtained by adjoining (that is, adding in) the additive inverses of the non-negative real numbers.

Remark 1.3. Some properties.

- It is important to recall, as discussed above, that for decimal expansions to be unique, you must exclude decimal expansions that end with an infinite string of nines. You do that by rounding up by one the digit that occurs just before the infinite string of nines, yielding a number with a finite decimal expansion.
- The rational numbers exhibit two key characteristics when we consider their decimal expansions:

- (a) **Terminating Decimal Expansion:** Some rational numbers, when converted into decimal form, have a finite

number of digits after the decimal point. This means the decimal expansion ends or “terminates” after a certain point. For instance, $\frac{1}{2} = 0.5$, $\frac{3}{4} = 0.75$, or $\frac{7}{5} = 1.4$. These are known as terminating decimals.

(b) **Repeating Decimal Expansion:** Other rational numbers, when converted into decimal form, have a sequence of one or more digits that repeat indefinitely. For example, $\frac{1}{3} = 0.333\dots$ where ‘3’ repeats indefinitely, or $\frac{2}{7} = 0.285714285714\dots$ where ‘285714’ repeats indefinitely. These are known as repeating or recurring decimals.

It is a fact that every rational number can be expressed either as a terminating decimal or a repeating decimal. This is not the case for irrational numbers, which neither terminate nor repeat in their decimal expansions.

Throughout the course, we’ll use the following Latin abbreviations, and, occasionally, one American abbreviation:

Abbreviations:

- “i.e.” stands for “id est,” which translates to “that is” or “in other words.” It is used to provide clarification or a more detailed explanation.
Example: “I enjoy playing board games, i.e., games that are played on a hard surface with pieces or cards.”
- “e.g.” stands for “exempli gratia,” which means “for example” or “such as”. It is used to provide examples to illustrate a point.
Example: “I enjoy playing outdoor sports, e.g., soccer, basketball, and tennis,” Which you could translate as “I enjoy playing outdoor sports, such as soccer, basketball, and tennis,” or as, “I enjoy playing outdoor sports, for example, soccer, basketball, and tennis.”
- “etc.” is an abbreviation for the Latin phrase “et cetera,” which translates to “and other things” or “and so forth” in English. It is used at the end of a list to indicate that other similar items could be listed but, for brevity’s sake, are not.
Example: “I love eating fruits like apples, oranges, bananas, etc.” Here, “etc.” indicates that the speaker enjoys eating other fruits too, not just the ones listed.
- “aka or AKA” is an abbreviation for “also known as.” It is commonly used to introduce an alternative name or alias for a person, place, thing, or concept. It is often used to provide additional information or clarification about a particular term or entity.

Examples:

- Marie Curie, aka “Madame Curie”
- Mount Everest, aka “Sagarmatha” or “Chomolungma”
- Sodium chloride, aka “table salt”.

In each case, the term “aka” is used to indicate that there is another name or alias associated with the mentioned term or entity.

1.3 The Approximation Principle: The Essence of Calculus through the Lens of Irrational Numbers

In the following, we illustrate how to compute some basic numbers that we use “all the time”, yet the numbers cannot be written as the ratio of two integers. In each case, we will provide a means of determining lower and upper bounds for an approximation of the number in the form

$$x^{\text{low}} \leq x \leq x^{\text{up}}, \quad (1.8)$$

where x is the quantity we seek to estimate or approximate, x^{low} is a lower bound, and x^{up} is an upper bound. If we can make the lower and upper bounds “converge” to one another, i.e., become arbitrarily close to one another, we’ll then have an effective means of computing the quantity x with arbitrary accuracy. **This is the essence of Calculus.**

When we write $x^{\text{low}} \leq x \leq x^{\text{up}}$, we are not giving an explicit approximation for x . This can be remedied by defining

$$\begin{aligned} x^{\text{est}} &:= \frac{x^{\text{up}} + x^{\text{low}}}{2} \\ x^{\text{error}} &:= \frac{x^{\text{up}} - x^{\text{low}}}{2}, \end{aligned} \quad (1.9)$$

and we can write

$$x = x^{\text{est}} \pm x^{\text{error}}, \quad (1.10)$$

which we read as x is equal to x^{est} plus or minus x^{error} . The plus-minus term provides our range of “confidence” in the approximation of x . You are familiar with the \pm sign from the quadratic formula. Here, it means something a bit different: it means the value of x can be anything in the **range given by the \pm sign**. For example, $x = 1 \pm 0.2$ means that x can be anything from $1 - 0.2 = 0.8$ to $1 + 0.2 = 1.2$, i.e., anything in the closed interval $[0.8, 1.2]$.

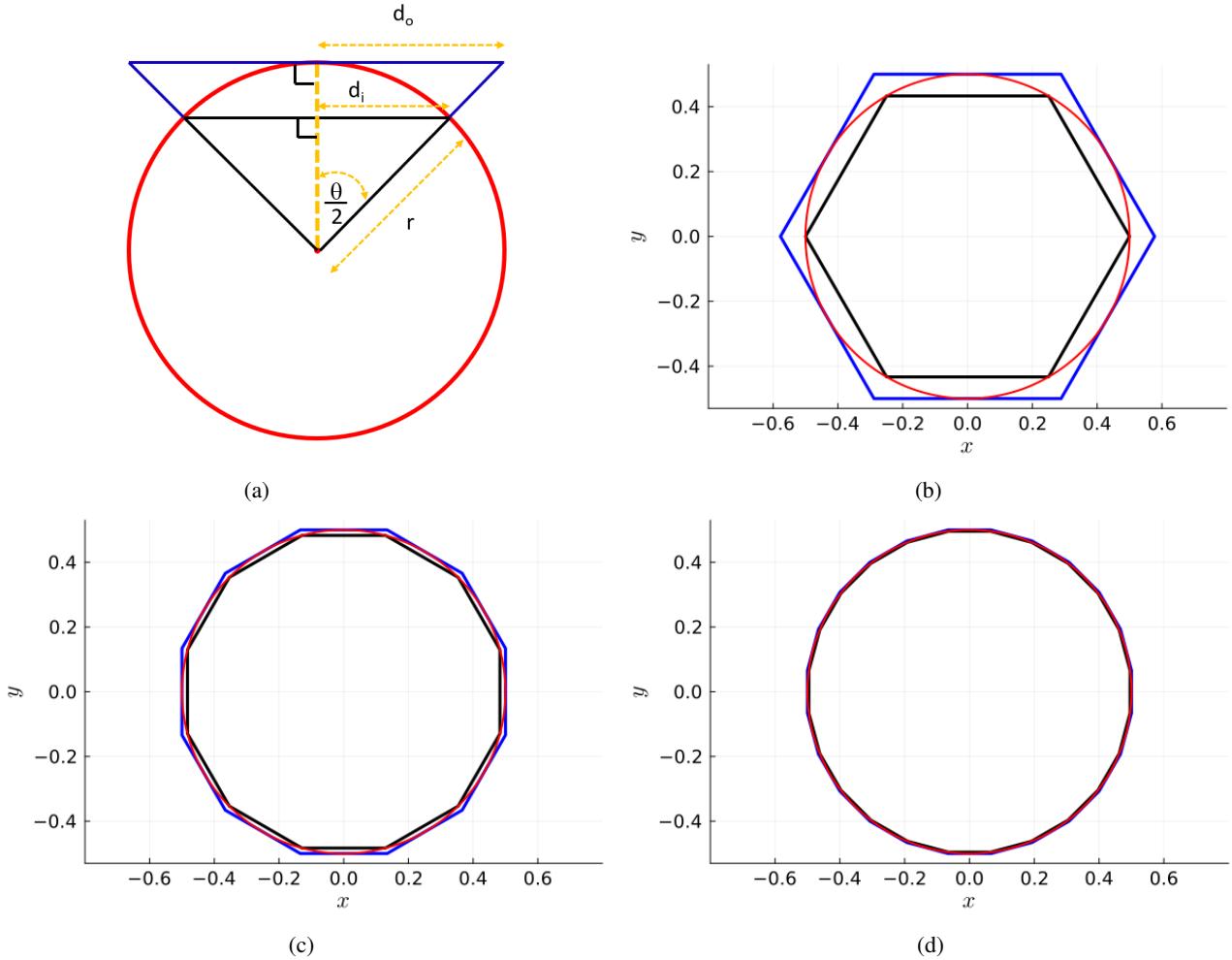


Figure 1.1: Archimedes’ approach to estimating π by approximating a circle with a polygon. (a) The radius of the circle is r , $d_i = r \sin(\theta/2)$ is half the length of the inner triangle’s base, and $d_o = r \cdot \tan(\theta/2)$ is half the length of the outer triangle’s base. The total length of the sides of the inner polygon with n sides is then $2nd_i$, which provides a **lower bound** on the circle’s circumference, while the total length of the sides of the outer polygon is $2nd_o$ and provides an **upper bound** on the circumference. As the number of sides of the polygon increases from six in (b), to 12 in (c), and 24 in (d), the lower and upper bounds become closer to one another and the true circumference of the circle. **Archimedes understood that as the number of sides tends to infinity (aka, an arbitrarily large number), the upper and lower bounds converge to one another, providing a means to compute the true circumference of the circle with arbitrary accuracy.** Not bad for 2,200 years ago!

1.3.1 Approximating π as Archimedes Did Around 200 BCE

We recall that π was **defined** by the Ancient Greeks to be the ratio of the circumference of a circle to its diameter,

$$\pi := C/D. \quad (1.11)$$

That's wonderful as a definition, but what is the value of π ? Figure 1.1 illustrates how **Archimedes** approximated it. Measuring the diameter of a circle is straightforward. It's the circumference that is hard to compute. The ancients understood right triangles. Archimedes had the brilliant idea to inscribe² a set of equally sized isosceles triangles in a circle so that he could add up their outer edges as a **lower bound** on the circumference of the circle. He then circumscribed a set of equally sized isosceles triangles and added up their outer edges as an **upper bound** on the circumference. He understood that as he used more and more triangles, he would get better and better approximations of the circumference. Moreover, he could then normalize by the diameter of the circle to obtain lower and upper bounds on an approximation of π . The method is purely geometric, and in 200 BCE, Figure 1.1 constituted proof of its well-foundedness.

Using trigonometry and the half-angle formulas for sine and tangent, Archimedes could compute everything he needed! Indeed, given a circle with radius $\frac{1}{2}$ and diameter 1, and using the trigonometric values for $\frac{\pi}{6}$, the side length of the inscribed polygon (inner approximation) and the side length of the circumscribed polygon (outer approximation) could be calculated as follows:

- For $\frac{\theta}{2} = \frac{\pi}{6} = 30$ (degrees), the exact values of sine and tangent are: $\sin\left(\frac{\pi}{6}\right) = \frac{1}{2}$ and $\tan\left(\frac{\pi}{6}\right) = \frac{\sqrt{3}}{3}$, and hence
 - **Inner Approximation:** $d_i = r \cdot \sin\left(\frac{\theta}{2}\right) = \frac{1}{2} \cdot \frac{1}{2} \implies C_i = 12 \cdot \frac{1}{4} = 3$
 - **Outer Approximation:** $d_o = r \cdot \tan\left(\frac{\theta}{2}\right) = \frac{1}{2} \cdot \frac{\sqrt{3}}{3} \implies C_o = 12 \cdot \frac{\sqrt{3}}{6} = 2\sqrt{3} \approx 3.464$
 - $\pi^{\text{est}} \approx \frac{3.464+3}{2} = 3.232 \pm 0.232$
- For $\frac{\theta}{2} = \frac{\pi}{12} = 15$ (degrees), $\frac{\pi}{12}$ can be represented as $\frac{\pi}{3} - \frac{\pi}{4}$, Archimedes could use the sum and difference formulas for sine and tangent to compute exact values, namely, $\sin\left(\frac{\pi}{12}\right) = \frac{\sqrt{6}-\sqrt{2}}{4}$ and $\tan\left(\frac{\pi}{12}\right) = 2 - \sqrt{3}$. Consequently,
 - **Inner Approximation:** $d_i = r \cdot \sin\left(\frac{\theta}{2}\right) = \frac{1}{2} \cdot \left(\frac{\sqrt{6}-\sqrt{2}}{4}\right) \implies C_i = 24 \cdot \frac{\sqrt{6}-\sqrt{2}}{8} \approx 3.106$
 - **Outer Approximation:** $d_o = r \cdot \tan\left(\frac{\theta}{2}\right) = \frac{1}{2} \cdot (2 - \sqrt{3}) \implies C_o = 24 \cdot \left(1 - \frac{\sqrt{3}}{2}\right) \approx 3.215$
 - $\pi^{\text{est}} \approx \frac{3.215+3.106}{2} = 3.161 \pm 0.0545$

Below is a Julia implementation of Archimedes' method, with illustrative results presented in Table 1.4. You are encouraged to play with the code. Using 1,000 triangles, we have estimated $\pi = 3.1415874858795636 \pm 7.75158829635636e-6$. That's not too bad! To be clear, this algorithm "converges" very slowly and is not used to compute π to billions of decimal places. To do that requires algorithms based on Calculus! Nevertheless, Archimedes' Algorithm is a valid means to compute lower and upper bounds for π . A crucial point is that we know the accuracy of the approximation. For us engineers, that can be important.

```

1 function archimedes_pi(n, radius=0.5)
2     # Circle with radius one half has a diameter of one
3     diameter = 2*radius
4     n = floor(Int, n) # force n to be an integer
5     if n < 3
6         n = 3
7     end
8     theta = 360/n
9     #
10    circumInnerApprox = radius*2*n*sind(theta/2)
11    circumOuterApprox = radius*2*n*tand(theta/2)
12    #
13    piLowerBound = circumInnerApprox/diameter
14    piUpperBound = circumOuterApprox/diameter
15    println(" $piLowerBound < pi < $piUpperBound for n = $n")

```

²Inscribe means to place inside. Circumscribe means to place outside.

```

16 piApprox = (piUpperBound + piLowerBound) / 2
17 piErrorBound = (piUpperBound - piLowerBound) / 2
18 return (piEst=piApprox, piErr=piErrorBound, piL=piLowerBound, piU=piUpperBound)
19 end
20
21 # Estimate pi using polygons with n sides
22 F = archimedes_pi(1000)

```

Output

3.141587485879563 < pi < 3.141602989056156 for n = 1000

(piEst = 3.1415952374678593, piErr = 7.75158829635636e-6,
piL = 3.141587485879563, piU = 3.141602989056156)

n	π^{low}	π^{est}	π^{up}	$\pm \pi^{\text{error}}$
6	3.0000	3.2320	3.4641	± 0.2320
24	3.1326	3.1461	3.1597	± 0.0135
128	3.1416	3.1416	3.1416	$\pm 1.3\text{e-}5$
512	3.1416	3.1416	3.1416	$\pm 8.2\text{e-}7$

Table 1.4: Archimedes estimated π in 200 BCE. By the time one reaches 24 sides on the polygon, the estimate is approximately 22/7. At 128 sides, the upper and lower bounds agree to four decimal places.

Here is a video on Archimedes' method [How Archimedes Trapped Pi](#) by Dubious Insights; you'll do something similar to this in Julia HW01. It is amazingly hard to prove that π is irrational. An animated proof can be found on [Mathologer](#).

(Optional Read:) A few Historical Remarks The ancient Indian mathematician [Madhava of Sangamagrama \(c. 1340 – c. 1425\)](#) is credited with discovering a series approximation for π , much before the invention of Calculus in Europe. His pioneering work laid the foundations for later developments in mathematical analysis. Here's a breakdown of his contributions for further exploration:

1. **Pi Approximation:** Madhava is known for deriving a series expansion for calculating the value of π , which is now known as the Madhava-Leibniz series. This series is given by:

$$\frac{\pi}{4} = 1 - \frac{1}{3} + \frac{1}{5} - \frac{1}{7} + \frac{1}{9} - \dots$$

2. **Infinite Series:** Madhava developed various infinite series expressions for trigonometric functions like sine, cosine, and arc-tangent. [His work showcased a deep understanding of the concept of infinity](#). He was the first to use infinite series approximations for a range of trigonometric functions, a precursor to the development of calculus.

3. **Legacy:** Madhava's discoveries laid the foundation for later mathematicians both within the Kerala School and eventually in Europe, where his series for π was rediscovered by Leibniz. The Madhava-Leibniz series is named in honor of both Madhava's and Leibniz's contributions to the development of Calculus.

This last point, “rediscovery”, is a constant battle in Science, even today with Google Scholar at our fingertips: results that are already known are “rediscovered” and published as new. Leibniz (1646 - 1716) did his work 300 years after Madhava; he did not have Google Scholar, so we give him a pass! During the so-called “Early Modern Period” (15th - 17th centuries), Europe and India had extensive trade relations. It seems unclear how much scientific exchange took place. Imagine where we might be today if revolutionary discoveries like Madhava's had made it to Europe during his lifetime.

1.3.2 Approximating $\sqrt{2}$ as the Greeks only Dreamed of Doing

The Ancient Greeks puzzled over the length of the hypotenuse of a right triangle when the other two sides had length one. We know from the Pythagorean Theorem that $a^2 + b^2 = c^2$ and hence $c = \sqrt{a^2 + b^2}$; well, Pythagoras was a Greek philosopher/mathematician, so they eventually figured out that the hypotenuse squared equaled two, but they hotly debated whether or not $\sqrt{2}$ could be expressed

as the ratio of two counting numbers, in other words, whether $\sqrt{2}$ was rational or not. Eventually, Euclid, another Greek philosopher/mathematician, formally proved, by a method we now call “Proof by Contradiction”, that $\sqrt{2}$ is **irrational**, meaning, it is not a **rational** number. We give the proof at the end of the Chapter in Prop. 1.90, in case you are curious about how it goes. You are not expected to know the proof.

The proof aside, how does one compute $\sqrt{2}$? (I know, it's a key on your calculator, just like π . Ha ha!) Well, literally, $\sqrt{2}$ is that number x , such that, when it is squared, you get two. In other words, $\sqrt{2}$ is the (positive) solution to

$$x^2 - 2 = 0. \quad (1.12)$$

In ROB 101 *Computational Linear Algebra*, you learned two methods for computing the roots of an equation, $f(x) = 0$. The simplest one was the Bisection Algorithm, where one starts with two values a and b that “bracket” the root (i.e., a is a lower bound for the root and b is an upper bound). You then compute the midpoint, $c := \frac{a+b}{2}$, and use it to update a or b depending on the sign of $f(a) \cdot f(c)$. The details are given in the Julia code below. A key point is that at **each step of the algorithm**, the bracketing values a and b provide **lower and upper bounds** on the root, respectively, so you actually know how accurately the root has been computed. In other words, we have an effective means for computing an approximation of $\sqrt{2}$ along with a bound on the accuracy of the estimate! (And yes, the Greeks had an algorithm for square roots, but they did not have Julia and the Bisection Algorithm.)

```

1 function bisection(f, a, b, tol)
2     if f(a) * f(b) > 0
3         error("f(a) and f(b) must have opposite signs")
4         return NaN
5     end
6     kmax = 1e5 # max number of loops so that the algorithm terminates
7     k=1
8     c=NaN #define c out of the loop so we can access it
9     while ( (b - a) / 2 > tol || abs(f(c))> tol ) && k < kmax
10        c = (a + b) / 2
11        #
12        if f(a) * f(c) < 0
13            b = c
14        else
15            a = c
16        end
17        k = k+1
18    end
19    root_ErrorBound = (b - a)/2
20    return (est=(a + b) / 2, low=a, up=b, error=root_ErrorBound)
21 end
22
23 # Define the function whose root we want to find
24 f(x) = x^2 - 2
25
26 # Call the bisection method on this function
27 F = bisection(f, 1.0, 2.0, 1e-10)
28 @show F
29
30 println(" ")
31 println("The root is approximately ", F.est)
32 println("The square of the root is ", F.est^2)

```

Output

```
F = (est = 1.414213562355144, low = 1.414213562355144, up = 1.4142135623842478,
error = 1.4551915228366852e-11)
```

The root is approximately 1.414213562355144
The square of the root is 1.999999999492266

From the above, we see that $\sqrt{2} = 1.414213562355144 \pm 1.4551915228366852e-11$. Moreover, in principle, we can determine its value to as many (finite) digits as we want.

n	$\sqrt{2}^{\text{low}}$	$\sqrt{2}^{\text{est}}$	$\sqrt{2}^{\text{up}}$	$\pm\sqrt{2}^{\text{error}}$
5	1.406250	1.421875	1.437500	± 0.01562500
10	1.414551	1.414063	1.415039	± 0.00048828
50	1.414214	1.414214	1.414214	$\pm 4.4409e-16$

Table 1.5: The Ancient Greeks discovered the existence of irrational numbers, such as $\sqrt{2}$. Here, n is the number of iterations used in the Bisection Algorithm. To generate this table, simply change the “while loop” into a “for loop”.

Remark 1.6. The Ancient Greek mathematician *Eudoxus of Cnidus* (408-355 BCE) is sometimes attributed with the first use of the Bisection Algorithm³, as part of his “method of exhaustion for calculating areas and volumes.” However, its systematic development and formalization in the context of algorithmic and numerical methods took place much later with the advent of modern computing. It’s important to note that while the method itself is ancient, the terminology “bisection method” or “algorithm” is much more recent, linked to the development of computer science and numerical methods in the 20th century. Because the Bisection Algorithm halves the distance between the bracketing values at each iteration, the lower and upper bounds will converge to one another, thereby providing an arbitrarily good approximation to $\sqrt{2}$.

```

1 function bisection2(f, a, b, n)
2     if f(a) * f(b) > 0
3         error ("f(a) and f(b) must have opposite signs")
4         return NaN
5     end
6     c=NaN #define c out of the loop so we can access it
7     for k = 1:n
8         c = (a + b) / 2
9         #
10        if f(a) * f(c) < 0
11            b = c
12        else
13            a = c
14        end
15    end
16    root_ErrorBound = (b - a) / 2
17    return (est=(a + b) / 2, low=a, up=b, error=root_ErrorBound)
18 end
19
20 # Define the function whose root we want to find
21 f(x) = x^2 - 2
22
23 # Call the bisection method on this function
24 F = bisection2(f, 1.0, 2.0, 50)
25 @show F
26
27 println(" ")
28 println("The root is approximately ", F.est)
29 println("The square of the root is ", F.est^2)
```

³Knorr, W. R. (1986). “The Ancient Tradition of Geometric Problems.” New York: Dover. This book includes a comprehensive review of the geometric techniques used by ancient Greek mathematicians, including Eudoxus.

1.3.3 Defining and Approximating Euler's Constant, e

Wikipedia recounts how in 1683, Jacob Bernoulli discovered the constant $e \approx 2.718281828\dots$ by asking a question regarding compound interest in a bank account. It is fun to see how math is discovered! Yes, facts in math are often discoveries. Bernoulli imagined he had a bank account that started with \$1.00 and paid 100 percent interest per year. Sign us all up! If the interest were credited once at year's end, the account's value would be \$2.00.

Bernoulli asked, what would happen if the interest were computed and credited more frequently during the year?

If the interest is credited twice a year, in the second half of the year, you will earn interest on the interest paid in the first half of the year in addition to the interest on your initial deposit! How does that work? Well, the interest rate for every 6 months will be 50%, so the initial \$1 is multiplied by 1.5 at the end of six months, and that amount is multiplied by 1.5 again at the end of the year. Hence, your \$1 gets multiplied by 1.5 twice, yielding $\$1.00 \times (1 + \frac{1}{2})^2 = \2.25 at the end of the year. Compounding quarterly yields $\$1.00 \times (1 + \frac{1}{4})^4 = \2.44140625 , and compounding monthly yields $\$1.00 \times (1 + \frac{1}{12})^{12} = \$2.613035\dots$. If there are n compounding intervals, the interest for each interval will be $100/n\%$ and the value at the end of the year will be $\$1.00 \times (1 + \frac{1}{n})^n$.

Bernoulli noticed that the sequence of numbers, $(1 + \frac{1}{n})^n$, approaches a limiting value with each larger value for n and, thus, more frequent, smaller compounding intervals. Compounding weekly ($n = 52$) yields \$2.692596..., while compounding daily ($n = 365$) yields \$2.714567... (approximately two cents more than when compounding weekly). The limit as n grows arbitrarily large is the number that came to be known as e . With continuous compounding, the account value would reach \$2.718281828... after one year.

$$e \approx \left(1 + \frac{1}{n}\right)^n \text{ for } n \text{ very large,} \quad (1.13)$$

with the approximation becoming better as n grows.

Remark 1.7. You may wonder why the constant is denoted “ e ” and not “ b ” since Bernoulli discovered it. Well, life is not fair, as you well know. Euler later found much more efficient algorithms for computing “Bernoulli’s constant”, and he denoted it by “ e ”. Euler’s notation stuck, and that was the end of the story.

The proof that the upper and lower bounds converge to one another is more difficult in this case than in the previous two. Archimedes' method is “geometrically intuitive”, the Bisection Algorithm” is “mathematically intuitive”, while “continuous compounding” only becomes “intuitively obvious” when you code it up, because “**Coding is Believing!**” (Prof. Chad Jenkins, circa summer 2020).

```

1 function bernoulli_e(n)
2     n = floor(Int,n) # force n to be an integer
3     if n < 2
4         n = 2
5     end
6     e_LowerApprox = (1.0 + 1.0/n)^n
7     e_UpperApprox = 1.0
8     for k = 1:n
9         e_UpperApprox = e_UpperApprox + 1.0/factorial(big(k))
10    end
11    e_UpperApprox = e_UpperApprox + 1.0/(n*factorial(big(n)))
12    e_Approx = (e_UpperApprox + e_LowerApprox)/2.0
13    e_ErrorBound = (e_UpperApprox - e_LowerApprox)/2.0
14    e_Approx = Float64(e_Approx)
15    e_ErrorBound = Float64(e_ErrorBound)
16    @show e_UpperApprox
17    @show e_LowerApprox
18    return (est=e_Approx, low=e_LowerApprox, up=e_UpperApprox, error=e_ErrorBound)
19 end

```

20
21 F = bernoulli_e(52*7*24) #Compounding hourly!

Output

```
e_UpperApprox = 2.7182818284590452353602874713526624977572
e_LowerApprox = 2.7181262654676885

(est = 2.718204046963367, low = 2.7181262654676885, up = 2.71828182845904523536,
error = 7.778149567834747e-5)
```

Hence, $e \approx 2.7182 \pm 7.77815e-5$. More information on the “many faces” of e is given by [Michael Penn](#) on his YouTube channel. A proof that e is irrational is given by [Numberphile](#), in song and verse by [TheGermanFox](#), and in what [BriTheMathGuy](#) calls the “The Most Beautiful Proof”. Intrigued enough to take a look?

n	e^{low}	e^{est}	e^{up}	$\pm e^{\text{error}}$
2	2.250000	2.500000	2.750000	± 0.250000
4	2.441406	2.580078	2.718750	± 0.138672
12	2.613035	2.665659	2.718283	± 0.052623
52	2.658970	2.688626	2.718282	± 0.029656
365	2.714567	2.716425	2.718282	± 0.001857
10^5	2.718126	2.718204	2.718282	$\pm 7.7782e-5$

Table 1.8: Bernoulli discovered the constant used in the natural logarithm, but Euler got to name it “ e ”.

1.3.4 Future Icons in Mathematics

Because Algebra was invented more than 2000 years ago and Calculus was invented in the mid 17th century (aka, mid 1600’s), you may find it hard to relate to its creators. **In case you would like to learn about people who have done amazing mathematical work in this century, you may enjoy investigating a few of the Fields Medal⁴ winners from 2022:**

- **Hugo Duminil-Copin** Lifted from the Comments: “[making mistakes is just an important component of the creative process]; I wish this poster was up there in every school!”
- **Maryna Viazovska** Lifted from the Comments: “People like you fill my heart with hope for tomorrow. Good going, Professor. Hope the scars of war will be over soon and you will be able to return to Kyiv.”
- **June Huh** Lifted from the Comments: “[Your] story is truly an inspiration to a struggling undergrad like myself, congratulations Dr. June Huh!
- **James Maynard** Lifted from the Comments: “Thank you so much for this video! It’s really helpful for me as a secondary school teacher trying to get teenagers to see maths as an attractive career choice.”

1.4 Algebraic Manipulation and Inequalities

A quick review of key properties when manipulating inequalities. We need these properties when we want to deeply understand upper and lower bounds for calculations at the heart of Calculus.

Proposition 1.9: Triangle and Reverse Triangle Inequalities, Plus a Few Other Handy Rules

For all real numbers x and y , the following hold

- (a) $|x + y| \leq |x| + |y|$ (Triangle Inequality).

⁴The [Fields Medal](#) is a prestigious award given to two to four mathematicians under 40 years old at the International Congress of the International Mathematical Union (IMU), held every four years. Named in honor of the Canadian mathematician John Charles Fields, it is often likened to the Nobel Prize of Mathematics but has significant differences, such as age limits, frequency of award, and criteria.

(b) $|x - y| \geq | |x| - |y| | = \text{abs}(|x| - |y|)$ (Reverse Triangle Inequality).

(c) $|x| \leq y \iff -y \leq x \leq y$ (No Name).

(d) For $a > 0$, $x < y \iff a \cdot x < a \cdot y$.

(e) For $a < 0$, $x < y \iff a \cdot x > a \cdot y$.

(f) For $a > 0$, $a \cdot x < y \iff x < \frac{y}{a}$ and $x < a \cdot y \iff \frac{x}{a} < y$.

(g) For $a < 0$, $a \cdot x < y \iff x > \frac{y}{a}$ and $x < a \cdot y \iff \frac{x}{a} > y$.

Proof: We'll prove the first two properties. The others are very standard, though it is always good to review them!

(a) $|x + y| \leq |x| + |y|$. Proof: If one or more of x and y is zero, the result is immediate. If x and y are both non-zero and have the same sign, then $|x + y| = |x| + |y|$. If x and y are both non-zero and have opposite signs, then $|x + y| < |x| + |y|$.

(b) $|x - y| \geq \text{abs}(|x| - |y|)$. Proof. By the triangle inequality,

$$\begin{aligned}|x| &= |x - y + y| \leq |x - y| + |y|, \text{ and} \\ |y| &= |y - x + x| \leq |y - x| + |x|,\end{aligned}$$

which yields,

$$\begin{aligned}|x| - |y| &\leq |x - y|, \text{ and} \\ |y| - |x| &\leq |y - x|.\end{aligned}$$

Because $|y - x| = |x - y|$ and $|y| - |x| = -(|x| - |y|)$, we have

$$\begin{aligned}|x| - |y| &\leq |x - y|, \text{ and} \\ -(|x| - |y|) &\leq |x - y|,\end{aligned}$$

which proves the result once (c) is established or accepted.

Remark 1.10. How to recall all these properties and ensure you use them correctly? Check special cases: $2 < 3 \implies \frac{1}{2} > \frac{1}{3}$ and $2 < 3 \implies -2 > -3$.

Example 1.11. Show the following properties dealing with inequalities.

(a) Show that if $x > 0$, then $x < y \iff 1 < \frac{y}{x}$.

(b) Show that if $x < 0$, then $x < y \iff 1 > \frac{y}{x}$.

(c) Show that if $x > 0$, then $x < y \iff \frac{1}{x} > \frac{1}{y}$. Note: Both x and y are positive.

(d) Show that if $y < 0$, then $x < y \iff \frac{1}{x} > \frac{1}{y}$. Note: Both x and y are negative.

(e) Show that if $x < 0$ and $y > 0$, then $x < y \iff \frac{1}{x} < \frac{1}{y}$. Note: x and y have opposite signs.

(f) Show that for all $x > 0$, $x + \frac{1}{x} \geq 2$.

Solutions:

(a) **Claim:** If $x > 0$, then $x < y \iff 1 < \frac{y}{x}$.

Proof: Because $x > 0$, in Prop. 1.9-(d), we can take $a := \frac{1}{x} > 0$, and deduce that $x < y \iff \frac{x}{a} < \frac{y}{a}$, which proves the result.

(b) **Claim:** If $x < 0$, then $x < y \iff 1 > \frac{y}{x}$.

Proof: Because $x < 0$, in Prop. 1.9-(e), we can take $a := \frac{1}{x} < 0$, and deduce that $x < y \iff \frac{x}{x} > \frac{y}{x}$, which proves the result.

(c) **Claim:** If $x > 0$, then $x < y \iff \frac{1}{x} > \frac{1}{y}$.

Proof: We note that $0 < x < y$ and hence both x and y are positive. In Prop. 1.9-(d), we can take $a := \frac{1}{x \cdot y} > 0$, and deduce that $x < y \iff \frac{x}{x \cdot y} < \frac{y}{x \cdot y}$, which proves the result.

(d) **Claim:** If $y < 0$, then $x < y \iff \frac{1}{x} > \frac{1}{y}$.

Proof: We note that $x < y < 0$ and hence both x and y are negative. In Prop. 1.9-(d), we can take $a := \frac{1}{x \cdot y} > 0$, and deduce that $x < y \iff \frac{x}{x \cdot y} < \frac{y}{x \cdot y}$, which proves the result.

(e) **Claim:** If $x < 0$ and $y > 0$, then $x < y \iff \frac{1}{x} < \frac{1}{y}$.

Because x and y are nonzero and have opposite signs, in Prop. 1.9-(e), we can take $a := \frac{1}{x \cdot y} < 0$, and deduce that $x < y \iff \frac{x}{x \cdot y} > \frac{y}{x \cdot y}$, which proves the result.

(f) **Claim:** If $x > 0$, then $x + \frac{1}{x} \geq 2$.

Proof: See [Sum of positive number and its reciprocal](#) by @MathVisualProofs.

■

Example 1.12. Solve for all $x \in \mathbb{R}$ satisfying $12 \geq 6 - 3x \geq -3$. The problem comes from [TC Math Academy](#).

Solution:

$$\begin{aligned} 12 &\geq 6 - 3x \geq -3 \\ &\Downarrow \\ 6 &\geq -3x \geq -9 \quad (\text{Subtracted 6 from all parts of the equation}) \\ &\Downarrow \\ -2 &\leq x \leq 3 \quad (\text{Divide by } -3, \text{ which flips the direction of the inequalities}). \end{aligned}$$

Hence, the answer is $-2 \leq x \leq 3$.

■

Example 1.13. Solve the following problems dealing with inequalities, which we will encounter in the study of “limits”.

(a) For $\epsilon > 0$, find all $x > 0$ such that $\frac{1}{x} < \epsilon$.

(b) For $\epsilon > 0$, solve for all $\delta > 0$, if any values exist, such that $|\frac{1}{3} - \frac{1}{\delta+3}| < \epsilon$.

(c) For $\epsilon > 0$, solve for all $x > 0$, if any values exist, such that $|\frac{1}{3} - \frac{1}{x^2+3}| < \epsilon$.

Solution:

(a) **Ans.** $\{x \in \mathbb{R} \mid x > 0, \frac{1}{x} < \epsilon\} = \{x \in \mathbb{R} \mid x > \frac{1}{\epsilon}\}$ by applying Example 1.11-(c); for positive quantities, when we take reciprocals, we have to flip the inequality.

(b) **Ans.** For $\epsilon > 0$, $\left\{ \delta > 0 \mid \left| \frac{1}{3} - \frac{1}{\delta+3} \right| < \epsilon \right\} = \left\{ \delta > 0 \mid \delta < \frac{9\epsilon}{1-3\epsilon} \right\}$. For this set to be nonempty, you need $\epsilon < \frac{1}{3}$ so that $1 - 3\epsilon > 0$.

We first note that for all $\delta > 0$, we have $\frac{1}{\delta+3} < \frac{1}{3}$ and hence $\left| \frac{1}{3} - \frac{1}{\delta+3} \right| = \frac{1}{3} - \frac{1}{\delta+3}$. So we can drop the absolute value, which simplifies the problem immensely. Next, it's just careful, though tedious, algebraic manipulation.

$$\begin{aligned}
 0 &< \frac{1}{3} - \frac{1}{\delta+3} < \epsilon \\
 &\Updownarrow \\
 0 &< \frac{\delta+3}{3(\delta+3)} - \frac{3}{3(\delta+3)} < \epsilon \text{ (put over common denominator)} \\
 &\Updownarrow \\
 \frac{\delta}{\delta+3} &< 3\epsilon \text{ (algebra)} \\
 &\Updownarrow \\
 \frac{\delta+3}{\delta} &> \frac{1}{3\epsilon} \text{ (positive reciprocals)} \\
 &\Updownarrow \\
 1 + \frac{3}{\delta} &> \frac{1}{3\epsilon} \text{ (simplify)} \\
 &\Updownarrow \\
 \frac{3}{\delta} &> \frac{1}{3\epsilon} - 1 \text{ (algebra)} \\
 &\Updownarrow \\
 \frac{1}{\delta} &> \frac{1-3\epsilon}{9\epsilon} \text{ (more algebra).}
 \end{aligned}$$

Hence, if $\epsilon \geq \frac{1}{3}$, then there is no solution $\delta > 0$ because $1 - 3\epsilon < 0$. For $0 < \epsilon < \frac{1}{3}$, we have

$$0 < \delta < \frac{9\epsilon}{1-3\epsilon}$$

satisfies $\frac{1}{3} - \frac{1}{\delta+3} < \epsilon$.

(c) **Ans.** This problem is nearly the same as the previous one, with δ replaced by $\frac{1}{x^2}$. Hence, we can re-use our previous work and obtain: For $\epsilon > 0$,

$$\begin{aligned}
 \left\{ x > 0 \mid \left| \frac{1}{3} - \frac{1}{\frac{1}{x^2} + 3} \right| < \epsilon \right\} &= \left\{ x > 0 \mid \frac{1}{x^2} < \frac{9\epsilon}{1-3\epsilon} \right\} \\
 &= \left\{ x > 0 \mid x^2 > \frac{1-3\epsilon}{9\epsilon} \right\} \\
 &= \left\{ x > 0 \mid x > \sqrt{\frac{1-3\epsilon}{9\epsilon}} \right\},
 \end{aligned}$$

where, for this set to be nonempty, you need $\epsilon < \frac{1}{3}$ so that $1 - 3\epsilon > 0$.

■

1.5 Functions, Domains, Ranges, Inverses, and Compositions

The topics of domains, ranges, and inverses become crucial when we compose functions and when we compute inverses of functions. To get us started, it is assumed that you are familiar with the following types of functions defined for a real number x :

- **monomials** x^k , for k a counting number.
- **polynomials** $a_nx^n + a_{n-1}x^{n-1} + \dots + a_1x + a_0$, for n a counting number and a_i real numbers.
- **basic trigonometric functions** $\sin(x)$ $\cos(x)$, and $\tan(x)$ (where $\tan(x)$ is not defined at odd multiples of $\frac{\pi}{2}$, because it “blows up” there!).

In addition, you have certainly seen **rational functions**, $\frac{n(x)}{d(x)}$, where $n(x)$, the numerator, is a polynomial and $d(x)$, the denominator, is a (non-zero) polynomial. And, you may recall that for values of x where $d(x) = 0$, a rational function is not well defined because we do not attach a meaning to the division of a real number⁵ by zero. We will also be careful about the values of x when we define $\sqrt[n]{x}$ because, unless we are willing to work with complex numbers, the even roots are not defined for negative numbers. The values of x for which a function is defined are called its **domain** and the values y can take on when $y = f(x)$ and x is in the domain are called the **range** of the function. What you may not have seen before is the term **codomain**; look for that in the following. Don’t worry; we’ll get all of this straightened out. You can also go to the [Math is Fun](#) website for additional information.

1.5.1 Intuitive Notion of a Function and Related Concepts

What is a Function: Imagine you have a machine that takes an input, does something to it, and produces an output. In mathematics and computer science, this machine is called a **function**.

- **Domain:** Think of the **domain** as a basket containing all possible inputs you can feed into the machine. Every piece in this basket is something the function can “handle”. If the machine makes soccer balls, it can likely “handle” rubber, leather, thread, and glue as inputs. If you tried to input aluminum, it would break the machine. Similarly, the function arcsine can “handle” any number $y \in [-1, 1]$, but numbers larger than one in magnitude are the function’s kryptonite and are not in its domain.
- **Codomain:** The **codomain** is like imagining every possible outcome that could ever come out of the machine, even if some outcomes don’t actually happen with the inputs you have available.
- **Function Rule:** The function itself is a specific rule that tells you how to transform an input from the domain into an output. This rule must work for every input in the domain, and for each input, there’s exactly one output.
- **Range:** The **range** is what actually comes out of the machine after you’ve fed it all the inputs in the domain. It’s a collection of the actual outcomes, which might be less than all the possibilities you imagined in the codomain. For the soccer ball machine, you may have originally set the codomain as the set of all balls (not just soccer balls), but in the end, the machine produces only high-quality soccer balls for use in professional sports, which would be its range.

Conceptual Example: Imagine a machine (function) where you put in a number (input part), and it gives back twice that number (output part).

- **Domain:** (Input bin to the machine) Any subset of the real numbers (you can double any real number); to be definite, we might take the domain as $[0, 1 \times 10^{307}]$, close to the largest number that can be represented in double-precision floating point.
- **Codomain:** (Output bin to the machine) Can be taken as large as all real numbers (doubling any real number gives a real number), or as $[0, 1.8 \times 10^{308}]$, based on the largest number that can be represented in double-precision floating point. The domain, however, cannot be smaller than $[0, 2 \times 10^{307}]$ because the codomain (output bin) must be large enough to contain the outputs of the machine (function). The machine (function) is not obliged to completely fill the output bin (codomain).
- **Function Rule:** Multiply by 2.
- **Range:** $[0, 2 \times 10^{307}]$ the part of the output bin (codomain) that is actually filled with products of the machine (numbers produced by the function).

Summary:

- The **domain** is where the function is defined, the set of points on which it is allowed to act.

⁵To be clear, not even zero divided by zero is given a definition. Eventually, we will learn “to compute a limiting value for it” under specific assumptions.

- The **codomain** can include every conceivable output value, not just the ones that actually happen.
- The **range** is where the function actually goes, the true results you see when applying the function to all points in its domain.
- Fact 1.17 illustrates why we don't "just simplify our lives" and always make the codomain equal to the range.

Intuitive Notion of Continuity: For now, we will say very informally that a function $f : \mathbb{R} \rightarrow \mathbb{R}$ is **continuous** if you can draw the graph of $y = f(x)$ on a sheet of paper without lifting your pencil (from the paper)! Figure 1.2-(a) clearly passes this test while Fig. 1.2-(b) does not.

Figure 1.2-(c) "seems" to pass the "without lifting your pencil test," but the graph does not represent a function! In High School, you may have learned the **vertical line test**: "The vertical line test is a quick way to check if a graph represents a function. If any vertical line touches the graph at more than one point, then the graph isn't a function. This is because a function should give only one output for each input." In Fig. 1.2-(c), for $x = 0$, we have $f(x) = y$ for all $y \in [-1, 2]$, which makes it not a function. What about the part of the graph where the "non-function" is constant? Does that also make it not a function? No, it is fine for the single value of $y = -1.0$ to be associated with many values of x (constant functions are functions); it's the other way around that is a problem. Functions map points to points and not points to something bigger, such as the interval $[-1, 2]$.

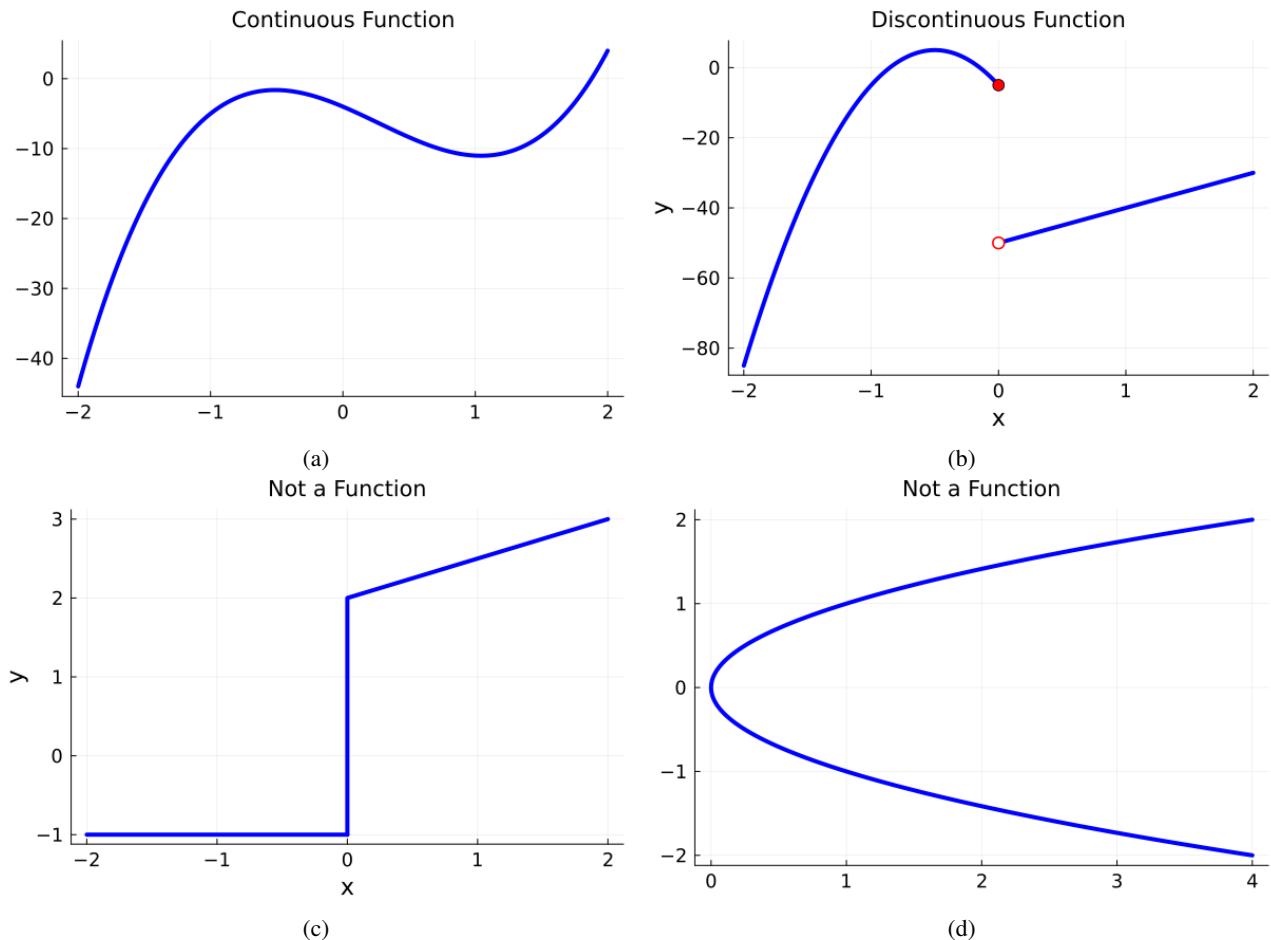


Figure 1.2: Examples of a continuous function, a discontinuous function, and two graphs that are not functions. Observe that in (b), at the jump, the value of the function has been indicated by a solid dot. Note that in (c), the point $x = 0$ is mapped to the interval $[-1, 2]$, while in (d), each $0 < x \leq 4$ is mapped to two points. To be a function, each point in the domain can only be mapped to a single point in the range.

1.5.2 Formal Notions of Domain, Codomain, and Range of a Function

Defining Functions in a Careful Manner

Definition 1.14. What is a function?

- (a) For A and B arbitrary sets, $f : A \rightarrow B$ denotes a **function from A to B** if, for all $a \in A$, $f(a) \in B$. Two things are important: (1) $f(a)$ is defined for every point $a \in A$, and (2), $f(a)$ is a single element of B , meaning it cannot contain two or more elements.
- (b) A is called the function's **domain**. It is the set of values for which the function is defined or is being defined. You can think of the domain as the "input values" of the function.
- (c) B is called its **codomain**. It is the set in which the function takes its values; for all $x \in A$, we have $f(x) \in B$. You can think of the codomain as the POTENTIAL "output values" of the function.
- (d) $C \subset B$ is the **range** of the function f if two things hold
 - (a) $\forall x \in A$, it is true that $f(x) \in C$ (C contains all of the values of f).
 - (b) $\forall y \in C$, there exists $x \in A$ such that $f(x) = y$ (every element of C is an output of f).

You can think of the range as the ACTUAL "output values" of the function because every point in C is an output of the function. Alternatively, the range is the set of points in B that are "landing points from A under the action of the function f ". If f does not "land on" $b \in B$, then b is not in the range of f .

Note: The range of $f : A \rightarrow B$ is the smallest possible codomain of the function. In other symbols, the range is

$$f(A) := \{f(a) \mid a \in A\}.$$

Looking Ahead: The difference between range and codomain becomes important when we consider the inverse of a function. For example, $\arcsin(x)$, the inverse function of $\sin(x)$, is not defined⁶ for $|x| > 1$. When considering the inverse of the sine function, we must be careful about how we define the function's domain and codomain. The same for x^2 ; its inverse function is not defined for $x < 0$. In other situations, it is fine to be much more relaxed about specifying the domain and codomain of a function.

Potentially helpful videos:

- [How to Find the Domain of Any Function](#) by Nancy Pi.
- [Domain and Range – Get Ready to Understand!](#) by TabletClass Math; easy to understand lesson on what the domain and range of a function represent.

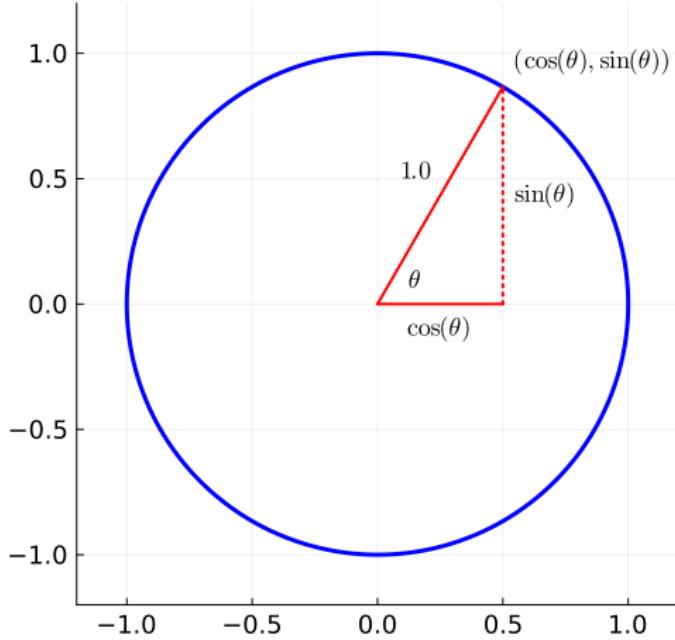
Example 1.15. We can write $\sin : \mathbb{R} \rightarrow \mathbb{R}$, and in this case, the sine function's domain and codomain have both been defined to be the set of real numbers, \mathbb{R} . What is the function's range?

Solution: We know that for all $x \in \mathbb{R}$,

$$-1 \leq \sin(x) \leq 1.$$

Moreover, from the actual definition of the sine function illustrated in the image below,

⁶Unless one allows complex numbers, $-1 \leq \sin(x) \leq 1$.



it follows that for every $y \in [-1, 1]$, there is a real number $x \in [-\frac{\pi}{2}, \frac{\pi}{2}]$ such that $\sin(x) = y$. In particular,

- $\sin(\frac{\pi}{2}) = 1$
- $\sin(-\frac{\pi}{2}) = -1$, and
- by selecting $-\frac{\pi}{2} < x < \frac{\pi}{2}$, all the values between -1 and 1 are obtained by $\sin(x)$.

The set $[-1, 1]$ is therefore the **range** of the function $\sin(x)$.

The set $[-1, 1]$ is the smallest subset of the codomain that contains all the values of $\sin(x)$. Hence, even though $[-2, 4] \subset \mathbb{R}$, and $\sin(x) \in [-2, 4]$ for all $x \in \mathbb{R}$, $[-2, 4]$ is **not the range of the sine function** because there does not exist any value of $x \in \mathbb{R}$ such that $\sin(x) = -2$. Moreover, $(0, 1]$ is **not the range** of the sine function because $\sin(-\frac{\pi}{2}) = -1$, and $-1 \notin (0, 1]$. ■

Bottom line: The range is a **Goldilocks notion**: it cannot be too big, and it cannot be too small; it cannot contain extra points on which the function f never lands, and it must be big enough to contain all of the values of f . ■

Example 1.16. As another example, consider the function, x^2 , which is defined for all real numbers. Hence, its domain can be taken as \mathbb{R} . Its values are real numbers, so we can write $(\bullet)^2 : \mathbb{R} \rightarrow \mathbb{R}$ (note that the “bullet point or dot” is just indicating a placeholder where you can plug in a value to be squared). We could also write

$$(\bullet)^2 : \mathbb{R} \rightarrow [0, \infty).$$

Both are correct. In the second case, we have used the function’s range as its codomain, which is really beautiful. ■

Fact 1.17: Range vs Codomain

Why don’t we always write functions as $f : A \rightarrow B$, where B has to be the **range** of the function?

The simple reason is that we **often do not know the range of a given function**. Consider $f : \mathbb{R} \rightarrow \mathbb{R}$ by

$$f(x) := \text{atan}\left(\sin(x^3 + 4) - 2^{-1-x^2}\right), \quad (1.14)$$

where $\text{atan}(x)$ is the inverse of the tangent function. What is the range of f ? It would take a lot of work to find out the

true answer, though a plot as in Fig. 1.3 can be very helpful to get a rough idea. Often, it is much easier just to say that the function takes values in the real numbers, and let it go at that, as a matter of convenience!

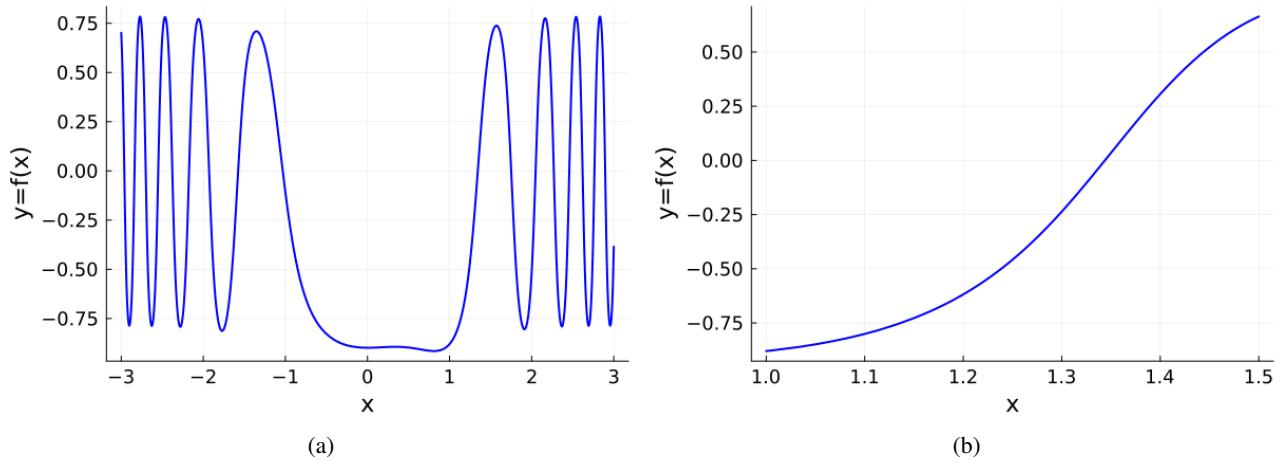


Figure 1.3: A graph of the function $f : A \rightarrow \mathbb{R}$ for $f(x) := \text{atan}(\sin(x^3 + 4) - 2^{-1-x^2})$ for two choices of its domain. In (a), the domain is $A = [-3.0, 3.0]$, while in (b), the domain is $A = [1.0, 1.5]$.

Example 1.18. Create a valid domain for the function $\tan(x) := \frac{\sin(x)}{\cos(x)}$.

Solution: Let $A := \{x \in \mathbb{R} \mid x \neq \frac{\pi}{2} \bmod \pi\}$, that is, $x \neq \frac{\pi}{2} + k\pi, \forall k \in \mathbb{Z}$, or said yet another way, x is not an odd multiple of $\frac{\pi}{2}$. Then $\tan : A \rightarrow \mathbb{R}$ is a function, because $\cos(x) \neq 0, \forall x \in A$. We have avoided division by zero. Yeah! The real numbers, by the way, are both a valid codomain and range for the function $\tan(x)$. ■

1.5.3 Composing Functions

Function composition is a fundamental concept in mathematics and computer programming. It refers to the process of applying one function to the result of another. In essence, we are “chaining” functions together so that the output of one becomes the input of the next. In a “transformer neural network”, such as the ones used in ChatGPT (2022), the chains are up to 96 units long.

The Composition of Two or More Functions

Definition 1.19. Given a “first function, $f : A \rightarrow B$, and a “second” function $g : C \rightarrow D$, their **composition** is denoted by $(g \circ f)(x)$ or simply $g(f(x))$. For the composition $g \circ f$ to be defined (which means to make sense, which in turn means not to be nonsense), **the range of f must be contained within the domain of g** .

In simple words, every output of f must be a valid input for g . This is crucial. If f produces a single value that g cannot accept, then their composition is undefined. In Julia, you are lucky; an undefined composition will (eventually) generate an error. In normal math, you are far less fortunate: you simply generate nonsense until someone finds the error. You and others continue treating the formula with respect when it is a lie!

In other symbols, given functions $f : A \rightarrow B$ and $g : C \rightarrow D$,

- (a) their composition is defined if $f(A) \subset C$, that is, the range of f is contained in C , the domain of g .
- (b) The function composition is $g \circ f : A \rightarrow D$, meaning $g \circ f$ has domain A and codomain D .

Note: Suppose $f_i : A_i \rightarrow B_i$, for $1 \leq i \leq N$. Then $f_N \circ f_{N-1} \circ \dots \circ f_1 : A_1 \rightarrow B_N$ as long as $f(A_i) \subset A_{i+1}$ for $1 \leq i \leq N-1$. Trying to write this out using the parenthesis-based notation for composition would be exciting, to say the least!

Example 1.20. Consider two functions:

- $f : (0, \infty) \rightarrow \mathbb{R}$ by $f(x) = \frac{1}{x}$, and
- $g : \mathbb{R} \rightarrow \mathbb{R}$ by $g(x) = x^2$.

(a) Does f composed with g make sense, and if it does, what is the resulting function?

(b) Does g composed with f make sense, and if it does, what is the resulting function?

Solution:

(a) The range of f is $(0, \infty)$. The function g is defined for all real numbers. Hence, their composition, $g \circ f$, makes sense.

The composition $g \circ f : (0, \infty) \rightarrow \mathbb{R}$ by $(g \circ f)(x) = g(f(x)) = g\left(\frac{1}{x}\right) = \frac{1}{x^2}$.

(b) The range of g is $\mathbb{R} = (-\infty, \infty)$, which contains zero. The function f is undefined for $x = 0$. Hence, their composition, $f \circ g$, does NOT make sense. The solution stops here because we do not need to compute something that is not defined!

Remark 1.21. Key Takeaways on Function Composition:

- Function composition allows us to create new functions by chaining existing ones.
- The order matters! In general, $g \circ f$ is not the same as $f \circ g$. In fact, one can be defined, and the other not.
- For the composition to be defined, the range of the INNERMOST function must be contained within the domain of the NEXT function. This way of thinking about it avoids the difficult task of saying which function is FIRST when writing $(g \circ f)(x) = g(f(x))$ and $(f \circ g)(x) = f(g(x))$.
- The best way to remember it is: we need g to accept $f(x)$ for $g(f(x))$ to make sense, and we need f to accept $g(x)$ for $f(g(x))$ to make sense. A function only accepts values in its domain.

Function composition is a powerful tool in mathematics and programming, enabling us to build complex operations from simpler ones.

1.5.4 Inverse Functions

Inverse of a Function

Definition 1.22. Let A and B be arbitrary sets and $f : A \rightarrow B$ a function. A function $g : B \rightarrow A$ is the **inverse of f** if,

- for all $a \in A$, $g(f(a)) = a$, and
- for all $b \in B$, $f(g(b)) = b$.

Both properties must hold.

Because it can be shown that if an inverse function exists, it is necessarily unique (meaning there cannot be more than one of them), we use the notation $f^{-1} : B \rightarrow A$. We say that $f : A \rightarrow B$ is **invertible** if it has an inverse function. We note that

$f : A \rightarrow B$ is the inverse of $f^{-1} : B \rightarrow A$, i.e., $(f^{-1})^{-1} = f$, as long as f^{-1} exists.

Note: The math community is counting on you to join them in NOT interpreting $f^{-1}(x)$ as $\frac{1}{f(x)}$. In case you are curious or just have an ornery streak, there is no (differentiable) function $f : \mathbb{R} \rightarrow \mathbb{R}$, where the function's inverse and reciprocal are the same function, meaning, very rarely, if you make this mistake, will it all work out in the end!

The following example relates the computation of inverse functions to root finding, a problem we tackled in Chapter 11 of ROB 101 *Computational Linear Algebra* (as highlighted earlier in this Chapter when we computed $\sqrt{2}$).

Example 1.23. Consider the function $f : [1.0, 1.5] \rightarrow \mathbb{R}$ by $f(x) := \text{atan}(\sin(x^3 + 4) - 2^{-1-x^2})$, with plot given in Fig. 1.3-(b). Compute $f^{-1}(y)$ for $y = -0.5$.

Solution: Wow! That looks nearly impossible to do. Analytically, yes, but with algorithms, it's quite easy. Our problem is to find x such that $f(x) = y$, which is equivalent to $f(x) - y = 0$. In other words, computing the inverse $f^{-1}(y)$ is the same as computing a root of $f(x) - y$. From Fig. 1.3-(b), we can see that $y = -0.5$ when $x \approx 1.25$, in other symbols, $f^{-1}(-0.5) \approx 1.25$. Below, we use the Bisection Algorithm to compute a more precise value than we can eyeball from the graph, and obtain that $f^{-1}(-0.5) = 1.2386 \pm 1.907e-6$.

```

1 f(x) = atan(sin(x^3 + 4) - 2.0^(-1 - x^2))
2
3 # Call the bisection method on this function
4 y = -0.5
5 h(x) = f(x) - y
6
7 # for
8 xmin = 1
9 xmax = 1.5
10
11 # One line of code
12 F = bisection(h, xmin, xmax, 1e-5)
13 @show F
14
15 y_est = f(F.est)
16
17 # Print relevant information
18 using IJulia
19 latex_string = "The value of \$f^{-1}(y = 0.5) = \$ (F.est) \$"
20 display("text/latex", latex_string)
21
22 latex_string = "The quality of the inverse is \$\\left(f^{-1}(y = 0.5)\\right) = \$ (y_est) \$ is pretty
23 display("text/latex", latex_string)

```

Output

```
F = (est = 1.2386131286621094, low = 1.2386093139648438, up = 1.2386131286621094, error = 1.907e-6)
```

The value of $f^{-1}(y = 0.5) = 1.2386112213134766$. The quality of the inverse is pretty good: $f(1.2386112213134766)) = -0.5000045415576979$; ideally, $f(f^{-1}(0.5)) = 0.5$

In Algebra II, you likely learned conditions for an inverse to exist. Let's first review them.

More Function Vocabulary

Definition 1.24. Let A and B be arbitrary subsets and $f : A \rightarrow B$ a function.

- f is **one-to-one or 1:1** if for all $a_1, a_2 \in A$, $a_1 \neq a_2 \implies f(a_1) \neq f(a_2)$. One also says the function is **injective**.
- f is **onto** if for all $b \in B$, there is some $a \in A$ such that $f(a) = b$. One also says the function is **surjective**.

Notes: When both A and B are subsets of the real numbers, you may also have been taught the **horizontal line test** for a function to be injective (i.e., 1:1). It goes like this: Let $A \subset \mathbb{R}$ and $B \subset \mathbb{R}$ be subsets of \mathbb{R} and let $f : A \rightarrow B$ be a function represented by its graph. The function f is one-to-one (or injective) if and only if no horizontal line intersects the graph of f more than once. And while, it is a perfectly valid test when A and B are subsets of \mathbb{R} , it fails in more general situations.

These definitions can be stated in different ways, some of which may be helpful to you, some of which may be confusing to you. **As long as you understand either Def. 1.24 or one of the following equivalent ways to express it, you're good to go.**

- A function $f : A \rightarrow B$ is **surjective or onto** if for each $b \in B$, the set $\{a \in A \mid f(a) = b\}$ (the set of all points in A that get mapped onto b) is nonempty. There could be multiple points $a \in A$ that get mapped onto the same $b \in B$, and the function would still be surjective, but there cannot exist a point $b \in B$ for which there are no points in A that get mapped onto it.
- **Said yet another way, the function's codomain is also its range.**

Example 1.25. Determine whether the following are surjective functions.

- (a) $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = \cos(x)$.
- (b) $f : \mathbb{R} \rightarrow [-1, 1]$ by $f(x) = \cos(x)$.
- (c) $f : (-\pi, \pi] \rightarrow [-1, 1]$ by $f(x) = \cos(x)$.
- (d) $f : [-\pi/2, 0] \rightarrow [-1, 1]$ by $f(x) = \cos(x)$.
- (e) $f : (-\pi, \pi] \rightarrow [0, 1]$ by $f(x) = \cos(x)$.

Solution: We make use of the fact that you “know well” the cosine function. Otherwise, we’d have to draw triangles to define the cosine function and find values for it!

- (a) $\cos : \mathbb{R} \rightarrow \mathbb{R}$ is not surjective (not onto) because there is no value $x \in \mathbb{R}$ such that $\cos(x) = 2$.
- (b) $\cos : \mathbb{R} \rightarrow [-1, 1]$ is onto (surjective) because $\cos(-\pi/2) = -1$, $\cos(0) = 1$, and for all $-1 < y < 1$, there is an $-\pi/2 < x < 0$ such that $\cos(x) = y$.
- (c) $\cos : (-\pi, \pi] \rightarrow [-1, 1]$ is onto (surjective) by the same reasoning as above.
- (d) $\cos : [-\pi/2, 0] \rightarrow [-1, 1]$ is onto (surjective) by the same reasoning as item (b)
- (e) $\cos : (-\pi, \pi] \rightarrow [0, 1]$ is not a function and hence it cannot be surjective (onto). Why is it not a function? Because $-1 = \cos(-\pi/2) \notin [0, 1]$. In other words, $[0, 1]$ is not a valid codomain for cosine.

Just as we gave alternative ways to say a function is surjective (onto), we do the same for injective (1:1).

- A function $f : A \rightarrow B$ is **injective or 1:1** if for each $b \in B$, the set $\{a \in A \mid f(a) = b\}$ (the set of all points in A that get mapped onto b) is either empty or a singleton. There could be no points $a \in A$ that get mapped onto a given $b \in B$, but if a point gets mapped onto b , there cannot exist a second point distinct from the first point.
- **Said another way, for an injective function, if there exists $a_1, a_2 \in A$ such that $f(a_1) = b$ and $f(a_2) = b$, then $a_1 = a_2$.**

You are free to ignore these additional ways to think about injectivity.

Example 1.26. Check if the following are injective functions or not.

- (a) $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = \cos(x)$.
- (b) $f : [-\pi/2, \pi/2] \rightarrow [-1, 1]$ by $f(x) = \cos(x)$.
- (c) $f : (-\pi, 0] \rightarrow [-1, 1]$ by $f(x) = \cos(x)$.
- (d) $f : [-\pi/2, 0] \rightarrow [-1, 1]$ by $f(x) = \cos(x)$.
- (e) $f : (-\pi, \pi] \rightarrow [0, 1]$ by $f(x) = \cos(x)$.

Solution: We once again assume that you “know well” the cosine function. We do not apply the horizontal line test, but you could in your solution. Just go to Julia and do the plots.

- (a) $\cos : \mathbb{R} \rightarrow \mathbb{R}$ is not injective (not 1:1) because $\cos(0) = 1$ and $\cos(2\pi) = 1$. To conclude **not injective**, all we need do is find a single pair of distinct points that get mapped to a common point. The fact that there are many other pairs of points that also get mapped to a common point is irrelevant to showing the definition of injectivity fails (aka, is not satisfied).
- (b) $\cos : [-\pi/2, \pi/2] \rightarrow [-1, 1]$ is not injective (not 1:1) because $\cos(-\pi/2) = 0$ and $\cos(\pi/2) = 0$.
- (c) $\cos : (-\pi, 0] \rightarrow [-1, 1]$ is 1:1 (injective) because for all $x, y \in (-\pi, 0]$, $x < y \implies \cos(x) < \cos(y)$ and hence distinct values in the domain produce distinct values in the range.
- (d) $\cos : [-\pi/2, 0] \rightarrow [-1, 1]$ is injective (1:1) because we already showed that $\cos : (-\pi, 0] \rightarrow [-1, 1]$ is injective and $[-\pi/2, 0] \subset (-\pi, 0]$. Alternatively, you can use the same reasoning as in (c).
- (e) $\cos : (-\pi, \pi] \rightarrow [0, 1]$ is not a function and hence it cannot be injective (1:1). Why is it not a function? Because $-1 = \cos(-\pi/2) \notin [0, 1]$.

■

Suggestion: Ask an LLM to create similar problems for $f(x) = x^2$.

Proposition 1.27: (Existence of Inverse Functions)

A function $f : A \rightarrow B$ has an inverse if, and only if it, is both injective (1:1) and surjective (onto).

Note: Such functions have a name: they are called **bijective**.

1.5.5 Strictly Monotonic Functions and Relation to Existence of Inverse Functions

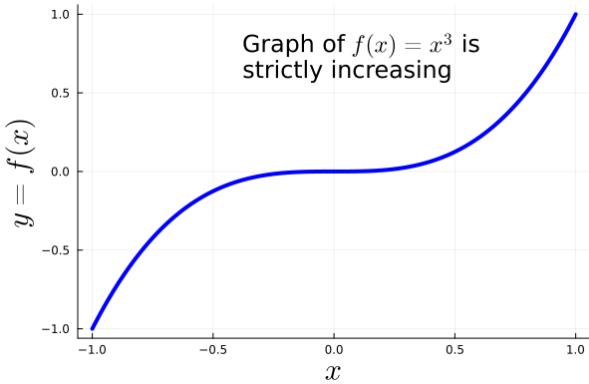
While it is great to have necessary and sufficient conditions for the existence of an inverse, in most cases, checking the conditions for the existence of an inverse is “nearly impossible”. There are a few easy cases; in that light, we introduce the following notions.

Strictly Monotonic Functions

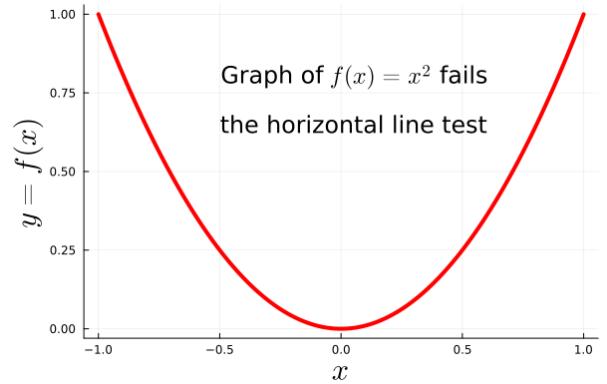
Definition 1.28. Let A and B be subsets of the real numbers and $f : A \rightarrow B$ a function.

- f is **strictly increasing** if for all $x, y \in A$, $y > x \implies f(y) > f(x)$.
- f is **strictly decreasing** if for all $x, y \in A$, $y > x \implies f(y) < f(x)$.

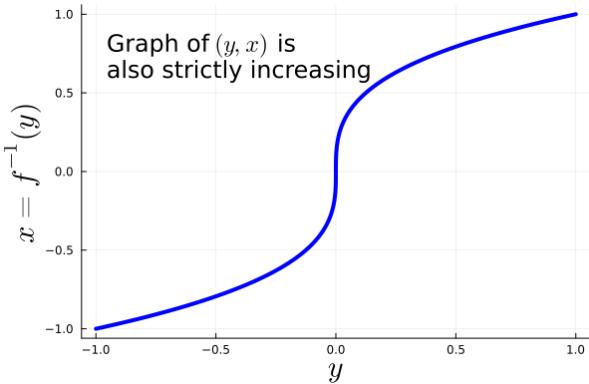
When the inequalities used above are **not strict**, we say that f is **increasing** or **decreasing**, respectively. Increasing and decreasing functions are also called **monotonically increasing** or **decreasing** to further emphasize that they never reverse direction; they can, however, remain “flat or constant” over a set of points.



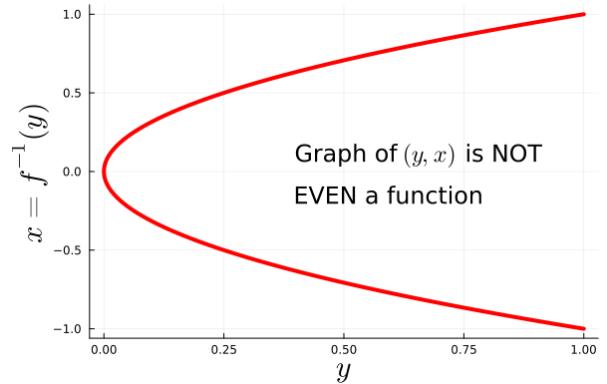
(a)



(b)



(c)



(d)

Figure 1.4: Comparison of “strictly monotonic” vs. “non-monotonic” functions and how that affects the existence of an inverse. In (a), the graph of $f : [-1, 1] \rightarrow \mathbb{R}$ by $f(x) = x^3$, and right below it in (c) is plotted (y, x) , that is, x and y have been swapped. Because $f : [-1, 1] \rightarrow \mathbb{R}$ by $f(x) = x^3$ is strictly increasing, it passes the horizontal line test and swapping the roles of x and y in the graph automatically produces the inverse of f . In (b), the graph of $f : [-1, 1] \rightarrow \mathbb{R}$ by $f(x) = x^2$, and right below it in (d) is plotted (y, x) (again, x and y have been swapped). When $f : [-1, 1] \rightarrow \mathbb{R}$ by $f(x) = x^2$ fails the horizontal line test, swapping the roles of x and y in the graph does NOT even produce a function, much less an inverse for f . The graph in (d) fails the vertical line test.

Remark 1.29. Among the powerful results of Calculus is an effective means to check if a function is (strictly) increasing or decreasing! For now, it is suggested that you graph the function and apply the horizontal line test. At least in Julia, that is not too painful.

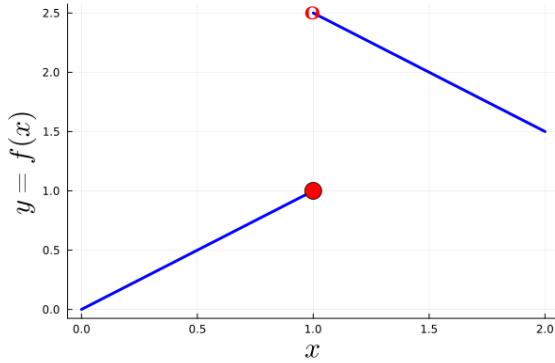
Proposition 1.30: Strict Monotonicity Helps with Invertibility

Suppose that A and B are subset of the real numbers.

- (a) If $f : A \rightarrow B$ is strictly monotonic (increasing or decreasing), then it is injective (one-to-one).
- (b) If $f : A \rightarrow B$ is strictly monotonic (increasing or decreasing) and surjective (onto), then it is invertible.

Remark 1.31. Suppose that A and B are subsets of the real numbers and $f : A \rightarrow B$ is a function.

- (a) When $f : A \rightarrow B$ is continuous as in Fig. 1.4, then the graph of f passes the horizontal line test if, and only if, $f : A \rightarrow B$ is EITHER strictly increasing or strictly decreasing.
- (b) When $f : A \rightarrow B$ is discontinuous, as in the plot below, then all bets are off!



For example,

$$y = f(x) := \begin{cases} x & 0 \leq x \leq 1 \\ 3.5 - x & 1 < x < 2 \end{cases}$$

has a jump at $x = 1$, and hence, is discontinuous, and yet it passes the horizontal line test, so it is injective (1:1). However, it is not monotonic. Indeed, if we look at the function on the first sub-interval, $[0, 1]$, it is strictly increasing, and on the second sub-interval, $(1, 2]$, it is strictly decreasing, and hence, overall, it is not monotonic.

- (c) When $A = [a, b]$ a **closed bounded interval** and $f : A \rightarrow B$ is **strictly monotonically increasing or decreasing**, then $f(A)$, the range of f is $f(A) = [c, d]$, where $c = \min\{f(a), f(b)\}$ and $d = \max\{f(a), f(b)\}$. In this case, $f : A \rightarrow B$ is surjective (onto) if, and only if, $B = [c, d]$. The result is false for many other choices for A . for example, it is false for $A = (a, b)$.

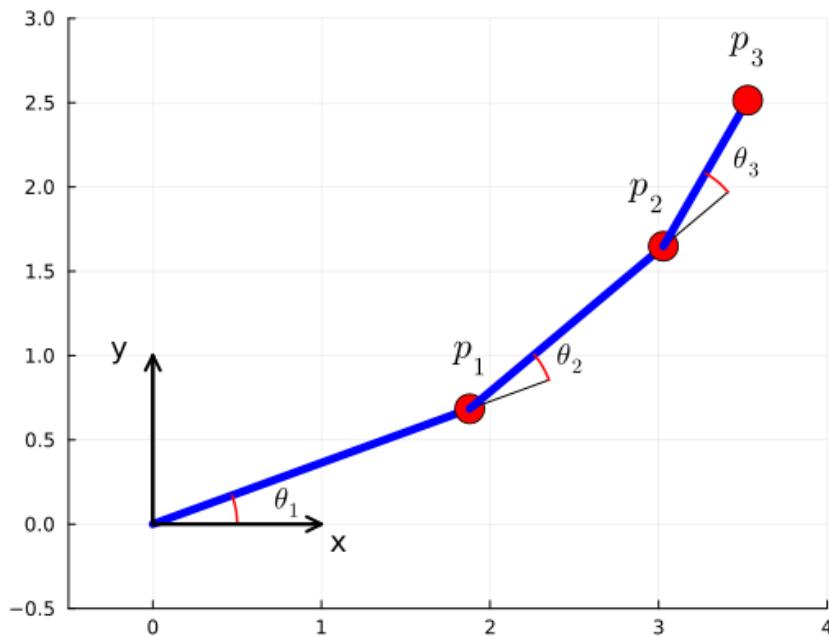


Figure 1.5: This is a robot kinematic chain, along with a Cartesian coordinate system, (x, y) . The angle θ_1 is measured counterclockwise (aka, CCW) from the x -axis, while θ_2 is measured CCW with respect to link-1 and θ_3 is measured CCW with respect to link-2. The text explains why the angles are defined this way: we can measure all of the indicated angles with practical instruments. The term “chain” just means one body is attached to another to form a chain. But what about “kinematic”?

1.6 Trigonometric and Inverse Trigonometric Functions

Figure 1.5 shows an idealized diagram of a robot with three links that we'll use to illustrate how Roboticists use trigonometry. In Robotics, we call such a diagram “a kinematic diagram of the robot” or a “kinematic chain representation of the robot”. Because the word “kinematic” is probably new to you, let's describe it.

Kinematics delves into the study of the motion of rigid bodies, focusing on the positions (not the speeds, which is dynamics) of points on a body as it moves through space, that is, \mathbb{R}^2 or \mathbb{R}^3 . In the context of robotics, such as the robot with three links in Fig 1.5, kinematics provides insights into the possible movements and configurations of the robot without considering the forces and physical properties influencing these movements (no masses, no inertias, no motor torques, which we'll study later). By emphasizing the geometric and algebraic aspects of motion, kinematics offers a mathematical framework to describe and analyze the geometry of motion, allowing Roboticists to achieve precise and coordinated movements in robotic systems.

1.6.1 Trig Functions

The main trigonometric functions we need are $\sin(x)$, $\cos(x)$, and $\tan(x)$. The others,

- $\sec(x) := \frac{1}{\cos(x)}$, called the **secant** of x ,
- $\csc(x) := \frac{1}{\sin(x)}$, called the **cosecant** of x , and
- $\cot(x) := \frac{\cos(x)}{\sin(x)} = \frac{1}{\tan(x)}$, called the **cotangent** of x ,

are reciprocals of the main trigonometric functions. They show up when we do trigonometric substitutions in Chapter 7.3.6. It is strongly suggested you watch this video, [Mathematical Visual Proofs: Trig Visualized](#) as review. An informative comment left by a viewer: “I've never seen this diagram. Even a static version is quite informative, but the animation knocks it out of the park.”

We do need the inverse $\sin(x)$, $\cos(x)$, and $\tan(x)$ functions. Those are treated separately because they can be “tricky”.

1.6.2 Trig Identities

There are “innumerable” trig identities. So, which ones to learn and remember? The ones we list here show up in the book in at least one location.

- $\cos^2(x) + \sin^2(x) = 1$.
- $\sec^2(x) = 1 + \tan^2(x)$.
- $\sin(2x) = 2 \sin(x) \cdot \cos(x)$.
- $\cos(2x) = 2 \cos^2(x) - 1$.
- $\cos^2(x) = \frac{1}{2} + \frac{1}{2} \cos(2x)$.
- $\sin^2(x) = \frac{1}{2} - \frac{1}{2} \cos(2x)$.
- $\cos(2x) = \cos^2(x) - \sin^2(x)$.

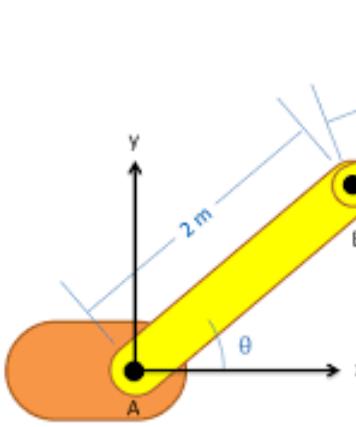
Even though several of the above identities are “redundant” in the sense that they can be easily derived from other listed identities, it's more efficient to have them handy in the “heat of battle”.

1.6.3 A 3-Link Robot Kinematic Chain in the Plane

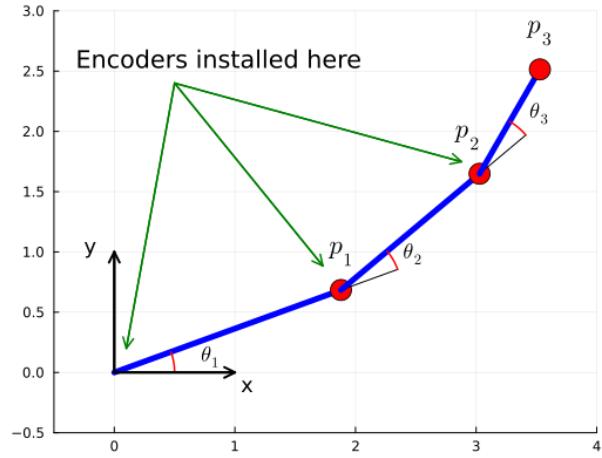
Let's define the lengths of the three links to be L_1 , L_2 , and L_3 . In Fig. 1.5, the lengths are 2, 1.5, and 1, respectively. Let's define p_0 to be the origin,

$$p_0 = \begin{bmatrix} 0.0 \\ 0.0 \end{bmatrix},$$

and we assume that the base of the first link is attached at the origin with a **revolute joint**. This means that its position is fixed at p_0 , but the link can rotate in a circle, that is 2π radians, or 360° . The angle in a revolute joint is measured with a **joint encoder**, which is



(a)



(b)

Figure 1.6: Absolute angles are those referenced to an agreed-upon stationary coordinate system (called a world frame in Mechanics and Robotics). In (a), both θ and ϕ are absolute angles; each is measured with respect to a fixed horizontal axis. **Relative angles** are those referenced to a coordinate system that moves with respect to the agreed-upon world frame. Typically, the angle is measured *relative* to another point on the robot's body. In (b), θ_1 is the absolute angle of the first link, while θ_2 measures the relative angle of the second link with respect to the first. The absolute angle of the second link is $\theta_1 + \theta_2$, the sum of the absolute angle of the first link and the relative angle between the second link and the first link. Similarly, the absolute angle of the third link is $\theta_1 + \theta_2 + \theta_3$. Source for (a): [Mechanics Map Open Textbook Project](#)

typically shortened to simply **encoder**. You have to specify a zero point (origin) for the encoder and the positive direction. We will say the positive x -axis is the zero point (0°), and counterclockwise (CCW) as the positive direction. With this convention, the y -axis is at $\frac{\pi}{2}$ or 90° .

With this data in our pocket, we can compute the Cartesian position of the end of link-1, labeled p_1 in Fig. 1.5. We obtain,

$$p_1 = p_0 + \begin{bmatrix} L_1 \cdot \cos(\theta_1) \\ L_1 \cdot \sin(\theta_1) \end{bmatrix}. \quad (1.15)$$

Because p_0 is the origin, we could leave it out of the equation, but including it (i) costs us nothing, (ii) it allows us to mount the base of the robot at any point we wish, **AND**, as you will see, (iii) it gives us a recursive structure for writing down the equations. That's a lot of benefits.

Continuing, we have

$$p_2 = p_1 + \begin{bmatrix} L_2 \cdot \cos(\theta_1 + \theta_2) \\ L_2 \cdot \sin(\theta_1 + \theta_2) \end{bmatrix}.$$

This looks just like (1.15), except we have $\theta_1 + \theta_2$. Why the sum of the angles? As illustrated in Fig. 1.6-(b), the angle θ_2 is measured with an encoder installed at the connecting point of links 1 and 2 of the robot. The encoder measures the angle of rotation of link 2 with respect to link 1, or in the vocabulary of Fig. 1.6, the relative angle between the second and first link. We need the absolute angle⁷ of the second link to calculate the position of its tip in the Cartesian coordinate system (x, y) . That angle is the sum of the absolute angle of the first link and the angle of the second link relative to the first link, in symbols, $\theta_1 + \theta_2$.

Based on the above discussion, computing the Cartesian position of the end of the third link is a breeze,

$$p_3 = p_2 + \begin{bmatrix} L_3 \cdot \cos(\theta_1 + \theta_2 + \theta_3) \\ L_3 \cdot \sin(\theta_1 + \theta_2 + \theta_3) \end{bmatrix}. \quad (1.16)$$

But wait a minute, do we now have to do all of the substitutions by hand to find p_2 and p_3 ? No, we have computers to do the work for us⁸. While in this case, we could easily do the math by hand, for more realistic robots, that would not be any fun. We illustrate this next. First, a function to compute p_1 , p_2 , and p_3 .

⁷While it is “theoretically possible” to measure the absolute angle of the second link directly, by using a laser, for example, the only practical way to do it is via an encoder attached at the joints between successive links. In any case, that's how we do it in Robotics!

⁸See (1.17) for the general formula for p_3 , the hardest one. You can do p_2 .

```

1 # functions for
2
3 function modelParameters()
4     # Place the parameters in a single location so that
5     # they can be updated at a moment's notice
6     L1, L2, L3 = [2, 1.5, 1]
7     return (L1=L1, L2=L2, L3=L3)
8 end
9
10 function linkPositions (th1, th2, th3)
11     params = modelParameters()
12     p0=[0.0; 0.0]
13     p1 = p0 + [ params.L1 * cos(th1); params.L1 * sin(th1) ]
14     p2 = p1 + [ params.L2 * cos(th1+th2); params.L2 * sin(th1+th2) ]
15     p3 = p2 + [ params.L3 * cos(th1+th2+th3); params.L3 * sin(th1+th2+th3) ]
16     return (p0=p0, p1=p1, p2=p2, p3=p3)
17 end
18
19 th1, th2, th3 = (2*pi/360) * [20, 20, 20]
20 positions = linkPositions (th1, th2, th3)

```

Output

```
(p0 = [0.0, 0.0], p1 = [1.8793852415718169, 0.6840402866513374],
p2 = [3.028451906250284, 1.6482217011811464],
p3 = [3.528451906250284, 2.514247104965585])
```

Next, we exercise the above function by computing p_1 , p_2 , and p_3 for various angles θ_1 , θ_2 , and θ_3 , and then plotting the links of the robot. Once we write the code, we can run it as many times as we want, without pain. We can also put in nominal angles, such as all angles equal to zero or all angles equal to $\frac{\pi}{2}$, and check that we obtain what we know has to be the correct configuration of the mechanism. We did that, of course, but do not show the plots here.

```

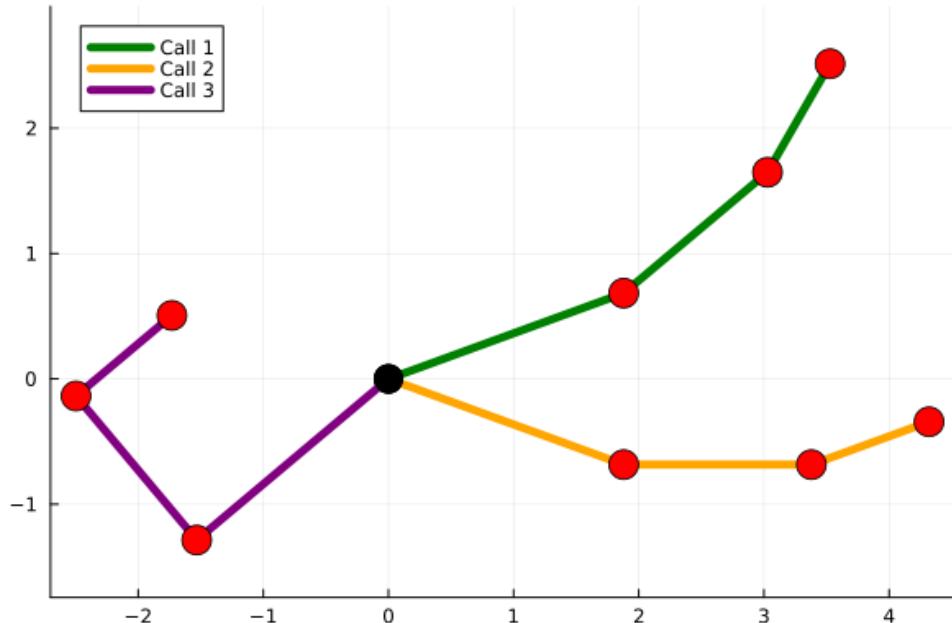
1 using Plots
2
3 # Global variable to keep track of the number of times the function has been called
4 global call_count = 0
5
6 function plot_points(positions; line_thickness=5, ball_size=10)
7     global call_count
8     p0 = positions.p0
9     p1 = positions.p1
10    p2 = positions.p2
11    p3 = positions.p3
12
13    # Colors for the lines for each call
14    line_colors = [:blue, :green, :orange, :purple, :cyan, :magenta, :yellow, :black]
15
16    # Increment the call count and determine the color for this call
17    call_count += 1
18    current_color = line_colors[mod(call_count, length(line_colors)) + 1]
19
20    # Plot the lines
21    plot!([p0[1], p1[1], p2[1], p3[1]], [p0[2], p1[2], p2[2], p3[2]],
22          linewidth=line_thickness, color=current_color, label="Call $(call_count)")
23
```

```

24     # Plot the balls
25     scatter! ([p0[1]], [p0[2]], color=:black, markersize=ball_size, label=nothing)
26     scatter! ([p1[1], p2[1], p3[1]], [p1[2], p2[2], p3[2]], color=:red,
27               markersize=ball_size, label=nothing)
28
29     display(current())
30 end
31
32 # Initialize the plot
33 #plot(; legend=:topright, aspect_ratio=:equal)
34 fig = plot(; legend=:topleft, aspect_ratio=:equal)
35
36 # Example usage
37 th1, th2, th3 = (2*pi/360) * [20, 20, 20]
38 positions = linkPositions(th1, th2, th3)
39 plot_points(positions)
40
41 # Call again with different angles to see the effect
42 th1, th2, th3 = (2*pi/360) * [-20, 20, 20]
43 positions = linkPositions(th1, th2, th3)
44 plot_points(positions)
45
46 # Call again with different angles to see the effect
47 th1, th2, th3 = (2*pi/360) * [-140, -90, -90]
48 positions = linkPositions(th1, th2, th3)
49 plot_points(positions)
50
51 # Make a png for inclusion in the textbook
52 png(fig, "robotKinematicChainWithOurNewFunction")

```

Output This cool image:



As promised, we do give the general formula for the position of the end effector, that is, the end of link 3. From this, you can deduce

p_2 very easily! We already gave p_1 in (1.15).

$$\begin{bmatrix} p_{3,x} \\ p_{3,y} \end{bmatrix} = \begin{bmatrix} L1 \cdot \cos(\theta_1) + L2 \cdot \cos(\theta_1 + \theta_2) + L3 \cdot \cos(\theta_1 + \theta_2 + \theta_3) \\ L1 \cdot \sin(\theta_1) + L2 \cdot \sin(\theta_1 + \theta_2) + L3 \cdot \sin(\theta_1 + \theta_2 + \theta_3) \end{bmatrix} \quad (1.17)$$

1.6.4 Notation for Inverse Trigonometric Functions

The trigonometric functions, sine, cosine, and tangent, have **inverse functions** if their domains and codomains are properly defined. The inverse functions help us find **angles in radians** when we know the trigonometric values. There are a few ways to denote the inverse functions, and we go through them below.

- The most popular way to write them is with an “arc” in front. For example:
 - $\arcsin(x)$ for inverse sine.
 - $\arccos(x)$ for inverse cosine.
 - $\arctan(x)$ for inverse tangent.
- In Julia and many other programming languages, they are written a bit shorter, namely
 - asin for inverse sine.
 - acos for inverse cosine.
 - atan for inverse tangent.
- Another way, which was introduced in the 1800s by John Herschel, is to write them with a “-1” as a power, like this,
 - $\sin^{-1}(x)$ for inverse sine.
 - $\cos^{-1}(x)$ for inverse cosine.
 - $\tan^{-1}(x)$ for inverse tangent.

It is recommended that you avoid this notation because it can be confusing. Really? Yes! If we write $\sin^2(x)$ for $(\sin(x))^2$, then how would we write $\frac{1}{\sin(x)}$? Oh yeah, we might write $\sin^{-1}(x)$, confusing the reciprocal and the inverse, which can lead to a mess. You will see it used in some courses, which is why we include this “awkward” notation here, even if we will avoid using it.

Remark 1.32. *In Julia, the trigonometric functions `sin`, `cos`, and `tan`, as well as their inverses, operate in radians by default. However, if you want to use degrees, Julia provides specific functions for that. `sind(x)`, `cosd(x)`, and `tand(x)` compute the sine, cosine, and tangent of x , where x is in degrees. For the inverse trigonometric functions in degrees, the functions `asind(x)`, `acosd(x)`, and `atand(x)` compute the arcsine, arccosine, and arctangent of x , returning the result in degrees.* ■

1.6.5 Understanding Inverse Trigonometric Functions

What messes everyone up when dealing with inverse trig functions, even the pros, is remembering the domains of the functions! In Figures 1.7-(a), 1.8-(a), and 1.9-(a), we show the functions as you would normally plot them. In each case, the functions fail the horizontal line test for being injective (1:1), and hence they do not have inverses. If you properly restrict their domains, as shown in the middle plots, then they do become injective (1:1) and have an inverse, as shown in the plot to the right, that is, Figures 1.7-(c), 1.8-(c), and 1.9-(c).

So what is so hard about that? Well, consider the sine function. It is injective (1:1) and surjective (onto) for all of the following domains

- $\sin : [-\frac{\pi}{2}, \frac{\pi}{2}] \rightarrow [-1, 1]$.
- $\sin : [\frac{\pi}{2}, \frac{3\pi}{2}] \rightarrow [-1, 1]$.
- $\sin : [-\frac{3\pi}{2}, -\frac{\pi}{2}] \rightarrow [-1, 1]$.
- $\sin : [-\frac{\pi}{2} + k\pi, -\frac{\pi}{2} + k\pi] \rightarrow [-1, 1], k \in \mathbb{Z}$.

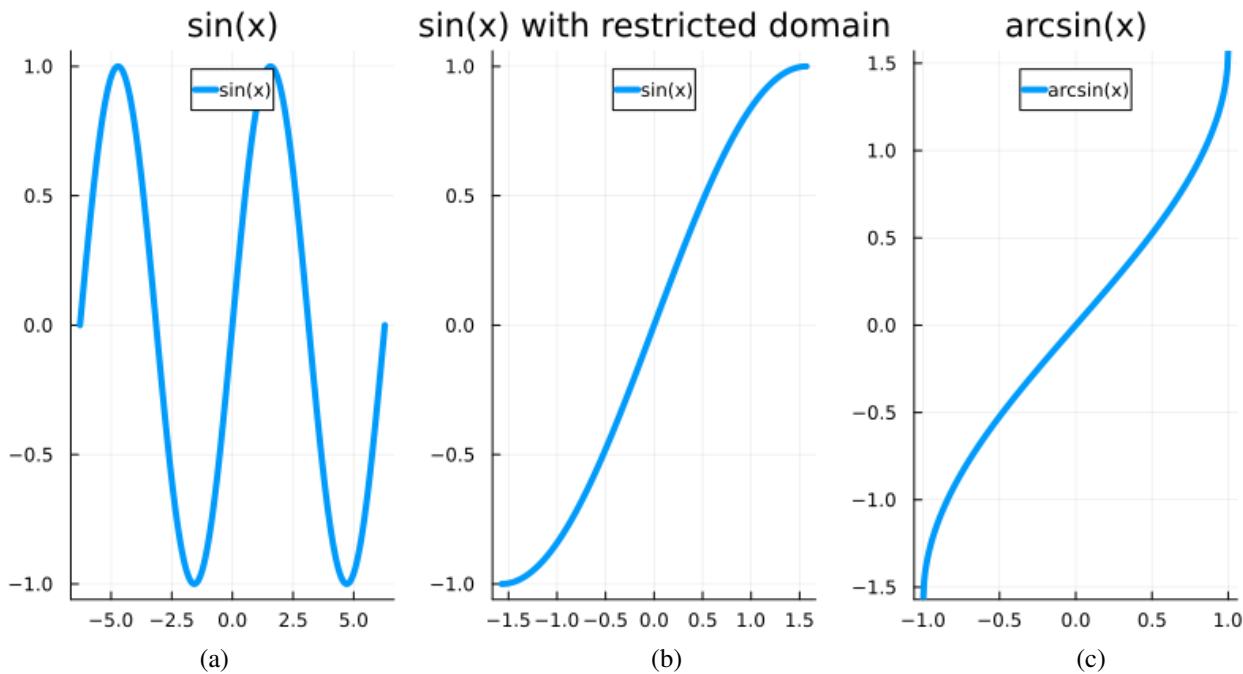


Figure 1.7: The sine function and its inverse. The sine function $\sin : \mathbb{R} \rightarrow [-1, 1]$ in (a) is surjective (onto) $[-1, 1]$ but does not have an inverse over its entire domain. You can see that it fails the “horizontal line test”. However, when restricted to $[-\frac{\pi}{2}, \frac{\pi}{2}]$, as in (b), it becomes injective (1:1) and has an inverse, shown on the right in (c).

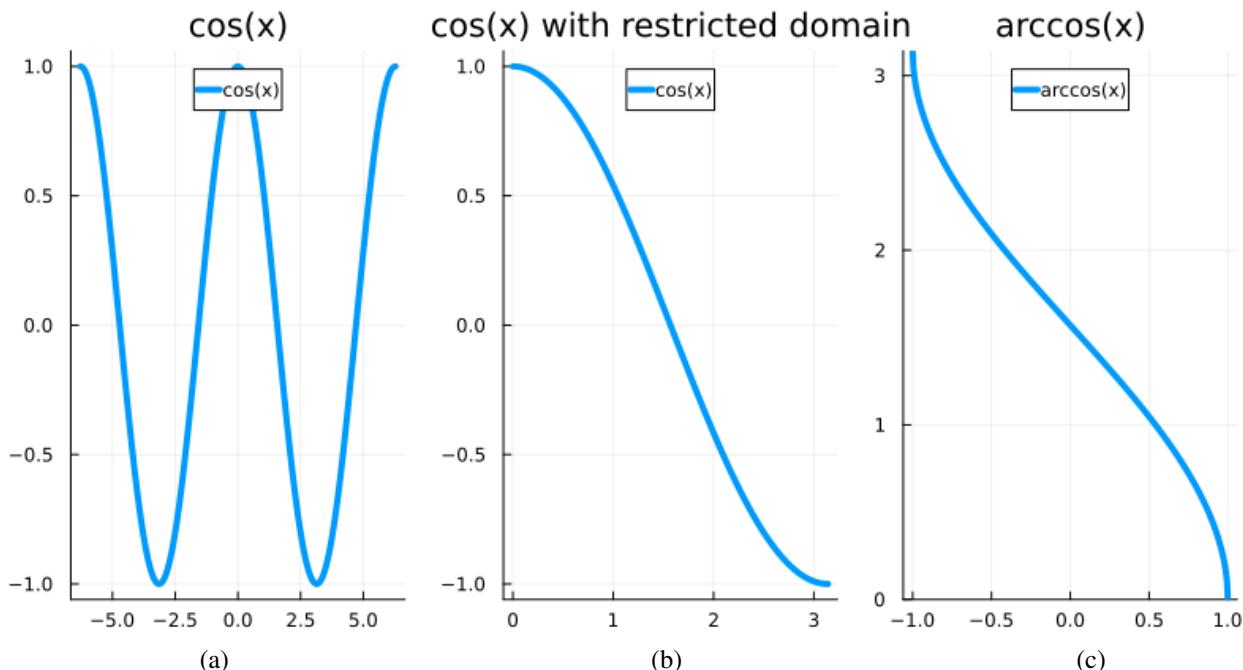


Figure 1.8: The cosine function and its inverse. The cosine function $\cos : \mathbb{R} \rightarrow [-1, 1]$ in (a) is surjective (onto) $[-1, 1]$ but does not have an inverse over its entire domain. You can see that it fails the “horizontal line test”. However, when restricted to $[0, \pi]$, as in (b), it becomes injective (1:1) and has an inverse, shown on the right in (c).

For all of these domains, $\sin(x)$ satisfies the “horizontal line test”, and hence, an inverse can be defined in all of these situations. **Therefore, you have to make a choice of domain, and making that choice is where most people do a face plant** (Steven Universe Faceplants - “I Knew You Were Trouble” to really get your attention).

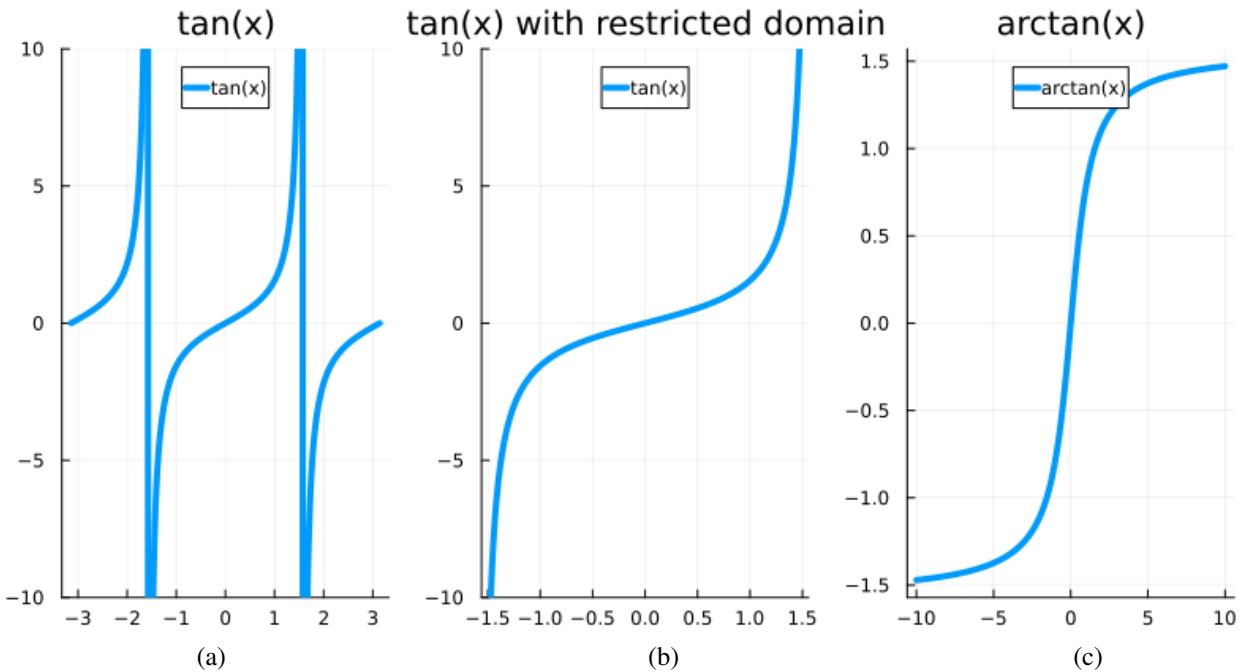


Figure 1.9: The tangent function and its inverse. The tangent function as shown in (a) does not have an inverse over \mathbb{R} , even after setting $\tan(x)$ to zero at odd multiples of $\frac{\pi}{2}$, where it blows up. However, when restricted to the open interval $(-\frac{\pi}{2}, \frac{\pi}{2})$, as shown in (b), it becomes injective (1:1) and has an inverse, shown on the right in (c).

The domain choices made in Table 1.33 and Figures 1.7, 1.8, and 1.9 are called the **Principal Domains**. Paraphrasing⁹, the principal¹⁰ domain of a trigonometric function is the largest interval near the origin (positive or symmetric about zero, if possible), where the function is one-to-one, strictly increasing or decreasing, and takes on all possible values. It is also the domain principally used in STEM fields. Ha ha!

We summarize this information in Table 1.33 and suggest the following videos if you are still having problems:

- Evaluating Inverse Trigonometric Functions (arcsin, arccos, arctan) Using Unit Circle by Mario's Math Tutoring.
- Inverse Trigonometric Functions by Professor Dave Explains.
- Inverse trig functions: arcsin | Trigonometry by the Khan Academy.

	$\sin(x)$	$\cos(x)$	$\tan(x)$	$\arcsin(x)$	$\arccos(x)$	$\arctan(x)$
Domain	$[-\frac{\pi}{2}, \frac{\pi}{2}]$	$[0, \pi]$	$(-\frac{\pi}{2}, \frac{\pi}{2})$	$[-1, 1]$	$[-1, 1]$	$(-\infty, \infty)$
Codomain	\mathbb{R}	\mathbb{R}	\mathbb{R}	\mathbb{R}	\mathbb{R}	\mathbb{R}
Range	$[-1, 1]$	$[-1, 1]$	$(-\infty, \infty)$	$[-\frac{\pi}{2}, \frac{\pi}{2}]$	$[0, \pi]$	$(-\frac{\pi}{2}, \frac{\pi}{2})$

Table 1.33: Domains, Codomains, and Ranges of Trigonometric and Inverse Trigonometric Functions

We illustrate the importance of principal domains with the following code. We evaluate \sin , \cos , and \tan for x in their principal domains and outside them. Then we apply the inverse functions and see if we recover the original x values or not!

⁹Stewart, James. Precalculus: Mathematics for Calculus. 7th ed. Cengage Learning, 2016. You can find a link to the book on the [Cengage Learning website](#).

¹⁰“Principle” means an idea or concept, whereas as “principal” means main or primary. More formally, A “principle” is a rule, a law, a guideline, or a fact. A “principal” is the headmaster of a school or a person who’s in charge of certain things in a company. “Principal” is also an adjective that means original, first, or most important.

```

1 # A few examples illustrating principal domains
2 # we take values of sin, cos, and tan within their principal
3 # domains and outside them. Then we apply the inverse functions
4 # and see if we recover the original values or not!
5
6 using Printf
7
8 x = pi*[-1.3, -1, -0.4, -0.25, 0, 0.25, 0.4, 1, 1.3]
9 function_names = ["sin(x)", "cos(x)", "tan(x)"]
10
11 for k = 1:length(x)
12     xk = x[k]
13     InData = [sin(xk) cos(xk) tan(xk)]
14     OutData = [asin(sin(xk)) acos(cos(xk)) atan(tan(xk))]
15     OriginalXk = [xk xk xk]
16     FlagDomainProblem = (abs.(OutData - OriginalXk) .> 1e-8) # Using a small
17                                     #threshold to account for floating-point errors
18
19     println("For x = $(@sprintf("%.10f", xk)):")
20     println("Trig function values: ", InData)
21     println("After applying inverse: ", OutData)
22     println("Hoping for: ", OriginalXk)
23     for j = 1:length(FlagDomainProblem)
24         if FlagDomainProblem[j]
25             println("x = $xk is outside the principal domain for $(function_names[j]) .")
26         else
27             println("x = $xk is inside the principal domain for $(function_names[j]) .")
28         end
29     end
30
31     println("-----")
32 end

```

Output (after truncating numbers to 4 decimal places)

```

For x = -4.0841:
Trig function values:  0.8090, -0.5878, -1.3764
After applying inverse: 0.9425, 2.1991, -0.9425
Hoping for:           -4.0841, -4.0841, -4.0841
x = -4.0841 is outside the principal domain for sin(x).
x = -4.0841 is outside the principal domain for cos(x).
x = -4.0841 is outside the principal domain for tan(x).
-----
For x = -3.1416:
Trig function values:  -0.0000, -1.0000, 0.0000
After applying inverse: -0.0000, 3.1416, 0.0000
Hoping for:           -3.1416, -3.1416, -3.1416
x = -3.1416 is outside the principal domain for sin(x).
x = -3.1416 is outside the principal domain for cos(x).
x = -3.1416 is outside the principal domain for tan(x).
-----
For x = -1.2566:
Trig function values:  -0.9511, 0.3090, -3.0777
After applying inverse: -1.2566, 1.2566, -1.2566
Hoping for:           -1.2566, -1.2566, -1.2566
x = -1.2566 is inside the principal domain for sin(x).

```

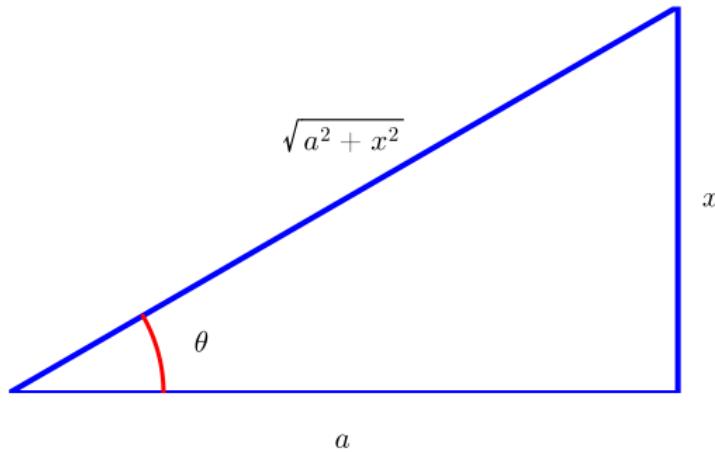
```

x = -1.2566 is outside the principal domain for cos(x) .
x = -1.2566 is inside the principal domain for tan(x) .
-----
For x = -0.7854:
Trig function values: -0.7071, 0.7071, -1.0000
After applying inverse: -0.7854, 0.7854, -0.7854
Hoping for: -0.7854, -0.7854, -0.7854
x = -0.7854 is inside the principal domain for sin(x) .
x = -0.7854 is outside the principal domain for cos(x) .
x = -0.7854 is inside the principal domain for tan(x) .
-----
For x = 0.0000:
Trig function values: 0.0000, 1.0000, 0.0000
After applying inverse: 0.0000, 0.0000, 0.0000
Hoping for: 0.0000, 0.0000, 0.0000
x = 0.0000 is inside the principal domain for sin(x) .
x = 0.0000 is inside the principal domain for cos(x) .
x = 0.0000 is inside the principal domain for tan(x) .
-----
For x = 0.7854:
Trig function values: 0.7071, 0.7071, 1.0000
After applying inverse: 0.7854, 0.7854, 0.7854
Hoping for: 0.7854, 0.7854, 0.7854
x = 0.7854 is inside the principal domain for sin(x) .
x = 0.7854 is inside the principal domain for cos(x) .
x = 0.7854 is inside the principal domain for tan(x) .
-----
For x = 1.2566:
Trig function values: 0.9511, 0.3090, 3.0777
After applying inverse: 1.2566, 1.2566, 1.2566
Hoping for: 1.2566, 1.2566, 1.2566
x = 1.2566 is inside the principal domain for sin(x) .
x = 1.2566 is inside the principal domain for cos(x) .
x = 1.2566 is inside the principal domain for tan(x) .
-----
For x = 3.1416:
Trig function values: 0.0000, -1.0000, -0.0000
After applying inverse: 0.0000, 3.1416, -0.0000
Hoping for: 3.1416, 3.1416, 3.1416
x = 3.1416 is outside the principal domain for sin(x) .
x = 3.1416 is inside the principal domain for cos(x) .
x = 3.1416 is outside the principal domain for tan(x) .
-----
For x = 4.0841:
Trig function values: -0.8090, -0.5878, 1.3764
After applying inverse: -0.9425, 2.1991, 0.9425
Hoping for: 4.0841, 4.0841, 4.0841
x = 4.0841 is outside the principal domain for sin(x) .
x = 4.0841 is outside the principal domain for cos(x) .
x = 4.0841 is outside the principal domain for tan(x) .
-----
```

1.6.6 Composing Trig Functions and Inverse Trig Functions

The video [Pre-Calculus Algebraic Expression Inverse Trig and Double Angle](#) by Malmos Math shows how to use right triangles and basic trig identities to derive algebraic formulas for expressions such as

$$\cos\left(2 \arctan\left(\frac{x}{a}\right)\right). \quad (1.18)$$



Let $\theta := \arctan\left(\frac{x}{a}\right) \iff \tan(\theta) = \frac{x}{a}$. Because tangent equals opposite over adjacent, θ is as represented in the above right triangle. This is the key step! From our handy list of trig identities, we have $\cos(2\theta) = 2\cos^2(\theta) - 1$. From our right triangle, because cosine equals adjacent over hypotenuse, we have that $\cos(\theta) = \frac{a}{\sqrt{a^2+x^2}}$. Now, we just have some algebra to do,

$$\begin{aligned} \cos\left(2 \arctan\left(\frac{x}{a}\right)\right) &= \cos(2\theta) \\ &= 2\cos^2(\theta) - 1 \\ &= 2\left(\frac{a}{\sqrt{a^2+x^2}}\right)^2 - 1 \\ &= \frac{2a^2}{a^2+x^2} - \frac{a^2+x^2}{a^2+x^2} \\ &= \frac{a^2-x^2}{a^2+x^2}. \end{aligned}$$

1.7 Powers and Roots: Common, but Often Less Understood

You have likely used the function x^y , where y is not an integer. For example, $x^{3/2}$. What does it really mean? We want it to mean $\sqrt{x^3}$, but then we have to wrestle with the fact that, as long as we work within the real numbers, the expression only makes sense for $x \geq 0$. Moreover, when we do calculations on a computer, we want it to return a single value, even though we know that square roots can be negative. So, let's define once and for all what we mean by whole-number roots, $\sqrt[n]{x}$, $n \geq 1$, which are often denoted as $x^{1/n}$, and then rational powers, $x^{p/q}$, for integers p and q . Finally, we'll broach the subject of x^y for a real number y .

1.7.1 n -th Roots, aka Integer Roots

Defining Roots in a Careful Manner and the Notation $x^{1/n}$

Definition 1.34. n -th Roots

- (a) For $n \geq 1$ an **odd integer** and x any real number, we define $\sqrt[n]{x}$ to be the unique real number y such that $y^n = x$. We note that x and y will have the same sign.
- (b) For $n \geq 1$ an **even integer** and x a non-negative real number, we define $\sqrt[n]{x}$ to be the unique non-negative real number y such that $y^n = x$.
- (c) For all **positive integers** n , we define $x^{1/n}$ to be the same as $\sqrt[n]{x}$. Hence, if n is odd, $x^{1/n}$ is defined for all real x , and if n is even, it is only defined for non-negative x . **The zeroth root of x is not defined.**
- (d) For all **positive integers** n and $x \neq 0$, we define $x^{-1/n}$ to be the same as $\frac{1}{x^{1/n}}$. Hence, if n is odd, $x^{-1/n}$ is defined for non-zero x , and if n is even, it is only defined for positive x .

Even though it is not a root, it is worth recalling that, for x any non-zero real number, we define $x^0 := 1$. This explains why a zeroth root is not defined (the zeroth root of 1 could be anything). In Calculus, it is important to leave 0^0 as **undefined**. For more on zero to power zero, check out this [link](#). There are plenty of sources that can lead you astray, such as [Math is Fun](#). When we cite an online source, we have checked it out, for the most part!

Example 1.35. Use the Bisection Algorithm to compute $\sqrt[7]{x}$ for $x = 3$ and $x = -5$. Compute the answer to a tolerance of $1e-5$

Solution:

The Bisection Algorithm implements the literal definition of a root: $\sqrt[n]{x}$ is the unique real number y such that $y^n = x$. Here, it is solving $y^7 - x = 0$.

```

1 function bisection(f, a, b, tol)
2     if f(a) * f(b) > 0
3         error("f(a) and f(b) must have opposite signs")
4         return NaN
5     end
6     kmax = 1e5 # max number of loops so that the algorithm terminates
7     k=1
8     c=NaN #define c out of the loop so we can access it
9     while ( (b - a) / 2 > tol || abs(f(c)) > tol ) && k < kmax
10        c = (a + b) / 2
11        #
12        if f(a) * f(c) < 0
13            b = c
14        else
15            a = c
16        end
17        k = k+1
18    end
19    root_ErrorBound = (b - a) / 2
20    return (est=c, low=a, up=b, error=root_ErrorBound)
21 end

```

```

1 # Define the function whose root we want to find
2 x=3
3 f(y) = y^7 - x
4

```

```

5 # Call the bisection method on this function
6 F = bisection(f, 1.0, 2.0, 1e-5)
7 @show F
8
9 println(" ")
10 println("The 7-th root is approximately ", F.est, )
11 println("The 7-th power of the root is ", F.est^7)

```

Output

F = (est = 1.1699304580688477, low = 1.1699304580688477, up = 1.169931411743164,
error = 4.76837158203125e-7)

The 7-th root is approximately 1.1699304580688477
The 7-th power of the root is 2.9999936334014214

We call the function with a different set of bracketing points and $x = -5$.

```

1 # Define the function whose root we want to find
2 x=-5
3 f(y) = y^7 - x
4
5 # Call the bisection method on this function
6 F = bisection(f, -3.0, 0.0, 1e-5)
7 @show F
8
9 println(" ")
10 println("The 7-th root is approximately ", F.est, )
11 println("The 7-th power of the root is ", F.est^7)

```

Output

F = (est = -1.2584989070892334, low = -1.2584996223449707, up = -1.2584989070892334,
error = 3.5762786865234375e-7)

The 7-th root is approximately -1.2584989070892334
The 7-th power of the root is -4.999998788762909

Example 1.36. Compute the following roots for $x = 3$ and $x = -5$. When defined, give the answers to 3 decimal points. You can use a calculator or the Julia function for x^y .

- (a) $\sqrt[3]{x}$
- (b) $\sqrt[3]{x}$
- (c) $\sqrt[9]{x}$
- (d) $\sqrt[1]{x}$
- (e) $\sqrt[5]{x}$

Solution: We use a built-in Julia function to do the hard part of the calculations.

- (a) $\sqrt[3]{3} = 1.442$ and $\sqrt[3]{-5} = -1.710$. You have to use the command `cbrt(x)` when computing the cubed root of a negative number in Julia. Otherwise, Julia returns a domain error with $(-5)^{1/3}$; you should check this out. The reason is explained when we cover rational powers. A second method is shown below.
- (b) $\sqrt[2]{3} = 1.732$ while $\sqrt[2]{-5}$ is undefined (when working with the real numbers). We note that $\sqrt{x} := \sqrt[2]{x}$ because we have previously defined the notation $\sqrt[2]{x}$, whereas, when we previously used the symbol \sqrt{x} , we were counting on your recall of how it was defined.

- (c) $\sqrt[0]{3}$ and $\sqrt[0]{-5}$ are both undefined.
- (d) $\sqrt[1]{3} = 3$ and $\sqrt[1]{-5} = -5$.
- (e) $\sqrt[5]{3} = 1.246$ and $\sqrt[5]{-5} = -1.380$. In Julia, you can use the command `sign(x) * (abs(x)^(1/5))`, where

$$\text{sign}(x) = \begin{cases} +1 & x > 0 \\ -1 & x < 0 \\ 0 & x = 0 \end{cases}$$

By taking the absolute value, Julia has no trouble computing the root and you restore the correct sign with the sign function¹¹.

```

1 function myRoot (x : :Float64, n : :Int)
2   if n % 2 == 0
3     println(n, " is even .")
4     if x < 0
5       println(n, "Root is not defined for n even and x negative .")
6       root = NaN
7     end
8   else
9     println(n, " is odd .")
10   end
11   root = sign(x) * ( abs(x)^(1/n) )
12   return root
13 end
14
15 myRoot (-5.0, 5)
```

Output

-1.379729661461215

Proposition 1.37: (Root of a Product Equals the Product of Roots)

For all natural numbers $n \geq 1$ and real numbers $x > 0$ and $y > 0$,

$$\sqrt[n]{x \cdot y} = \sqrt[n]{x} \cdot \sqrt[n]{y}. \quad (1.19)$$

Proof: This proof, like many proofs in early math courses, is a “simple application” of the definitions. Let’s see if you agree! Let $z_1 := \sqrt[n]{x}$ and $z_2 := \sqrt[n]{y}$. Then, by the definition of an n -th root, z_1 and z_2 must satisfy

$$(z_1)^n = x \text{ and } (z_2)^n = y.$$

Hence,

$$x \cdot y = (z_1)^n \cdot (z_2)^n = (z_1 \cdot z_2)^n$$

by standard rules for multiplication. Indeed,

$$\underbrace{z_1 \cdot z_1 \cdots z_1}_{n \text{ times}} \cdot \underbrace{z_2 \cdot z_2 \cdots z_2}_{n \text{ times}} = \underbrace{(z_1 \cdot z_2) \cdots (z_1 \cdot z_2)}_{n \text{ times}} = (z_1 \cdot z_2)^n.$$

Because $(z_1 \cdot z_2)^n = x \cdot y$, we can conclude that $z_1 \cdot z_2 = \sqrt[n]{x \cdot y}$, that is, $\sqrt[n]{x} \cdot \sqrt[n]{y} = \sqrt[n]{x \cdot y}$, proving the result. ■

Remark 1.38. What if we have $\sqrt[n]{x^m}$ for n and m positive integers? You can successively apply Prop. 1.37 to obtain

$$\sqrt[n]{x^m} = \sqrt[n]{x} \cdot \sqrt[n]{x^{m-1}} = \sqrt[n]{x} \cdot \sqrt[n]{x} \cdot \sqrt[n]{x^{m-2}} = \underbrace{(\sqrt[n]{x}) \cdot (\sqrt[n]{x}) \cdots (\sqrt[n]{x})}_{m \text{ times}} = (\sqrt[n]{x})^m.$$

¹¹The sign function is sometimes called the signum function because “signum” is Latin for “sign”

For the first factorization of the n -th root, we take $x = x$ and $y = x^{m-1}$ to factor $\sqrt[n]{x^m}$ as $\sqrt[n]{x} \cdot \sqrt[n]{x^{m-1}}$; for the second factorization, we take $x = x$ and $y = x^{m-2}$ to factor $\sqrt[n]{x^{m-1}}$ as $\sqrt[n]{x} \cdot \sqrt[n]{x^{m-2}}$; etc., until the last factorization, when we take $x = x$ and $y = x$ to factor $\sqrt[n]{x^2}$ as $\sqrt[n]{x} \cdot \sqrt[n]{x}$. ■

Proposition 1.39: (Nested Radicals (aka, Nested Roots) Can be Combined)

For all natural numbers $n \geq 1$, $m \geq 1$ and real numbers $x > 0$,

$$\sqrt[n]{\sqrt[m]{x}} = \sqrt[n \cdot m]{x}.$$

In other symbols, $(x^{1/m})^{1/n} = (x)^{\frac{1}{m \cdot n}}$.

Important: $\sqrt[n]{\sqrt[m]{x}} = \sqrt[n \cdot m]{x}$ and NOT $\sqrt[n+m]{x}$. This is analogous to $(x^n)^m = x^{n \cdot m}$ and NOT x^{n+m} , because the exponential of an exponential means to multiply the exponents, NOT add them. See Prop. 1.49 for general “Power Towers”.

Proof: Let $z_1 := \sqrt[n]{x}$ so that $(z_1)^n = x$. Next, set $z_2 = \sqrt[m]{z_1}$ so that $(z_2)^m = z_1$. Combining these two results, we have

$$z_1 = (z_2)^n \implies \underbrace{(z_1)^m}_{x} = ((z_2)^n)^m = (z_2)^{n \cdot m} \implies x = (z_2)^{n \cdot m}.$$

Hence, we have that

$$x = (z_2)^{n \cdot m},$$

and therefore, $z_2 = \sqrt[n \cdot m]{x}$, that is, $\sqrt[n]{\sqrt[m]{x}} = \sqrt[n \cdot m]{x}$. ■

In a similar fashion, for any finite number of nested (integer) roots of $x > 0$, you can combine them into a single (integer) root of x by multiplying the integers specifying the nested roots. A more general procedure is called “denesting of radicals”. It’s pretty rad, so be careful if you use it.

Proposition 1.40: (Roots and Powers are Inverse Operations)

For all natural numbers $n \geq 1$ and real numbers $x > 0$,

$$\sqrt[n]{x^n} = x \text{ and } (\sqrt[n]{x})^n = x.$$

In other symbols, $(x^n)^{1/n} = x$ and $(x^{1/n})^n = x$. You can also write this as $(x^n)^{1/n} = x^{\frac{n}{n}} = x^1 = x$ and $(x^{1/n})^n = x^{\frac{n}{n}} = x^1 = x$.

Proof: Let $z_1 := \sqrt[n]{x^n}$. Then, by the definition of the n -th root, this holds if, and only if, $(z_1)^n = x^n$. Hence, $z_1 = x$ by standard properties of multiplication¹². This proves the first statement.

Next, let $z_2 := \sqrt[n]{x}$. Then, once again by the definition of the n -th root, this holds if, and only if, $(z_2)^n = x$. Hence, $(\sqrt[n]{x})^n = x$, which is what we wanted to show. ■

1.7.2 Rational Exponents or Powers

In the following, we seek to make sense of $x^{p/q}$. A key observation is that we want to obtain the same answer, whether we use $\sqrt[q]{x^p}$ or $(\sqrt[q]{x})^p$ as the definition. Just as with roots, for negative values of x , we have to be careful.

¹²Suppose that x and y are both positive real numbers. Then $x = y \iff x^n = y^n$ for all positive integers, n . Indeed, if $x = y$, then $x^n = y^n$ is immediate. Suppose, therefore, that $x \neq y$. Without loss of generality, we can assume $0 < x < y$. Then, $x \cdot x < x \cdot y < y \cdot y$, which proves that $x^2 < y^2$. Other powers are similar. Hence, $x^n = y^n$ and $x > 0, y > 0$ together imply that $x = y$.

Definition 1.41. *Rational Powers*

- (a) For $x \geq 0$ any non-negative real number, and for p and q any positive integers, then both $\sqrt[q]{x^p}$ and $(\sqrt[q]{x})^p$ return the same answer^a. Hence, we can define $x^{p/q}$ using either one. For definiteness, we define $x^{p/q} := \sqrt[q]{x^p}$. Of course, you are free to use $(\sqrt[q]{x})^p$.
- (b) For $x > 0$ any positive real number, and for p and q any positive integers, we define $x^{-p/q} := \frac{1}{x^{p/q}}$.
- (c) For $x < 0$ any negative real number, and for all p and q integers, we are leaving $x^{p/q}$ as **undefined**. The reason is indicated below: we do not want to deal with complex roots at this time.

If you find this treatment inadequate, [Wikipedia](#) provides more information on the meaning of rational exponents.

^aThis follows from Prop. 1.37

Remark 1.42. When p and q are both **odd integers**, it is OK to define $x^{p/q} := \text{sign}(x) \sqrt[q]{|x|^p}$. Some textbooks do this. For $p = 1$ and $q = 3$, this is consistent with how we defined the cube root of -1 . However, how would you define $x^{0.6}$ for $x < 0$? We are leaving it as **undefined** when $x < 0$ because 0.6 could be interpreted as $3/5$ or $6/10$, which is problematic, as we show next.

- $\sqrt[5]{(-1)^3} = \sqrt[5]{-1} = -1$ (when all of the powers and roots are odd, it all works), whereas
- $\sqrt[10]{(-1)^6} = \sqrt[10]{1} = +1$, where we added the plus sign for emphasis.

To avoid these inconsistencies, we recommend you adopt $x^{p/q}$ as being undefined for $x < 0$.

Note: Some programming languages, such as Julia, throw a TYPE¹³ error, when you try to compute $(-1)^{0.6}$, while others, such as MATLAB, return a complex number.

Example 1.43. Compute the following rational powers for $x = 3$ and $x = -5$:

(a) $x^{3/2}$

(b) $x^{2/3}$

(c) $x^{-3/2}$

(d) $x^{-2/3}$

(e) $x^{8/12}$

(f) $x^{3/5}$

Solution:

¹³The TYPE error in Julia arises when you are using real numbers (Float64 and Int64) in the expression $(-1)^{0.6}$, and yet the answer is a complex number. We show later that if you use numbers that are of complex TYPE, Julia returns a complex answer, and there is no longer a TYPE error. This is a FEATURE and not a BUG! There is no better way to screw up a computation than to throw in Complex Numbers when you had no intention of doing so. Julia warns you, while MATLAB assumes you know what you are doing. **Which one do you prefer?**

- | | |
|--|--|
| (a) $3^{3/2} = \sqrt{3^3} \approx 5.196$ | $(-5)^{3/2} = \text{undefined}$ |
| (b) $3^{2/3} = \sqrt[3]{3^2} \approx 2.080$ | $(-5)^{2/3} = \text{undefined}$ |
| (c) $3^{-3/2} = \frac{1}{3^{3/2}} \approx 0.192$ | $(-5)^{-3/2} = \text{undefined}$ |
| (d) $3^{-2/3} = \frac{1}{3^{2/3}} \approx 0.481$ | $(-5)^{-2/3} = \text{undefined}$ |
| (e) $3^{8/12} = \sqrt[12]{3^8} = \sqrt[3]{3^2}$ | $(-5)^{8/12} = \text{undefined}$ |
| (f) $3^{3/5} = \sqrt[5]{3^3} \approx 1.933$ | $(-5)^{3/5} = \text{undefined}$ |
| (f') $3^{3/5} = \sqrt[5]{3^3} \approx 1.933$ | $(-5)^{3/5} = -\sqrt[5]{5^3} \approx -2.627$ |

Note: While the solution in (f'), $(-5)^{3/5} = -\sqrt[5]{5^3}$, is **not wrong, neither is it a recommended solution.** Why? Because, if instead, you wrote $0.6 = 6/10$ for $\frac{3}{5}$, which seems totally legit, you would then likely define $(-5)^{6/10}$ as $\sqrt[10]{(-5)^6}$, which gives a positive number due to $(-5)^6 = ((-5)^2)^3 = 25^3 > 0$. When $0.6 = \frac{6}{10}$ is inconsistent with writing $0.6 = \frac{3}{5}$, you can have all sorts of confusion on hand. You have to be careful! ■

Julia has a command for x^y and not one for $x^{p/q}$, with p and q called out separately. Many programming languages are like this. You could create your own function, of course. Just out of curiosity, what happens if you use Julia's power command $^$ to solve the above problem?

```

1 p = 3
2 q = 2
3 a = p/q
4 @show a
5
6 x = 3
7 y = -5
8 @show x^a
9 @show y^a

```

Output

```

a = 1.5
x ^ a = 5.196152422706632

DomainError with -5.0:
Exponentiation yielding a complex result requires a complex argument.
Replace x^y with (x+0im)^y, Complex(x)^y, or similar.

```

Note: `im` is the Julia symbol for **i**, the imaginary number satisfying $(\mathbf{i})^2 = -1$. You learned in Algebra II how to use complex numbers when computing square roots of a negative number. Most of you have not learned how to use complex numbers to compute $\sqrt[3]{-1}$ or $\sqrt[8]{-1}$ for example, because, when you allow complex solutions¹⁴, the n -th root, $\sqrt[n]{-1}$, involves n numbers, with at most one of them real.

Example 1.44. Solve the following *SAT problem* for p/q :

$$\frac{\sqrt{x^5}}{\sqrt[3]{x^4}} = x^{\frac{p}{q}}.$$

¹⁴Wikipedia for the n-roots of unity.

Solution:

$$\begin{aligned}
 \sqrt{x^5} &= x^{5/2} \\
 \sqrt[3]{x^4} &= x^{4/3} \\
 &\Downarrow \\
 \frac{\sqrt{x^5}}{\sqrt[3]{x^4}} &= x^{5/2} \cdot x^{-4/3} \\
 &= x^{5/2-4/3} \\
 &= x^{7/6}.
 \end{aligned}$$

Hence, $p/q = 7/6$. ■

Example 1.45. Solve the following *Math Olympiad problem* for $K > 0$,

$$\sqrt{K \cdot \sqrt{K \cdot \sqrt{K}}} = 3.$$

Solution:

Squaring both sides gives $K \cdot \sqrt{K \cdot \sqrt{K}} = 3^2$. Then, using our notation, $x^{\frac{1}{2}} = \sqrt{x}$, and the two facts $x^a \cdot x^b = x^{a+b}$ and $(x^a)^b = x^{a \cdot b}$, we obtain

$$\begin{aligned}
 K \cdot \left(K \cdot K^{\frac{1}{2}}\right)^{\frac{1}{2}} &= 3^2 \\
 &\Updownarrow \\
 K \cdot \left(K^{\frac{3}{2}}\right)^{\frac{1}{2}} &= 3^2 \\
 &\Updownarrow \\
 K \cdot K^{\frac{3}{4}} &= 3^2 \\
 &\Updownarrow \\
 K^{\frac{7}{4}} &= 3^2 \\
 &\Updownarrow \\
 K^7 &= 3^8 \\
 &\Updownarrow \\
 K &= 3^{\frac{8}{7}}.
 \end{aligned}$$

Hence, $K \approx 3.509792$.

An alternative solution is to deduce $\sqrt{K \cdot \sqrt{K \cdot \sqrt{K}}} = 3 \iff \left(K \cdot \left(K \cdot K^{\frac{1}{2}}\right)^{\frac{1}{2}}\right)^{\frac{1}{2}} = 3 \iff K^{\frac{7}{8}} = 3$, and then to conclude $K = \left(K^{\frac{7}{8}}\right)^{\frac{8}{7}} = 3^{\frac{8}{7}}$. ■

1.7.3 Real Exponents

What is x^y , when the exponent, y , is a real number? It is traditional to define x^y using the natural logarithm, as introduced in the next section; it can also be defined using “limits” of rational powers, which can be reduced to integer roots and integer powers. The set of real numbers is roughly defined as “everything you can approximate arbitrarily closely with a rational number”. The University of Michigan’s Math 451 course makes this precise using the concepts of Cauchy sequences and equivalence classes, which are both outside the scope of this course. Let’s, nevertheless, consider **conceptually how to define x^π** , i.e., let’s try to have a general idea of the process. If we approximate the real number π by a rational number, $r = p/q$, then we can use that number in place of π , to compute

$$x^\pi \approx x^{p/q} := \sqrt[q]{x^p}.$$

Then, while holding x constant, we take successively “better” rational approximations of π (means, $|\pi - \frac{p}{q}|$ becomes closer to zero) and check that the result $x^{p/q}$ “tends” to a fixed value (means the values are only changing in more and more negative powers of ten). The same process works for an arbitrary real number y in place of π .

Of course, when we form a rational approximation of a real number y , p/q will not necessarily be a ratio of odd numbers, so we must only define x^y for $x \geq 0$.

What is x^y , when the exponent, y , is a real number?

Definition 1.46. For $x \geq 0$, x^y is the “limiting” value you obtain as you approximate y more and more closely by a rational number, such as a decimal expansion.

Table 1.47 illustrates this process for 2^π , that is, $x = 2$ and $y = \pi$. In HW, you will simply use Julia’s native power command, $^$. [Wikipedia](#) provides more information on real exponents.

Approximation of π	Approximation of 2^π
3	8.00000
$\frac{22}{7}$	8.51396
3.14	8.81524
3.14159	8.82498
3.141592653589793	8.824977827076287

Table 1.47: Computation of 2^π for different rational approximations of π .

Example 1.48. Compute the following powers using Julia.

(a) $(\sqrt{2})^\pi$

(b) $(\sqrt{2})^{-\pi}$

(c) $e^{2/3}$

(d) π^e

(e) $(-3)^e$

Solution:

(a) $(\sqrt{2})^\pi \approx 2.97069$

(b) $(\sqrt{2})^{-\pi} \approx 0.336623$

(c) $e^{2/3} \approx 1.94773$

(d) $\pi^e \approx 22.4592$

(e) $(-3)^e$ is **undefined**.

```

1 x = [sqrt(2) sqrt(2) exp(1) pi]
2 y = [pi -pi 2/3 exp(1)]
3 x.^y

```

Output

```
The 1×4 Matrix{Float64}:
2.97069  0.336623  1.94773  22.4592
```

```
1 # Two ways to access Euler's constant in Julia
2 @show e # obtained by \euler followed by the Tab key
3 println("The value of Euler's constant is ", e) # leaves it as the symbol e
4 @show exp(1);
5 println("The value of Euler's constant is ", exp(1))
6 @show e^e # definitely uses the value of Euler's constant
7 @show exp(1)^exp(1);
```

Output

```
e = e
The value of Euler's constant is e
exp(1) = 2.718281828459045
The value of Euler's constant is 2.718281828459045
e ^ e = 15.154262241479262
exp(1) ^ exp(1) = 15.15426224147926W
```

```
1 (-3)^exp(1)
```

Output

```
DomainError with -3.0:
Exponentiation yielding a complex result requires a complex argument.
Replace x^y with (x+0im)^y, Complex(x)^y, or similar.
```

Only shown to appease your curiosity!

```
1 (-3 + 0*im)^exp(1)
```

Output

```
-12.546688358893388 + 15.334119260562304im
```

Proposition 1.49: (Power Towers)

For all $x > 0$, and all $y, z \in \mathbb{R}$

$$(x^y)^z = x^{y \cdot z}. \quad (1.20)$$

From Prop. 1.49, we have that $(x^y)^z = (x^z)^y$ because both are equal to $x^{y \cdot z}$. Moreover, you can use this fact to reduce power towers of any (finite) length to a single power, such as, for all $x > 0$,

$$\left(((x^{y_1})^{y_2})^{y_3} \cdots \right)^{y_n} = x^{y_1 \cdot y_2 \cdot \cdots \cdot y_n}.$$

The proof of Prop. 1.49 is given in Chapter 1.19.

Proposition 1.50: Other Useful Properties of Real Exponents

Suppose that $x > 0$ is a positive real number. Then the following hold:

- (a) For all $y \in \mathbb{R}$, $x^y > 0$.
- (b) For all $y \in \mathbb{R}$, $x^{-y} = 1/x^y$.

- (c) For all $y, z \in \mathbb{R}$, $x^{y+z} = x^y \cdot x^z$.
- (d) For all $y > 0$, and $z \in \mathbb{R}$, $x^z \cdot y^z = (x \cdot y)^z$.

If you wanted to prove these properties, you would replace the exponents with rational approximations and then apply properties of rational powers.

Example 1.51. (*Optional Read:*) Is it possible for an irrational number raised to an irrational power to be a rational number, or even wilder, an integer?

Solution: If you got into engineering because you think math is fun, this question is for you. If you are taking math just to survive in engineering, then the logic of this problem's solution may appeal to you. See the solution to what **BriTheMathGuy** says “[...] Will Be Your Favorite Simple [Math] Problem” because it emphasizes the amazing power of basic logic. ■

1.8 Exponentials and Logarithms

When working with x^y , which variable are we changing: x , y , or both? Exponential functions remove this vagueness by fixing x to a constant value and then varying the exponent. The constant value is called a base, and following Descartes' notation, we should use a lowercase a , b , or c to denote it, right?

Exponential Function of Base a

Definition 1.52. For a real number $a > 0$, the **exponential function of base a** is $\exp_a : \mathbb{R} \rightarrow (0, \infty)$ by $\exp_a(x) := a^x$.

Remark 1.53. We do not define the exponential function for bases $a < 0$. See the discussion of rational powers, $x^{p/q}$. And because we will define a new “inverse operation” called the “logarithm”, we also do not define the exponential function for $a = 0$.

Remark 1.54. The domain of $\exp_a(x)$ is all real numbers, and its codomain is all positive real numbers. In fact, its range is also all positive real numbers.

Proposition 1.55: Properties of the Exponential Function

Suppose that $a > 0$ is a positive real number. Then the following hold:

- (a) For all $x \in \mathbb{R}$, $\exp_a(x) > 0$.
- (b) If $a > 1$, then $\exp_a(x)$ is strictly increasing, that is, for $y > x$, $\exp_a(y) > \exp_a(x)$.
- (c) If $0 < a < 1$, then $\exp_a(x)$ is strictly decreasing, that is, for $y > x$, $\exp_a(y) < \exp_a(x)$.
- (d) If $a = 1$, then $\exp_a(x) := 1^x = 1$ for all $x \in \mathbb{R}$.
- (e) For all $x \in \mathbb{R}$, $\exp_a(-x) = 1/\exp_a(x)$.
- (f) For all $x, y \in \mathbb{R}$, $\exp_a(x + y) = \exp_a(x) \cdot \exp_a(y)$.
- (g) For all $x, y \in \mathbb{R}$, $(\exp_a(x))^y = \exp_a(x \cdot y)$.

Remark 1.56. Proposition 1.55-(d) shows that $\exp_1(x)$ is a constant function and is not very useful! Also, because the function is constant, it cannot have an inverse.

Because $0 < a < 1 \iff 1 < \frac{1}{a} = a^{-1} < \infty$, it follows that for all $x \in \mathbb{R}$ and $0 < a < 1$, $\exp_a(x) = \exp_{\frac{1}{a}}(-x)$. Hence, it is not really necessary to use exponentials for bases $0 < a < 1$ as they can be easily replaced with exponentials of a base greater than one. This motivates some textbooks to only define the logarithm for bases larger than one.

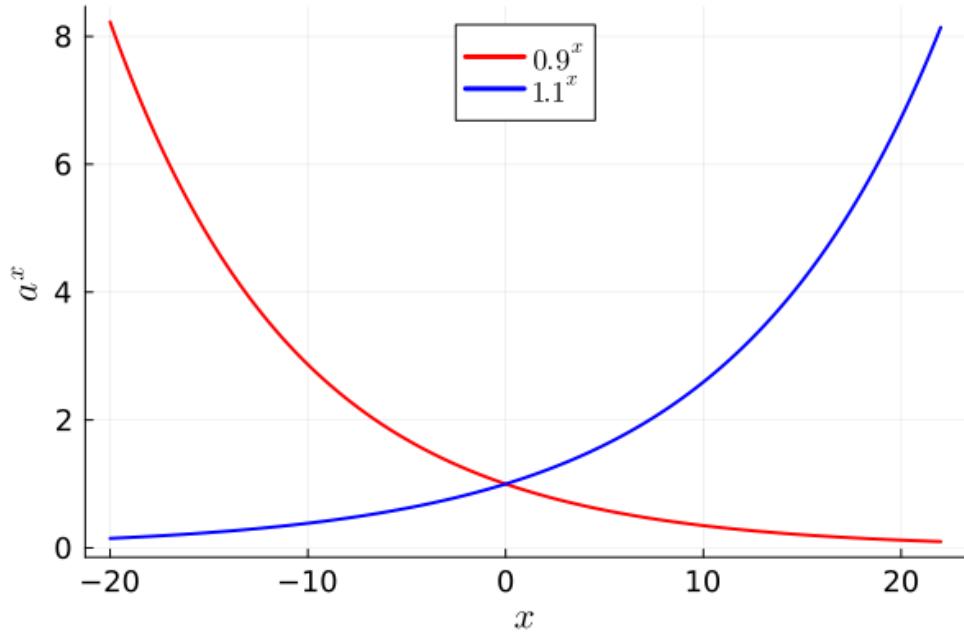


Figure 1.10: Shown are plots of $\exp_a(x) := a^x$ for $a = 0.9$ and $a = 1.1$. We can observe the trends. For $0 < a < 1$, a^x is strictly decreasing; in particular, $a^x \rightarrow \infty$ as $x \rightarrow -\infty$, and $a^x \rightarrow 0$ as $x \rightarrow \infty$. For $a > 1$, however, a^x is strictly increasing, and the trends are flipped, namely $a^x \rightarrow 0$ as $x \rightarrow -\infty$ and $a^x \rightarrow \infty$ as $x \rightarrow \infty$. Here, you can read the symbol \rightarrow as “tends to” or “approaches”.

Logarithm Base a

Definition 1.57. For real numbers $a > 0$, $a \neq 1$, and $x > 0$,

- $y \in \mathbb{R}$ is the **logarithm base a** of x if $a^y = x$, that is, $\exp_a(y) = x$.
- The notation is $y = \log_a(x)$, and is read, y equals log base a of x .
- The **logarithm function of base a** is $\log_a : (0, \infty) \rightarrow \mathbb{R}$, that is, it has domain the positive real numbers and range the real numbers.

Note: For $a = 1$ and $x \neq 1$, there is no $y \in \mathbb{R}$ such that $a^y = x$. That is why $a = 1$ has to be excluded. From the definition $y = \log_a(x) \iff a^y = x$ and our power rules, it follows that $\log_{\frac{1}{a}}(x) = -\log_a(x)$.

Remark 1.58. If you want to see an example worked out for logarithms to base $\frac{1}{3}$, you may enjoy this video by the **Math Sorcerer**, aka, **Daniel Gabilondo**, a former math professor at Seminole College^a. His YouTube channel is very down-to-earth. This video also talks about how math can be hard for anyone. Prof. Gabilondo traveled after High School and then restarted with math at age 24. Today, we would call that taking six gap years.

^aHere's a comment on Prof. Gabilondo from one of his students: “His class alone, his teaching style, turned everything I thought I knew about math on its head and I not only [under]stood the math, but why it was useful and how it was derived. I excelled, got an A, and have since gone on to self-study (through MIT open courses) linear algebra as well as complex analysis. [...] I found that I'm super hard on myself and if I didn't understand something, I would just quit. But that doesn't help anything in your life.”

Note how the definition of the logarithm follows a pattern similar to the definition of the n -th root. They are both defined “indirectly”,

$$y = \sqrt[n]{x} \iff y^n = x \quad \text{and} \quad y = \log_a(x) \iff a^y = x; \quad (1.21)$$

for roots, y is “downstairs” while for logarithms, y is “upstairs”. Specifically, for y to be an n -th root, we need to **raise y to power n** to recover the original number, whereas, for y to be a logarithm base a , we need to **raise a to power y** to recover the original number. Similar to the n -th root, we can compute logarithms via the Bisection Algorithm.

Example 1.59. For $x = 10$, compute its logarithm base 2, base e , base π , and base 10 using the Bisection Algorithm.

Solution: We gave the Bisection Algorithm in Chapter 1.7. Applying as shown below yields

a (base)	x	$\log_a(x)$	Error
2.0	10.0	2.302585	$\pm 3.337860 \times 10^{-6}$
e	10.0	3.321927	$\pm 3.337860 \times 10^{-6}$
π	10.0	2.011467	$\pm 1.668930 \times 10^{-6}$
10.0	10.0	1.000000	$\pm 8.344650 \times 10^{-7}$

Table 1.60: Computation of $\log_a(x)$ using the Bisection Algorithm. Note that rounding 0.9999997615814209 to six decimal points gave the 1.0 for $\log_{10}(10)$.

```

1 bases = [2.0 exp(1) π 10.0]
2 x = 10.0
3
4 for k = 1:length(bases)
5   a = bases[k]
6   # Define the function whose root we want to find
7   f(y) = a^y - x
8
9   # Call the bisection method on this function
10  F = bisection(f, 0.5, 4.0, 1e-5)
11  @show F
12
13
14  println("The log base $a is approximately ", F.est, )
15  println("The a to power loga(x) is approximately ", a^F.est)
16  println(" ")
17 end

```

Output

```

F = (est = 3.321927070617676, low = 3.321927070617676, up = 3.3219337463378906,
error = 3.337860107421875e-6)
The log base 2.0 is approximately 3.321927070617676
The a to power loga(x) is approximately 9.999992900306067

F = (est = 2.302584648132324, low = 2.302584648132324, up = 2.302591323852539,
error = 3.337860107421875e-6)
The log base 2.718281828459045 is approximately 2.302584648132324
The a to power loga(x) is approximately 9.999995551383774

F = (est = 2.0114665031433105, low = 2.011463165283203, up = 2.0114665031433105,
error = 1.6689300537109375e-6)
The log base 3.141592653589793 is approximately 2.0114665031433105
The a to power loga(x) is approximately 10.000007275393529

F = (est = 0.9999997615814209, low = 0.9999997615814209, up = 1.0000014305114746,
error = 8.344650268554688e-7)
The log base 10.0 is approximately 0.9999997615814209
The a to power loga(x) is approximately 9.999994510210845

```

Proposition 1.61: Properties of Logarithms (aka, Logarithm Rules)

For all bases $a > 0, a \neq 1$,

$$\begin{aligned}
&a^{\log_a(x)} = x, \text{ for all } x > 0. \\
&\log_a(a^x) = x, \text{ for all } x \in \mathbb{R} \\
&\log_a(a) = 1 \\
&\log_a(1) = 0 \\
&\log_a(x \cdot y) = \log_a(x) + \log_a(y), \text{ for all } x > 0, y > 0 \\
&\log_a\left(\frac{x}{y}\right) = \log_a(x) - \log_a(y), \text{ for all } x > 0, y > 0 \\
&\log_a(x^y) = y \log_a(x), \text{ for all } x > 0, y \in \mathbb{R}.
\end{aligned} \tag{1.22}$$

In addition,

$$\begin{aligned}
a > 1 \implies \log_a(x) \text{ is strictly increasing for all } x > 0 \\
0 < a < 1 \implies \log_a(x) \text{ is strictly decreasing for all } x > 0.
\end{aligned}$$

Remark 1.62. The first two lines of Prop. 1.61 can be written as

$$\begin{aligned}
\exp_a(\log_a(x)) &= x, \text{ for all } x > 0, \text{ and} \\
\log_a(\exp_a(x)) &= x, \text{ for all } x \in \mathbb{R},
\end{aligned} \tag{1.23}$$

which underlines that $\log_a(x)$ is the inverse function of $\exp_a(x)$ and vice versa. ■

Remark 1.63. Suppose that $f : \mathbb{R} \rightarrow \mathbb{R}$ is a function and for some $y \in \mathbb{R}$, $f(y) > 0$. Then $a^{\log_a(f(y))} = f(y)$. This follows from “ $a^{\log_a(x)} = x$ ” by taking $x = f(y)$. You might call this a “trick”, but it’s not really a trick. It’s just basic math. Similarly, taking $x = f(y)$, we have that $\log_a(a^{f(y)}) = f(y)$ for all $f(y)$. ■

Common logarithms are \log_{10} (log base 10), \log_2 (log base two), and \log_e (log base e , Euler’s constant, the irrational number 2.71828182846...). In Julia, the notation is

- `log(x)` is used for $\log_e(x)$
- `log10(x)` is used for $\log_{10}(x)$
- `log2(x)` is used for $\log_2(x)$
- you can generate Euler’s constant, e , by typing “`\euler`” (without the quotes) and then pressing the TAB key. You can also use “`e=exp(1)`”; the latter has the advantage of being understood by almost everyone.

Example 1.64. Solve for x in the following problem from a *Chinese Math Olympiad*:

$$9^x - 6^x = 4^x.$$

Solution: We’ll give two solutions, one analytical and one numerical. First, the analytical solution.

$$\begin{aligned}
9^x - 6^x - 4^x &= 0 \quad (\text{express numbers as products of primes}) \\
&\Downarrow \\
3^{2x} - 2^x \cdot 3^x - 2^{2x} &= 0 \quad (\text{divide both sides by } 3^{2x}) \\
&\Downarrow \\
1 - \left(\frac{2}{3}\right)^x - \left(\frac{2}{3}\right)^{2x} &= 0 \quad (\text{look for structure in the equation}).
\end{aligned} \tag{1.24}$$

Letting $y = \left(\frac{2}{3}\right)^x$ yields a quadratic equation

$$\begin{aligned} 1 - y - y^2 &= 0 \\ \Downarrow \\ y^2 + y - 1 &= 0 \\ \Downarrow \\ y &= \frac{-1 \pm \sqrt{5}}{2}. \end{aligned} \tag{1.25}$$

Because $y := \left(\frac{2}{3}\right)^x > 0$, we have that

$$\begin{aligned} \left(\frac{2}{3}\right)^x &= \frac{-1 + \sqrt{5}}{2} \quad (\text{take log on both sides, can use any base}) \\ \Downarrow \\ x \log_{10} \left(\frac{2}{3}\right) &= \log_{10} \left(\frac{-1 + \sqrt{5}}{2}\right) \quad (\text{logarithm rules bring } x \text{ down so that we can solve for it}) \\ \Downarrow \\ x &= \frac{\log_{10} \left(\frac{-1 + \sqrt{5}}{2}\right)}{\log_{10} \left(\frac{2}{3}\right)} \quad (\text{can stop here; it is a valid answer, or}) \end{aligned} \tag{1.26}$$

≈ 1.1868 (do the indicated operations, in Julia or on calculator, to get a decimal approx.)

Second, a numerical solution using the Bisection Algorithm.

```

1 f(x) = 9.0^x - 6.0^x - 4.0^x
2 F = bisection(f, 1.0, 4.0, 1e-9)

```

Output

```
(est = 1.1868143903557211, low = 1.1868143896572292, up = 1.1868143903557211,
error = 3.4924596548080444e-10)
```

Which solution approach did you prefer? The one where you are more likely to obtain a correct answer? The Bisection Algorithm? In High School, they did not let you choose your preferred method. Engineering is a practical science!

Example 1.65. Solve the following exponential Diophantine equation

$$2^x + x^3 = 4.$$

Solution: We observe that $2^1 + 1^3 = 3 < 4$ and $2^2 + 2^3 = 12 > 4$. Hence, we can use $a = 1$ and $b = 2$ as bracketing points in the Bisection Algorithm. The result is $x = 1.1956916 \pm 2.384185e-7$. ■

Example 1.66. Practice with logarithm rules, roots, and quadratic equations. Solve the following problem from *Vijay Maths*, which in turn comes from one of the Math Olympiads:

(a) Find all real x satisfying

$$x^{16 \cdot (\log_5(x))^3 - 68 \cdot \log_5(x)} = 5^{-16}. \tag{1.27}$$

(b) Compute the product of all solutions to the above equation.

Hints: Useful logarithm properties: for any $b > 0$, $b \neq 1$,

$$\begin{aligned}\log_b(x^{f(x)}) &= f(x) \cdot \log_b(x) \\ \log_b(b^c) &= c.\end{aligned}\tag{1.28}$$

Solution: We first take the log base 5 of both sides of (1.27), yielding

$$(16 \cdot (\log_5(x))^3 - 68 \cdot \log_5(x)) \cdot \log_5(x) = -16.$$

Next, we make the substitution $a = \log_5(x)$, giving us

$$(16 \cdot a^3 - 68 \cdot a) \cdot a = -16,$$

which, after diving both sides by -4 and rearranging terms, gives

$$4a^4 - 17a^2 + 4 = 0.$$

The above is a quadratic equation in a^2 . Indeed, if we let $y = a^2$, we have

$$4y^2 - 17y + 4 = 0,$$

and the quadratic formula provides us with

$$y = \frac{17 \pm \sqrt{17^2 - 4^3}}{8} = \frac{17 \pm \sqrt{225}}{8} = \frac{17 \pm 15}{8} = 4 \text{ and } \frac{1}{4}.$$

Recalling that $y = a^2$ and that for $c > 0$, $a^2 = c \iff a = \pm\sqrt{c}$, gives us

$$a \in \{-2, 2, -\frac{1}{2}, \frac{1}{2}\}.$$

Then, because $a = \log_5(x) \iff x = 5^a$, the answer to part (a) is

$$x \in \{5^{-2}, 5^2, 5^{-\frac{1}{2}}, 5^{\frac{1}{2}}\}.$$

Finally, the answer to part (b) is the product of these solutions, namely

$$\text{Ans} = 5^{-2} \cdot 5^2 \cdot 5^{-\frac{1}{2}} \cdot 5^{\frac{1}{2}} = 5^{-2+2-\frac{1}{2}+\frac{1}{2}} = 5^0 = 1.$$

While the roots of $x^{16 \cdot (\log_5(x))^3 - 68 \cdot \log_5(x)} - 5^{-16} = 0$ can be determined numerically via the Bisection Algorithm, it is challenging to find the correct bracketing points for each of the four solutions. In many engineering problems, a combination of analytical and numerical prowess is required. ■

Remark on Logarithms: In High School, logarithms terrify many students. The whole idea of being given a base a and a number x and being asked to find y such that $a^y = x$, somehow, seems very difficult and even more pointless. In Engineering, $\log_a(x)$ is a function with really useful properties. We employ logarithms frequently, though we will never compute the logarithm of a number using its definition; we have computers to do the dirty work for us. **What we need to know as engineers are the properties in Prop. 1.61; bookmark it, and you'll be fine. Use the properties enough, and you won't need the bookmark anymore.**

1.9 One Logarithm to Rule Them All

In engineering practice, we mostly use the so-called **natural logarithm**, which is log base e , Euler's constant. In Mathematics, the natural logarithm is denoted $\ln(x)$, where the "n" is for natural. The corresponding exponential function base e is written e^x or $\exp(x)$, without a subscript. The terminology "natural logarithm" comes from Calculus because the logarithm of base e is the only

Log Base a Rules	Natural Log Rules
$\log_a(x \cdot y) = \log_a(x) + \log_a(y)$	$\ln(x \cdot y) = \ln(x) + \ln(y)$
$\log_a\left(\frac{x}{y}\right) = \log_a(x) - \log_a(y)$	$\ln\left(\frac{x}{y}\right) = \ln(x) - \ln(y)$
$\log_a(x^y) = y \cdot \log_a(x)$	$\ln(x^y) = y \cdot \ln(x)$
$\log_a(a^x) = x$	$\ln(e^x) = x$
$a^{\log_a(x)} = x$	$e^{\ln(x)} = x$

Table 1.67: Side-by-side Comparison of Logarithm Rules and Natural Logarithm Rules (they are identical).

logarithm with a “slope” equal to 1 at $x = 0$ and its exponential function e^x is the only one of the exponential functions, a^x , that is equal to its **own derivative**. For now, we’ll paraphrase Shakespeare and say, “A logarithm by any other justification smells as natural”. The logarithm rules make it easy to prove the result on “power towers”, Prop. 1.49. In other words, does the equation

$$(x^y)^z = x^{yz} \quad (1.29)$$

hold for $x > 0$? If one of y or z is zero, then the result is immediate because $x^0 = 1$ on both sides. Hence, we assume both x and y are positive. Taking the natural logarithm of the left side of (1.29) and applying the logarithm rules give us

$$\ln((x^y)^z) = z \cdot \ln(x^y) = z \cdot y \cdot \ln(x).$$

Next, taking the natural logarithm of the right side of (1.29) and applying the logarithm rules give us

$$\ln(x^{yz}) = y \cdot z \cdot \ln(x).$$

Hence, when y and z are nonzero,

$$\ln((x^y)^z) = \ln(x^{yz}).$$

Taking the exponential of both sides of the above equation gives

$$(x^y)^z = x^{yz},$$

and, hence, we have proved that (1.29) is true. ■

It is sometimes useful to know how to compute all logs from just the natural log. From $x = a^{\log_a(x)}$, taking the natural log of both sides yields,

$$\ln(x) = \log_a(x) \cdot \ln(a).$$

Rearranging gives the following result.

Proposition 1.68: All Logs from the Natural Log

The natural logarithm can be used to compute any logarithm by the **change of base formula**: for $a > 0$, $a \neq 1$, and $x > 0$,

$$\log_a(x) = \frac{\ln(x)}{\ln(a)}. \quad (1.30)$$

A similar formula is true for any base $b > 0$, $b \neq 1$ replacing base e , though it is almost never used.

Example 1.69. Compute $\log_a(\pi)$ for $a \in \{1.1, 2.0, e, \pi, 10.0\}$ using (1.30).

Solution:

```

1 function myLog (x: :Real, a: :Real)
2   if ! (x>0)
3     println ("x must be positive")
4     return NaN
5   end
6   if ! (a>0)
7     println ("the base-a must be positive")

```

x	a	$\log_a(x) = \ln(x)/\ln(a)$	$a^{\log_a(x)}$
3.14159	1.10000	12.0106	3.14159
3.14159	2.00000	1.6515	3.14159
3.14159	2.71828	1.1447	3.14159
3.14159	3.14159	1.0000	3.14159
3.14159	10.0000	0.4972	3.14159

Table 1.70: Computing logarithms to an arbitrary base using the natural logarithm and (1.30).

```

8   return NaN
9   end
10  return log(x) /log(a)
11 end
12
13 a = [1.1, 2.0, exp(1), pi, 10.0]
14 x = pi
15 y = 0*a
16
17 for i = 1:length(a)
18     y[i] = myLog(x, a[i])
19 end
20
21 [x*ones(length(a), 1) a y a.^y]
```

Output

```
5×4 Matrix{Float64}:
3.14159    1.1      12.0106    3.14159
3.14159    2.0      1.6515     3.14159
3.14159    2.71828  1.14473   3.14159
3.14159    3.14159  1.0        3.14159
3.14159    10.0     0.49715   3.14159
```

■

1.10 Set Notation and How to Make Sense of It

In college mathematics, precise notation allows us to express complex ideas clearly and concisely. However, it can feel intimidating at first. Don't worry! With practice, you'll become comfortable. Let's break it down through some specific examples.

What Does This Notation Mean?

Consider the following set:

$$S_A := \{x \in (0, \infty) \mid |\sin(x)| = 1.0\}.$$

Step-by-step explanation:

- 1. Domain restriction:** The symbol $x \in (0, \infty)$ means we are only considering positive real numbers. This is similar to specifying the domain of a function.
- 2. The vertical bar “|”:** This symbol means “such that.” The set includes only those x values that satisfy the condition after the bar, namely, $|\sin(x)| = 1.0$.
- 3. Interpretation:** This means we are selecting all x in $(0, \infty)$ such that the absolute value of $\sin(x)$ is exactly 1.

Enumerating the Set

It is sometimes possible and useful to list all of the elements in a set to make things more concrete. To do this for S_A , recall key properties of $\sin(x)$:

- $\sin(x) = 1$ at $x = \frac{\pi}{2} + 2k\pi$ for $k \in \mathbb{Z}$.
- $\sin(x) = -1$ at $x = \frac{3\pi}{2} + 2k\pi$ for $k \in \mathbb{Z}$.

The absolute value, $|\sin(x)|$, is periodic with period π . Therefore,

$$|\sin(x)| = 1 \implies x = \frac{\pi}{2} + k\pi, \quad k = 0, 1, 2, \dots$$

Since we are restricted to $x > 0$, negative values of k are excluded. Thus, we can write,

$$S_A = \left\{ \frac{\pi}{2} + k\pi \mid k = 0, 1, 2, \dots \right\}.$$

This step of **enumerating the elements of the set** is central to understanding **countability**—a property that will arise later.

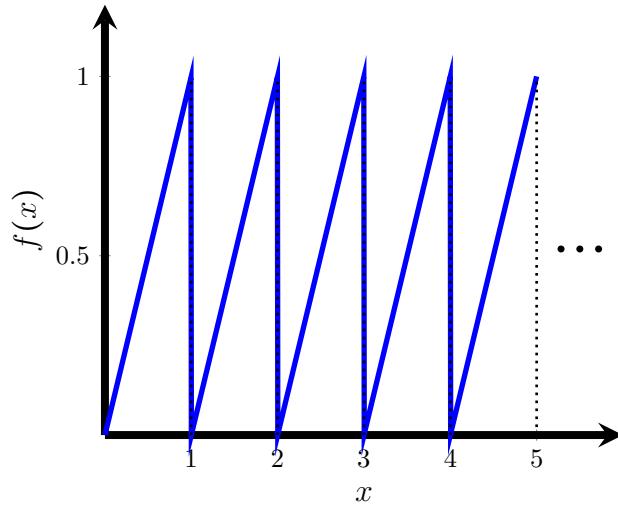


Figure 1.11: A **sawtooth function with period one**. It is plotted in the usual manner, with a vertical line connecting the points $(k, 1)$ and $(k, 0)$ for $k = 1, 2, 3, \dots$, which, in theory, makes it not a function, but in reality, we ignore the vertical lines because they are “plotting artifacts” left to guide our eyes.

Working With More Complex Sets

Let’s build on this with another example, inspired by the **sawtooth wave** function, $f(x)$, defined as,

$$f(x) = \begin{cases} x, & 0 \leq x < 1, \\ x \bmod 1, & x \geq 1. \end{cases}$$

The goal is to understand sets defined by conditions such as,

$$S_B := \{x \in (0, \infty) \mid f(x) = \frac{0.9}{i}, i = 1, 2, 3, \dots\}.$$

Breaking it down: For simplicity, let’s first analyze subsets of S_B based on intervals of x . For example:

- Points in $(0, 1]$ that satisfy $f(x) = \frac{0.9}{i}$ are:

$$S_{B,1} = \{x \in (0, 1] \text{ such that } x = \frac{0.9}{i}, i = 1, 2, 3, \dots\}.$$

- Points in $(1, 2]$ that satisfy $f(x) = \frac{0.9}{i}$ are:

$$S_{B,2} = \{x \in (1, 2] \text{ such that } x = 1 + \frac{0.9}{i}, i = 1, 2, 3, \dots\}.$$

- Similarly, for $(j - 1, j]$, we can write:

$$S_{B,j} = \{x \in (j - 1, j] \text{ such that } x = (j - 1) + \frac{0.9}{i}, i = 1, 2, 3, \dots\}.$$

Key Takeaways

1. **Set-builder notation** uses conditions to define a set, like restricting a domain or specifying a property.
2. If you are confused about a set, writing out many of the “mysterious symbols” in your favorite spoken language can be helpful.

Example 1.71. Here are some practice problems.

- (a) Write in set notation: The set of all positive integers divisible by 3 (same as being a multiple of 3).
- (b) Enumerate the elements of the set $T := \{x \in (0, 1] \mid x = \frac{1}{i}, i = 2, 3, 4, \dots\}$.
- (c) Explain what is meant by $\{n \in \mathbb{N} \mid n \bmod 3 = 0\}$ and then, how is the set different from $\{n \in \mathbb{Z} \mid n \bmod 3 = 0\}$.

Solution:

- (a) Let's define the set $S := \{x \in \mathbb{N} \mid x = 3k, k \in \mathbb{N}\}$. This means that S contains all positive integers that can be written as a multiple of 3, where k is a positive integer. Enumerating a few elements of S , we find

$$S = \{3, 6, 9, 12, \dots\}.$$

- (b) The condition $x = \frac{1}{i}$ means that x is the reciprocal of a positive integer $i \geq 2$. Substituting successive values of i , we get:

$$T = \left\{ \frac{1}{2}, \frac{1}{3}, \frac{1}{4}, \frac{1}{5}, \dots \right\}.$$

Observe that T is infinite, but we can systematically list all its elements, helping us to understand the set.

- (c) The set $\{n \in \mathbb{N} \mid n \bmod 3 = 0\}$ consists of all natural numbers (positive integers) that are divisible by 3. In terms of set notation, this can be written as:

$$\{n \in \mathbb{N} \mid n \bmod 3 = 0\} = \{3, 6, 9, 12, \dots\}.$$

This is the set of positive multiples of 3, starting from 3, since \mathbb{N} only includes positive integers.

The set $\{n \in \mathbb{Z} \mid n \bmod 3 = 0\}$, however, includes all integers divisible by 3, zero as well as positive and negative integers. In terms of set notation, this is,

$$\{n \in \mathbb{Z} \mid n \bmod 3 = 0\} = \{\dots, -9, -6, -3, 0, 3, 6, 9, \dots\}.$$

The key difference is that the first set (with “domain” \mathbb{N}) only contains positive integers, while the second set (with “domain” \mathbb{Z}) includes negative integers, zero, and positive integers, representing all integer multiples of 3.

1.11 Euler's Formula

He did a lot of stuff, Euler! Somehow, in popular culture, we know some of the famous physicists, such as Newton, Curie, and Einstein, but know little about the mathematicians. Euler is one to know. You can learn more [here](#), if you wish.

Many of us do not learn Euler's Formula in High School because it comes from Calculus. However, Euler's Formula makes working with complicated trig expressions a breeze, and thus we will cover it now by introducing it as a definition. When we cover Taylor Expansions (McLaurin Expansions), we'll be able to replace the definition with a theorem.

Euler's Formula

Definition 1.72. Let \mathbf{i} denote $\sqrt{-1}$ (that is, $\mathbf{i}^2 = -1$). Then, for $x \in \mathbb{R}$,

$$e^{\mathbf{i}x} := \cos(x) + \mathbf{i}\sin(x). \quad (1.31)$$

Let's note that $e^{-\mathbf{i}x} = e^{\mathbf{i}(-x)} := \cos(-x) + \mathbf{i}\sin(-x) = \cos(x) - \mathbf{i}\sin(x)$, because cosine is an even function (that is, $\cos(-x) = \cos(x)$) and sine is an odd function (that is, $\sin(-x) = -\sin(x)$). Hence,

$$e^{\mathbf{i}x} + e^{-\mathbf{i}x} = 2\cos(x),$$

yielding

$$\cos(x) = \frac{e^{\mathbf{i}x} + e^{-\mathbf{i}x}}{2}. \quad (1.32)$$

Similarly,

$$e^{\mathbf{i}x} - e^{-\mathbf{i}x} = 2\mathbf{i}\sin(x),$$

yielding

$$\sin(x) = \frac{e^{\mathbf{i}x} - e^{-\mathbf{i}x}}{2\mathbf{i}}. \quad (1.33)$$

At first blush, these formulas seem a needless complication of two perfectly manageable functions, $\sin(x)$ and $\cos(x)$. We get it. Do you recall, however, those crazy trig reduction formulas you had to derive and memorize? They can be derived much more easily through Euler's formula. Here is our first example.

Example 1.73. Derive a simpler formula for $\cos^2(x) := (\cos(x))^2$.

Solution:

$$\begin{aligned} \cos^2(x) &= \left(\frac{e^{\mathbf{i}x} + e^{-\mathbf{i}x}}{2} \right)^2 \\ &= \frac{e^{\mathbf{i}2x} + 2 + e^{-\mathbf{i}2x}}{4} \\ &= \frac{1}{2} + \frac{1}{2} \frac{e^{\mathbf{i}2x} + e^{-\mathbf{i}2x}}{2} \\ &= \frac{1}{2} + \frac{1}{2} \cos(2x). \end{aligned} \quad (1.34)$$

The key “trick” is to use the usual exponential rules, so that $e^{\mathbf{i}x} \cdot e^{\mathbf{i}x} = e^{\mathbf{i}2x}$ and $e^{\mathbf{i}x} \cdot e^{-\mathbf{i}x} = e^0 = 1$, and so on. ■

Here is a more complicated example.

Example 1.74. Use Euler's formula to derive a simpler expression for $\sin^2(x)\cos(x)$.

Solution: We start by writing out $\sin(x)$ and $\cos(x)$ in terms of exponential functions and substituting these expressions into

$\sin^2(x) \cos(x)$, yielding,

$$\begin{aligned}
\sin^2(x) \cos(x) &= \left(\frac{e^{ix} - e^{-ix}}{2i} \right)^2 \left(\frac{e^{ix} + e^{-ix}}{2} \right) \text{ (square the first term)} \\
&= \frac{(e^{i2x} - 2 + e^{-i2x})(e^{ix} + e^{-ix})}{(-4)(2)} \text{ (multiply out the terms)} \\
&= \frac{e^{i3x} - 2e^{ix} + e^{-ix} + e^{ix} - 2e^{-ix} + e^{-i3x}}{-8} \text{ (add like terms and group)} \\
&= \frac{(e^{ix} + e^{-ix}) - (e^{i3x} + e^{-i3x})}{8} \\
&= \frac{1}{4} \frac{(e^{ix} + e^{-ix})}{2} - \frac{1}{4} \frac{(e^{i3x} + e^{-i3x})}{2} \\
&= \frac{1}{4} \cos(x) - \frac{1}{4} \cos(3x).
\end{aligned} \tag{1.35}$$

Fact 1.75: Consequence of Euler's Formula

Any (finite) product of powers of sines and cosines,

$$\prod_{k=1}^N \sin^{n_k}(\omega_k x) \cdot \cos^{m_k}(\omega_k x), \tag{1.36}$$

with n_k and m_k positive integers, can be expressed as a linear combination of sines and cosines to the first power,

$$a_0 + \sum_{i=1}^J a_i \sin(\nu_i x) + \sum_{i=1}^K b_i \cos(\nu_i x), \tag{1.37}$$

plus a constant, where the ν_i 's are sums of $\pm \omega_k$, $1 \leq k \leq N$. We recall from ROB 101 that linear combinations are always finite.

Note: In (1.35), $N = 1$, $\omega_1 = 1$, $\nu_1 = \omega_1$ and $\nu_2 = \omega_1 + \omega_1 + \omega_1 = 3\omega_1$.

Note: See Table 1.89 and Chapter 1.18.2 for a convenient set of power reductions for sine and cosine.

You've earned a bit of hilarity because only you, a denizen of Calculus, can even begin to appreciate this video: All The Math References Frame by Frame From Animation vs. Math A link to the un-commented video is in the description. Even your author found the running commentary helpful! The video is deep and a stitch at the same time. Warning: Video Game style of violence, stick figures vs math symbols, no blood. Includes extra dimensions and time travel.

1.12 Hyperbolic Trig Functions and Relation to Euler's Formula

The ordinary trigonometric functions arise by studying the unit circle. The **hyperbolic trigonometric functions** arise from the study of the unit hyperbola. A really great explanation is given in [Hyperbolic functions and the unit hyperbola](#), from Khan Academy's pre-calculus series. If you become an Aero, Civil, or Mechanical Engineer, you may encounter them daily, as explained in [The applications of hyperbolic trig: Why do we even care about these things?](#) by Zach Star. We introduce them here due to their parallels with Euler's formula in Def. 1.72.

Hyperbolic Trig Functions

Definition 1.76. For all $x \in \mathbb{R}$,

$$\begin{aligned}\cosh(x) &:= \frac{e^x + e^{-x}}{2} \\ \sinh(x) &:= \frac{e^x - e^{-x}}{2} \\ \tanh(x) &:= \frac{\sinh(x)}{\cosh(x)} \\ \text{etchn} &:= \dots\end{aligned}\tag{1.38}$$

The functions are called hyperbolic cosine, hyperbolic sine, hyperbolic tangent, and hyperbolic etc. (of course!).

The hyperbolic trig functions play a minor role in this course. We introduce them here to help you digest Euler's Formula. We note that

$$\begin{aligned}\cos(x) &= \cosh(ix) = \frac{e^{ix} + e^{-ix}}{2} \\ \sin(x) &= \sinh(ix) = \frac{e^{ix} - e^{-ix}}{2},\end{aligned}\tag{1.39}$$

even though the hyperbolic functions were discovered before Euler's formula.

We would be remiss not to mention how the hyperbolic functions show up in Architecture, in the form of the **Catenary Arch**, which is a mainstay of the beautiful buildings by **Gaudi**. Also, you know that $\sin^2(x) + \cos^2(x) = 1$. The hyperbolic version is

$$\cosh^2(x) \underset{\uparrow}{-} \sinh^2(x) = 1,$$

(not a typo) as you can multiply out and check for yourself.

Function	Principal Value	Domain
$\operatorname{acosh}(x)$	$\operatorname{acosh}(x) = \ln(x + \sqrt{x^2 - 1})$	$x \geq 1$
$\operatorname{asinh}(x)$	$\operatorname{asinh}(x) = \ln(x + \sqrt{x^2 + 1})$	All $x \in \mathbb{R}$

Table 1.77: Principal domains for $\operatorname{acosh}(x)$ and $\operatorname{asinh}(x)$.

In case you want to learn more:

- [Hyperbolic Trig Identities](#) by the The Organic Chemistry Tutor.
- [Why hyperbolic functions are actually really nice](#) by Dr. Trefor Bazett.

1.13 Summing Symbol and Change of Indices

The summing symbol $\sum_{i=1}^n x_i$ is shorthand for $x_1 + x_2 + \dots + x_n$. In Julia, it can be implemented with a `for` loop or with the command `sum(x)`, if x is a vector with components x_i .

It is sometimes useful to perform a **change of index** in a sum. For example,

$$\sum_{i=1}^n (i-1) = \sum_{j=0}^{n-1} j = \sum_{j=1}^{n-1} j\tag{1.40}$$

by defining the change of index $j = i - 1$. Hence, when $i = 1$, its lower limit, $j = 0$, and when $i = n$, its upper limit, $j = (n - 1)$. The simplification in the third sum is because $j = 0$ does not contribute to the sum, so one can remove it.

1.14 Binomial Theorem

The Binomial Theorem typically covered in Algebra II provides a method for expanding $(x+y)^n$, where $x+y$ and $(x+y)^n$ are called **binomials**. This theorem enables the expansion of the polynomial $(x+y)^n$ into a sum composed of terms of the form $x^{n-k} \cdot y^k$, where $0 \leq k \leq n$ is an integer. As an example, for $n = 4$, the expansion yields: $(x+y)^4 = x^4 + 4x^3y + 6x^2y^2 + 4xy^3 + y^4$. The coefficient in front of each term $x^{n-k}y^k$ is known as the **binomial coefficient**, often denoted as $\binom{n}{k}$.

Some of you may find this video helpful: [Unveiling the Joy of Math: Inventing the Binomial Theorem](#) by Quick QEDs.

1.14.1 Binomial Coefficients

In mathematics, binomial coefficients are positive integers that appear as coefficients in the Binomial Theorem. They are denoted by $\binom{n}{k}$, where n and k are integers such that $n \geq k \geq 0$. These coefficients represent the x^k term in the expansion of $(1+x)^n$. The binomial coefficients are defined by

$$\binom{n}{k} := \frac{n!}{k!(n-k)!}, \quad (1.41)$$

where ! denotes factorial; for example, $5! = 5 \cdot 4 \cdot 3 \cdot 2 \cdot 1 = 120$. Equivalently, this formula can be written as

$$\binom{n}{k} = \frac{n \cdot (n-1) \cdots (n-k+1)}{k \cdot (k-1) \cdots 1} = \prod_{\ell=1}^k \frac{n-\ell+1}{\ell} = \prod_{\ell=0}^{k-1} \frac{n-\ell}{k-\ell}, \quad (1.42)$$

each of which is more efficient for doing the actual calculation.

Remark 1.78. By convention, $0! := 1$ and hence $\binom{n}{0} = 1$ and $\binom{n}{n} = 1$.

In the field of combinatorics, these coefficients denote the number of **distinct ways** that k elements can be selected from a set of n elements. Hence, $\binom{n}{k}$ is frequently pronounced as “ n choose k ”. As an example, if you have a group of five people ($n = 5$) and you want to select a team of two people ($k = 2$) from this group, the number of different teams you can form is given by “5 choose 2” or $\binom{5}{2}$, which equals 10 according to (1.41). It is important to note that $\binom{n}{k}$ doesn’t take into account the order of the elements chosen. If the order were important (for instance, if you were choosing a president and a vice president from a group), you would use a different formula (called a permutation) to determine the number of possible distinct combinations of a president and a vice president.

Proposition 1.79: Useful Lower Bound for Studying Euler’s Number

For $n > 2k$,

$$\binom{n}{k} := \frac{n!}{k!(n-k)!} = \frac{n \cdot (n-1) \cdots (n-k+1)}{k!} > \frac{n^k}{2^k k!}. \quad (1.43)$$

Note: This is a **lower bound** because $\frac{n^k}{2^k k!} < \binom{n}{k} := \frac{n!}{k!(n-k)!}$. This simple bound will allow us in Chapter 2.6.3 to show that exponentials of the form a^x , with $a > 1$, **grow faster than any polynomial function** as x tends to infinity.

Proof: From $n > 2k$, we have $\frac{n}{2} > k$ and $-k > -\frac{n}{2}$. Thus, $(n-k+1) > (n-\frac{n}{2}+1) > \frac{n}{2}$. Hence, $n > (n-1) > \cdots > (n-k+1) > \frac{n}{2}$. Therefore,

$$\underbrace{n \cdot (n-1) \cdots (n-k+1)}_{k\text{-times}} > \frac{n^k}{2^k} \implies \frac{n \cdot (n-1) \cdots (n-k+1)}{k!} > \frac{n^k}{2^k k!}.$$

Proposition 1.80: Useful Upper Bounds for Studying Euler’s number

For all $n \geq 1$ and $0 \leq k \leq n$,

$$\binom{n}{k} := \frac{n!}{k!(n-k)!} = \frac{n \cdot (n-1) \cdots (n-k+1)}{k!} \leq \frac{n^k}{k!}, \quad (1.44)$$

and for all $0 < k < n$, the inequality is strict, namely,

$$\binom{n}{k} := \frac{n!}{k!(n-k)!} = \frac{n \cdot (n-1) \cdots (n-k+1)}{k!} < \frac{n^k}{k!}. \quad (1.45)$$

Note: This even simpler bound will help us to approximate Euler's number.

Proof: Equation (1.45) holds because $(n-j) \leq n$ for all integers $j \geq 0$. Concerning the strict inequality, for $k=1$, the result is immediate. For $2 \leq k < n$, we have

$$\underbrace{(n-k+1) < (n-k+2) < \cdots < (n-1)}_{k \text{ terms}} < n,$$

and thus $n \cdot (n-1) \cdots (n-k+1) < n^k$. Substituting this into (1.45) proves the result. ■

Proposition 1.81: A Second Useful Upper Bound for Studying Euler's Number

For all $n \geq 1$ and $0 \leq k \leq n$,

$$\binom{n}{k} \frac{1}{n^k} \leq \binom{n+1}{k} \frac{1}{(n+1)^k}, \quad (1.46)$$

and for $n \geq 2$ and $2 \leq k \leq n$, the inequality is strict, namely,

$$\binom{n}{k} \frac{1}{n^k} < \binom{n+1}{k} \frac{1}{(n+1)^k}. \quad (1.47)$$

Note: These bounds will help us show that $(1 + \frac{1}{n})^n < (1 + \frac{1}{n+1})^{n+1}$, a property that Bernoulli noted when he discovered e .

Proof: We only show (1.47) and leave the other for the interested learner. For all $n \geq 2$ and $2 \leq k \leq n$, we have

$$\begin{aligned} \binom{n}{k} \frac{1}{n^k} &= \frac{n \cdot (n-1) \cdot (n-2) \cdots (n-k+1)}{k! n^k} \\ \binom{n+1}{k} \frac{1}{(n+1)^k} &= \frac{(n+1) \cdot (n+1-1) \cdot (n+1-2) \cdots (n+1-k+1)}{k! (n+1)^k}. \end{aligned} \quad (1.48)$$

We can cancel the common term, $k!$, in the denominators. Then, (1.47) holds if, and only if,

$$\begin{aligned} \frac{n \cdot (n-1) \cdot (n-2) \cdots (n-k+1)}{n^k} &< \frac{(n+1) \cdot (n+1-1) \cdot (n+1-2) \cdots (n+1-k+1)}{(n+1)^k} \\ \frac{n}{n} \cdot \frac{n-1}{n} \cdot \frac{n-2}{n} \cdots \frac{n-k+1}{n} &\stackrel{\uparrow}{<} \underbrace{\frac{(n+1)}{(n+1)} \cdot \frac{(n+1)-1}{(n+1)} \cdot \frac{(n+1)-2}{(n+1)} \cdots \frac{(n+1)-k+1}{(n+1)}}_{k \text{ terms}}. \end{aligned}$$

But this last inequality is true because, for all $n \geq 1$ and $m \geq 1$,

$$\frac{n-m}{n} = 1 - \frac{m}{n} < 1 - \frac{m}{(n+1)} = \frac{(n+1)-m}{(n+1)}.$$

■

1.14.2 Expansion of $(x+y)^n$

Binomial	Expansion
$(x+y)^0$	1
$(x+y)^1$	$x+y$
$(x+y)^2$	$x^2 + 2xy + y^2$
$(x+y)^3$	$x^3 + 3x^2y + 3xy^2 + y^3$
$(x+y)^4$	$x^4 + 4x^3y + 6x^2y^2 + 4xy^3 + y^4$
$(x+y)^5$	$x^5 + 5x^4y + 10x^3y^2 + 10x^2y^3 + 5xy^4 + y^5$
$(x+y)^6$	$x^6 + 6x^5y + 15x^4y^2 + 20x^3y^3 + 15x^2y^4 + 6xy^5 + y^6$
$(x+y)^7$	$x^7 + 7x^6y + 21x^5y^2 + 35x^4y^3 + 35x^3y^4 + 21x^2y^5 + 7xy^6 + y^7$
$(x+y)^8$	$x^8 + 8x^7y + 28x^6y^2 + 56x^5y^3 + 70x^4y^4 + 56x^3y^5 + 28x^2y^6 + 8xy^7 + y^8$

Table 1.82: A handy list of binomial expansions. Note the symmetry in the coefficients about the middle of each formula. This is because $\binom{n}{k} = \binom{n}{n-k}$.

Theorem 1.83: Binomial Theorem

For all integers $n \geq 1$,

$$\begin{aligned}(x+y)^n &= \sum_{k=0}^n \binom{n}{k} x^{n-k} y^k \\ &= \sum_{k=0}^n \binom{n}{k} x^k y^{n-k}.\end{aligned}\tag{1.49}$$

Early on, the “summing notation” can be a bit much to digest. We note that the above can also be written as

$$(x+y)^n = \binom{n}{0} x^n y^0 + \binom{n}{1} x^{n-1} y^1 + \binom{n}{2} x^{n-2} y^2 + \cdots + \binom{n}{n-1} x^1 y^{n-1} + \binom{n}{n} x^0 y^n.$$

Table 1.82 provides explicit expansions for $0 \leq n \leq 8$. A proof by induction is given by [Wikipedia](#), [Khan Academy](#), [Math is Fun](#), and many other sources.

Example 1.84. Use the Binomial Theorem to determine the expansion of $(1+x)^n$.

Solution:

$$\begin{aligned}(1+x)^n &= (x+1)^n = \sum_{k=0}^n \binom{n}{k} x^k \\ &= \binom{n}{0} x^0 + \binom{n}{1} x^1 + \binom{n}{2} x^2 + \binom{n}{3} x^3 + \cdots + \binom{n}{n-1} x^{n-1} + \binom{n}{n} x^n \\ &= 1 + nx + \frac{n(n-1)}{2} x^2 + \frac{n(n-1)(n-2)}{6} x^3 + \cdots + nx^{n-1} + x^n.\end{aligned}$$

1.15 Some Special Functions

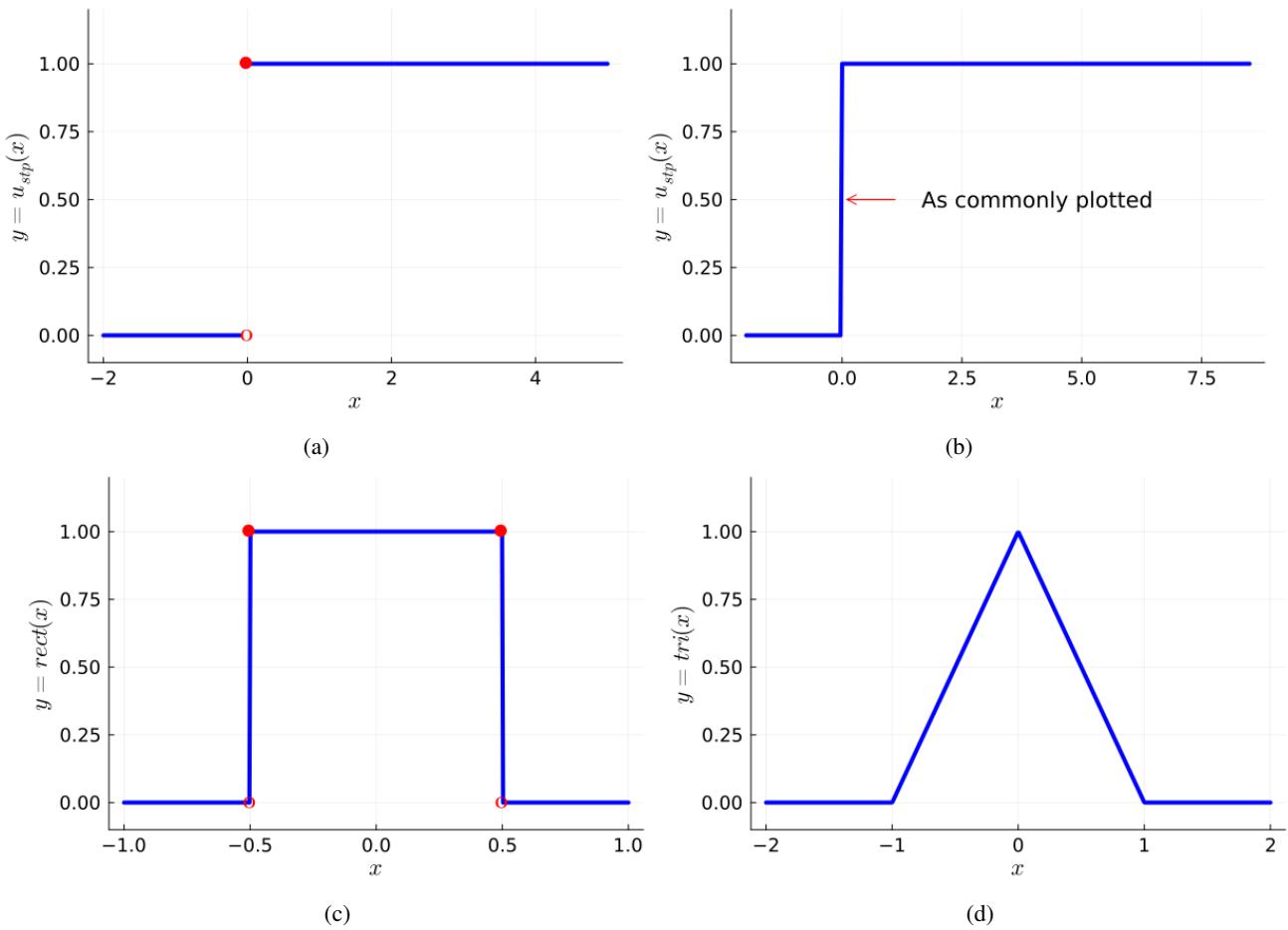


Figure 1.12: Some useful functions. (a) The unit step function. Here, the graph is done “by the book”, defining clearly a value of the function at the point of discontinuity, the origin. (b) This is how the unit step function is normally plotted (graphed). You are expected to understand that while you must choose a SINGLE value for the function at the origin, it is usually not important what value you assign at the point of discontinuity. (c) The unit rectangle function is plotted in a “hybrid manner” so that its rectangular shape is evident. Normally, one does not care what value is assigned at each of the two points of discontinuity. (d) The unit triangle function is continuous everywhere (it has no jumps). Note the base of the unit rectangle function has a width of one, while the base of the unit triangle function has a width of two. **This is done so that each function’s “shape” encloses “unit area”, meaning the area is one.**

Figure 1.12 shows three functions, the **Unit Step Function** (aka, Heaviside function, or Heaviside unit step function) in (a) and (b), the **Unit Rectangle Function** (aka, rectangular function) in (c), and the **Unit Triangle Function** (aka, triangular function, hat function, or tent function) in (d). The functions are defined as follows:

Unit Step:

$$u_{\text{stp}}(x) := \begin{cases} 1 & x \geq 0 \\ 0 & x < 0. \end{cases} \quad (1.50)$$

Unit Rectangle:

$$\text{rect}(x) := \begin{cases} 1 & |x| \leq \frac{1}{2} \\ 0 & \text{otherwise;} \end{cases} \quad (1.51)$$

the area under the function is equal to one because the base and height are both equal to one.

Unit Triangle:

$$\text{tri}(x) := \begin{cases} 1 - |x| & |x| \leq 1 \\ 0 & \text{otherwise;} \end{cases} \quad (1.52)$$

the area under the function is also equal to one because the base of the triangle is equal to two, and the height is equal to one.

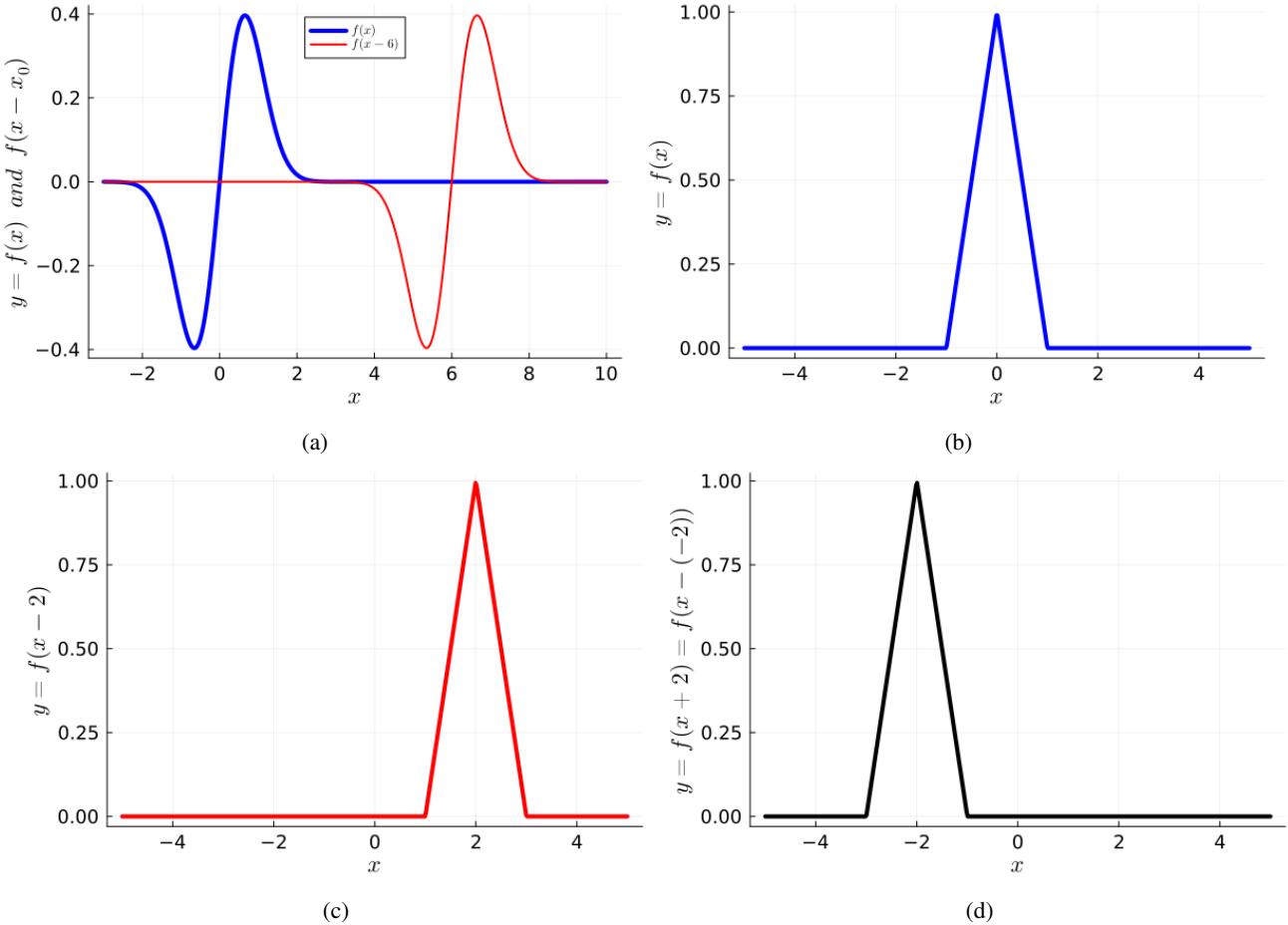


Figure 1.13: Shifted functions. (a) In blue is shown the nominal function and in red, the function shifted by $x_0 = 6$. Because $x_0 > 0$, the function's graph is shifted to the right. The shape of the function is preserved. (b) A unit triangle function. (c) The unit triangle function has been shifted to the right by $x_0 = 2 > 0$. (d) The unit triangle function has been shifted to the left by two because $x_0 = -2 < 0$. Note that the area under the triangles has not changed under shifting because their shape remains the same: their bases are the same, and their heights are the same. The same is true for rectangles.

1.16 Shifting and Scaling

Shifting and scaling are common ways to create new functions from existing functions. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be some function that you already have in hand. And let x_0 and $c \neq 0$ be two real constants.

The shift of $f(x)$ by x_0 is denoted $f(x - x_0)$. As illustrated in Fig. 1.13, the shift operation applied to a function shifts its graph to the right by x_0 , when $x_0 > 0$ and to the left by $|x_0|$, when $x_0 < 0$.

Example 1.85. Show that with the exception of $x = \pm \frac{1}{2}$, the two points of discontinuity of the rectangle function, we have

$$\text{rect}(x) = u_{\text{stp}}(x + \frac{1}{2}) - u_{\text{stp}}(x - \frac{1}{2}).$$

Solution: For $x < -\frac{1}{2}$, we have $x + \frac{1}{2} < 0$ and $x - \frac{1}{2} < 0$ and thus $u_{\text{stp}}(x + \frac{1}{2}) - u_{\text{stp}}(x - \frac{1}{2}) = 0 - 0 = 0$, which is the value of $\text{rect}(x)$. Similarly, for $x > \frac{1}{2}$, we have $x + \frac{1}{2} > 0$ and $x - \frac{1}{2} > 0$ and thus $u_{\text{stp}}(x + \frac{1}{2}) - u_{\text{stp}}(x - \frac{1}{2}) = 1 - 1 = 0$, which is also the value of $\text{rect}(x)$. Finally, we consider $-\frac{1}{2} < x < \frac{1}{2}$. Then, we have $x + \frac{1}{2} > 0$ and $x - \frac{1}{2} < 0$ and thus $u_{\text{stp}}(x + \frac{1}{2}) - u_{\text{stp}}(x - \frac{1}{2}) = 1 - 0 = 1$, which is the value of $\text{rect}(x)$.

Note: For $x = -\frac{1}{2}$, we have $u_{\text{stp}}(x + \frac{1}{2}) - u_{\text{stp}}(x - \frac{1}{2}) = u_{\text{stp}}(0) - u_{\text{stp}}(-1) = 1 - 0 = 1$, which is the value of $\text{rect}(-\frac{1}{2})$, but for $x = \frac{1}{2}$, we have $u_{\text{stp}}(x + \frac{1}{2}) - u_{\text{stp}}(x - \frac{1}{2}) = u_{\text{stp}}(1) - u_{\text{stp}}(0) = 1 - 1 = 0$, which is NOT the value of $\text{rect}(\frac{1}{2})$. ■

The scale of $f(x)$ by $c \neq 0$ is denoted $f(c \cdot x)$. As illustrated in Fig. 1.14, the scale operation applied to a function “squeezes” its graph when $c > 1$ and expands its graph, when $0 < c < 1$. We are not defining here the scaling operation for negative values of c because it also “flips” the function about the y -axis.

Example 1.86. Compute the area under the functions,

- (a) $\text{rect}(x)$
- (b) $\text{rect}(3x)$
- (c) $\text{rect}(0.2x)$
- (d) $2 \cdot \text{tri}(x)$
- (e) $2 \cdot \text{tri}(11x)$
- (f) $2 \cdot \text{tri}(\frac{x}{3})$

Solutions:

- (a) $\text{rect}(x)$ has a base of one and a height of one, and thus its area is also one, i.e., 1.0.
- (b) $\text{rect}(3x)$ has a base of one-third and a height of one, and thus its area is one-third, i.e., $\frac{1}{3}$.
- (c) $\text{rect}(0.2x)$ has a base of five and a height of one, and thus its area is five, i.e., 5.0.
- (d) $2 \cdot \text{tri}(x)$ has a base of two and a height of two, and thus its area is two (one half base times height), i.e., 2.0.
- (e) $2 \cdot \text{tri}(11x)$ has a base of two-elevenths and a height of two, and thus its area is $\frac{2}{11}$.
- (f) $2 \cdot \text{tri}(\frac{x}{3})$ has a base of six and a height of two, and thus its area is six, i.e., 6.0.

1.17 A Peek at Feedback Control and a Partial Roadmap to this Textbook

Why feedback control and why in a Chapter on Pre-Calculus:

- Feedback control is a pervasive technology. In Michigan Engineering, feedback control courses can be found in approximately half of its departments. Its widespread representation rivals Physics and Mathematics.
- The robot arms in Fig. 1.6 are not worth much if you cannot command them to reach an object at a given position. In many cases, that comes down to regulating each of the joints in the robot arm to a particular value, called a **setpoint**. Commanding robot joints to setpoints is the bread-and-butter of feedback control. Because some of the key ideas only involve basic algebra, pre-calculus knowledge is enough to glean key insight into how a feedback controller works.
- The description of a feedback control system will also allow us to highlight where Calculus is used in Robotics, partially giving you a roadmap to this textbook.

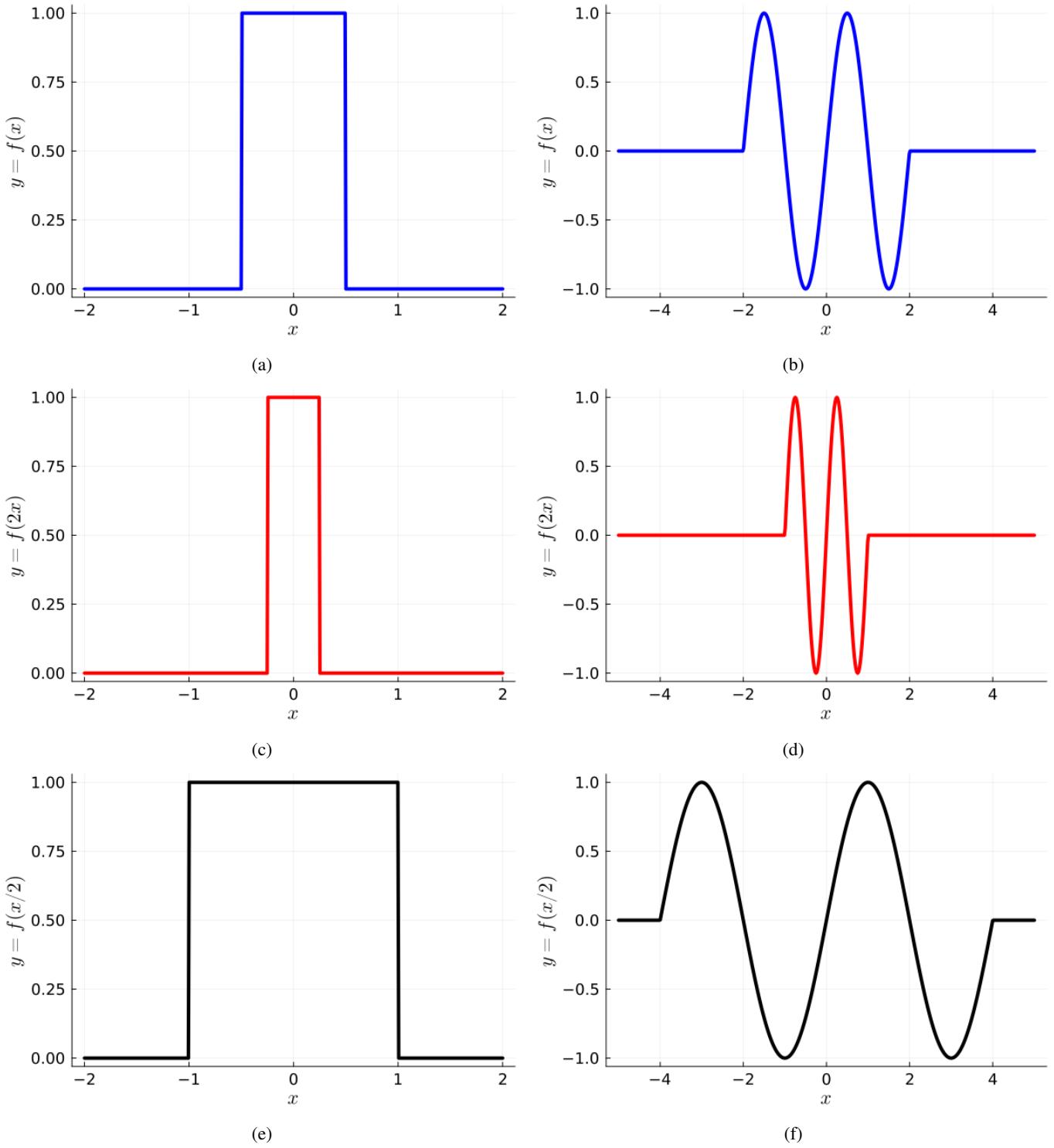


Figure 1.14: Scaled functions. (a) Unit rectangle function in blue. (c) Unit rectangle function scaled by $c = 2$ in red. The area of the rectangle has been decreased by a factor of two. (e) Unit rectangle function scaled by $c = \frac{1}{2}$ in black. The area of the rectangle has been increased by a factor of two. (b) Two periods of the function $\sin(x)$ in blue. (d) Two periods of the function $\sin(2x)$ in red. The frequency has been doubled. (f) Two periods of the function $\sin(\frac{x}{2})$. The frequency has been halved.

Consider Fig. 1.15-(b). The diagram represents the algebraic equation

$$y = P \cdot u + d, \quad (1.53)$$

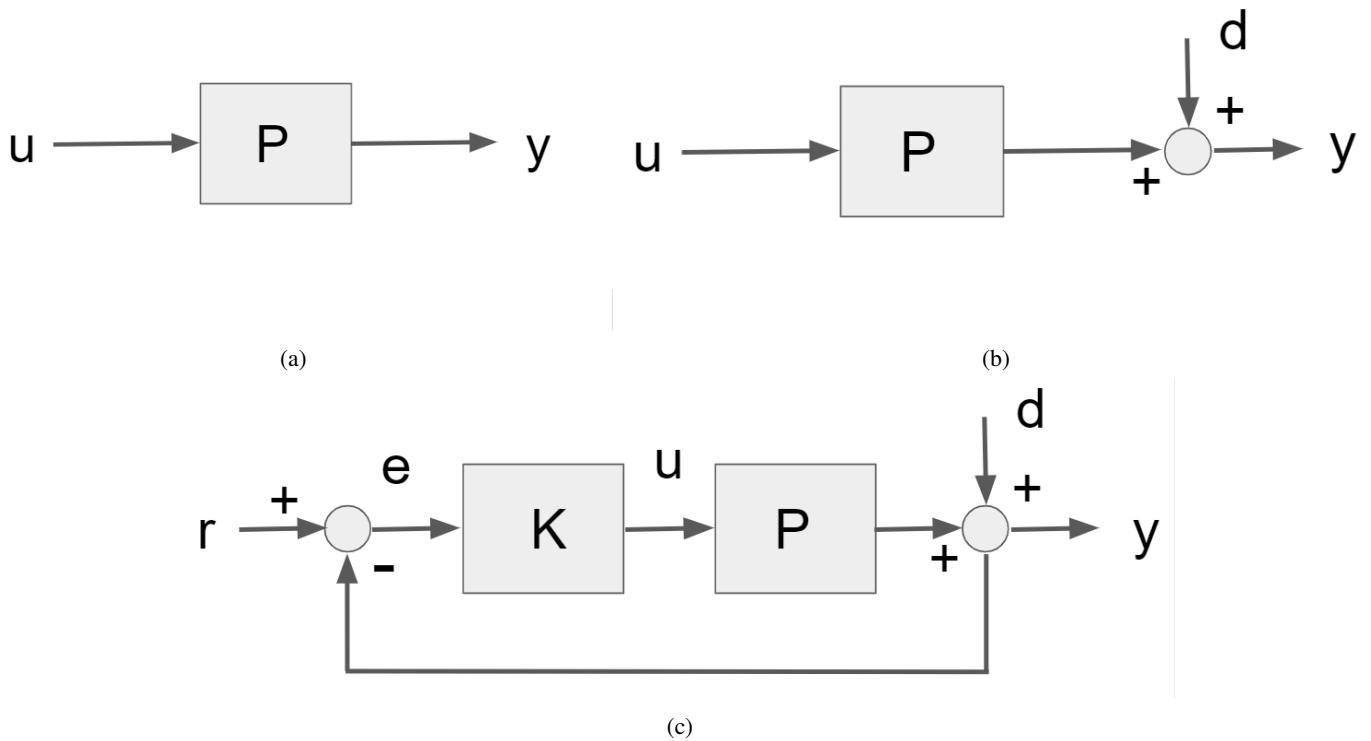


Figure 1.15: Feedback diagrams, where P stands for plant, which is the generic name for the system that is to be controlled (e.g., motor with a link connected to it), y is the measured output (e.g. joint angle), u is the input to the plant (e.g., motor torque), and d is a disturbance signal, also called a perturbation (e.g., friction in the joint). (a) An open-loop system without a disturbance. (b) An open-loop system with a disturbance. (c) A **closed-loop system**, where K is the **controller gain**, r is the reference signal (what we want y to become), and e is the error term, the difference between the reference signal and the measured output. In (a) and (b), there is no loop, hence the term “open-loop system”, while in (c), a loop is present, hence the term “closed-loop control”. **Credit for this idea: Prof. Semyon Meerkov, University of Michigan, ECE Department.**

where

- $P > 0$ is a positive constant representing the **gain** (aka, multiplying factor) of the system to be controlled or regulated; let's imagine that our plant is one of the links in Fig. 1.5 attached to a DC motor. Our objective is to regulate the joint angle to a desired value, which we will denote by r because, in the control community, we call it a **reference signal**.
- u is the **input or control signal** to the plant, a signal we can manipulate to achieve our goal; let's imagine that u is the torque of the DC motor.
- y is the **measured output** of the system; let's imagine it is the output of the encoder at the joint angle.
- d is our nemesis, called a **disturbance signal** or **perturbation** in control parlance, which is universally present in real life; let's imagine that d is friction in the joint of the robot. In industrial robots, unfortunately, d can be very large. Yes, manufacturers try to minimize it, but alas, that is hard and expensive.

How to Read the Diagrams in Fig. 1.15: These diagrams represent linear equations. In (a), $y = P \cdot u$, whatever signal comes into the box is multiplied by the system gain P . In (b), the circle with the plus signs is called a **summing junction** aka **summer**; the output of the summing junction is equal to the sum of all the signals that enter it, with each input signal weighted by the sign nearest to it. Because in (b), they are all plus signs, $y = P \cdot u + d$. In (c), we have our first encounter with a feedback system, aka a closed-loop system. The first summer represents $e = r - y$, and the second summer represents $y = P \cdot u + d$, as in (b). These two equations are connected by $u = K e$. We will solve these simultaneous equations shortly, but not now!

Objective: Regulate y to r , a desired angle, despite us not being able to measure the disturbance, d .

Open-loop Approach to the Problem

We know the plant model is $y = P \cdot u + d$, and we cannot measure d . So we propose

$$u = K \cdot r \quad (1.54)$$

and seek to select K so that $y \approx r$, its desired value. We write this as

$$r \approx y = P \cdot u + d = P \cdot K \cdot r + d. \quad (1.55)$$

Because we know nothing about d , we assume it is as likely to be positive as negative, with a mean value of zero. This leads to

$$r \approx y \iff P \cdot K = 1 \iff K = \frac{1}{P}, \quad (1.56)$$

that is, the control gain is set to the reciprocal of the plant gain. Because $P \cdot K = 1$, our open-loop control system results in

$$y = r + d. \quad (1.57)$$

There are two issues with this solution:

- (a) the disturbance passes directly through to the output, causing error in the precision of the control system; and
- (b) we may not know the plant gain exactly, and hence our control gain is really $K = \frac{1}{P_{\text{nom}}}$, a nominal value for the plant gain, and the real output is then

$$y = \left(\frac{P}{P_{\text{nom}}} \right) \cdot r + d.$$

There must be a better way to handle the uncertainty in the model, namely, the unmeasurable disturbance and any variations in the plant gain, P . Right?

Closed-loop Approach to the Problem

The feedback diagram in Fig. 1.15-(c) yields

$$\begin{aligned} y &= P \cdot u + d \\ e &= r - y \\ u &= K \cdot e, \end{aligned} \quad (1.58)$$

where K is a feedback gain. We get to adjust it so as to best meet our objective of making the output as close to the commanded reference value as possible, that is, $y \approx r$. Combining the bottom two equations of (1.58), we have

$$u = K \cdot e = K \cdot (r - y).$$

Plugging this into the top equation yields

$$y = P \cdot K \cdot e = P \cdot K \cdot (r - y) + d,$$

a scalar linear equation. Solving for y gives us

$$(1 + P \cdot K) \cdot y = P \cdot K \cdot r + d,$$

and thus,

$$y = \frac{P \cdot K}{1 + P \cdot K} \cdot r + \frac{1}{1 + P \cdot K} \cdot d, \quad (1.59)$$

a linear equation relating the reference signal, r , and the disturbance signal, d , to the output, y . **This is the most important equation in feedback control** because it tells us everything about how the free parameter K in the controller relates to the multiplying factor in front of the reference command, namely, $\frac{P \cdot K}{1 + P \cdot K}$, and in front of the disturbance, namely $\frac{1}{1 + P \cdot K}$.

In a “perfect world”, we would have $\frac{P \cdot K}{1 + P \cdot K} = 1$ and $\frac{1}{1 + P \cdot K} = 0$, so that $y = r$. Perfection is hard to achieve, so we seek to make the following approximations

$$\frac{P \cdot K}{1 + P \cdot K} \approx 1 \text{ and } \frac{1}{1 + P \cdot K} \approx 0$$

hold as closely as we can.

- Suppose that P has a nominal value of 2, and we select $K = 25$. Then, plugging in values gives us

$$y = \frac{50}{51} \cdot r + \frac{1}{51} \cdot d \approx 0.98r + 0.0196d.$$

The disturbance is greatly attenuated, and y is within 2% of its desired value, r . Not bad, huh?

- Suppose that P could vary between $1 \leq P \leq 3$, and we keep $K = 25$. Then we have

$$\begin{aligned} y &\approx 0.96r + 0.0385d \quad P = 1 \\ y &\approx 0.99r + 0.0136d \quad P = 3, \end{aligned}$$

which is also great!

Note: Suppose you could select $K = 100$. How would the above numbers change?

Taking Stock

Closed-loop control is a game changer. It allows one to “track” reference commands, that is, drive y to its desired value r , even in the face of

- uncertainty in the mathematical model of the system being controlled, and
- nefarious disturbances trying to prevent us from achieving our tracking objective.

Where's the Calculus?

Nowhere to be seen. Yet! In the above, we've assumed the system to be controlled can be represented by a constant gain, P . You know full well that “ $F = ma$ ” has to enter into the picture when dealing with a robot. That is called **dynamics**, and that's where Calculus first comes in.

- In Chapter 3.4, we will draw on your knowledge of Physics to compute the displacement of an object from its speed of motion, using definite integration. Further applications of integration allow us to compute the mass of a robot link from its density and shape, and its “center of gravity”, aka “center of mass”, the point at which a single point mass can summarize the overall effect of gravity on the object.
- In Chapter 5, we'll study derivatives and realize that velocity is the derivative of position with respect to time. In a similar manner, acceleration is the derivative of velocity with respect to time.
- Deriving dynamical models of robots via $F = ma$ is very tedious. Lagrange invented another method for creating dynamical models of collections of rigid bodies, such as our 3-link robot. It involves potential energy and kinetic energy. Chapter 6.4 explores Lagrange's Equations of Motion.
- In Chapter 9, we will learn about ordinary differential equations (ODEs for short) and how to solve them. We'll also learn how to develop linear approximations of complicated nonlinear differential equations, opening up their analysis to undergraduate students.
- In Chapter 10, we study the Laplace transform, a means to understand the stability of ODEs. The Laplace Transform is a kind of superpower: it transforms linear ODEs into algebraic equations! Yes, really. By looking at the roots of the resulting algebraic equations, we can determine if small perturbations introduced into an ODE model of a robot eventually dissipate (this is what stability is all about) or cause the system to explode (this is the essence of instability).
- In Chapter 10, we'll use all of this knowledge to design feedback control laws for linear ODEs, culminating in us stabilizing the **BallBot**. If you really like this material, at Michigan, you would follow up with one of EECS 460 Control Systems Analysis and Design, ME 461 Automatic Control, or AEROSP 345 Flight Dynamics and Control. And then there are the **graduate-level courses!** Their cardinality is a good approximation of infinity, which may sound intimidating, but it really means that in any given semester, you have a smorgasbord of interesting courses to consider.

1.18 (Optional Read:) Binomial Theorem meets Euler

On the surface, the Binomial Theorem, which expands $(x + y)^n$, seems to be limited to helping us with polynomials. However, Euler's number involves $\left(1 + \frac{1}{n}\right)^n$ and Euler's formula allows us to write $\cos^n(x) = \left(\frac{e^{ix} + e^{-ix}}{2}\right)^n$ and $\sin^n(x) = \left(\frac{e^{ix} - e^{-ix}}{2i}\right)^n$. Here, we apply the Binomial Theorem to understand these quantities.

1.18.1 Binomial Theorem meets Euler's Constant

Applying the Binomial Theorem to the expression $\left(1 + \frac{1}{n}\right)^n$ yields another means to compute Euler's number, e . Indeed,

$$\left(1 + \frac{1}{n}\right)^n = \sum_{k=0}^n \binom{n}{k} \frac{1}{n^k}. \quad (1.60)$$

The k th term of this sum is

$$\begin{aligned} \binom{n}{k} \frac{1}{n^k} &= \frac{n!}{k!(n-k)!} \cdot \frac{1}{n^k} \\ &= \frac{1}{k!} \cdot \frac{n \cdot (n-1) \cdot (n-2) \cdots (n-k+1)}{n^k} \\ &= \frac{1}{k!} \cdot \frac{n}{n} \cdot \frac{n-1}{n} \cdot \frac{n-2}{n} \cdots \frac{n-k+1}{n} \\ &= \frac{1}{k!} \cdot \left(1 - \frac{1}{n}\right) \cdot \left(1 - \frac{2}{n}\right) \cdots \left(1 - \frac{k+1}{n}\right), \end{aligned} \quad (1.61)$$

where we noted that $\frac{n}{n} = 1$, $\frac{n-1}{n} = 1 - \frac{1}{n}$, $\frac{n-2}{n} = 1 - \frac{2}{n}$, all the way to $\frac{n-k+1}{n} = 1 - \frac{k+1}{n}$.

We will now make some approximations in the “spirit of Calculus” that will become more rigorous in the next Chapter; you can also see [Wikipedia](#).

For k fixed and “ n much larger than k ”, $1 - \frac{1}{n} \approx 1$, $1 - \frac{2}{n} \approx 1$, all the way to $1 - \frac{k+1}{n} \approx 1$. If you are with us so far, then you won't mind the next approximation, namely, the product of terms that are close to one is also close to one,

$$\left(1 - \frac{1}{n}\right) \cdot \left(1 - \frac{2}{n}\right) \cdots \left(1 - \frac{k+1}{n}\right) \approx 1.$$

Therefore, for k constant and “ n large”,

$$\binom{n}{k} \frac{1}{n^k} \approx \frac{1}{k!}.$$

This indicates that e can be approximated as

$$e \approx \left(1 + \frac{1}{n}\right)^n = \sum_{k=0}^n \binom{n}{k} \frac{1}{n^k} \approx \sum_{k=0}^n \frac{1}{k!} = \frac{1}{0!} + \frac{1}{1!} + \frac{1}{2!} + \frac{1}{3!} + \cdots + \frac{1}{n!}, \quad (1.62)$$

with the approximation becoming better and better as n gets larger and larger. The notion of approximations improving for large n will lead us to the concept of limits in Chapter 2.

Now, after all these approximations, what do we have? Table 1.87 compares our two approximations for Euler's number. We see that the sum, $\sum_{k=0}^n \frac{1}{k!}$, “converges” to e much faster than the product, $\left(1 + \frac{1}{n}\right)^n$. Euler's name is attached to the constant e because he proved rigorously¹⁵, among other things, that as n tends to infinity, the sum converges to e .

¹⁵Here, “rigorously” means that Euler filled in all the holes in the argument we sketched out. In mathematics, rigor refers to a level of thoroughness and consistency in reasoning that ensures the validity and correctness of results. It involves precise definitions, clear assumptions, and sound logical techniques to establish the truth of mathematical statements. Rigor helps in avoiding errors, clarifying concepts, and ensuring that the conclusions reached are reliable and trustworthy. Through rigorous processes, mathematicians aim to build a solid foundation for the development of mathematical theories and the advancement of the field. Engineering is built upon these solid foundations.

n	$(1 + \frac{1}{n})^n$	$\sum_{k=0}^n \frac{1}{k!}$
5	2.48832000	2.71666667
10	2.59374246	2.71828180
10^2	2.70481383	2.71828183
10^3	2.71692393	2.71828183

Table 1.87: Comparison of the values of $e \approx 2.71828183 \dots$, the product $(1 + \frac{1}{n})^n$, and the sum $\sum_{k=0}^n \frac{1}{k!}$ for various n . Values are shown to 8 significant digits. The sum seems to converge much faster, meaning that it provides a more practical way of computing e .

```

1 n=100
2
3 e_approxA = (1+1/n)^n
4
5 e_approxB = 0.0
6 for k = 0:n
7     e_approxB = e_approxB + 1/factorial(big(k))
8     # big(k) needed for 21! and beyond
9 end
10 [e_approxA; e_approxB; exp(1)]

```

Output

```

3-element Array{BigFloat,1}:
2.704813829421528481589120929129421710968017578125
2.718281828459045235360287471352662497757247093699959574966967627724076630353416
2.718281828459045090795598298427648842334747314453125

```

1.18.2 Binomial Theorem and Powers of Sine and Cosine

This section shows how powers of sine and cosine can be simplified through the Binomial Theorem. We first use Euler's formula to write

$$\cos^n(\theta) = \left(\frac{e^{i\theta} + e^{-i\theta}}{2} \right)^n = \frac{1}{2^n} \cdot (e^{i\theta} + e^{-i\theta})^n = \frac{1}{2^n} \cdot (x + x^{-1})^n, \quad (1.63)$$

for $x = e^{i\theta}$. The binomial expansion of $(x + x^{-1})^n$ for any integer n is given by the Binomial Theorem:

$$(x + x^{-1})^n = \sum_{k=0}^n \binom{n}{k} x^k \cdot (x^{-1})^{n-k} = \sum_{k=0}^n \binom{n}{k} x^k \cdot x^{k-n} = \sum_{k=0}^n \binom{n}{k} x^{2k-n}. \quad (1.64)$$

Because the sum starts at zero and ends at n , there are $(n + 1)$ terms in the expansion. When n is odd, $(n + 1)$ is even, and we can pair the first term with the last term, the second term with the second-to-last term, and in general, the k -th term from the start with the $(n - k)$ -th term from the end to obtain,

$$\binom{n}{k} x^{2k-n} + \binom{n}{n-k} x^{n-2k}. \quad (1.65)$$

Since $\binom{n}{k} = \binom{n}{n-k}$, the expansion simplifies to

$$(x + x^{-1})^n = \sum_{k=0}^{\frac{n-1}{2}} \binom{n}{k} (x^{n-2k} + x^{-(n-2k)}). \quad (1.66)$$

Looping back to (1.63) and substituting in $x = e^{i\theta}$, for n odd, we have

$$\begin{aligned}\cos^n(\theta) &= \frac{1}{2^n} \cdot \sum_{k=0}^{\frac{n-1}{2}} \binom{n}{k} (e^{i(n-2k)\theta} + e^{-i(n-2k)\theta}) \\ &= \frac{1}{2^{n-1}} \cdot \sum_{k=0}^{\frac{n-1}{2}} \binom{n}{k} \left(\frac{e^{i(n-2k)\theta} + e^{-i(n-2k)\theta}}{2} \right) \\ &= \frac{1}{2^{n-1}} \cdot \sum_{k=0}^{\frac{n-1}{2}} \binom{n}{k} \cos((n-2k)\theta).\end{aligned}$$

Similar steps yield the result for n even. The main difference is, when n is even, $(n+1)$ is odd. We can still do the same pairing up of terms, except the middle term remains unpaired, giving

$$(x + x^{-1})^n = \sum_{k=0}^{\frac{n-2}{2}} \binom{n}{k} (x^{n-2k} + x^{-(n-2k)}) + \binom{n}{\frac{n}{2}}. \quad (1.67)$$

The rest is left to the learner. The result for $\sin^n(x)$ is developed in Chapter 1.19.

Proposition 1.88: Reducing Powers of Cosine and Sine

For all integers $n > 1$,

$$\cos^n(\theta) = \begin{cases} \frac{1}{2^{n-1}} \cdot \sum_{k=0}^{\frac{n-1}{2}} \binom{n}{k} \cos((n-2k)\theta) & n \text{ odd} \\ \frac{1}{2^n} \binom{n}{\frac{n}{2}} + \frac{1}{2^{n-1}} \cdot \sum_{k=0}^{\frac{n-2}{2}} \binom{n}{k} \cos((n-2k)\theta) & n \text{ even} \end{cases} \quad (1.68)$$

and

$$\sin^n(\theta) = \begin{cases} \frac{(-1)^{\frac{n-1}{2}}}{2^{n-1}} \cdot \sum_{k=0}^{\frac{n-1}{2}} \binom{n}{k} (-1)^k \sin((n-2k)\theta) & n \text{ odd} \\ \frac{1}{2^n} \binom{n}{\frac{n}{2}} + \frac{(-1)^{\frac{n}{2}}}{2^{n-1}} \cdot \sum_{k=0}^{\frac{n-2}{2}} \binom{n}{k} (-1)^k \cos((n-2k)\theta) & n \text{ even} \end{cases} \quad (1.69)$$

Note: It's not a typo that for n even, $\sin^n(\theta)$ expands in terms of $\cos((n-2k)\theta)$. If you are curious as to how that happens, all of the steps are given in Chapter 1.19 or you can consult [Wolfram World](#), in addition to many other sources on the web. While the formulas may appear complicated, they are easy to code up in Julia, and Voila !, you have all of your power rules at your fingertips. Powers two through seven are given in Table 1.89.

```

1 using SpecialFunctions: binomial
2
3 function cos_power_n(theta, n)
4     if n % 2 == 0 # n is even
5         sum = binomial(n, n ÷ 2) / 2^n
6         for k in 0:(n ÷ 2 - 1)
7             sum += binomial(n, k) * cos((n - 2k) * theta) / 2^(n-1)
8         end
9         return sum
10    else # n is odd
11        sum = 0.0
12        for k in 0:(n - 1) ÷ 2)

```

n	$\cos^n(\theta)$	$\sin^n(\theta)$
2	$\frac{1}{2} + \frac{1}{2} \cos(2\theta)$	$\frac{1}{2} - \frac{1}{2} \cos(2\theta)$
3	$\frac{3}{4} \cos(\theta) + \frac{1}{4} \cos(3\theta)$	$\frac{3}{4} \sin(\theta) - \frac{1}{4} \sin(3\theta)$
4	$\frac{3}{8} + \frac{1}{2} \cos(2\theta) + \frac{1}{8} \cos(4\theta)$	$\frac{3}{8} - \frac{1}{2} \cos(2\theta) + \frac{1}{8} \cos(4\theta)$
5	$\frac{5}{8} \cos(\theta) + \frac{5}{16} \cos(3\theta) + \frac{1}{16} \cos(5\theta)$	$\frac{5}{8} \sin(\theta) - \frac{5}{16} \sin(3\theta) + \frac{1}{16} \sin(5\theta)$
6	$\frac{5}{16} + \frac{15}{32} \cos(2\theta) + \frac{3}{16} \cos(4\theta) + \frac{1}{32} \cos(6\theta)$	$\frac{5}{16} - \frac{15}{32} \cos(2\theta) + \frac{3}{16} \cos(4\theta) - \frac{1}{32} \cos(6\theta)$
7	$\frac{35}{64} \cos(\theta) + \frac{21}{64} \cos(3\theta) + \frac{7}{64} \cos(5\theta) + \frac{1}{64} \cos(7\theta)$	$\frac{35}{64} \sin(\theta) - \frac{21}{64} \sin(3\theta) + \frac{7}{64} \sin(5\theta) - \frac{1}{64} \sin(7\theta)$

Table 1.89: Powers of Sine and Cosine

```

13     sum += binomial(n, k) * cos((n - 2k) * theta) / 2^(n-1)
14 end
15 return sum
16 end
17 end
18
19 function sin_power_n(theta, n)
20   if n % 2 == 0 # n is even
21     sum = binomial(n, n / 2) / 2^n
22     for k in 0:(n - 2) / 2
23       sum += (-1)^(n / 2 + k) * binomial(n, k) * cos((n - 2k) * theta) / 2^(n-1)
24     end
25   return sum
26 else # n is odd
27   sum = 0.0
28   for k in 0:(n - 1) / 2
29     sum += (-1)^((n - 1) / 2 + k) * binomial(n, k) * sin((n - 2k) * theta) / 2^(n-1)
30   end
31   return sum
32 end
33 end

```

Output

sin_power_n (generic function with 1 method)

```

1 using Symbolics
2 @variables theta
3 n=7
4 myAns = cos_power_n(theta, n); display(myAns)
5 myAns = sin_power_n(theta, n); display(myAns)

```

Output

$$\frac{21}{64} \cos(3\theta) + \frac{1}{64} \cos(7\theta) + \frac{35}{64} \cos(\theta) + \frac{7}{64} \cos(5\theta)$$

$$\frac{35}{64} \sin(\theta) + \frac{7}{64} \sin(5\theta) - \frac{21}{64} \sin(3\theta) - \frac{1}{64} \sin(7\theta)$$

1.19 (Optional Read:) Proofs Associated with the Chapter

For those who want to know why some of the things we covered are actually true, we provide the proofs. While an instructor may find them useful for review, students may want to study the proofs for a reason highlighted in [Why Everyone should learn Math in school - Neil deGrasse](#) (video thanks to LearnwithJaspal). Despite Dr. deGrasse's endorsement, we do not anticipate that most students of the course will want to read them. They truly are optional readings.

Proposition 1.90: Are There Any Irrational Numbers?

$\sqrt{2}$ is an irrational number.

Proof: Here is the 10,000 feet view of a proof by contradiction. We first define a logical statement such as

$$P : \text{"}\sqrt{2}\text{ is an irrational number".}$$

We then assume its logical negative, namely,

$$\neg P : \text{"}\sqrt{2}\text{ is a rational number".}$$

We seek to show that assuming $\neg P$ is true leads to an “absurdity”; specifically, we seek to deduce from $\neg P$ a second logical statement R that is both true and false! In logic, statements that are both true and false are called **contradictions**. Moreover, *the only way to arrive at a contradiction by properly applying the rules of logic¹⁶ is to start with a false premise*. If a premise is false, then its negative is true, and vice versa.

The Logic of Euclid's Proof by Contradiction

One assumes $\neg P$: “ $\sqrt{2}$ is a rational number” is **T** and deduces from that a contradiction. The only way out of the contradiction is for $\neg P$ to be **F**. Therefore, P : “ $\sqrt{2}$ is an irrational number” is **T**.

Note: Not bad for 2,000 years ago!

For us, the logical statement R will come from us starting with pair of integers “that do not have any common factors,” and yet, if their ratio is equal to $\sqrt{2}$, we'll prove that “they both must be even,” and hence have two as a common factor, yielding a contradiction! We'll then know that our premise was false, and hence its negation is true.

In the proof, we'll need the following fact about an integer m : **Fact:** m^2 is even if, and only if, m is even. In other words, if you square an integer and obtain an even number, you cannot have started with an odd¹⁷ integer! Recall that by definition, an integer m is even if there exists an integer k such that $m = 2k$. Similarly, m is odd if there exists an integer k such that $m = 2k + 1$.

If $\sqrt{2}$ is rational, then by the definition of the rational numbers, there must exist natural numbers m and n such that

- m and n have no common factors,
- $n \neq 0$, and

$$\sqrt{2} = \frac{m}{n}. \tag{1.70}$$

So far, all we have done is apply the definition of a rational number. Next, we square both sides of (1.70) to arrive at

$$\left(2 = \frac{m^2}{n^2}\right) \implies (2n^2 = m^2) \implies (m^2 \text{ is even}),$$

because it is the product of two and the natural number n^2 . From our fact about squares of integers, we deduce that m must be even, and hence there must exist an integer k such that $m = 2k$.

¹⁶It's also easy to produce something that seems to be a contradiction by messing up the rules of logic, but you already guessed that to be true!

¹⁷A counting number m is **odd** if there exists $k \in \mathbb{N}$ such that $m = (2k + 1)$. We square the number and obtain $m^2 = (2k + 1)^2 = 2(2k^2 + k) + 1$, which is an odd integer because $(2k^2 + k)$ is an integer. Hence, if you square a number and obtain an even number, you cannot have started with an odd number. You must have started with an even number.

From $2n^2 = m^2$, we deduce that

$$(2n^2 = (2k)^2) \implies (2n^2 = 4k^2) \implies (n^2 = 2k^2) \implies n^2 \text{ is even.}$$

Once again appealing to our fact about squares of integers, we deduce that n must be even, and hence there must exist an integer j such that $n = 2j$.

Because both m and n are even, they have 2 as a common factor, which contradicts that m and n have no common factors.

Because we arrived at this contradiction from the statement (premise) “ $\sqrt{2}$ is rational”, we deduce that our premise “ $\sqrt{2}$ is rational” must be false. Hence, $\sqrt{2}$ is NOT RATIONAL, which means that $\sqrt{2}$ IS IRRATIONAL. Ta-da! ■

Remark 1.91. Estermann offers an alternative proof that is pretty cool. Here is its presentation by Michael Penn.

Proposition 1.92: Bounds for Studying Euler's number

For $n \geq 1$, define $a_n := \left(1 + \frac{1}{n}\right)^n$. Then for all $n \geq 1$ and $m \geq 0$,

- (a) $a_n \leq a_{n+1}$ (the terms are monotonically increasing),
- (b) $a_n \leq \sum_{k=0}^n \frac{1}{k!}$, (we have an upper bound on the n -th term) and, for all $m \geq 1$,
- (c) $a_{n+m} < \sum_{k=0}^n \frac{1}{k!} + \frac{1}{n \cdot n!}$ (we have an upper bound for all terms).
- (d) $a_m < 3$ (by taking $n = 3$, the above property gives $a_{3+m} < 3$ for all m ; checking individually a_1 , a_2 , and a_3 gives the general result).

Proof: From the Binomial Theorem, for any $n \geq 1$,

$$a_n := \left(1 + \frac{1}{n}\right)^n = \sum_{k=0}^n \binom{n}{k} \left(\frac{1}{n}\right)^k = \sum_{k=0}^n \binom{n}{k} \cdot \frac{1}{n^k}. \quad (1.71)$$

Everything will follow from this equation.

(a) We consider the binomial expansions of a_n and a_{n+1}

$$\begin{aligned} a_n &= \sum_{k=0}^n \binom{n}{k} \cdot \frac{1}{n^k} \\ a_{n+1} &= \sum_{k=0}^n \binom{n+1}{k} \cdot \frac{1}{(n+1)^k} + \frac{1}{(n+1)^{n+1}}. \end{aligned}$$

From Prop. 1.81, for all $0 \leq k \leq n$,

$$\binom{n}{k} \cdot \frac{1}{n^k} \leq \binom{n+1}{k} \cdot \frac{1}{(n+1)^k},$$

and thus $a_{n+1} - a_n \geq \frac{1}{(n+1)^{n+1}}$.

(b) For the upper bound, we can use Prop. 1.80 stating that the binomial coefficient $\binom{n}{k}$ is upper bounded by $\frac{n^k}{k!}$. This gives us

$$\left(1 + \frac{1}{n}\right)^n = \sum_{k=0}^n \binom{n}{k} \cdot \frac{1}{n^k} \leq \sum_{k=0}^n \frac{n^k}{k!} \cdot \frac{1}{n^k} = \sum_{k=0}^n \frac{1}{k!}, \quad (1.72)$$

proving (b).

(c) Continuing from (b), for all $m \geq 1$,

$$\left(1 + \frac{1}{n+m}\right)^{n+m} \leq \sum_{k=0}^n \frac{1}{k!} + \sum_{j=1}^m \frac{1}{(n+j)!}. \quad (1.73)$$

We note that

$$(n+j)! = \underbrace{(n+j) \cdot (n+j-1) \cdots (n+1)}_{\geq (n+1)^j} \cdot n! \geq (n+1)^j \cdot n!$$

Hence,

$$\left(1 + \frac{1}{n+m}\right)^{n+m} \leq \sum_{k=0}^n \frac{1}{k!} + \frac{1}{n!} \sum_{j=1}^m \frac{1}{(n+1)^j}. \quad (1.74)$$

The term “ $\sum_{j=1}^m \frac{1}{(n+1)^j}$ ” is nearly a geometric series, and we’ll show later that for all $m \geq 1$, $\sum_{j=1}^m \frac{1}{(n+1)^j} < \frac{1}{n}$. Hence, for all $n \geq 1$ and $m \geq 1$,

$$\left(1 + \frac{1}{n+m}\right)^{n+m} < \sum_{k=0}^n \frac{1}{k!} + \frac{1}{n \cdot n!}. \quad (1.75)$$

■

Remark 1.93. Because the right side of (1.75) does not depend on m , we have that for all $n \geq 1$, Euler’s number is sandwiched between

$$\left(1 + \frac{1}{n}\right)^n \leq e \leq \sum_{k=0}^n \frac{1}{k!} + \frac{1}{n \cdot n!},$$

giving easily computable upper and lower bounds. Moreover, the upper and lower bounds approach one another as n becomes larger and larger. We do not have the tools at this point to prove this last statement, but you can verify it numerically. Later, we will be able to make sense of and write

$$e := \lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n, \text{ and,}$$

$$e = \sum_{k=0}^{\infty} \frac{1}{k!} := \lim_{n \rightarrow \infty} \sum_{k=0}^n \frac{1}{k!}. \quad (1.76)$$

Proof of Prop. 1.49 (Power Towers) that Avoids the use of Logarithms

This kind of proof is called “First Principles” because we only use the basic properties of integer powers and integer roots that have been introduced in the Chapter.

Let $x > 0$, and let $y \approx p/q$ and $z \approx r/s$ be rational approximations, with the denominators q and s positive integers. We seek to show that

$$(x^{p/q})^{r/s} = x^{\frac{pr}{qs}},$$

in other words, the exponents multiply. Indeed, by applying the definitions

$$\begin{aligned} (x^{p/q})^{r/s} &= \sqrt[s]{(x^{p/q})^r} && \text{r-th power and s-th root of } x^{p/q} \\ &= \sqrt[s]{\underbrace{(x^{p/q}) \cdots (x^{p/q})}_{r \text{ times}}} && \text{work out what it means to be the r-th power} \\ &= \sqrt[s]{\underbrace{(\sqrt[q]{x^p}) \cdots (\sqrt[q]{x^p})}_{r \text{ times}}} && \text{substitute in now for } x^{p/q} \text{ and try to group terms.} \end{aligned} \quad (1.77)$$

By Prop. 1.37, a product of roots equals the root of the product. Hence,

$$\underbrace{(\sqrt[q]{x^p}) \cdots (\sqrt[q]{x^p})}_{r \text{ times}} = \sqrt[r]{x^{rp}}.$$

The above result plus Prop. 1.39 on nested roots yields,

$$\begin{aligned} (x^{p/q})^{r/s} &= \sqrt[s]{\sqrt[q]{x^{rp}}} \\ &= \sqrt[s \cdot q]{x^{rp}} \\ &= x^{\frac{rp}{sq}}, \end{aligned} \quad (1.78)$$

which is amazing! The last step is to use Math 451 to conclude that because (1.78) holds for all rational approximations of the real numbers y and z , it must also hold for y and z themselves. ■

Proof of Theorem 1.83 (Binomial Theorem)

The Binomial Theorem states that for any scalar variables x and y , and for any counting number n , the following equation holds:

$$(x + y)^n = \sum_{k=0}^n \binom{n}{k} x^{n-k} y^k,$$

where $\binom{n}{k} := \frac{n!}{k!(n-k)!}$ is the binomial coefficient. We will prove this by induction.

Base Case (n=1): When $n = 1$, the left side of the equation is $(x + y)^1 = x + y$. The right side of the equation is $\binom{1}{0}x^1y^0 + \binom{1}{1}x^0y^1 = x + y$, because $\binom{1}{0} = 1$, $\binom{1}{1} = 1$, $x^0 = 1$, and $y^0 = 1$. So, the hypothesis holds for $n = 1$. Yeah.

Inductive Step: Assume that the theorem holds for some n , that is:

$$(x + y)^n = \sum_{k=0}^n \binom{n}{k} x^{n-k} y^k$$

We need to show that the theorem holds for $n + 1$, that is:

$$(x + y)^{n+1} = \sum_{k=0}^{n+1} \binom{n+1}{k} x^{n+1-k} y^k$$

We start with the left side of the equation,

$$(x + y)^{n+1} = (x + y)(x + y)^n.$$

Using the inductive hypothesis, we can expand $(x + y)^n$ as

$$(x + y)^{n+1} = (x + y) \sum_{k=0}^n \binom{n}{k} x^{n-k} y^k.$$

Distribute the $(x + y)$ to obtain

$$(x + y)^{n+1} = \sum_{k=0}^n \binom{n}{k} x^{n+1-k} y^k + \sum_{k=0}^n \binom{n}{k} x^{n-k} y^{k+1}.$$

We can shift the index of the second sum to align with the first, via

$$(x + y)^{n+1} = \sum_{k=0}^n \binom{n}{k} x^{n+1-k} y^k + \sum_{k=1}^{n+1} \binom{n}{k-1} x^{n+1-k} y^k.$$

Now, we can combine the two sums, per

$$(x + y)^{n+1} = \sum_{k=0}^{n+1} \left[\binom{n}{k} x^{n+1-k} y^k + \binom{n}{k-1} x^{n+1-k} y^k \right]$$

Finally, using a property of binomial coefficients $\binom{n}{k} + \binom{n}{k-1} = \binom{n+1}{k}$, yields,

$$(x + y)^{n+1} = \sum_{k=0}^{n+1} \binom{n+1}{k} x^{n+1-k} y^k.$$

Proof of Prop. 1.88 Reducing Powers of Sine and Cosine

For all integers $n > 1$,

$$\cos^n(\theta) = \begin{cases} \frac{1}{2^{n-1}} \cdot \sum_{k=0}^{\frac{n-1}{2}} \binom{n}{k} \cos((n-2k)\theta) & n \text{ odd} \\ \frac{1}{2^n} \binom{n}{\frac{n}{2}} + \frac{1}{2^{n-1}} \cdot \sum_{k=0}^{\frac{n-2}{2}} \binom{n}{k} \cos((n-2k)\theta) & n \text{ even} \end{cases} \quad (1.79)$$

and

$$\sin^n(\theta) = \begin{cases} \frac{(-1)^{\frac{n-1}{2}}}{2^{n-1}} \cdot \sum_{k=0}^{\frac{n-1}{2}} \binom{n}{k} (-1)^k \sin((n-2k)\theta) & n \text{ odd} \\ \frac{1}{2^n} \binom{n}{\frac{n}{2}} + \frac{(-1)^{\frac{n}{2}}}{2^{n-1}} \cdot \sum_{k=0}^{\frac{n-2}{2}} \binom{n}{k} (-1)^k \cos((n-2k)\theta) & n \text{ even} \end{cases} \quad (1.80)$$

The derivation for $\cos^n(\theta)$ was already given. The derivation for $\sin^n(\theta)$ follows the same reasoning as the one for $\cos^n(\theta)$, with the key differences being the expansion for $(x - x^{-1})^n$ instead of $(x + x^{-1})^n$, and factoring out $\frac{1}{(2i)^n}$ instead of $\frac{1}{2^n}$ when considering $\sin^n(\theta)$. Keeping track of all the minus signs is painful, and thus we provide every step of the derivation.

Starting with Euler's formula for $\sin^n(\theta)$, we have:

$$\sin^n(\theta) = \left(\frac{e^{i\theta} - e^{-i\theta}}{2i} \right)^n = \frac{1}{(2i)^n} \cdot (e^{i\theta} - e^{-i\theta})^n = \frac{1}{(2i)^n} \cdot (x - x^{-1})^n, \quad (1.81)$$

for $x = e^{i\theta}$. The binomial expansion of $(x - x^{-1})^n$ is:

$$(x - x^{-1})^n = \sum_{k=0}^n \binom{n}{k} x^{n-k} (-x^{-1})^k = \sum_{k=0}^n \binom{n}{k} (-1)^k x^{2k-n}. \quad (1.82)$$

For n odd, the sum has an even number of terms, allowing us to pair the first term with the last term, and so on, to obtain:

$$(x - x^{-1})^n = \sum_{k=0}^{\frac{n-1}{2}} \left(\binom{n}{k} (-1)^k x^{n-2k} + \binom{n}{n-k} (-1)^{n-k} x^{-(n-2k)} \right) = \sum_{k=0}^{\frac{n-1}{2}} \binom{n}{k} (-1)^k (x^{n-2k} - x^{-(n-2k)}), \quad (1.83)$$

where we used $\binom{n}{n-k} = \binom{n}{k}$, and $(-1)^{n-k} = (-1)^n \cdot (-1)^{-k} = -(-1)^k$ (because n is odd and because $(-1)^{-k} = (-1)^k$).

Substituting back into (1.81) for n odd, gives,

$$\begin{aligned}
\sin^n(\theta) &= \frac{1}{(2i)^n} \cdot (x - x^{-1})^n \Big|_{x=e^{i\theta}} \\
&= \frac{1}{(2i)^n} \cdot \sum_{k=0}^{\frac{n-1}{2}} \binom{n}{k} (-1)^k (x^{n-2k} - x^{-(n-2k)}) \Big|_{x=e^{i\theta}} \\
&= \frac{1}{(2i)^{n-1}} \cdot \sum_{k=0}^{\frac{n-1}{2}} \binom{n}{k} (-1)^k \left(\frac{(e^{i\theta})^{n-2k} - (e^{i\theta})^{-(n-2k)}}{2i} \right) \\
&= \frac{1}{(2i)^{n-1}} \cdot \sum_{k=0}^{\frac{n-1}{2}} \binom{n}{k} (-1)^k \left(\frac{e^{i(n-2k)\theta} - e^{-(n-2k)\theta}}{2i} \right) \\
&= \frac{1}{(2i)^{n-1}} \cdot \sum_{k=0}^{\frac{n-1}{2}} \binom{n}{k} (-1)^k \sin((n-2k)\theta) \\
&= \frac{(-1)^{\frac{n-1}{2}}}{2^{n-1}} \cdot \sum_{k=0}^{\frac{n-1}{2}} \binom{n}{k} (-1)^k \sin((n-2k)\theta),
\end{aligned}$$

where on the last step, we used the fact that n odd implies $(n-1)$ is even, and hence $(i)^{n-1} = (i^2)^{\frac{n-1}{2}} = (-1)^{\frac{n-1}{2}}$.

For n even, the sum has an odd number of terms and hence the middle term remains unpaired.

$$\begin{aligned}
(x - x^{-1})^n &= \sum_{k=0}^{\frac{n-2}{2}} \left(\binom{n}{k} (-1)^k x^{n-2k} + \binom{n}{n-k} (-1)^{n-k} x^{-(n-2k)} \right) + (-1)^{\frac{n}{2}} \binom{n}{\frac{n}{2}} \\
&= \sum_{k=0}^{\frac{n-2}{2}} \binom{n}{k} (-1)^k (x^{n-2k} + x^{-(n-2k)}) + (-1)^{\frac{n}{2}} \binom{n}{\frac{n}{2}}.
\end{aligned}$$

where we used again $\binom{n}{n-k} = \binom{n}{k}$, and $(-1)^{n-k} = (-1)^n \cdot 1^{-k} = 1^{-k} = (-1)^k$ (because n is even and because $(-1)^{-k} = (-1)^k$).

Substituting back into (1.81) for n even, gives,

$$\begin{aligned}
\sin^n(\theta) &= \frac{1}{(2i)^n} \cdot (x - x^{-1})^n \Big|_{x=e^{i\theta}} \\
&= (-1)^{\frac{n}{2}} \cdot \frac{1}{(2i)^n} \binom{n}{\frac{n}{2}} + \frac{1}{(2i)^n} \cdot \sum_{k=0}^{\frac{n-2}{2}} \binom{n}{k} (-1)^k (x^{n-2k} + x^{-(n-2k)}) \Big|_{x=e^{i\theta}} \\
&= (-1)^{\frac{n}{2}} \cdot \frac{1}{(2i)^n} \binom{n}{\frac{n}{2}} + \frac{1}{2^{n-1}} \cdot \frac{1}{i^n} \cdot \sum_{k=0}^{\frac{n-2}{2}} \binom{n}{k} (-1)^k \left(\frac{(e^{i\theta})^{n-2k} + (e^{i\theta})^{-(n-2k)}}{2} \right) \\
&= (-1)^{\frac{n}{2}} \cdot \frac{1}{(2i)^n} \binom{n}{\frac{n}{2}} + \frac{1}{2^{n-1}} \cdot \frac{1}{i^n} \cdot \sum_{k=0}^{\frac{n-2}{2}} \binom{n}{k} (-1)^k \left(\frac{e^{i(n-2k)\theta} + e^{-(n-2k)\theta}}{2} \right) \\
&= (-1)^{\frac{n}{2}} \cdot \frac{1}{(2i)^n} \binom{n}{\frac{n}{2}} + \frac{1}{2^{n-1}} \cdot \frac{1}{i^n} \cdot \sum_{k=0}^{\frac{n-2}{2}} \binom{n}{k} (-1)^k \cos((n-2k)\theta) \\
&= \frac{1}{2^n} \binom{n}{\frac{n}{2}} + \frac{(-1)^{\frac{n}{2}}}{2^{n-1}} \cdot \sum_{k=0}^{\frac{n-1}{2}} \binom{n}{k} (-1)^k \cos((n-2k)\theta),
\end{aligned}$$

where on the last step, we used the fact that n even $(i)^n = (i^2)^{\frac{n}{2}} = (-1)^{\frac{n}{2}}$ and $(-1)^{\frac{n}{2}} \cdot (-1)^{\frac{n}{2}} = (-1)^n = 1$. ■

Chapter 2

Calculus Foundations: Proofs, Finite Sums, Limits at Infinity, and Geometric Sums

Learning Objectives

By the end of this chapter, the student should be able to:

- Understand the significance of mathematical proofs and articulate their importance in Calculus and its applications.
- Master the technique of proof by induction and apply it to demonstrate the validity of mathematical statements.
- Comprehend the concept of limits and see some of their initial uses in Calculus.

Outcomes

Upon successful completion of this chapter, students will be able to:

- Apply logical rules to construct (simple) mathematical arguments.
- Learn there are at least two kinds of infinity.
- Analyze sums of powers of integers to build a foundational understanding of integration.
- Calculate limits at infinity for functions that are significant in Calculus.
- Acquire additional knowledge about Euler's number and its unique properties.

2.1 Setting Expectations

While this section could have been placed in the introduction to Chapter 1, there's a good chance that many of you would have skipped it.

2.1.1 How to Use the Textbook

To make the most of this textbook and your learning experience, start by scanning the relevant parts of the textbook before attending (or watching) the lecture. This quick scan, which should take 15 minutes or less, will help you anticipate the main topics and identify potentially challenging points. During the lecture, focus on taking essential notes rather than copying everything verbatim—trust that the textbook has detailed explanations for most material (and besides, the instructor's live notes are posted that day). After the lecture, but before the next one, take time to rewrite and condense your notes onto one side of a single page. The act of summarizing forces you to engage with the material thoughtfully, creating a third encounter with the content after scanning and attending (watching) the lecture. By this point, you should have a solid grasp of the material. If not, this process will help you pinpoint specific questions to ask during Office Hours or study sessions. Remember, understanding comes from active engagement, not just passively reading or watching.

To further enhance your learning, review your condensed notes briefly every few days to reinforce key ideas—spaced practice improves retention. Collaborating with peers can also help you refine your understanding and tackle challenging problems together. Finally, focus on solving problems rather than memorizing formulas; calculus is best learned by doing.

If a topic in lecture went by too fast, leaving you feeling unprepared—such as understanding how to write a proof—additional help is available quickly and easily. You can use UMGPT, our version of ChatGPT, or Maizey, an AI tool indexed on the course textbook and lecture audio transcripts. To make the most of these tools, practice writing specific prompts. A broad request like “Tell me about proofs” may not yield useful guidance, but a focused prompt such as, “I am struggling with writing a proof by induction and I'd like examples with clearly delineated steps and the overall thought process to help me learn,” will provide much more relevant and actionable assistance.

When working on the homework problems, you may discover gaps in your understanding. Sometimes this happens because you zoned out during that part of the lecture—don't worry, it happens! Other times, it's intentional: we deliberately leave certain topics for you to learn on your own. By now, most of you have been in school for over a dozen years, so you've learned how to learn, right? If not, it's never too late to take ownership of your education. Embrace these moments as opportunities to build your skills and independence—they're as valuable as mastering the content itself.

2.1.2 The Role of Lectures

In this course, lectures are designed to be to the textbook what *CliffsNotes* are to WAR AND PEACE or SHAKESPEARE: a condensed, high-level summary that captures the key ideas while providing a framework for deeper understanding. Lectures offer a concise presentation of the most important topics, allowing you to grasp the core concepts quickly and focus on how they connect to broader ideas.

Much like *CliffsNotes*, lectures aim to be:

- **Quick Summaries:** They provide a streamlined overview of complex topics, which can be especially useful for orienting yourself to the material and preparing for deeper study.
- **Accessible:** The straightforward structure and focused content make lectures approachable, particularly for students who may initially struggle with the density of the textbook.
- **Provide Supplemental Learning:** Lectures are not a replacement for reading the textbook. Instead, they complement it by helping you identify key sections and giving you the context needed to fully engage with the more detailed explanations in the book.

After attending a lecture, you are expected to read the corresponding sections of the textbook to gain a fuller understanding of the ideas introduced. The textbook provides the depth, nuance, and detailed examples necessary for mastering the material, which cannot be fully conveyed in a lecture's condensed format.

To make the most of this approach, treat lectures as an orientation tool: focus on understanding the main points and taking essential notes. Use these notes as a guide for your textbook reading, allowing you to delve into the details with clarity and purpose. By combining the broad overview provided in lectures with the depth of the textbook, you will build a well-rounded understanding of the material and be better prepared to tackle challenging problems in homework, projects, and your engineering courses.

2.1.3 Is this a Programming Course?

Absolutely not! The textbook integrates programming with mathematics so that realistic problems can be solved with relative ease. Part of Calculus dread comes from the length of the algebraic manipulations that come with the application of Calculus principles, where, with each new line of calculation, comes a new opportunity for a sign error. Coding in Julia allows us to circumvent much of the drudgery and error sources.

The coding assignments are designed to strike a balance between showcasing the power of Calculus and allowing students to complete the HW assignments within the expected 4- to 6-hour time period for a 200-level course. To achieve this, the textbook provides well-structured code templates that focus on the concepts and methods central to the problem. These templates serve as scaffolding to guide students through the logic of the computation while removing unnecessary coding overhead that could detract from the mathematical insights.

Critically, these assignments aim to marry the elegance of Calculus with the efficiency of Julia. For example, problems might involve using symbolic differentiation to derive gradients, solving systems of nonlinear equations, or numerically integrating trajectories of dynamic systems. With the provided templates, students focus on interpreting and implementing the mathematical models rather than wrestling with syntax or boilerplate code.

Some students may feel that the templates make the coding assignments too easy. However, these assignments are deliberately designed to ensure that time spent on coding enhances understanding of the mathematical principles, not debugging unrelated programming errors. The integration of written assignments alongside coding tasks ensures a well-rounded experience that highlights the theoretical underpinnings of Calculus while demonstrating its application to solving exciting and complex engineering problems.

The goal is not to transform students into programmers but to empower them with tools that extend the reach of their mathematical abilities. By reducing the cognitive load associated with low-level programming, students can engage with real-world problems that would otherwise be inaccessible within the scope of a 200-level course. Julia serves as the bridge between theoretical understanding and practical application, allowing students to see the beauty and utility of Calculus in action.

2.1.4 Why Have Written Homework at All?

The inclusion of written homework ensures that students develop the fluency needed to solve problems by hand, a skill essential for their follow-on engineering courses where manual computations will often be expected. While coding assignments demonstrate the power of computational tools like Julia, written problems serve a complementary role by emphasizing the fundamental principles of Calculus.

These written exercises provide a space for students to reflect on core concepts, strengthen their problem-solving techniques, and gain a deeper understanding of the mathematical structures that underpin their computational solutions. Together, the written and coding tasks create a balanced learning experience, showcasing the elegance and utility of Calculus while equipping students with the practical and conceptual tools necessary for their future studies.

2.1.5 What if You Don't Like Julia?

You have access to the code files on Vocareum. You can paste them into your favorite LLM and ask for them to be translated into Python, MATLAB, C++, etc., or even cleaner Julia! Coding languages come and go. Ultimately, they all provide for the manipulation of arrays, the definition of new functions, the creation of loops, and branching via if-then-else. We chose Julia over MATLAB because Julia is open-source. Moreover, it can be implemented on robots for real-time control as shown in [adubred/KinodynamicFabrics.jl](#) and [Exploring Kinodynamic Fabrics for Reactive Whole-Body Control of Underactuated Humanoid Robots](#), Alphonsus Adu-Bredu, Grant Gibson, and JG.

2.2 What is a Mathematical Proof? And Why are Mathematical Proofs Important?

We ask an expert, **ChatGPT**:

A *mathematical proof* is a logical argument that establishes the truth of a mathematical statement. It is a demonstration that, if certain fundamental statements (called axioms or premises) are assumed to be true, then another statement necessarily follows from these assumptions.

Proofs use the rules of logic combined with the axioms of the mathematical system in question to make their arguments. There are different types of proofs, such as direct proof, proof by contradiction, and proof by induction.

A *direct proof* begins with the assumptions and applies the rules of logic and other mathematical truths that have already been established to arrive at the conclusion.

Proof by induction is a method often used to prove statements about all natural numbers. It involves proving a base case (usually for the number 1 or 0) and then proving that if the statement is true for an arbitrary natural number $k > 0$, it's also true for the number $k + 1$.

Proof by contradiction (also known as *reductio ad absurdum*) assumes that the statement to be proved is actually false, and then shows that this assumption leads to an absurdity. This then means that the statement cannot be false, and hence must actually be true because we are using “binary logic” (any statement is either **T** or **F**; there is nothing in between).

Mathematical proofs are crucial for a number of reasons:

- (a) *Certainty*: Unlike in many other fields, in mathematics, knowledge can be established beyond a shadow of a doubt. A proof guarantees that, given the acceptance of the axioms, the conclusion must be true.
- (b) *Building Blocks*: Each proven statement becomes a tool that mathematicians can use in further proofs. Mathematics is built up from layers upon layers of proven statements.
- (c) *Understanding*: The process of creating a proof often leads to a deeper understanding of the mathematical structure being studied.
- (d) *Discovery*: Sometimes, the process of proving a statement can lead to the discovery of new mathematical facts or even whole areas of mathematics.
- (e) *Rigor*: Proofs provide a standard of rigor that helps prevent errors in mathematical reasoning. This rigor is crucial in areas such as computer programming and algorithm design, where a small error in reasoning can lead to significant problems.

Mathematical Proofs Promote Understanding

In short, mathematical proofs are the cornerstone of mathematical knowledge, providing certainty, depth of understanding, and a rigorous basis for further discovery. We'll add one more reason they are important: some things that seem “obviously false”, even to the larger Mathematical Community, can turn out to be true, and vice versa! In the next Section, we'll use a famous example to accomplish two things: (1) learn about countable and uncountable sets and (2) underline that you can be right even when everyone tells you that you are wrong. A sad example of the latter is [Georg Cantor](#), whose surprising discovery about the nature of infinity led to public opposition and personal attacks. Mathematics being as controversial as [evolution](#)? You cannot make this stuff up! Please read on.

Related Videos:

- [Don't let it fool you](#) by [3Blue1Brown](#) (moral of the story: patterns often break).
- [Proofs Are Awesome! Here's Why!](#) Math Time With Professor Prime (a short low-key pep talk).
- [Is math discovered or invented?](#) a TED-ED talk by Jeff Dekofsky (asks the question, “does math exist independent of humans”).

2.3 Vocabulary and Helpful Notation

The following definitions borrow liberally from [Math 299](#), Michigan State University. In this course, you will rarely be asked to write proofs yourself. You will be asked to read them from time to time. In lectures and [quizzes](#), we will use the terminology, vocabulary, and shorthand notation called out below.

Terminology:

- *Definition*: A precise explanation of the mathematical meaning of a term or concept.
- *Axiom*: A fundamental principle or assumption that is accepted as true without proof in a mathematical system.
- *Theorem*: A mathematical statement that has been proven to be true.
- *Lemma*: A mathematical statement that is proven true and used as a stepping stone in the proof of other statements.

- *Claim*: A mathematical statement that is made as a step towards proving a theorem or to highlight an important property of a mathematical object. When used as a stepping stone to prove something else, Claims are essentially equivalent to Lemmas. When a Claim is used to highlight an important property, it is a “baby” Proposition.
- *Proposition*: A mathematical statement that is less central but nonetheless true and interesting.
- *Corollary*: A mathematical statement that is a straightforward deduction from a theorem or proposition.
- *Proof*: The logical reasoning that demonstrates why a statement is true.
- *Conjecture*: A mathematical statement that is believed to be true, but for which no proof has been established. Or, a communication of a potentially true statement that the author hopes to prove or that the mathematical community seeks to prove. **Fermat’s Last Theorem** was a Conjecture until it was proven by **Wiles** in 1995.

We will mostly use just definitions and propositions so that we do not get tied up in the hierarchy of theorems, lemmas, etc. We will not take this to extremes, by calling the very well-known “Pythagorean Theorem” the “Pythagorean Proposition”, or awkwardly referring to the “Binomial Proposition”.

Crib Notes for Proofs: where P and Q are logical statements. They can be either true (denoted **T**) or false (denoted **F**). Because we use binary logic, there is no in-between value, such as “maybe”. Statement are either **T** or **F**.

In the following, you might have in mind, $P : \det(A) \neq 0$ and $Q : A$ is a square invertible matrix, or, $P :$ the rows of A are linearly independent and $Q : Ax = b$ has a minimum norm solution.

- \neg or \sim denotes the logical operation “NOT”, where $\neg\mathbf{T} = \mathbf{F}$ and $\neg\mathbf{F} = \mathbf{T}$. Both symbols will be used interchangeably in these notes.
- $P \implies Q$ represents “if statement P is true, then statement Q is also true”.
- $P \iff Q$ indicates that “ P is true if and only if Q is true”. Although P iff Q is a valid expression, it will be seldom used in these notes.
- $P \iff Q$ is logically equivalent to:
 - (a) $P \implies Q$ and
 - (b) $Q \implies P$.
- The *contrapositive* of $P \implies Q$ is denoted as $\neg Q \implies \neg P$.
- The *converse* of $P \implies Q$ is expressed as $Q \implies P$. It is crucial to note that $P \implies Q$ does not generally imply $Q \implies P$ and vice versa. If they were equivalent, there would be no need for $P \iff Q$.
- The relation $(P \implies Q) \iff (\neg Q \implies \neg P)$ indicates that these two statements are logically equivalent. Although they are equivalent, one may be easier to use in a proof than the other.
- *Logical AND*: $P_1 \wedge P_2$ is read as “ P_1 AND P_2 ”. It holds true when both P_1 and P_2 are true.
- *Logical OR*: $P_1 \vee P_2$ is interpreted as “ P_1 OR P_2 ”. It is true when at least one of P_1 or P_2 is true. In this course, we use the “inclusive OR”, which means $\mathbf{T} \vee \mathbf{T} = \mathbf{T} \vee \mathbf{F} = \mathbf{F} \vee \mathbf{T} = \mathbf{T}$, where **T** stands for true and **F** stands for false.
- Q.E.D. is an abbreviation of the Latin phrase “quod erat demonstrandum”, which translates to “what was to be demonstrated”. It is used to indicate the completion of a proof. Nowadays, you more frequently see \square or \blacksquare instead of Q.E.D.

Warning: In the beginning, it is quite frequent that students confuse the meanings of *contrapositive* and *converse*. Just be careful. With practice, it becomes second nature. The fact that they both start with “con-” is not helpful!

(Optional Read:) Symbology of Mathematics:

Definition 2.1. Mathematical Shorthand:

- (a) \forall means “**for every**”, “**for all**”, or “**for each**” as in $\forall x \in \mathbb{R}, x^0 := 1$, which means that for all real numbers x , we define x^0 to stand for the number one. Continuing in this vein, the notation $\forall k \in \mathbb{N}$ and $\forall x \in \mathbb{R}, x^k := x \cdot x^{k-1}$, which means that for all counting numbers, and for all real numbers x , we define x^k to stand for the product of x times x^{k-1} . We note that x^1 makes sense because $x^1 := x \cdot x^0 = x \cdot 1$, and we know how to form the product $x \cdot 1$ within the set of real numbers.
- (b) \exists means “**there exist(s)**”, “**there is/are**”, “**for some**”, “**for at least one**” as in $\exists x \in [0, \pi/2]$ such that $\sin(x) = \frac{\sqrt{2}}{2}$, which reads, there exists at least one point x in the interval from zero to $\pi/2$ such that $\sin(x) = \frac{\sqrt{2}}{2}$. Or let A and B be sets. Then $f : A \rightarrow B$ is a function if, $\forall a \in A, \exists b \in B$, such that $f(a) = b$, that is, f is a function from A to B if for all elements a in the set A there is some element b in the set B such that $f(a) = b$. Or a final example, $\forall x \in \mathbb{R}$ and $\forall \epsilon > 0, \exists r \in \mathbb{Q}$ such that $|x - r| < \epsilon$, which means: for every real number x and every positive number ϵ no matter how small, there exists a rational number r that approximates x with an error that is less than or equal to ϵ (or you could say, approximates x within a tolerance of ϵ). When we do numerical computations on a computer, we use rational approximations of real numbers. In that case, the smallest positive number that, when added to 1, yields a result different from 1 in the machine’s arithmetic is called machine epsilon. For single-precision IEEE 754 standard floating-point format (often used in the float data type in programming languages), the machine epsilon is 2^{-23} , or approximately 1.19×10^{-7} , and for double-precision IEEE 754 standard floating-point format (often used in the double data type in programming languages), the machine epsilon is 2^{-52} , or approximately 2.22×10^{-16} .
- (c) \nexists means “**there do/does not exist**”, “**there is no element**”, or “**there are no elements**”, as in $\nexists x \in \mathbb{R}$ such that $x \cdot 0 = 1$, which is how you would write that the number zero does not have a multiplicative inverse. Or another example, $\tan(x)$ is not a function from \mathbb{R} to \mathbb{R} because $\pi/2$ is in the domain \mathbb{R} , but $\nexists y \in \mathbb{R}$ (used here as the codomain) such that $\tan(\pi/2) = y$. Recall that a function $f : A \rightarrow B$ **must assign** to each element of the domain, A , an element in the codomain, B .

Initially, these symbols might appear as unnecessary complexities. **Yet, they aim to serve as a concise language**, enabling precise and efficient communication of concepts that we encounter regularly. With practice, you will come to appreciate them just as much as you appreciate Descartes’ notation for unknowns.

2.4 Proofs by Induction, Some Finite Sums, and Their Applications

The only way to learn proofs is by doing them.

2.4.1 Two Finite Sums via Induction

Logic Principle 2.2: First Principle of Induction (Standard Induction)

Let $P(n)$ denote a statement about the counting numbers with the following properties:

- (a) **Base case:** $P(1)$ is true
(b) **Induction hypothesis:** For all $k \geq 1$, if $P(k)$ is true, then $P(k + 1)$ is true.

Then $P(n)$ is true for all $n \geq 1$.

Remark 2.3. Suppose the base case involves an integer $k_0 \neq 1$. Then you can re-index things and reduce it to the base case having $k_0 = 1$. Alternatively, you assume that $P(k)$ true for $k \geq k_0$ implies $P(k + 1)$ is true, and then you get $P(n)$ is true for all $n \geq k_0$. A common mistake is not to use the correct base case. For example, you can read about “how to prove by induction” that all **horses are the same color**.

Example 2.4. Let's prove the **Claim:** For all $n \geq 1$,

$$1 + 2 + 3 + \cdots + n = \frac{n(n+1)}{2}. \quad (2.1)$$

Proof:

- **Step 0:** Write down $P(n)$: $1 + 2 + 3 + \cdots + n = n(n+1)/2$, the statement that is to be proved.
- **Step 1:** Check the base case, $P(1)$: For $n = 1$, we have that is $1 = 1$ for the left side (the sum of the first one integer), and for the right side, $n(n+1)/2 = 1 \cdot (1+1)/2 = 2/2 = 1$. Because $1 = 1$, the base case is true.
- **Step 2:** Show the induction hypothesis is true. That is, assuming that $P(k)$ is true, show that $P(k+1)$ is true. For us,

$$P(k+1) : (1 + 2 + \cdots + k) + (k+1) = (k+1) \cdot (k+2)/2.$$

Often, showing that $P(k) = \mathbf{T} \implies P(k+1) = \mathbf{T}$ involves re-writing the terms in $P(k+1)$ as a sum of terms that show up in $P(k)$, plus another term.

For the induction step, we assume that

$$P(k) : 1 + 2 + 3 + \cdots + k = k \cdot (k+1)/2$$

is true. Next, we note there are terms in $P(k+1)$ that are common to $P(k)$, namely

$$(1 + 2 + \cdots + k) + (k+1) = \underbrace{(1 + 2 + \cdots + k)}_{k \cdot (k+1)/2} + (k+1) = k \cdot (k+1)/2 + (k+1),$$

because we assumed that $P(k)$ holds. Then, by the accepted rules of arithmetic,

$$k(k+1)/2 + (k+1) = (k+1) \cdot (k/2 + 1) = (k+1) \cdot (k+2)/2,$$

and we deduce that

$$(1 + 2 + \cdots + k) + (k+1) = (k+1) \cdot (k+2)/2,$$

and therefore, $P(k+1)$ is true. Because we have shown that $P(1)$ is true and for all $k \geq 1$, $P(k) \implies P(k+1)$, by the Principle of Induction, we conclude that for all $n \geq 1$,

$$1 + 2 + 3 + \cdots + n = n(n+1)/2,$$

that is, (2.1) holds. ■

Remark 2.5. If you run into a case where base case $P(1)$ seems so totally trivial that you are unsure whether there is anything to show, it's OK to check $P(2)$, just to convince yourself that it is true too. In our case, you would have checked $P(2) : 1 + 2 = 2 \cdot 3/2$; since this is easily established to be true, you may now be more confident of your proof. If you do one more, $P(3) : 1 + 2 + 3 = 3 \cdot 4/2$, you are now on a roll and ready to attack the general case by induction. Bottom line: Initially, it's natural to be tentative when you write a proof. **It takes practice to learn the art of proving things. In this course, writing proofs is not a core skill we are trying to develop. Instead, we want you to respect that proofs are important because they explain why something is true.** They turn a mathematical result from "black box" (opaque, unclear, huh?, pure manipulation of symbols) to "white box" (clear, transparent, understandable, Oh, that's why!, feeling like you know what you are doing)¹.

¹The terms "black box" and "white box" are commonly used in various fields, including systems engineering, computing, and business, with their exact meaning depending on the context.

- **Black Box Testing/Modeling:** In black box testing or modeling, the internal workings of the system or model are not known or considered. It's called a "black box" because the system is opaque; we can't see what's happening inside. Instead, we focus on the input and output. We provide certain inputs, and then test to see if the output is as expected. In the context of AI, a black box model is a system where we can observe the input and output, but not the internal decision-making process. For example, deep learning models are often considered black boxes because, although they can make highly accurate predictions, it's difficult to understand how they arrived at a particular decision.
- **White Box Testing/Modeling:** In white box testing or modeling, we do have access to the internal workings of the system or model. This type of testing is sometimes called "clear box" or "glass box" testing. We not only examine input and output but also the processes that lead from one to the other. For example, a decision tree can be considered a white box model because it's possible to follow the path from beginning to end to understand exactly how a decision is made.

In a similar sense, a mathematical proof provides the steps that show us how to go from the hypotheses (i.e., input) to the claimed result (i.e., output). Without proofs, we are taking results on faith.

Example 2.6. Let's prove the **Claim:** For all $n \geq 1$,

$$1^2 + 2^2 + \dots + n^2 = \frac{n(n+1)(2n+1)}{6} \quad (2.2)$$

Proof:

Proofs can be written in many different ways. Let's approach this proof a bit differently by defining

$$\begin{aligned} L(n) &:= 1^2 + 2^2 + \dots + n^2, \text{ the left side of the Claim,} \\ R(n) &:= n(n+1)(2n+1)/6, \text{ the right side of the Claim.} \end{aligned} \quad (2.3)$$

We want to show for all counting numbers $n \geq 1$,

$$P(n) \text{ is true, that is, } L(n) = R(n).$$

- **Step 0:** Write down $P(k)$: $L(k) = R(k)$.
- **Step 1:** Check the base case, $P(1)$: For $k = 1$, we have that $L(1) = 1$ and $R(1) = 1 \cdot (1+1) \cdot (2+1)/6 = 1$ and hence the base case is true.
- **Step 2:** Let's check one more just so we get the hang of it, $P(2)$: we have that $L(2) = 1^2 + 2^2 = 5$ and $R(2) = 2 \cdot (2+1)(2 \cdot 2+1)/6 = 2 \cdot (15)/6 = 5$, and hence $L(2) = R(2)$, showing that $P(2)$ holds. Our confidence is higher!
- **Step 3:** Show the induction hypothesis is true. That is, using the fact that $P(k)$ is true for some $k \geq 1$, show that $P(k+1)$ is true.

When we assume that $P(k)$ is true, we are assuming that

$$L(k) = R(k).$$

To see if this implies $L(k+1) = R(k+1)$, we first evaluate

$$L(k+1) = \underbrace{1^2 + 2^2 + \dots + k^2}_{L(k)} + (k+1)^2 = L(k) + (k+1)^2.$$

Using the induction hypothesis, namely, $L(k) = R(k) = k(k+1)(2k+1)/6$, we compute

$$\begin{aligned} L(k+1) &= L(k) + (k+1)^2 \\ &= \underbrace{k(k+1)(2k+1)/6}_{L(k)} + (k+1)^2 \\ &= \underbrace{(2k^3 + 3k^2 + k)/6}_{k(k+1)(2k+1)/6} + \underbrace{(k^2 + 2k + 1)}_{(k+1)^2} \\ &= (k^3/3 + k^2/2 + k/6) + (k^2 + 2k + 1) \\ &= k^3/3 + 3k^2/2 + 13k/6 + 1. \end{aligned} \quad (2.4)$$

There was some “messy” algebra, but nothing too taxing.

Next, we plug $(k+1)$ into the definition of $R(k)$ to evaluate

$$\begin{aligned} R(k+1) &= (k+1) \cdot ((k+1)+1) \cdot (2(k+1)+1)/6 \\ &= (k+1) \cdot (k+2) \cdot (2k+3)/6 \\ &= (k+1) \cdot (2k^2 + 7k + 6)/6 \\ &= k^3/3 + 3k^2/2 + 13k/6 + 1. \end{aligned} \quad (2.5)$$

Once again, there was some “messy” algebra, but nothing too taxing.

Comparing (2.4) and (2.5), we see that indeed, $L(k+1) = R(k+1)$, and therefore $P(k+1)$ is true. Because we have shown that $P(1)$ is true, and for all $k \geq 1$, $P(k) \implies P(k+1)$, by the Principle of Induction, we conclude that for all $n \geq 1$,

$$1^2 + 2^2 + \dots + n^2 = n(n+1)(2n+1)/6.$$

■

This example, from [Prove that \$n^3 + 11n\$ is divisible by 6](#) by PrimeNewtons, is included to show a different style of proof by induction, always with the goal of enhancing your understanding. A counting number (positive integer) $i \in \mathbb{N}$ is divisible by 6 if there exists another integer $m \in \mathbb{N}$ such that $i = 6 \cdot m$; in other words, $\frac{i}{6}$ is also a counting number.

Exercise 2.7. Use induction to show that for all $n \in \mathbb{N}$, $n^3 + 11n$ is divisible by 6.

Solution: We want to show for all counting numbers $n \geq 1$, $P(n) : n^3 + 11n = 6 \cdot m$, for some counting number m .

Base Case: $n = 1$: $1^3 + 11 \cdot 1 = 12 = 2 \cdot 6$, and hence the base case is true.

Induction Step: We show that if $P(k)$ is true for some $k \geq 1$, then $P(k+1)$ is true.

We first use basic algebra to arrive at

$$\begin{aligned} (k+1)^3 + 11 \cdot (k+1) &= (k^3 + 3k^2 + 3k + 1) + (11k + 11) \\ &= (k^3 + 11k) + 3 \cdot (k^2 + k + 4). \end{aligned} \tag{2.6}$$

Because $P(k)$ holds, there exists $m_1 \in \mathbb{N}$ such that $(k^3 + 11k) = 6m_1$. We prove below that $(k^2 + k + 4)$ is even for all counting numbers k , and hence there exists m_2 such that $(k^2 + k + 4) = 2m_2$. Substituting these facts into (2.6) gives

$$\begin{aligned} (k+1)^3 + 11 \cdot (k+1) &= (k^3 + 11k) + 3 \cdot (k^2 + k + 4) \\ &= 6m_1 + 3 \cdot 2m_2 \\ &= 6 \cdot (m_1 + m_2). \end{aligned}$$

This completes the proof because, $m_1 + m_2$, the sum of two counting numbers, is another counting number.

Claim: $k^2 + k + 4$ is even for all counting numbers k .

Proof: Uses basic properties of counting numbers, whose proofs you can either provide or seek online. If k is odd, then so is k^2 . Because the sum of two odd numbers is even, we have that $k^2 + k$ is even. Finally, adding the even number 4 to the even number $k^2 + k$ gives an even number. Next, if k is even, so is k^2 . Because the sum of two even numbers is even, we have that $k^2 + k$ is even. Finally, adding the even number 4 to the even number $k^2 + k$ gives an even number.

Note: The information we learned from $P(k)$, namely that $k^3 + 11k$ is divisible by 6, needed to be augmented with an additional fact about counting numbers, namely that $k^2 + k + 4$ is always even.

■

Related Videos:

- Basic: [Proof by induction | Sequences, series and induction | Precalculus](#) by Khan Academy.
- Basic: [Mathematical Induction Practice Problems](#) by the Organic Chemistry Tutor.
- More Advanced: [Prove by induction, Sum of the first n cubes, \$1^3 + 2^3 + 3^3 + \dots + n^3\$](#) by BlackPenRedPen.
- Quite Advanced: [Introduction to Proof by Induction \[Discrete Math Class\]](#) by Mathematical Visual Proofs.

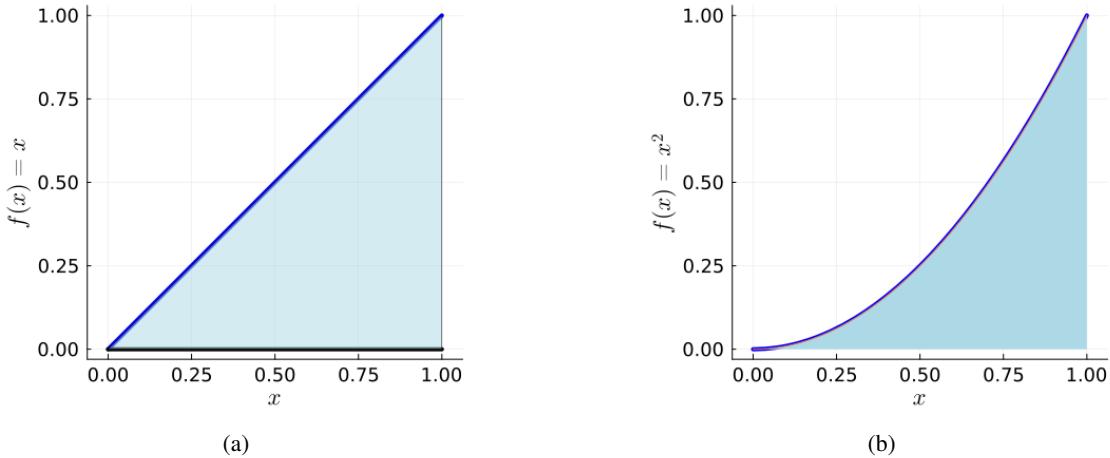


Figure 2.1: Areas under curves. In blue, the areas under two curves between $x = 0$ and $x = 1$ are shown. (a) For a triangle, we know the formula for its area, $\frac{1}{2}b \cdot h$, one-half base times height! (b) For the parabola $y = x^2$, we do not have a simple formula for its area. We'll develop a method to estimate lower and upper bounds for the computation of areas. We'll first show how our estimate works for the triangle, and then apply it to the parabola. Sounds like a plan?

2.4.2 The Amazing Power of the Rectangle: Finding Area under a Curve

In this section, we will use what we have learned so far to lay the foundations for computing areas under the simple “curves” shown in Fig. 2.1, namely a line and a parabola, **assuming only knowledge of the area of a rectangle!** How can we do this? . We will follow Archimedes’ basic plan: create both a lower bound and an upper bound for our desired quantity (for Archimedes, it was the circumference of the circle, here, it is area under the curve), and then show that we can make the lower and upper bounds approach one another. You may wish to refer to Fig. 1.1 and the text around it to refresh your memory.

The curves in Fig. 2.1 are monotonically increasing, making it particularly easy to formulate upper and lower bounds for their areas by summing the areas of carefully chosen rectangles. We already know a formula for the area of a triangle, namely, one-half base times height, $\frac{1}{2}b \cdot h$, giving us 0.5 for the area in Fig. 2.1-(a). What is the area for the parabola in Fig. 2.1-(b)?

Method 2.8: Lower and Upper Bounds for Area Under an Increasing Function

Approximating Area: We seek to estimate the area under an increasing function $f : [a, b] \rightarrow \mathbb{R}$, where $a < b$. We first divide the interval $[a, b]$ into $n > 1$ evenly spaced subintervals, $[x_i, x_{i+1}]$, where

$$\begin{aligned}\Delta x &:= \frac{b - a}{n} \\ x_i &:= a + (i - 1) \cdot \Delta x,\end{aligned}\tag{2.7}$$

so that

$$a = x_1 < x_2 < \dots < x_n < x_{n+1} = b\tag{2.8}$$

$$x_{i+1} = x_i + \Delta x, \quad 1 \leq i \leq n\tag{2.9}$$

$$f(x_1) \leq f(x_2) \leq \dots \leq f(x_n) \leq f(x_{n+1}).\tag{2.10}$$

- We underestimate the area between $[x_i, x_{i+1}]$ by a rectangle of height $h_i^{\text{LOW}} := f(x_i)$ and base $b_i := \Delta x$. This is an **underapproximation** because $h_i^{\text{LOW}} := f(x_i) \leq f(x)$ for all $x \in [x_i, x_{i+1}]$. We underestimate the total area between $[a, b]$ by summing up the underapproximations,

$$\text{Area}_n^{\text{Low}} := \sum_{i=1}^n h_i^{\text{LOW}} \cdot b_i = \sum_{i=1}^n f(x_i) \cdot \Delta x.\tag{2.11}$$

- We overestimate the area between $[x_i, x_{i+1}]$ by a rectangle of height $h_i^{\text{UP}} := f(x_{i+1})$ and base $b_i := \Delta x$. This is an **overapproximation** because $h_i^{\text{UP}} := f(x_{i+1}) \geq f(x)$ for all $x \in [x_i, x_{i+1}]$. We overestimate the total area between

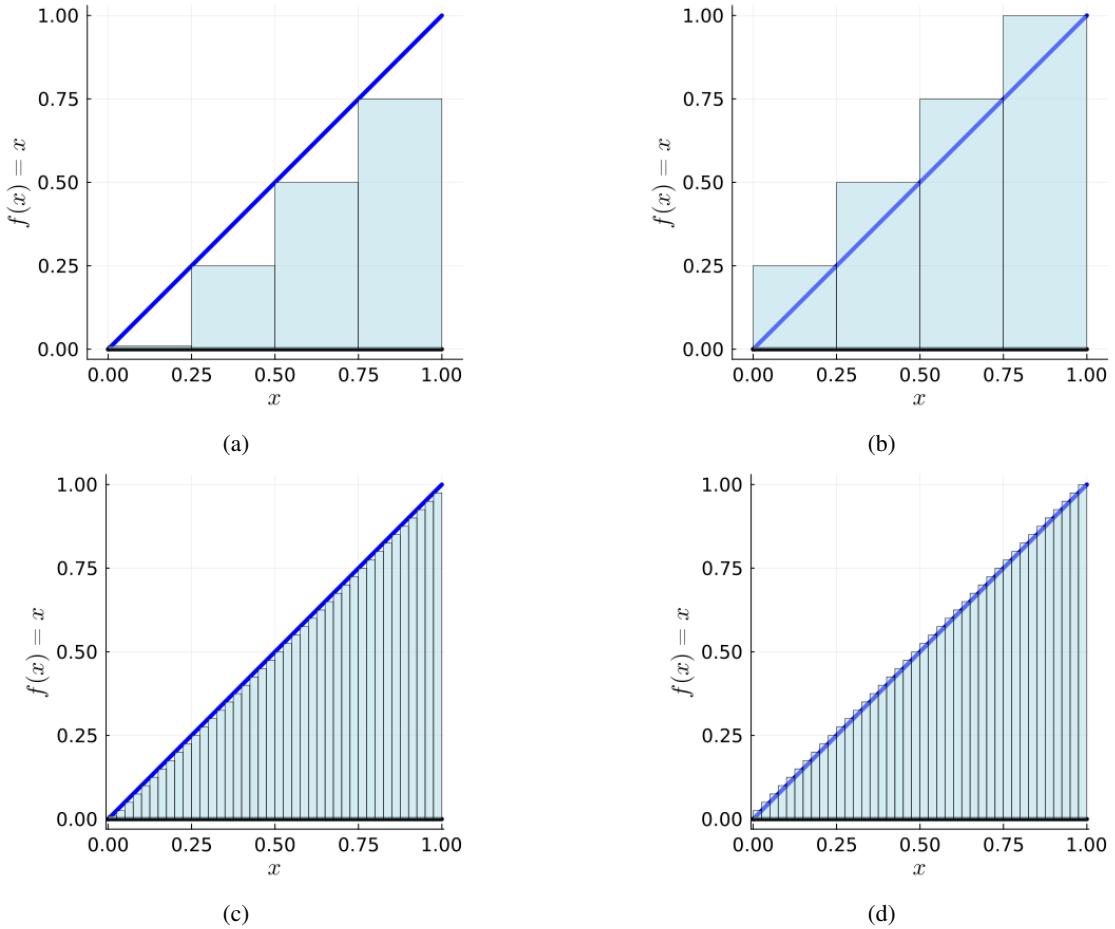


Figure 2.2: Under- and overapproximations for the area of a triangle defined by a line of slope one through the origin. (a) An underapproximation for the triangle defined by $f(x) = x$, $0 \leq x \leq 1$. (b) An overapproximation for the triangle defined by $f(x) = x$, $0 \leq x \leq 1$. For (a) and (b), $n = 4$ and thus $\Delta x = 0.25$, while for (c) and (d), $n = 40$ and $\Delta x = 0.025$.

$[a, b]$ by summing up the overapproximations,

$$\text{Area}_n^{\text{Up}} := \sum_{i=1}^n h_i^{\text{Up}} \cdot b_i = \sum_{i=1}^n f(x_{i+1}) \cdot \Delta x. \quad (2.12)$$

- See Fig. 2.2 and Fig. 2.3.

Then, just like with Archimedes, we have that for all $n > 1$,

$$\text{Area}_n^{\text{Low}} \leq \text{Area} \leq \text{Area}_n^{\text{Up}},$$

where Area is the true area under the curve between the values $x = a$ and $x = b$. Moreover, as we take n larger and larger, the lower and upper bounds will approach one another, giving us the ability to approximate the area with arbitrary accuracy.

Let's carry out this method for $f(x) = x$, $a = 0$ and $b = 1$. Then with $n > 0$, we have $\Delta x = \frac{1}{n}$, $x_i = (i-1)\Delta x$, $f(x_i) = (i-1)\Delta x$, and

$$\begin{aligned}
\text{Area}_n^{\text{Low}} &:= \sum_{i=1}^n f(x_i) \cdot \Delta x \\
&= \sum_{i=1}^n ((i-1)\Delta x) \cdot \Delta x \\
&= \sum_{i=1}^n \left((i-1) \frac{1}{n} \right) \cdot \frac{1}{n} \\
&= \frac{1}{n^2} \cdot \sum_{i=1}^n (i-1) \quad (\text{because } n \text{ is a constant, we take it out of the sum}) \\
&= \frac{1}{n^2} \cdot \sum_{i=1}^{n-1} i \quad (\text{change of index}) \\
&= \frac{1}{n^2} \cdot \frac{(n-1)(n)}{2} \quad (\text{by Example 2.4}) \\
&= \frac{1}{2} - \frac{1}{2n},
\end{aligned} \tag{2.13}$$

Similarly,

$$\begin{aligned}
\text{Area}_n^{\text{Up}} &:= \sum_{i=1}^n f(x_{i+1}) \cdot \Delta x \\
&= \sum_{i=1}^n (i\Delta x) \cdot \Delta x \\
&= \sum_{i=1}^n \left(i \frac{1}{n} \right) \cdot \frac{1}{n} \\
&= \frac{1}{n^2} \cdot \sum_{i=1}^n i \quad (\text{because } n \text{ is a constant, we take it out of the sum}) \\
&= \frac{1}{n^2} \cdot \frac{n(n+1)}{2} \quad (\text{by Example 2.4}) \\
&= \frac{1}{2} + \frac{1}{2n}.
\end{aligned} \tag{2.14}$$

Hence, for all $n > 1$,

$$\text{Area}_n^{\text{Low}} = \frac{1}{2} - \frac{1}{2n} \leq \text{Area} \leq \frac{1}{2} + \frac{1}{2n} = \text{Area}_n^{\text{Up}}. \tag{2.15}$$

This is pretty nice, because we know that the true area of the triangle is $\frac{1}{2}$, and we easily understand that we can make $\frac{1}{2n}$ as small as we like by taking n sufficiently large. Because $\frac{\text{Area}_n^{\text{Up}} + \text{Area}_n^{\text{Low}}}{2} = \frac{1}{2}$ and $\frac{\text{Area}_n^{\text{Up}} - \text{Area}_n^{\text{Low}}}{2} = \frac{1}{n}$, we have

$\text{Area}_n^{\text{est}} = \frac{1}{2} \pm \frac{1}{n}.$

(2.16)

Area Under a Parabola: Next, let's carry out our Archimedes-inspired approximation method for $f(x) = x^2$, $a = 0$, and $b = 1$.

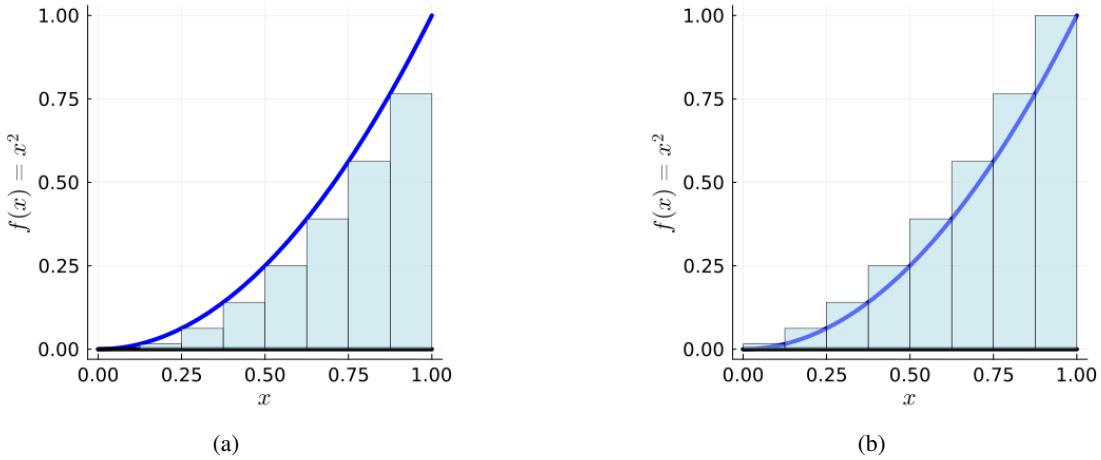


Figure 2.3: Under- and overapproximations for the area under a parabola. (a) An underapproximation for a parabola defined by $f(x) = x^2$, $0 \leq x \leq 1$. (b) An overapproximation for a parabola defined by $f(x) = x^2$, $0 \leq x \leq 1$. For (a) and (b), $n = 8$ and thus $\Delta x = 0.125$.

Again, with $n > 0$, we have $\Delta x = \frac{1}{n}$, $x_i = (i - 1)\Delta x$, $f(x_i) = ((i - 1)\Delta x)^2$, and

$$\begin{aligned}
\text{Area}_n^{\text{LOW}} &:= \sum_{i=1}^n f(x_i) \cdot \Delta x \\
&= \sum_{i=1}^n ((i - 1)\Delta x)^2 \cdot \Delta x \\
&= \sum_{i=1}^n \left((i - 1) \frac{1}{n} \right)^2 \cdot \frac{1}{n} \\
&= \frac{1}{n^3} \cdot \sum_{i=1}^n (i - 1)^2 \quad (\text{because } n \text{ is a constant, we take it out of the sum}) \\
&= \frac{1}{n^3} \cdot \sum_{i=1}^{n-1} i^2 \quad (\text{change of index}) \\
&= \frac{1}{n^3} \cdot \frac{(n-1)(n-1+1)(2(n-1)+1)}{6} \quad (\text{by Example 2.6}) \\
&= \frac{1}{n^3} \cdot \frac{2n^3 - 3n^2 + n}{6} \\
&= \frac{1}{3} - \frac{1}{2n} + \frac{1}{6n^2}.
\end{aligned} \tag{2.17}$$

Similarly,

$$\begin{aligned}
\text{Area}_n^{\text{UP}} &:= \sum_{i=1}^n f(x_{i+1}) \cdot \Delta x \\
&= \sum_{i=1}^n (i\Delta x)^2 \cdot \Delta x \\
&= \sum_{i=1}^n \left(i \frac{1}{n} \right)^2 \cdot \frac{1}{n} \\
&= \frac{1}{n^3} \cdot \sum_{i=1}^n i^2 \quad (\text{because } n \text{ is a constant, we take it out of the sum}) \\
&= \frac{1}{n^3} \cdot \frac{n(n+1)(2n+1)}{6} \quad (\text{by Example 2.6}) \\
&= \frac{1}{n^3} \cdot \frac{2n^3 + 3n^2 + n}{6} \\
&= \frac{1}{3} + \frac{1}{2n} + \frac{1}{6n^2}.
\end{aligned} \tag{2.18}$$

Hence, for all $n > 1$,

$$\text{Area}_n^{\text{Low}} = \frac{1}{3} - \frac{1}{2n} + \frac{1}{6n^2} \leq \text{Area} \leq \frac{1}{3} + \frac{1}{2n} + \frac{1}{6n^2} = \text{Area}_n^{\text{Up}}. \quad (2.19)$$

Because $\frac{\text{Area}_n^{\text{Up}} + \text{Area}_n^{\text{Low}}}{2} = \frac{1}{3} + \frac{1}{6n^2}$ and $\frac{\text{Area}_n^{\text{Up}} - \text{Area}_n^{\text{Low}}}{2} = \frac{1}{n}$, we have

$$\boxed{\text{Area}_n^{\text{est}} = \left(\frac{1}{3} + \frac{1}{6n^2} \right) \pm \frac{1}{n}.} \quad (2.20)$$

Moreover, as n gets larger and larger, both $\frac{1}{6n^2}$ and $\frac{1}{n}$ become smaller and smaller, leading us to the conclusion that the area under the parabola is $\frac{1}{3}$.

Summary: To estimate the areas, all we did was add up the areas of a carefully selected set of rectangles. If we were doing the entire thing in Julia, we would not have needed any analytical results at all. We'd have naturally taken n large and obtained lower and upper bounds on our desired quantity. In order to obtain analytical results, all we needed were Examples 2.4 and 2.6, which address the sum of the first n integers raised to power “1” and power “2”, respectively. That's a pretty low bar, you have to admit.

Area Under under a Triangle and a Parabola, Generalized just a Smidgen

Instead of computing the areas under the curves over the interval $[0, 1]$, let's consider $[0, x]$, for $x > 0$. We won't repeat the steps for the triangle, because the formula $\frac{1}{2}b \cdot h$ still applies, giving us $x^2/2$ for the area. What about the parabola? For the lower bound, with $n > 1$ rectangles, we have $\Delta x = \frac{x}{n}$, $x_i = (i-1)\Delta x$, $f(x_i) = ((i-1)\Delta x)^2$, and

$$\begin{aligned} \text{Area}_n^{\text{Low}} &:= \sum_{i=1}^n f(x_i) \cdot \Delta x \\ &= \sum_{i=1}^n ((i-1)\Delta x)^2 \cdot \Delta x \\ &= \sum_{i=1}^n \left((i-1) \frac{x}{n} \right)^2 \cdot \frac{x}{n} \\ &= \frac{x^3}{n^3} \cdot \sum_{i=1}^n (i-1)^2 \quad (\text{because both } n, x \text{ are independent of } i, \text{ we take them out of the sum}) \\ &= \frac{x^3}{n^3} \cdot \sum_{i=1}^{n-1} i^2 \quad (\text{change of index}) \\ &= \frac{x^3}{n^3} \cdot \frac{(n-1)(n-1+1)(2(n-1)+1)}{6} \quad (\text{by Example 2.6}) \\ &= \frac{x^3}{n^3} \cdot \frac{2n^3 - 3n^2 + n}{6} \\ &= \frac{x^3}{3} - \frac{x^3}{2n} + \frac{x^3}{6n^2}. \end{aligned} \quad (2.21)$$

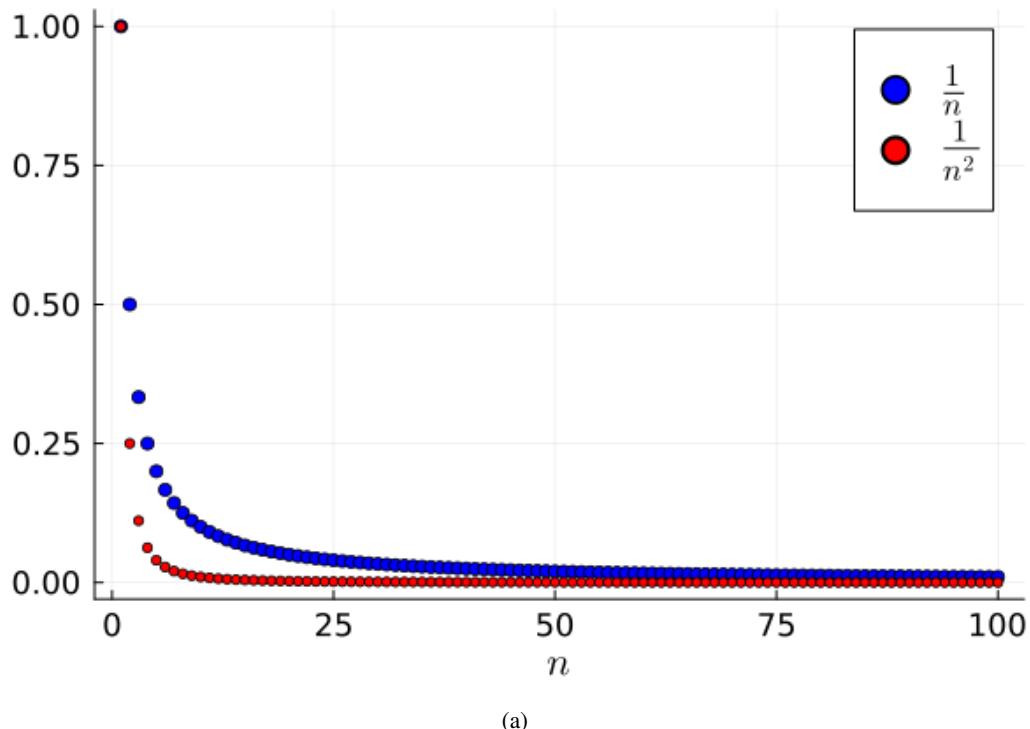
When the learner repeats the calculation for the upper bound, they will obtain $\text{Area}_n^{\text{Up}} = \frac{x^3}{3} + \frac{x^3}{2n} + \frac{x^3}{6n^2}$. Combining the lower and upper bounds gives,

$$\boxed{\text{Area}_n^{\text{est}} = \left(\frac{x^3}{3} + \frac{x^3}{6n^2} \right) \pm \frac{x^3}{2n}.} \quad (2.22)$$

Moreover, for x fixed (i.e., held constant), as n gets larger and larger, the learner can verify numerically, or via Fig. 2.4, that both $\frac{x^3}{6n^2}$ and $\frac{x^3}{2n}$ become smaller and smaller, leading us to the conclusion that the area under the parabola is $\frac{x^3}{3}$. For $x > 0$, the area under the parabola $f(x) = x^2$ for the interval $[0, x]$ is equal to $\frac{x^3}{3}$. In other words, we have a closed-form solution for the area for a general interval $[0, x]$, with very little effort. That's remarkable.

Function	Interval	Area
$f(x) = x$	$[0, 1]$	$1/2$
$f(x) = x$	$[0, x]$	$x^2/2$
$f(x) = x^2$	$[0, 1]$	$1/3$
$f(x) = x^2$	$[0, x]$	$x^3/3$

Table 2.9: Areas under two simple curves, following a method inspired by Archimedes. The only “leap of faith” required so far is that terms of the form $\frac{\text{constant}}{n}$ and $\frac{\text{constant}}{n^2}$ become arbitrarily small (i.e., negligible) as n becomes very large, as illustrated in Fig. 2.4.



(a)

Figure 2.4: Graphs of $1/n$ and $1/n^2$ versus n , for $1 \leq n \leq 100$.

Observation 2.10: Ratios of Polynomials of n as n tends to Infinity

When approximating our areas with a sum of n rectangles, we were faced with terms of the form

- Triangle ($f(x) = x$): $\frac{1+2+\dots+n}{n^2} = \frac{an^2+bn}{n^2}$
- Parabola ($f(x) = x^2$): $\frac{1^2+2^2+\dots+n^2}{n^3} = \frac{an^3+bn^2+cn}{n^3}$

and we can imagine that for the area under x^k , $k > 2$, we'd be faced with a term of the form

$$\frac{1^k + 2^k + \dots + n^k}{n^{k+1}} = \frac{an^{k+1} + bn^k + \dots + \text{lower powers of } n}{n^{k+1}}.$$

In each case, we want to let n become large so that our scheme of under- and over-approximations through adding up rectangles becomes more accurate. In the limit, we'd like to take an “unbounded” number of “infinitely thin” rectangles, which brings us to the notion of “limits at infinity”, the topic of Chapters 2.5 and 2.6. Before we do that, let's address sums of powers of integers for arbitrary $k \geq 1$.

Proposition 2.11: Power Sums

For all $n \geq 1$ and $k \geq 1$, the following holds,

$$\begin{aligned} \sum_{i=1}^n i^k &:= 1^k + 2^k + 3^k + \dots + n^k \\ &= \frac{n^{k+1}}{k+1} + \frac{n^k}{2} + O(n, k-1), \end{aligned} \tag{2.23}$$

where $O(n, k-1)$ stands for a polynomial in the variable n with degree less than or equal to $k-1$ (i.e., all of the lower-order terms after n^k).

Remark 2.12. If you want to see the explicit values of the sums, here are the first few:

$$\begin{aligned} \sum_{i=1}^n i &= \frac{1}{2}(n^2 + n) = \frac{n^2}{2} + \frac{n}{2} \\ \sum_{i=1}^n i^2 &= \frac{1}{6}(2n^3 + 3n^2 + n) = \frac{n^3}{3} + \frac{n^2}{2} + O(n, 1) \\ \sum_{i=1}^n i^3 &= \frac{1}{4}(n^4 + 2n^3 + n^2) = \frac{n^4}{4} + \frac{n^3}{2} + O(n, 2) \\ \sum_{i=1}^n i^4 &= \frac{1}{30}(6n^5 + 15n^4 + 10n^3 - n) = \frac{n^5}{5} + \frac{n^4}{2} + O(n, 3) \\ \sum_{i=1}^n i^5 &= \frac{1}{12}(2n^6 + 6n^5 + 5n^4 - n^2) = \frac{n^6}{6} + \frac{n^5}{2} + O(n, 4). \end{aligned} \tag{2.24}$$

You can find longer lists online. A general formula was given by [Faulhaber](#).

2.5 Limits at Infinity: The Hard Way

Here, we take our first cut at the concept of a limit. We'll first lay out the concept using “everyday language”, and then we'll formalize it. To make sure we understand the formalization, the examples will be worked in “excruciating detail”. How's that for encouragement? In compensation, once we really understand what is going on, we'll be able to compute most limits “by inspection”.

Finite Limits at ∞

Conceptual Idea: Suppose $f : (0, \infty) \rightarrow \mathbb{R}$ is a function and $L \in \mathbb{R}$ is a (finite) real number. Then we say **the limit of $f(x)$ as x approaches infinity equals L** , denoted,

$$\lim_{x \rightarrow \infty} f(x) = L,$$

if we can make $f(x)$ stay “arbitrarily close” to L by taking x “sufficiently large in the positive direction.”

The formal definition makes precise what we mean by “arbitrarily close” and “sufficiently large.”

Definition 2.13. (*Finite Limit at Infinity*) Suppose $f : (0, \infty) \rightarrow \mathbb{R}$ is a function and $L \in \mathbb{R}$ is a (finite) real number. Then we say

$$\lim_{x \rightarrow \infty} f(x) = L$$

if, for all $\epsilon > 0$ (no matter how small), there exists $N < \infty$, such that, for all $x \geq N$, it holds that $|f(x) - L| \leq \epsilon$.

Notes:

- for all $x \geq N$ defines “ x sufficiently large”.
- $|f(x) - L| \leq \epsilon$ defines “arbitrarily close”. You can think of $\epsilon > 0$ as specifying a level of precision. How precisely is $f(x)$ approximated by the constant L for x sufficiently large.
- The definition can also be written as: for all $\epsilon > 0$, there exists $N < \infty$, such that, $x \geq N \implies |f(x) - L| \leq \epsilon$.
- It is emphasized that if there is a single point $x \geq N$ such that $|f(x) - L| > \epsilon$, then the “closeness (aka, precision)” criterion has failed. We have found a point x where $f(x)$ is not well approximated by the constant L .

Additional Remarks:

- We note that as we take ϵ smaller and smaller (aka, our precision requirement becomes finer and finer), we typically need to take N larger and larger. In other words, N typically depends on ϵ . It is sometimes useful to express this dependence as $N(\epsilon)$ or N_ϵ .
- Another way to interpret a limit at infinity is that as we take N larger and larger, the function restricted to $[N, \infty)$ becomes closer and closer to having a constant value equal to L .
- We will treat the definition of a limit as x approaches $-\infty$ later.

Example 2.14. Compute $\lim_{x \rightarrow \infty} f(x)$, if it exists, for the following functions.

(a) $f : (0, \infty) \rightarrow \mathbb{R}$ by $f(x) = \frac{1}{x}$.

(b) $f : (0, \infty) \rightarrow \mathbb{R}$ by $f(x) = \frac{1}{x^k}$.

(c) $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = \frac{x^2}{1+3x^2}$.

(d) $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = \frac{x}{1+3x^2}$.

Solutions:

(a) $f(x) = \frac{1}{x}$ for $x > 0$. **Ans.** $\lim_{x \rightarrow \infty} f(x) = 0.0$

To show this, set $L = 0.0$ and choose $\epsilon > 0$ but otherwise arbitrary. Then for all $x > \frac{1}{\epsilon}$, we have

$$|f(x) - L| = \left| \frac{1}{x} - 0.0 \right| = \frac{1}{x} < \frac{1}{\epsilon} = \epsilon,$$

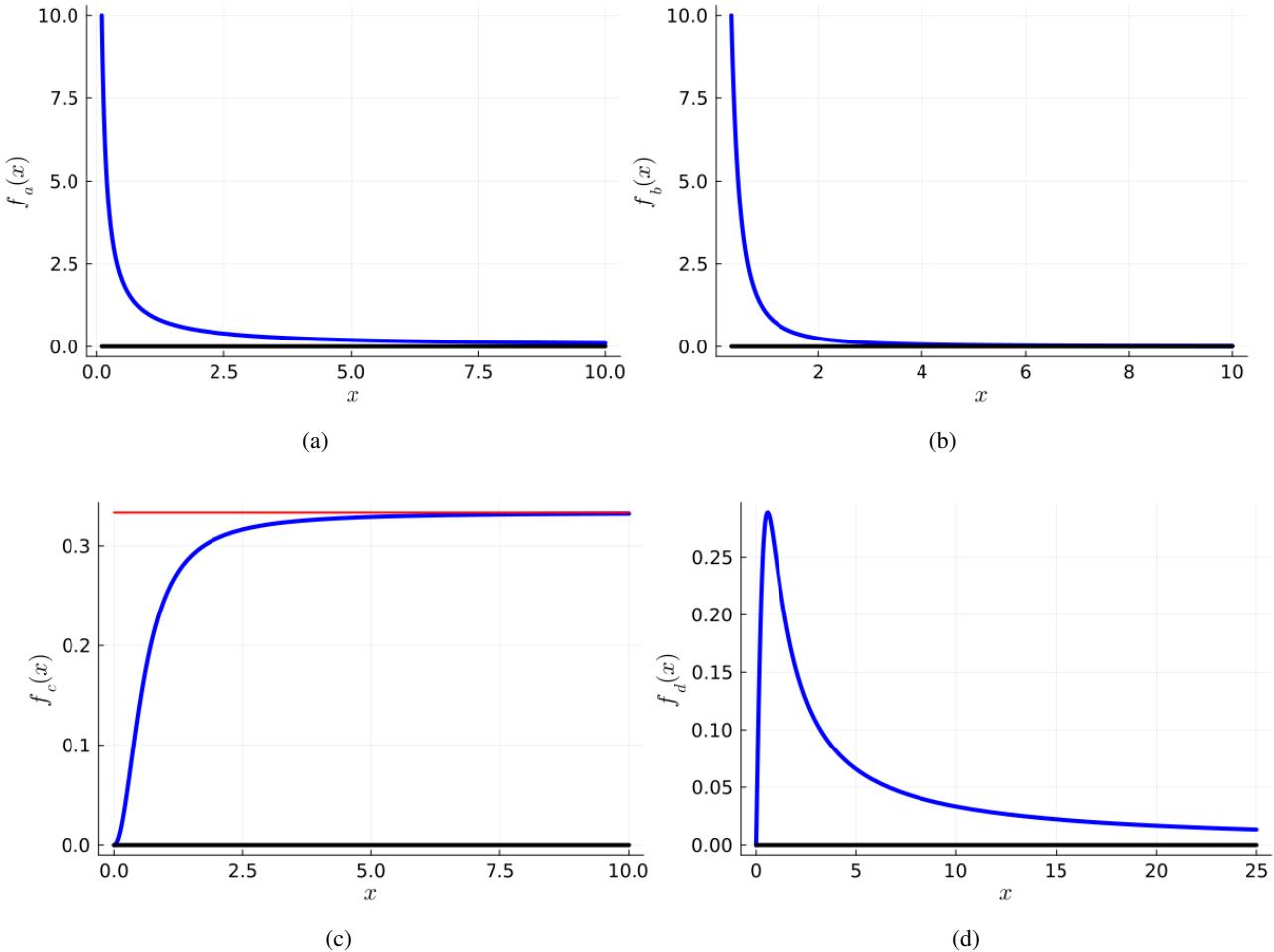


Figure 2.5: Graphs of the four functions in Example 2.14. The red line in (c) shows the non-zero value to which the function is converging.

where we used $\frac{1}{x} = \frac{1}{\text{a bigger number}} < \frac{1}{\text{a smaller number}} = \frac{1}{\epsilon}$, when $x > \frac{1}{\epsilon}$.

Note: Because $N(\epsilon) := \frac{1}{\epsilon}$, when ϵ becomes small to bound $|f(x) - L|$, we have at the same time $N(\epsilon) = \frac{1}{\epsilon}$ becoming large, imposing that $x > N$ is large.

$$(b) f(x) = \frac{1}{x^k} \text{ for } x > 0. \text{ Ans. } \lim_{x \rightarrow \infty} f(x) = 0.0.$$

To show this, set $L = 0.0$ and choose $\epsilon > 0$ but otherwise arbitrary. Then for all $x > \frac{1}{\sqrt[k]{\epsilon}}$, we have

$$|f(x) - L| = \left| \frac{1}{x^k} - 0.0 \right| = \frac{1}{x^k} < \frac{1}{\left(\frac{1}{\sqrt[k]{\epsilon}} \right)^k} = \frac{1}{\frac{1}{\epsilon}} = \epsilon.$$

Notes: The remarks from the previous problem apply here, with the exception that $N(\epsilon) := \frac{1}{\sqrt[k]{\epsilon}}$. It is still true that ϵ small forces $N(\epsilon)$ to be large.

$$(c) f : \mathbb{R} \rightarrow \mathbb{R} \text{ by } f(x) = \frac{x^2}{1+3x^2}. \text{ Ans. } \lim_{x \rightarrow \infty} f(x) = \frac{1}{3}.$$

To show this, set $L = 1/3$ and choose $\epsilon > 0$ but otherwise arbitrary. We rewrite the function as, for all $x \neq 0$,

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

Then,

$$|f(x) - L| = \left| \frac{1}{\frac{1}{x^2} + 3} - \frac{1}{3} \right| = \frac{1}{3} - \frac{1}{\frac{1}{x^2} + 3}.$$

Now, let's set $\delta = \frac{1}{x^2}$ and solve for $\delta > 0$ such that $\frac{1}{3} - \frac{1}{\frac{1}{\delta} + 3} < \epsilon$. Hence, after some algebra (nothing to do with Calculus), we work out

$$\frac{1}{3} - \frac{1}{\delta + 3} = \frac{\delta}{3(\delta + 3)} < \epsilon \iff \delta < 3\epsilon(\delta + 3) \iff (1 - 3\epsilon)\delta < 9\epsilon.$$

Because $(1 - 3\epsilon) < 1$ for all $\epsilon > 0$, we have that $(1 - 3\epsilon)\delta < \delta$, and

$$\delta < 9\epsilon \implies \left| \frac{1}{\frac{1}{\delta} + 3} - \frac{1}{3} \right| < \epsilon.$$

Hence, replacing δ by $\frac{1}{x^2}$, we have that $\frac{1}{x^2} < 9\epsilon \implies \left| \frac{1}{\frac{1}{x^2} + 3} - \frac{1}{3} \right| < \epsilon$, from which we deduce that

$$x > \frac{1}{3\sqrt{\epsilon}} \implies \left| \frac{1}{\frac{1}{x^2} + 3} - \frac{1}{3} \right| < \epsilon,$$

and hence we have shown that $\lim_{x \rightarrow \infty} f(x) = \frac{1}{3}$.

Notes:

- The above is what computing limits by hand typically looks like: **a blaze of symbols, multiple lines of computations, each of which is a likely opportunity for making errors.**
- While in this example, the messy algebra and bounding of terms was actually done back in Example 1.13, there is no way you would have remembered that!
- Bottom line:** Finding situations where we could write down the limits by inspection, while avoiding the messy algebra, would be nice. We will do that shortly.

(d) $f(x) = \frac{x}{1+3x^2}$ **Ans.** $\lim_{x \rightarrow \infty} f(x) = 0.0$.

To show this, set $L = 0.0$ and choose $\epsilon > 0$ but otherwise arbitrary. We rewrite the function as, for all $x \neq 0$,

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

Then, for all $x > 0$,

$$|f(x) - L| = \left| \frac{1}{\frac{1}{x} + 3x} - 0.0 \right| < \frac{1}{0+3x}.$$

We observe that

$$\frac{1}{3x} < \epsilon \iff 3x > \frac{1}{\epsilon} \iff x > \frac{1}{3\epsilon}.$$

Hence, here, $N(\epsilon) = \frac{1}{3\epsilon}$.

There is more than one way to make these estimates. For example, if we had instead written, for all $x > 0$,

$$f(x) = \frac{x}{1+3x^2} \cdot \frac{\frac{1}{x^2}}{\frac{1}{x^2}} = \frac{\frac{1}{x}}{\frac{1}{x^2} + 3},$$

then, for all $x > 1$,

$$|f(x) - L| = \left| \frac{\frac{1}{x}}{\frac{1}{x^2} + 3} - 0.0 \right| < \frac{\frac{1}{x}}{0+3} = \frac{1}{3x},$$

yielding the same result we had before. ■

Strategy for Bounding Ratios of Functions: To find an upper bound for a fraction, $\frac{n(x)}{d(x)}$, you replace the numerator with an upper bound and the denominator with a lower bound for the inequality to hold. In symbols, suppose that $0 \leq n(x) \leq N(x)$. Then,

$$\frac{n(x)}{d(x)} \leq \frac{N(x)}{d(x)},$$

because having a larger number in the numerator gives a larger fraction. Next, suppose that $0 < D(x) \leq d(x)$. Then,

$$\frac{n(x)}{d(x)} \leq \frac{N(x)}{d(x)} \leq \frac{N(x)}{D(x)},$$

because having a smaller number in the denominator yields a larger fraction. **While these are “obvious” facts, applying them successfully when you are first learning limits can be hard to master!**

Conversely, to find a lower bound for a fraction, $\frac{n(x)}{d(x)}$, you replace the numerator with a lower bound and the denominator with an upper bound for the inequality to hold.

In other words, if $n(x) \geq N(x)$ and $0 < d(x) \leq D(x)$, then $\frac{n(x)}{d(x)} \geq \frac{N(x)}{D(x)}$. This is useful for showing a limit is equal to infinity.

Infinite Limit at ∞

Conceptual Idea: A function $f : (0, \infty) \rightarrow \mathbb{R}$ converges to infinity as x tends to ∞ if for all finite values of L , “ $f(x)$ eventually equals or exceeds L and does not drop below L ”.

The formal definition makes precise what we mean by “eventually equals or exceeds L and does not drop below it.”

Definition 2.15. (Unbounded Limits at Infinity) The function $f : (0, \infty) \rightarrow \mathbb{R}$ has an **unbounded positive limit at infinity** if, for every finite $L > 0$, there exists $N < \infty$ such that,

$$x \geq N \implies f(x) \geq L,$$

that is, for all $x \in [N, \infty)$, $f(x) \geq L$. We emphasize that there cannot be any $x \in [N, \infty)$ such that $f(x) < L$; the function has to exceed or equal L and not drop below it.

The function $f : (0, \infty) \rightarrow \mathbb{R}$ has an **unbounded negative limit at infinity** if, for every finite $L < 0$, there exists $N < \infty$ such that,

$$x \geq N \implies f(x) \leq L,$$

that is, for all $x \in [N, \infty)$, $f(x) \leq L$. We emphasize that there cannot be any $x \in [N, \infty)$ such that $f(x) > L$; the function has to be less than or equal to L and not rise above it.

We denote these limits by $\lim_{x \rightarrow \infty} f(x) = \infty$ and $\lim_{x \rightarrow \infty} f(x) = -\infty$, respectively.

Note: For the definition of $\lim_{x \rightarrow \infty} f(x) = -\infty$, you may be worried that $L < 0$ does not “look big and negative”. However, $L = -10^9$ is big and negative. We allow L to be any negative real number.

Example 2.16. Compute $\lim_{x \rightarrow \infty} f(x)$, if it exists, for the following functions.

(a) $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = \frac{20x^3 + 11x^2 + 14}{1+x^2}$.

(b) $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = -x^3$.

(c) $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = x \cdot \sin(x)$.

$$(d) f : \mathbb{R} \rightarrow \mathbb{R} \text{ by } f(x) = \frac{-4x^3 + 11x + 14}{1+x^4}.$$

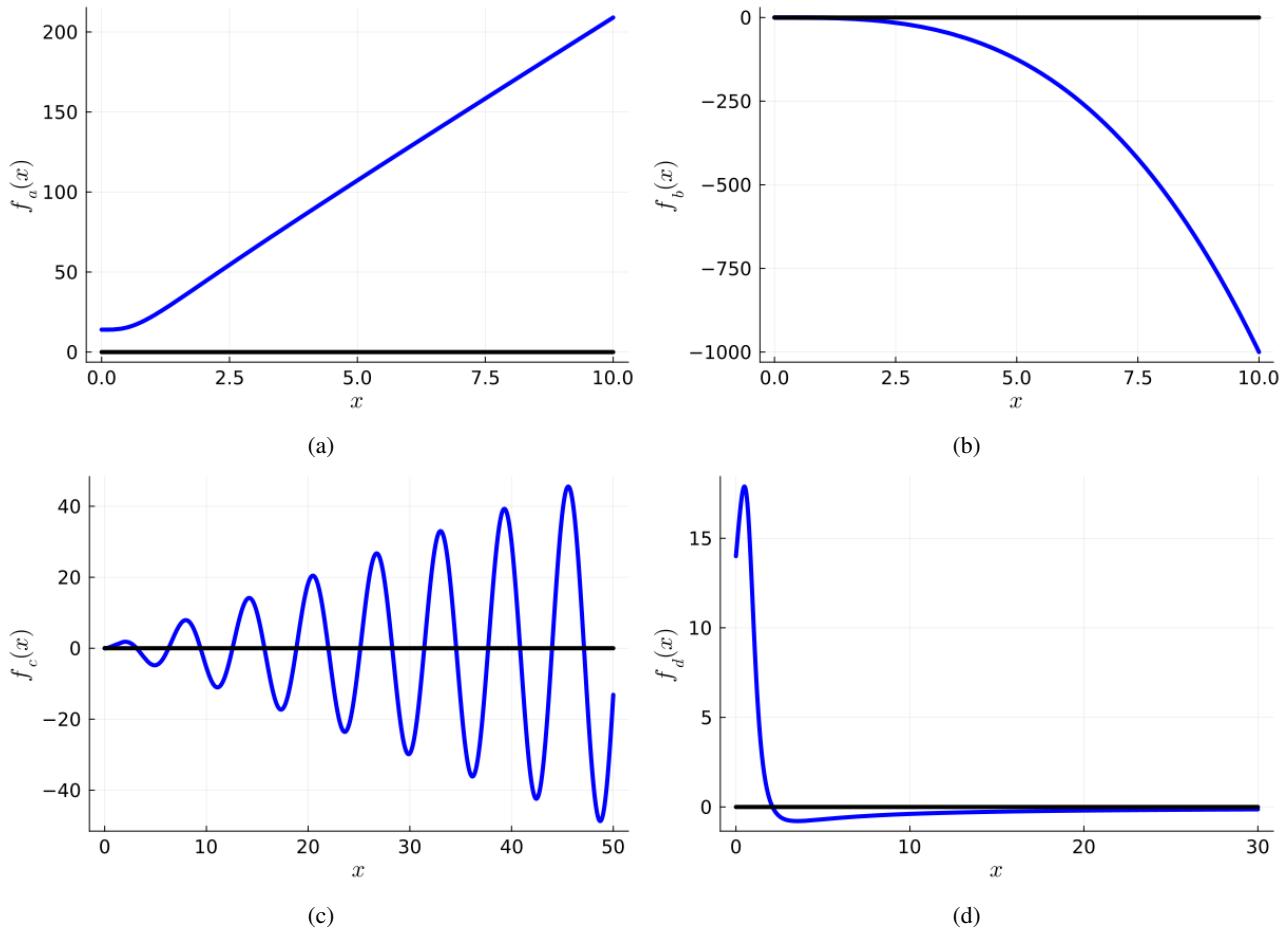


Figure 2.6: Graphs of the four functions in Example 2.16. The plot of $x \cdot \sin(x)$ given in (c) shows that the function becomes unbounded without converging to anything, not even infinity.

Solutions:

$$(a) f(x) = \frac{20x^3 + 11x^2 + 14}{1+x^2}. \text{ Ans. } \lim_{x \rightarrow \infty} f(x) = \infty.$$

To show this, we first rewrite the function as we have done several times before, for $x > 1$,

$$f(x) = \frac{20x^3 + 11x^2 + 14}{1+x^2} = \frac{20x^3 + 11x^2 + 14}{1+x^2} \cdot \frac{\frac{1}{x^2}}{\frac{1}{x^2}} = \frac{20x + 11 + \frac{14}{x^2}}{\frac{1}{x^2} + 1} \geq \frac{20x + 11}{1 + \frac{1}{x^2}} \geq \frac{20x + 11}{1 + 1} \geq 10x.$$

Next, let $L > 0$ be otherwise arbitrary and let $N \geq L/10$. Then, for all $x \geq N$, $f(x) \geq L$.

$$(b) f(x) = -x^3. \text{ Ans. } \lim_{x \rightarrow \infty} f(x) = -\infty.$$

To show this, let $L < -1$ be otherwise arbitrary and set $N = |L|$. Then $x > N \implies f(x) = -x^3 < -N^3 < -|L|^3 < L$. You could also have taken $N = \sqrt[3]{|L|}$ and made similar estimates. All that matters is that you show that you find a region of $[N, \infty) \subset \mathbb{R}$ such that $f(x) \leq L$ for $x \in [N, \infty)$.

$$(c) f(x) = x \cdot \sin(x). \text{ Ans. } \lim_{x \rightarrow \infty} f(x) \text{ does not exist.}$$

The issue is that $x \cdot \sin(x)$ becomes arbitrarily large in both the positive direction and the negative direction as $x \rightarrow \infty$. To show this, we first note that

$$\sin(x) = \begin{cases} +1 & x = \frac{\pi}{2} + 2k\pi, k \in \mathbb{N} \\ -1 & x = -\frac{\pi}{2} + 2k\pi, k \in \mathbb{N}, \end{cases}$$

where we used the $+1$ to emphasize the positive value. Hence, for arbitrary L and $N > 0$, taking $k = \max\{|L|, N\}$ results in $f(\frac{\pi}{2} + 2k\pi) = \frac{\pi}{2} + 2k\pi > L$ and $f(-\frac{\pi}{2} + 2k\pi) = \frac{\pi}{2} - 2k\pi < L$. We cannot make the function stay above L or below L for all $x > N$.

(d) $f(x) = \frac{-4x^3 + 11x + 14}{1+x^4}$. **Ans.** $\lim_{x \rightarrow \infty} f(x) = 0$.

To show this, we first rewrite the function as we have done several times before, for $x > 1$,

$$f(x) = \frac{-4x^3 + 11x + 14}{1+x^4} = \frac{-4x^3 + 11x + 14}{1+x^4} \cdot \frac{\frac{1}{x^4}}{\frac{1}{x^4}} = \frac{-4\frac{1}{x} + 11\frac{1}{x^3} + 14\frac{1}{x^4}}{\frac{1}{x^4} + 1}.$$

Next, let $L = 0$ and for all $x > 1$,

$$|f(x) - L| = |f(x) - 0| = \frac{|-4\frac{1}{x} + 11\frac{1}{x^3} + 14\frac{1}{x^4}|}{|\frac{1}{x^4} + 1|} \leq \frac{1}{x} \frac{|4 + 11 + 14|}{0 + 1} \leq 29 \frac{1}{x}.$$

Hence, for all $x > \frac{29}{\epsilon} > 1$, $|f(x) - 0| \leq \epsilon$, which shows that the limit is zero. ■

Neither Finite nor Infinite Limit at ∞ Exists

We need to work out what it means to fail a definition. If you are super good at logic, all you have to do is negate the definition! At this stage of your mathematics education, that is easier said than done, so we'll take it step by step. The main mental hurdle is that a definition typically requires several conditions to hold, and hence it fails if any ONE OF THEM does not hold. Logically, this arises from the negation of "and" being "or". For example, if we have **MyDef is True** if **P1** and **P2** (both) hold, then **MyDef is False** if either of **P1** or **P2** does not hold.

Conceptual Ideas: We say that $f : (0, \infty) \rightarrow \mathbb{R}$ **does not have a finite limit** as x tends to infinity if, for all candidate finite limits L , it is impossible to ensure that $f(x)$ remains near L by taking x sufficiently large. Similarly, $f(x)$ does not converge to infinity if for some finite $L > 0$, the function keeps "dipping" below L .

The formal anti-definitions clarify what "impossible to keep $f(x)$ near L " or "dipping below L " really mean!

Definition 2.17. (No Finite Limit at Infinity) Suppose $f : (0, \infty) \rightarrow \mathbb{R}$ is a function. Then $f(x)$ does not have a finite limit at infinity if for every choice of $L \in \mathbb{R}$, there exists some $\epsilon > 0$ such that one of the following hold:

- there does not exist $N < \infty$ so that $x > N \implies |f(x) - L| < \epsilon$.
- for all $N < \infty$ no matter how large, it is not true that $x > N \implies |f(x) - L| < \epsilon$.
- for all $N < \infty$ no matter how large, there (always) exists some $x > N$ such that $|f(x) - L| \geq \epsilon$.

The last three statements are all equivalent. We have simply "unwrapped" the negations one by one. They all capture the idea that "for every choice of $L \in \mathbb{R}$, we are unable to keep the function's values near L by taking x sufficiently large". The notion of " $f(x)$ is not near L " is codified by the existence of some positive value ϵ being used to bound the "distance of $f(x)$ from L ", via $|f(x) - L| \geq \epsilon$.

Definition 2.18. (No Infinite Limit at Infinity) Suppose $f : (0, \infty) \rightarrow \mathbb{R}$ is a function. Then $f(x)$ does not have an unbounded limit at infinity in the positive direction if for some choice of $L > 0$ and every choice of $0 < N_L < \infty$ (no matter how large), there exists $x \geq N_L$ such that $f(x) \leq L$. Similarly, $f(x)$ does not have an unbounded limit at infinity in the negative direction if for some choice of $L < 0$ and every choice of $0 < N_L < \infty$ (no matter how large), there exists $x \geq N_L$ such that $f(x) \geq L$.

Classic Example is shown in Fig. 2.6-(c). Yes, the function grows without bound, but it fails to stay large in the positive direction (or alternatively, the negative direction), as do the curves in Fig. 2.6-(a) and -(b).

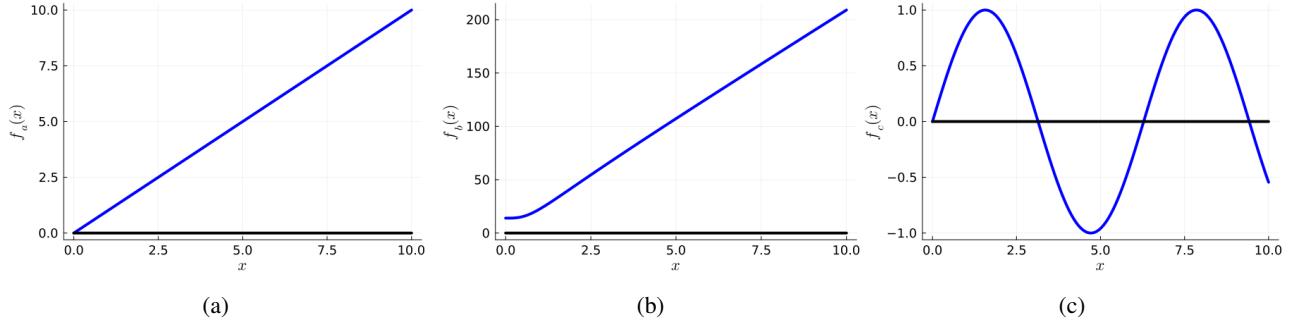


Figure 2.7: Graphs of the three functions in Example 2.19.

Example 2.19. Compute $\lim_{x \rightarrow \infty} f(x)$, if a finite limit exists for the following functions.

- (a) $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = x$.
- (b) $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = \frac{20x^3 + 11x + 14}{1 + x^2}$.
- (c) $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = \sin(x)$.

Solutions:

(a) $f(x) = x$ does not have a finite limit as x becomes arbitrarily large. While this seems patently obvious, how would one actually show it?

- Solution 1: We use Condition 1 of our “anti-definition”, namely, for every choice of $L \in \mathbb{R}$, and for every choice of $N < \infty$, there is always an $x > N$ such that $f(x) > L$. This is easy to show because we can let $L \in \mathbb{R}$ and $0 < N < \infty$ be arbitrary. Then for $x > \max\{L, N\}$, $f(x) = x > L$.
- Solution 2: We use Condition 4 of our “anti-definition”, namely, for every choice of $L \in \mathbb{R}$, there exists some $\epsilon > 0$ such that, for all $N < \infty$ no matter how large, there (always) exists some $x > N$ such that $|f(x) - L| \geq \epsilon$. To show this, let L be arbitrary and choose $\epsilon = 1$. Then, for all $x > L + \epsilon$,

$$|x - L| = x - L > L + \epsilon - L = \epsilon.$$

Hence, no matter what L we choose, and for all N finite, we have $x > \max\{N, L + \epsilon\} \implies |f(x) - L| \geq \epsilon$.

Note: Yes, it seems clear that $\lim_{x \rightarrow \infty} x = \infty$ and we will formalize that shortly.

(b) $f(x) = \frac{20x^3 + 11x + 14}{1 + x^2}$ does not have a finite limit as x becomes arbitrarily large. We observe that the numerator is growing at “rate” x^3 , while the denominator is growing at rate x^2 . Hence, the numerator should “win”. How do we formalize this? Let’s rewrite, for all $x \neq 0$,

$$\frac{20x^3 + 11x + 14}{1 + x^2} = \frac{20x^3 + 11x + 14}{1 + x^2} \cdot \frac{\frac{1}{x^2}}{\frac{1}{x^2}} = \frac{20x + \frac{11}{x} + \frac{14}{x^2}}{\frac{1}{x^2} + 1}.$$

Then, for all $x > 1$,

$$f(x) = \frac{20x + \frac{11}{x} + \frac{14}{x^2}}{\frac{1}{x^2} + 1} > \frac{20x}{1 + 1} \geq 10x,$$

which grows unbounded as x tends to infinity, meaning, there is no finite limit.

- (c) $f(x) = \sin(x)$ has no limit as $x \rightarrow \infty$ because no matter how large N becomes, $\sin : [N, \infty) \rightarrow \mathbb{R}$ still oscillates between -1 and 1 . While that really is “enough said”, let’s be stubborn and spell out all the details.
- Case 1: The candidate limit $L \notin [-1, 1]$ and hence $|L| > 1$. Let $\delta = |L| - 1 > 0$. For all $y \in [-1, 1]$, it follows that $|y| \leq 1$ and hence $||L| - |y|| \geq \delta$. By the reverse triangle inequality in Prop. 1.9, we have $|L - y| \geq \delta$. Hence, for all $\epsilon \leq \delta$ and $x \in \mathbb{R}$, $|L - \sin(x)| = |\sin(x) - L| > \delta \geq \epsilon$, and thus L cannot be a limit.
 - Case 2: The candidate limit $L = 0$. Then for $\epsilon = 0.5$ and for all $N < \infty$, it follows that $x = N^{\frac{\pi}{2}}$ satisfies $x > N$ and $|\sin(x) - L| = |1 - 0| = 1 > 0.5 = \epsilon$, and hence L is not a limit.
 - Case 3: The candidate limit $L \in [-1, 1]$ and $L \neq 0$. Then there exists $x \in [-\pi, \pi]$ such that $f(x) = L$ and $f(-x) = -L$. Moreover, because $\sin(-x + 2k\pi) = \sin(-x) = -\sin(x)$, we have that $f(-x + 2k\pi) = -L$. Hence, we can take $\epsilon = |L| > 0$ and note that for all $1 < N < \infty$, $-x + 2N\pi > N$ and $|f(-x + 2N\pi) - L| = |\sin(-x + 2N\pi) - L| = |\sin(-x) - L| = 2|L| > \epsilon$.

■

Secrets of the Arcane 2.20: Why Do We Care About Limits Anyway?

Limits seem hard. Can’t Calculus exist without them?

- Limits are the Building Blocks of Calculus:** Calculus without limits is like a castle floating in the air. Limits are the bedrock of this mathematical fortress.
- Infinity Demystified:** Limits turn the concept of infinity from a perplexing puzzle into a manageable marvel, allowing us to reason about the infinitely small and large.
- Solving the Unsolvable:** Limits provide a method to solve problems that might otherwise be unsolvable. For example, they can be used to determine the value of a function at a point where it’s undefined. To see what this might mean, what values should we assign at ...

- $x = 1$ for $\frac{x^3 - 6x^2 + 11x - 6}{x-1}$ (were we not taught to NEVER divide by zero)?
- $t = 0$ for $\frac{\sin(t)}{t}$ (shows up in Signal Processing when building low-pass filters)?
- $h = 0$ for $\frac{e^h - 1}{h}$ (rate of change of the natural exponential at the origin)?
- $n \rightarrow \infty$ for $\frac{n^2 + n}{2n^2}$ (area under a unit parabola between zero and one)?
- $t \rightarrow \infty$ for $t \cdot e^{-t} \cdot \sin(3t)$ (infinity times zero times an oscillation)? Shows up in feedback control!

2.6 Limits at Infinity: The Easy Way

Now that we’ve computed limits by the book, we can make a few observations and profit from them.

2.6.1 Rational Functions

In Examples 2.14, 2.19, and 2.16, we had many **rational functions**, that is, functions consisting of the ratio of two polynomials. Recall that the **degree of a polynomial (aka, “order” of the polynomial)** is the highest power of the variable in the polynomial expression. The coefficient multiplying the term of the highest degree is called the **leading coefficient**; by definition, it cannot be zero. In all of our examples, the nature of the limit of a rational function could be determined by degrees of the numerator and denominator polynomials, and their corresponding leading coefficients. **Goodbye, tedious calculations!**

Intuition for the limit of a rational function of x : With a little algebra, we can rewrite a general rational function of x as:

$$\frac{b_m x^m + b_{m-1} x^{m-1} + \dots + b_0}{a_n x^n + a_{n-1} x^{n-1} + \dots + a_0} = \left[\frac{x^m}{x^n} \right] \cdot \frac{b_m + b_{m-1} \frac{1}{x} + \dots + b_0 \frac{1}{x^m}}{a_n + a_{n-1} \frac{1}{x} + \dots + a_0 \frac{1}{x^n}}.$$

Because all terms of the form $\frac{1}{x^k}$ go to zero as $|x|$ tends to infinity, the “messy fraction” on the right has a simple limit, namely,

$$\lim_{x \rightarrow \pm\infty} \frac{b_m + b_{m-1}\frac{1}{x} + \cdots + b_0\frac{1}{x^m}}{a_n + a_{n-1}\frac{1}{x} + \cdots + a_0\frac{1}{x^n}} = \boxed{\frac{b_m}{a_n}}.$$

Hopefully, you see it. If not, play with the limit in Julia. Once we have replaced the messy term,

$$\frac{b_m + b_{m-1}\frac{1}{x} + \cdots + b_0\frac{1}{x^m}}{a_n + a_{n-1}\frac{1}{x} + \cdots + a_0\frac{1}{x^n}},$$

with the constant $\frac{b_m}{a_n}$, all we need to do is understand

$$\lim_{x \rightarrow \pm\infty} \frac{b_m}{a_n} \cdot \frac{x^m}{x^n} = \frac{b_m}{a_n} \cdot \lim_{x \rightarrow \pm\infty} \frac{x^m}{x^n},$$

which seems “infinitely easier”! **This is Prop. 2.21 in a nutshell.**

Proposition 2.21: Limit of a Rational Function as x tends to Plus or Minus Infinity

Consider a rational function

$$\frac{b_m x^m + b_{m-1} x^{m-1} + \cdots + b_0}{a_n x^n + a_{n-1} x^{n-1} + \cdots + a_0},$$

where b_m and a_n are non-zero. Then the limit as x tends to plus or minus infinity is equal to

$$\lim_{x \rightarrow \pm\infty} \frac{b_m x^m + b_{m-1} x^{m-1} + \cdots + b_0}{a_n x^n + a_{n-1} x^{n-1} + \cdots + a_0} = \frac{b_m}{a_n} \cdot \lim_{x \rightarrow \pm\infty} x^{m-n}. \quad (2.25)$$

In particular,

$$\lim_{x \rightarrow \infty} \frac{b_m x^m + b_{m-1} x^{m-1} + \cdots + b_0}{a_n x^n + a_{n-1} x^{n-1} + \cdots + a_0} = \begin{cases} \text{sign}\left(\frac{b_m}{a_n}\right) \cdot \infty & m > n \\ \frac{b_m}{a_n} & m = n \\ 0 & m < n. \end{cases} \quad (2.26)$$

and

$$\lim_{x \rightarrow -\infty} \frac{b_m x^m + b_{m-1} x^{m-1} + \cdots + b_0}{a_n x^n + a_{n-1} x^{n-1} + \cdots + a_0} = \begin{cases} \text{sign}\left(\frac{b_m}{a_n}\right) \cdot (-1)^{m-n} \cdot \infty & m > n \\ \frac{b_m}{a_n} & m = n \\ 0 & m < n. \end{cases} \quad (2.27)$$

Note: If one of m and n is odd and the other even, then $(-1)^{m-n} = -1$; otherwise, it equals $+1$. The easiest thing to remember is (2.25). It states that the **ratio of the highest degree terms in the numerator and denominator determines the limit**.

$$\lim_{x \rightarrow \pm\infty} \frac{b_m x^m + b_{m-1} x^{m-1} + \cdots + b_1 x + b_0}{a_n x^n + a_{n-1} x^{n-1} + \cdots + a_1 x + a_0} = \lim_{x \rightarrow \pm\infty} \frac{x^m \cdot (b_m + b_{m-1}\frac{1}{x} + \cdots + b_1\frac{1}{x^{m-1}} + b_0\frac{1}{x^m})}{x^n \cdot (a_n + a_{n-1}\frac{1}{x} + \cdots + a_1\frac{1}{x^{n-1}} + a_0\frac{1}{x^n})} = \lim_{x \rightarrow \pm\infty} \frac{b_m x^m}{a_n x^n}.$$

You may enjoy the video by [Nancy Pi, MIT graduate](#), “How to Find the Limit at Infinity.” She earns great comments, such as “I usually never comment on YouTube videos but I just had to say how much this has helped me. THANK YOU SO MUCH, and never stop making math videos!! :)" If you want even more practice problems, please check out the video by [Jennifer McCracken](#).

Example 2.22. Where applicable, use Prop. 2.21 to compute $\lim_{x \rightarrow \pm\infty} f(x)$ for the following functions.

(a) $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = \frac{20x^2+11x+14}{1+4x^2}$.

(b) $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = -x^3$.

(c) $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = x \cdot \sin(x)$.

(d) $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = \frac{-4x^3+11x+14}{1+x^4}$.

(e) $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = \frac{4x^4+11x+14}{1+2x^3}$.

Solutions:

(a) $f(x) = \frac{20x^2+11x^2+14}{1+4x^2}$. **Ans.** $\lim_{x \rightarrow \infty} f(x) = 5$ and $\lim_{x \rightarrow -\infty} f(x) = 5$.

$$\bullet \lim_{x \rightarrow \infty} f(x) = \frac{20}{4} \cdot \lim_{x \rightarrow \infty} x^{2-2} = 5.$$

$$\bullet \lim_{x \rightarrow -\infty} f(x) = \frac{20}{4} \cdot \lim_{x \rightarrow -\infty} x^{2-2} = 5.$$

(b) $f(x) = -x^3$. **Ans.** $\lim_{x \rightarrow \infty} f(x) = -\infty$ and $\lim_{x \rightarrow -\infty} f(x) = +\infty$.

$$\bullet \lim_{x \rightarrow \infty} f(x) = \frac{-1}{1} \cdot \lim_{x \rightarrow \infty} x^{3-0} = -\infty.$$

$$\bullet \lim_{x \rightarrow -\infty} f(x) = \frac{-1}{1} \cdot \lim_{x \rightarrow -\infty} x^{3-0} = +\infty.$$

(c) $f(x) = x \cdot \sin(x)$. **Ans.** Prop. 2.21 is not applicable because $f(x)$ is not a rational function.

(d) $f(x) = \frac{-4x^3+11x+14}{1+x^4}$. **Ans.** $\lim_{x \rightarrow \infty} f(x) = 0$ and $\lim_{x \rightarrow -\infty} f(x) = 0$.

$$\bullet \lim_{x \rightarrow \infty} f(x) = \frac{-4}{1} \cdot \lim_{x \rightarrow \infty} x^{3-4} = \frac{-4}{1} \cdot \lim_{x \rightarrow \infty} \frac{1}{x} = 0.$$

$$\bullet \lim_{x \rightarrow -\infty} f(x) = \frac{-4}{1} \cdot \lim_{x \rightarrow -\infty} x^{3-4} = \frac{-4}{1} \cdot \lim_{x \rightarrow -\infty} \frac{1}{x} = 0.$$

(e) $f(x) = \frac{4x^4+11x+14}{1+2x^3}$. **Ans.** $\lim_{x \rightarrow \infty} f(x) = +\infty$ and $\lim_{x \rightarrow -\infty} f(x) = -\infty$.

$$\bullet \lim_{x \rightarrow \infty} f(x) = \frac{4}{2} \cdot \lim_{x \rightarrow \infty} x^{4-3} = \frac{4}{2} \cdot \lim_{x \rightarrow \infty} x = +\infty.$$

$$\bullet \lim_{x \rightarrow -\infty} f(x) = \frac{4}{2} \cdot \lim_{x \rightarrow -\infty} x^{4-3} = \frac{4}{2} \cdot \lim_{x \rightarrow -\infty} x = -\infty.$$

Sometimes, a change of variable can turn a hard problem into an easy one. We illustrate this next.

Example 2.23. Evaluate² $\lim_{x \rightarrow \infty} \frac{2+x}{x^{1/3}-4}$.

Solution: We introduce a change of variable $y := x^{1/3}$. We do this because then

$$\frac{2+x}{x^{1/3}-4} = \frac{2+y^3}{y-4},$$

a simple rational function of y . We note that $x \rightarrow \infty \iff y \rightarrow \infty$, and hence,

$$\lim_{x \rightarrow \infty} \frac{2+x}{x^{1/3}-4} = \lim_{y \rightarrow \infty} \frac{2+y^3}{y-4} = \infty.$$

Transformations/substitutions/changes of variables are used a lot in Calculus to simplify problems. ■

²The problem is motivated by All Things Mathematics. Here, we take the limit as x tends to infinity. The YouTube video is looking at finite limits, which we treat in Chapter 4.

2.6.2 Limit as an Integer Tends to Infinity

Our original motivation for studying limits in the first place, our area computations in Chapter 2.4.2, resulted in functions that depended on a counting number, $n \in \mathbb{N}$. All of the definitions and results of Chapter 2.5 and 2.6.1, on limits for a real variable x , apply equally well to limits as an integer $n \in \mathbb{Z}$ tends to plus or minus infinity. We summarize them briefly here.

Taking limits as $n \rightarrow \pm\infty$ is no different than taking limits as $x \rightarrow \pm\infty$. While we formalize this next, all we are doing is changing variable names to protect the guilty. If you are good with limits as $x \rightarrow \pm\infty$, you can skim a few of the examples and then move on.

Fact 2.24: Limits of Functions of Integers

Definition 2.25. (Limits at Infinity) Suppose $f : \mathbb{Z} \rightarrow \mathbb{R}$ is a function and $L \in \mathbb{R}$ is a (finite) real number. Then

- (a) $\lim_{n \rightarrow \infty} f(n) = L$ if, for all $\epsilon > 0$ (no matter how small), there exists $N < \infty$, such that, for all $n \geq N$, it holds that $|f(n) - L| \leq \epsilon$.
- (b) $\lim_{n \rightarrow -\infty} f(n) = L$ if, for all $\epsilon > 0$ (no matter how small), there exists $N > -\infty$, such that, for all $n \leq N$, it holds that $|f(n) - L| \leq \epsilon$.
- (c) $\lim_{n \rightarrow \infty} f(n) = \infty$ if, for all $M > 0$ (no matter how large), there exists $N < \infty$, such that, for all $n \geq N$, it holds that $f(n) > M$.
- (d) $\lim_{n \rightarrow -\infty} f(n) = \infty$ if, for all $M > 0$ (no matter how large), there exists $N > -\infty$, such that, for all $n \leq N$, it holds that $f(n) > M$.
- (e) $\lim_{n \rightarrow \infty} f(n) = -\infty$ if, for all $M < 0$ (no matter how large in the negative direction), there exists $N < \infty$, such that, for all $n \geq N$, it holds that $f(n) < M$.
- (f) $\lim_{n \rightarrow -\infty} f(n) = -\infty$ if, for all $M < 0$ (no matter how large in the negative direction), there exists $N > -\infty$, such that, for all $n \leq N$, it holds that $f(n) < M$.

Note: Just as with limits of a real variable x tending to infinity, we look at the behavior of the function on sets of the form $\{n \in \mathbb{N} \mid n \geq N\}$ or $\{n \in \mathbb{Z} \mid n \leq N\}$. We want to know whether the function becomes essentially constant as N becomes arbitrarily large in the positive (or, negative) direction, giving us a finite limit, or can we make the function always stay above (or, below) an arbitrarily large positive (or, negative) number, giving us an infinite limit. Of course, sometimes neither of these conditions hold, in which case a limit as n tends to infinity does not exist. Rational functions will always have either a finite limit or an infinite limit. You can play with $\sin(n)$ as $n \rightarrow \infty$ to see that not all functions have limits.

Example 2.26. Where applicable, use Fact 2.24 and Prop. 2.21 to compute $\lim_{n \rightarrow \infty} f(n)$ for the following functions.

- (a) $f : \mathbb{Z} \rightarrow \mathbb{R}$ by $f(n) = \frac{1}{3} + \frac{1}{6n^2}$.
- (b) $f : \mathbb{Z} \rightarrow \mathbb{R}$ by $f(n) = \frac{1}{n^2} \cdot \frac{n(n+1)}{2}$.
- (c) $f : \mathbb{Z} \rightarrow \mathbb{R}$ by $f(n) = \frac{x^3}{n^3} \cdot \frac{2n^3 - 3n^2 + n}{6}$, where x is a constant not depending upon n .
- (d) $f : \mathbb{Z} \rightarrow \mathbb{R}$ by $f(n) = \frac{1}{n^{k+1}} \cdot \sum_{i=1}^n i^k$, for $k \geq 1$.

Solutions:

(a) $f(n) = \frac{1}{3} + \frac{1}{6n^2}$ **Ans.** $\lim_{n \rightarrow \infty} f(n) = \frac{1}{3}$.

In principle, to apply Prop. 2.21, we should rewrite $f(n) = \frac{2n^2+1}{6n^2}$ and then compute $\lim_{n \rightarrow \infty} f(n) = \frac{2}{6} = \frac{1}{3}$. However, you probably already have guessed that the finite limit of a sum is the sum of the finite limits, and hence we can also compute the result as $\lim_{n \rightarrow \infty} f(x) = \frac{1}{3} + \lim_{n \rightarrow \infty} \frac{1}{6n^2} = \frac{1}{3} + 0$.

(b) $f(n) = \frac{1}{n^2} \cdot \frac{n(n+1)}{2}$. **Ans.** $\lim_{n \rightarrow \infty} f(n) = \frac{1}{2}$.

$$f(n) = \frac{n^2+n}{2n^2} \text{ and hence } \lim_{n \rightarrow \infty} f(n) = \frac{1}{2}.$$

(c) $f(n) = \frac{x^3}{n^3} \cdot \frac{2n^3-3n^2+n}{6}$. **Ans.** $\lim_{n \rightarrow \infty} f(n) = \frac{1}{k+1}$.

$$f(n) = \frac{x^3}{n^3} \cdot \frac{2n^3-3n^2+n}{6} = \frac{2x^3 n^3 - 3x^3 n^2 + x^3 n}{6n^3} \text{ and hence } \lim_{n \rightarrow \infty} f(x) = \frac{2x^3}{6} = \frac{x^3}{3}.$$

(d) $f(n) = \frac{1}{n^{k+1}} \cdot \sum_{i=1}^n i^k$, for $k \geq 1$. **Ans.** $\lim_{n \rightarrow \infty} f(x) = \frac{1}{k+1}$.

$f(n) = \frac{1}{n^{k+1}} \cdot \sum_{i=1}^n i^k = \frac{1}{n^{k+1}} \cdot \left(\frac{n^{k+1}}{k+1} + \frac{n^k}{2} + O(n, k-1) \right)$ by Prop. 2.11. Given what we know so far, we should work the problem as follows: $O(n, k-1)$ is some polynomial of degree $(k-1)$ and hence, $O(n, k-1) = a_{k-1}n^{k-1} + a_{k-2}n^{k-2} + \dots + a_1n + a_0$. Hence, we have that

$$f(n) = \frac{\frac{n^{k+1}}{k+1} + \frac{n^k}{2} + a_{k-1}n^{k-1} + a_{k-2}n^{k-2} + \dots + a_1n + a_0}{n^{k+1}},$$

and thus, $\lim_{n \rightarrow \infty} f(x) = \frac{\frac{1}{k+1}}{1} = \frac{1}{k+1}$.

Note: Once again, you are probably thinking that we should be able to take limits term-by-term and not have to rewrite everything as one giant rational function before taking the limit. **You are mostly correct**, as we will clarify in Chapter 2.6.5. ■

Example 2.27. Compute the following limits.

(a) $\lim_{N \rightarrow \infty} \frac{(x + \frac{1}{N})^2 - x^2}{\frac{1}{N}}$.

(b) $\lim_{N \rightarrow \infty} \frac{(x + \frac{1}{N})^3 - x^3}{\frac{1}{N}}$.

(c) $\lim_{N \rightarrow \infty} \frac{(x + \frac{1}{N})^n - x^n}{\frac{1}{N}}$, for $n \in \mathbb{N}, n \geq 3$.

(d) $\lim_{N \rightarrow \infty} A_N^{\text{Low}}$ for the strictly increasing function $f(y) = y^4$ over the interval $[0, x]$, where $x > 0$; see Chapter 2.4.2.

(e) $\lim_{N \rightarrow \infty} A_N^{\text{Up}}$ for the strictly increasing function $f(y) = y^4$ over the interval $[0, x]$, where $x > 0$; see Chapter 2.4.2.

Solution:

(a) **Ans.** $\lim_{N \rightarrow \infty} \frac{(x + \frac{1}{N})^2 - x^2}{\frac{1}{N}} = 2x$.

$$\begin{aligned} \frac{(x + \frac{1}{N})^2 - x^2}{\frac{1}{N}} &= \frac{(x^2 + \frac{2x}{N} + \frac{1}{N^2}) - x^2}{\frac{1}{N}} \\ &= N \cdot \left(\frac{2x}{N} + \frac{1}{N^2} \right) \\ &= \frac{2xN + 1}{N}. \end{aligned}$$

Hence,

$$\begin{aligned} \lim_{N \rightarrow \infty} \frac{(x + \frac{1}{N})^2 - x^2}{\frac{1}{N}} &= \lim_{N \rightarrow \infty} \frac{2xN + 1}{N} \\ &= 2x. \end{aligned}$$

Note: $2xN$ is the leading term of the numerator if, and only if, $x \neq 0$. Hence, to be 100% totally correct, we really have two cases to consider here: (i) $x \neq 0$ and (ii) $x = 0$. In Case (i), we have already computed the answer to be $2x$. In Case (ii), the leading coefficient of the numerator is 1, and we have $\lim_{N \rightarrow \infty} \frac{1}{N} = 0$. Because $2x = 0$ when $x = 0$, being “super careful” versus “not so careful” has not changed the answer. In the following, if you care to do so, you can check that being super careful does not change the answers we give below.

(b) **Ans.** $\lim_{N \rightarrow \infty} \frac{(x + \frac{1}{N})^3 - x^3}{\frac{1}{N}} = 3x^2$.

From Table 1.82, we have

$$\begin{aligned} \frac{(x + \frac{1}{N})^3 - x^3}{\frac{1}{N}} &= \frac{(x^3 + 3\frac{x^2}{N} + 3\frac{x}{N^2} + \frac{1}{N^3}) - x^3}{\frac{1}{N}} \\ &= N \cdot \left(3\frac{x^2}{N} + 3\frac{x}{N^2} + \frac{1}{N^3} \right) \\ &= \frac{3x^2N^2 + 3xN + 1}{N^2}. \end{aligned}$$

Hence,

$$\begin{aligned} \lim_{N \rightarrow \infty} \frac{(x + \frac{1}{N})^3 - x^3}{\frac{1}{N}} &= \lim_{N \rightarrow \infty} \frac{3x^2N^2 + 3xN + 1}{N^2} \\ &= 3x^2. \end{aligned}$$

(c) **Ans.** $\lim_{N \rightarrow \infty} \frac{(x + \frac{1}{N})^n - x^n}{\frac{1}{N}} = nx^{n-1}$.

For this problem, we need the full-blown Binomial Theorem 1.83. So that we do not lose track of the term, $-x^n$, we move it to the front and compute

$$\begin{aligned} -x^n + (x + \frac{1}{N})^n &= -x^n + \binom{n}{0}x^n \cdot \left(\frac{1}{N}\right)^0 + \binom{n}{1}x^{n-1} \cdot \left(\frac{1}{N}\right)^1 + \sum_{k=2}^n \binom{n}{k}x^{n-k} \cdot \left(\frac{1}{N}\right)^k \\ &= -x^n + x^n + nx^{n-1} \cdot \left(\frac{1}{N}\right)^1 + \sum_{k=2}^n \binom{n}{k}x^{n-k} \cdot \left(\frac{1}{N}\right)^k \\ &= nx^{n-1} \cdot N^{-1} + \sum_{k=2}^n \binom{n}{k}x^{n-k} \cdot N^{-k} \\ &= nx^{n-1} \cdot N^{-1} + \sum_{k=2}^{n-1} \binom{n}{k}x^{n-k} \cdot N^{-k} + x^0 \cdot N^{-n} \text{ multiply by } \frac{N^n}{N^n} \\ &= \frac{nx^{n-1}N^{n-1} + \sum_{k=2}^{n-1} \binom{n}{k}x^{n-k} \cdot N^{n-k} + 1}{N^n}. \end{aligned}$$

Hence,

$$\begin{aligned} \lim_{N \rightarrow \infty} \frac{\left(x + \frac{1}{N}\right)^n - x^n}{\frac{1}{N}} &= \lim_{N \rightarrow \infty} N \cdot \frac{nx^{n-1}N^{n-1} + \sum_{k=2}^{n-1} \binom{n}{k}x^{n-k} \cdot N^{n-k} + 1}{N^n} \\ &= \lim_{N \rightarrow \infty} \frac{nx^{n-1}N^{n-1} + \sum_{k=2}^{n-1} \binom{n}{k}x^{n-k} \cdot N^{n-k} + 1}{N^{n-1}} \\ &= nx^{n-1}. \end{aligned}$$

(d) **Ans.** $\lim_{N \rightarrow \infty} A_N^{\text{Low}} = \lim_{N \rightarrow \infty} \frac{x^5(6(N-1)^5 + 15(N-1)^4 + 10(N-1)^3 - (N-1))}{30N^5} = \frac{x^5}{5}.$

For $N \geq 2$, define $\Delta x = \frac{x}{N}$, $x_i = (i-1)\Delta x$, $f(x_i) = ((i-1)\Delta x)^4$,

$$\begin{aligned} \text{Area}_N^{\text{Low}} &:= \sum_{i=1}^N f(x_i) \cdot \Delta x \\ &= \sum_{i=1}^N ((i-1)\Delta x)^4 \cdot \Delta x \\ &= \sum_{i=1}^N \left((i-1)\frac{x}{N}\right)^4 \cdot \frac{x}{N} \\ &= \frac{x^5}{N^5} \cdot \sum_{i=1}^N (i-1)^4 \quad (\text{because both } n, x \text{ are independent of } i, \text{ we take them out of the sum}) \\ &= \frac{x^5}{N^5} \cdot \sum_{i=1}^{N-1} i^4 \quad (\text{change of index}) \\ &= \frac{x^5}{N^5} \cdot \frac{6(N-1)^5 + 15(N-1)^4 + 10(N-1)^3 - (N-1)}{30} \quad (\text{by Remark 2.12}) \\ &= \frac{x^5(6(N-1)^5 + 15(N-1)^4 + 10(N-1)^3 - (N-1))}{30N^5}. \end{aligned}$$

Taking the limit as $N \rightarrow \infty$ completes the result, or does it? Do we need to expand $(N-1)^5$, $(N-1)^4$, and $(N-1)^3$ using the Binomial Theorem so as to identify the leading coefficient carefully? The answer is that with a bit of experience, NO, we do not need to do that. Why? The leading coefficient will come from $(N-1)^5$, so we can ignore the other two. Moreover, from the Binomial Theorem,

$$\begin{aligned} (N-1)^5 &= \binom{5}{0}N^5 + \text{terms of degree 4 and lower in } N \\ &= N^5 + O(N, 4). \end{aligned}$$

Hence, the limit is equal to $\lim_{N \rightarrow \infty} \frac{6x^5N^5}{30N^5} = \frac{x^5}{5}$.

(e) **Ans.** $\lim_{N \rightarrow \infty} A_N^{\text{Up}} = \lim_{N \rightarrow \infty} \frac{x^5(6N^5 + 15N^4 + 10N^3 - N)}{30N^5} = \frac{x^5}{5}.$

For $N \geq 2$, define $\Delta x = \frac{x}{N}$, $x_i = (i-1)\Delta x$, $f(x_{i+1}) = (i\Delta x)^4$,

$$\begin{aligned}
\text{Area}_N^{\text{Up}} &:= \sum_{i=1}^N f(x_{i+1}) \cdot \Delta x \\
&= \sum_{i=1}^N (i\Delta x)^4 \cdot \Delta x \\
&= \sum_{i=1}^N \left(i \frac{x}{N}\right)^4 \cdot \frac{x}{N} \\
&= \frac{x^5}{N^5} \cdot \sum_{i=1}^N i^4 \quad (\text{because both } N, x \text{ are independent of } i, \text{ we take them out of the sum}) \\
&= \frac{x^5}{N^5} \cdot \frac{6N^5 + 15N^4 + 10N^3 - N}{30} \quad (\text{by Remark 2.12}) \\
&= \frac{x^5(6N^5 + 15N^4 + 10N^3 - N)}{30N^5}.
\end{aligned}$$

Taking the limit as $N \rightarrow \infty$ completes the result. ■

2.6.3 Products and Ratios of Exponential and Monomial Terms

It is also easy to determine the limits of products and ratios of exponentials (a^x) and monomials (x^k), where k can be positive or negative.

Intuition for the limits of $x^k \cdot a^x$ and $\frac{a^x}{x^k}$: The picture to have in mind is that exponentials a^x (for $a > 1, x > 0$) **grow infinitely faster** than any monomial, and hence, the exponential term dominates the monomial even while the monomial in the denominator is exploding to infinity. (It's pretty clear that the product of an exploding exponential and an exploding monomial is going to infinity).

Similarly, exponentials $a^{-x} = \left(\frac{1}{a}\right)^x$ (for $a > 1, x > 0$) **decay infinitely faster** than any monomial, and hence, the exponential term dominates the monomial even when the monomial in the numerator is exploding to infinity. (It's pretty clear that the product of a decaying exponential and a decaying monomial will be zero).

While the exponential always dominates the tug of war, **there is a subtlety with negative limits**, that is, when you look at limits as $x \rightarrow -\infty$: then the monomial can change the sign. That's the only catch. Otherwise, the same story holds: the exponential wins, but the monomial determines the sign.

The above is Prop. 2.28, Corollary 2.29, and Prop. 2.30 in a nutshell.

Proposition 2.28: Exponentials Grow Faster than any Monomial Function

The following limits hold for real variables x and for integers n replacing x :

- (a) For all $a > 1$ and integers $0 \leq m < \infty$, $\lim_{x \rightarrow \infty} \frac{a^x}{x^m} = \infty$.
- (b) For all $a > 1$ and integers $0 \leq m < \infty$, $\lim_{x \rightarrow \infty} \frac{x^m}{a^x} = 0$.
- (c) For all $a > 1$ and integers $0 \leq m < \infty$, $\lim_{x \rightarrow \infty} x^m \cdot a^x = \infty$.

Note: For the interested learner, the proof is given in Chapter 2.11. When using this result, don't forget that $x^{-m} := \frac{1}{x^m}$ and $a^{-x} := \frac{1}{a^x}$, which allow you to rewrite a function with negative exponents (powers) into one with positive exponents (powers).

By noting that $0 < b < 1$ implies that $a := \frac{1}{b} > 1$, and that $a^x = \frac{1}{b^x}$, the limits in Prop. 2.28 can also be applied to exponentials with bases between zero and one. For consistency and ease of use, we restate them for $0 < a < 1$.

Corollary 2.29: Exponentials Also Decay Faster than any Monomial Function

The following limits hold for real variables x and for integers n replacing x :

$$(a) \text{ For all } 0 < a < 1 \text{ and integers } 0 \leq m < \infty, \lim_{x \rightarrow \infty} x^m \cdot a^x = 0.$$

$$(b) \text{ For all } 0 < a < 1 \text{ and integers } 0 \leq m < \infty, \lim_{x \rightarrow \infty} \frac{x^m}{a^x} = \infty.$$

$$(c) \text{ For all } 0 < a < 1 \text{ and integers } 0 \leq m < \infty, \lim_{x \rightarrow \infty} \frac{1}{x^m \cdot a^x} = \infty.$$

Proposition 2.30: Relating Limits at Negative Infinity to Limits at (Positive) Infinity

The following relations hold for real variables x and for integers n replacing x :

$$(a) \text{ For all } a > 0, a \neq 1 \text{ and integers } 0 \leq m < \infty, \lim_{x \rightarrow -\infty} x^m \cdot a^x = (-1)^m \cdot \lim_{x \rightarrow +\infty} \frac{x^m}{a^x}.$$

$$(b) \text{ For all } a > 0, a \neq 1 \text{ and integers } 0 \leq m < \infty, \lim_{x \rightarrow -\infty} \frac{x^m}{a^x} = (-1)^m \cdot \lim_{x \rightarrow +\infty} x^m \cdot a^x.$$

$$(c) \text{ For all } a > 0, a \neq 1 \text{ and integers } 0 \leq m < \infty, \lim_{x \rightarrow -\infty} \frac{a^x}{x^m} = (-1)^m \cdot \lim_{x \rightarrow +\infty} \frac{1}{x^m \cdot a^x}.$$

Note: Because $1^x = 1$ and $1^n = 1$ for all $x \in \mathbb{R}$ and $n \in \mathbb{Z}$, it is not strictly necessary to make an exception for $a = 1$ in the above limits. It is still true that

$$\lim_{x \rightarrow -\infty} x^m = (-1)^m \cdot \lim_{x \rightarrow +\infty} x^m \text{ and } \lim_{x \rightarrow -\infty} \frac{1}{x^m} = (-1)^m \cdot \lim_{x \rightarrow +\infty} \frac{1}{x^m}.$$

However, because $a = 1$ is excluded from Prop. 2.28 and Corollary 2.29, we wanted to avoid misapplication of those results.

Example 2.31. Where applicable, use Prop. 2.28, Corollary 2.29, and Prop. 2.30 to compute $\lim_{x \rightarrow \pm\infty} f(x)$ for the following functions.

$$(a) f : \mathbb{R} \rightarrow \mathbb{R} \text{ by } f(x) = x^3 \cdot (0.5)^x.$$

$$(b) f : \mathbb{R} \rightarrow \mathbb{R} \text{ by } f(x) = \frac{2^x}{x^{27}}.$$

$$(c) f : \mathbb{R} \rightarrow \mathbb{R} \text{ by } f(x) = e^x \cdot \sin(x).$$

$$(d) f : \mathbb{R} \rightarrow \mathbb{R} \text{ by } f(x) = (1.01)^x \cdot \frac{-4x^3 + 11x + 14}{1+x^4}.$$

$$(e) f : \mathbb{R} \rightarrow \mathbb{R} \text{ by } f(x) = \frac{x^4}{5^{\frac{x}{2}}}.$$

Solutions:

$$(a) f(x) = x^3 \cdot (0.5)^x. \text{ Ans. } \lim_{x \rightarrow \infty} f(x) = 0 \text{ and } \lim_{x \rightarrow -\infty} f(x) = -\infty.$$

- $a = 0.5 \implies 0 < a < 1$. Hence, by Corollary 2.29-(a), $\lim_{x \rightarrow \infty} f(x) = 0$.

- $a = 0.5 \implies 0 < a < 1$ and m is odd. Hence, by Corollary 2.30-(a), $\lim_{x \rightarrow -\infty} f(x) = (-1)^3 \lim_{x \rightarrow +\infty} \frac{x^3}{(0.5)^x} = -\infty$.

(b) $f(x) = \frac{2^x}{x^{27}}$. **Ans.** $\lim_{x \rightarrow \infty} f(x) = \infty$ and $\lim_{x \rightarrow -\infty} f(x) = 0$.

- $a = 2 \implies a > 1$. Hence, by Prop. 2.28-(a), $\lim_{x \rightarrow \infty} f(x) = \infty$.

- $a = 2 \implies a > 1$ and m is odd. Hence, by Corollary 2.30-(c), $\lim_{x \rightarrow -\infty} f(x) = (-1)^{27} \lim_{x \rightarrow +\infty} \frac{1}{x^{27} \cdot 2^x} = 0$.

(c) $f(x) = e^x \cdot \sin(x)$. **Ans.** Due to the $\sin(x)$ term, Prop. 2.28, Corollary 2.29, and Prop. 2.30 are not applicable. Later, we will learn that $\lim_{x \rightarrow -\infty} e^x \cdot \sin(x) = 0$, while $\lim_{x \rightarrow +\infty} e^x \cdot \sin(x)$ does not exist.

(d) $f(x) = (1.01)^x \cdot \frac{-4x^3 + 11x + 14}{1+x^4}$. **Ans.** Because the exponential term $(1.01)^x$ is multiplying a rational function instead of a monomial, Prop. 2.28, Corollary 2.29, and Prop. 2.30 are not (directly) applicable. This is a fine answer. If you want to dig a bit deeper, as $|x| \rightarrow \infty$, Prop. 2.21 encourages us to approximate a rational function by

$$\frac{b_m x^m + b_{m-1} x^{m-1} + \dots + b_0}{a_n x^n + a_{n-1} x^{n-1} + \dots + a_0} \approx \frac{b_m}{a_n} \cdot x^{m-n} \text{ for } |x| \text{ sufficiently large.}$$

Doing so for $\frac{-4x^3 + 11x + 14}{1+x^4}$ results in, for $|x|$ large, $f(x)$ behaves like $g(x) := -4 \frac{(1.01)^x}{x}$. For this function, we can apply Prop. 2.28, Corollary 2.29, and Prop. 2.30 to obtain

- $\lim_{x \rightarrow \infty} (1.01)^x \cdot \frac{-4x^3 + 11x + 14}{1+x^4} = \lim_{x \rightarrow \infty} g(x) = \lim_{x \rightarrow \infty} -4 \frac{(1.01)^x}{x} = -\infty$.
- $\lim_{x \rightarrow -\infty} (1.01)^x \cdot \frac{-4x^3 + 11x + 14}{1+x^4} = \lim_{x \rightarrow -\infty} g(x) = \lim_{x \rightarrow -\infty} -4 \frac{(1.01)^x}{x} = 0$.

(e) $f(x) = \frac{x^4}{5^{\frac{x}{2}}}$ **Ans.** $\lim_{x \rightarrow \infty} f(x) = 0$ and $\lim_{x \rightarrow -\infty} f(x) = \infty$.

We rewrite $f(x) = \frac{x^4}{5^{\frac{x}{2}}} = \frac{x^4}{(\sqrt{5})^x}$

- $a = \sqrt{5} \implies a > 1$. Hence, by Prop. 2.28-(b), $\lim_{x \rightarrow \infty} f(x) = 0$.

- $a = \sqrt{5} \implies a > 1$ and m is even. Hence, by Prop. 2.30-(b), $\lim_{x \rightarrow -\infty} f(x) = (-1)^4 \lim_{x \rightarrow +\infty} x^4 \cdot \left(\frac{\sqrt{5}}{5}\right)^x = 0$, because $\frac{\sqrt{5}}{5} < 1$.

■

Example 2.32. From *Math Curiosity*, solve for $a \in \mathbb{N}$, if one exists, such that

$$\lim_{n \rightarrow \infty} \frac{n^{2022}}{n^a - (n-1)^a} = \frac{1}{2023}.$$

Hint: Think about the Binomial Theorem.

Solution By the Binomial Theorem,

$$\begin{aligned} n^a - (n-1)^a &= n^a - \left(n^a + \binom{a}{1} n^{a-1} (-1)^1 + O(n, a-2) \right) \\ &= an^{a-1} - O(n, a-2). \end{aligned}$$

Hence,

$$\frac{n^{2022}}{n^a - (n-1)^a} = \frac{n^{2022}}{an^{a-1} - O(n, a-2)}.$$

Applying Prop. 2.21 and Fact 2.24, we have a finite limit if, and only if, $a \geq 2023$, and we have a non-zero limit if, and only if, $a \leq 2023$. Hence, our only hope is $a = 2023$. Indeed, we then have

$$\lim_{n \rightarrow \infty} \frac{n^{2022}}{n^a - (n-1)^a} = \lim_{n \rightarrow \infty} \frac{n^{2022}}{n^{2023} - (n-1)^{2023}} = \lim_{n \rightarrow \infty} \frac{n^{2022}}{2023n^{2022} - O(n, 2021)} = \frac{1}{2023}.$$

■

2.6.4 Algebraic Operations with Infinity: Determinate and Indeterminate Forms

In Definition 1.1, we emphasized that the symbol “ ∞ ” is a concept and not a number. In particular,

- ∞ represents that which is larger than any real number, meaning for all $\alpha \in \mathbb{R}$, $\infty > \alpha$ (aka, $\alpha < \infty$), while
- $-\infty$ represents that which is smaller than any real number, meaning for all $\alpha \in \mathbb{R}$, $-\infty < \alpha$ (aka, $\alpha > -\infty$).

Nevertheless, as with numbers, some **algebraic operations** make sense³. For example,

- for all $\alpha \in \mathbb{R}$, $\alpha + \infty = \infty$ (a finite quantity plus an infinite quantity is still infinite). Its negative is,
- for all $\alpha \in \mathbb{R}$, $\alpha - \infty = -\infty$, which can also be written as $-\infty + \alpha = -\infty$.
- The above can be extended by $\infty + \infty = \infty$ (the sum of two infinitely “positive” quantities is infinitely positive, and the negative of this,
- $-\infty - \infty = -\infty + (-\infty) = -\infty$ (the sum of two infinitely negative quantities).

What about multiplication and division? Similar rules can be defined:

- For all $\alpha > 0$, $\alpha \cdot \infty = \infty$ (an infinite quantity scaled by a positive number is still infinite).
- For all $\alpha < 0$, $\alpha \cdot \infty = -\infty$ (an infinite quantity scaled by a negative number yields minus infinity). This has to hold if we accept that $(-1) \cdot \infty := -\infty$. Moreover,
- For all $\alpha \in \mathbb{R}$, $\frac{\alpha}{\pm\infty} = 0$.

In the above, we have carefully excluded the following cases, called **indeterminate forms**, because assigning a fixed value to the operations will lead to logical fallacies.

Indeterminate Forms

Definition 2.33. *The following are **indeterminate forms** (aka, **operations that are undefined and must remain undefined**):*

(a) $\infty - \infty$

(b) $\frac{\infty}{\infty}$

(c) $0 \cdot \infty$

(d) $\frac{0}{0}$

Though less important to us in this course, for completeness, we mention three additional indeterminate forms

(e) 1^∞

(f) 0^0

(g) ∞^0

In the following, we illustrate why the above are “indeterminate”.

- $\infty - \infty$ can arrive through $\lim_{x \rightarrow \infty} (x - x^3) = -\infty$, $\lim_{x \rightarrow \infty} (x^3 - x) = \infty$, and $\lim_{x \rightarrow \infty} (x - x) = 0$. Without knowing “the rate at which each infinity” arises, it is impossible to know which value to assign!
- $\frac{\infty}{\infty}$ can arrive through $\lim_{x \rightarrow \infty} \frac{x^3}{x+1} = \infty$, $\lim_{x \rightarrow \infty} \frac{x}{x^3+1} = 0$, and $\lim_{x \rightarrow \infty} \frac{x}{x+1} = 1$. Without knowing “the rate at which each infinity” arises, it is impossible to know which value to assign!

³Determinate is the opposite of indeterminate. Legitimate algebraic operations with infinity are sometimes called **determinate** forms.

- $0 \cdot \infty$ can arrive from $\lim_{x \rightarrow \infty} e^{-x} \cdot x = 0$, $\lim_{x \rightarrow \infty} 2^{-x} \cdot 3^x = \infty$, and $\lim_{x \rightarrow \infty} \frac{1}{1+x^2} \cdot x^2 = 1$. Without knowing “the rate at which zero and infinity” arise, it is impossible to know which value to assign!
- $\frac{0}{0}$ can arrive from $\lim_{x \rightarrow \infty} \frac{e^{-x^2}}{e^{-x}} = 0$, $\lim_{x \rightarrow \infty} \frac{2^{-x}}{3^{-x}} = \infty$, and $\lim_{x \rightarrow \infty} \frac{1 - e^{-\frac{1}{x}}}{\frac{1}{x}} = 1$. Without knowing “the rate at which each zero” arises, it is impossible to know which value to assign!

Associated Video: Indeterminate: the hidden power of 0 divided by 0 by Mathologer. Be aware that his videos are meant for those who are deeply curious about the why and how of Calculus. He has almost 1 M subscribers, which is pretty outrageous for a math channel.

2.6.5 Limit Algebra

It would seem much more efficient to take limits term by term instead of what we did in the solutions to Example 2.26. In fact, what you would hope to be true is true: **as long as you avoid the indeterminate forms listed above, Calculus is not evil :-).**

Proposition 2.34: Limits of Linear Combinations

Suppose that for all $1 \leq i \leq n$, $f_i : (0, \infty) \rightarrow \mathbb{R}$ has a finite limit at infinity, that is, $\lim_{x \rightarrow \infty} f_i(x)$ exists and is finite. Then, for all coefficients $\alpha_i \in \mathbb{R}$,

$$\lim_{x \rightarrow \infty} \left(\sum_{i=1}^n \alpha_i f_i(x) \right) = \sum_{i=1}^n \left(\alpha_i \lim_{x \rightarrow \infty} f_i(x) \right); \quad (2.28)$$

the limit of the linear combination exists and equals the linear combination of the limits.

Note: The assumption of the individual limits existing and being finite allows us to avoid the dreaded “ $\infty - \infty$ ”.

Once again, the above shows that Calculus is not an intuition-defying monster. Finite limits of products and ratios can be also computed sequentially. It is even straightforward to allow one of the limits to be unbounded and still draw a conclusion. Below, we avoid the dreaded $\frac{\infty}{\infty}$, $0 \cdot \infty$, and $\frac{0}{0}$ by forcing $g(x)$ to have a finite non-zero limit.

Proposition 2.35: Limits of Products and Ratios

Suppose that $g : (0, \infty) \rightarrow \mathbb{R}$ has a limit at infinity of one, i.e., $\lim_{x \rightarrow \infty} g(x) = 1.0$. Then for any function $f : (0, \infty) \rightarrow \mathbb{R}$,

$$\begin{aligned} \lim_{x \rightarrow \infty} f(x) \cdot g(x) &= \lim_{x \rightarrow \infty} f(x), \text{ and} \\ \lim_{x \rightarrow \infty} \frac{f(x)}{g(x)} &= \lim_{x \rightarrow \infty} f(x). \end{aligned} \quad (2.29)$$

Useful special cases: in particular,

- if $\lim_{x \rightarrow \infty} f(x) = L$, for $L \in \mathbb{R}$, then $\lim_{x \rightarrow \infty} f(x) \cdot g(x) = L$ and $\lim_{x \rightarrow \infty} \frac{f(x)}{g(x)} = L$;
- if $\lim_{x \rightarrow \infty} f(x) = \pm\infty$, then $\lim_{x \rightarrow \infty} f(x) \cdot g(x) = \pm\infty$ and $\lim_{x \rightarrow \infty} \frac{f(x)}{g(x)} = \pm\infty$; and
- if $\lim_{x \rightarrow \infty} f(x)$ does not exist, then $\lim_{x \rightarrow \infty} f(x) \cdot g(x)$ and $\lim_{x \rightarrow \infty} \frac{f(x)}{g(x)}$ do not exist.

Note: If $\lim_{x \rightarrow \infty} g(x) = M \neq 1$, and $M \neq 0$, then $\frac{g(x)}{M}$ has limit one. Hence, using $f(x) \cdot g(x) = (Mf(x)) \cdot \frac{g(x)}{M}$, and $\frac{f(x)}{g(x)} = \frac{\frac{1}{M}f(x)}{\frac{1}{M}g(x)}$ reduces the general case of the problem to the given (special) conditions.

If you remember (2.29), then the special cases follow immediately.

Example 2.36. Compute the indicated limits using Propositions 2.34 and 2.35:

$$(a) \lim_{x \rightarrow \infty} \frac{x^2}{3-x^3} + \frac{4+x+3x^2}{x^2-6x+11}.$$

$$(b) \lim_{x \rightarrow \infty} \frac{x^3}{3-x^3} \cdot \frac{4+x+3x^2}{x^2-6x+11}.$$

$$(c) \lim_{x \rightarrow \infty} \frac{\frac{x^2}{3-x^3}}{\frac{4+x+3x^2}{x^2-6x+11}}.$$

$$(d) \lim_{x \rightarrow \infty} \frac{x \cdot \sin(x)}{1 + e^{-x} + e^{-2x}}.$$

Solution:

$$(a) \text{ Ans. } \lim_{x \rightarrow \infty} \frac{x^2}{3-x^3} + \frac{4+x+3x^2}{x^2-6x+11} = 3.0$$

Why: By Prop. 2.21, $\lim_{x \rightarrow \infty} \frac{x^2}{3-x^3} = 0$ and $\lim_{x \rightarrow \infty} \frac{4+x+3x^2}{x^2-6x+11} = 3.0$. Hence, by Prop. 2.34, the limits add, and we obtain the given answer.

$$(b) \text{ Ans. } \lim_{x \rightarrow \infty} \frac{x^2}{3-x^3} \cdot \frac{4+x+3x^2}{x^2-6x+11} = 0.0$$

Why: By Prop. 2.21, $\lim_{x \rightarrow \infty} \frac{x^2}{3-x^3} = 0$ and $\lim_{x \rightarrow \infty} \frac{4+x+3x^2}{x^2-6x+11} = 3.0$. Hence, by Prop. 2.35, the limits multiply, and we obtain the given answer.

$$(c) \text{ Ans. } \lim_{x \rightarrow \infty} \frac{\frac{x^3}{3-x^3}}{\frac{4+x+3x^2}{x^2-6x+11}} = -\frac{1}{3}.$$

Why: By Prop. 2.21, $\lim_{x \rightarrow \infty} \frac{x^3}{3-x^3} = -1.0$ and $\lim_{x \rightarrow \infty} \frac{4+x+3x^2}{x^2-6x+11} = 3.0$. Hence, by Prop. 2.35, the limits divide, and we obtain the given answer.

$$(d) \text{ Ans. } \lim_{x \rightarrow \infty} \frac{x \cdot \sin(x)}{1 + e^{-x} + e^{-2x}} \text{ does not exist.}$$

Why: By Corollary 2.29 and Prop. 2.34, the denominator limit, $\lim_{x \rightarrow \infty} 1 + e^{-x} + e^{-2x} = 1.0$, while the numerator limit, $\lim_{x \rightarrow \infty} x \cdot \sin(x)$, does not exist due to the oscillating sinusoid. Hence, by Prop. 2.35, the limit of the ratio does not exist. ■

2.7 Geometric Sums

A geometric sequence is a sequence of numbers where each term after the first is found by multiplying the previous term by a fixed, non-zero number r called the **common ratio**. The general form of a **geometric sequence** is:

$$a_0 = a, a_1 = a \cdot r, a_2 = a \cdot r^2, a_3 = a \cdot r^3, \dots, a_i = a \cdot r^i,$$

where a is the first term, and r is the **common ratio**. We are more interested in a **geometric sum**, which is the finite sum of the form

$$S_n := \sum_{i=0}^n a \cdot r^i = a + a \cdot r + a \cdot r^2 + \dots + a \cdot r^n. \quad (2.30)$$

Intuition for Geometric Sums (aka, geometric series) is uncommonly visual:

- Three Geometric Series in an Equilateral Triangle (visual proof without words) by @MathVisualProofs.
- Playlist for Geometric Sums by @MathVisualProofs.

For later use, we note that $S_{n+1} = a + r \cdot S_n$.

Proposition 2.37: Geometric Sums

For $r \neq 1$,

$$S_n = \frac{a(1 - r^{n+1})}{1 - r}, \quad (2.31)$$

and if $|r| < 1$,

$$\lim_{n \rightarrow \infty} S_n = \frac{a}{1 - r}. \quad (2.32)$$

Note: (i) If $r = 1$, then $S_n = a(n + 1)$ and if $a \neq 0$, its limit is unbounded as $n \rightarrow \infty$. (ii) $\sum_{i=0}^{\infty} ar^i := \lim_{n \rightarrow \infty} \sum_{i=0}^n ar^i$ is called a **geometric series**.

Proof: The proof can be given by induction. We'll take $k = 1$ as the base case, though we could have taken $k = 0$.

Base Case: $k = 1$: $S_1 = a + a \cdot r = a \cdot (1 + r) = \frac{a \cdot (1+r) \cdot (1-r)}{1-r} = \frac{a \cdot (1-r^2)}{1-r}$. As is commonly the case, it is quite straightforward.

Induction Step: We define $L_k := S_k$ and $R_k := \frac{a(1-r^{k+1})}{1-r}$ and seek to show that $L_k = R_k$ being true for some $k \geq 1$ implies it is also true for $k + 1$.

$$\begin{aligned} L_{k+1} &= S_{k+1} \\ &= a + r \cdot S_k \quad (\text{by the note just below (2.30)}) \\ &= a + r \cdot \frac{a(1 - r^{k+1})}{1 - r} \quad (\text{by the induction hypothesis}) \\ &= a + \frac{a(r - r^{k+2})}{1 - r} \quad (\text{algebra, multiply by } r) \\ &= a \cdot \left(1 + \frac{r - r^{k+2}}{1 - r}\right) \quad (\text{algebra, factor out } a) \\ &= a \cdot \left(\frac{1 - r}{1 - r} + \frac{r - r^{k+2}}{1 - r}\right) \quad (\text{rewrite the number 1 as a ratio}) \\ &= a \cdot \left(\frac{1 - r^{k+2}}{1 - r}\right) \quad (\text{algebra}). \end{aligned}$$

By its definition,

$$R_{k+1} = \frac{a(1 - r^{k+2})}{1 - r},$$

which equals L_{k+1} , so the proof by induction is complete.

Now, we look at using Prop. 2.29 to evaluate the limit. We first note that if $r < 0$, then $r = (-1) \cdot |r|$ and therefore, $r^n = (-1)^n \cdot |r|^n$. Hence, by Prop. 2.29, if $|r| < 1$, then $\lim_{n \rightarrow \infty} r^{n+1} = r \cdot \lim_{n \rightarrow \infty} r^n = 0$, independent of the sign of r . Therefore,

$$\lim_{n \rightarrow \infty} S_n = \lim_{n \rightarrow \infty} \frac{a(1 - r^{n+1})}{1 - r} = \frac{a}{1 - r}.$$

Example 2.38. Determine which of the following sums have the form of (2.30), and hence are a **geometric sum**. For those cases that are geometric sums, use Prop. 2.37 to compute S_n , and if it exists, $s_\infty := \lim_{n \rightarrow \infty} S_n$.

- (a) $S_n = \sum_{i=0}^n \frac{1}{2^i}$.
- (b) $S_n = 1 + \frac{1}{2} + \frac{1}{2^2} + \cdots + \frac{1}{2^n}$.
- (c) $S_n = \sum_{i=0}^n (-1)^i$.
- (d) $S_n = 1 + \frac{1}{2} + \frac{1}{3} + \cdots + \frac{1}{n}$.
- (e) $S_n = 3 + 3e^x + 3e^{2x} + \cdots + 3e^{nx}$.
- (f) $S_n = 2 + \frac{1}{3^2} + \frac{1}{3^3} + \frac{1}{3^4} + \cdots + \frac{1}{3^n}$.

Solutions:

- (a) **Ans.** : $\sum_{i=0}^n \frac{1}{2^i} = \sum_{i=0}^n \left(\frac{1}{2}\right)^i$ is a geometric sum, with common ratio $r = \frac{1}{2}$ and $a = 1$. Hence,

$$S_n = \frac{a \cdot (1 - r^{n+1})}{1 - r} = \frac{1 \cdot (1 - (\frac{1}{2})^{n+1})}{1 - \frac{1}{2}} = 2 \cdot (1 - 2^{-(n+1)}) = 2 - 2^{-n},$$

and therefore, $S_\infty := \lim_{n \rightarrow \infty} S_n = 2$.

- (b) **Ans.** : $S_n = 1 + \frac{1}{2} + \frac{1}{2^2} + \cdots + \frac{1}{2^n} = \sum_{i=0}^n \frac{1}{2^i}$, the sum given in (a). Hence, the answer is the same.

- (c) **Ans.** : $S_n = \sum_{i=0}^n (-1)^i$ is a geometric series, with common ratio $r = -1$ and $a = 1$. Hence,

$$S_n = \frac{a \cdot (1 - r^{n+1})}{1 - r} = \frac{1 \cdot (1 - (-1)^{n+1})}{1 - (-1)} = \frac{1 - (-1)^{n+1}}{2} = \begin{cases} 0 & n \text{ odd} \\ 1 & n \text{ even.} \end{cases}$$

Because $|r| = 1$, $\lim_{n \rightarrow \infty} S_n$ does not exist; that is hopefully clear because S_n alternates between 0 and 1, which does not converge.

- (d) **Ans.** : $S_n = 1 + \frac{1}{2} + \frac{1}{3} + \cdots + \frac{1}{n}$ is NOT a geometric sum. If it were to be a geometric series, from the first term, we would have that $a = 1$, and from the second term, we would deduce that $r = \frac{1}{2}$. The problem comes with the third term because

$$\frac{1}{3} \neq r^2 = \left(\frac{1}{2}\right)^2 = \frac{1}{4}.$$

- (e) **Ans.** : $S_n = 3 + 3e^x + 3e^{2x} + \cdots + 3e^{nx} = \sum_{i=0}^n 3 \cdot (e^x)^i$ is a geometric sum, with common ratio $r = e^x$ and $a = 3$. Hence,

$$S_n = \frac{a \cdot (1 - r^{n+1})}{1 - r} = \frac{3 \cdot (1 - (e^x)^{n+1})}{1 - e^x} = 3 \frac{1 - e^{(n+1)x}}{1 - e^x}.$$

For $x < 0$, $|r| = |e^x| = e^x < 1$, while for $x \geq 0$, $|r| > 1$. Therefore,

$$S_\infty := \lim_{n \rightarrow \infty} S_n = \begin{cases} \frac{1}{1-e^x} & x < 0 \\ \text{undefined} & x \geq 0. \end{cases}$$

- (f) $S_n = 2 + \frac{1}{3^2} + \frac{1}{3^3} + \frac{1}{3^4} + \cdots + \frac{1}{3^n}$ is NOT a geometric sum. If it were to be a geometric sum, from the first term, we would have that $a = 2$, and from the second term, we would deduce that $r = \frac{1}{3^2} = \frac{1}{9}$. The problem comes with the third term because

$$\frac{1}{3^3} = \frac{1}{27} \neq a \cdot r^2 = 2 \left(\frac{1}{9}\right)^2 = \frac{2}{81} = \frac{1}{2} \cdot \frac{1}{81}.$$

However, the given sum can be rewritten as a constant plus a geometric sum, and hence we can still compute it with relative ease. Indeed, let's define

$$\bar{S}_n := 1 + \frac{1}{3} + \frac{1}{3^2} + \frac{1}{3^3} + \frac{1}{3^4} + \cdots + \frac{1}{3^n} = -1 + \frac{1}{3} + S_n.$$

\overline{S}_n is a geometric sum with common ratio $r = \frac{1}{3}$ and $a = 1$. Hence,

$$\overline{S}_n = \frac{a \cdot (1 - r^{n+1})}{1 - r} = \frac{1 \cdot (1 - (\frac{1}{3})^{n+1})}{1 - \frac{1}{3}} = \frac{3}{2} \cdot (1 - 3^{-(n+1)}),$$

and therefore, $\overline{S}_\infty := \lim_{n \rightarrow \infty} \overline{S}_n = \frac{3}{2}$.

It follows that,

$$S_n = \overline{S}_n - \frac{4}{3} = \frac{3}{2} \cdot (1 - 3^{-(n+1)}) + \frac{2}{3},$$

and $S_\infty = \frac{3}{2} + \frac{2}{3} = \frac{9+4}{6} = \frac{13}{6} = 2 + \frac{1}{6}$.

Example 2.38-(f) illustrates that sometimes a sum can be modified slightly so that it becomes a geometric sum. That should be considered a standard “trick of the trade”. Real problems do not come pre-packaged for directly applying schoolwork in a rote manner. Some creativity is required!

2.8 Countable Sets

Our goal here is to illustrate why proofs are important. Moreover, we wish to do that in a setting where the logic of the proof (its roadmap, if you will) may be hard to invent on your own, yet, once it is spelled out to you, it is quite beautiful and can be understood without a lot of experience with proofs. Even more, after reading through the material, there is a chance, albeit a small one, that you might even find it edifying (i.e, enlightening, educational, informative, or illuminating).

To accomplish all of these goals, we'll present a non-intuitive, *even shocking*, result that **Georg Cantor** proved in 1874. The following is a loosely paraphrased excerpt from Wikipedia.

Secrets of the Arcane 2.39: Renegade Mathematician and Corrupter of Youth ? You be the judge.

In the heat of the mathematical revolution, Georg Cantor stirred up controversy with his mind-bending theory of transfinite numbers. This audacious idea was met with the equivalent of academic side-eye, even triggering the ire of mathematical heavyweights such as Leopold Kronecker and Henri Poincaré, not to mention the philosophical doubts raised by Ludwig Wittgenstein. The vitriol Cantor faced was far from mild. Kronecker went so far as to lambaste Cantor publicly, **branding him a “scientific charlatan”, a “renegade”, and a “corrupter of youth.”** The ridicule and condemnation of the mathematical community eventually led Cantor to depression and admission to a mental healthcare facility.

While Cantor's treatment by his peers was unjust, he had his issues as well. “Cantor was paradoxically opposed to theories of infinitesimals of his contemporaries Otto Stolz and Paul du Bois-Reymond, describing [Stolz and du Bois-Reymond] as “**an abomination” and “a cholera bacillus of mathematics.**” When was the last time you called someone “a cholera bacillus,” a cholera bacterium?

In a certain sense, the concept we will study is quite simple: how many elements are contained in a set? For example, if we can list the elements of a set, such as

$$S = \{s_1, s_2, s_3, s_4, s_5\},$$

where $s_i \neq s_j$ for $i \neq j$, then we'd naturally say the set has five elements. In mathematics, the number of elements in a set is called its **cardinality**. If you can index the elements of the set using a finite number of counting numbers, then it has finite cardinality. The listing does not have to use consecutive integers for us to conclude that the set's cardinality is finite; for example, we can consider

$$T := \{2, 4, 11, 10^9\}$$

as already being listed by the numbers 2, 3, 11, and 10^9 ; because the largest element, 10^9 , is finite, there is no need to assign $t_1 := 2, t_2 := 4, t_3 := 11, t_4 := 10^9$. However, sets such as the natural numbers \mathbb{N} and the real numbers \mathbb{R} have an infinite number of elements, and hence their cardinality must be infinity. **Cantor asked the question: are the infinities associated with the natural numbers and the real numbers the same, or are they different in some fundamental way?** Huh? And that led him to be branded a “corrupter of youth?” To dig a little deeper, we must carefully define what we mean by “counting the number of elements in a set.”

What does it mean for a set to be countable?

Definition 2.40. A set S is **countable** if there exists an injective (i.e., 1:1) function $f : S \rightarrow \mathbb{N}$. This means that each element of the set is assigned a counting number. Moreover, distinct elements of the set are assigned distinct numbers so that we are not undercounting the number of elements.

S is **finitely countable** if the range of f is contained in a bounded subset of the natural numbers. Otherwise, S is called **countably infinite**. Said in a different way: in a finitely countable set, the process of counting comes to an end, whereas in a countably infinite set, it does not, just as with the counting numbers themselves.

Notes:

- (a) The **empty set** has no elements and hence, by tradition, is considered finitely countable.
- (b) If there is no injective function from S to the counting numbers, the set is called **uncountable**.
- (c) If an injective function $f : S \rightarrow \mathbb{N}$ exists and its range is contained in a bounded subset of the natural numbers, then S is called **countably finite** or **finitely countable** (both are used). Otherwise, S is called **countably infinite**. To be extra clear, this means that the range of $f : S \rightarrow \mathbb{N}$ contains an unbounded number of natural numbers (equivalently, the range does not have a largest element).

Remark 2.41. Some people prefer the terminology **listable** over **countable**. The term “listable” is less formal and not universally used in mathematics. However, when it is used, it typically implies the same concept as **countable**. In other words, if you can devise a method to list all the elements of a set – with no element being missed and each being listed only once – then the set is “listable” or **countable**.

Example 2.42. Determine which of the following sets are countably finite and which are countably infinite. In each case, exhibit a function from the set to a **finite** subset of the counting numbers or the full set of natural numbers.

- (a) $S_1 := \{\text{bus, submarine, carrot}\}$
- (b) $S_2 := \text{the set of even prime numbers.}$
- (c) $S_3 := \{2k \mid k \in \mathbb{N}\}$, the set of even numbers.
- (d) $S_4 := \text{the set of all } 3 \times 3 \text{ matrices with entries from the set } \{0, 1, 2, \dots, 9\}.$
- (e) $S_5 := \text{the set of all } \mathbf{invertible} 2 \times 2 \text{ matrices with entries from the set } \{1, 2, 3, 4, 5\}.$
- (f) $S_6 := \{x \in \mathbb{R} \mid x > 0, \log_{10}(x) \in \mathbb{N}, \text{ and } \log_{10}(x) > \pi\}$ (deliberately made to appear absurdly complicated).

Solution:

- (a) S_1 is finitely countable. We define $f : S_1 \rightarrow \{1, 2, 3\}$ by $f(\text{bus}) = 1$, $f(\text{submarine}) = 2$, and $f(\text{carrot}) = 3$. f is clearly injective (1:1).

Note: Our mapping does not have to be surjective (onto). Hence, another solution is $f : S_1 \rightarrow \{1, 2, 3, 4, 11, 17\}$ by $f(\text{bus}) = 1$, $f(\text{submarine}) = 2$, and $f(\text{carrot}) = 3$, and yet another solution is $f : S_1 \rightarrow \{1, 2, 3, 4, 11, 17\}$ by $f(\text{bus}) = 11$, $f(\text{submarine}) = 3$, and $f(\text{carrot}) = 17$. We can even define $f : S_1 \rightarrow \mathbb{N}$ by $f(\text{bus}) = 10$, $f(\text{submarine}) = 200$, and $f(\text{carrot}) = 300$, and then note that the range of f is contained in the finite set $\{1, 2, \dots, 500\}$.

In the same vein, a non-solution is $f : S_1 \rightarrow \{1, 2, 3\}$ by $f(\text{bus}) = 3$, $f(\text{submarine}) = 2$, and $f(\text{carrot}) = 2$ because f is not injective (1:1) with both “carrot” and “submarine” being mapped onto the number two. In this case, the undercounting is not a big deal, but I guess you would agree that the map $f : \mathbb{N} \rightarrow \mathbb{N}$ by $f(k) = 1$ for all k would result in severe undercounting of the natural numbers! That is why injectivity is important.

- (b) S_2 is finitely countable because there is only one even prime number, 2.

- (c) S_3 is countably infinite. We can define $f : S_3 \rightarrow \mathbb{N}$ by $f(2k) = k$. The function is one-to-one because $f(s) = k \iff s = 2k$. Hence, there is only one element of S_3 assigned to any element of \mathbb{N} . We could also have defined $f : S_3 \rightarrow \mathbb{N}$ by $f(s) = 10^s$. There is no restriction that the numbers be consecutive, i.e., the function can “skip over” numbers in \mathbb{N} . You can check that the second function is also injective.
- (d) S_4 is finitely countable. A 3×3 matrix has nine entries, and 3×3 matrices in S_4 have entries drawn from the set $\{0, 1, 2, \dots, 9\}$. Hence, there are 10^9 distinct matrices in S_4 . Building a map from $f : S_4 \rightarrow \{1, 2, 3, \dots, 10^9\}$ can be done using basic arithmetic. Let a_1, a_2, \dots, a_9 denote the entries of the 3×3 matrix, enumerated in any order you like. We define $f : S_4 \rightarrow \{1, 2, 3, \dots, 10^9\}$ by $f(A) = 1 + a_1 + a_2 \cdot 10 + \dots + a_8 \cdot 10^8 + a_9 \cdot 10^9$, which associates a unique base-10 counting number to each element of the set. (The “1” is added so that when f is applied to the matrix with all zero entries, it evaluates to a counting number, namely 1, instead of zero.)

Note: Suppose the matrices in S_4 had n -entries and their values were drawn from the set $\{0, 1, \dots, M - 1\}$, where n and M are counting numbers, $M \geq 2$. Then there are M^n distinct matrices, and you can define $f(A) = 1 + \sum_{i=1}^n a_i M^{i-1}$, which associates a unique counting number to each element of the set, and injectivity (one-to-one-ness) is easy to prove if you think of the number as being in base M .

- (e) S_5 is finitely countable, but why? For all of the other solutions, we constructed an explicit injective mapping from the set to the natural numbers or a subset of the natural numbers. To do that here, it appears that we’d have to identify those matrices that result in a zero determinant and eliminate them, and then start assigning numbers to the remaining matrices. Right?

Well, no, we do not have to do that. We see easily that the set of all 2×2 matrices with entries from the set $\{1, 2, 3, 4, 5\}$ has $5^4 = 625$ distinct elements, which is finite. The set S_5 has fewer elements than that because the matrices with a zero determinant are excluded (e.g., the matrix with ones in all four entries is not invertible). A subset of a finite set is finite because it has at most the same number of elements as the “larger, or parent” set. We formalize this observation below in Prop. 2.43. Hence, S_5 has at most 625 elements and is therefore finitely countable.

- (f) S_6 is countably infinite. The first natural number greater than π is four. Hence, $x \in S_6 \iff x = 10^k$, with $k \in \{4, 5, 6, \dots\}$. We can define $f : S_6 \rightarrow \mathbb{N}$ by

$$f(x) = k,$$

where $x = 10^k$; in other words, $f(x) = \log_{10}(x)$. You could also define $f(10^k) = k - 3$, for all $k \geq 4$, so that f would also be onto (surjective), but as we discussed above, that is not a requirement.

Proposition 2.43: Subsets of Countable Sets

A subset of a countable set is countable. In particular, a subset of a finitely countable set is finitely countable, while a subset of a countably infinite set can be either finitely countable or countably infinite.

Proof: Suppose that $f : S \rightarrow \mathbb{N}$ is injective and that $A \subset S$. Then because $a \in A \implies a \in S$, we have that $f(a) \in \mathbb{N}$. Moreover, if $a_1, a_2 \in A$ then $a_1, a_2 \in S$. Hence, if $a_1 \neq a_2$, the injectivity of $f : S \rightarrow \mathbb{N}$ yields $f(a_1) \neq f(a_2)$, so $f : A \rightarrow \mathbb{N}$ is injective whenever $f : S \rightarrow \mathbb{N}$ is injective. The rest of the proof is left to the learner. ■

Remark 2.44. When we restrict a function to a subset of its original domain, we are technically creating a new function. Why? Because a function is defined by its domain, codomain, and values. All three are important. In High School mathematics, we mostly think of a function as being defined by its values, whereas in College, we have to keep the fuller picture in mind. When the distinction is not all that important, we often do not change the function’s name, just as we used the notation f above when its domain was S and when its domain was $A \subset S$. When we want to call attention to the “new function’s (new) domain”, we use the notation $f|_A$, which is read “ f restricted to A ” or the “restriction of f to A ”. We can also write out the full function, $f : A \rightarrow \mathbb{N}$.

All of this has seemed rather tame, so, again, what got Cantor into so much trouble with the math community? Thank you for asking! The notion of countability led Cantor to the idea that there is more than one kind of infinity, and the mathematicians of his day found that heretical. To consider some easy cases first, ask yourself,

- How many counting numbers are there? (Hint, it is not a number.)
- How many even numbers are there? Let’s denote the even numbers by \mathbb{E} .
- How many odd numbers are there? Let’s denote the odd numbers by \mathbb{O} .

- Then wait, is the infinity of the even numbers half that of the counting numbers? And the same for the odd numbers? No, they are all the same. We have

- $f : \mathbb{N} \rightarrow \mathbb{N}$ by $f(k) = k$,
- $f : \mathbb{E} \rightarrow \mathbb{N}$ by $f(2k) = k$, for all $k \geq 1$,
- $f : \mathbb{O} \rightarrow \mathbb{N}$ by $f(2k - 1) = k$, for all $k \geq 1$, and
- $\mathbb{E} \cup \mathbb{O} = \mathbb{N}$ and $\mathbb{E} \cap \mathbb{O} = \emptyset$.

Each set is actually in one-to-one correspondence with the natural numbers. The same is true for all multiples of three or five, for example. Okay, so $\infty + \infty = \infty$ and $\frac{\infty}{2} = \infty$? Yes, that's an acceptable interpretation of the above, so long as you recall that you are doing arithmetic on a concept and not a number. Now, what about a matrix with an infinite number of rows and columns? Would that have ∞^2 elements? And is ∞^2 , assuming we can make sense out of it, different than ∞ ? Cantor really did stir up a hornet's nest!

2.8.1 The Rational Numbers are Countable

Because all counting numbers are rational numbers with denominators equal to one, there are at least as many rational numbers as there are counting numbers. But, it seems like there should be many more. Maybe infinitely more? Recall that we defined the rational numbers by,

$$\mathbb{Q} := \left\{ \frac{p}{q} \mid p \in \mathbb{Z}, q \in \mathbb{N}, \text{no common factors other than one} \right\}. \quad (2.33)$$

Having to check that the gcd (greatest common divisor) of two integers equals one would complicate the counting process, so Cantor employed Prop. 2.43, just as we did in Example 2.42 (dealing with invertible matrices with integer coefficients): he said, if he could prove that the set of all ratios of counting numbers is countable, with $1/2, 2/4, \dots, \frac{k}{2k}, \dots$ **all counted as distinct elements**⁴, then the rational numbers must be countable as well. So that is what he did, and we will too. Well, almost, we'll focus on the positive rational numbers, because if we can show that they are countable, then all rational numbers are countable following the same reasoning applied to the even and odd numbers.

Consider $S := \left\{ \frac{i}{j} \mid i \in \mathbb{N}, j \in \mathbb{N} \right\}$, which we can envision as an infinite version of Table 2.45. We switched from p and q to i and j because you are used to using ij as indices in a matrix, something we do in Table 2.45 and you will also do in Lab.

The first problem faced by Cantor was, how to count these numbers. If you try to count row by row, you hit an immediate snag; the first row never terminates, so you cannot come back and start counting the elements in the second row. The same issue arises if you try to count by columns. So what is a young Cantor to do? He counts the elements using the reverse diagonals, as shown in this [YouTube video](#). We do something similar in Table 2.46. You can even write down an explicit function, called the Cantor Pairing Function

$$f : S \rightarrow \mathbb{N}, \text{ by } f(i, j) := \frac{(i+j)(i+j-1)}{2} + j. \quad (2.34)$$

In Julia, you can easily verify that this formula allows you to associate a unique counting number to each pair (i, j) of counting numbers! Moreover, you can even play with the inverse of the Cantor Pairing Function. Here is a ChatGPT prompt that will produce the Julia code for you.

Write two functions in Julia that handle Cantor's pairing function and its inverse, using 1-based indexing for counting numbers:

1. **Cantor Pairing Function:** Implement a function “cantor_pair(i, j)” that takes two positive integers $i \geq 1$ and $j \geq 1$ and returns the unique Cantor pairing number associated with the pair (i, j) . Ensure that the function works with 1-based indexing, which means that i and j are both positive integers starting from 1.
2. **Inverse Cantor Function:** Implement a function “inverse_cantor(z)” that takes a Cantor pairing number $z \geq 1$ and returns the

⁴Mathematically, it is better to think of the “non-reduced” ratios p/q as pairs of numbers, (p, q) . This avoids you even thinking about 120/240 as being distinct from 1/2 because the pair of integers (120, 240) is clearly distinct from the pair (2, 4) and the pair (1, 2), while any pair of integers (p, q) generates a rational number through division, p/q .

$i \setminus j$	1	2	3	4	5	6	7	8	9	10	...
1	1/1	1/2	1/3	1/4	1/5	1/6	1/7	1/8	1/9	1/10	...
2	2/1	2/2	2/3	2/4	2/5	2/6	2/7	2/8	2/9	2/10	...
3	3/1	3/2	3/3	3/4	3/5	3/6	3/7	3/8	3/9	3/10	...
4	4/1	4/2	4/3	4/4	4/5	4/6	4/7	4/8	4/9	4/10	...
5	5/1	5/2	5/3	5/4	5/5	5/6	5/7	5/8	5/9	5/10	...
6	6/1	6/2	6/3	6/4	6/5	6/6	6/7	6/8	6/9	6/10	...
7	7/1	7/2	7/3	7/4	7/5	7/6	7/7	7/8	7/9	7/10	...
8	8/1	8/2	8/3	8/4	8/5	8/6	8/7	8/8	8/9	8/10	...
9	9/1	9/2	9/3	9/4	9/5	9/6	9/7	9/8	9/9	9/10	...
10	10/1	10/2	10/3	10/4	10/5	10/6	10/7	10/8	10/9	10/10	...
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\ddots

Table 2.45: The set of all positive ratios of natural numbers, now with each row having a constant numerator and the denominator increasing. You can think of it as a matrix with infinite rows and columns.

original pair (i, j) of positive integers. The function should correctly handle the inverse calculation using 1-based indexing.

Hence, the rational numbers are countable because they are a subset of the countable set consisting of all ratios of natural numbers, where every fraction occurs an infinite number of times! Indeed, if p/q is a fraction in reduced form, then $\frac{kp}{kq}$ shows up in the list of all ratios for all $k \geq 1$ (equivalently, the pair of numbers (kp, kq) shows up in Table 2.45).

Remark 2.47. Tables 2.45 and 2.46 were generated by ChatGPT4 under prompts created by the author.

2.8.2 The Real Numbers are Uncountable

Cantor employed a proof by contradiction to show that the real numbers are not countable. If you read how Euclid proved that $\sqrt{2}$ is not rational, then, yes, Cantor employed a similar method. But of course, there is a super cool twist, and it earned the moniker, **Cantor's Diagonal Argument**.

Because every real number $y \in \mathbb{R}$ can be expressed as $y = k + x$, for $k \in \mathbb{Z}$ and $x \in [0, 1)$, it is enough to show that the interval $[0, 1)$ is uncountable. We recall that

$$x \in [0, 1) \iff x = 0.d_1d_2d_3\dots, \quad (2.35)$$

and the decimal expansion does not terminate in an infinite string of nines.

Armed with these facts, we assume that the real numbers in $[0, 1)$ ARE COUNTABLE (i.e., listable). Hence, we can enumerate them as in Table 2.49, that is,

$$[0, 1) = \{x_1, x_2, x_3, \dots\}.$$

We emphasize that if the real numbers between zero and one are indeed countable, then we can *completely list* (i.e., without missing a single one of them) all the real numbers in the interval $[0, 1)$. Cantor next shows that this list must be incomplete by producing a new real number in the interval $[0, 1)$ that is not in the list, contradicting the countability of $[0, 1)$, and hence, of \mathbb{R} .

Let $d_i[n]$ denote the i -th digit of the decimal expansion of the n -th real number, x_n . The digits $d_n[n]$ are on the diagonal of Table 2.49. Cantor uses these diagonal elements to construct a new real number that is nowhere to be found in the original list! The new number has the decimal expansion,

$$\sigma := \sum_{n=1}^{\infty} \sigma_n 10^{-n}, \text{ where, } \sigma_n := \begin{cases} 0 & \text{if } d_n[n] \in \{1, 2, \dots, 8\} \\ 1 & \text{otherwise, i.e., } d_n[n] \text{ equals 0 or 9} \end{cases};$$

$i \setminus j$	1	2	3	4	5	6	7	8	9	10
1	1	3	6	10	15	21	28	36	45	55
2	2	5	9	14	20	27	35	44	54	65
3	4	8	13	19	26	34	43	53	64	76
4	7	12	18	25	33	42	52	63	75	88
5	11	17	24	32	41	51	62	74	87	101
6	16	23	31	40	50	61	73	86	100	115
7	22	30	39	49	60	72	85	99	114	130
8	29	38	48	59	71	84	98	113	129	146
9	37	47	58	70	83	97	112	128	145	163
10	46	57	69	82	96	111	127	144	162	181

Table 2.46: Counting ratios via a “broken” Cantor Snake! We assign one to the first element of the first row. Then we move to the second row and assign two to the first element of the second row, and then, we work our way up the reverse diagonal, assigning three to the second element of the first row. We then go to the third row and repeat the process for the next reverse diagonal. This allows the entries of every row and column to be visited exactly once. **This blew the minds of Cantor’s colleagues.** What about yours? **Note: we are not showing numbers outside of the 10×10 “footprint”.** You can do more with Julia. For example, using the Cantor Pairing Function in (2.34), you can compute that $(i, j) = (11, 1) \rightarrow 56$.

in other words, he takes all the digits on the diagonal of his infinite matrix and constructs a new number that differs from the n -th number in its n -th digit. By construction, this new number⁵ is nowhere in the original list, because it is different from the first number in the first digit, it is different from the second number in the second digit, etc. The existence of a number that is not in the list contradicts the list containing all real numbers in $[0, 1]$. Because the new number is built from the diagonal elements of the table, the method is called “**Cantor’s Diagonal Argument**”. It’s pretty clever!

Secrets of the Arcane 2.48: Infinity Comes in Multiple Sizes!

By Cantor showing that real numbers are not countable, the world was informed that the “infinity” of the real numbers is **Strictly Greater Than** the “infinity” of the counting numbers. While we’ll stop here, Cantor went on to show that there is an infinite chain^a of “infinities”, each one strictly more “numerous” than the previous one. That is mind-bending stuff, and yeah, perhaps capable of corrupting young minds!

^aIn case you are wondering, the collection of all “infinities” (technically, the collection of all infinite cardinal numbers), does not form a set. ChatGPT can tell you more about that!

While the notion of countable sets is not crucial for the remainder of the course, the idea that mathematics needs to be done carefully is paramount. You got a taste of this in Chapter 1, where we tried to express concepts you learned in High School in a careful manner.

2.9 (Optional Read:) How to Really Define Euler’s Number

We want to define

$$e := \lim_{n \rightarrow \infty} \underbrace{\left(1 + \frac{1}{n}\right)^n}_{a_n}, \quad (2.36)$$

but we don’t have enough tools to make sense of the limit. Yet, the numbers $a_n := \left(1 + \frac{1}{n}\right)^n$ that Bernoulli discovered through his study of compound interest have some very nice properties: for each integer $n \geq 0$,

1. a_n is rational,

⁵Based on Table 2.49, the first ten digits of Cantor’s new number are $\sigma = 0.00010000\dots$

$d_i[n] \setminus x_n$	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}
	$\pi/10$	$\phi/10$	$\sqrt{2}/10$	$e/10$	$\sqrt{3}/10$	$\ln(2)/10$	$\sqrt{5}/10$	$\log_{10}(2)/10$	$\delta/10$	$\gamma/10$
d_1	3	1	1	2	1	0	2	3	4	5
d_2	1	6	4	7	7	6	2	0	6	7
d_3	4	1	1	1	3	9	3	4	9	6
d_4	1	8	2	8	2	3	6	3	1	1
d_5	5	0	1	2	0	1	0	0	2	2
d_6	9	3	3	8	5	4	4	1	3	0
d_7	2	3	5	1	0	7	7	4	1	7
d_8	6	9	6	8	7	8	9	6	4	8
d_9	5	8	2	1	2	5	8	9	2	5
d_{10}	3	8	3	8	0	5	2	4	8	6
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\ddots

Table 2.49: A hypothesized listing of the real numbers between zero and one. Shown are the first ten digits of some interesting real numbers in the interval $[0, 1)$, with the digits on the diagonal, $d_n[n]$, highlighted in red. You may have heard about the golden ratio, $\phi := (1 + \sqrt{5})/2 \approx 1.618033988749895\dots$ You can learn about δ and γ , the Feigenbaum constant and the Euler–Mascheroni constant, respectively, if you wish. They will not appear later in the course.

2. $a_n \leq a_{n+1}$ (the terms are monotonically increasing),
3. $a_n < 3$ (the terms are bounded from above).

The first property is true because $1 + \frac{1}{n} = \frac{n+1}{n}$ is rational and $(1 + \frac{1}{n})^n$ is therefore the product of n rational numbers. The second and third properties were given in Prop. 1.92 and proved at the end of the Chapter.

Because the terms are monotonically increasing and bounded from above, they have to “bunch up.” Formally, one can show that for all $\epsilon > 0$, there exists $N_\epsilon < \infty$ such that $|a_n - a_m| < \epsilon$ for all $n, m \geq N_\epsilon$. If they are “bunching up,” are they, in turn, “clustering around something” that has all the properties we want numbers to have (we can add and subtract them, multiply them, and even divide by them when they are non-zero)? The answer is YES. We’ve mentioned a few times that the real numbers, \mathbb{R} , are sufficiently complicated that they are typically treated in 400-level math courses, such as Michigan’s Math 451 Advanced Calculus I. The complication arises because one wants to construct the real numbers from the rational numbers in as convincing a fashion as the rational numbers are constructed from the integers. One of the constructions can be paraphrased as follows:

Secrets of the Arcane 2.50: What is a Real Number Anyway?

Every real number $x \in \mathbb{R}$ can be written as $x = \lim_{n \rightarrow \infty} a_n$, where for all $n \geq 1$,

1. a_n is rational (we are building new numbers from existing numbers),
2. $a_n \leq a_{n+1}$ (the terms are monotonically increasing),
3. there is a finite integer K (independent of n) such that $a_n \leq K$ (the terms are bounded from above).

Note: Hence, (2.36) is how Euler’s number is constructed! It is literally built from the ground up from the rational numbers. As long as we start the bisection algorithm with rational lower and upper bounds (bracketing points), all updates will have rational lower and upper bounds; hence, the construction of $\sqrt{2}$ can be achieved through two monotonic sequences of rational numbers, one that is monotonically increasing and the other decreasing. Archimedes’ construction of π cannot be executed using rational numbers. However, it can be done using “algebraic numbers”, such as $\frac{\sqrt{2}}{2}$, which are built up from the rational

numbers. **Being able to understand constructions like these is a huge benefit of an undergraduate course on real analysis.**

2.10 (Optional Read:) Where to Learn Proofs at Michigan

- Demonstration: Proof Beyond the Possibility of Doubt are the notes for the 1-credit Pass-Fail course, **Math 201 Introduction to Mathematical Writing**, a laboratory course in writing mathematics, and in writing down proofs. The laboratory focuses on learning the fundamentals of mathematical writing, rather than on the mathematics itself. It uses worksheets, which are given in the course text. All work is done during the class time. This course was designed in major part by Prof. Stephen DeBacker, developed since 2009.
- **EECS 203 Discrete Mathematics:** Introduction to the mathematical foundations of computer science. Topics covered include: **propositional and predicate logic, proof methods, set theory, functions, sequences and summations**, introduction to algorithms, growth of functions and asymptotic notation, **mathematical induction**, counting, pigeonhole principle, permutations and combinations, relations and their properties, closure of relations, partial orderings, and introduction to graph theory.
- **Math 217 - Linear Algebra (with proofs):** The topics covered include: systems of linear equations; matrix algebra; vectors, vector spaces, and subspaces; geometry of R^n ; linear dependence, bases, and dimension; linear transformations; Eigenvalues and Eigenvectors; diagonalization; inner products. **Throughout there will be an emphasis on the concepts, logic, and methods of theoretical mathematics.**

2.11 (Optional Read:) Proofs Associated with the Chapter

Proof of Prop. 2.28 (Exponentials Grow Faster than any Monomial Function)

The following limits hold for real variables x and for integers n replacing x :

- For all $a > 1$ and integers $0 \leq m < \infty$, $\lim_{x \rightarrow \infty} \frac{a^x}{x^m} = \infty$.
- For all $a > 1$ and integers $0 \leq m < \infty$, $\lim_{x \rightarrow \infty} \frac{x^m}{a^x} = 0$.
- For all $a > 1$ and integers $0 \leq m < \infty$, $\lim_{x \rightarrow \infty} x^m \cdot a^x = \infty$.

We only prove the second statement because the first one is the reciprocal of the second, and the third one is clear because, for $x > 1$, $x^m \cdot a^x \geq a^x$, and hence, $\lim_{x \rightarrow \infty} x^m \cdot a^x \geq \lim_{x \rightarrow \infty} a^x = \infty$.

Because $a > 1$, we can write $a = 1 + b$ for some $b > 0$. Then, for all $k > m$, and all $(n+1) > x > n > 2k$, the Binomial Theorem 1.83

and Prop. 1.79 gives

$$\begin{aligned}
0 &\leq \frac{x^m}{a^x} \\
&= \frac{x^m}{(1+b)^x} \\
&\leq \frac{x^m}{(1+b)^n} \quad (\text{because } n < x \text{ and } 1+b \geq 1) \\
&\leq \frac{x^m}{\sum_{k=0}^n \binom{n}{k} b^k} \quad (\text{applying the Binomial Theorem 1.83}) \\
&< \frac{(n+1)^m}{\binom{n}{k} b^k} \quad (\text{keeping only one term, and everything dropped is strictly } > 0) \\
&\leq \frac{x^m}{\frac{n^k}{2^k k!} \cdot b^k} \quad (\text{applying (1.43) from Prop. 1.79}) \\
&\leq \frac{2^k k!}{b^k} \cdot \frac{x^m}{n^k} \quad (\text{grouping terms independent of } n \text{ and those dependent on } n) \\
&\leq \frac{2^k k!}{b^k} \cdot \frac{(n+1)^m}{n^k} \quad (\text{because } (n+1) > x).
\end{aligned} \tag{2.37}$$

Because $k > m$, by Prop. 2.21,

$$\lim_{n \rightarrow \infty} \frac{(n+1)^m}{n^k} = 0,$$

which completes the proof. ■

Sketch of the Proof of Prop. 2.11 (Power Sums)

For all $n \geq 1$ and $k \geq 1$, the following holds,

$$\begin{aligned}
\sum_{i=1}^n i^k &:= 1^k + 2^k + 3^k + \dots + n^k \\
&= \frac{n^{k+1}}{k+1} + \frac{n^k}{2} + O(n, k-1),
\end{aligned} \tag{2.38}$$

where $O(n, k-1)$ stands for a polynomial in the variable n with degree less than or equal to $k-1$ (i.e., all of the lower-order terms).

The proof is by induction on n with $k \geq 1$ a fixed counting number. We define

$$\begin{aligned}
L(n) &:= \sum_{i=1}^n i^k \quad \text{and} \\
R(n) &:= \frac{n^{k+1}}{k+1} + \frac{n^k}{2}.
\end{aligned} \tag{2.39}$$

Our goal is to show that if for some $n \geq 1$, $L(n) = R(n) + O(n, k-1)$, then it holds that

$$L(n+1) = R(n+1) + O(n+1, k-1). \tag{2.40}$$

We first evaluate $L(n+1)$. By its definition and the induction hypothesis,

$$\begin{aligned}
L(n+1) &= L(n) + (n+1)^k \quad (\text{by the definition of } L(n)) \\
&= R(n) + O(n, k-1) + (n+1)^k \quad (\text{because, we assume } L(n) = R(n) + O(n, k-1)) \\
&= \frac{n^{k+1}}{k+1} + \frac{n^k}{2} + O(n, k-1) + (n+1)^k \quad (\text{using the definition of } R(n)) \\
&= \frac{n^{k+1}}{k+1} + \frac{n^k}{2} + O(n, k-1) + n^k + kn^{k-1} + O(n, k-2) \quad (\text{applying the Binomial Theorem}) \\
&= \frac{n^{k+1}}{k+1} + \frac{3n^k}{2} + O(n, k-1),
\end{aligned}$$

where we noted that $kn^{k-1} + O(n, k-1) + O(n, k-2) = O(n, k-1)$.

Next, we evaluate $R(n+1)$. Using its definition and expanding out terms,

$$\begin{aligned}
R(n+1) &:= \frac{(n+1)^{k+1}}{k+1} + \frac{(n+1)^k}{2} \\
&= \frac{n^{k+1} + (k+1)n^k + O(n, k-1)}{k+1} + \frac{n^k + kn^{k-1} + O(n, k-2)}{2} \quad (\text{applying the Binomial Theorem twice}) \\
&= \frac{n^{k+1}}{k+1} + \frac{3n^k}{2} + O(n, k-1),
\end{aligned}$$

where we noted that

$$\frac{O(n, k-1)}{k+1} + \frac{kn^{k-1} + O(n, k-2)}{2} = O(n, k-1).$$

Because⁶ $O(n+1, k-1) = O(n, k-1)$, we have shown that (2.40) holds and the proof is complete. ■

Proof of Prop. 2.34 (Limits of Linear Combinations)

Suppose that for all $1 \leq i \leq n$, $f_i : (0, \infty) \rightarrow \mathbb{R}$ has a finite limit at infinity, that is, $\lim_{x \rightarrow \infty} f_i(x)$ exists and is finite. Then, for all coefficients $\alpha_i \in \mathbb{R}$,

$$\lim_{x \rightarrow \infty} \left(\sum_{i=1}^n \alpha_i f_i(x) \right) = \sum_{i=1}^n \left(\alpha_i \lim_{x \rightarrow \infty} f_i(x) \right); \quad (2.41)$$

the limit of the linear combination exists and equals the linear combination of the limits.

Proof: We can assume without loss of generality that for all $1 \leq i \leq n$, $|\alpha_i| > 0$, because otherwise, the corresponding term(s) can be dropped from the sum. Let $\epsilon > 0$ be given and define $\epsilon_i := \frac{\epsilon}{|\alpha_i|^n} > 0$, $L_i := \lim_{x \rightarrow \infty} f_i(x)$. Then, because the individual limits exist, there exists N_i finite such that, for all $x \geq N_i$, $|f_i(x) - L_i| \leq \epsilon_i$.

⁶The Binomial Theorem implies that $(n+1)^{k-1}$ is a polynomial in the variable n with highest power $k-1$. For the proof to be rock solid, we need to note that the terms from the Binomial Theorem that are being added to $O(n, k-1)$ have coefficients that depend on k , which stays fixed, and not n , which can get arbitrarily large. In fact, this follows from $\binom{k-1}{j} := \frac{(k-1)!}{j!(k-1-j)!}$, $0 \leq j \leq k-1$.

Next, for $N := \max\{N_1, \dots, N_n\}$ and $x \geq N$, it follows that

$$\begin{aligned}
\left| \sum_{i=1}^n \alpha_i f_i(x) - \sum_{i=1}^n \alpha_i L_i \right| &= \left| \sum_{i=1}^n \alpha_i (f_i(x) - L_i) \right| \\
&\leq \sum_{i=1}^n |\alpha_i| |f_i(x) - L_i| \\
&\leq \sum_{i=1}^n |\alpha_i| \epsilon_i \\
&\leq \sum_{i=1}^n |\alpha_i| \frac{\epsilon}{|\alpha_i| n} \\
&\leq \epsilon,
\end{aligned}$$

which proves the result. ■

Sketch of Proof of Prop. 2.35 (Limits of Products and Ratios)

Suppose that $g : (0, \infty) \rightarrow \mathbb{R}$ has a limit at infinity of one, i.e., $\lim_{x \rightarrow \infty} g(x) = 1.0$. Then for any function $f : (0, \infty) \rightarrow \mathbb{R}$,

$$\begin{aligned}
\lim_{x \rightarrow \infty} f(x) \cdot g(x) &= \lim_{x \rightarrow \infty} f(x), \text{ and} \\
\lim_{x \rightarrow \infty} \frac{f(x)}{g(x)} &= \lim_{x \rightarrow \infty} f(x).
\end{aligned}$$

Proof: We will prove that case when both limits exist and are finite, leaving the other cases to the learner. We define $L := \lim_{x \rightarrow \infty} f(x)$ and $M := \lim_{x \rightarrow \infty} g(x)$. Then, the key inequality is

$$\begin{aligned}
|f(x) \cdot g(x) - L \cdot M| &= |(f(x) - L) \cdot (g(x) - L) + L(g(x) - M) + M(f(x) - L)| \\
&\leq |(f(x) - L) \cdot (g(x) - L)| + |L| \cdot |g(x) - M| + |M| \cdot |f(x) - L| \\
&\leq |f(x) - L| \cdot |g(x) - L| + |L| \cdot |g(x) - M| + |M| \cdot |f(x) - L|
\end{aligned}$$

Let $\epsilon > 0$ and choose $N_\epsilon < \infty$ such that, for all $x \geq N_\epsilon$

$$|f(x) - L| \leq \sqrt{\frac{\epsilon}{3}}$$

$$|g(x) - M| \leq \sqrt{\frac{\epsilon}{3}}$$

$$|f(x) - L| \leq \frac{|M|}{3}$$

$$|g(x) - M| \leq \frac{|L|}{3}.$$

This is possible because of the assumptions on the limits of f and g existing and being finite. Then, by the key inequality, for all

$x \geq N_\epsilon$

$$\begin{aligned} |f(x) \cdot g(x) - L \cdot M| &\leq |f(x) - L| \cdot |g(x) - L| + |L| \cdot |g(x) - M| + |M| \cdot |f(x) - L| \\ &\leq \frac{\epsilon}{3} + \frac{\epsilon}{3} + \frac{\epsilon}{3} \\ &\leq \epsilon, \end{aligned}$$

proving the first result.

To prove the second result, we must assume $|M| > 0$. The key inequality is then,

$$\begin{aligned} \left| \frac{f(x)}{g(x)} - \frac{L}{M} \right| &= \left| \frac{M f(x) - L g(x)}{M g(x)} \right| \\ &= \left| \frac{M (f(x) - L) + M L - L (g(x) - M) - M L}{M g(x)} \right| \\ &= \left| \frac{M (f(x) - L) - L (g(x) - M)}{M g(x)} \right| \\ &\leq \left| \frac{f(x) - L}{g(x)} \right| + \left| \frac{L}{M} \right| \left| \frac{g(x) - M}{g(x)} \right|. \end{aligned}$$

Let $0 < \epsilon < \left| \frac{M}{2} \right|$ and choose $N_\epsilon < \infty$ such that, for all $x \geq N_\epsilon$

$$|f(x) - L| \leq \frac{\epsilon |M|}{4}$$

$$|g(x) - M| \leq \left| \frac{M}{L} \right| \frac{\epsilon |M|}{4}.$$

Then, from the second main inequality,

$$\begin{aligned} \left| \frac{f(x)}{g(x)} - \frac{L}{M} \right| &\leq \left| \frac{f(x) - L}{g(x)} \right| + \left| \frac{L}{M} \right| \left| \frac{g(x) - M}{g(x)} \right| \\ &\leq \left| \frac{f(x) - L}{M/2} \right| + \left| \frac{L}{M} \right| \left| \frac{g(x) - M}{M/2} \right| \\ &\leq \frac{\epsilon}{2} + \frac{\epsilon}{2} \\ &\leq \epsilon, \end{aligned}$$

proving the result.

Analytical Background for Secrets of the Arcane 2.50: Why Euler's Number Exists

Suppose that x_k is a monotonically increasing sequence of real numbers that is bounded from above. Then, for every $\epsilon > 0$, there exists $N_\epsilon < \infty$ such that $m, n \geq N_\epsilon \implies |x_m - x_n| < \epsilon$.

In plain language, the terms at the tail of the sequence cluster arbitrarily closely to one another. $m, n \geq N_\epsilon$ defines the terms in the “tail” and $|x_m - x_n| < \epsilon$ defines what it means to cluster together.

Proof: We prove by induction the contrapositive⁷: if $\{x_k\}$ is a sequence of real numbers that is monotonically increasing, and there exists an $\epsilon > 0$ such that, for any natural number n , there exists a second natural number $m > n$ such that $x_m > x_n + \epsilon$, then the

⁷The logical proposition that $(A \implies B) \iff (\neg B \implies \neg A)$.

sequence $\{x_k\}$ is unbounded.

Assumption 1: $\{x_k\}$ is a monotonically increasing sequence. This means that for all $k \in \mathbb{N}$, $x_{k+1} \geq x_k$.

Assumption 2: For some $\epsilon > 0$, for all $n \in \mathbb{N}$, there exists an $m \in \mathbb{N}$ with $m > n$ such that $x_m > x_n + \epsilon$.

Induction Hypothesis: Fix $\epsilon > 0$. For every $n \in \mathbb{N}$, there exists an index $m > n$ such that $x_m > x_1 + n\epsilon$.

Base Case: For $n = 1$, Assumption 2 directly gives us an $m > 1$ such that $x_m > x_1 + \epsilon$.

Inductive Step: Assume that for some $n \in \mathbb{N}$, there exists an index k such that $x_k > x_1 + n\epsilon$. By our assumption, there must exist an index $m > k$ such that $x_m > x_k + \epsilon$.

From the induction hypothesis and the inductive step, we have:

$$x_m > x_k + \epsilon > (x_1 + n\epsilon) + \epsilon = x_1 + (n+1)\epsilon.$$

This shows that for $n+1$, we have found an index m such that $x_m > x_1 + (n+1)\epsilon$, completing the inductive step.

Since for every n , we can find an m such that $x_m > x_1 + n\epsilon$, the sequence $\{x_k\}$ must be unbounded. For any proposed bound B , choose n such that $x_1 + n\epsilon > B$. The corresponding m from our inductive step guarantees that $x_m > B$, thus $\{x_k\}$ exceeds the bound. ■

Hence, if the sequence is monotonically increasing and bounded from above, the terms must cluster together. The real numbers can be constructed from rational sequences that satisfy these two properties. The existence of Euler's number, e , is a result of this construction.

Chapter 3

Definite Integration as the Signed Area Under a Curve

Learning Objectives

By the end of this chapter, the student should be able to:

- Define and explain the concept of a definite integral within the context of calculus.
- Execute accurate computations of definite integrals using appropriate mathematical techniques.
- Describe and apply the basic properties of definite integrals to solve problems.
- Recognize and demonstrate the applications of definite integrals in various engineering scenarios.

Outcomes

Upon successful completion of this chapter, students will be able to:

- Construct and calculate Riemann lower and upper sums to approximate definite integrals.
- Acquire insight into what kinds of functions can be integrated.
- Employ numerical algorithms used in engineering practice to compute definite integrals.
- Utilize integration techniques to infer changes in position from a given velocity function, particularly in robotic applications.
- Identify the mathematical origins of parabolic trajectories in ballistic motion.
- Determine the area enclosed between two functions.
- Compute essential parameters for robotic models, such as total mass and center of mass.
- Learn how to trick a single-variable integral into computing the volume of a solid of revolution.

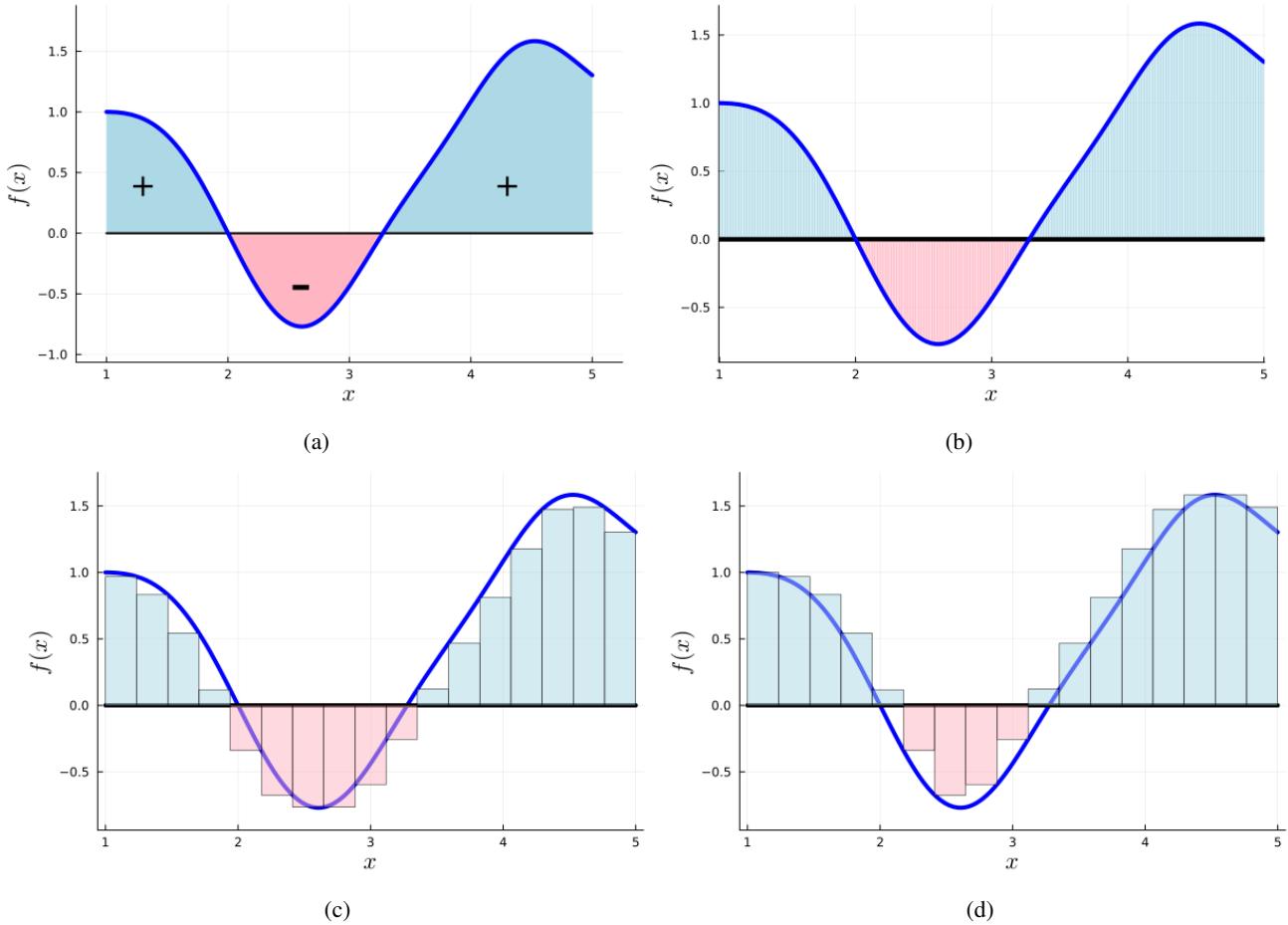


Figure 3.1: The Riemann Integral computes the area under a (decently, well-behaved¹) curve. (a) Shows that where the function is positive, the area is assigned a positive value, and where the function is negative, it is assigned a negative value. The Riemann Integral is the sum of the positive and negative areas. (c) Shows an underapproximation via rectangles, (d) shows an overapproximation, while in (b), we've used so many rectangles that you cannot tell if it is an under- or over-approximation. Ha!

The Riemann Integral is the result of applying the Approximation Principle to compute the “signed area under a curve”, as illustrated in Fig. 3.1. The term “signed area” is used to highlight that we assign a positive number to the area “under” positive values of the function and a negative number to the area “above” negative values of the function. Moving forward, we’ll use the generic terminology “AREA UNDER the CURVE” to cover both situations. The process of approximating the area should look familiar to you because we carried out a similar process in Chapter 2.4.2, for monotonically increasing functions, and we only dealt with “positive area”.

3.1 The Riemann Integral (aka, Riemann-Darboux Integral)

In this Chapter, we will address the more general situation of functions that are non-monotonic (aka, “wiggle around”) and take on both positive and negative values. For us to even define what we mean by the area under the curve, we’ll have to impose some conditions on the functions we consider. We will not address the “most general conditions known to humankind”, but we’ll build toward a general enough class of functions for “most engineering applications”. The fundamental idea remains the same we used for monotonically increasing functions in Chapter 2.4.2: we sandwich the area we wish to compute between an underapproximation of the area and an overapproximation of the area. We then seek to analyze when we can make the under- and over-approximations converge to one another in the limit as we take finer and finer rectangular slices. The limiting value is defined to be the area under the curve, that is, the Riemann Integral.

3.1.1 Simple Version of Riemann Lower and Upper Sums

As in Chapter 2.4.2, for a function $f : [a, b] \rightarrow \mathbb{R}$, where $a < b$, we divide the interval $[a, b]$ into $n > 1$ **evenly spaced subintervals**, $[x_i, x_{i+1}]$, where

$$\begin{aligned}\Delta x &:= \frac{b-a}{n} \\ x_i &:= a + (i-1) \cdot \Delta x,\end{aligned}\tag{3.1}$$

so that

$$a = x_1 < x_2 < \dots < x_n < x_{n+1} = b\tag{3.2}$$

$$x_{i+1} = x_i + \Delta x, \quad 1 \leq i \leq n.\tag{3.3}$$

We still want to define

$$\text{Area}_n^{\text{Low}} := \sum_{i=1}^n h_i^{\text{Low}} \cdot b_i \text{ and } \text{Area}_n^{\text{Up}} := \sum_{i=1}^n h_i^{\text{Up}} \cdot b_i,\tag{3.4}$$

where $b_i = \Delta x$ is the base of the rectangle and h_i^{Low} and h_i^{Up} are the heights. To form the “best” underapproximation we can, we want to assign to h_i^{Low} the **maximum** (aka, largest) value we can so that the rectangle fits under the function, similarly, to form the “best” overapproximation we can, we want to assign to h_i^{Up} the **minimum** (aka, smallest) value we can so that the rectangle fits over the function.

For monotonically increasing continuous functions, this was a triviality because we could take $h_i^{\text{Low}} = f(x_i)$, the function evaluated at the lower end of the interval because $f(x_i) \leq f(x)$ for all $x \in [x_i, x_{i+1}]$, and we could take $h_i^{\text{Up}} = f(x_{i+1})$, the function evaluated at the upper end of the interval because $f(x_{i+1}) \geq f(x)$ for all $x \in [x_i, x_{i+1}]$. In other symbols, for monotonically increasing continuous functions, we used

$$h_i^{\text{Low}} := \min_{x \in [x_i, x_{i+1}]} f(x) = f(x_i) \quad \text{and} \quad h_i^{\text{Up}} := \max_{x \in [x_i, x_{i+1}]} f(x) = f(x_{i+1}).\tag{3.5}$$

When the function wiggles up and down, we have to search for (or compute) the minimum and maximum values of $f : [x_i, x_{i+1}] \rightarrow \mathbb{R}$ instead of just using the function values at the endpoints. But with this one exception, the process remains the same!

Riemann Lower and Upper Sums

Definition 3.1. (Riemann (Approximating) Sums) For $a < b$ (finite) real numbers, suppose that $f : [a, b] \rightarrow \mathbb{R}$ is continuous. For $n \geq 1$, define $\Delta x := \frac{b-a}{n}$ and $x_i := a + (i-1) \cdot \Delta x$, for $1 \leq i \leq n+1$. Then,

(a) $\text{Area}_n^{\text{Low}} := \sum_{i=1}^n h_i^{\text{Low}} \cdot \Delta x$, for $h_i^{\text{Low}} := \min_{x \in [x_i, x_{i+1}]} f(x)$ is a **Riemann lower sum**, and

(b) $\text{Area}_n^{\text{Up}} := \sum_{i=1}^n h_i^{\text{Up}} \cdot \Delta x$, for $h_i^{\text{Up}} := \max_{x \in [x_i, x_{i+1}]} f(x)$ is a **Riemann upper sum**.

Notes: Riemann’s original formulation was more complicated and more general than the above. We’ll get to a more general formulation in due time. The simplified formulation we have followed mirrors the formulation of Darboux, though, as a Calculus novice, you “could (and should) care less!”

The division of $[a, b]$ as $a = x_1 < x_2 < \dots < x_n < x_{n+1} = b$ is called a **partition**. We have been using a **uniform partition** where $x_{i+1} - x_i = \Delta x$, for all $1 \leq i \leq n$. Darboux allowed for non-uniform partitions by requiring that

$$\max_{1 \leq i \leq n} (x_{i+1} - x_i) \xrightarrow{n \rightarrow \infty} 0.$$

This is immediate for a uniform partition because $\max_{1 \leq i \leq n} (x_{i+1} - x_i) = \Delta x = \frac{b-a}{n} \xrightarrow{n \rightarrow \infty} 0$. Some theorems about Riemann Integrals only work when non-uniform partitions are allowed. At the University of Michigan, MATH 451 teaches this level of detail. You will not find it in the 100- and 200-level courses.

Below is a Julia implementation to estimate the Riemann Lower and Upper Sums. They are estimations because we are NOT computing the exact minima and maxima of the function $\sin[x_i, x_{i+1}] \rightarrow \mathbb{R}$. Instead, we introduce `nRefine` additional points in

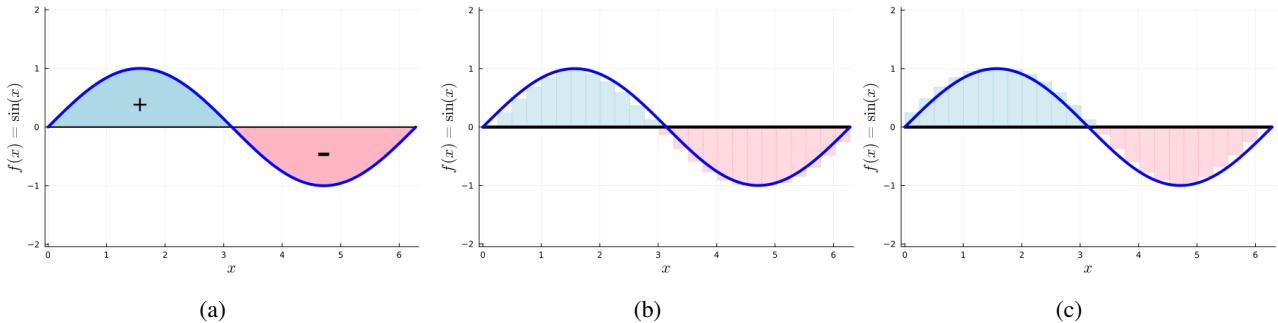


Figure 3.2: An illustration of Riemann Lower and Upper Sums for $\sin : [0, 2\pi] \rightarrow \mathbb{R}$. (a) Shows that the area we are approximating should equal to zero, because the area in pink appears to be the negative of the area in blue. A lower sum is shown in (b) and an upper sum is shown in (c). The number of intervals is $n = 25$.

between $[x_i, x_{i+1}]$, evaluate the function at the array of points

```
xRefine = range(x[i], x[i+1], length=ceil(Int64, nRefine)),
```

and then apply Julia's minimum and maximum functions to the vector $f(xRefine)$, which contains the function's values at the two endpoints, x_i and x_{i+1} , plus any additional points specified by $nRefine \geq 2$.

Be sure to review (1.9) for how the estimate and the error bound are defined!

```

1 function riemannLowerUpperSums(f, a, b, N=100, nRefine=4, aTol=1e-6)
2     if abs(a-b) < aTol
3         println("abs(a-b) < $aTol. Estimated integral is zero.")
4         return 0.0
5     end
6     if nRefine < 2
7         nRefine = 2
8         println("nRefine was set to its minimum value of $nRefine")
9     end
10    # make sure N is an integer greater or equal to 2
11    N=max(ceil(Int64,N),2)
12    # Define \Delta x
13    Deltax = (b-a)/N
14    # Both sums start at zero
15    LowerSum = 0.0
16    UpperSum = 0.0
17    # Initialize x_i
18    xi=a
19    for i = 1:N
20        # x_{i+1}
21        xip1 = xi + Deltax
22        # add in extra points so that we can find an estimate of the max and min
23        xRefine = range(xi, xip1, length=ceil(Int64, nRefine))
24        fRefined = f.(xRefine)
25        fMax = maximum(fRefined)
26        fMin = minimum(fRefined)
27        # Now that we have max and min for f, we can compute upper and lower sums
28        UpperSum = UpperSum + fMax*Deltax
29        LowerSum = LowerSum + fMin*Deltax
30        # Update x_i

```

```

31     xi = xip1
32
33     # @show b-xi # for debugging
34     estIntegral = (UpperSum + LowerSum)/2.0
35     pmError = (UpperSum - LowerSum)/2
36     return (LowerSum=LowerSum, UpperSum=UpperSum, estIntegral = estIntegral, pmError = pmError)
37 end

```

Output

riemannLowerUpperSums (generic function with 4 methods)

Two illustrations follow.

```

1 f(x) = sin(x)
2 a = 0.0
3 b = 2*pi
4 N=25
5 F=riemannLowerUpperSums(f, a, b, N)
6 @show F.lowerSum
7 @show F.upperSum;

```

Output

```
F.lowerSum = -0.5021037670817627
F.upperSum = 0.502103767081769
```

Note the lower and upper Riemann sums are not close to one another.

```

1 f(x) = sin(x)
2 a = 0.0
3 b = 2*pi
4 N=250
5 F=riemannLowerUpperSums(f, a, b, N)
6 @show F.lowerSum
7 @show F.upperSum;

```

Output

```
F.lowerSum = -0.05046634766507299
F.upperSum = 0.05046634766507339
```

The lower and upper Riemann sums are now closer to one another. With even more rectangles, $N = 2,500$, we obtain $F.lowerSum \approx -0.00503$ and $F.upperSum \approx 0.00503$.

3.1.2 Riemann Integral of a Continuous Function over a Bounded Interval

Riemann Integral

Definition 3.2. (Riemann Integral for a Continuous Function over a (Bounded) Interval $[a, b]$, $a < b$, Finite) If $\lim_{n \rightarrow \infty} \text{Area}_n^{\text{Low}}$ and $\lim_{n \rightarrow \infty} \text{Area}_n^{\text{Up}}$ both exist, are finite, and equal one another, then the **Riemann Integral** of $f : [a, b] \rightarrow \mathbb{R}$ exists and is equal to the common limit. In other symbols,

$$\int_a^b f(x) dx := \begin{cases} \lim_{n \rightarrow \infty} \text{Area}_n^{\text{Low}} \\ \lim_{n \rightarrow \infty} \text{Area}_n^{\text{Up}} \end{cases} \quad \text{provided both limits exist, are finite, and are equal.} \quad (3.6)$$

Notes:

- The symbol \int is an old-fashioned “lazy S ” and comes from the fact that an integral is a sum.
- The symbol \int_a^b stands for “the sum from a to b ” and is read “**the integral from a to b** ”.
- a is called the **lower limit** and b the **upper limit**. Together, they are called **limits of integration**.
- The function being integrated is called the **integrand**.
- The Riemann integral exists for bounded continuous functions on bounded intervals.
- The bounded intervals can be open (a, b) “with endpoints excluded”, semi-open $(a, b]$ or $[a, b)$ “with one endpoint excluded, or closed $[a, b]$ “with endpoints included”. In all cases, $\int_a^b f(x) dx$ gives the same value.
- More general functions admitting a Riemann Integral are discussed in Chapter 4.7 “When does a Function have a Definite Riemann Integral?”.

```
1 f(x) = sin(x)
2 a = 0.0
3 b = 2*pi
4 N=250
5 F=riemannLowerUpperSums(f, a, b, N)
```

Output

```
(lowerSum = -0.05026327759307434, upperSum = 0.05026327759309964,
estIntegral = 1.2649603586822877e-14, pmError = 0.05026327759308699)
```

```
1 f(x) = sin(x)
2 a = 0.0
3 b = 2*pi
4 N=2500
5 F=riemannLowerUpperSums(f, a, b, N)
```

Output

```
(lowerSum = -0.005026548246017541, upperSum = 0.0050265482454683935,
estIntegral = -2.74573766501085e-13, pmError = 0.005026548245742967)
```

In case you did not review (1.9) for how the estimate and the error bounds are defined, they are defined by

- $\text{estIntegral} = (\text{upperSum} + \text{lowerSum})/2$
- $\text{pmError} = (\text{upperSum} - \text{lowerSum})/2$.

Even though the estimated integral `estIntegral` may look amazingly good, you must look at the error bound `pmError` to understand how much faith to put in the answer. The error bound tells you that with 250 rectangles, $\text{pmError} \approx \pm 0.05$, which is not so great, while with 2,500 rectangles, the error bound is approximately ± 0.005 .

It is straightforward to numerically integrate complicated functions that you could never integrate by hand.

```

1 f(x) = cos(exp(-x)*sin(x)/(1 + x + 5x^2)) + log2((1 + x)/cos(x) + 5x^2))
2 a = 0.0
3 b = 10
4 N = 5000
5 F = riemannLowerUpperSums(f, a, b, N)

```

Output

(`lowerSum` = -28.837938067489223, `upperSum` = -28.826646469106663,
`estIntegral` = -28.83229226829794, `pmError` = 0.005645799191279721)

$$\int_0^{10} \cos\left(\frac{\exp(-x)\sin(x)}{1+x+5x^2}\right) + \log_2\left(\frac{1+x}{\cos(x)+5x^2}\right) dx \approx -28.832292 \pm 0.005646. \quad (3.7)$$

Note: Here is what ChaptGPT4 with the Wolfram plugin gives: “I’m sorry, but the computation for the integral took too long and timed out. This is likely due to the complexity of the function you’re trying to integrate. In such cases, it might be helpful to try a numerical approximation method like Simpson’s rule, or to use a specialized numerical integration package. However, please note that these methods might not always provide an exact result, especially for complex functions. If you’re working in a programming environment like Python or Julia, you might want to use a numerical integration function provided by a library like SciPy (for Python) or QuadGK (for Julia). If you’re doing this by hand or in a simpler calculator, you might need to simplify the function or break the integral into smaller parts that can be computed separately.”

3.1.3 Illustration of the Riemann Indefinite Integral of a Monomial

If we replace one of the limits of integration with a variable, say x , meaning we are integrating over $[a, x]$ or $[x, b]$, the resulting integral produces a function of x . Because x is not a “definite real number”, the integral is called an **indefinite integral**. In other words, if $f : [a, b] \rightarrow \mathbb{R}$ is Riemann integrable, then $g : [a, b] \rightarrow \mathbb{R}$ by

$$g(x) := \int_a^x f(y) dy$$

is called an **indefinite integral**.

Proposition 3.3: Integrating x^k

For all counting numbers k and for all $x > 0$ finite,

$$\int_0^x y^k dy = \frac{x^{k+1}}{k+1}. \quad (3.8)$$

Note: For the moment, we only know how to integrate functions that can be built from monomials. In Prop. 4.24, we’ll have enough tools to show $\int_0^x e^y dy = e^x - 1$.

Proof: We will work this one out in full detail so as to further ingrain the **Approximation Principle** that underlies most of Calculus. We let $n > 1$ be an integer and define $\Delta y := \frac{x-0}{n}$ along with $y_i := a + (i-1)\Delta y$, yielding

$$a = y_1 < y_2 < \cdots < y_n < y_{n+1} = b.$$

Because $f(y) = y^k$ is increasing, $f(y_i) \leq f(y) \leq f(y_{i+1})$ for $y \in [y_i, y_{i+1}]$. This makes computing the min and max of the function over subintervals very easy.

Riemann Lower Sum: We underestimate the total area between $[0, x]$ by summing up the underapproximations, $f(y_i) \cdot \Delta y$, for each interval $[y_i, y_{i+1}]$,

$$\begin{aligned} \text{Area}_n^{\text{Low}} &:= \sum_{i=1}^n f(y_i) \cdot \Delta y \\ &= \sum_{i=1}^n [(i-1)\Delta y]^k \cdot \Delta y \quad (\text{substituting in terms}) \\ &= \sum_{i=1}^n (i-1)^k \cdot (\Delta y)^k \cdot \Delta y \quad (\text{algebra}) \\ &= (\Delta y)^{k+1} \cdot \sum_{i=1}^n (i-1)^k \quad (\text{take terms that do not depend on } i \text{ outside the sum}) \\ &= \left(\frac{x}{n} \right)^{k+1} \cdot \sum_{i=1}^{n-1} i^k \quad (\text{change of index and algebra}) \\ &= \frac{x^{k+1}}{n^{k+1}} \cdot \left(\frac{(n-1)^{k+1}}{k+1} + \frac{(n-1)^k}{2} + O(n-1, k-1) \right) \quad (\text{Prop. 2.11}). \end{aligned} \tag{3.9}$$

From our results in Chapter 2.6.1, we have that

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{x^{k+1}}{n^{k+1}} \cdot \left(\frac{(n-1)^{k+1}}{k+1} + \frac{(n-1)^k}{2} + O(n-1, k-1) \right) &= x^{k+1} \cdot \lim_{n \rightarrow \infty} \left(\frac{(n-1)^{k+1}}{(k+1) \cdot n^{k+1}} + \frac{(n-1)^k}{2 \cdot n^{k+1}} + \frac{O(n-1, k-1)}{n^{k+1}} \right) \\ &= \frac{x^{k+1}}{k+1}. \end{aligned} \tag{3.10}$$

We next repeat this process.

Riemann Upper Sum: We overestimate the total area between $[0, x]$ by summing up the overapproximations, $f(y_{i+1}) \cdot \Delta y$, for each interval $[y_i, y_{i+1}]$,

$$\begin{aligned} \text{Area}_n^{\text{Up}} &:= \sum_{i=1}^n f(y_{i+1}) \cdot \Delta y \\ &= \sum_{i=1}^n [i\Delta y]^k \cdot \Delta y \quad (\text{substituting in terms}) \\ &= \sum_{i=1}^n i^k \cdot (\Delta y)^k \cdot \Delta y \quad (\text{algebra}) \\ &= (\Delta y)^{k+1} \cdot \sum_{i=1}^n i^k \quad (\text{take terms that do not depend on } i \text{ outside the sum}) \\ &= \left(\frac{x}{n} \right)^{k+1} \cdot \sum_{i=1}^n i^k \quad (\text{algebra}) \\ &= \frac{x^{k+1}}{n^{k+1}} \cdot \left(\frac{n^{k+1}}{k+1} + \frac{n^k}{2} + O(n, k-1) \right) \quad (\text{Prop. 2.11}). \end{aligned} \tag{3.11}$$

From our results in Chapter 2.6.1, we have that

$$\lim_{n \rightarrow \infty} \frac{x^{k+1}}{n^{k+1}} \cdot \left(\frac{n^{k+1}}{k+1} + \frac{n^k}{2} + O(n, k-1) \right) = \frac{x^{k+1}}{k+1}. \tag{3.12}$$

Because the limits of the Upper and Lower Riemann Sums both exist and are equal, by Def. 3.2, the Riemann Integral of y^k over the interval $[0, x]$ is defined and equals $\frac{x^{k+1}}{k+1}$; that is

$$\int_0^x y^k dy = \frac{x^{k+1}}{k+1}. \quad (3.13)$$

■

3.1.4 Are all Functions Riemann Integrable?

The answer is NO! Consider this wild function, $f : [0, 1] \rightarrow \mathbb{R}$ by

$$f(x) := \begin{cases} 1 & x \text{ is rational} \\ 0 & x \text{ is irrational.} \end{cases} \quad (3.14)$$

Then, for any interval $0 \leq x_i < x_{i+1} \leq 1$,

$$\min_{x \in [x_i, x_{i+1}]} f(x) = 0.0 \text{ and } \max_{x \in [x_i, x_{i+1}]} f(x) = 1.0 \quad (3.15)$$

It follows that for all n , the Riemann Lower Sum, $\text{Area}_n^{\text{Low}} = 0.0$ (is identically equal to zero), while for all n , the Riemann Upper Sum, $\text{Area}_n^{\text{Up}} = 1.0$. (is identically one). Hence, even though both limits exist, they are not equal. The function is, therefore, not Riemann Integrable. In the next Chapter, we'll nail down a better understanding of what types of functions have a Riemann Integral.

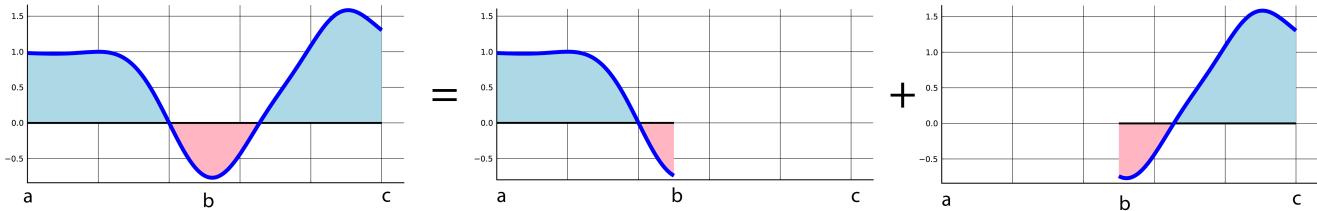


Figure 3.3: The signed area under the curve from $x = a$ to $x = c$ is equal to the signed area from $x = a$ to $x = b$ plus the signed area from $x = b$ to $x = c$. Hence, when the Riemann Integral exists, it can be computed as $\int_a^c f(x) dx = \int_a^b f(x) dx + \int_b^c f(x) dx$.

3.2 A Few Properties of the Riemann Integral

Here are a few properties that are often helpful when evaluating Riemann Integrals.

3.2.1 A Basic Additivity Property of Area Under a Curve

Figure 3.3 illustrates a general additivity property of area under a curve: if the interval is broken up into disjoint pieces, then the individual signed areas can be added up to get the total signed area. This translates into a nice property of the Riemann Integral.

Proposition 3.4: First Additivity Property of the Riemann Integral

Suppose that $a < b < c$ and $f : [a, c] \rightarrow \mathbb{R}$ is continuous. Then the three Riemann Integrals $\int_a^c f(x) dx$, $\int_a^b f(x) dx$, and $\int_b^c f(x) dx$ all exist and satisfy

$$\int_a^c f(x) dx = \int_a^b f(x) dx + \int_b^c f(x) dx. \quad (3.16)$$

Remark 3.5. There is nothing to stop you from using a different number of rectangles when approximating $\int_a^b f(x) dx$ versus $\int_b^a f(x) dx$. In fact, “professional” packages for numerically approximating integrals use “adaptive refinement” where they selectively refine some subintervals $[x_i, x_{i+1}]$ more than others when approximating area under a curve. We discuss later in (3.31) a criterion that can be used for deciding where to use a finer partition.

3.2.2 Dummy Variables of Integration

A variable that appears in a calculation only as a placeholder and disappears completely in the final result is called a **Dummy Variable**. As an example, in (2.20) we computed an estimate of the area under a parabola $y = x^2$ for the interval $[0, 1]$ using n rectangles and obtained $\text{Area}_n^{\text{est}} = \left(\frac{1}{3} + \frac{1}{6n^2}\right) \pm \frac{1}{n}$. Nowhere does the variable x appear in the final answer. Hence, it is a dummy variable. We could have similarly computed the area under the parabola $\beta = \alpha^2$, and we would have obtained the same answer, $\text{Area}_n^{\text{est}} = \left(\frac{1}{3} + \frac{1}{6n^2}\right) \pm \frac{1}{n}$. Similarly, when we computed the Riemann Integral of y^k over the interval $[0, x]$, we obtained an answer that depended on x , the upper limit, but not on the variable y .

Proposition 3.6: Dummy Variables

Suppose that $a < b$ and $f : [a, b] \rightarrow \mathbb{R}$ is continuous. Then,

$$\int_a^b f(x) dx = \int_a^b f(y) dy = \int_a^b f(\lambda) d\lambda, \quad (3.17)$$

for any choice of real variables x, y , and λ . The Riemann Integral depends on the function $f : [a, b] \rightarrow \mathbb{R}$ and the lower and upper limits. The integral does not depend on what we name the variable we are plugging into f . In other words, the variable under the integral sign is a **dummy variable**.

Note: You may enjoy the video Why the “dummy variable” of a definite integral doesn’t matter by **BlackPenRedPen**.

3.2.3 Integrating Linear Combinations

Proposition 3.7: Linearity of the Riemann Integral

Suppose that $a < b$ and each of the functions $1 \leq i \leq N$, $f_i : [a, b] \rightarrow \mathbb{R}$, is Riemann integrable (i.e., their Riemann Integral exists and is finite). Then for constants $\alpha_i \in \mathbb{R}$, the function

$$f(x) := \alpha_1 \cdot f_1(x) + \alpha_2 \cdot f_2(x) + \cdots + \alpha_N \cdot f_N(x),$$

defined as a “linear combination of the functions $\{f_1, f_2, \dots, f_N\}$ ”, is Riemann integrable and, moreover,

$$\int_a^b f(x) dx = \alpha_1 \int_a^b f_1(x) dx + \alpha_2 \int_a^b f_2(x) dx + \cdots + \alpha_N \int_a^b f_N(x) dx. \quad (3.18)$$

Note: Because we know how to integrate monomials and linear combinations, we now know how to integrate any polynomial.

3.2.4 Making Sense of an Integral when its Lower Limit is Greater than its Upper Limit

How to make sense of $\int_a^b f(x) dx$ when $a > b$? The interval $[a, b]$ does not make sense because $\{x \in \mathbb{R} \mid a \leq x \leq b\} = \emptyset$; so what gives? By convention, when $a > b$, the Riemann Integral from a to b is defined as the negative of the integral from b to a . This is compatible with our approach to integration via Riemann sums because when $a > b$, $\Delta x := \frac{b-a}{n} = -\frac{a-b}{n} < 0$, and $x_i := a + (i-1)\Delta x$ satisfy $b = x_{n+1} < x_n < \cdots < x_2 < x_1 = a$. In symbols, we have, if $b < a$ and $f : [b, a] \rightarrow \mathbb{R}$ is continuous, then

$$\int_a^b f(x) dx := - \int_b^a f(x) dx. \quad (3.19)$$

3.2.5 Generalized First Additivity Property of Integrals

We build on Prop. 3.4, which we restate as: for $a < c < b$ and $f : [a, b] \rightarrow \mathbb{R}$ continuous,

$$\int_a^b f(x) dx = \int_a^c f(x) dx + \int_c^b f(x) dx.$$

It's useful to know that the relative order of a , b , and c is immaterial (aka, does not matter).

Corollary 3.8: Generalized First Additivity Property

For a , b , and c finite real numbers in no particular order, define $\alpha := \min\{a, b, c\}$ to be the smallest of the three numbers, and $\beta := \max\{a, b, c\}$ the largest. If $f : [\alpha, \beta] \rightarrow \mathbb{R}$ is continuous, then

$$\int_a^b f(x) dx = \int_a^c f(x) dx + \int_c^b f(x) dx. \quad (3.20)$$

Example 3.9. For a and b both finite real numbers, determine $\int_a^b y^k dy$.

Solutions: Because $(\bullet)^k : \mathbb{R} \rightarrow \mathbb{R}$ is continuous, Corollary 3.8 yields,

$$\begin{aligned} \int_a^b y^k dy &= \int_a^0 y^k dy + \int_0^b y^k dy \\ &= - \int_0^a y^k dy + \int_0^b y^k dy \\ &= -\frac{a^{k+1}}{k+1} + \frac{b^{k+1}}{k+1} = \frac{b^{k+1} - a^{k+1}}{k+1}. \end{aligned}$$

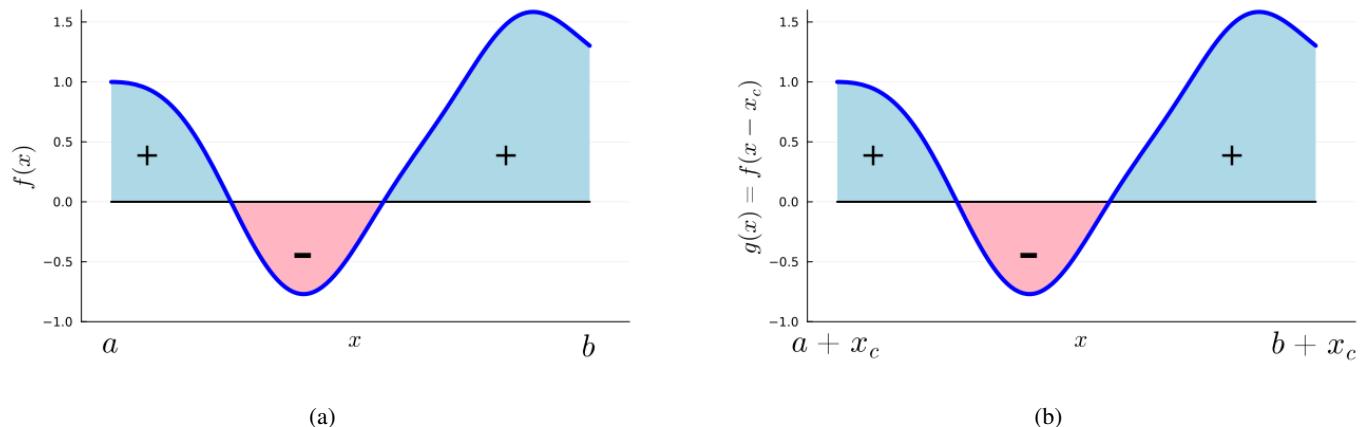


Figure 3.4: (**Shifting does not Change Area:**) A plot of $f : [a, b] \rightarrow \mathbb{R}$ is shown in (a) while (b) shows a plot of a shifted version of f , namely, $g : [a + x_c, b + x_c] \rightarrow \mathbb{R}$, where $g(x) := f(x - x_c)$. Clearly, the area under each function is the same, **when the limits of integration shift with the function**, as shown in the plots. Note that $g(a + x_c) = f(a + x_c - x_c) = f(a)$; similarly, $g(b + x_c) = f(b)$.

3.2.6 Shifting and Integration

Chapter 1.16 introduced the idea of shifting a function along the x -axis. It also noted that the area enclosed by rectangles and triangles does not change under shifts (said another way, their area is **invariant** under the operation of shifting). Because the Riemann-Darboux integral is based on the area of rectangles, it inherits this invariance property, as illustrated in Fig. 3.4.

Proposition 3.10: Shift Property

Suppose that $a < b$, $f : [a, b] \rightarrow \mathbb{R}$ is continuous, and $x_c \in \mathbb{R}$. Then $g : [a + x_c, b + x_c] \rightarrow \mathbb{R}$ by $g(x) := f(x - x_c)$ is also continuous and $\int_a^b f(x) dx = \int_{a+x_c}^{b+x_c} g(x) dx$. That is,

$$\int_a^b f(x) dx = \int_{a+x_c}^{b+x_c} f(x - x_c) dx. \quad (3.21)$$

Note: By defining $\alpha := a + x_c$ and $\beta := b + x_c$, the above result can be rewritten as

$$\boxed{\int_{\alpha}^{\beta} f(x - x_c) dx = \int_{\alpha-x_c}^{\beta-x_c} f(x) dx.} \quad (3.22)$$

The latter form is usually easier to apply.

Example 3.11. Evaluate the following definite integrals.

$$(a) \int_0^1 (x - 2)^2 dx$$

$$(b) \int_0^1 (x - x_c)^2 dx$$

$$(c) \int_0^1 \text{rect}(x - 3) dx$$

$$(d) \int_3^5 \text{tri}(x - 4) dx$$

Solution:

$$(a) \text{ Ans. } \int_0^1 (x - 2)^2 dx = 2\frac{1}{3}.$$

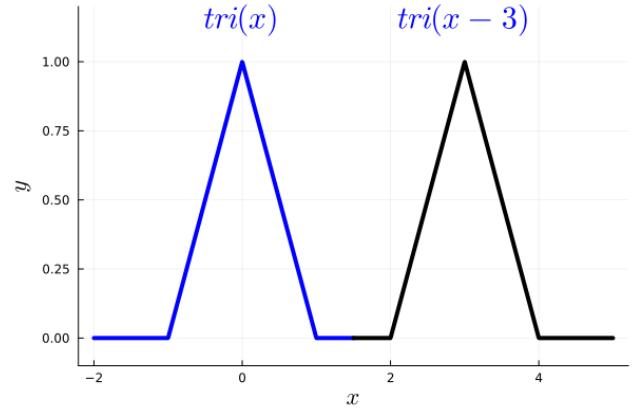
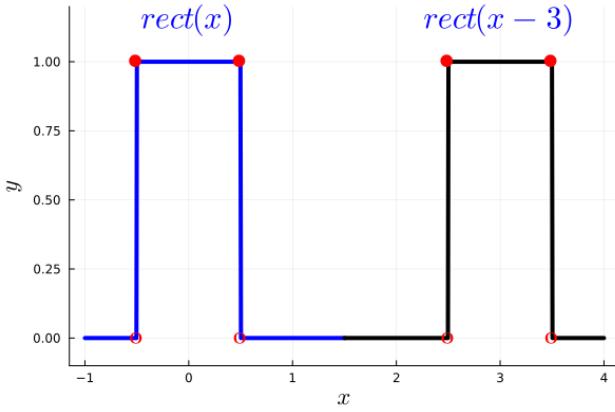
Applying Prop. 3.10,

$$\begin{aligned} \int_0^1 (x - 2)^2 dx &= \int_{0-2}^{1-2} x^2 dx \\ &= \int_{-2}^{-1} x^2 dx \\ &= \left. \frac{x^3}{3} \right|_{-2}^{-1} \\ &= \frac{(-1)^3}{3} - \frac{(-2)^3}{3} \\ &= -\frac{1}{3} + \frac{8}{3} \\ &= 2\frac{1}{3}. \end{aligned}$$

The solution is long because we gave every step. Normally, two or three lines is all one would give!

$$(b) \text{ Ans. } \int_0^1 (x - x_c)^2 dx = \frac{1}{3} - x_c + x_c^2$$

$$\begin{aligned}
\int_0^1 (x - x_c)^2 dx &= \int_{-x_c}^{1-x_c} x^2 dx \\
&= \frac{x^3}{3} \Big|_{-x_c}^{1-x_c} \\
&= \frac{(1-x_c)^3}{3} - \frac{(-x_c)^3}{3} \\
&= \frac{1}{3} - x_c + x_c^2.
\end{aligned}$$



(c) **Ans.** $\int_0^1 \text{rect}(x-3) dx = 0$

A plot of the shifted rectangle function is shown above. It is zero over the range of integration, $[0, 1]$, and therefore, its integral is zero.

Applying Prop. 3.10,

$$\begin{aligned}
\int_0^1 \text{rect}(x-3) dx &= \int_{-3}^{-2} \text{rect}(x) dx \\
&= 0
\end{aligned}$$

because the function $\text{rect}(x)$ is zero over the range of integration, $[-3, -2]$

(d) **Ans.** $\int_3^5 \text{tri}(x-3) dx = \frac{1}{2}$

A plot of the shifted (unit) triangle function is shown above. The range of integration, $[3, 5]$, captures the right half of the triangle, and hence its integral is one half.

Applying Prop. 3.10,

$$\begin{aligned}
\int_3^5 \text{tri}(x-3) dx &= \int_0^2 \text{tri}(x) dx \\
&= 0.5.
\end{aligned}$$

Because the function $\text{tri}(x)$ is centered at the origin, the (shifted) range of integration, $[0, 2]$, excludes the left half of the triangle and includes its right half. ■

3.2.7 Scaling and Integration

Scaling, introduced in Chapter 1.16, causes most learners and (some instructors) more “issues” than shifting. This is because area scales inversely proportional to the scale factor, as long as the limits of integration are scaled as well. Figures 3.6 and 3.7 illustrate how $f(2x)$ encloses half the area of $f(x)$, and similarly, $f(\frac{x}{2})$ encloses double the area. The easiest way to keep this straight is to sketch the first half-periods of $\sin(x)$ and $\sin(2x)$. You are unlikely to forget that $\sin(2x)$ has twice the frequency and hence, half the period, of $\sin(x)$, and hence encloses half the area.

Proposition 3.12: Scaling Property

Suppose that $a < b$, $f : [a, b] \rightarrow \mathbb{R}$ is continuous, and $\omega_0 > 0$ is a real number. Then $g : [\frac{a}{\omega_0}, \frac{b}{\omega_0}] \rightarrow \mathbb{R}$ by $g(x) := f(\omega_0 \cdot x)$ is also continuous and

$$\int_{\frac{a}{\omega_0}}^{\frac{b}{\omega_0}} g(x) dx = \frac{1}{\omega_0} \cdot \int_a^b f(x) dx. \quad (3.23)$$

Similarly, if $\omega_0 < 0$, Then $g : [\frac{b}{\omega_0}, \frac{a}{\omega_0}] \rightarrow \mathbb{R}$ by $g(x) := f(\omega_0 \cdot x)$ is also continuous and

$$\int_{\frac{b}{\omega_0}}^{\frac{a}{\omega_0}} g(x) dx = \frac{-1}{\omega_0} \cdot \int_a^b f(x) dx = \frac{1}{|\omega_0|} \cdot \int_a^b f(x) dx. \quad (3.24)$$

Note: Assume $\omega_0 \neq 0$. By defining $\alpha := \frac{a}{\omega_0}$ and $\beta := \frac{b}{\omega_0}$, the above results can be rewritten as

$$\int_{\alpha}^{\beta} f(\omega_0 \cdot x) dx = \frac{1}{\omega_0} \cdot \int_{\alpha \cdot \omega_0}^{\beta \cdot \omega_0} f(x) dx. \quad (3.25)$$

This latter form is often easier to apply and is valid for $\omega_0 < 0$ and $\omega_0 > 0$.

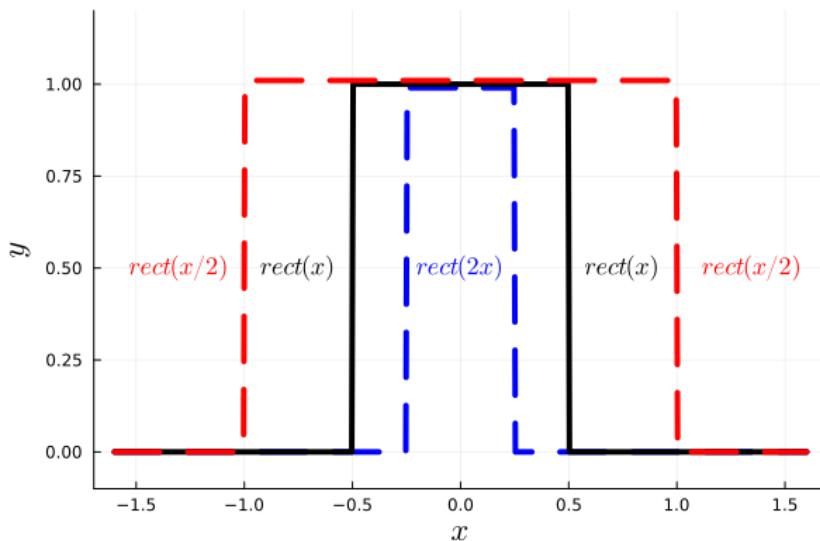


Figure 3.6: (**Scaling does Change Area:**) Plots for three rectangles are shown each with a different scaling. The area scales inversely with the scale factor used in the function when the integration limits are appropriately adjusted. A unit rectangle has area 1.0. A unit rectangle scaled by $2x$ has area 0.5, while a unit rectangle scaled by $x/2$ has area 2.0.

Example 3.13. Evaluate the following definite integrals.

$$(a) \int_0^1 (2x)^2 dx$$

$$(b) \int_0^2 \text{rect}(\frac{x}{2}) dx$$

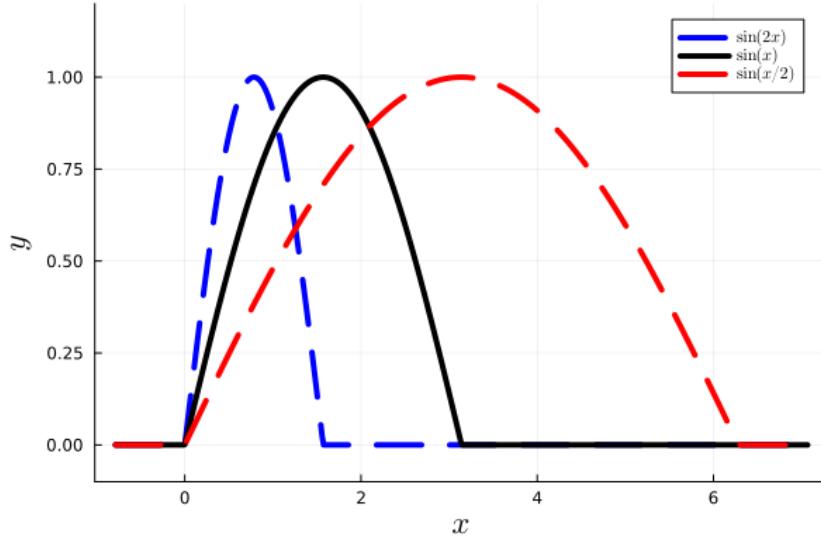


Figure 3.7: (**Scaling does Change Area, Take 2:**) Similarly, a half-period of $\sin(x)$ (meaning integrated over $[0, \pi]$) has area 2.0, whereas a half-period of $\sin(2x)$ (meaning integrated over $[0, \frac{\pi}{2}]$) has area 1.0 and a half-period of $\sin(\frac{x}{2})$ (meaning integrated over $[0, 2\pi]$) has area 4.0. **If we had integrated the sinusoids over a full period, all of the answers would have been zero. Doubling and halving zero is not so informative.**

$$(c) \int_0^1 \text{rect}(\omega_0 x) dx, \text{ for } \omega_0 > 0$$

$$(d) \int_3^5 \text{tri}\left(\frac{x-2}{2}\right) dx \text{ (has both a shift and a scale).}$$

Solution:

$$(a) \text{ Ans. } \int_0^1 (2x)^2 dx = \frac{4}{3}$$

The most obvious way to solve this problem is to note $(2x)^2 = 4x^2$ and hence

$$\int_0^1 (2x)^2 dx = 4 \int_0^1 x^2 dx = 4 \left. \frac{x^3}{3} \right|_0^1 = 4 \left(\frac{1}{3} - 0 \right) = \frac{4}{3}.$$

Applying Prop. 3.12 gives

$$\begin{aligned} \int_0^1 (2x)^2 dx &= \frac{1}{2} \int_0^2 x^2 dx \\ &= \frac{1}{2} \left. \frac{x^3}{3} \right|_0^2 \\ &= \frac{1}{2} \left(\frac{8}{3} - 0 \right) \\ &= \frac{4}{3}. \end{aligned}$$

$$(b) \text{ Ans. } \int_0^2 \text{rect}\left(\frac{x}{2}\right) dx = 1.0$$

Recall that $\text{rect}(x) = \begin{cases} 1 & |x| \leq 0.5 \\ 0 & \text{otherwise.} \end{cases}$. Applying Prop. 3.12 gives

$$\begin{aligned} \int_0^2 \text{rect}\left(\frac{x}{2}\right) dx &= 2 \int_0^1 \text{rect}(x) dx \\ &= 2 \int_0^{0.5} 1 dx + 2 \int_{0.5}^1 0 dx \\ &= 2 \cdot \frac{1}{2} = 1. \end{aligned}$$

(c) **Ans.** For $\omega_0 > 0$, $\int_0^1 \text{rect}(\omega_0 x) dx = \frac{\min\{0.5, \omega_0\}}{\omega_0} = \begin{cases} 1 & 0 < \omega_0 \leq 0.5 \\ \frac{1}{2\omega_0} & \omega_0 > 0.5. \end{cases}$

Applying Prop. 3.12 gives

$$\begin{aligned} \int_0^1 \text{rect}(\omega_0 x) dx &= \frac{1}{\omega_0} \int_0^{\omega_0} \text{rect}(x) dx \\ &= \frac{1}{\omega_0} \int_0^{\min\{0.5, \omega_0\}} 1 dx + \frac{1}{\omega_0} \int_{\min\{0.5, \omega_0\}}^{\omega_0} 0 dx \\ &= \frac{1}{\omega_0} \cdot \min\{0.5, \omega_0\}. \end{aligned}$$

(d) **Ans.** $\int_3^5 \text{tri}\left(\frac{x-2}{2}\right) dx = \frac{3}{4}$

Recall that $\text{tri}(x) = \begin{cases} |x| & |x| \leq 1 \\ 0 & \text{otherwise.} \end{cases}$. Applying Propositions 3.10 and 3.12 gives

$$\begin{aligned} \int_3^5 \text{tri}\left(\frac{x-2}{2}\right) dx &= \int_1^3 \text{tri}\left(\frac{x}{2}\right) dx \quad (\text{remove the shift first}) \\ &= 2 \int_{1/2}^{3/2} \text{tri}(x) dx \quad (\text{then take care of the scale}) \\ &= 2 \int_{1/2}^1 x dx + 2 \int_1^{3/2} 0 dx \quad (\text{apply def. of triangle function}) \\ &= 2 \left. \frac{x^2}{2} \right|_{1/2}^1 + 0 \\ &= 2 \left(\frac{1}{2} - \frac{1}{8} \right) \\ &= 2 \cdot \frac{3}{8} = \frac{3}{4}. \end{aligned}$$

■

3.3 Numerical Methods for Approximating Riemann Integrals

You may have been thinking that rectangles are not very efficient for estimating area and that other geometric shapes might better “fill in” the area. One popular shape is a trapezoid, as shown in Fig. 3.8, which caps a rectangle with a right triangle; this gives rise to the **Trapezoidal Rule** for approximating the area under a curve. A trapezoid results from approximating a function over a “small” interval $[x_i, x_{i+1}]$ by a (straight) line segment connecting $f(x_i)$ and $f(x_{i+1})$. Another popular idea is to approximate a function over $[x_i, x_{i+1}]$ by a quadratic; this gives rise to **Simpson’s Rule**.

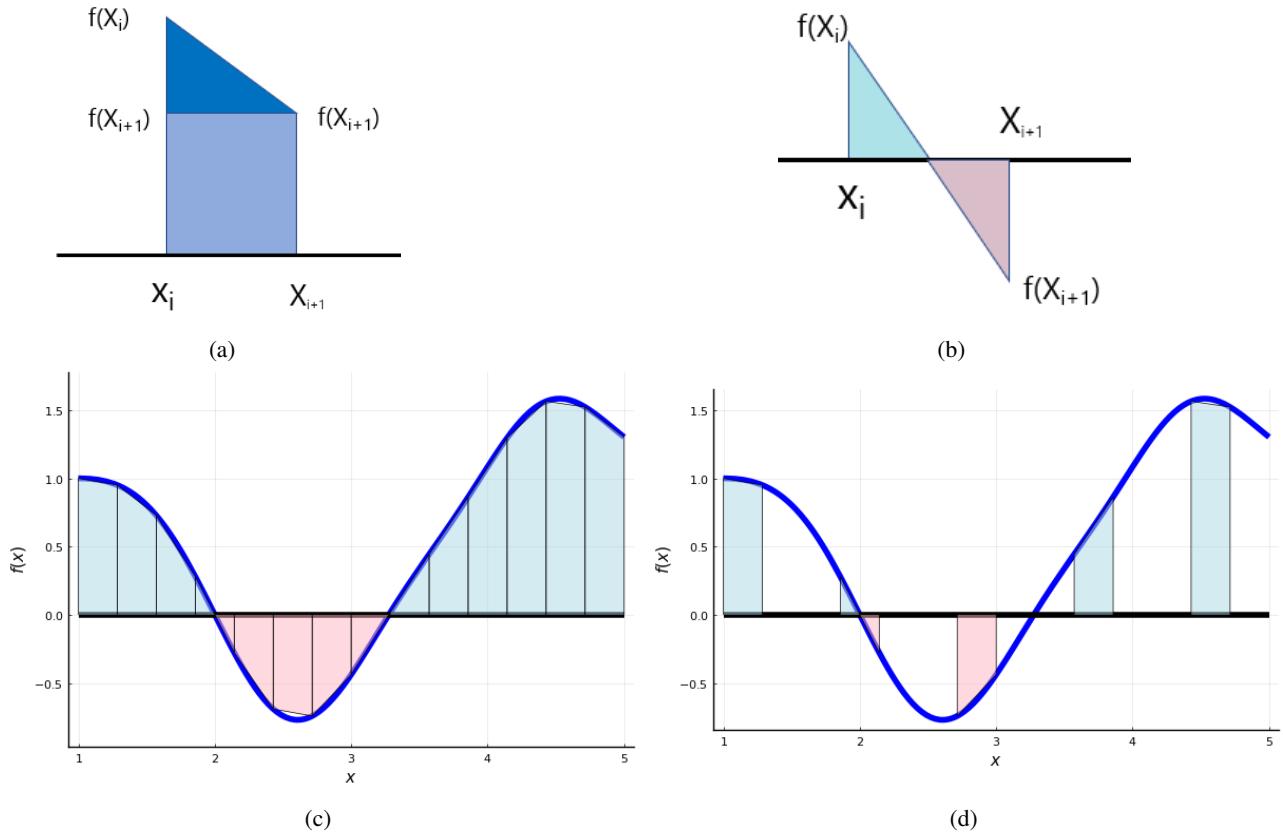


Figure 3.8: Using trapezoids to approximate the area of a function $f : [1, 5] \rightarrow \mathbb{R}$. Trapezoids are polygons with four corners, parallel sides, and a flat base, like a rectangle, but where the two sides can have different heights. (a) Illustration that a trapezoid can also be viewed as a rectangle plus a right triangle, where the rectangle is determined by the smallest height, the base of the right triangle is the top of the rectangle, and the side of the right triangle is determined by the difference in the two heights. (b) When $f(x_i)$ and $f(x_{i+1})$ have opposite signs, the trapezoid becomes two right triangles, one “pointing up” and one “pointing down”. (c) Shows 14 trapezoids approximating the area under the curve defined by f . (d) Shows every third trapezoid so that the individual shapes can be more clearly seen.

3.3.1 Trapezoidal Rule

The area under the trapezoids shown in Fig. 3.8-(a) and -(b) is given by

$$\text{trapArea}_i = \frac{f(x_i) + f(x_{i+1})}{2} \cdot (x_{i+1} - x_i), \quad (3.26)$$

which is the average (or, mean) of the two rectangles used in the Riemann lower and upper sums, when the function is increasing. For a partition, $a = x_1 < x_2 < \dots < x_n < x_{n+1} = b$, the area under the curve over the interval $[a, b]$ is approximated by

$$\text{Area}_n^{\text{Trap}} := \sum_{i=1}^n \text{trapArea}_i = \sum_{i=1}^n \frac{f(x_i) + f(x_{i+1})}{2} \cdot (x_{i+1} - x_i). \quad (3.27)$$

Proposition 3.14: Trapezoidal Rule for Computing a Riemann Integral

Suppose that $a < b$ and $f : [a, b] \rightarrow \mathbb{R}$ is continuous. Then the Riemann Integral $\int_a^b f(x) dx$ exists and it can be evaluated by the **Trapezoidal Rule**. That is,

$$\int_a^b f(x) dx = \lim_{n \rightarrow \infty} \text{Area}_n^{\text{Trap}} = \lim_{n \rightarrow \infty} \sum_{i=1}^n \frac{f(x_i) + f(x_{i+1})}{2} \cdot (x_{i+1} - x_i) = \lim_{n \rightarrow \infty} \sum_{i=1}^n \frac{f(x_i) + f(x_{i+1})}{2} \cdot \Delta x, \quad (3.28)$$

where, $\Delta x := \frac{b-a}{n}$ and $x_i := a + (i-1)\Delta x$. In the above, the limits mean that for all $\epsilon > 0$, there exists $N < \infty$ such that $n \geq N \implies \left| \int_a^b f(x) dx - \sum_{i=1}^n \frac{f(x_i) + f(x_{i+1})}{2} \cdot (x_{i+1} - x_i) \right| < \epsilon$.

Remark 3.15. There is no obvious “Lower Trapezoidal Sum” or “Upper Trapezoidal Sum” that can be used to bound the error in the approximation for a given, finite value of n , as we had with the method of rectangles. A common “hack” is to form the sum for an initial value of n , and then recompute it for $2n$, $3n$, etc., until either the absolute change or the relative change in the calculation is less than a given aTol , that is, for $k = 1, 2, \dots$, stop when

$$\left| \text{Area}_{(k+1)\cdot n}^{\text{Trap}} - \text{Area}_{k\cdot n}^{\text{Trap}} \right| \leq \text{aTol} \quad (\text{absolute change, or})$$

$$\left| \frac{\text{Area}_{(k+1)\cdot n}^{\text{Trap}} - \text{Area}_{k\cdot n}^{\text{Trap}}}{\text{Area}_{k\cdot n}^{\text{Trap}}} \right| \leq \text{aTol} \quad (\text{relative change}).$$

Neither of these hacks provides the firm bound we had previously, $|\text{Area}^{\text{True}} - \text{Area}_n^{\text{Approx}}| \leq \text{aTol}$.

Is there something more that can be done? Yes.

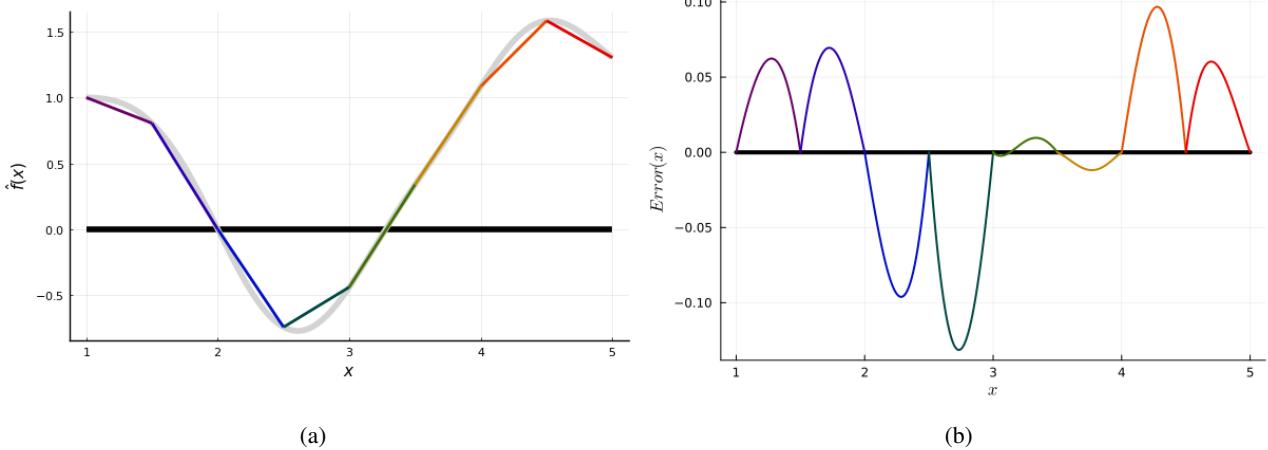


Figure 3.9: The Trapezoidal Rule can be viewed as a piecewise linear approximation of the function $f : [1, 5] \rightarrow \mathbb{R}$ over intervals $[x_i, x_{i+1}]$. (a) Shows an approximation of f using eight linear segments, one for each of the intervals $[1, 1.5]$, $[1.5, 2]$, etc., to $[4.5, 5]$. (b) Shows the error in the fit. The peak values in the error over each interval $[x_i, x_{i+1}]$ can be used to bound the total error in the approximated integral.

We formally define a **piecewise linear approximation** to the function $f : [1, 5] \rightarrow \mathbb{R}$ by

$$\widehat{f}(x) := m_i(x - x_i) + f(x_i) \text{ for } x \in [x_i, x_{i+1}), \quad (3.29)$$

where

$$m_i := \frac{f(x_{i+1}) - f(x_i)}{x_{i+1} - x_i} \quad (\text{rise over run}). \quad (3.30)$$

Figure 3.9-(a) shows a plot of a typical \hat{f} . The **approximation uses the term “piecewise”** because a different approximation (in this case, linear) is used for each subinterval $[x_i, x_{i+1}]$. The Trapezoidal Rule integrates exactly the linear approximation. The error in the Riemann Integral of the original function f arises from the difference, $\text{Error}(x) := f(x) - \hat{f}(x)$. A typical error plot is shown in Fig. 3.9-(b). We can use this to form an estimate of the error in the Riemann Integral by

$$\text{Error}_i^{\max} := \max_{x_i \leq x \leq x_{i+1}} |f(x) - \hat{f}(x)|. \quad (3.31)$$

The error in the computation of the Riemann Integral can then be upper bounded by

$$\text{Error}^{\text{RI}} \leq \sum_{i=1}^n \text{Error}_i^{\max} \cdot \Delta x = \Delta x \cdot \sum_{i=1}^n \max_{x_i \leq x \leq x_{i+1}} |f(x) - \hat{f}(x)|. \quad (3.32)$$

For any interval $[x_i, x_{i+1}]$ where Error_i^{\max} in (3.31) is larger than some preset tolerance, you are free to subdivide the interval and attempt to reduce the error.

```

1 function trapezoidalRule (f, a, b, N, aTol=1e-6)
2   if abs (a-b) < aTol
3     println ("abs (a-b) < $aTol. Estimated integral is zero .")
4     return 0.0
5   end
6   # make sure N is an integer greater or equal to 2
7   N=max (ceil (Int64, N), 2)
8   # Define \Delta x
9   Deltax = (b-a) / N
10  # Initialize Sum
11  estIntegralTrapezoidal = 0.0
12  # Initialize x_i
13  xi=a
14  for i = 1 : N
15    # define x_{i+1}
16    xip1 = xi + Deltax
17    estIntegralTrapezoidal = estIntegralTrapezoidal + Deltax* (f (xi) + f (xip1)) / 2.0
18    # update xi
19    xi = xip1
20  end
21  return estIntegralTrapezoidal
22 end
```

Output

```
trapezoidalRule (generic function with 2 methods)
```

```

1 f (x) = sin (x)
2 a = 0.0
3 b = 2*pi
4 N = 25
5 trapezoidalRule (f, a, b, N)
```

Output

```
2.373101715136272e-15
```

Recall that we know that the area under the function $\sin : [0, 2\pi] \rightarrow \mathbb{R}$ is zero. The Trapezoidal Rule has returned a value very close to zero, 2.37×10^{-15} , but how accurate is it? To find out, we implement the error estimate in (3.32).

```

1 function trapezoidalRuleWithErrorBounds (f, a, b, N=100, aTol=1e-6)
2     if abs(a-b) < aTol
3         println ("abs (a-b) < $aTol. Estimated integral is zero .")
4         return 0.0
5     end
6     # make sure N is an integer greater or equal to 2
7     N=max (ceil(Int64,N),2)
8     # Define \Delta x
9     Deltax = (b-a)/N
10    # Initialize Sum
11    estIntegralTrapezoidal = 0.0
12    # Initialize Error Estimate
13    pmError = 0.0
14    nRefine = 10
15    # Initialize x_i
16    xi=a
17    for i = 1:N
18        # define x_{i+1}
19        xip1 = xi + Deltax
20        y1 = f(xi)
21        y2 = f(xip1)
22        estIntegralTrapezoidal = estIntegralTrapezoidal + Deltax*(y1 + y2)/2.0
23        xRefined = range(xi, xip1, length=nRefine)
24        fRefined = f.(xRefined)
25        g(x) = y1 + (x-xi)*(y2-y1)/Deltax
26        gRefined = g.(xRefined)
27        errorLocal = fRefined-gRefined
28        errorChoice = "max"
29        if errorChoice == "max"
30            maxError = maximum(abs.(fRefined-gRefined))
31            pmError = pmError + Deltax * maxError
32        elseif errorChoice == "avg"
33            avgError = sum(abs.(fRefined-gRefined))/length(errorLocal)
34            pmError = pmError + Deltax * avgError
35        else # mean error
36            meanError = sum(errorLocal)/length(errorLocal)
37            pmError = pmError + Deltax * meanError
38        end
39        # update xi
40        xi = xip1
41    end
42    return (estIntegralTrapezoidal=estIntegralTrapezoidal, pmError=pmError)
43 end

```

Output

trapezoidalRuleWithErrorBounds (generic function with 3 methods)

```

1 f(x) = sin(x)
2 a = 0.0
3 b = 2*pi
4 N = 25
5 G=trapezoidalRuleWithErrorBounds(f, a, b, N)
6 display(G)

```

Output

```
(estIntegralTrapezoidal = 2.373101715136272e-15, pmError = 0.03127822283578252)
```

If we increase the number of line segments, the result improves, as we would expect.

```
1 f(x) = sin(x)
2 a = 0.0
3 b = 2*pi
4 N = 250
5 G=trapezoidalRuleWithErrorBounds(f, a, b, N)
6 display(G)
```

Output

```
(estIntegralTrapezoidal = 1.2034589904480475e-14, pmError = 0.0003120774612516434)
```

We can also test it on the complicated function in (3.7).

```
1 f(x) = cos(exp(-x)*sin(x)/(1+x+5x^2)) + log2((1+x)/(cos(x)+5x^2))
2 a = 0.0
3 b = 10
4 N = 5000
5 G=trapezoidalRuleWithErrorBounds(f, a, b, N)
6 display(G)
```

Output

```
(estIntegralTrapezoidal = -28.832292268297902, pmError = 2.701204661411364e-6)
```

The error is now more than a thousand times smaller than with the Riemann Upper and Lower Sums, using the same number of subintervals,

$$\int_0^{10} \cos\left(\frac{\exp(-x)\sin(x)}{1+x+5x^2}\right) + \log_2\left(\frac{1+x}{\cos(x)+5x^2}\right) dx \approx -28.832292 \pm 2.701205e-6, \quad (3.33)$$

versus the previous error bound of $\approx \pm 5e-3$.

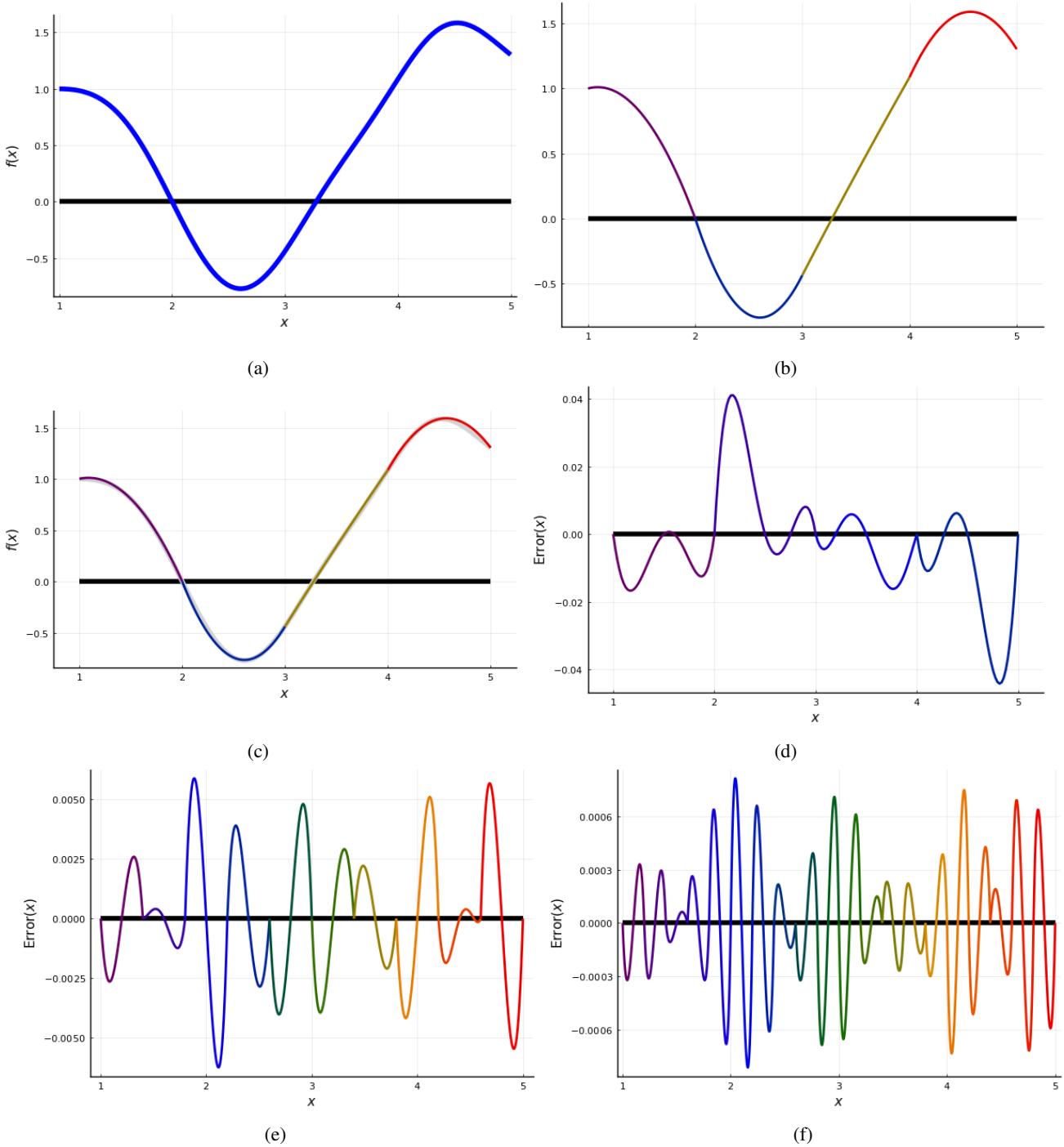


Figure 3.10: Simpson's Rule uses piecewise quadratic functions to approximate the function $f : [1, 5] \rightarrow \mathbb{R}$ over intervals $[x_i, x_{i+1}]$. (a) The function f that is to be integrated. (b) Shows an approximation of f using only four quadratics, one for each of the intervals $[1, 2]$, $[2, 3]$, $[3, 4]$, and $[4, 5]$. See Prop. 3.17 for how to design the quadratic functions. (c) Shows an overlay of the approximation via four quadratics onto the original function; it's hard to tell the difference. (d) Shows the error in the fit, once again, using only four quadratics. (e) Shows the error in the fit using ten quadratics. (f) Shows the error in the fit using twenty quadratics. You need to look at the scales to appreciate that the magnitude of the error is being reduced by a factor of ten each time.

3.3.2 Simpson's Rule

Simpson's Rule is inspired by the thought, "if one can make headway by piecewise approximating a function by lines, why not make even more headway by piecewise approximating a function by quadratics?" From Prop. 3.3 and Prop. 3.4, we obtain that

$$\begin{aligned}\int_a^b 1 dx &= b - a \\ \int_a^b x dx &= \frac{b^2 - a^2}{2} \\ \int_a^b x^2 dx &= \frac{b^3 - a^3}{3}.\end{aligned}\tag{3.34}$$

When combined with Prop. 3.7, we know how to integrate exactly a general quadratic function as the sum of the integrals of the functions 1, x , and x^2 .

Interpolating a Quadratic Function through a Set of Three Points

Definition 3.16. Consider an interval $[x_i, x_{i+1}] \subset \mathbb{R}$ and a function $f : [x_i, x_{i+1}] \rightarrow \mathbb{R}$. Let $x_c := \frac{x_i + x_{i+1}}{2}$ be the midpoint of the interval and $\Delta x := x_{i+1} - x_i > 0$ be its length. The quadratic function

$$q(x) := \alpha(x - x_c)^2 + \beta(x - x_c) + \gamma\tag{3.35}$$

is said to be **centered** about x_c . The function q is said to **interpolate the function** f if the two functions agree on the three points $\{x_i, x_c, x_{i+1}\}$, that is

$$\begin{aligned}q(x_i) &= f(x_i) \\ q(x_c) &= f(x_c) \\ q(x_{i+1}) &= f(x_{i+1}).\end{aligned}$$

Intuition for Simpson's Rule: The quadratic function $q(x) := \alpha(x - x_c)^2 + \beta(x - x_c) + \gamma$ has three unknowns, α , β , and γ . One needs, therefore, three equations to determine the unknowns. The equations come from matching the quadratic to $f : [x_i, x_{i+1}] \rightarrow \mathbb{R}$ at the two endpoints x_i and x_{i+1} , and the center point $x_c = \frac{x_{i+1} + x_i}{2}$. These are the interpolation conditions given in the above definition.

Matching $q(x)$ and $f(x)$ at the midpoint gives $\gamma = f(x_c)$. The algebra to do the match at the endpoints is given in the proof of Prop. 3.17. Integrating the quadratic, linear, and constant terms in $q(x)$ over the interval $[x_i, x_{i+1}]$ then gives the exact area under the quadratic function, $q : [x_i, x_{i+1}] \rightarrow \mathbb{R}$, and a "good approximation" of the area under the given function, $f : [x_i, x_{i+1}] \rightarrow \mathbb{R}$. There is nothing magical about using a quadratic function. You can use any function where solving for its unknown coefficients is straightforward, and you know how to compute its integral in closed form.

Proposition 3.17: Basis of Simpson's Rule: Integral of an Interpolated Quadratic Function

The quadratic function (3.35) centered about the midpoint $x_c = x_i + \frac{\Delta x}{2}$ interpolates $f : [x_i, x_{i+1}] \rightarrow \mathbb{R}$ at the points $x_i < x_c < x_{i+1}$ if, and only if,

$$\alpha = \frac{2}{(\Delta x)^2} (f(x_i) + f(x_{i+1}) - 2f(x_c))$$

$$\beta = \frac{1}{\Delta x} (f(x_{i+1}) - f(x_i))\tag{3.36}$$

$$\gamma = f(x_c).$$

Moreover, the Riemann integral of the interpolating quadratic function is equal to

$$\begin{aligned}\int_{x_i}^{x_{i+1}} q(x) dx &= \alpha \cdot \frac{(\Delta x)^3}{12} + \gamma \cdot \Delta x \\ &= \frac{\Delta x}{6} \cdot (f(x_i) + 4 f(x_c) + f(x_{i+1})),\end{aligned}\tag{3.37}$$

where $x_c := \frac{x_i+x_{i+1}}{2}$.

Proof: Because $q(x_c) = \gamma$, we immediately obtain $\gamma = f(x_c)$. The other two coefficients are obtained from the following 2×2 system of linear equations

$$\begin{bmatrix} \left(\frac{\Delta x}{2}\right)^2 & -\frac{\Delta x}{2} \\ \left(\frac{\Delta x}{2}\right)^2 & \frac{\Delta x}{2} \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} = \begin{bmatrix} f(x_i) - f(x_c) \\ f(x_{i+1}) - f(x_c) \end{bmatrix},$$

where we have used $\Delta x = x_{i+1} - x_i$.

To obtain Simpson's formula, we apply Prop. 3.10 to $q(x)$ to obtain

$$\int_{x_i}^{x_{i+1}} q(x) dx = \int_{x_i-x_c}^{x_{i+1}-x_c} q(x+x_c) dx = \int_{-\frac{\Delta x}{2}}^{\frac{\Delta x}{2}} (\alpha x^2 + \beta x + \gamma) dx = \alpha \cdot \frac{(\Delta x)^3}{12} + \gamma \cdot \Delta x.$$

■

Remark 3.18. The “integration” formula is amazingly simple: two-thirds of the function’s value at the midpoint of the interval $[x_i, x_{i+1}]$, plus one-sixth for each of the two endpoints; we note that $\frac{2}{3} + 2 \cdot \frac{1}{6} = 1.0$. Let’s keep in mind the road we have traveled: Riemann Sums use one value of the function for each interval to approximate area via a rectangle, the Trapezoidal Rule uses two values, $f(x_i)$ and $f(x_{i+1})$ with weights of one half, and Simpson’s Rule uses three values, $f(x_i)$, $f(x_{i+1})$, with weight one-sixth and $f(x_c)$ with weight two thirds. If you approximate with a cubic, you’ll use four function values and different weights. In general, you can fit with an n -th degree polynomial, but eventually, you are faced with **Diminishing Returns**.

```

1 function simpsonRule(f, a, b, N=100, aTol=1e-6)
2     # Uses a centered quadratic
3     if abs(a-b) < aTol
4         println("abs(a-b) < $aTol. Estimated integral is zero.")
5         return 0.0
6     end
7     # make sure N is an integer greater or equal to 2
8     N=max(ceil(Int64,N),2)
9     # Define \Delta x
10    Deltax = (b-a)/N
11    # Initialize Sum
12    estIntegralSimpson = 0.0
13    # Initialize x_i
14    xi=a
15    for i = 1:N
16        xip1 = xi + Deltax
17        a1 = f(xi)
18        a2 = f(xip1)
19        ac = f((xi + xip1)/2)
20        estIntegralSimpson = estIntegralSimpson + Deltax*(a1 + a2)/6.0 + 2.0*Deltax*ac/3.0
21        # update xi
22        xi = xip1
23    end
24    return estIntegralSimpson
25 end

```

Output

```
simpsonRule (generic function with 3 methods)
```

Here is a version with error estimation.

```
1 function simpsonRuleWithErrorBounds(f, a, b, N=100, aTol=1e-6)
2     # Uses a centered quadratic
3     if abs(a-b) < aTol
4         println("abs(a-b) < $aTol. Estimated integral is zero.")
5         return 0.0
6     end
7     # make sure N is an integer greater or equal to 2
8     N=max(ceil(Int64,N),2)
9     # Define \Delta x
10    Deltax = (b-a)/N
11    # Initialize Sum
12    estIntegralSimpson = 0.0
13    # Initialize total error
14    pmError = 0.0
15    # Initialize x_i
16    xi=a
17    for i = 1:N
18        xip1 = xi + Deltax
19        xc = (xi + xip1)/2
20        a1 = f(xi)
21        a2 = f(xip1)
22        ac = f(xc)
23        estIntegralSimpson = estIntegralSimpson + Deltax*(a1 + a2)/6.0 + 2.0*Deltax*ac/3.0
24        nRefine = 50
25        xRefined = range(xi, xip1, length=nRefine)
26        fRefined = f.(xRefined)
27        # Coefficients of the centered quadratic
28        alpha = (2/Deltax^2)*(a1 + a2 - 2*ac)
29        beta = (1/Deltax)*(a2-a1)
30        gamma = ac
31        # centered quadratic
32        g(x) = alpha*(x-xc)^2 + beta*(x-xc) + gamma
33        gRefined = g.(xRefined)
34        errorLocal = fRefined-gRefined
35        errorChoice = "max"
36        if errorChoice == "max"
37            maxError = maximum(abs.(fRefined-gRefined))
38            pmError = pmError + Deltax * maxError
39        elseif errorChoice == "avg"
40            avgError = sum(abs.(fRefined-gRefined))/length(errorLocal)
41            pmError = pmError + Deltax * avgError
42        else # mean error
43            meanError = sum(errorLocal)/length(errorLocal)
44            pmError = pmError + Deltax * meanError
45        end
46        # update xi
47        xi = xip1
48    end
49    return (estIntegralSimpson=estIntegralSimpson, pmError=pmError)
50 end
```

Output

```
simpsonRuleWithErrorBounds (generic function with 3 methods)
```

```
1 f(x) = sin(x)
2 a = 0.0
3 b = 2*pi
4 N = 25
5 G = simpsonRuleWithErrorBounds(f, a, b, N)
6 display(G)
```

Output

```
(estIntegralSimpson = 3.021888295151598e-15, pmError = 0.0005179561770644536)
```

```
1 f(x) = sin(x)
2 a = 0.0
3 b = 2*pi
4 N = 250
5 G = simpsonRuleWithErrorBounds(f, a, b, N)
6 display(G)
```

Output

```
(estIntegralSimpson = 1.3069054302303207e-14, pmError = 5.096276651101236e-7)
```

```
1 f(x) = cos(exp(-x)*sin(x)/(1+x+5x^2)) + log2((1+x)/(cos(x)+5x^2))
2 a = 0.0
3 b = 10
4 N = 5000
5 G = simpsonRuleWithErrorBounds(f, a, b, N)
6 display(G)
```

Output

```
(estIntegralSimpson = -28.832291734252568, pmError = 1.1358825205876213e-9)
```

The error is now more than a thousand times smaller than with the Trapezoidal Rule, using the same number of subintervals,

$$\int_0^{10} \cos\left(\frac{\exp(-x)\sin(x)}{1+x+5x^2}\right) + \log_2\left(\frac{1+x}{\cos(x)+5x^2}\right) dx \approx -28.832292 \pm 1.135883e-9, \quad (3.38)$$

versus the previous error bound of $\approx \pm 2.701205e-6$. And the error is two and a half million times (2.5×10^6) smaller than when we estimated the area with rectangles.

Even though, for many real problems, such accuracy is overkill, it's still fun to see how it all works out. In real life, you'd be happy with an error of 10^{-3} and would use a far smaller n to achieve it with Simpson's Rule than with Riemann Lower and Upper Sums.

3.3.3 Julia Packages

QuadGK is a numerical integration function found in languages like MATLAB and Julia, which uses the **Gauss-Kronrod** quadrature rule to calculate definite integrals. Recall that Simpson's Rule,

$$\int_{x_i}^{x_i+\Delta x} f(x) dx \approx \frac{\Delta x}{6} \cdot (f(x_i) + f(x_{i+1})) + \frac{2 \cdot \Delta x}{3} \cdot f(x_c),$$

works by forming a weighted sum of $f : [x_i, x_{i+1}] \rightarrow \mathbb{R}$ evaluated at x_i , $x_c := \frac{x_i+x_{i+1}}{2}$, and x_{i+1} , and then adding up all of these terms to form an approximation of area. QuadGK is similar, except it uses a more sophisticated choice of points and weights discovered by Gauss and improved on by [Kronrod](#). The estimates over each interval $[x_i, x_{i+1}]$ are then summed for the final result. Moreover, the package QuadGK refines estimates iteratively until they're within a threshold, making it adaptive and able to handle various functions and difficult areas like discontinuities. See the [documentation](#) for more information.

Example 3.19. Use the function `quadgk` to evaluate the following definite integrals taken from YouTube. Each integral is supposed to be uncommonly difficult for students to integrate.

(a) $\int_1^e \frac{1+x^2 \cdot \ln(x)}{x+x^2 \cdot \ln(x)} dx$ from [Cipher](#).

(b) $\int_0^\pi \frac{2 \cdot x \cdot \sin(x)}{3+\cos(2 \cdot x)} dx$ from [Lets Solve Math Problems](#), 2018 Stanford Math Tournament.

(c) $\int_{0.1}^{0.5} \left(\sqrt{1-x^2} + \sqrt{x^2+1} + x^2 \cdot \ln(x) + \frac{\arctan(x)}{x^2} + x^{\frac{\sqrt{x}}{\ln(x)}} + \frac{x^2+1}{x^3-2x} + \sin(x)^4 \cdot \cos(x)^2 + \frac{1}{2+\cos(x)} + \tan(x)^{10} + x^x \right)^2 dx$ from “when the [Calculus teacher promises there will be only one integral on the test](#)” by [BlackPenRedPen](#). Yikes!

(d) $\int_{-\frac{1}{2}}^{\frac{1}{2}} \sqrt{x^2 + 1 + \sqrt{x^4 + x^2 + 1}} dx$ from [MathTuts](#), solving integrals from the 2023 MIT Integration Bee Finals: Q3.

Solutions:

(a)

```

1 using QuadGK
2 f(x) = ( 1 + x^2 * log(x) ) / ( x + x^2 * log(x) )
3 area, error = quadgk(f, 1, exp(1))

```

Output

(1.4050201409408223, 2.707664759071804e-9)

(b)

```

1 f(x) = ( 2.0*x*sin(x) ) / ( 3 + cos(2*x) )
2 area, error = quadgk(f, 0, pi)

```

Output

(2.4674011002723395, 4.898666333685853e-10)

(c)

```

1 f(x) = ( sqrt(1 - x^2) + sqrt(x^2+1) + x^2 * log(x) + atan(x)/x^2 +
2      x^(sqrt(x)/log(x)) + (x^2+1)/(x^3-2x) + sin(x)^4 * cos(x)^2 +
3      1/(2 + cos(x)) + tan(x)^10 + x^x )^2
4 F = quadgk(f, 0.1, 0.5)
5 G = simpsonRuleWithErrorBounds(f, 0.1, 0.5, 3000)
6 display(F)
7 display(G);

```

Output

(16.53669833564563, 6.188844992038867e-10)

(estIntegralSimpson = 16.5366983356456,
pmError = 4.499301236175278e-10)

(d)

```

1 f(x) = sqrt( x^2 + 1 + sqrt( x^4 + x^2 + 1 ) )
2 F = quadgk(f, -0.5, 0.5)
3 G = simpsonRuleWithErrorBounds(f, -0.5, 0.5, 3000)
4 display(F)
5 display(G);

```

Output

(1.4586629142874117, 7.910228028151778e-11)

(estIntegralSimpson = 1.458662914287415, pmError = 3.0790407260876096e-14)

(Optional Read): ApproxFun is a second interesting package for numerical integration. The syntax is hard to grasp at first, but it even allows the solution of “integral equations²”, a generalization of the Ordinary Differential Equations we study in Chapter 9. This package first takes a function $f : [a, b] \rightarrow \mathbb{R}$ defined on a bounded interval $[a, b]$, $-\infty < a < b < \infty$, and transforms it to an equivalent function on the interval $[-1, 1]$ through shifting and scaling. It then approximates the function by a linear combination of the form,

$$\widehat{f}(x) = \sum_{k=0}^N f_k \cdot \cos(k \cdot \arccos(x)),$$

where the functions $\cos(k \cdot \arccos(x))$ are **orthogonal** on $-1 \leq x \leq 1$ with respect to the weight $\frac{1}{\sqrt{1-x^2}}$. The orthogonality condition is

$$\int_{-1}^1 \cos(j \cdot \arccos(x)) \cdot \cos(k \cdot \arccos(x)) \frac{1}{\sqrt{1-x^2}} dx = 0,$$

for all $j \neq k$. The number of terms N is adjusted so that the approximation error $f(x) - \widehat{f}(x)$ is close to “machine epsilon”. Hence, the results are pretty accurate.

The downside of the ApproxFun package is that you need more than basic Calculus to really understand what it is doing, whereas the concepts behind QuadGK are not dissimilar to Simpson’s Rule.

We’ll repeat two of the previous examples so that you can see how the syntax goes when using ApproxFun.

```

1 using SpecialFunctions, ApproxFun
2 f(x) = ( 1 + x^2 * log(x) ) / ( x + x^2 * log(x) )
3
4 # Define the approximating function
5 g = Fun(f, 1..exp(1))
6 # Define the operation of computing a Definite Integral
7 Sigma = DefiniteIntegral()
8 # Apply the Definite Integral operator to g
9 Sigma*g

```

Output

1.4050201409408223 anywhere

```

1 f(x) = ( sqrt(1 - x^2) + sqrt(x^2+1) + x^2 * log(x) + atan(x)/x^2 +
2      x^( sqrt(x)/log(x) ) + (x^2+1)/(x^3-2x) + sin(x)^4 * cos(x)^2 +
3      1/( 2 + cos(x) ) + tan(x)^10 + x^x )^2
4 F = quadgk(f, 0.1, 0.5)
5 G = simpsonRuleWithErrorBounds(f, 0.1, 0.5, 3000)

```

²Think differential equations, but replace derivatives with integrals.

```

6 display (F)
7 display (G)
8
9 # Define the approximating function
10 g = Fun(f, 0.1..0.5)
11 # Define the operation of computing a Definite Integral
12 Sigma = DefiniteIntegral()
13 # Apply the Definite Integral operator to g
14 Sigma*g

```

Output

```
(16.53669833564563, 6.188844992038867e-10)

(estIntegralSimpson = 16.5366983356456, pmError = 4.499301236175278e-10)

16.536698335645635 anywhere
```

3.3.4 (Optional Read:) Even and Odd Functions as Great Test Cases

When using a numerical integration package on a challenging function, it is always good to be able to generate “test cases”, whereby we mean functions that are related to an important problem you are trying to solve, and yet, you know ahead of time what the integral of the function is, exactly. If this seems useful, then the following may be of interest to you.

Even and Odd Functions, and Parts of Functions

Definition 3.20.

- (a) $f : \mathbb{R} \rightarrow \mathbb{R}$ is **even** if $f(-x) = f(x)$.
- (b) $f : \mathbb{R} \rightarrow \mathbb{R}$ is **odd** if $f(-x) = -f(x)$.

Notes:

- Even functions are modeled on $f(x) = x^2$, while odd functions are modeled on $f(x) = x^3$.
- Common even functions are x^{2k} for $k \geq 0$, $\cos(x)$, and $\cosh(x)$.
- Common odd functions are x^{2k+1} for $k \geq 0$, $\sin(x)$, and $\sinh(x)$.
- More examples can be found at [Math is Fun](#) and [Khan Academy](#).

Most functions are neither even nor odd. That is what makes this class of functions very useful to us.

Definition 3.21. Consider $f : \mathbb{R} \rightarrow \mathbb{R}$.

- (a) $f_e : \mathbb{R} \rightarrow \mathbb{R}$ by $f_e(x) := \frac{f(x)+f(-x)}{2}$ is called the **even part of f** .
- (b) $f_o : \mathbb{R} \rightarrow \mathbb{R}$ by $f_o(x) := \frac{f(x)-f(-x)}{2}$ is called the **odd part of f** .

Notes:

- $f_e(-x) := \frac{f(-x)+f(-(-x))}{2} = \frac{f(-x)+f(x)}{2} = \frac{f(x)+f(-x)}{2} = f_e(x)$, and hence f_e is an even function.
- $f_o(-x) := \frac{f(-x)-f(-(-x))}{2} = \frac{f(-x)-f(x)}{2} = -\frac{f(x)-f(-x)}{2} = f_o(x)$, and hence f_o is an odd function.
- $f_e(x) + f_o(x) = \frac{f(x)+f(-x)}{2} + \frac{f(x)-f(-x)}{2} = \frac{f(x)+f(x)}{2} + \frac{f(-x)-f(-x)}{2} = f(x)$
- **Hence, every function is the sum of its even and odd parts.**

Example 3.22. Compute the even and odd parts of e^x .

Solutions:

We define $f(x) := e^x$. Then,

$$f_e(x) = \frac{f(x) + f(-x)}{2}$$

$$= \frac{e^x + e^{-x}}{2}$$

$$= \cosh(x).$$

Similarly,

$$f_o(x) = \frac{f(x) - f(-x)}{2}$$

$$= \frac{e^x - e^{-x}}{2}$$

$$= \sinh(x).$$

■

Proposition 3.23: Integrals of Even and Odd Functions over Symmetric Intervals

Suppose that $f : \mathbb{R} \rightarrow \mathbb{R}$ is a function and that for $a > 0$, the definite integrals $\int_{-a}^0 f(x) dx$ and $\int_0^a f(x) dx$ both exist and are finite. Then the following hold:

(a) If f is an even function,

- $\int_{-a}^0 f(x) dx = \int_0^a f(x) dx$, and
- $\int_{-a}^a f(x) dx = 2 \int_0^a f(x) dx$.

(b) If f is an odd function,

- $\int_{-a}^0 f(x) dx = -\int_0^a f(x) dx$, and
- $\int_{-a}^a f(x) dx = 0$.

(c) For all $[b, c] \subset [-a, a]$,

- $\int_b^c f(x) dx = \int_b^c f_e(x) dx + \int_b^c f_o(x) dx$.

Example 3.24. Figure 3.11 shows a complicated function to integrate. We are concerned about the ability of QuadGK to integrate the function. One means of checking is

- Integrate $f(x)$.
- Integrate the odd part of $f(x)$, $f_o(x)$.
- Integrate the even part of $f(x)$, $f_e(x)$.
- Check that $\int_a^b f(x) dx = \int_a^b f_0(x) dx + \int_a^b f_e(x) dx$.

They match to 4.385×10^{-15} , which is either outstanding or QuadGK actually integrated the even and odd parts separately and added them together!

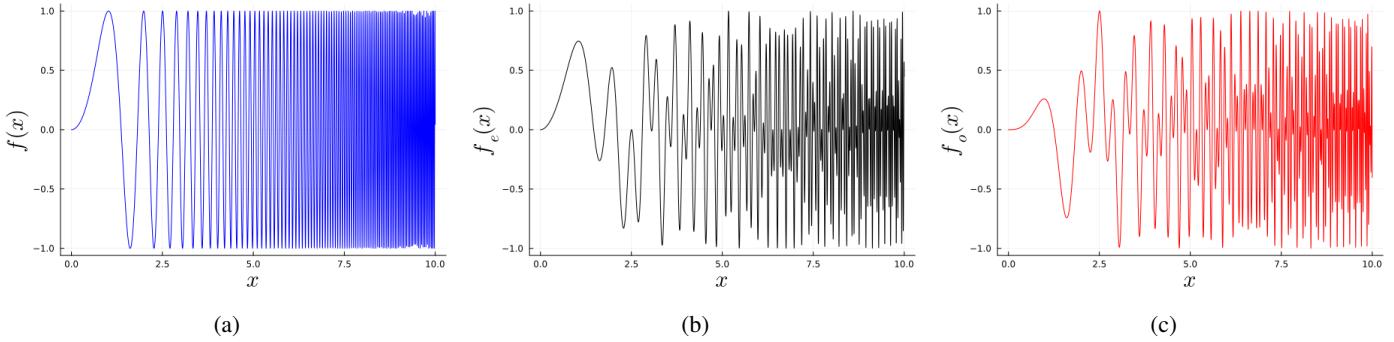


Figure 3.11: (**Even and Odd Parts to the Rescue:**) (a) shows a complicated function that may be hard to integrate. (b) and (c) show the even and odd parts, respectively. While they still look hard to integrate, Prop. 3.23 shows how they can be used to check the accuracy of an integration routine.

Solution:

```

1 using QuadGK
2
3 f(x) = sin(x^2 + .5*x^3)
4 fo(x) = (f(x) - f(-x))/2.0
5 fe(x) = (f(x) + f(-x))/2.0
6
7 xmax = 10.0
8 xmin = 0
9
10 area, error = quadgk(f, xmin, xmax);           @show area, error
11 area_o, error_o = quadgk(fo, xmin, xmax);     @show area_o, error_o
12 area_e, error_e = quadgk(fe, xmin, xmax);     @show area_e, error_e
13
14 area = (area_o + area_e)

```

Output

```

(area, error) = (0.4817407024672528, 6.724935180040881e-9)
(area_o, error_o) = (0.08683542238431138, 1.235326109499035e-9)
(area_e, error_e) = (0.39490528008293707, 5.824010571660494e-9)

4.385380947269368e-15

```

3.4 Applications of the Definite Integral

Now that we've delved into the theoretical underpinnings of the definite integral, it's time to witness its practical prowess. This cornerstone of calculus isn't just an abstract mathematical concept; it's a versatile tool that brings clarity to complex problems across diverse fields. In robotics and mechanical design, its value becomes immediately apparent. Take, for instance, a mobile robot traversing a landscape: by integrating its velocity, we can pinpoint its exact change in position, offering a detailed account of its journey. When we're faced with determining the area sandwiched between two functions, the definite integral steps in, quantifying the shared space with precision. This is especially vital in robot link design, where discerning properties like the center of mass is paramount for ensuring both stability and accurate movement. Venturing further, the definite integral plays a starring role in the realm of solids of revolution, allowing us to both visualize and calculate the volume of shapes birthed from rotating curves around an axis. As we work through these various applications, the definite integral's depth and adaptability will shine through, underscoring why it was invented at the precise time in human history when we were trying to decode the intricacies of the world and the solar system.



Figure 3.12: A small mobile robot that is used in various Michigan Robotics courses. The user commands wheel speed. If both wheels move at the same speed, the robot moves in a straight line. If one wheel turns faster than the other, the robot turns toward the side of the slower wheel.

3.4.1 From Speed to Change in Position

When we drive a car, fly an airplane, or pilot a boat, with the goal of going from point A to point B, we typically don't set new position destinations (i.e., waypoints) every second along a planned path. Instead, we command speed and direction. Here, we'll explore the relationship between speed and change in position in the context of a small mobile ground robot, such as those used in ROB 330 Localization, Mapping, and Navigation; see Fig. 3.12.

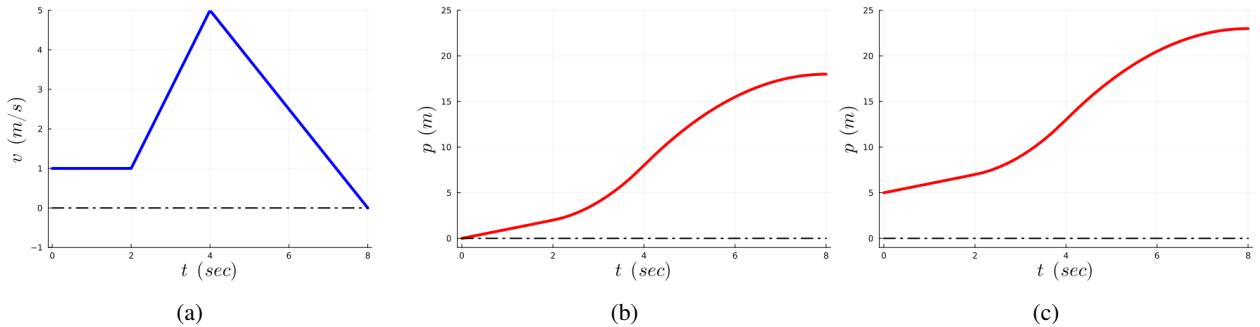


Figure 3.13: How is position related to speed of a moving body? (a) Shows a simple speed profile given by (3.39). (b) Shows the resulting position profile when the body starts at the point $p_0 = 0.0$, while (c) shows the position profile when the body starts at a point, $p_0 = 5.0$ meters, instead of the origin. Varying p_0 simply shifts the curve up or down.

Consider the speed profile,

$$v(t) = \begin{cases} 1 & 0 \leq t < 2 \\ 1 + 2 \cdot (t - 2) & 2 \leq t < 4 \\ 5 - 1.25 \cdot (t - 4) & 4 \leq t \leq 8, \end{cases} \quad (3.39)$$

as plotted in Fig. 3.13-(a). For the first two seconds, the robot is moving at 1 m/s (i.e., meters per second). Hence, after 1 s, it has traveled 1 m, and after 2 s, it has traveled 2 m. That's pretty simple to comprehend because the speed is constant over the first two seconds. What about the next two seconds, where the speed is ramping up with a slope of 2 m/s^2 ? Well, we could compute the change in position by approximating the speed over a small interval of time as being constant, then multiply the constant by the time interval, and finally, add the displacement to the previous speed position. But wait a minute, that's the same procedure we followed

when defining the Riemann Integral! Yes, indeed! It follows that **position is related to speed** by³

$$p(t) = p_0 + \int_{t_0}^t v(\tau) d\tau, \quad (3.40)$$

where t_0 is the initial time, $p_0 := p(t_0)$ is the initial position of the moving body at time t_0 , and τ is a dummy variable of integration. Fig. 3.13-(b) and -(c) show the position profile for two different initial conditions, namely, $p(0) = 0.0$ m and $p(0) = 5.0$ m. The integral itself computes the **change in position** or **relative displacement** of the robot. We have to add to the relative displacement the initial position of the robot to obtain the true position of the robot.

In the following, we show an inefficient, but effective means of computing $p_0 + \int_{t_0}^t v(t) dt$, as a function of time, for $t \geq t_0$.

```

1 # Speed profile
2 function v(t)
3     if t<2
4         v = 1
5     elseif t < 4
6         v = 1 + 2*(t-2)
7     else
8         v = 5 - 1.25*(t-4)
9     end
10    return v
11 end
12
13 # How to integrate a generic function f from t0 to t
14 function integral_t(f, t0, t)
15     result, err = quadgk(f, t0, t)
16     return result
17 end
18
19 # Set an initial position
20 p0=5.0
21
22 # Define the integral as a function so that we can plot it
23 int(t) = p0 + integral_t(v, 0, t)
24
25 p1 = plot(int, 0, 8, linewidth=4, color=:red, label=false,
26           xlabel=L"\$t \sim (\{\rm sec\}) \$", ylabel=L"\$p \sim (\rm m) \$",
27           yguidefontsize=18, xguidefontsize=18)
28 plot!(g, 0, 8, linewidth=4, color=:black, label=false)
29 xlims!(p1, (-0.1, 8.1))
30 ylims!(p1, (-1, 25))
31 display(p1)
32 png(p1, "SimpleDisplacementProfilOffset")

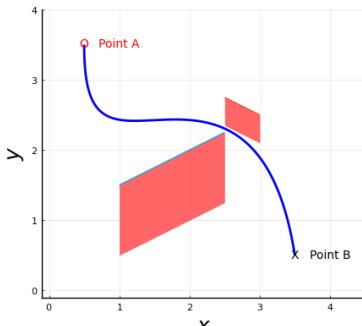
```

Output The plots in Fig. 3.13-(b) and -(c).

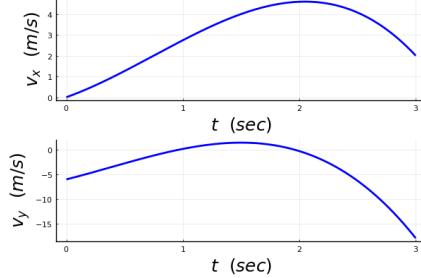
Let's now do something a bit wilder and consider a robot moving in the plane \mathbb{R}^2 , instead of along a line. We'll model a velocity profile as a vector, $[v_x(t), v_y(t)]$. Because each component of the velocity depends on a single variable, t , we can perform the Riemann Integral separately for each component, as in

$$\begin{aligned} p(t) = \begin{bmatrix} p_x(t) \\ p_y(t) \end{bmatrix} &= \begin{bmatrix} p_x(t_0) + \int_{t_0}^t v_x(\tau) d\tau \\ p_y(t_0) + \int_{t_0}^t v_y(\tau) d\tau \end{bmatrix}. \end{aligned} \quad (3.41)$$

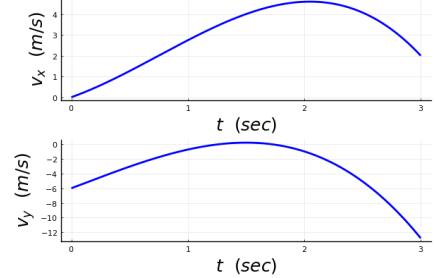
³So as not to confuse the variable being integrated with the time variable t , we used τ as the “dummy variable”.



(a)



(b)



(c)

Figure 3.14: We want to steer from Point A to Point B, without hitting the obstacles in (a). Among the speed profiles shown in (b) and (c), which one will get the job done? They look pretty similar.

Example 3.25. Figure 3.14 shows a patch of \mathbb{R}^2 with some obstacles in it. We want to steer our robot so that it starts at Point A, goes to Point B, and avoids the obstacles shown by the colored sections. Which of the following velocity profiles will work?

$$(i) \begin{bmatrix} v_y(t) \\ v_x(t) \end{bmatrix} = \begin{bmatrix} 0.00807179 \\ -1.97588 \end{bmatrix} + \begin{bmatrix} 0.677259 \\ 2.01576 \end{bmatrix} \cdot t + \begin{bmatrix} 0.437326 \\ 0.645312 \end{bmatrix} \cdot t^2 + \begin{bmatrix} -0.197531 \\ -0.592593 \end{bmatrix} \cdot t^3; \text{ see Fig. 3.14-(b).}$$

$$(ii) \begin{bmatrix} v_y(t) \\ v_x(t) \end{bmatrix} = \begin{bmatrix} 0.00807179 \\ -1.97596 \end{bmatrix} + \begin{bmatrix} 0.677259 \\ 2.00199 \end{bmatrix} \cdot t + \begin{bmatrix} 0.437326 \\ 0.0767879 \end{bmatrix} \cdot t^2 + \begin{bmatrix} -0.197531 \\ -0.335802 \end{bmatrix} \cdot t^3; \text{ see Fig. 3.14-(c)}$$

Solution

We use the velocity profile from (i).

```

1 # Show the parameters being used
2 @show P0
3 @show P4
4 @show beta
5
6 # Define the proposed velocity profiles and so that we can integrate them
7 Velx(t) = beta[1,1] + t*beta[2,1] + t^2*beta[3,1] + t^3*beta[4,1]
8 Vely(t) = beta[1,2] + t*beta[2,2] + t^2*beta[3,2] + t^3*beta[4,2]
9
10 # How to integrate a generic function f from t0 to t
11 function integral_t(f, t0, t)
12     # Could just as easily use our version of Simpson's Rule
13     result, err = quadgk(f, t0, t)
14     return result
15 end
16
17 # Set an initial position
18 px0 = P0[1]
19 py0 = P0[2]
20
21 # Define the integral as a function so that we can plot it
22 intx(t) = px0 + integral_t(Velx, 0, t)
23 inty(t) = py0 + integral_t(Vely, 0, t)
24
25 # Use integration to compute the Robot's Path in the Plane
26 N=1000
27 t = range(0, 3, length=N)
28 posRobotx = intx.(t)
29 posRoboty = inty.(t)

```

```

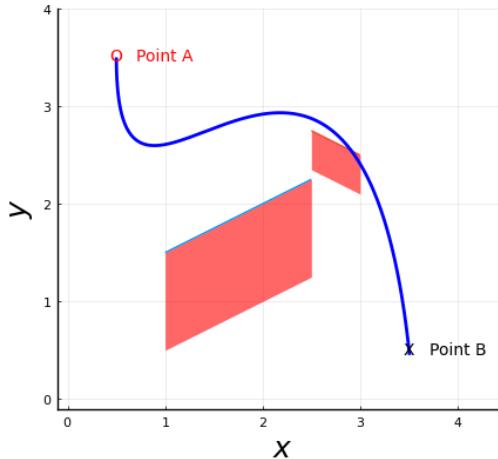
30
31 # Plot the Obstacles
32 ob1(x)=1 + x/2
33 ob2(x)=0.0 + x/2
34 x=range(1, 2.5, length=N)
35 p1=plot(x, ob1.(x), fillrange=ob2.(x), fillalpha=0.6, fillcolor=:red, label=false)
36 ob1(x)=4 - x/2
37 ob2(x)=3.6 - x/2
38 x=range(2.5, 3.0, length=N)
39 p1=plot!(x, ob1.(x), fillrange=ob2.(x), fillalpha=0.6, fillcolor=:red, label=false)
40
41 # Plot the Robot's path in the Plane with Obstacles
42 plot!(posRobotx, posRoboty, legend=false, linewidth=2, color=:blue, label=false,
43 aspect_ratio=1, xlabel=L"$x$ ", ylabel=L"$y$ ",
44 yguidefontsize=18, xguidefontsize=18)
45 annotate!(P0[1], P0[2], text("O", :red, :center, 10))
46 annotate!(P0[1] + 0.5, P0[2], text("Point A", :red, :center, 10))
47 annotate!(P4[1], P4[2], text("X", :black, :center, 10))
48 annotate!(P4[1] + 0.5, P4[2], text("Point B", :black, :center, 10))
49 xlims!(p1, (-0.1, 4.5))
50 ylims!(p1, (-0.1, 4))

```

Output

```
P0 = [0.5, 3.5]
P4 = [3.5, 0.5]
```

```
beta = [0.008071793840493424 -1.9758808107720176; 0.6772585064197526 2.015759503265411;
0.4373262151028052 0.6453119786307883; -0.19753086419795485 -0.592592592589968]
```



A plot of the path corresponding to the velocity profile in (i) is shown above. While it ALMOST misses the upper obstacle, it does, in fact, hit it, as shown in the plot. The correct answer is (ii), which produces the path shown in Fig. 3.14-(a).

3.4.2 Ballistic Motion

Ballistic motion, also known as projectile motion, refers to the motion of objects that are launched, thrown, or projected into the air and subject to only the force of gravity. This type of motion combines both vertical and horizontal movements. When an object is



Figure 3.15: (**Ballistic Motion:**) All Hail San Antonio's Coyote, Inventor Of The T-shirt Cannon? A normally welcome example of ballistic motion in everyday life.

launched, it follows a curved path known as a parabola, assuming there is no air resistance. The analysis of ballistic motion allows us to predict the trajectory, range, maximum height, and time of flight of the projectile, making it an important concept in various applications, from sports to engineering.

Similar to how (3.40) allows us to go from velocity, $v(t)$, to position, $p(t)$, we can also go from acceleration, $a(t)$, to velocity, $v(t)$ via

$$v(t) = v_0 + \int_{t_0}^t a(\tau) d\tau, \quad (3.42)$$

and then apply (3.40) to obtain position. For **ballistic motion** in the absence of air resistance, the acceleration is constant,

$$a(t) = \begin{bmatrix} 0.0 \\ -g \end{bmatrix}, \quad (3.43)$$

making it quite easy to obtain velocity,

$$v(t) = \begin{bmatrix} v_x(t) \\ v_y(t) \end{bmatrix} = \begin{bmatrix} v_x(t_0) + \int_{t_0}^t a_x(\tau) d\tau \\ v_y(t_0) + \int_{t_0}^t a_y(\tau) d\tau \end{bmatrix} = \begin{bmatrix} v_x(t_0) + \int_{t_0}^t 0.0 d\tau \\ v_y(t_0) - \int_{t_0}^t g d\tau \end{bmatrix} = \begin{bmatrix} v_x(t_0) \\ v_y(t_0) - g \cdot (t - t_0) \end{bmatrix}. \quad (3.44)$$

Integrating once more to obtain position, yields, the well-known quadratic equations

$$p(t) = \begin{bmatrix} p_x(t_0) + \int_{t_0}^t v_x(\tau) d\tau \\ p_y(t_0) + \int_{t_0}^t v_y(\tau) d\tau \end{bmatrix} = \begin{bmatrix} p_x(t_0) + \int_{t_0}^t v_x(t_0) d\tau \\ p_y(t_0) + \int_{t_0}^t (v_y(t_0) - g \cdot (\tau - t_0)) d\tau \end{bmatrix} = \begin{bmatrix} p_x(t_0) + v_x(t_0) \cdot (t - t_0) \\ p_y(t_0) + v_y(t_0) \cdot (t - t_0) - \frac{1}{2} g \cdot (t - t_0)^2 \end{bmatrix}. \quad (3.45)$$

In the above equation, t_0 , $p(t_0)$, and $v(t_0)$ are initial conditions for the motion, whereas g is a physical parameter.

Example 3.26. Jalen Rose of *Fab Five* fame wants to pop a t-shirt to Jimmy King, whom he has spotted hiding up in the stands. If Jimmy is at a position 20 meters longitudinally and 10 meters vertically from the t-shirt cannon, at what angle should Jalen tilt the cannon, if the exit velocity of the t-shirt is 35 m/s (meters/second)?

Additional Data: The mass of a t-shirt is 140 g or 0.14 Kg; zero degrees corresponds to the cannon's barrel being parallel to the ground; and $g \approx 9.81 \text{ m/s}^2$.



(a)

(b)

Solution: **Ans.** $\approx 0.6395 \text{ rad} \approx 37^\circ$.

We can assume $t_0 = 0$. And because we are given the desired end position of the t-shirt relative to its exit point from the cannon, we have $p_x(t_0) = p_y(t_0) = 0.0$. Next, we resolve the t-shirt's exit velocity from the cannon barrel into its x and y components,

$$\begin{bmatrix} v_x(0) \\ v_y(0) \end{bmatrix} = \begin{bmatrix} 25 \cos(\theta) \\ 25 \sin(\theta) \end{bmatrix}.$$

Let T be the time of flight (the time the t-shirt lands in Jimmy King's lap, a quantity we need to solve for). The ballistic equation (3.45) simplifies to

$$p(T) = \begin{bmatrix} 25 \cos(\theta) \cdot T \\ 25 \sin(\theta) \cdot T - \frac{1}{2}(9.81)T^2 \end{bmatrix} = \begin{bmatrix} 20 \\ 10 \end{bmatrix}. \quad (3.46)$$

Solving the first equation for T yields $T = \frac{20}{25 \cos(\theta)}$, which when substituted into the second equation gives

$$10 = 20 \tan(\theta) - 4.905 \cdot \left(\frac{4}{5 \cos(\theta)} \right)^2 \implies \theta \approx 0.6395 \text{ rad} \approx 36.6^\circ,$$

which we solve with the bisection algorithm. From the Julia code, we also obtain that it takes 0.997 s for the t-shirt to arrive in Jimmy's lap.

```

1
2 # make the algorithm handy
3 function bisection(f, a, b, tol)
4     if f(a) * f(b) > 0
5         error("f(a) and f(b) must have opposite signs")
6         return NaN
7     end
8     kmax = 1e5 # max number of loops so that the algorithm terminates
9     k=1
10    c=NaN #define c out of the loop so we can access it
11    while ( (b - a) / 2 > tol || abs(f(c)) > tol ) && k < kmax
12        c = (a + b) / 2
13        #
14        if f(a) * f(c) < 0
15            b = c
16        else
17            a = c
18        end
19        k = k+1
20    end
21    pmErrorBound = (b - a) / 2

```

```

22   println("Number of iterations = ", k)
23   return (est=(a + b) / 2, low=a, up=b, error=pmErrorBound)
24 end
25
26 # Set up the problem
27 px0 = 0 # meters
28 py0 = 0 # meters
29
30 v0 = 25
31 0 # m/s
32
33 pxT = 20 # meters
34 pyT = 10 # meters
35
36 gG = 9.81
37
38 # pxT = v0*cos(theta)*T
39 T(theta) = pxT/(v0*cos(theta))
40 # py = v0*sin(theta)*T
41 px(theta) = v0*cos(theta)*pxT/(v0*cos(theta))
42
43 #f(theta) = -pyT + v0*sin(theta)*T - 0.5*gG*T^2
44 f(theta) = -pyT + v0*sin(theta)*pxT/(v0*cos(theta)) - 0.5*gG*(pxT/(v0*cos(theta)))^2
45
46 thetaStar = bisection(f, pi/16, pi/4, 1e-6)
47
48 @show theta_rad = thetaStar.est
49
50 @show theta_deg = theta_rad*180/pi
51
52 Tstar = T(thetaStar.est)
53 println("Duration of the tShirts toss = $Tstar seconds")

```

Output

```

Number of iterations = 22
theta_rad = thetaStar.est = 0.6395180791349615
theta_deg = (theta_rad * 180) / pi = 36.64168685674669
Duration of the tShirts toss = 0.9970295179594787 seconds

```

3.4.3 Area between Two Functions

When an area of interest is bounded by two functions instead of one, the calculation for area changes slightly, as illustrated in Fig. 3.17. The area formula becomes

$$\text{Signed Area}_g^f := \int_a^b (f(x) - g(x)) dx. \quad (3.47)$$

As described in the figure, this is the signed area, just as with the normal Riemann Integral. If you want the pink-shaded regions to be counted as positive areas, then you would need to evaluate

$$\text{Unsigned Area}_g^f := \int_a^b |f(x) - g(x)| dx. \quad (3.48)$$

It is a fact that if $f : [a, b] \rightarrow \mathbb{R}$ and $g : [a, b] \rightarrow \mathbb{R}$ are both continuous, then so are $f - g : [a, b] \rightarrow \mathbb{R}$ and $|f - g| : [a, b] \rightarrow \mathbb{R}$, and thus their Riemann Integrals exist.

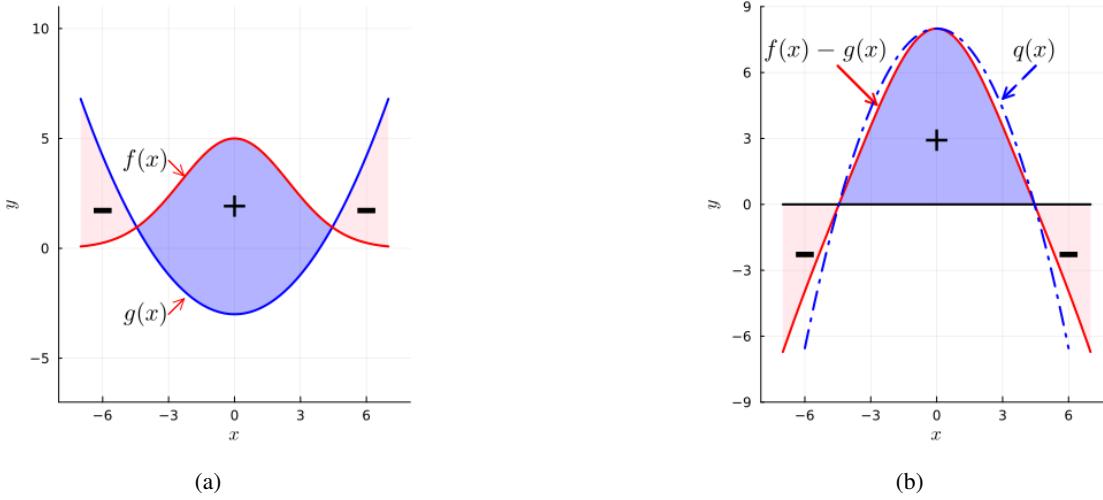


Figure 3.17: The two images have the same scale, the left going from -7.0 to +11.0, and the right from -9.0 to +9.0. We defined the Riemann Integral as the (signed) area under a curve. In reality, we were looking at the area between a function $f(x)$ and the x -axis. When the function was positive, i.e., above the x -axis, the area was positive, and when the function was below the x -axis, the area was negative. More generally, it is common for a shape to be defined by more than one function. In (a), we treat the case of two functions, $f(x)$ and $g(x)$. The (signed) area is the Riemann Integral of $f(x) - g(x)$ over the region of interest. Here, we emphasize that when $f(x) - g(x) < 0$, the area “between the curves” is negative. If you want to count that as a positive area, then you would integrate $|f(x) - g(x)|$, the absolute value of $f(x) - g(x)$. In (b), we plot the function $f(x) - g(x)$. The corresponding shaded areas in each figure have the same area. Moreover, the plot in (b) looks like an inverted parabola, but it is not. The dashed line is the function $q(x) = -0.404022x^2 + 8.0$, which matches $f(x) - g(x)$ at its peak and its two zero crossings. Because a parabola is determined by its values at three distinct points, this shows that $f(x) - g(x)$ is not a parabola.

Example 3.27. Compute the area bounded by the two functions in Fig. 3.18, defined by

$$f(x) = \begin{cases} 2 - |x| & |x| \leq 2 \\ 0 & \text{otherwise} \end{cases} \quad \text{and} \quad g(x) = \begin{cases} |x| - 1 & |x| \leq 1 \\ 0 & \text{otherwise.} \end{cases} \quad (3.49)$$

Solution: The area in Fig. 3.18-(a) is the sum of two triangles,

$$\text{Area}_a = \frac{1}{2} \cdot 4 \cdot 2 + \frac{1}{2} \cdot 2 \cdot 1 = 5.0$$

The area in Fig. 3.18-(b) is the sum of three triangles and a rectangle,

$$\text{Area}_b = \underbrace{\frac{1}{2} \cdot 2 \cdot 2}_{\text{top triangle}} + \underbrace{2 \cdot \frac{1}{2} \cdot 1 \cdot 1}_{\text{two side triangles}} + \underbrace{2 \cdot 1}_{\text{rectangle}} = 5.0$$

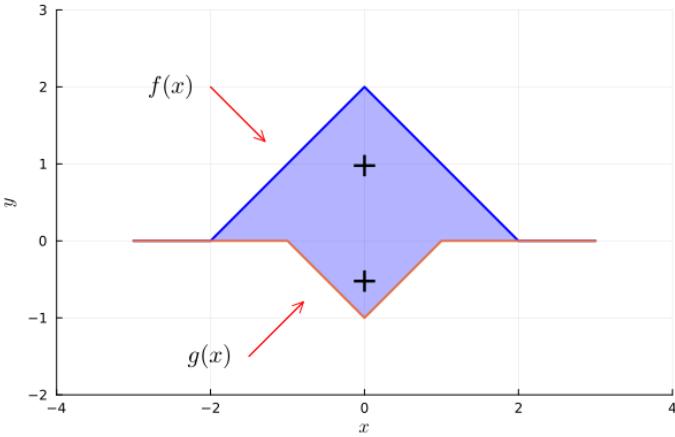
Hence, the two areas are equal. ■

Example 3.28. For the functions $f(x) := 5 \cdot \exp(-x^2/12)$ and $g(x) := -3 + 0.2 \cdot x^2$ plotted in Fig. 3.17, compute the following areas:

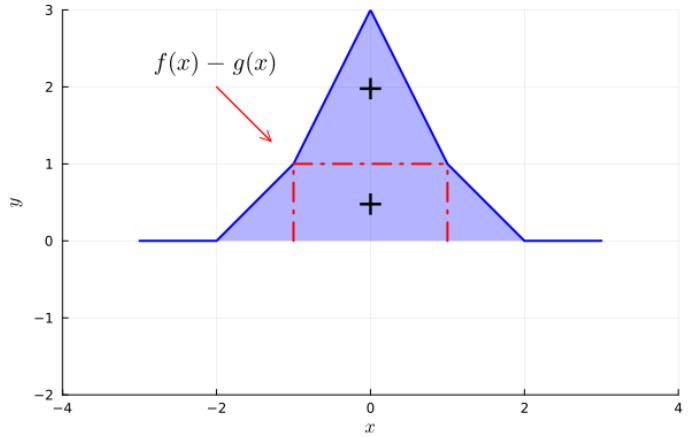
- (a) Signed area between $f(x)$ and $g(x)$ for the interval $[a, b] := [-6, 6]$.
- (b) Unsigned area between $f(x)$ and $g(x)$ for the interval $[a, b] := [-6, 6]$.
- (c) The area indicated by the blue patch in Fig. 3.17.

Solutions:

- (a) **Ans.** 37.461.



(a)



(b)

Figure 3.18: Triangles provide a simpler situation in which to understand that the area bounded by two functions is the same as the area defined by the difference of the two functions. The functions are defined in (3.49). (a) From the image, we compute the area of the top triangle is 4, and the bottom triangle is 1, for a total area of 5 units squared. (b) Decomposing the shape into a top triangle, a rectangle, and two side triangles, we arrive at its area as $1 + 2 + 1 + 1 = 5$. They are the same, whether they look it or not!

(b) **Ans.** 49.587

(c) **Ans.** 43.524

```

1 # Problem Solutions
2
3 f(x) = 5*exp(-x^2/12)
4 g(x) = -3 + 0.2*x^2
5
6 h(x) = f(x) - g(x)
7 absh(x) = abs(h(x))
8 a = -6
9 b = 6
10
11 # For bounding the blue-shaded region, we use the
12 # Bisection Algorithm to find where h(x)=0
13 xRoot1 = bisection(h, -5, -4, 1e-6)
14 xRoot1 = xRoot1.est
15 xRoot2 = bisection(h, 2.5, 5, 1e-6)
16 xRoot2 = xRoot2.est
17
18 # We compute the various areas for parts (a), (b), and (c)
19
20 @show Ans_a, error_a = quadgk(h, a, b)
21 println(" ")
22 @show Ans_b, error_b = quadgk(absh, a, b)
23 println(" ")
24 @show [xRoot1, xRoot2]
25 @show Ans_c, error_c = quadgk(h, xRoot1, xRoot2);

```

Output

(Ans_a, error_a) = quadgk(h, a, b) = (37.46061361388632, 1.2664415294239006e-8)

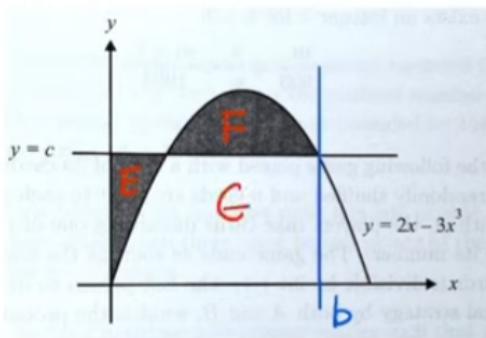
```
(Ans_b, error_b) = quadgk(absh, a, b) = (49.58734405494167, 1.984501950419177e-7)
[xRoot1, xRoot2] = [-4.4498186111450195, 4.449818730354309]
(Ans_c, error_c) = quadgk(h, xRoot1, xRoot2) = (43.52397876865702, 1.0885337076160795e-9)
```

■

Example 3.29. Use numerical methods to solve this problem from a *Putnam Exam*: Find c so that the two areas defined by the curves $y = 2x - 3x^3$ and $y = c$ are equal.



- A1. The horizontal line $y = c$ intersects the curve $y = 2x - 3x^3$ in the first quadrant as in the figure. Find c so that the areas of the two shaded regions are equal.



Solution:

The peak of the curve $y = 2x - 3x^3$ occurs at $x = \frac{\sqrt{2}}{3}$. For $c \in \mathbb{R}$ such that the horizontal line $y = c$ intersects the curve $y = 2x - 3x^3$ in two points, let b be the rightmost intersection. Then, areas E and F are equal if, and only if, $\int_0^b (\underbrace{2x - 3x^3}_f(x) - \underbrace{c}_g(x)) dx = 0$. We use the Bisection Algorithm to solve for b and c .

```

1 f(x) = 2x - 3x^3
2
3 function Int(c)
4     h(x) = f(x)-c
5     # For c given, solve for b such that f(b)-c = 0.
6     # Bisection is great for finding roots
7     Ans = bisection(h, sqrt(2)/3, 2, 1e-6)
8     b = Ans.est
9     # Compute the integral from zero to b. We seek the value of c
10    # such that the integral from 0 to b equals zero
11    # Hence, we can use Bisection for that too.
12    value, error = quadgk(h, 0, b)
13    return Area=value
14 end

```

Output

```
Int (generic function with 1 method)
```

```

1 Ans = bisection(Int, .3, .5, 1e-9)
2 @show c = Ans.est
3
4 # let's find b, even though it is not part of the answer
5 h(x) = f(x)-c
6 Ans = bisection(h, sqrt(2)/3, 2, 1e-6)
7 @show b = Ans.est;

```

Output

```

c = Ans.est = 0.4444444440305233
b = Ans.est = 0.6666664352015328

```

Note: $\frac{4}{9} = 0.444\dots$ and $\frac{2}{3} = 0.666\dots$ Hence, we guess the analytical answers, should we care, are $c = \frac{4}{9}$ and $b = \frac{2}{3}$. ■

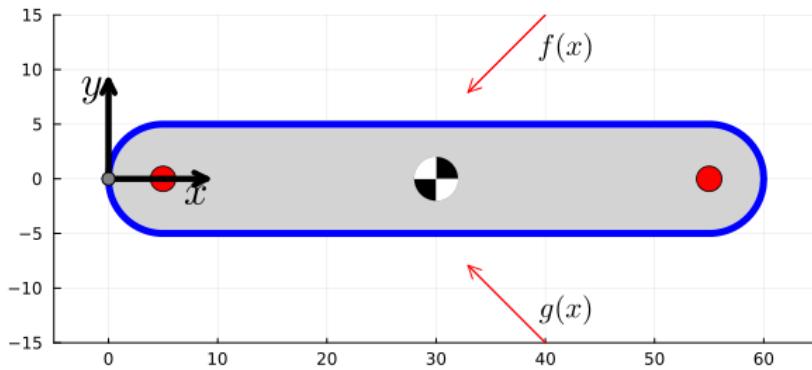


Figure 3.19: A planar version of a very basic robot link in units of centimeters (cm). A Cartesian coordinate system is attached to the link, with its origin at the left end of the link. In these coordinates, the leftmost point on the link is at $(0, 0)$, and the rightmost point is $(60, 0)$. The red dots indicate the location of two revolute joints, located at $(0, 5)$ and $(0, 55)$. The symbol in the middle indicates the location of the link's **center of mass**, a term that will be defined in the text.

Helpful Videos:

- [Area Between Curves by Integration \(Easy Way\)](#) by Quoc Dat Phung has high-quality drawings to support the discussion.
- [Area Between Curves - Integrating with Respect to Y](#) by Quoc Dat Phung has high-quality drawings to support the discussion.

3.4.4 Using the Power of the Rectangle to Find the Center of Mass of a (Planar) Robot Link

In physics and robotics, we are constantly leaning upon foundational principles to explore new frontiers. One of those is Newton's second law, $F = ma$, a seemingly simple equation that elegantly describes the motion of a (single) point mass. Yet, we quickly realize that a real robot link is not just a solitary point but an entire body comprised of infinitely many point masses, all moving in tandem. This realization is not just a curiosity of nature but a profound insight that echoes the essence of Calculus. By viewing a body as an infinite collection of points, we unveil deeper layers of understanding. The **Total Mass** becomes more than just a summation; it's an integral, a continuous accumulation of infinitesimal masses. The **Center of Mass** transforms from a mere geometric point to a gravitational anchor, summarizing the collective pull of gravity on every particle, as seen in this [video](#). And, the **Moment of Inertia** discussed a bit later, it encapsulates a body's resistance to change in its rotational state, just as total mass encapsulates a body's

resistance to change in its horizontal and vertical position. The more mass a body has, the harder it is to push, and similarly, the more inertia a body has, the more effort it takes to change its angle of rotation.

Because this is not a mechanics course, we will make two simplifying assumptions:

- each of our robot's links has a constant thickness in the z -direction, which will not be shown in our link drawings as they only give the outline in the x - y -plane; and
- the density of a link's material is constant throughout the link.

We'll still do awesome things, don't worry. Here is some useful notation to get us going.

Notation:

- A is the (total) area of a link.
- h is the thickness (or height in the z -direction) of a link.
- V is the (total) volume of a link, and hence equals $A \cdot h$.
- ρ is the material density of a link. Because the volume of a typical robot link in m^3 is tiny, we use units of grams per centimeter cubed for the density, that is, g/cm^3 . For reference, the density of aluminum is $\rho_0 := 2.7\ g/cm^3$.
- M_{Total} , or M_T for short, is the total mass of a link. Hence,

$$M_T = \rho \cdot V = \rho \cdot h \cdot A,$$

because of our assumptions of uniform link thickness and density. Hence, for us, computing total mass is tantamount to computing the area, something we know how to do.

Total Mass:

From Chapter 3.4.3, if $f : [x_{\min}, x_{\max}] \rightarrow \mathbb{R}$ and $g : [x_{\min}, x_{\max}] \rightarrow \mathbb{R}$ provide the upper and lower boundaries of an object, such as the robot link shown in Fig. 3.19, then,

$$M_T = \rho \cdot h \cdot \int_{x_{\min}}^{x_{\max}} (f(x) - g(x)) \, dx. \quad (3.50)$$

This is very straightforward due to our assumptions of constant density and thickness. We next move on to an object's center of mass.

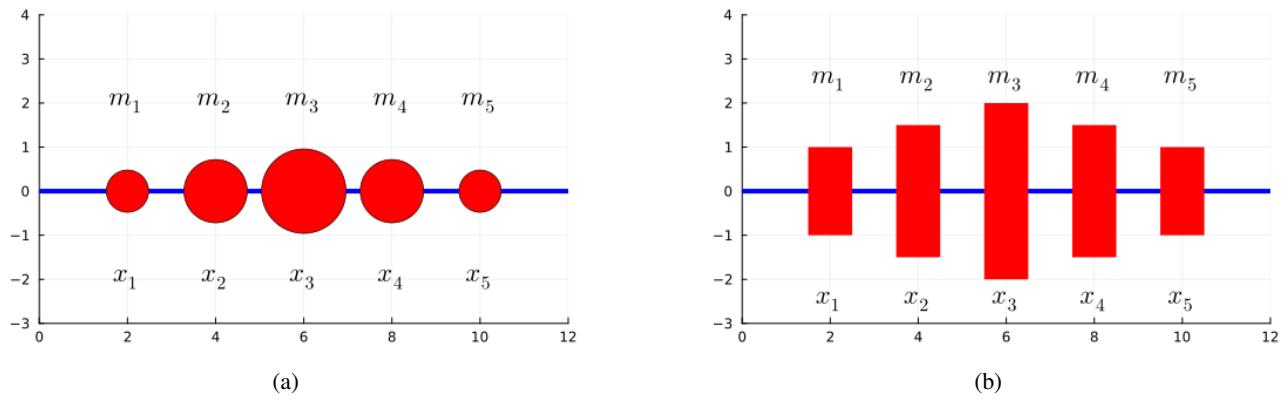


Figure 3.20: The x -center of mass is the point where the collection of masses would balance, taking into account both the position and magnitude of each mass. (a) Representation of a system of discrete point masses m_i located at positions x_i along the x -axis. The size of each red ball is proportional to its corresponding mass. (b) Visualization of the same system of discrete masses as rectangles. Each rectangle's height corresponds to the mass m_i and is centered at position x_i . This representation bridges the concept from discrete point masses to continuous mass distributions, emphasizing the idea of the center of mass as the system's balance point in a continuous setting. The depiction with rectangles should inspire you to recall that the Riemann integral is, after all, a sum of infinitely many, infinitely thin rectangles.

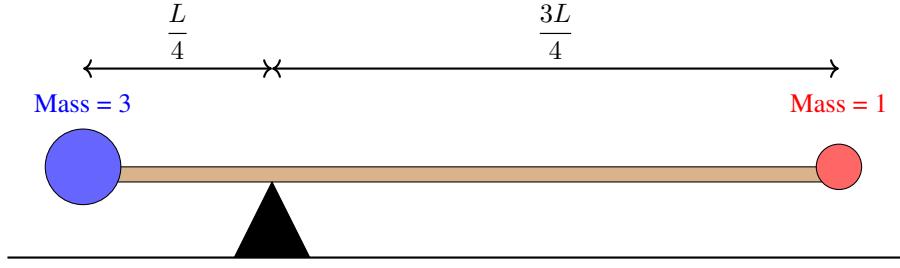


Figure 3.21: To be in balance, the sum of the moments about the fulcrum must be equal. If x_c is the position of the fulcrum and x_1, x_2 denote the positions of the masses m_1 and m_2 , then the condition for balance is $m_1 \cdot (x_1 - x_c) + m_2 \cdot (x_2 - x_c) = 0$.

Center of Mass, First for Point Masses:

You may have encountered this concept in a Physics course for a discrete set of masses as shown in Figs. 3.20. Given the symmetric distribution of the masses in the figure, the center of mass is computed by eye to be the position of the middle mass, x_3 . But how does one determine it in general, such as in Fig. 3.21?

Consider a collection of N point masses m_1, m_2, \dots, m_N of various sizes located at various positions x_1, x_2, \dots, x_N respectively along the x -axis. The **total mass** is

$$M_T := \sum_{i=1}^N m_i.$$

For the case of point masses lying on the x -axis (i.e., their y -coordinate is identically zero) the moment about the origin and the moment about the y -axis are the same. Because we will consider soon masses with non-zero y -coordinates, we use the term **Moment about the y -axis**, denoted M_y .

For the i -th mass, its contribution to M_y is the product of its mass and its distance from the y -axis, namely

$$M_{y,i} := m_i \cdot x_i.$$

The **Total Moment about the y -axis** is the sum of the individual moments, giving

$$M_y := \sum_{i=1}^N M_{y,i} = \sum_{i=1}^N m_i \cdot x_i.$$

The **x -coordinate of the center of mass**, denoted x_c , is the total y -moment divided by the total mass of the system,

$$x_c := \frac{M_y}{M_T} = \frac{\sum_{i=1}^N m_i \cdot x_i}{\sum_{i=1}^N m_i}.$$

Interpretation: The formula for x_c is the weighted average of the x -positions of the point masses, where each position is weighted by its corresponding mass. The center of mass is the “balance point” of the system, and the formula reflects this by taking into account both the position and mass of each point in the system. If you need more information, you may enjoy [Equation for Center of Mass](#) by Khan Academy.

Source of the Equation for Center of Mass. From Physics, for a “plank” (or teeter-totter) to be in equilibrium the sum of the “moment arms” must be zero, that is

$$\sum_{i=1}^N m_i \cdot (x_i - x_c) = 0.$$

This yields,

$$\sum_{i=1}^N m_i \cdot x_i = \sum_{i=1}^N m_i \cdot x_c = x_c \cdot \sum_{i=1}^N m_i.$$

Solving this equation for x_c yields the above formulas.

Example 3.30. For Fig. 3.20, the masses and x -locations are

$$(m_1, m_2, m_3, m_4, m_5) = (2, 3, 4, 3, 2)$$

and

$$(x_1, x_2, x_3, x_4, x_5) = (2, 4, 6, 8, 10).$$

Compute the center of mass, x_c .

Solution: The total mass is $M_T = 14$ and the total y -moment equals 84. Hence,

$$x_c := \frac{84}{14} = 6 = x_3.$$

To be in equilibrium, a fulcrum would be placed exactly in the middle at x_3 . Hopefully, this makes physical sense to you, perhaps from your experience⁴ with a teeter-totter?

```

1 x_positions = [2, 4, 6, 8, 10]
2 masses = [2, 3, 4, 3, 2]
3 @show M_T = sum(masses)
4 @show M_y = sum(masses.*x_positions)
5 @show x_c = M_y / M_T;

```

Output

```

M_T = sum(masses) = 14
M_y = sum(masses .* x_positions) = 84
x_c = M_y / M_T = 6.0

```

Example 3.31. Suppose we make the mass at the far end truly massive,

$$(m_1, m_2, m_3, m_4, m_5) = (2, 3, 4, 3, 200)$$

and leave the x values the same,

$$(x_1, x_2, x_3, x_4, x_5) = (2, 4, 6, 8, 10).$$

Compute the new center of mass, x_c .

Solution: The total mass is $M_T = 212$ and the total y -moment equals 2064. Hence,

$$x_c := \frac{2064}{212} = 9.736 \approx x_5.$$

To be in equilibrium, a fulcrum would have to be placed almost directly under the huge mass at x_5 . Hopefully, this makes physical sense to you.

```

1 x_positions = [2, 4, 6, 8, 10]
2 masses = [2, 3, 4, 3, 200]
3 @show M_T = sum(masses)
4 @show totalMomentx = sum(masses.*x_positions)
5 @show x_c = totalMomentx / M_T;

```

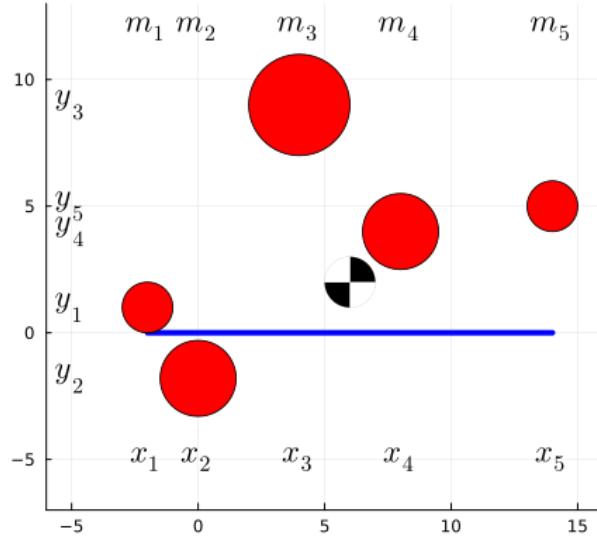
Output

```

M_T = sum(masses) = 212
totalMomentx = sum(masses .* x_positions) = 2064
x_c = totalMomentx / M_T = 9.735849056603774

```

⁴If not you, maybe your goat?



Example 3.32. As in the image above, suppose the masses also have non-zero y -locations in addition to their x -locations. Compute both the x and y components of the center of mass. Note that it is harder to eyeball its location.

Data:

$$(m_1, m_2, m_3, m_4, m_5) = (2, 3, 4, 3, 2)$$

$$(x_1, x_2, x_3, x_4, x_5) = (-2, 0, 4, 8, 14),$$

and

$$(y_1, y_2, y_3, y_4, y_5) = (1, -1.8, 9, 4, 5).$$

Compute the center of mass, (x_c, y_c) .

Solution: The total mass is $M_T = 14$ and the total y -moment equals 64. Hence,

$$x_c := \frac{64}{14} = 4.571.$$

The x -moment equals 54.6. Hence,

$$y_c := \frac{54.6}{14} = 3.9$$

It follows that if the masses were attached to a massless support and we were to drill a finger hole precisely at

$$(x_c, y_c) = (4.571, 3.9),$$

the overall object would balance perfectly on our fingertip without any tendency to rotate. Hopefully, this also makes physical sense to you. If not, this video may help, [Center of Mass along x y plane](#) by Ian Page.

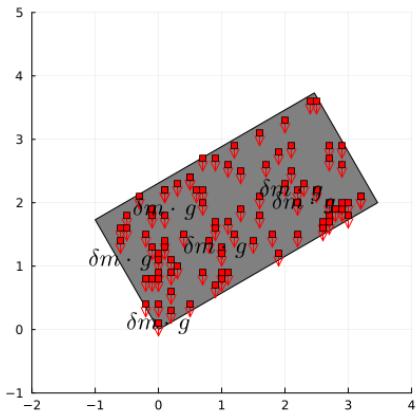
Our closing remark is that the center of mass of a “curved” object does not have to lie within the object itself; consider, for example, an annulus or a doughnut.

```

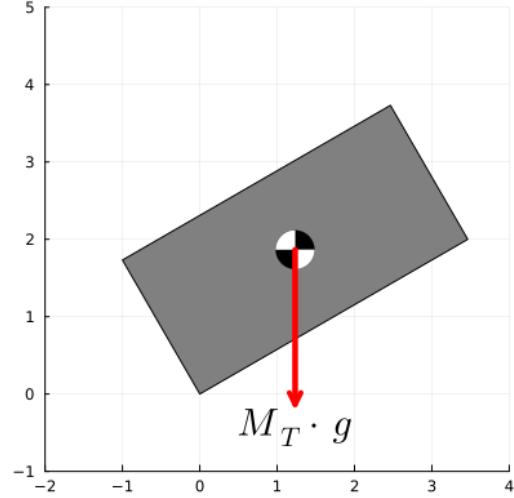
1 # Balls with non-zero y_positions
2 x_positions = [-2, 0, 4, 8, 14]
3 y_positions = [1, -1.8, 9, 4, 5]
4 masses = [2, 3, 4, 3, 2]
5
6 @show M_T = sum(masses)
7 @show M_y = sum(masses.*x_positions)
8 @show M_x = sum(masses.*y_positions)
9 @show x_c = M_y / M_T
10 @show y_c = M_x / M_T;
```

Output

```
M_T = sum(masses) = 14
M_y = sum(masses .* x_positions) = 64
M_x = sum(masses .* y_positions) = 54.6
x_c = M_y / M_T = 4.571428571428571
y_c = M_x / M_T = 3.9
```



(a)



(b)

Figure 3.22: Physical Significance of the Center of Mass: Calculus is what allows us to correctly apply Newton's Laws to macroscopic objects. Image (a) illustrates that a physical object is comprised of a myriad of tiny particles of mass δm , each one undergoing the effects of gravity. Image (b) shows that once the **Center of Mass** is computed, we recover our beloved Second Law of Newton, $F = M_T \cdot g$. **There is one more important thing about the center of mass:** if the object has zero initial rotation (spin), placing the force of gravity at the center of mass will not produce any “torque” on the object, and hence it will not rotate under the action of gravity, which is what Newton observed (his famous apple was not spinning as it fell from the tree; that would have complicated his explanation immensely!). If we were to place $M_T \cdot g$ at a corner of the rectangle, it would cause the rectangle to rotate, which would go counter to experimental evidence. M_y and M_x are capturing that effect.

Center of Mass for a (Continuous) Planar Robot Link

Center of Mass in the Plane

In the following, we assume that $f : [x_{\min}, x_{\max}] \rightarrow \mathbb{R}$ and $g : [x_{\min}, x_{\max}] \rightarrow \mathbb{R}$ provide the upper and lower boundaries of a planar object, such as the robot link shown in Fig. 3.19.

Definition 3.33. *The following mechanical parameters are defined for a planar object:*

(a) **Total Moment about the y-axis:**

$$M_y = \rho \cdot h \cdot \int_{x_{\min}}^{x_{\max}} x \cdot (f(x) - g(x)) \, dx$$

(b) **Total Moment about the x-axis:**

$$M_x = \rho \cdot h \cdot \int_{x_{\min}}^{x_{\max}} \frac{1}{2} \cdot (f^2(x) - g^2(x)) \, dx$$

The formula for M_y is directly analogous to what we had for point masses, while the formula for M_x seems totally different. The justification is given in an optional reading section further down. Once we get over the hurdle for M_x , the definition of

the Cartesian center of mass is straightforward, namely

(c) ***x*-Coordinate of Center of Mass**

$$x_c := \frac{M_y}{M_T} = \frac{\int_{x_{\min}}^{x_{\max}} x \cdot (f(x) - g(x)) dx}{\int_{x_{\min}}^{x_{\max}} (f(x) - g(x)) dx} \quad (3.51)$$

(d) ***y*-Coordinate of Center of Mass**

$$y_c := \frac{M_x}{M_T} = \frac{\int_{x_{\min}}^{x_{\max}} \frac{1}{2} \cdot (f^2(x) - g^2(x)) dx}{\int_{x_{\min}}^{x_{\max}} (f(x) - g(x)) dx} \quad (3.52)$$

where the term $\rho \cdot h$ cancels in the numerators and denominators. Because of this, (x_c, y_c) is related to the **centroid** of an object. The two notions are ONLY the same when the mass is uniformly distributed throughout the object (e.g., ρ is a constant, independent of (x, y)).

Notes: $x \cdot (f(x) - g(x))$ and $(f^2(x) - g^2(x))$ have the same units, namely distance squared. This is because x , $f(x)$, and $g(x)$ are each measured in meters or centimeters. Consequently, the units work out fine for M_x and M_y , and most importantly for us, the units on x_c and y_c are distance. If anything else were true, we'd have a major problem.

A good source of SIMPLE examples is [How to find the centroid of simple composite shapes](#) by **Engineer4Free**.

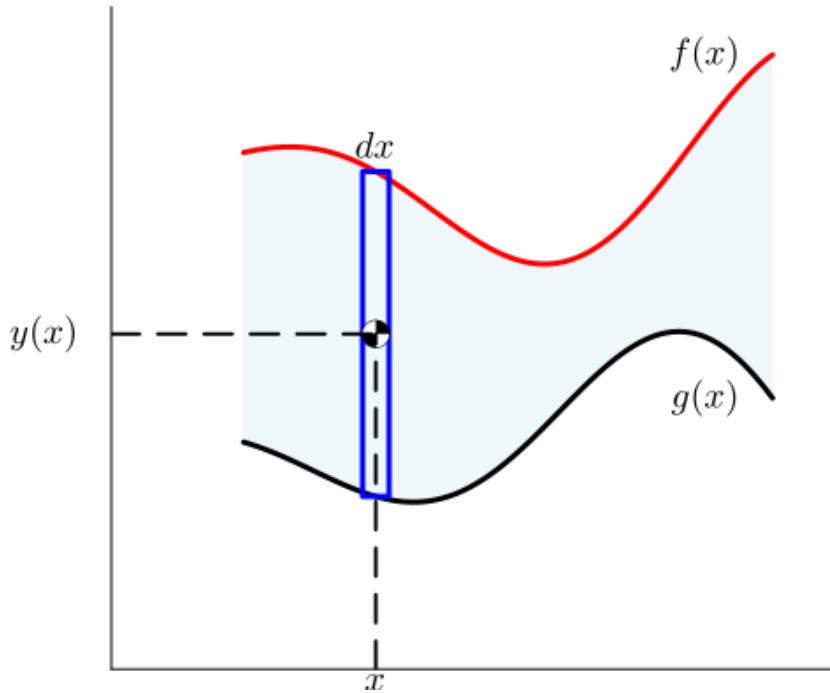


Figure 3.23: The surprising power of the rectangle is on full display. The (differential) area of the rectangle is height \times width or $dA(x) = (f(x) - g(x)) dx$. Its center (of mass) is located at $(x, \frac{f(x)+g(x)}{2})$, where the y -value is the mean of the top and bottom of the rectangle. Calculus tells us how to use an integral to add up the values associated with each rectangle to compute the object's overall mass and center of mass. The details are given in the text. The same process can be followed for moments of inertia.

Source of the Equation for Center of Mass. Figure 3.23 shows a body defined by functions two functions, $f : [a, b] \rightarrow \mathbb{R}$ and $g : [a, b] \rightarrow \mathbb{R}$; moreover, for all $x \in [a, b]$ we have $f(x) > g(x)$. Let (x_c, y_c) denote the Cartesian position of the center of mass.

Moreover, define

- **Differential Element and Mass:** Choose a vertical differential element of width dx at a distance x from the y -axis. The mass of this strip is

$$dm = \rho \cdot h \cdot (f(x) - g(x)) dx,$$

(density times z-thickness times height times differential width).

- **Condition of Balance from Physics for x_c :**

$$\int_a^b (x - x_c) \cdot \rho \cdot h \cdot (f(x) - g(x)) dx = 0,$$

which is equivalent to

$$\int_a^b x \cdot \rho \cdot h \cdot (f(x) - g(x)) dx = \int_a^b x_c \cdot \rho \cdot h \cdot (f(x) - g(x)) dx = x_c \cdot \int_a^b \rho \cdot h \cdot (f(x) - g(x)) dx$$

because x_c is a constant and can be taken out of the integral. Solving for x_c gives the answer below.

Similarly, the balance condition for y_c is

- **Condition of Balance from Physics for y_c :**

$$\int_a^b (y(x) - y_c) \cdot \rho \cdot h \cdot (f(x) - g(x)) dx = 0,$$

which is equivalent to

$$\int_a^b y(x) \cdot \rho \cdot h \cdot (f(x) - g(x)) dx = \int_a^b y_c \cdot \rho \cdot h \cdot (f(x) - g(x)) dx = y_c \cdot \int_a^b \rho \cdot h \cdot (f(x) - g(x)) dx$$

because y_c is a constant and can be taken out of the integral. Things become a bit more interesting when we note that the $y(x)$ -center of mass position of the infinitesimally small rectangle is its midpoint, namely, $y(x) = \frac{f(x)+g(x)}{2}$, as shown below.

Derivation of x_c :

- **Moment of the Differential Element along the x -axis:** This is given by:

$$dM_y = \rho \cdot h \cdot x \cdot (f(x) - g(x)) dx$$

- **Total Moment about the y -axis:**

$$M_y = \int_a^b x \cdot \rho \cdot h \cdot (f(x) - g(x)) dx$$

- **Total Mass of the Body:**

$$M_T = \int_a^b \rho \cdot h \cdot (f(x) - g(x)) dx$$

- **x -coordinate of the Center of Mass:**

$$x_c := \frac{M_y}{M_T} = \frac{\int_a^b x \cdot (f(x) - g(x)) dx}{\int_a^b (f(x) - g(x)) dx}$$

Interpretation: The formula for x_c is the weighted average of the x -positions of the differential elements, where each position is weighted by its corresponding mass. The center of mass is the “balance point” of the planar body.

Derivation of y_c :

- **Differential Element and Mass:** We use the same vertical differential element as before; it has width dx at a distance x from the y -axis. The height of this strip is $f(x) - g(x)$ and its differential mass is

$$dm = \rho \cdot h \cdot (f(x) - g(x)) dx.$$

- **y-Center of Mass for the Differential Element:** is the middle of the rectangle,

$$y(x) = \frac{f(x) + g(x)}{2}.$$

- **Moment of the Differential Element about the x -axis:**

$$\begin{aligned} dM_x &= \rho \cdot h \cdot y(x) \cdot (f(x) - g(x)) dx \\ &= \rho \cdot h \cdot \frac{f(x) + g(x)}{2} \cdot (f(x) - g(x)) dx \\ &= \rho \cdot h \cdot \frac{(f(x) - g(x))^2}{2} dx \end{aligned}$$

- **Total Moment about the x -axis:**

$$M_x = \int_a^b \rho \cdot h \cdot \frac{(f(x) - g(x))^2}{2} dx$$

- **Total Mass of the Body:**

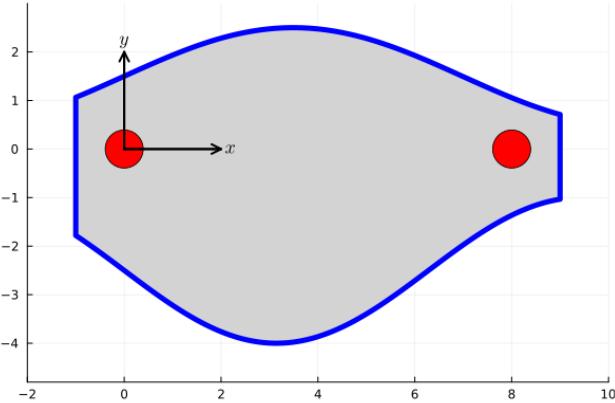
$$M_T = \int_a^b \rho \cdot h \cdot (f(x) - g(x)) dx$$

- **y -coordinate of the Center of Mass:**

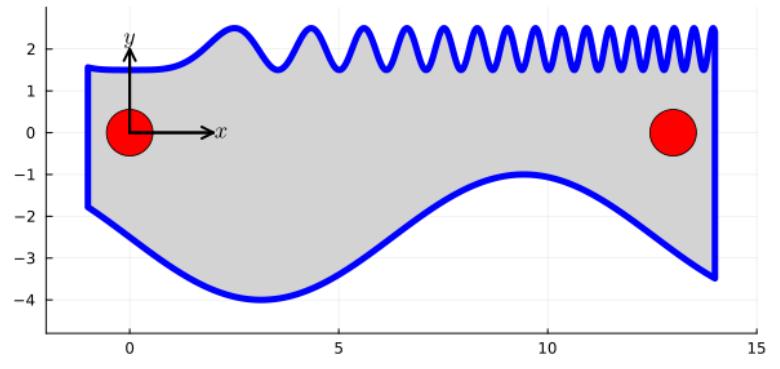
$$y_c = \frac{M_x}{M_T} = \frac{1}{2} \frac{\int_a^b (f(x) - g(x))^2 dx}{\int_a^b (f(x) - g(x)) dx},$$

where, because the density ρ and thickness h are constant, they cancel out.

There are many derivations available on YouTube. Here is one by **BlackPenRedPen**: [Centroid formulas of a region bounded by two curves](#).



(a)



(b)

Figure 3.24: We will compute the center of mass of these two planar objects, where (a) could be a “hip” and (b) could be an “old-fashioned saw”? Our point is that doing the calculations by hand would be extremely tedious, while in Julia, they are cake. For today’s robots, the mechanical parameters are determined numerically.

Example 3.34. Confirm the position of the center of mass for the standard link in Fig. 3.19 and determine the center of mass for the objects in Fig. 3.24. For all of the calculations, use Julia and quadgk.

Note: For the basic link,

$$f(x) = \begin{cases} \sqrt{|5^2 - (x - 5)^2|} & 0 \leq x \leq 5 \\ 5 & 5 < x \leq 55 \\ \sqrt{|5^2 - (55 - x)^2|} & 55 < x \leq 60 \end{cases}$$

and $g(x) = -f(x)$.

Solutions:

Basic Robot Link:

```

1 # Basic Link Top
2 function f(x)
3     #
4     # Basic Link
5     #
6     # Recall origin is placed at left end of the link
7     if x < 5
8         # Left end semicircle
9         # (x-r)^2 + y^2 = r^2
10        y = sqrt(abs(5^2 - (x-5)^2))
11    elseif x < 55
12        y = 5
13    else
14        # Right end semicircle
15        # (r + distancebetweenJoints - x)^2 + y^2 = r^2
16        y = sqrt(abs(5^2 - (5 + 50 - x)^2))
17    end
18    return y
19 end
20 g(x) = -f(x)
21
22 xmin = 0.0
23 xmax = 60.0
24
25 (xc, yc) = computeCoM(f, g, xmin, xmax)

```

Output

(30.000000027228815, 0.0)

Hip:

```

1 # Hip
2 #
3 # Bounding functions having sinusoidal waves
4 f(x) = sin(0.45x) + 1.75
5 g(x) = -1.5sin(0.50x) - 2.5
6 L = 10.0
7
8 xmin = 0.0
9 xmax = L
10

```

```

11 (xc, yc) = computeCoM(f, g, xmin, xmax)
12 center_of_mass_symbol((xc,yc), .5, 0.2)

```

Output

(4.1845164408139786, -0.4255811707111333)

Saw:

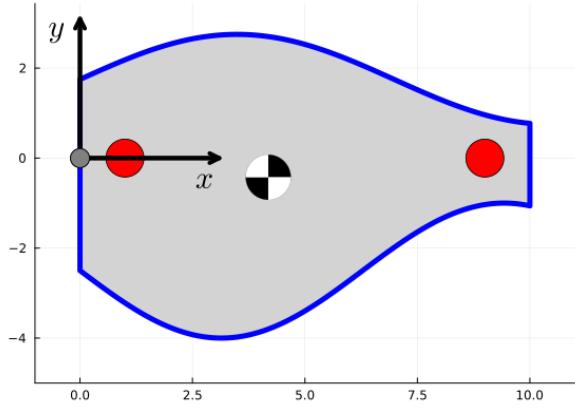
```

1 # Saw
2 #
3 # Functions having crazier sinusoidal waves
4 f(x) = sin(0.25x^2)^2 + 1.5
5 g(x) = -1.5sin(0.50x) - 2.5
6 L = 15.0
7
8 xmin = 0.0
9 xmax = L
10
11 (xc, yc) = computeCoM(f, g, xmin, xmax)

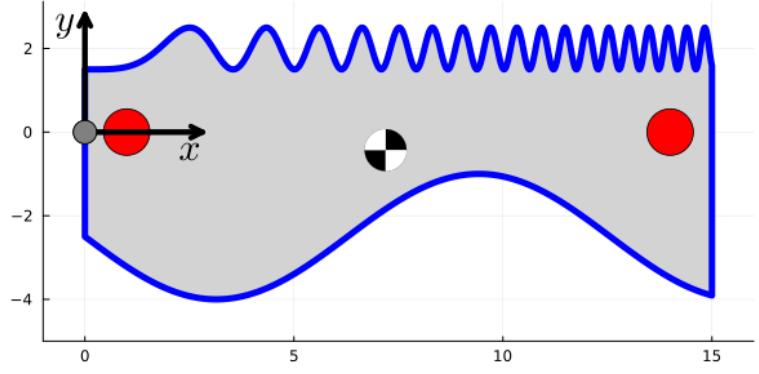
```

Output

(7.192793764557787, -0.43051124387170386)



(a)



(b)

```

1 # Compute center of mass
2 using QuadGK
3
4 function computeCoM(f, g, xmin, xmax)
5     #
6     # Note: we are leaving out the constants for density and thickness
7     # because they cancel in the end. Hence, mass is replaced with area.
8     #
9     # f = upper bounding function
10    # g = lower bouding function
11    # [xmin, xmax] = domain of f and g

```

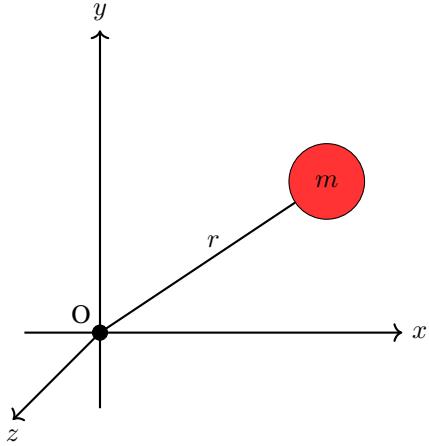
```

12
13 # Set up functions for calling quadgk
14 areaIntegrand(x) = f(x) - g(x)
15 MxIntegrand(x) = x*areaIntegrand(x)
16 MyIntegrand(x) = 0.5*(f(x)^2 - g(x)^2)
17
18 # Apply quadgk
19 A,      = quadgk(areaIntegrand, xmin, xmax)      # Area
20 Mx,     = quadgk(MxIntegrand, xmin, xmax)        # Total x Moment
21 My,     = quadgk(MyIntegrand, xmin, xmax)        # Total y Moment
22
23 # (xc, yc) of object
24 xc = Mx/A
25 yc = My/A
26
27 # return values
28 return (xc, yc)
29 end

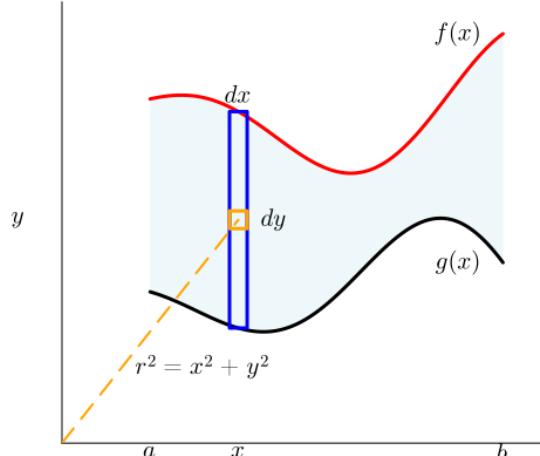
```

Output

computeCoM (generic function with 1 method)



(a) Physics tells us the moment of inertia for a point mass.



(b) Computing the moment of inertia for a planar area uses Calculus.

Figure 3.26: Illustrations of the moment of inertia: (a) for a point mass about the z -axis and (b) for a planar area using differential elements. The amazing power of the rectangle is on even fuller display for the computation of the moment of inertia. We want to compute the moment of inertia about the z -axis for a continuous body, based on what we know about the moment of inertia of a point mass, m , namely, $I_z = mr^2$, where r is the distance of the point mass from the origin. We apply this knowledge to compute the moment of inertia of the pseudo point mass defined by the differential orange rectangle centered at (x, y) , where $\rho \cdot h \cdot dx \cdot dy$ plays the role of a point mass and $\sqrt{x^2 + y^2}$ is the distance from the origin. We use Calculus to sum up the differential inertias along the y -direction, yielding a formula for the differential inertia of a very thin rectangle (its width is dx). We use Calculus a second time to add up the inertia associated with each infinitesimally thin rectangle to compute the object's overall moment of inertia. Incredible!

3.4.5 (Optional Read:) Using Again the Power of the Rectangle to Compute the Inertia of a (Planar) Robot Link

This is the first and only time in the textbook that we “more-or-less” do a two-dimensional integral. It is perhaps the ultimate illustration of a Superpower of Calculus: adding up Infinitesimal Objects in 2d instead of 1D! Let’s get to work. We define the

differential mass contained in the orange rectangle of Fig. 3.26b by

$$dM(x, y) := \rho \cdot h \cdot dy \cdot dx, \quad (3.53)$$

which is density, times thickness, times height, times width. The distance squared of the differential mass from the origin is

$$r^2 = x^2 + y^2. \quad (3.54)$$

To compute the differential inertia dI_z of the entire rectangle of height $(f(x) - g(x))$, we need to do an integral in the y -direction, namely

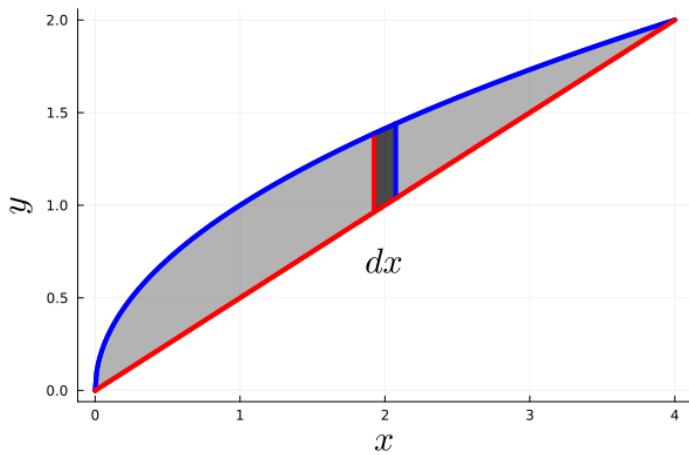
$$\begin{aligned} dI_x &= \int_{g(x)}^{f(x)} r^2 \cdot dM(x, y) \\ &= \int_{g(x)}^{f(x)} (x^2 + y^2) \cdot \rho \cdot h \cdot dy \cdot dx \text{ (factor out constants)} \\ &= \rho \cdot h \cdot dx \cdot \int_{g(x)}^{f(x)} (x^2 + y^2) \cdot dy \text{ (}dx\text{ is treated as a constant when integrating with respect to }y\text{)} \\ &= \rho \cdot h \cdot dx \cdot \left(x^2 \cdot y + \frac{1}{3} y^3 \right) \Big|_{g(x)}^{f(x)} \text{ (this time, }x\text{ is treated as a constant when integrating with respect to }y\text{)} \\ &= \rho \cdot h \cdot \left[x^2 \cdot (f(x) - g(x)) + \frac{1}{3} \cdot (f^3(x) - g^3(x)) \right] \cdot dx \text{ (moved }dx\text{ to the end).} \end{aligned} \quad (3.55)$$

Wow! That's one big impressive formula for an infinitesimally thin rectangle. Indeed. Now, to get the entire moment of inertia about the z -axis, we use Calculus to sum up along the x -direction all of the differential moments of inertia, yielding,

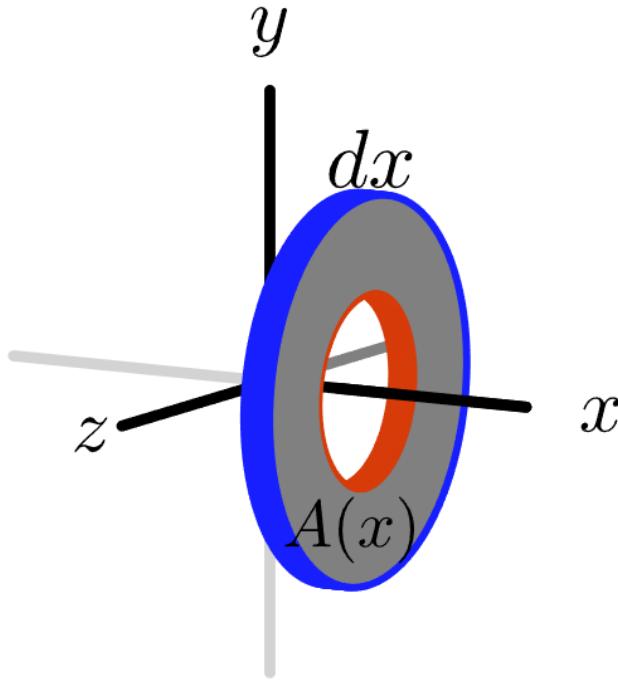
$$I_z = \int_a^b dI_z = \rho \cdot h \cdot \int_a^b \left[x^2 \cdot (f(x) - g(x)) + \frac{1}{3} \cdot (f^3(x) - g^3(x)) \right] dx. \quad (3.56)$$

The integral can be evaluated in closed form in simple toy examples. However, for most robot links of interest, numerical methods are used to evaluate the integral. QuadGK to the rescue!

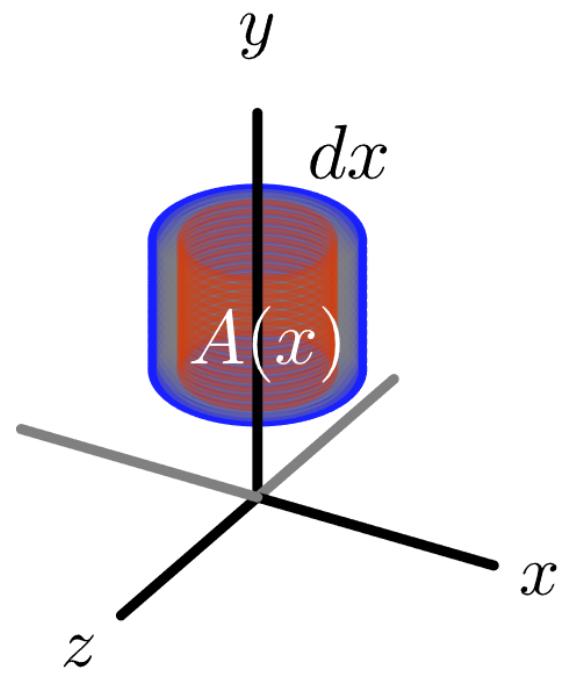
P.S. We really just did a simple form of two-dimensional integration, a mainstay of Michigan's Math 215 Multivariable Calculus



(a)



(b)



(c)

Figure 3.27: Differential volumes, dV , in the form of disk/washers and shells. The area upper bounded by $f : [a, b] \rightarrow \mathbb{R}$ in blue and lower bounded by $g : [a, b] \rightarrow \mathbb{R}$ in red is shown in (a), for $[a, b] = [0, 4]$. The area enclosed is shaded gray. A solid of revolution results by rotating the indicated area about either the x -axis or the y -axis. Can you visualize it? Through the magic of Calculus, we can compute the volume of the solid of rotation, whether or not we can visualize the resulting solid in our mind. The key point is the following: The Riemann integral comes down to understanding narrow rectangles. Within the area enclosed in (a) is a dark gray “rectangle” of width dx , the same rectangle we would use to compute the area enclosed by f and g . The result of rotating the narrow rectangle about the x -axis is the **washer** shown in (b); a washer is a thin disk with a hole in it, having thickness dx . Because its thickness is non-zero, the washer has a small, but non-zero volume, dV . The same dark gray rectangle is rotated about the y -axis in (c). The blue cylinder is traced out by the right side of the rectangle, and the red cylinder is traced out by the left side of the rectangle. The result is a **shell**, that is, a cylinder with thickness dx . Due to its non-zero thickness, the shell has a small but non-zero volume, dV . Calculus tells us how to add up the infinitesimal volumes dV to get the overall volume.

3.4.6 Solids of Revolution or More on the Power of the Rectangle

Consider a shape in the plane enclosing a non-zero area, as illustrated in Fig. 3.27. A **solid of revolution** is the result of rotating the shape about either the x -axis or the y -axis. Can you visualize the two shapes? They are not shown in the Figure. When rotating about the x -axis, can you imagine a blue parabolic-shaped vase with a red cone-shaped interior that is laying on its side? Maybe. When rotating about the y -axis, can you imagine an upright red cone-shaped object with an anti-parabolic-(square-root)-shaped blue interior? Maybe. To compute the volume of these two shapes, you don't need to be able to sketch them! It's nice if you can sketch the shapes, but that might just confuse you when computing their respective volumes. As shown in Fig. 3.27, rotating a rectangle of width dx about the x -axis gives a washer (aka, a disk with a hole in the middle) with thickness dx , while rotating that same rectangle about the y -axis produces a shell, that is, a cylinder with thickness dx . Each of these objects has a differential volume equal to $A(x)$ times dx . The total volume is then

$$V_{\text{solid}} := \int_a^b dV(x) := \int_a^b A(x) dx. \quad (3.57)$$

Computing the Volume

Without even sketching the solid of revolution that results from rotating the area about an axis, we can understand how to compute its volume. Sounds crazy, right? The secret is that Riemann Integration is all about narrow rectangles. Rotating a rectangle is much easier to understand than rotating an arbitrary shape. As shown in Fig. 3.27, rotating a rectangle of width dx about the x -axis gives a washer (aka, a disk with a hole in the middle) with thickness dx , while rotating that same rectangle about the y -axis produces a shell, that is, a cylinder with thickness dx . Each of these objects has a differential volume equal to $A(x)$ times dx .

- **Area and Differential Volume of the Disk:** The outer radius is $f(x)$ and the inner radius is $g(x)$, and its thickness is dx . Hence,

$$A(x) := \pi (f^2(x) - g^2(x)) \text{ and } dV_{x\text{-axis}} := \pi (f^2(x) - g^2(x)) dx.$$

- **Area and Differential Volume of the Cylinder/Shell:** The top edge of the cylinder/shell is $f(x)$ and its bottom edge is $g(x)$, giving a height of $h := (f(x) - g(x))$. The radius is $r := x$, and its thickness is dx . Because the surface area of a cylinder is $2\pi \cdot r \cdot h$, we have that

$$A(x) := 2\pi x (f(x) - g(x)) \text{ and } dV_{y\text{-axis}} := 2\pi x (f(x) - g(x)) dx.$$

Proposition 3.35: Volume of a Solid of Revolution

Key Assumptions: The area lies totally in the first-quadrant of the (x, y) -plane, that is, the quadrant where $x \geq 0$ and $y \geq 0$. Moreover, the area is delimited by two functions $f : [a, b] \rightarrow \mathbb{R}$ and $g : [a, b] \rightarrow \mathbb{R}$, and for all $0 \leq a \leq x \leq b$, $f(x) \geq g(x)$. The last condition means that f is the upper bound on the area, g is the lower bound, and the two functions do not enclose any “negative area”.

Under the above assumptions, the volumes of the two standard solids of revolution are computed as follows:

$$V_{x\text{-axis}} := \int_a^b dV_{x\text{-axis}} = \int_a^b \pi (f^2(x) - g^2(x)) dx \quad (\text{called the } \mathbf{disc/washer method} \text{ for rotation about the } x\text{-axis}) \quad (3.58)$$

$$V_{y\text{-axis}} := \int_a^b dV_{y\text{-axis}} = \int_a^b 2\pi x (f(x) - g(x)) dx \quad (\text{called the } \mathbf{shell method} \text{ for rotation about the } y\text{-axis}). \quad (3.59)$$

The terms disc/washer and shell are defined in Fig. 3.27 and in the text below.

Remark 3.36. On YouTube, there are tons of videos that illustrate these two methods. You need to be careful of two points:

- We are applying the disc-washer method to a solid formed by rotating an area about the x -axis. If you apply the disc-washer method when rotating the area about the y -axis, you will have to use the inverse functions, $x = f^{-1}(y)$ and $x = g^{-1}(y)$. Watch out for this in the videos. Also, what do you do when the inverse functions do not exist? Ouch!
- We are applying the shell method to a solid formed by rotating an area about the y -axis. If you apply the shell method when rotating the area about the x -axis, you will have to use the inverse functions, $x = f^{-1}(y)$ and $x = g^{-1}(y)$. Watch out for this in the videos. Also, what do you do when the inverse functions do not exist? Ouch! You apply our method!
- Our suggestion: do not put yourself in a situation where you need inverse functions; they are typically hard to compute, and their calculation is easily messed up. If you apply the methods as we have presented them, you are good to go. No inverse functions are needed.
- (Optional Read:) A minor point: On rare occasions, an area is delimited by $x = F(y)$ and $x = G(y)$, and thus your functions are parameterized by y instead of x . This means the roles of the x and y axes have been flipped. In this case, you would draw the dark rectangle as wide in the x -direction and short in the y -direction, flipping the axes that produce a disc/washer versus a shell. Hence, you have to flip the formulas in (3.58) and (3.59) for the disc/washer and shell methods.
- (Optional Read:) A minor point: There are a few videos on YouTube that treat a solid of revolution about an arbitrary axis. Your author has not seen that show up in an engineering course. If it does, you can ask the instructor to teach you how to do it!

Justification of (3.58): Consider Fig. 3.27-(a). When the darkly shaded rectangle of width dx is rotated about the x -axis, it traces out the **washer** shown in Fig. 3.27-(b). The outside radius of the washer is $r_2 = f(x)$, the inside radius is $r_1 = g(x)$, and its width is dx . It follows that the area of the washer is $A(x) = \pi(r_2^2 - r_1^2) = \pi(f^2(x) - g^2(x))$. The washer is the solid of revolution traced out by the small rectangle when it is rotated about the x -axis. The **differential volume** of the washer is its area times width (or thickness), that is, $dV = A(x) dx$. Hence, the total volume of the solid of revolution about the x -axis is the sum of the differential volumes of revolution. That is, $V = \int_a^b dV = \int_a^b A(x) dx$, we simply add up the differential volumes $dV := A(x) dx$ of each washer from the lower limit of x to its upper limit.

Justification of (3.59): Consider again Fig. 3.27-(a). When the darkly shaded rectangle of width dx is rotated about the y -axis, it traces out the **shell** shown in Fig. 3.27-(c). A “shell” is a cylinder with a thin but non-zero thickness; in our case, the thickness is dx . The area of the cylinder is $2\pi \cdot r \cdot h$, where its radius $r = x$ and height $h = f(x) - g(x)$. Therefore, the surface area of the shell is $A(x) = 2\pi \cdot x \cdot (f(x) - g(x))$, and the differential volume of the resulting solid of revolution is $dV = A(x) dx = 2\pi \cdot x \cdot (f(x) - g(x)) dx$. Hence, the total volume of the solid of revolution about the y -axis is the sum of the differential volumes of revolution. That is, $V = \int_a^b dV = \int_a^b A(x) dx$, we simply add up the differential volumes $dV := A(x) dx$ of each shell from the lower limit of x to its upper limit.

Example 3.37. For the two functions $f : [0, 4] \rightarrow \mathbb{R}$ by $f(x) = \sqrt{x}$ and $g : [0, 4] \rightarrow \mathbb{R}$ by $g(x) = 0.5x$, compute the volumes for the resulting solids of revolution about the x - and y -axes. After this example, we'll learn how to sketch the volumes.

Rotation about x -axis:

```

1 # Volume of revolution about the x-axis
2 using QuadGK
3
4 f(x) = sqrt(x)
5 g(x) = 0.5x
6 xmin = 0.0
7 xmax = 4.0
8
9 # Area of disc
10 areaDiscWasher(x) = pi * ( (f(x))^2 - (g(x)^2) )
11
12 Volume_x, Error = quadgk(areaDiscWasher, xmin, xmax)

```

Output

(8.377580409572783, 0.0)

Rotation about y -axis:

```
1 # Volume of revolution about the y-axis
2 using QuadGK
3
4 f(x) = sqrt(x)
5 g(x) = 0.5x
6 xmin = 0.0
7 xmax = 4.0
8
9 # Area of disc
10 areaShell(x) = 2*pi * x * (f(x) - g(x))
11
12 Volume_y, Error = quadgk(areaShell, xmin, xmax)
```

Output

(13.40412865472173, 1.244413129087435e-7)

Note: The volumes are not the same.

How to Sketch the Volumes

Figure 3.28 walks you through the steps for sketching a solid of revolution about the x -axis. The same procedure works for revolving about the y -axis. To get the hang of the process, seeing “live” examples worked out step-by-step is extremely helpful.

Recommended Video: [Disk Washer Method to Find the Volume of Solids of Revolution](#) by Quoc Dat Phung, an engineering student in Canada. The video has excellent graphics to support the discussion. The first two examples, which take 10 minutes total at 1X, are adequate to get you going. Quoc Dat uses the term “antiderivative”. Don’t let it bother you; we will cover the concept in an upcoming Chapter. Roughly speaking, it’s another term for the “indefinite integral” that we covered in Prop. 3.3: said another way, the antiderivative of x^k is $\frac{x^{k+1}}{k+1}$. The first two examples worked by Quoc Dat are for rotations about the x -axis, and he applies the disk/washer method.

ChatGPT Prompt for Students: “Using the disk/washer method, determine the volume of the solid of revolution formed by rotating about the x -axis the region bounded by $y = \sqrt{x}$ and $y = \frac{1}{2}x$ for x between zero and four. In addition, using the shell method, determine the volume of the solid of revolution formed by rotating about the y -axis the region bounded by $y = \sqrt{x}$ and $y = \frac{1}{2}x$ for x from 0 to 4. Provide well-documented Julia code for both methods, explaining each step in its solution. Additionally, generate supporting plots to illustrate the methods explained in the code.”

Example 3.38. Volume of a Basketball: Use the disk/washer method or the shell method to compute the volume of an NBA basketball.

Solution: According to [Wikipedia](#), an NBA basketball has a circumference of 75 centimeters (cm), meaning its radius is $r = \frac{75}{2\pi} = 11.9366$ cm. We’ll define $f : [0, r] \rightarrow \mathbb{R}$ by $f(x) = \sqrt{r^2 - 1}$ and $g : [0, r] \rightarrow \mathbb{R}$ by $g(x) = 0$. The graph of f is a quarter of a circle, as shown in Fig. 3.31.

```
1 # Volume of revolution about the x-axis
2 using QuadGK
3
4 r = 75 / (2*pi)
5
6 f(x) = sqrt(r^2 - x^2)
7 g(x) = 0.0x
```

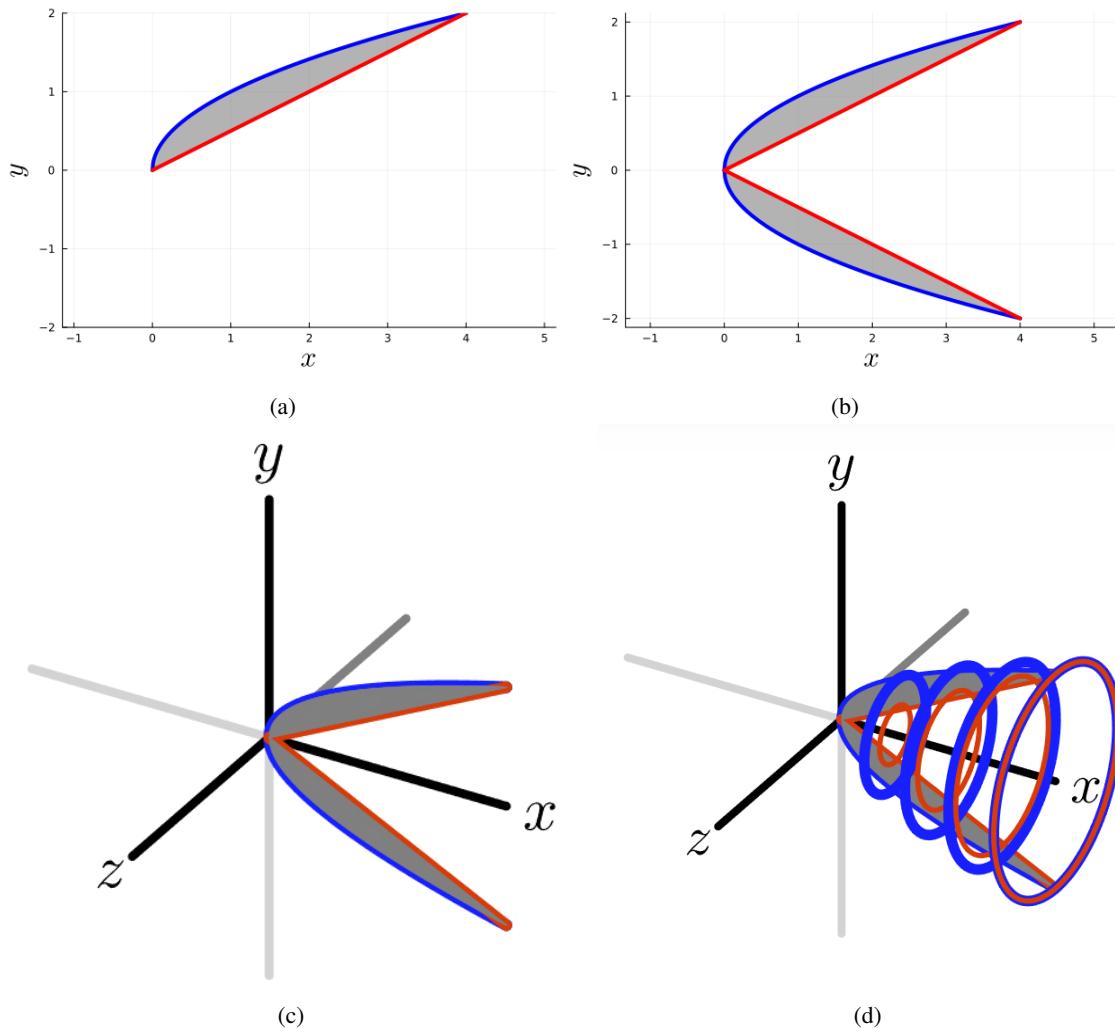


Figure 3.28: **How to Sketch a Solid of Revolution.** (a) shows the area of interest in the (x, y) -plane. To form a rotation about the x -axis, make a mirror image of the area below the x -axis, as in (b); this is really the area rotated by 180 degrees about the x -axis. Next, as in (c), add in a z -axis to form a 3D object. The final step is to realize that each point (x, y) in the original area forms a circle in the (y, z) -plane when rotated about the x -axis. The blue circles in (d) delineate the outer surface of the solid of revolution, and the red circles delineate the inner surface of the solid.

```

8 xmin = 0
9 xmax = r
10
11 # Area of disc
12 areaDiscWasher(x) = pi * ( (f(x) ^2 - (g(x) ^2) )
13
14 VolumeSemisphere_x, Error = quadgk(areaDiscWasher, xmin, xmax)
15
16 VolumeBasketball = 2 * VolumeSemisphere_x
```

Output

7124.145724851874

```

1 # Volume of revolution about the y-axis
2 using QuadGK
3
```

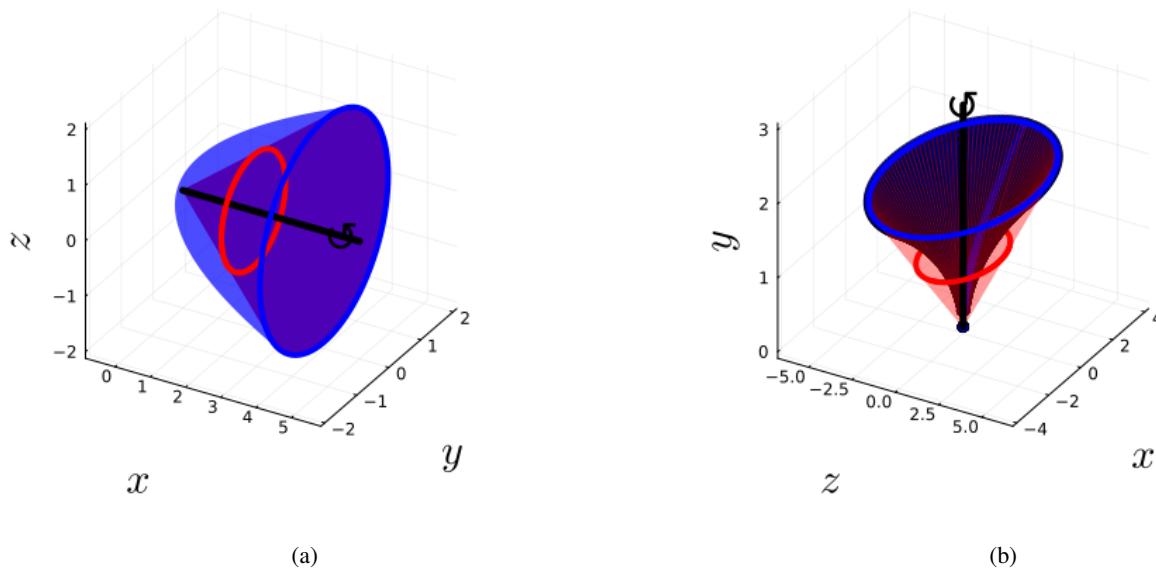


Figure 3.29: Solids of Revolution for Example 3.37 and Fig.3.28 with bounding functions $f : [0, 4] \rightarrow \mathbb{R}$ and $g : [0, 4] \rightarrow \mathbb{R}$, with $f(x) \geq g(x)$. (a) Rotation of the area about the x -axis. (b) Rotation of the area about the y -axis. f forms the outer boundary when the area is rotated about the x -axis, while it forms the inner boundary when the area is rotated about the y -axis. The opposite is true for g .

```

4 r = 75 / (2*pi)
5
6 f(x) = sqrt(r^2 - x^2)
7 g(x) = 0.0x
8 xmin = 0
9 xmax = r
10
11 # Area of disc
12 areaShell(x) = 2*pi * x * (f(x) - g(x))
13
14 VolumeSemisphere_y, Error = quadgk(areaShell, xmin, xmax)
15
16 VolumeBasketball = 2 * VolumeSemisphere_y

```

Output

7124.145729241054

Our estimate is, therefore, Volume = 7124.145727 cm³ ± 4.39e-06

Other Recommended Videos:

- [Cylindrical Shells to Find the Volume of Solids of Revolution](#) by Quoc Dat Phung has excellent graphics to support what he is teaching.
- [Disk/Washer vs. Cylindrical Shell...when to use which?](#) by Quoc Dat Phung has excellent graphics to support what he is teaching. **Note:** For us, the answer is clear for functions $y = f(x)$ and $y = g(x)$: when rotating about the x -axis, use the Disk/Washer Method, and when rotating about the y -axis, use the Shell Method. This avoids having to compute inverse functions.
- [Disc/Washer Method vs. Shell Method \(rotated about different lines\)](#) by **BlackPenRedPen** has clear, hand-drawn graphics so that you too can make nice sketches.

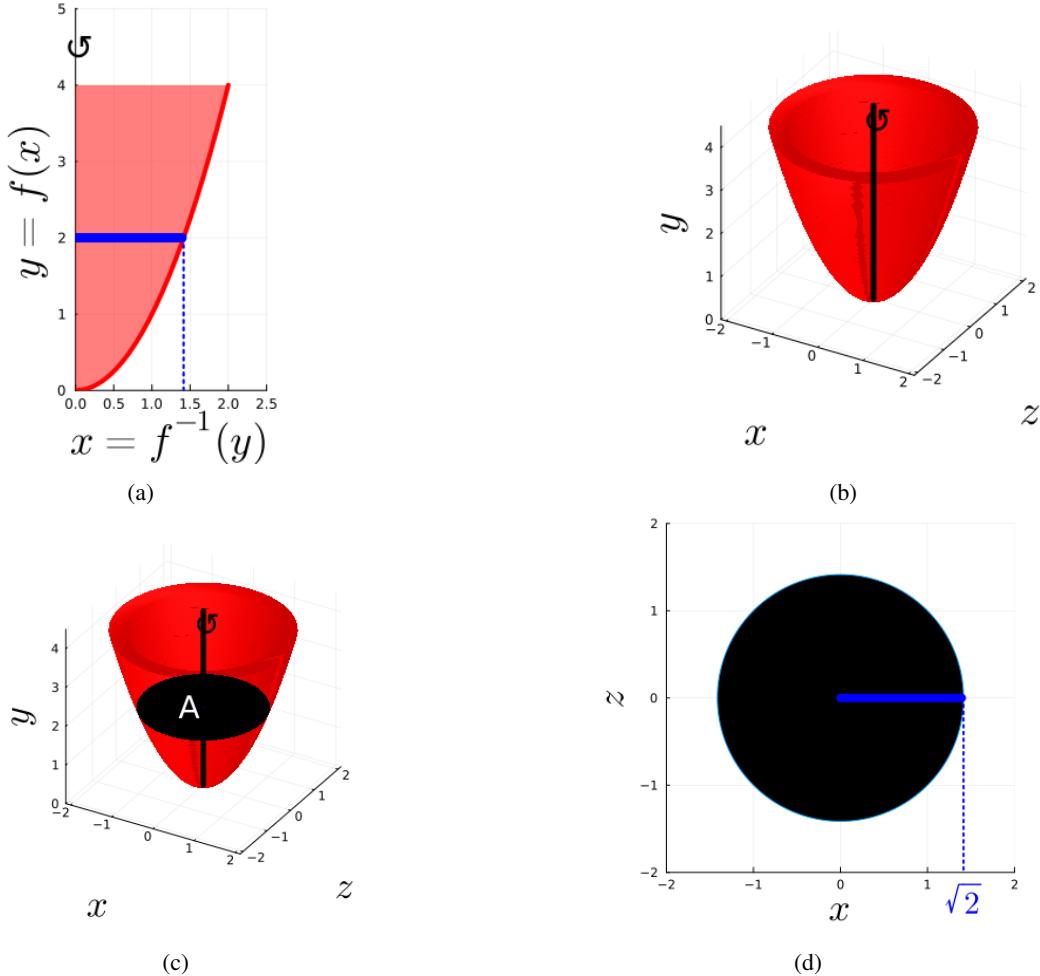


Figure 3.30: **Applying the disk/washer method to a solid of revolution about the y -axis requires inverse functions.** The area between the line $y = 4$ and $y = x^2$ has been shaded in (a). In (b), the shaded area has been rotated around the y -axis. This means that each point $(x = f^{-1}(y), y)$ in the shaded area will trace out a circle of radius $f^{-1}(y)$ in the (x, z) -plane at a height of $y = f(x)$. In particular, the blue bar in (a), which has thickness dy and radius $r(y) = f^{-1}(y)$, traces out the disk in (c). In (d), the disk has been projected to the (x, z) -plane, confirming it is a “perfect” circle. Hence, the area in (c) is $A = \pi r^2$, where $r = f^{-1}(y)$. You will be wise to apply the shell method in cases like this one to avoid having to compute inverses of functions.

- **Shell method for volume of revolution (rotated about different axis and lines)** by **BlackPenRedPen** has clear, hand-drawn graphics so that you too can make nice sketches.

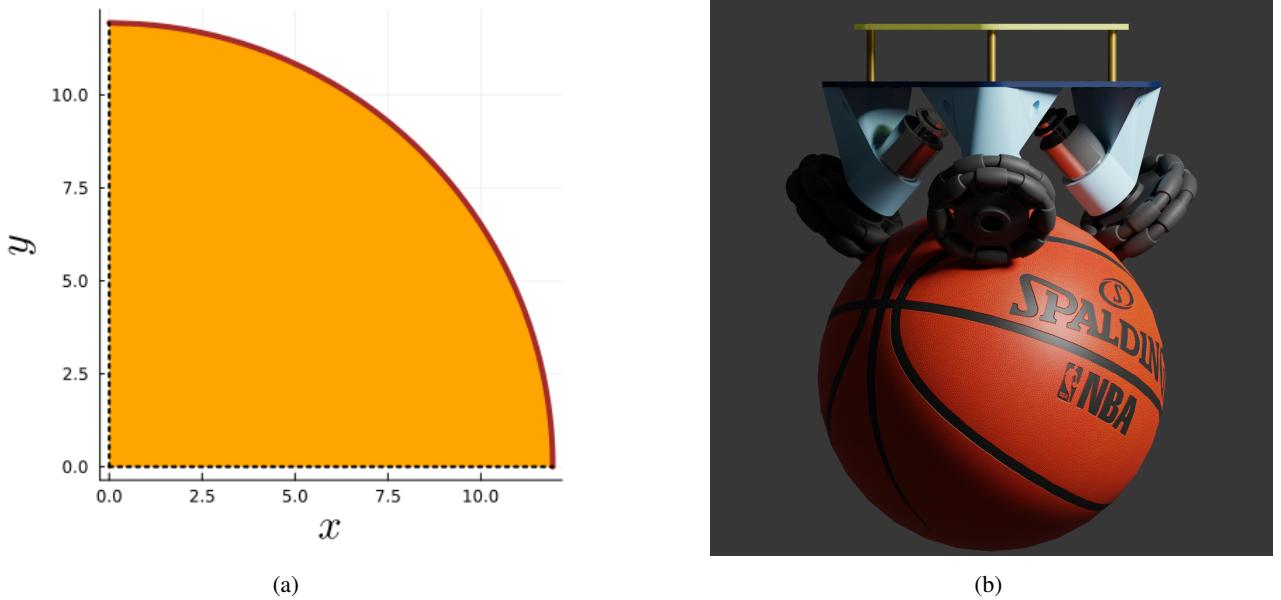


Figure 3.31: We can revolve a quarter circle around either axis to form a hemisphere. The volume of a sphere will be twice that of a hemisphere.

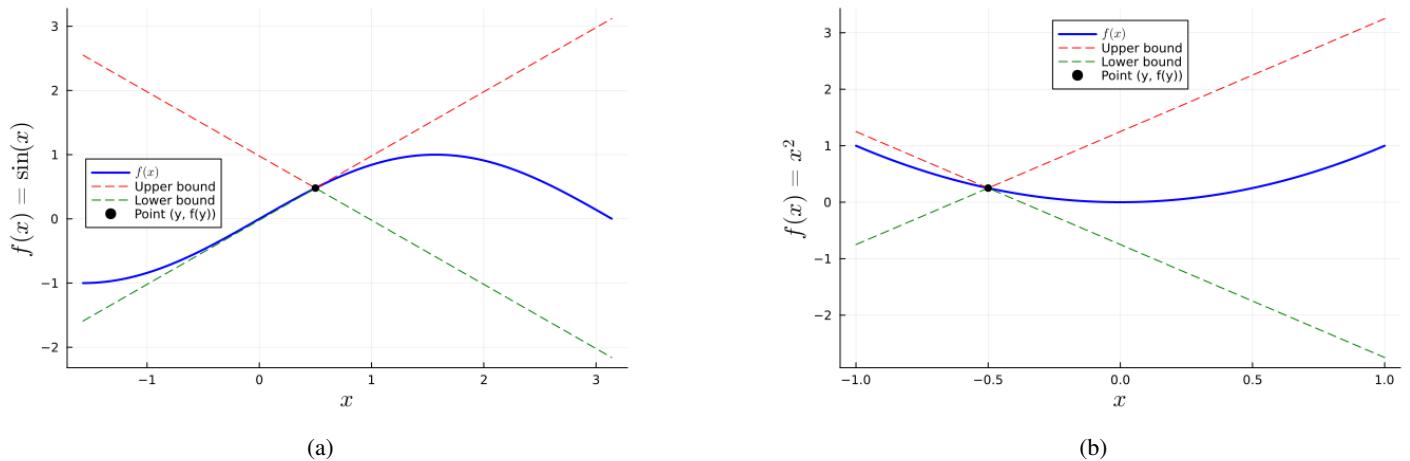


Figure 3.32: (**Lipschitz Continuity**) is more restrictive than “normal continuity”. You can think of Lipschitz continuity as a rule that puts a bound on how quickly a function can change as you move from one point to another along its graph, forcing the function to behave as if it were “locally linear”. This means the function doesn’t spike too steeply or drop too suddenly, ensuring a certain smoothness in its graph.

3.5 (Optional Read:) Proofs Associated with the Chapter

We open this section by providing proof of the existence of the Riemann-Darboux Integral for functions that satisfy a stronger condition than the “usual notion of continuity” as we’ve used that term so far in the course.

Lipschitz Continuity

Definition 3.39. A function $f : [a, b] \rightarrow \mathbb{R}$ is **Lipschitz continuous** if there exists a finite constant $L < \infty$ such that for all $x, y \in [a, b]$,

$$|f(x) - f(y)| \leq L \cdot |x - y|. \quad (3.60)$$

Note: Essentially, Lipschitz continuity means the function cannot grow faster than a linear function with slope L . More precisely, for any point $y \in [a, b]$ and all points $x \in [a, b]$

$$f(y) - L|x - y| \leq f(x) \leq f(y) + L|x - y|, \quad (3.61)$$

as illustrated in Fig. 3.32. Near each point $y \in [a, b]$, the local change in f as a function of $x - y$ is linear with slope L . Indeed, for $x > y$, $|x - y| = x - y$, and we have $f(x) \leq f(y) + L(x - y)$ (a constant plus a slope times the change of x with respect to y), while for $x < y$, $|x - y| = -(x - y)$, and we have $f(x) \geq f(y) + L(x - y)$ (one again, a constant plus a slope times the change of x with respect to y).

Note: From Prop. 1.9-(c), the “no-name property of inequalities”, (3.60) is equivalent to

$$-L \cdot |x - y| \leq f(x) - f(y) \leq L \cdot |x - y|.$$

Adding $f(y)$ to each term gives,

$$f(y) - L \cdot |x - y| \leq f(x) \leq f(y) + L \cdot |x - y|.$$

Note: Due to the inequality (3.61), a Lipschitz continuous function cannot have a “jump”, and hence every Lipschitz continuous function is also “continuous” in the sense we have been using the term in this textbook. Indeed, considering Fig. 1.2-(b) or -(c), drawing a discontinuous function requires a line with an infinite slope (i.e., a vertical line) or lifting your pen from the paper. In the next Chapter, we’ll be more precise about the notion of continuity. ■

Proposition 3.40: (Existence of the Riemann-Darboux Integral for Lipschitz Continuous Functions)

Suppose that $f : [a, b] \rightarrow \mathbb{R}$ is Lipschitz continuous. Then its Riemann-Darboux integral $\int_a^b f(x) dx$ exists and is finite.

Proof: We must check the conditions in Definition 3.2 are satisfied. To do that, we recall some of the details. For $a < b$ (finite) real numbers, and $n \geq 1$, define $\Delta x := \frac{b-a}{n}$ and $x_i := a + (i-1) \cdot \Delta x$, for $1 \leq i \leq n+1$. Then,

- (a) $\text{Area}_n^{\text{Low}} := \sum_{i=1}^n h_i^{\text{Low}} \cdot \Delta x$, for $h_i^{\text{Low}} := \min_{x \in [x_i, x_{i+1}]} f(x)$ is the **Riemann lower sum**, and
- (b) $\text{Area}_n^{\text{Up}} := \sum_{i=1}^n h_i^{\text{Up}} \cdot \Delta x$, for $h_i^{\text{Up}} := \max_{x \in [x_i, x_{i+1}]} f(x)$ is the **Riemann upper sum**.

We need to show that

- $\lim_{n \rightarrow \infty} \text{Area}_n^{\text{Low}}$ exists and is finite,
- $\lim_{n \rightarrow \infty} \text{Area}_n^{\text{Up}}$ exists and is finite, and
- the two limits are equal.

We first show the last property: if the two limits exist and are finite, then they converge to the same value.

The key step is to estimate the max and min of $f : [x_i, x_{i+1}] \rightarrow \mathbb{R}$. But from (3.61), for all $x \in [x_i, x_{i+1}]$

$$f(x_i) - L|x - x_i| \leq f(x),$$

and hence,

$$h_i^{\text{Low}} := \min_{x \in [x_i, x_{i+1}]} f(x) \geq \min_{x \in [x_i, x_{i+1}]} (f(x_i) - L|x - x_i|) = f(x_i) - L|x_{i+1} - x_i| = f(x_i) - L\Delta x.$$

Similarly, from (3.61), for all $x \in [x_i, x_{i+1}]$

$$f(x) \leq f(x_i) + L|x - x_i|,$$

and hence,

$$h_i^{\text{Up}} := \max_{x \in [x_i, x_{i+1}]} f(x) \leq \max_{x \in [x_i, x_{i+1}]} (f(x_i) + L|x - x_i|) = f(x_i) + L|x_{i+1} - x_i| = f(x_i) + L\Delta x.$$

In case the math in these steps is unclear, the bounds can be deduced from Fig. 3.32.

From these inequalities,

$$\begin{aligned}
\text{Area}_n^{\text{UP}} - \text{Area}_n^{\text{LOW}} &= \sum_{i=1}^n (h_i^{\text{UP}} - h_i^{\text{LOW}}) \Delta x \\
&\leq \sum_{i=1}^n (h_i^{\text{UP}} - h_i^{\text{LOW}}) \Delta x \\
&\leq \sum_{i=1}^n (f(x_i) + L\Delta x - f(x_i) + L\Delta x) \Delta x \\
&\leq \sum_{i=1}^n 2L\Delta x \cdot \Delta x \\
&\leq n2L \left(\frac{b-a}{n} \right)^2 \\
&\leq 2L \frac{(b-a)^2}{n}.
\end{aligned}$$

Hence,

$$\lim_{n \rightarrow \infty} (\text{Area}_n^{\text{UP}} - \text{Area}_n^{\text{LOW}}) \leq \lim_{n \rightarrow \infty} 2L \frac{(b-a)^2}{n} = 0.$$

Because $\text{Area}_n^{\text{UP}} - \text{Area}_n^{\text{LOW}} \geq 0$, we deduce that

$$\lim_{n \rightarrow \infty} (\text{Area}_n^{\text{UP}} - \text{Area}_n^{\text{LOW}}) = 0,$$

and therefore, if the two limits exist and are finite, then they converge to the same value.

To show the convergence of the Riemann Upper or Lower Sum requires a result from Real Analysis that is beyond the scope of this course: the **Monotone Convergence Theorem**, which is taught in Michigan's Math 451 *Advanced Calculus I*. The theorem states that if a sequence is monotonically increasing and bounded from above, then it has a limit, and the limit is finite; similarly, if a sequence is monotonically decreasing and bounded from below, then it has a limit, and the limit is finite.

For us, the sequences of interest are $\text{Area}_n^{\text{UP}}$ and $\text{Area}_n^{\text{LOW}}$, for n a power of two, that is, $n = 2^k$. **Why Powers of 2?** When n is chosen as a power of 2, each step from n to $2n$ involves splitting each subinterval into exactly two equal parts. This ensures every new subinterval's maximum cannot exceed the maximum of the “parent interval” in the previous partition. To see this, consider a subinterval $[x_i, x_i + \Delta x]$ and divide it into two equal halves,

$$[x_i, x_i + \Delta x] = [x_i, x_i + \frac{\Delta x}{2}] \cup [x_i + \frac{\Delta x}{2}, x_i + \Delta x].$$

Then, upon defining,

$$h_i^{\text{UP}} := \max_{x \in [x_i, x_i + \Delta x]} f(x)$$

$$\bar{h}_{2i}^{\text{UP}} := \max_{x \in [x_i, x_i + \frac{\Delta x}{2}]} f(x)$$

$$\bar{h}_{2i+1}^{\text{UP}} := \max_{x \in [x_i + \frac{\Delta x}{2}, x_i + \Delta x]} f(x),$$

we have that

$$h_i^{\text{UP}} = \max\{\bar{h}_{2i}^{\text{UP}}, \bar{h}_{2i+1}^{\text{UP}}\},$$

and therefore,

$$h_i^{\text{UP}} \cdot \Delta x \geq \bar{h}_{2i}^{\text{UP}} \cdot \frac{\Delta x}{2} + \bar{h}_{2i+1}^{\text{UP}} \cdot \frac{\Delta x}{2},$$

showing that the sequence, $\text{Area}_n^{\text{UP}}$, is monotonically decreasing. Using Lipschitz continuity, we have the sequence is lower bounded by

$$(f(a) - L|b-a|) \cdot (b-a),$$

showing that all of the conditions of the Monotone Convergence Theorem are met.

The learner can repeat the above argument with small changes and show that $\text{Area}_n^{\text{Low}}$ is monotonically increasing and bounded from above, and hence also satisfies the conditions of the Monotone Convergence Theorem. ■

Proof of Prop. 3.4 (Additivity Property of the Riemann Integral)

Suppose that $a < b < c$ and $f : [a, c] \rightarrow \mathbb{R}$ is continuous. Then the three Riemann Integrals $\int_a^c f(x) dx$, $\int_a^b f(x) dx$, and $\int_b^c f(x) dx$ all exist and satisfy

$$\int_a^c f(x) dx = \int_a^b f(x) dx + \int_b^c f(x) dx.$$

Proof: For many people, Fig. 3.3 is proof enough! If we were using the more general Darboux partitions, instead of the uniform partitions $x_i := a + (i-1) \cdot \Delta x$, then the proof would be immediate as well. Due to our use of simplified partitions, the best way to organize the proof is to wait for the integrability results in Chapter 4.7 on piecewise continuous functions, define

$$\begin{aligned} f_{ab}(x) &:= f(x), \quad a \leq x < b \\ f_{bc}(x) &:= f(x), \quad b \leq x \leq c \end{aligned}$$

and note that $f(x) = f_{ab}(x) + f_{bc}(x)$. Finally, you can then apply Prop. 3.7. ■

Proof of Prop. 3.7 (Linearity of the Riemann Integral)

Suppose that $a < b$ and each of the functions $1 \leq k \leq N$, $f_k : [a, b] \rightarrow \mathbb{R}$, is Riemann integrable (i.e., their Riemann Integral exists and is finite). Then for constants $\alpha_k \in \mathbb{R}$, the function

$$f(x) := \alpha_1 \cdot f_1(x) + \alpha_2 \cdot f_2(x) + \cdots + \alpha_N \cdot f_N(x),$$

defined as a “linear combination of the functions $\{f_1, f_2, \dots, f_N\}$ ”, is Riemann integrable and, moreover,

$$\int_a^b f(x) dx = \alpha_1 \int_a^b f_1(x) dx + \alpha_2 \int_a^b f_2(x) dx + \cdots + \alpha_N \int_a^b f_N(x) dx.$$

Proof: The proof uses the Riemann Lower and Upper Sums and needs a few key facts that are left to the learner, namely,

$$\max_{x \in [x_i, x_{i+1}]} (\alpha_1 \cdot f_1(x) + \cdots + \alpha_N \cdot f_N(x)) \leq \max_{x \in [x_i, x_{i+1}]} (\alpha_1 \cdot f_1(x)) + \cdots + \max_{x \in [x_i, x_{i+1}]} (\alpha_N \cdot f_N(x)) \quad (3.62)$$

$$\min_{x \in [x_i, x_{i+1}]} (\alpha_1 \cdot f_1(x) + \cdots + \alpha_N \cdot f_N(x)) \geq \min_{x \in [x_i, x_{i+1}]} (\alpha_1 \cdot f_1(x)) + \cdots + \min_{x \in [x_i, x_{i+1}]} (\alpha_N \cdot f_N(x));$$

in words, the max of a sum is less than or equal to the sum of the maxes, while the min of a sum is greater than or equal to the sum of the mins.

Given, this, we complete the proof. Without loss of generality, we can assume that the coefficients⁵ $\alpha_k \geq 0$. Indeed, if for some k , we have $\alpha_k < 0$, then the substitutions $f_k(x) \rightarrow -f_k(x)$ and $\alpha_k \rightarrow -\alpha_k$ result in the sum being unchanged and all of the coefficients

⁵If minus signs were involved, then we'd have to deal with $\max f(x) = -\min(-f(x))$, etc., making the argument clumsy. In math, whenever possible, we like to reduce things to a simpler case that still shows the result we seek.

being non-negative. With this condition, it is now true that for all $1 \leq k \leq n$,

$$\alpha_k \cdot \min_{x \in [x_i, x_{i+1}]} f_k(x) = \min_{x \in [x_i, x_{i+1}]} (\alpha_k \cdot f_k(x))$$

and

$$\alpha_k \cdot \max_{x \in [x_i, x_{i+1}]} f_k(x) = \max_{x \in [x_i, x_{i+1}]} (\alpha_k \cdot f_k(x)).$$

When combined with (3.62), we have,

$$\begin{aligned} \max_{x \in [x_i, x_{i+1}]} (\alpha_1 \cdot f_1(x) + \cdots + \alpha_N \cdot f_N(x)) &\leq \alpha_1 \cdot \max_{x \in [x_i, x_{i+1}]} f_1(x) + \cdots + \alpha_N \cdot \max_{x \in [x_i, x_{i+1}]} f_N(x) \\ &= h_i^{\text{Up}} \cdot \Delta x \\ \min_{x \in [x_i, x_{i+1}]} (\alpha_1 \cdot f_1(x) + \cdots + \alpha_N \cdot f_N(x)) &\geq \alpha_1 \cdot \min_{x \in [x_i, x_{i+1}]} f_1(x) + \cdots + \alpha_N \cdot \min_{x \in [x_i, x_{i+1}]} f_N(x). \end{aligned} \tag{3.63}$$

We now follow our general (simplified) procedure for Riemann Lower and Upper Sums from Chapter 3.1.1, let $n \geq 2$, and define

$$\begin{aligned} \Delta x &:= \frac{b-a}{n} \\ x_i &:= a + (i-1) \cdot \Delta x \\ h_i^{\text{Low}} &:= \min_{x \in [x_i, x_{i+1}]} f(x) \\ h_i^{\text{Up}} &:= \max_{x \in [x_i, x_{i+1}]} f(x) \\ h_{i,j}^{\text{Low}} &:= \min_{x \in [x_i, x_{i+1}]} f_j(x) \\ h_{i,j}^{\text{Up}} &:= \max_{x \in [x_i, x_{i+1}]} f_j(x) \\ \text{Area}_n^{\text{Low}} &:= \sum_{i=1}^n h_i^{\text{Low}} \cdot \Delta x \\ \text{Area}_n^{\text{Up}} &:= \sum_{i=1}^n h_i^{\text{Up}} \cdot \Delta x \\ \text{Area}_{n,j}^{\text{Low}} &:= \sum_{i=1}^n h_{i,j}^{\text{Low}} \cdot \Delta x \\ \text{Area}_{n,j}^{\text{Up}} &:= \sum_{i=1}^n h_{i,j}^{\text{Up}} \cdot \Delta x. \end{aligned} \tag{3.64}$$

Then, from (3.63), we have

$$\sum_{j=1}^N \alpha_j \text{Area}_{n,j}^{\text{Low}} \leq \text{Area}_n^{\text{Low}} \leq \text{Area}_n^{\text{Up}} \leq \sum_{j=1}^N \alpha_j \text{Area}_{n,j}^{\text{Up}}. \tag{3.65}$$

Then, because N is finite,

$$\begin{aligned} \lim_{n \rightarrow \infty} \sum_{j=1}^N \alpha_j \text{Area}_{n,j}^{\text{Low}} &= \sum_{j=1}^N \alpha_j \lim_{n \rightarrow \infty} \text{Area}_{n,j}^{\text{Low}} \\ \lim_{n \rightarrow \infty} \sum_{j=1}^N \alpha_j \text{Area}_{n,j}^{\text{Up}} &= \sum_{j=1}^N \alpha_j \lim_{n \rightarrow \infty} \text{Area}_{n,j}^{\text{Up}}. \end{aligned} \tag{3.66}$$

Because $\{f_1, \dots, f_N\}$ are Riemann Integrable, the above limits exist, are finite, and are equal. Hence, by the **Squeeze Theorem for Limits**,

$$\lim_{n \rightarrow \infty} \text{Area}_n^{\text{Low}} = \lim_{n \rightarrow \infty} \text{Area}_n^{\text{Up}} = \sum_{j=1}^N \alpha_j \lim_{n \rightarrow \infty} \text{Area}_{n,j}^{\text{Up}} = \sum_{j=1}^N \alpha_j \int_a^b f_j(x).$$

This proves the result. ■

Proof of Cor. 3.8 (Generalized First Additivity Property)

For a, b , and c finite real numbers in no particular order, define $\alpha := \min\{a, b, c\}$ to be the smallest of the three numbers, and $\beta := \max\{a, b, c\}$ the largest. If $f : [\alpha, \beta] \rightarrow \mathbb{R}$ is continuous, then

$$\int_a^b f(x) = \int_a^c f(x) + \int_c^b f(x). \quad (3.67)$$

Proof:

From Prop. 3.4, if $a < b < c$ and $f : [a, c] \rightarrow \mathbb{R}$ is continuous, then the three Riemann Integrals $\int_a^c f(x)$, $\int_a^b f(x)$, and $\int_b^c f(x)$ all exist and satisfy

$$\int_a^c f(x) = \int_a^b f(x) + \int_b^c f(x).$$

Hence, $\int_a^b f(x) = \int_a^c f(x) - \int_b^c f(x)$. If we apply (3.19) on swapping the order of the limits of the integral, we obtain

$$\int_a^b f(x) = \int_a^c f(x) + \int_c^b f(x),$$

which shows the additivity property holds when c is larger than both limits and not just for c between the two limits.

Similarly,

$$\int_b^c f(x) = - \int_a^b f(x) + \int_a^c f(x),$$

by the First (Basic) Additivity Property, and, if we apply (3.19), we obtain

$$\int_b^c f(x) = \int_b^a f(x) + \int_a^c f(x),$$

which shows the additivity property holds when c is smaller than both limits. This is half⁶ of what we need to show.

If we negate the First (Basic) Additivity Property, we obtain,

$$-\int_a^c f(x) = -\int_a^b f(x) - \int_b^c f(x),$$

which, upon flipping all three sets of limits of integration, yields

$$\int_c^a f(x) = \int_b^a f(x) + \int_c^b f(x).$$

From this, the remaining cases follow easily, and hence, the order of a, b , and c is immaterial. ■

Proof of Prop. 3.10 (Shift Property)

Suppose that $a < b$ and $f : [a, b] \rightarrow \mathbb{R}$ is continuous and let $x_c \in \mathbb{R}$. Then $g : [a+x_c, b+x_c] \rightarrow \mathbb{R}$ by $g(x) := f(x-x_c)$ is also continuous and $\int_a^b f(x) dx = \int_{a+x_c}^{b+x_c} g(x) dx$.

A graphical proof is provided by the [Khan Academy](#) and in Fig. 3.4.

Proof: If you want a formal proof with Riemann Sums, we can develop that too! For simplicity and clarity, we assume the function is monotonically increasing and we use the lower Riemann sum:

- $\Delta x := \frac{b-a}{n}$
- $x_i = a + (i-1)\Delta x$

⁶The integrals on the left of our three cases all have a lower limit that is smaller than the upper limit. What happens when the reverse is true?

- $\text{Area}_n^{\text{Low}} := \sum_{i=1}^n f(x_i) \cdot \Delta x$ is the lower Riemann sum.

For $g : [a+x_c, b+x_c] \rightarrow \mathbb{R}$ by $g(x) := f(x - x_c)$, we construct its lower Riemann sum, via

- $\Delta \bar{x} := \frac{(b+x_c) - (a+x_c)}{n} = \frac{b-a}{n} = \Delta x$
- $\bar{x}_i = a + x_c + (i-1)\Delta x = x_i + x_c$
- $\overline{\text{Area}}_n^{\text{Low}} := \sum_{i=1}^n g(\bar{x}_i) \cdot \Delta \bar{x} = \sum_{i=1}^n g(x_i + x_c) \cdot \Delta x$ is the lower Riemann sum for g .

Because $g(x_i + x_c) := f(x_i + x_c - x_c) = f(x_i)$, the two lower Riemann sums are identical, and hence, if the integrals exist, they are also identical. ■

Proof of Prop. 3.12 (Scaling Property)

Suppose that $a < b$, $f : [a, b] \rightarrow \mathbb{R}$ is continuous, and $\omega_0 > 0$ is a real number. Then $g : [\frac{a}{\omega_0}, \frac{b}{\omega_0}] \rightarrow \mathbb{R}$ by $g(x) := f(\omega_0 \cdot x)$ is also continuous and $\int_a^b f(x) dx = \omega_0 \cdot \int_{\frac{a}{\omega_0}}^{\frac{b}{\omega_0}} g(x) dx$.

Proof:

We reuse the proof strategy for the Shifting Property. For $g : [\frac{a}{\omega_0}, \frac{b}{\omega_0}] \rightarrow \mathbb{R}$ by $g(x) := f(\omega_0 \cdot x)$, we construct its lower Riemann sum, via

- $\Delta \bar{x} := \frac{\frac{b}{\omega_0} - \frac{a}{\omega_0}}{n} = \frac{1}{\omega_0} \cdot \frac{b-a}{n} = \frac{1}{\omega_0} \cdot \Delta x$
- $\bar{x}_i = \frac{a}{\omega_0} + (i-1)\Delta \bar{x} = \frac{a}{\omega_0} + (i-1)\frac{1}{\omega_0} \cdot \Delta x = \frac{x_i}{\omega_0}$
- $\overline{\text{Area}}_n^{\text{Low}} := \sum_{i=1}^n g(\bar{x}_i) \cdot \Delta \bar{x} = \sum_{i=1}^n g(\frac{x_i}{\omega_0}) \cdot \frac{1}{\omega_0} \cdot \Delta x$ is the lower Riemann sum for g .

Because $g(\frac{x_i}{\omega_0}) := f(\omega_0 \cdot \frac{x_i}{\omega_0}) = f(x_i)$, the two lower Riemann sums are identical, up to the scale factor of $\frac{1}{\omega_0}$, and hence, if the integrals exist, they are also identical up to the scale factor of $\frac{1}{\omega_0}$. ■

For the case $\omega_0 < 0$, the same proof is valid and is not repeated. ■

Proof of Prop. 3.14 (Trapezoidal Rule for Computing a Riemann Integral)

Suppose that $a < b$ and $f : [a, b] \rightarrow \mathbb{R}$ is continuous. Then the Riemann Integral $\int_a^b f(x) dx$ exists and it can be evaluated by the **Trapezoidal Rule**. That is,

$$\int_a^b f(x) dx = \lim_{n \rightarrow \infty} \text{Area}_n^{\text{Trap}} = \lim_{n \rightarrow \infty} \sum_{i=1}^n \frac{f(x_i) + f(x_{i+1})}{2} \cdot (x_{i+1} - x_i),$$

where, $\Delta x := \frac{b-a}{n}$ and $x_i := a + (i-1)\Delta x$. In the above, the limits mean that for all $\epsilon > 0$, there exists $N < \infty$ such that $n \geq N \implies \left| \int_a^b f(x) dx - \sum_{i=1}^n \frac{f(x_i) + f(x_{i+1})}{2} \cdot (x_{i+1} - x_i) \right| < \epsilon$.

Proof: We give the proof for functions that are monotonically increasing. For the general case, one must appeal to the uniform continuity property of continuous functions over closed bounded intervals.

If the integral $\int_a^b f(x) dx$ exists, then the lower and upper Riemann sums for f both exist and converge to the same number, which defines $\int_a^b f(x) dx$. By Prop. 2.34, the average of the lower and upper Riemann sums must also converge to $\int_a^b f(x) dx$. For monotonic functions, the trapezoidal rule is precisely the average of the upper and lower Riemann sums. Hence, the result is proved. ■.

Proof of Prop. 3.23 (Integrals of Even and Odd Functions over Symmetric Intervals)

Suppose that $f : \mathbb{R} \rightarrow \mathbb{R}$ is a function and that for $a > 0$, the definite integrals $\int_{-a}^0 f(x) dx$ and $\int_0^a f(x) dx$ both exist and are finite. Then the following hold:

(a) If f is an even function,

- $\int_{-a}^0 f(x) dx = \int_0^a f(x) dx$, and
- $\int_{-a}^a f(x) dx = 2 \int_0^a f(x) dx$.

(b) If f is an odd function,

- $\int_{-a}^0 f(x) dx = -\int_0^a f(x) dx$, and
- $\int_{-a}^a f(x) dx = 0$.

(c) For all $[b, c] \subset [-a, a]$,

- $\int_b^c f(x) dx = \int_b^c f_e(x) dx + \int_b^c f_o(x) dx$.

Proof:

- (a) If f is an even function, $\int_{-a}^0 f(x) dx = \int_0^a f(x) dx$ follows from Prop. 3.12 by taking $\omega_0 = -1$. The second statement is immediate from the first statement.
- (b) If f is an odd function, $\int_{-a}^0 f(x) dx = -\int_0^a f(x) dx$ follows from Prop. 3.12 by taking $\omega_0 = -1$. The second statement is immediate from the first statement.
- (c) For all $[b, c] \subset [-a, a]$, $\int_b^c f(x) dx = \int_b^c f_e(x) dx + \int_b^c f_o(x) dx$ follows from Prop. 3.7 because $f(x) = f_e(x) + f_o(x)$ by the construction of the odd and even parts.

Chapter 4

Properties of Functions: Left and Right Limits, Types of Continuity, Boundedness, and Generalizations of Max and Min

Remark: This Chapter reveals some of the mathematical backbone of Calculus. Traditionally, much of this material is placed earlier in a Calculus course. Because many of the topics are abstract and very technical, we delayed them until you've gotten the hang of combining programming and numerical calculations when learning mathematical concepts.

Learning Objectives

By the end of this chapter, the student should be able to:

- Appreciate the foundational role of calculus in mathematical modeling and problem-solving, integrating programming and numerical methods to solidify these concepts.
- Analyze the behavior of functions at specific points using the concept of one-sided limits.
- Thoroughly understand the nature of function continuity and discover that continuity comes in more than one flavor.
- Derive closed-form expressions for the integrals of exponential functions and apply these techniques to integrate trigonometric functions.
- Determine the boundedness of functions and understand the implications for mathematical analysis and applications.
- Recognize that while maximum and minimum values are important, they are not the full story.

Outcomes

Upon successful completion of this chapter, students will be able to:

- Acquire both intuitive and formal understandings of one-sided limits and their calculation.
- Apply numerical methods to estimate one-sided limits and assess the continuity of functions using the epsilon-delta definition.
- Recognize when it is permissible to take limits within functions and apply this knowledge to the integration of exponentials.
- Identify and analyze piecewise continuous functions and extend this understanding to a broader class of Riemann integrable functions.
- Review and apply the concepts of maximum and minimum function values, and explore alternative behaviors when these extrema do not exist.

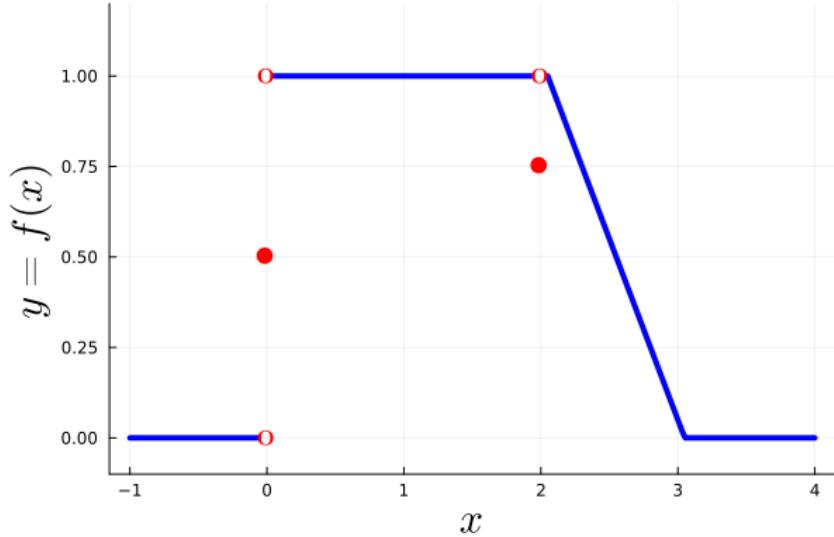


Figure 4.1: This oddly defined function will help us understand the notions of limit from the left and limit from the right at a point. The value of the function at $x = 0$ is 0.5, and its value at $x = 2$ is 0.75.

4.1 Intuition for Limits from the Left and the Right

To help with developing an intuitive feel for one-sided limits, consider the function

$$f(x) := \begin{cases} 0.0 & x < 0.0 \\ 0.5 & x = 0.0 \\ 1.0 & 0.0 < x < 2.0 \\ 0.75 & x = 2.0 \\ 1.0 - (x - 2.0) & 2.0 < x \leq 3.0 \\ 0.0 & x > 3.0 \end{cases} \quad (4.1)$$

as plotted in Fig. 4.1. We have defined the value of the function at the origin to be neither zero nor one; in fact, $f(0) = 0.5$. However, if we approach the origin from the left side, namely through negative values of x , the value of the function is constant and equal to zero. Hence, we say that the limit from the left at $x = 0.0$ equals zero, and we denote this by

$$\lim_{x \rightarrow 0^-} f(x) = 0. \quad (4.2)$$

Here, the minus sign as a superscript on the right side of a value of x , currently $x_0 = 0$, indicates the limit is taken with values of x that are strictly less than $x_0 = 0$, are never equal to x_0 , but become “arbitrarily close to x_0 ”. Similarly, if we approach the origin from the right side, namely through small positive values of x , the value of the function is constant and equal to one. Hence, we say that the limit from the right at $x_0 = 0.0$ equals one, and we denote this by

$$\lim_{x \rightarrow 0^+} f(x) = 1. \quad (4.3)$$

Here, the plus sign as a superscript on the right side of a value of x , say x_0 , indicates the limit is taken with values of x that are strictly greater than $x_0 = 0$, are never equal to x_0 , but become “arbitrarily close to x_0 ”. Before we move on, let’s note that the limit of the function from the left at the origin equals zero, its limit from the right equals one, and the value of the function at the origin is 0.5. The function is clearly discontinuous at the origin because there is a jump.

Let’s next consider the behavior of the function near $x_0 = 2$. If we approach $x_0 = 2$ from its left side, namely through values of x strictly less than two, the value of the function is constant and equal to one. Hence, we say that the limit from the left at $x = 2.0$ equals one, and we denote this by

$$\lim_{x \rightarrow 2^-} f(x) = 1. \quad (4.4)$$

Again, the minus sign as a superscript on the right side of a value of $x_0 = 2$ indicates the limit is taken with values of x that are strictly less than $x_0 = 2$, are never equal to x_0 , but become arbitrarily close to x_0 . Similarly, if we approach $x_0 = 2$ from the right side, namely through values of x strictly greater than two, the value of the function is less than one but becomes arbitrarily close to one as x approaches two. Hence, we say that the limit from the right at $x = 2$ equals one, and we denote this by

$$\lim_{x \rightarrow 2^+} f(x) = 1. \quad (4.5)$$

Again, the plus sign as a superscript on the right side of $x_0 = 2$ indicates the limit is taken with values of x that are strictly greater than x_0 , are never equal to x_0 , but become arbitrarily close to x_0 . Before we move on, let's note that the limit of the function from the left at $x_0 = 2$ equals one, its limit from the right equals one, and the value of the function at $x_0 = 2$ is equal to 0.75. The function is discontinuous at $x_0 = 2$ because its value jumps from 1.0 to 0.75 at the point $x_0 = 2$.

For our final case, let's consider the behavior of the function near $x_0 = 3$. If we approach $x_0 = 3$ from its left side, namely through values of x strictly less than three, the value of the function decreases to zero. Hence, we say that the limit from the left at $x = 3$ equals zero, and we denote this by

$$\lim_{x \rightarrow 3^-} f(x) = 0. \quad (4.6)$$

As before, the minus sign as a superscript on the right side of a value of $x_0 = 3$ indicates the limit is taken with values of x that are strictly less than $x_0 = 3$, are never equal to x_0 , but become arbitrarily close to x_0 . Similarly, if we approach $x_0 = 3$ from the right side, namely through values of x strictly greater than three, the value of the function is constant and equal to zero. Hence, we say that the limit from the right at $x = 3$ equals zero, and we denote this by

$$\lim_{x \rightarrow 3^+} f(x) = 0. \quad (4.7)$$

Again, the plus sign as a superscript on the right side of $x_0 = 3$ indicates the limit is taken with values of x that are strictly greater than $x_0 = 3$, are never equal to x_0 , but become arbitrarily close to x_0 . Before we move on, let's note that the limit of the function from the left at $x_0 = 3$ equals zero, its limit from the right equals zero, and the value of the function at $x_0 = 3$ is equal to zero. The function is continuous at $x_0 = 3$ because its value does not undergo a jump as x is varied near $x_0 = 3$. This is "mathematically certified" by the limit from the left, the limit from the right, and the function all having the same value at the point $x_0 = 3$. In fact, this property, namely, the two one-sided (left and right) limits existing and agreeing with the value of the function, holds for all $-1 < x \leq 0$, $0 < x < 2$, and $2 < x < 4$; in other symbols, for all $x_0 \in \{(-1, 4) \mid x_0 \neq 0 \text{ and } x_0 \neq 2\}$,

$$\lim_{x \rightarrow x_0^-} f(x) = f(x_0) \text{ and } \lim_{x \rightarrow x_0^+} f(x) = f(x_0),$$

and this analysis agrees with our intuitive notion of the continuity of the function at those points.

Fact 4.1: Relation to Limits at $\pm\infty$

We connect the one-sided limits to the work we did in Chapters 2.5 and 2.6. The limit from the right is the same as the following (positive) infinite limit

$$\lim_{x \rightarrow x_0^+} f(x) = \lim_{\eta \rightarrow +\infty} f(x_0 + \frac{1}{\eta})$$

and the limit from the left is the same as the corresponding negative infinite limit

$$\lim_{x \rightarrow x_0^-} f(x) = \lim_{\eta \rightarrow -\infty} f(x_0 + \frac{1}{\eta}).$$

The key points are:

- $\frac{1}{\eta}$ is never exactly equal to zero, so we never evaluate the function at x_0 ;
- as $\eta \rightarrow \infty$, $x_0 + \frac{1}{\eta}$ approaches x_0 from the **right** because $x_0 + \frac{1}{\eta} > x_0$ for all $\eta > 0$; and
- as $\eta \rightarrow -\infty$, $x_0 + \frac{1}{\eta}$ approaches x_0 from the **left** because $x_0 + \frac{1}{\eta} < x_0$ for all $\eta < 0$.

Note: If we make the substitution $h := \frac{1}{\eta}$, then $\eta \rightarrow +\infty \iff h \rightarrow 0^+$ (h decreases to zero through positive values, that is, h approaches zero from the right) and $\eta \rightarrow -\infty \iff h \rightarrow 0^-$ (because h increases to zero through negative values, that is, h approaches zero from the left).

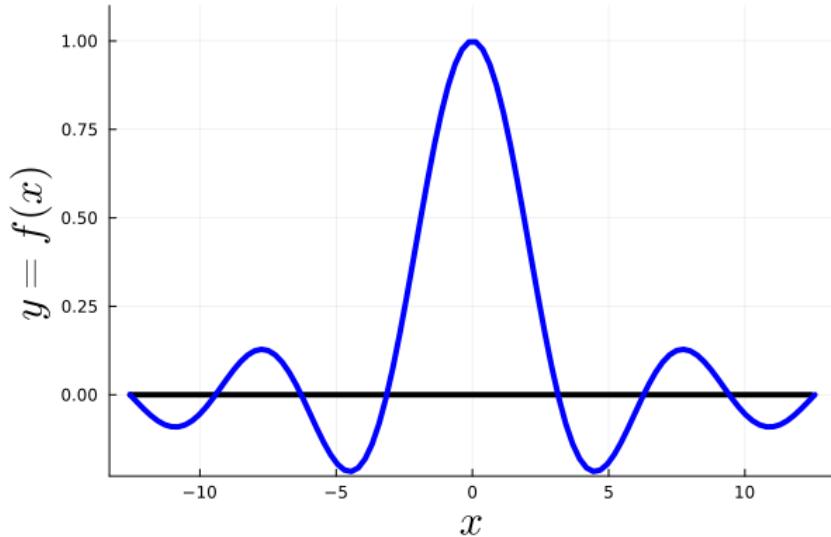


Figure 4.2: This is a plot of $\sin(x)/x$, a function that plays a starring role in Signal Processing. How is it possible to make sense of its value at the origin, where we have $0/0$? Inquiring minds want to know!

Secrets of the Arcane 4.2: How Might Limits Be Useful

The function $\text{sinc}(x) := \begin{cases} \frac{\sin(x)}{x} & x \neq 0 \\ 1 & x = 0 \end{cases}$ (pronounced “sink” or “sync”) plays a big role in Signal Processing because it can be used to create^a low-pass filters (i.e., filters that remove high-frequency “noise” from a data signal). Why is it defined to be one at the origin? **Limits give us a way to understand and analyze such conundrums as $0/0$.** The function’s graph is given in Fig. 4.2.

This Chapter puts us on the road to understanding how to deal with “ $0/0$,” while L’Hôpital’s Rule in the next Chapter should clear up any remaining doubts.

^aMichigan’s EECS 216 Signals and Systems covers this material. In particular, it covers the Fourier Transform, which is the key to unlocking how a sinc-function has anything to do with filtering noise from signals.

4.2 Formal Definition of Limits from the Left and the Right

We build off the observations presented in Fact 4.1.

One-sided Limits at Zero

Definition 4.3. Suppose that $f : A \rightarrow \mathbb{R}$ is a real-valued function, $A \subset \mathbb{R}$, and both L and R are (finite) real numbers.

- (a) The **limit of f from the left at 0** is equal to L if, for all $\epsilon > 0$ (no matter how small), there exists $\delta > 0$, such that the open interval $(-\delta, 0) \subset A$ (i.e., is contained in the domain of definition of the function), and $-\delta < h < 0 \implies |f(h) - L| \leq \epsilon$. The left limit is denoted as

$$\lim_{h \rightarrow 0^-} f(h) = L.$$

- (b) The **limit of f from the right at 0** is equal to R if, for all $\epsilon > 0$ (no matter how small), there exists $\delta > 0$, such that the open interval $(0, \delta) \subset A$ (i.e., is contained in the domain of definition of the function), and $0 < h < \delta \implies |f(h) - R| \leq \epsilon$. The right limit is denoted as

$$\lim_{h \rightarrow 0^+} f(h) = R.$$

Notes:

- In the above, h is a dummy variable. We can just as easily write $\lim_{x \rightarrow 0^-} f(x)$, $\lim_{y \rightarrow 0^+} f(y)$, etc.
- As we take ϵ smaller and smaller, we typically need to take δ smaller and smaller. In other words, δ typically depends on ϵ . It is sometimes useful to express this dependence as $\delta(\epsilon)$ or δ_ϵ , just as we did with $N(\epsilon)$ and N_ϵ when doing limits at infinity.
- Similar to limits as x tends to infinity, most one-sided limits in Engineering can be computed by inspection, as in Chapter 2.6, or with minor computations. This will become more clear when we get to L'Hôpital's Rule. For now, focus on the concept of how the function behaves as we approach the origin from the right or the left.**

Example 4.4. Compute the left and right limits at zero for the following functions.

$$(a) f(y) = y^2, y \in \mathbb{R}$$

$$(b) f(z) = \begin{cases} -1 + z^4 & z < 0 \\ 0 & z = 0 \\ 1 + z^2 & z > 0 \end{cases}$$

$$(c) f(h) = \begin{cases} \frac{h}{|h|} & h \neq 0 \\ \text{undefined} & h = 0 \end{cases}$$

$$(d) f(x) = \begin{cases} \sin(\frac{1}{x}) & x > 0 \\ 0 & x \leq 0 \end{cases}$$

Solutions: We first give solutions that use the epsilon-delta-definitions, and then provide a pointer for doing the limits by inspection.

$$(a) \text{ Ans. } \lim_{y \rightarrow 0^-} f(y) = 0 \text{ and } \lim_{y \rightarrow 0^+} f(y) = 0.$$

Let $\epsilon > 0$ be given and set $\delta := \sqrt{\epsilon}$. Then for all $-\delta < y < 0$, we have $|f(y) - L| = |y^2 - 0| = |y^2| = y^2 < (-\delta)^2 = \epsilon$, and hence the definition of the limit from the left is verified.

Similarly, for all $0 < y < \delta$, we have $|f(y) - R| = |f(y) - 0| = |y^2| = y^2 < (\delta)^2 = \epsilon$, and hence the definition of the limit from the right is verified. In this case, we could use the same δ for both sides of the limit.

$$(b) \text{ Ans. } \lim_{z \rightarrow 0^-} f(z) = -1 \text{ and } \lim_{z \rightarrow 0^+} f(z) = 1.$$

Let $\epsilon > 0$ be given and set $\delta := \sqrt[4]{\epsilon}$. Then for all $-\delta < z < 0$, we have $|f(z) - L| = |-1 + z^4 - (-1)| = |z^4| = z^4 < (-\delta)^4 = \epsilon$, and hence the definition of the limit from the left is verified.

Let $\epsilon > 0$ be given and set $\delta := \sqrt{\epsilon}$ for all $0 < z < \delta$, we have $|f(z) - R| = |1 + z^2 - 1| = |z^2| = z^2 < (\delta)^2 = \epsilon$, and hence the definition of the limit from the right is verified.

$$(c) \text{ Ans. } \lim_{h \rightarrow 0^-} f(h) = -1 \text{ and } \lim_{h \rightarrow 0^+} f(h) = 1.$$

Let $\epsilon > 0$ be given and set $\delta := 27$. Then for all $-\delta < h < 0$, we have $|f(h) - L| = |\frac{h}{|h|} - (-1)| = |-1 - (-1)| = 0 < \epsilon$, and hence the definition of the limit from the left is verified. In this case, we can take δ to be any positive real number. It did not have to be small. You will understand why when you see the graph of $f(h)$ in Fig. 4.3-(c).

Similarly, for all $0 < h < \delta$, we have $|f(h) - R| = |\frac{h}{|h|} - 1| = |1 - 1| = 0 < \epsilon$, and hence the definition of the limit from the right is verified.

(d) **Ans.** $\lim_{x \rightarrow 0^-} f(x) = 0$ and $\lim_{x \rightarrow 0^+} f(x)$ does not exist.

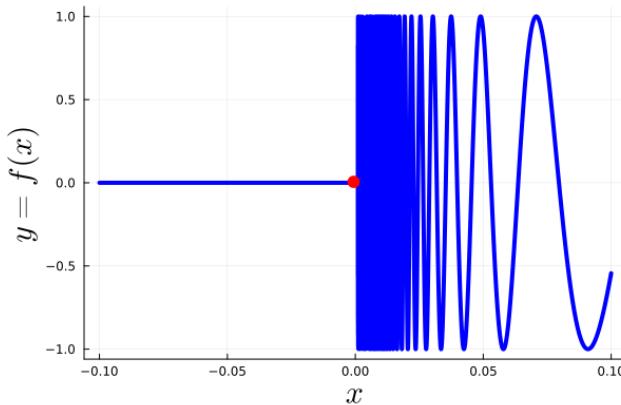
Here is the insight: $\lim_{x \rightarrow 0^+} \sin\left(\frac{1}{x}\right) = \lim_{y \rightarrow +\infty} \sin(y)$, and we know that $\sin(y)$ oscillates forever between -1 and $+1$ as $y \rightarrow \infty$; hence, the limit from the right at zero does not exist. Can we show this violates the existence of the limit, using epsilons and deltas? Of course!

Let $R \in \mathbb{R}$, the candidate right limit, be arbitrary. Set $\epsilon = 0.25$ and let $\delta > 0$ be arbitrary. We seek to show that there always exists $0 < x < \delta$ such that $|f(x) - R| > \epsilon$.

Let k be an even integer such that $0 < \frac{1}{k} < \delta$. Then $x_1 := \frac{1}{k\pi}$ and $x_2 := \frac{1}{(k+0.5)\pi}$ satisfy $0 < x_2 < x_1 < \delta$. These numbers were chosen because

$$f(x_1) = \sin(k\pi) = 0 \text{ and } f(x_2) = \sin((k+0.5)\pi) = \sin\left(\frac{\pi}{2}\right) = 1.$$

If $|f(x_1) - R| = |0 - R| = |R| > \epsilon$, then R is not the right limit. Hence, suppose that $|R| \leq \epsilon = 0.25$. Then, by the Reverse Triangle Inequality, $|f(x_2) - R| \geq |f(x_2) - |R|| \geq 1 - 0.25 = 0.75 > \epsilon$, and thus R is not the right limit. Because R was arbitrary, the right limit does not exist at zero. A zoom of the function's graph near the origin is shown below.



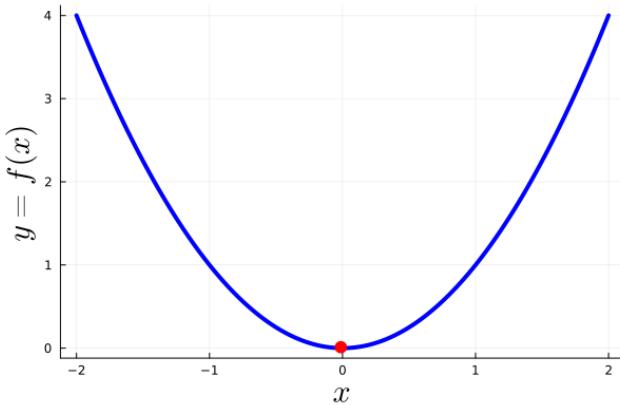
Because the limit from the left is “obvious”, this completes the analytical solutions. **The solutions by inspection** proceed as in Chapter 4.1, using Fig. 4.3. You can trace the graph of each function to the left and right of zero to determine the limits.

You can also estimate the limits numerically like this: define $h[i] = \frac{1}{i}$, $0 < i_{\min} \leq i \leq i_{\max}$, for example, and then compute $f(h)$ and $f(-h)$ on the vector of values specified by h . This is illustrated in the next code block.

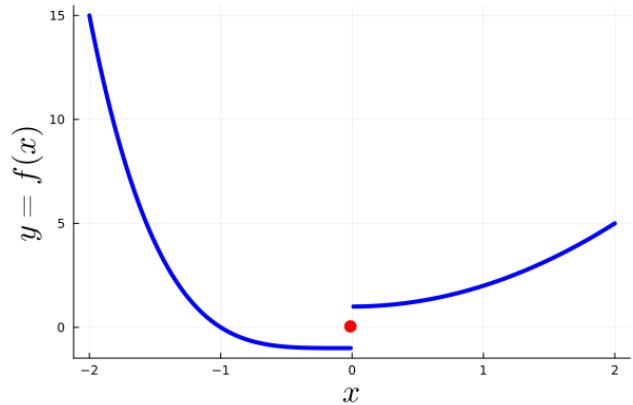
```

1 eta = 100:1:1e4
2 h_Right = 1.0 ./eta
3 k = 2
4 y_Right = f.(h_Right, k)
5 p1 = plot(h_Right, y_Right, linewidth=4, color=:blue, label=false, guidefont=20)
6 plot!( xlabel=L"\$x\$", ylabel=L"\$y=f(x)\$")
7 h_Left = -h_Right
8 y_Left = f.(h_Left, k)
9 p1 = plot!(h_Left, y_Left, linewidth=4, color=:green, label=false, guidefont=20)
10 if k != 3
11     annotate!(0, f(0, k), text(L"\$\bullet\$", :red, :center, 20))
12 end
13 ylims!(p1, (minimum(y_minus), maximum(y_Right)))

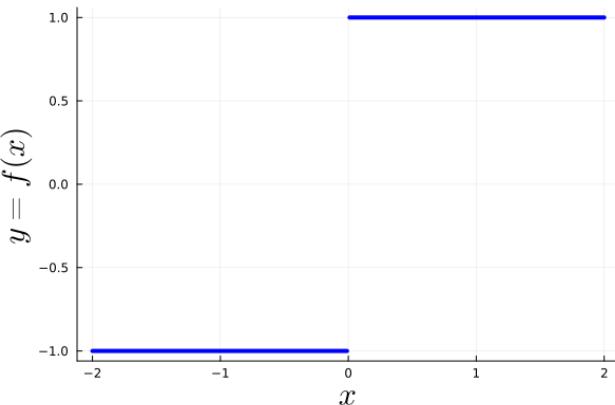
```



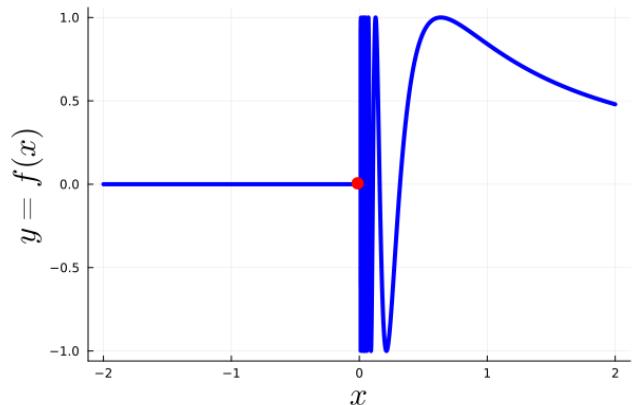
(a)



(b)



(c)



(d)

Figure 4.3: Graphs of the functions defined in Example 4.4, which treats left and right limits at zero. (a) The left and right limits are both equal to zero. (b) The left limit equals -1 and the right limit equals $+1$. (c) The left limit equals -1 and the right limit equals $+1$. (d) The left limit equals 0 and the right limit is undefined because, from the right side of zero, the function oscillates between -1 and $+1$ an infinite number of times for $0 < x < \delta$, for any $\delta > 0$, no matter how small.

```

14
15 display(p1)
16
17 @show LeftLimit = y_Left[end]
18 @show RightLimit = y_Right[end]
19
20 png(p1, "LeftRightLimitAtZeroBSecondMethod")

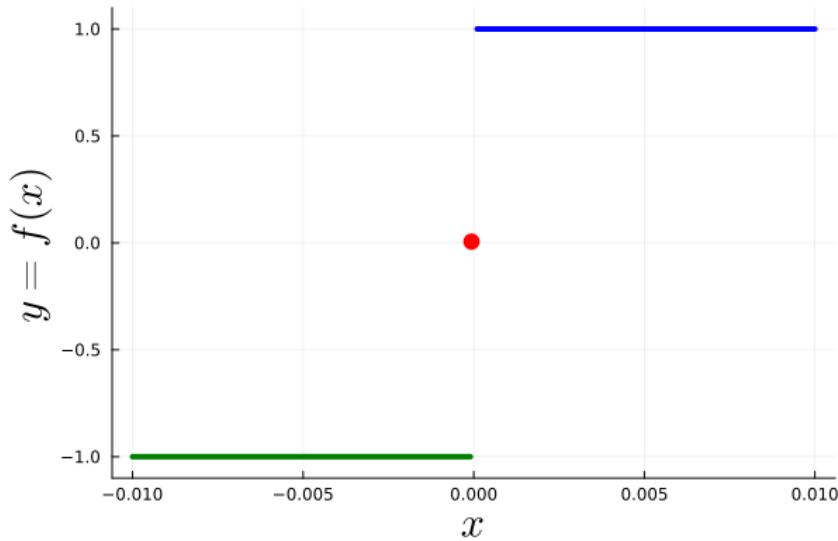
```

Output

```

LeftLimit = y_Left[end] = -0.9999999999999999
RightLimit = y_Right[end] = 1.00000001

```



If the function diverges to infinity as we approach zero from the left or right, then finite limits will certainly not exist. We take care of this, just as we did for limits at infinity in Definition 2.15.

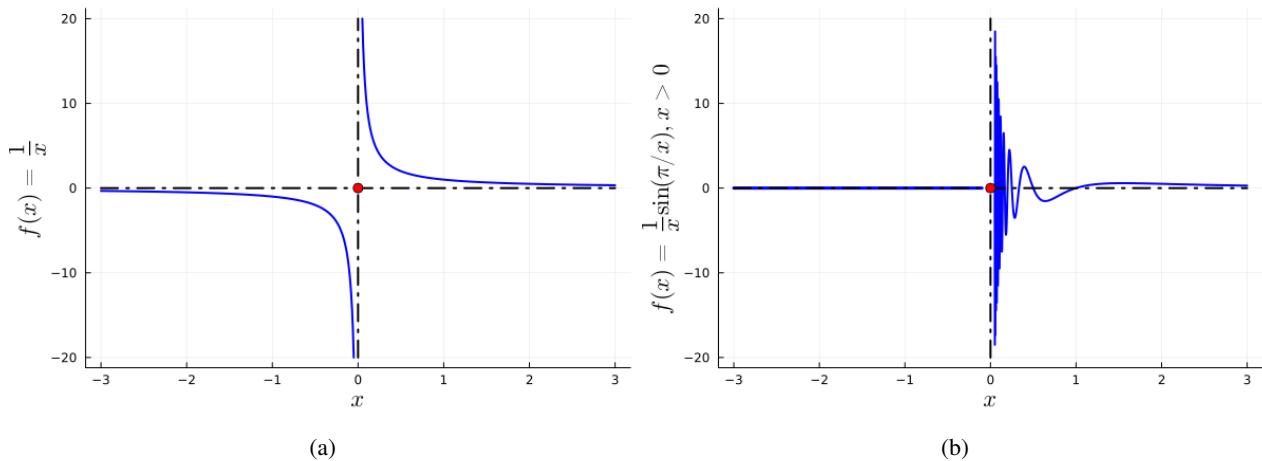


Figure 4.4: Functions can also blow up as the origin is approached from the left or the right. (a) Both one-sided limits of $1/x$ exist but are $\pm\infty$. (b) The limit from the left exists and is bounded, while the limit from the right does not exist for $(1/x) \cdot \sin(\pi/x)$, $x > 0$.

Unbounded One-sided Limits at 0

Definition 4.5. Suppose that $A \subset \mathbb{R}$ and $f : A \rightarrow \mathbb{R}$ is a real-valued function.

- (a) The **limit of f from the left at 0** is equal to ∞ if, for all $0 < N < \infty$ (no matter how large), there exists $\delta > 0$ such that the open interval $(-\delta, 0) \subset A$ (i.e., is contained in the domain of definition of the function), and $-\delta < h < 0 \implies f(h) > N$. The left limit is then denoted as

$$\lim_{h \rightarrow 0^-} f(h) = \infty.$$

Similarly, the **limit of f from the left at 0** is equal to $-\infty$ if, for all $0 < N < \infty$ (no matter how large), there exists $\delta > 0$ such that the open interval $(0, \delta) \subset A$ (i.e., is contained in the domain of definition of the function), and $-\delta < h < 0 \implies f(h) < -N$. The left limit is then denoted as

$$\lim_{h \rightarrow 0^-} f(h) = -\infty.$$

- (b) The **limit of f from the right at 0** is equal to ∞ if, for all $0 < N < \infty$ (no matter how large), there exists $\delta > 0$ such that the open interval $(0, \delta) \subset A$ (i.e., is contained in the domain of definition of the function), **and** $0 < h < \delta \implies f(h) > N$. The right limit is then denoted as

$$\lim_{h \rightarrow 0^+} f(h) = \infty.$$

Similarly, the **limit of f from the right at 0** is equal to $-\infty$ if, for all $0 < N < \infty$ (no matter how large), there exists $\delta > 0$ such that the open interval $(0, \delta) \subset A$ (i.e., is contained in the domain of definition of the function), **and** $0 < h < \delta \implies f(h) < -N$. The right limit is then denoted as

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

Example 4.6. For $f : \mathbb{R} \rightarrow \mathbb{R}$ by

$$f(x) = \begin{cases} \frac{1}{x} & x \neq 0 \\ 0 & x = 0, \end{cases}$$

compute the left and right limits at the origin.

Solution: We provide first an analytical solution and then a numerical solution.

Ans. $\lim_{h \rightarrow 0^-} f(h) = -\infty$. To prove this using the definition, let $0 < N < \infty$ be arbitrary. Then for $0 < \delta < \frac{1}{N}$,

$$-\delta < x < 0 \implies f(x) = \frac{1}{x} \leq -\frac{1}{\delta} \leq -\frac{1}{\frac{1}{N}} = -N.$$

Ans. $\lim_{h \rightarrow 0^+} f(h) = \infty$. To prove this using the definition, let $0 < N < \infty$ be arbitrary. Then for $0 < \delta < \frac{1}{N}$,

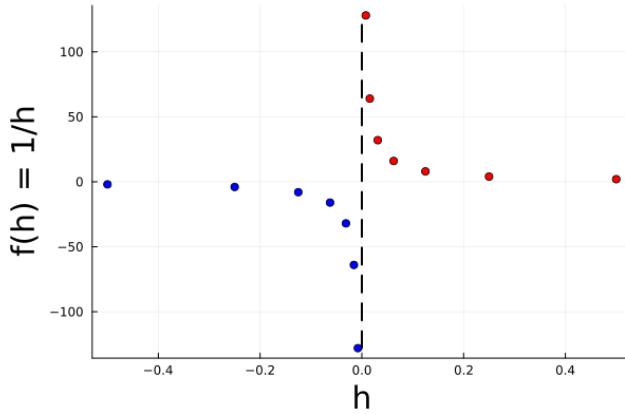
$$0 < x < \delta \implies f(x) = \frac{1}{x} \geq \frac{1}{\delta} \geq \frac{1}{\frac{1}{N}} = N.$$

1 using Plots, LaTeXStrings

```

2
3 # Create a logarithmic scale of h values
4 k = 7 #
5 h = 2.0 .^ range(-1, -k, step=-1)
6
7 f(x) = 1/x
8
9 # Limit from the right
10 p1 = scatter(h, f.(h), legend=false, markersize = 4, color=:red, guidefont=20)
11 plot!( xlabel = "h", ylabel = "f(h) = 1/h")
12 # Limit from the left
13 scatter!(-h, f.(-h), markersize = 4, color=:blue)
```

Output



One-sided Limits at a Point

Definition 4.7. Suppose that $f : A \rightarrow \mathbb{R}$ is a real-valued function, $A \subset \mathbb{R}$, $x_0 \in \mathbb{R}$, and there exists $\delta > 0$ such that the two open intervals $(x_0 - \delta, x_0)$ and $(x_0, x_0 + \delta)$ are contained in A , the domain of definition of the function.

(a) The **limit from the left at x_0** is

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

(b) The **limit from the right at x_0** is

$$\lim_{x \rightarrow x_0^+} f(x) := \lim_{h \rightarrow 0^+} f(x_0 + h).$$

Notes:

- The above is valid for finite limits and unbounded limits.
- $(x_0 - \delta, x_0) \subset A$ and $(x_0, x_0 + \delta) \subset A$ ensure that considering the limits from the left and right makes sense.

Example 4.8. Let's reuse the function that we studied in (4.1) and Fig. 4.1, namely,

$$f(x) := \begin{cases} 0.0 & x < 0.0 \\ 0.5 & x = 0.0 \\ 1.0 & 0.0 < x < 2.0 \\ 0.75 & x = 2.0 \\ 1.0 - (x - 2.0) & 2.0 < x \leq 3.0 \\ 0.0 & x > 3.0. \end{cases}$$

Determine the limits from the left and right for

- (a) $x_0 = 0$
- (b) $x_0 = 2$, and
- (c) $x_0 = 3$.

Solution: We will take a numerical approach. We could plot the function near the specified values of x_0 . We'll do something a bit more quantitative. We'll define a (finite) vector of h values that terminates in very small values, e.g., on the order of 10^{-6} . We'll also compute the **mean** and (sample) **standard deviation** of the function evaluated $x_0 \pm h$, namely

- $h[i] = \frac{1}{i+i_{\min}}$, $1 \leq i \leq N$, $i_{\min} > 0$ a constant.
- $y_R[i] := f(x_0 + h[i])$

- $\mu_R := \frac{1}{N} \cdot \sum_{i=1}^N y_R[i]$ (mean value)
- $\sigma_R := \sqrt{\frac{1}{N-1} \cdot \sum_{i=1}^N (y_R[i] - \mu_R)^2}$ (sample standard deviation)

and similarly for the left-sided limit. If the function has a right limit at x_0 , then y_R should be nearly a constant for N large. Hence, when

- $\mu_R \approx y_R[N]$ and
- $\sigma_R \approx 0$,

the numerical evidence suggests that $y_R[N]$ is a good estimate of the limit.

```

1  function f(x)
2    if x < 0
3      y = -1.0
4    elseif abs(x)<1e-12
5      y = 0.0
6    elseif (0 < x < 2)
7      y = 1.0
8    elseif abs(x-2.0)<1e-12
9      y = 0.75
10   elseif (2 < x <= 3.0)
11     y = 1.0 - (x - 2.0)
12   else
13     y = 0.0
14   end
15   return y
16 end
17 #
18 using Statistics
19 eta = 1e3:1e3:1e6
20 h_Right = (1.0). / eta
21 h_Left = -h_Right
22 x0Vec = [0.0 2.0 3.0]
23
24 delta = h_Right [end]
25
26 @show delta
27 println (" ")
28
29 for k = 1:length(x0Vec)
30   x0 = x0Vec [k]
31   y_Left = f. (x0 .+ h_Left)
32   y_Right = f. (x0 .+ h_Right)
33   #
34   @show x0
35   @show limitLeft = y_Left [end]
36   @show meanLeft = mean (y_Left)
37   @show stdDevLeft = std (y_Left)
38   println ("-----")
39   @show limitRight = y_Right [end]
40   @show meanRight = mean (y_Right)
41   @show stdDevRight = std (y_Right)
42   println (" ")
43   println (" ")
44   #

```

45 **end**

Output

```

delta = 1.0e-6

x0 = 0.0
limitLeft = y_Left[end] = -1.0
meanLeft = mean(y_Left) = -1.0
stdDevLeft = std(y_Left) = 0.0
-----
limitRight = y_Right[end] = 1.0
meanRight = mean(y_Right) = 1.0
stdDevRight = std(y_Right) = 0.0

x0 = 2.0
limitLeft = y_Left[end] = 1.0
meanLeft = mean(y_Left) = 1.0
stdDevLeft = std(y_Left) = 0.0
-----
limitRight = y_Right[end] = 0.9999989999999999
meanRight = mean(y_Right) = 0.9999925145291394
stdDevRight = std(y_Right) = 3.9868430925503734e-5

x0 = 3.0
limitLeft = y_Left[end] = 1.000000000139778e-6
meanLeft = mean(y_Left) = 7.4854708605545105e-6
stdDevLeft = std(y_Left) = 3.986843092550373e-5
-----
limitRight = y_Right[end] = 0.0
meanRight = mean(y_Right) = 0.0
stdDevRight = std(y_Right) = 0.0

```

We see that the estimated numerical limits correspond well to what we computed previously in Chapter 4.1, when all we knew to do was read values from a graph. In each case, $y[N] \approx \mu$ and $\sigma \approx 0$. ■

Let's now try the more challenging example of $\sin(\frac{1}{x})$, $x > 0$ and zero otherwise.

Example 4.9. Determine the limits from the left and right for $g(x) := \sin(\frac{1}{x})$, $x > 0$ and $g(x) = 0$, $x \leq 0$, for $x_0 \in \{0.0, \frac{1}{\pi}, \frac{2}{\pi}\}$.

Solution:

```

1 eta = 1e3:1e3:1e6
2 h_Right = (1.0)./ eta
3 h_Left = -h_Right
4 x0Vec = [0.0 1.0/pi 2.0/pi]
5
6 g(x) = x > 0 ? sin(1/x) : 0.0
7
8 delta = h_Right[end]
9
10 @show delta
11 println(" ")
12
13 for k = 1:length(x0Vec)
14     x0 = x0Vec[k]
15     y_Left = g.(x0 .+ h_Left)

```

```

16 y_Right = g.(x0 .+ h_Right)
17 #
18 @show x0
19 @show limitLeft = y_Left[end]
20 @show meanLeft = mean(y_Left)
21 @show stdDevLeft = std(y_Left)
22 println ("-----")
23 @show limitRight = y_Right[end]
24 @show meanRight = mean(y_Right)
25 @show stdDevRight = std(y_Right)
26 println (" ")
27 println (" ")
28
29 println (" ")
30 #
31 end

```

Output

```

delta = 1.0e-6

x0 = 0.0
limitLeft = y_Left[end] = 0.0
meanLeft = mean(y_Left) = 0.0
stdDevLeft = std(y_Left) = 0.0
-----
limitRight = y_Right[end] = -0.34999350217129294
meanRight = mean(y_Right) = -0.00011524369884409064
stdDevRight = std(y_Right) = 0.7075827861802105

x0 = 0.3183098861837907
limitLeft = y_Left[end] = -9.869635406884313e-6
meanLeft = mean(y_Left) = -7.392953161190493e-5
stdDevLeft = std(y_Left) = 0.0003944103587102789
-----
limitRight = y_Right[end] = 9.86957339451198e-6
meanRight = mean(y_Right) = 7.382758963882898e-5
stdDevRight = std(y_Right) = 0.0003925579978142143

x0 = 0.6366197723675814
limitLeft = y_Left[end] = 0.999999999996956
meanLeft = mean(y_Left) = 0.9999999949842888
stdDevLeft = std(y_Left) = 1.0037167987310087e-7
-----
limitRight = y_Right[end] = 0.999999999996956
meanRight = mean(y_Right) = 0.9999999950072798
stdDevRight = std(y_Right) = 9.976893600410653e-8

```

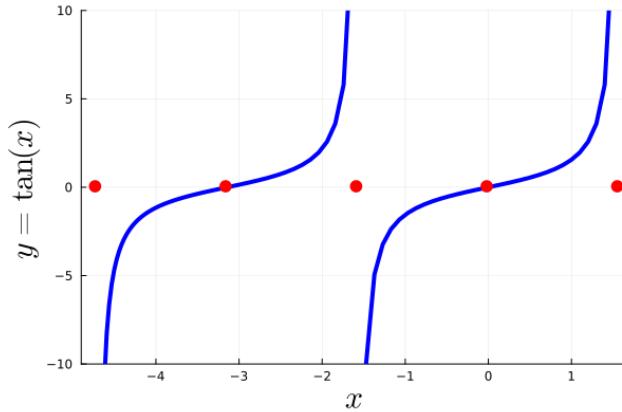
As we suspected from Example 4.4-(d), the left limit at the $x_0 = 0$ exists and is equal to zero, but the right limit does not exist. This shows up in the numerical calculations by

- $|\mu_R - y_R[N]| \approx 0.35$ is “large”; and
- $\sigma_R \approx 0.708$ is also large.

These observations are compatible with the extreme oscillations on the right side of the origin, as seen in Fig. 4.3-(d). All other one-sided limits exist. ■

We next illustrate unbounded one-sided limits at a point.

Example 4.10. Determine the limits from the left and right at $x_0 \in \{0.0, -\frac{\pi}{2}, -\pi\}$ for $g(x) := \tan(x)$, $(x \bmod \pi) \neq \frac{\pi}{2}$ (i.e., x is not an odd multiple of $\pi/2$) and $g(x) = 0$ otherwise.



Solution: Because $(x \bmod \pi) = \frac{\pi}{2} \implies \sin(x) = \pm 1$ and $\cos(x) = 0$, we expect unbounded limits at these points (numerator is non-zero and the denominator is equal to zero). Similarly, because $(x \bmod \pi) = 0 \implies \sin(x) = 0$ and $\cos(x) = \pm 1$, we expect the limits to be zero. Indeed, based on the plot, we can see the one-sided limits at $x_0 = 0$ and $x_0 = -\pi$ all exist and are equal to zero, which is confirmed by the numerical calculations below. Moreover, the limit from the left at $-\frac{\pi}{2}$ appears to be $+\infty$, while the limit from the right appears to be $-\infty$, consistent with the “analytical” study of the problem.

```

1 using Statistics
2 eta = 1e3:1e3:1e6
3 h_Right = (1.0) ./ eta
4 h_Left = -h_Right
5 x0Vec = [0.0 -pi/2 -pi]
6
7 g(x) = abs(x % pi) > 1e-10 ? tan(x) : 0.0
8
9 delta = h_Right [end]
10
11 @show delta
12 println (" ")
13
14 for k = 1:length(x0Vec)
15     x0 = x0Vec[k]
16     y_Left = g.(x0 .+ h_Left)
17     y_Right = g.(x0 .+ h_Right)
18     #
19     @show x0
20     @show limitLeft = y_Left [end]
21     @show meanLeft = mean(y_Left)
22     @show stdDevLeft = std(y_Left)
23     println ("-----")
24     @show limitRight = y_Right [end]
25     @show meanRight = mean(y_Right)
26     @show stdDevRight = std(y_Right)
27     println (" ")
28     println (" ")

```

```

29
30     println( "  " )
31     #
32 end

```

Output

```

delta = 1.0e-6

x0 = 0.0
limitLeft = y_Left[end] = -1.0000000000003333e-6
meanLeft = mean(y_Left) = -7.485471261235949e-6
stdDevLeft = std(y_Left) = 3.986843990838751e-5
-----
limitRight = y_Right[end] = 1.0000000000003333e-6
meanRight = mean(y_Right) = 7.485471261235949e-6
stdDevRight = std(y_Right) = 3.986843990838751e-5

x0 = -1.5707963267948966
limitLeft = y_Left[end] = 1.0000000001431657e6
meanLeft = mean(y_Left) = 500500.000018453
stdDevLeft = std(y_Left) = 288819.43611785874
-----
limitRight = y_Right[end] = -1.000000000020701e6
meanRight = mean(y_Right) = -500499.9999775702
stdDevRight = std(y_Right) = 288819.4360824532

x0 = -3.141592653589793
limitLeft = y_Left[end] = -1.0000000000176466e-6
meanLeft = mean(y_Left) = -7.485471261117652e-6
stdDevLeft = std(y_Left) = 3.986843990838485e-5
-----
limitRight = y_Right[end] = 1.000000000262576e-6
meanRight = mean(y_Right) = 7.485471261362585e-6
stdDevRight = std(y_Right) = 3.986843990838485e-5

```

This can be the Point where the Wheels Come off the Track!

Wheels off the track! Why? Because this Chapter is difficult. It speaks to the heart of Calculus. Once you master limits and continuity of functions, everything else in the book seems easier by comparison. It's maybe a good time to assess where you stand. Are your study habits really good enough? You have been in school for, like, forever. But have you really learned how to learn? While you should seek out resources on your own, here are a few suggestions:

- [The Biggest Myth In Education](#) by Veritasium. Speaks about learning “styles”.
- [Math is Hard for Everyone and Don't Be A Math Imposter](#) by the Math Sorcerer. **Scan the comments before watching the videos.**
- Basic things are important:
 - Go to Office Hours the moment you catch yourself falling behind or being confused by any topic. Avoiding office hours and spending 2-3 weeks to figure it out yourself will only hurt you. It's okay to ask for help or an extra

explanation!

- Read the relevant textbook material before attending (watching) the lecture (you can see what topics will be covered on the course schedule in the syllabus).
 - Review the homework solutions once they are posted to learn from your mistakes.
 - Don't wait until the last minute on a project. Complete small chunks of it over the span of 2 weeks.
- **Put in the time: (4 hours lectures and recitation) + (8 hours assignments and studying) = 12 hours per week for a typical student in the course. Are you typical? Lying to yourself is the worst.**
 - Your author struggled Y1 in College (lied to himself in semester one) and had to right the ship in semester two or lose essential scholarships and quit college. Here is what worked for him (one person out of the seven or eight billion on the planet).
 - Carefully review all extracurricular activities. Trim to only the ones you need to stay sane. Once your results match your expectations, you can slowly add more fun back in. This is called feedback control!
 - Scan the Chapter before attending (or watching) the lecture so that you know what to expect and where the challenging points are; should take 15 minutes or less.
 - Take essential notes during the lecture, keeping in mind that the instructor's notes are posted after each lecture.
 - Before the next lecture, rewrite and condense your notes to one side of a single page. The book has the full details on most material. You do not need all of that in your notes. Being thoughtful about condensing your notes provides a third pass on the material. Along with the first two passes, scanning the Chapter before the lecture and attending or watching the lecture, you should be good to go. If not, you know what questions to ask during Office Hours!

4.3 The Power of Continuity Done Right

Now that we've done the hard work of understanding limits, we have the tools to move beyond our folksy definition of "a function is continuous if we can draw its graph without lifting our pencil from the paper". **The real definition of continuity at a point will provide a major benefit to us when it comes to evaluating limits and computing integrals.**

4.3.1 Continuity at a Point and Continuous Everywhere on the Domain of Definition

Continuity at a Point and Continuous Everywhere

Definition 4.11. Suppose that $I \subset \mathbb{R}$ is an open interval, that is, I is either (a, b) , (a, ∞) , $(-\infty, b)$, or $(-\infty, \infty)$ as in Definition 1.1.

(a) A function $f : I \rightarrow \mathbb{R}$ is **continuous at $x_0 \in I$** if the two one-sided limits at x_0 exist and agree with the value of the function at x_0 . In other symbols,

$$\begin{aligned} \lim_{x \rightarrow x_0^-} f(x) &= f(x_0), \text{ and} \\ \lim_{x \rightarrow x_0^+} f(x) &= f(x_0). \end{aligned} \tag{4.8}$$

For intervals containing an endpoint, such as $[a, b]$, $[a, b)$, $[a, \infty)$, $(a, b]$, and $[b, \infty)$, we have to treat the endpoints differently, but all other points in the interval are treated as above.

(b) For I a closed or half-open interval as highlighted above and $x_0 \in I$ **not an endpoint**, f is continuous at x_0 if the conditions in (4.8) hold. At the endpoints of the interval, only one of the one-sided limits makes sense. Hence, f is

continuous at $x_0 = a$, a lower endpoint, if

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

and f is continuous at $x_0 = b$, an upper endpoint, if

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

(c) A function f that is not continuous at x_0 is said to be **discontinuous at x_0** or x_0 is said to be a **point of discontinuity** of the function.

(d) If f is continuous at each point $x_0 \in I$, its domain of definition, then f is simply said to be **continuous**. If you want to say that f is **continuous everywhere in (or on) its domain of definition**, that is fine too!

(e) If f is not continuous at all points of its domain of definition, then it is **discontinuous**. It only takes one point of discontinuity to make a function discontinuous.

Remark 4.12. For $x_0 \in \mathbb{R}$, it can NEVER be the case that $f(x_0) = \pm\infty$ because **the codomains of our functions contain only real numbers**. We recall that ∞ is a concept and not a number. Limits can equal ∞ because we allow functions to grow without bound, such as the function $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = x$ or $f : (0, \infty) \rightarrow (0, \infty)$ by $f(x) = \frac{1}{x}$.

Equation 4.8 can be rearranged as

$$\begin{aligned} \lim_{x \rightarrow x_0^-} (f(x) - f(x_0)) &= 0, \text{ and} \\ \lim_{x \rightarrow x_0^+} (f(x) - f(x_0)) &= 0. \end{aligned} \tag{4.9}$$

We'll use this equivalence in Example 4.13.

Example 4.13. For the following functions, determine, if any exist, their points of discontinuity. If there are no points of discontinuity, then state the function is continuous everywhere in its domain of definition.

(a) $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = x^3$.

(b) $f : [2, \infty) \rightarrow \mathbb{R}$ by $f(x) = x^2$.

(c) $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = \begin{cases} x^3 & x \geq 0 \\ x^2 & x < 0. \end{cases}$

(d) $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = \begin{cases} 1 + x & x \geq 0 \\ x^2 & x < 0. \end{cases}$

(e) $f : [0, \infty) \rightarrow \mathbb{R}$ by $f(x) = \begin{cases} \sin(\frac{1}{x}) & x > 0 \\ 0 & x = 0 \end{cases}$

(f) $f : (0, \infty) \rightarrow \mathbb{R}$ by $f(x) = \log_a(x)$, where $a > 1$.

Solutions:

(a) **Ans.** x^3 is continuous at all points of its domain of definition. We show this via the definition. Let $x_0 \in \mathbb{R}$ be arbitrary. Then,

$$\lim_{x \rightarrow x_0^-} x^3 := \lim_{h \rightarrow 0^-} (x_0 + h)^3 = \lim_{h \rightarrow 0^-} (x_0^3 + 3hx_0^2 + 3h^2x_0 + h^3) = x_0^3 + 0 + 0 + 0 = x_0^3$$

$$\lim_{x \rightarrow x_0^+} x^3 := \lim_{h \rightarrow 0^+} (x_0 + h)^3 = \lim_{h \rightarrow 0^+} (x_0^3 + 3hx_0^2 + 3h^2x_0 + h^3) = x_0^3 + 0 + 0 + 0 = x_0^3.$$

Because x_0 was arbitrary, x^3 is continuous at all points of its domain of definition. **Note that here, we did the limits by inspection after some simple algebra (namely, the Binomial Theorem).**

- (b) **Ans.** $(\bullet)^2 : [2, \infty) \rightarrow \mathbb{R}$ is continuous at all points of its domain of definition. We show this via the definition. We first consider $2 < x_0 < \infty$, but otherwise arbitrary. Then,

$$\lim_{x \rightarrow x_0^-} x^2 := \lim_{h \rightarrow 0^-} (x_0 + h)^2 = \lim_{h \rightarrow 0^-} (x_0^2 + 2hx_0 + h^2) = x_0^2 + 0 + 0 = x_0^2$$

$$\lim_{x \rightarrow x_0^+} x^2 := \lim_{h \rightarrow 0^+} (x_0 + h)^2 = \lim_{h \rightarrow 0^+} (x_0^2 + 2hx_0 + h^2) = x_0^2 + 0 + 0 = x_0^2.$$

So far, so good! Let's now let $x_0 = 2$. We only check the limit from the right, namely,

$$\lim_{x \rightarrow 2^+} x^2 := \lim_{h \rightarrow 0^+} (2 + h)^2 = \lim_{h \rightarrow 0^+} (4 + 4h + h^2) = 4 + 0 + 0 = 4 = x_0^2.$$

Because we checked continuity at all points of its domain, $(\bullet)^2 : [2, \infty) \rightarrow \mathbb{R}$ is continuous at all points of its domain of definition. **Once again, the limits themselves were computed by inspection.**

- (c) **Ans.** $f(x) = \begin{cases} x^3 & x \geq 0 \\ x^2 & x < 0. \end{cases}$ is continuous at all points of its domain of definition. We solve this one a bit differently, because continuity for $x > 0$ and $x < 0$ can be checked the same as we did for parts (a) and (b). The only point that seems a bit different is the origin, but because x^3 and x^2 both vanish at the origin, there should not be a problem; to be sure, we'll check it anyway!

$$\lim_{x \rightarrow 0^-} x^2 := \lim_{h \rightarrow 0^-} (0 + h)^2 = \lim_{h \rightarrow 0^-} h^2 = 0 = 0^2$$

$$\lim_{x \rightarrow 0^+} x^3 := \lim_{h \rightarrow 0^+} (0 + h)^2 = \lim_{h \rightarrow 0^+} h^3 = 0 = 0^3.$$

Hence, the definition of continuity at a point is met for all points of the domain of definition of the function.

- (d) **Ans.** $f(x) = \begin{cases} 1 + x & x \geq 0 \\ x^2 & x < 0. \end{cases}$ is discontinuous at $x_0 = 0$ and continuous everywhere else. Once again, checking the limits for all $x_0 \neq 0$ is straightforward, and we leave that to the learner. For $x_0 = 0$, we see that

$$\lim_{x \rightarrow 0^-} x^2 := \lim_{h \rightarrow 0^-} (0 + h)^2 = \lim_{h \rightarrow 0^-} h^2 = 0 \neq f(0) = 1$$

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

Because the left limit at zero is not equal to the value of the function at zero, the function is discontinuous at the origin. **It was not necessary to compute the limit from the right; we simply did it for extra practice.**

- (e) **Ans.** $f(x) = \begin{cases} \sin(\frac{1}{x}) & x > 0 \\ 0 & x = 0 \end{cases}$ is discontinuous at the origin, and continuous everywhere else in its domain of definition, meaning, it is continuous for all $x > 0$. The function is discontinuous at $x_0 = 0$ because we showed earlier, in Example 4.4, that $\lim_{x \rightarrow 0^+} \sin(1/x)$ does not exist. Hence, the limit from the right cannot equal the value of the function, which is zero.

For $x > 0$, set $y := \frac{1}{x}$. Then $y > 0 \iff x > 0$ reduces us to whether $\sin(y)$ is continuous or not for $y > 0$. But we "know" from experience that it is continuous, so we're done! Really? For now, in any case. A formal proof can be given using trig identities and the bounds given in Fig. 4.7.

- (f) **Ans.** $f(x) = \log_a(x)$, where $a > 1$ is continuous everywhere in its domain of definition, which we know is $(0, \infty)$. We first show continuity at the origin and then use the pattern established for the special case to handle the general case. Recall from Prop. 1.61, for all $a > 1$, $\log_a(1) = 0$. Hence, we need to show that both

$$\lim_{h \rightarrow 0^-} \log_a(1 + h) = 0, \text{ and}$$

$$\lim_{h \rightarrow 0^+} \log_a(1 + h) = 0.$$

This limit is less obvious than the previous limits; hence, we give an epsilon-delta proof: Let $\epsilon > 0$ be arbitrary. We seek $\delta > 0$ such that

$$-\delta < h < 0 \implies |\log_a(1 + h)| \leq \epsilon, \text{ and}$$

$$0 < \delta < h \implies |\log_a(1 + h)| \leq \epsilon = 0,$$

where we note that we can use different values for δ in the two cases. From Prop. 1.9-(c), we have that

$$|\log_a(1+h)| < \epsilon \iff -\epsilon < \log_a(1+h) < \epsilon.$$

Because, for $a > 1$, a^x is strictly monotonically increasing on $(-\infty, \infty)$, taking the exponential preserves the inequalities; hence,

$$-\epsilon < \log_a(1+h) < \epsilon \iff a^{-\epsilon} < 1+h < a^\epsilon.$$

Therefore¹,

$$|\log_a(1+h)| < \epsilon \iff \underbrace{a^{-\epsilon} - 1}_{<0} < h < \underbrace{a^\epsilon - 1}_{>0},$$

which implies that

$$\begin{aligned} -\delta < h < 0 &\implies |\log_a(1+h)| \leq \epsilon, \text{ when } \delta = 1 - a^{-\epsilon}, \text{ and} \\ 0 < h < \delta &\implies |\log_a(1+h)| \leq \epsilon, \text{ when } \delta = a^\epsilon - 1, \end{aligned}$$

establishing continuity at the origin. To show continuity at a general $x > 0$, we want to show that the left and right limits as h tends to zero of $\log_a(x+h) - \log_a(x) = 0$. But, using log properties,

$$\log_a(x+h) - \log_a(x) = \log_a\left(\frac{x+h}{x}\right) = \log_a\left(1 + \frac{h}{x}\right).$$

Because $x > 0$, if we define $\bar{h} := \frac{h}{x}$, then $\bar{h} \rightarrow 0^- \iff h \rightarrow 0^-$ and $\bar{h} \rightarrow 0^+ \iff h \rightarrow 0^+$. Hence,

$$\begin{aligned} \lim_{h \rightarrow 0^-} \log_a\left(1 + \frac{h}{x}\right) &= \lim_{\bar{h} \rightarrow 0^-} \log_a(1 + \bar{h}), \text{ and} \\ \lim_{h \rightarrow 0^+} \log_a\left(1 + \frac{h}{x}\right) &= \lim_{\bar{h} \rightarrow 0^+} \log_a(1 + \bar{h}), \end{aligned}$$

reducing the problem to the case we have already considered, namely, continuity at $x_0 = 1$.

■

4.3.2 Two-sided Limits and Continuous Functions

You might imagine that if the limit from the left and the right both exist at a point, then you can combine them into a single limit. And you can!

Two-sided Limit at a Point

To set the stage, consider a function $f : A \rightarrow \mathbb{R}$ where A is a subset of \mathbb{R} , $x_0 \in A$, and there exists $\delta > 0$ such that both $(x_0 - \delta, x_0) \subset A$ and $(x_0, x_0 + \delta) \subset A$. Then $x_0 + h \in A$ for all $0 < |h| < \delta$, meaning we can evaluate

$$\begin{aligned} \lim_{h \rightarrow 0^-} f(x_0 + h), \text{ and,} \\ \lim_{h \rightarrow 0^+} f(x_0 + h), \end{aligned}$$

the left and right limits.

Definition 4.14. Suppose that $f : A \rightarrow \mathbb{R}$, where A is a subset of \mathbb{R} and $x_0 \in A$. Then the limit as x approaches x_0 exists and is equal to L , a finite real number, if for all $\epsilon > 0$, there exists $\delta > 0$, such that $0 < |h| < \delta \implies |f(x_0 + h) - L| \leq \epsilon$. The limit is denoted as

$$\lim_{x \rightarrow x_0} f(x) = L.$$

“Two-sided limits” are typically called “limits,” with the adjective “two-sided” dropped. It’s only in the case of a one-sided limit that you need to call it out its **chirality** (“sidedness” or handedness of the limit).

¹For all $a > 1$, $a^x < 1 \iff x < 0$.

Fact 4.15: Equivalence between a Two-sided Limit and Two One-sided Limits

Two of one, versus one of the other: the following are equivalent statements for $|L| < \infty$:

(a) $\lim_{x \rightarrow x_0^-} f(x) = L$ and $\lim_{x \rightarrow x_0^+} f(x) = L$.

(b) $\lim_{x \rightarrow x_0} f(x) = L$.

Notes: Statements being equivalent means $(a) \iff (b)$, where we note that **(a) is true if, and only if, both one-sided limits exist and are equal to L** . Just as with a one-sided limit, we never evaluate the function at x_0 when evaluating a two-sided limit. **In theory, you can replace every “limit” with two one-sided limits. However, this is not common practice and some faculty or GSIs may mark you wrong for using two one-sided limits in place of a single (two-sided) limit!** People fear what they do not fully understand. You are expected to fall into line...ha!

Using the above Fact, continuity of a function at a point can be checked with a single (two-sided) limit instead of two one-sided limits. We'll leave that to you. What we'll do instead is give the classical epsilon-delta definition of a function being continuous at a point and then **relate all of the (commonly seen) definitions of continuity at a point**.

Classical Definition of Continuity at a Point

Definition 4.16. Suppose that $f : I \rightarrow \mathbb{R}$, where I is an interval and $x_0 \in A$. Then f is continuous at x_0 if for all $\epsilon > 0$, there exists $\delta > 0$, such that

$$\{x \mid |x - x_0| < \delta\} \subset I \text{ and } |x - x_0| < \delta \implies |f(x) - f(x_0)| \leq \epsilon.$$

Note: If $\{x \in \mathbb{R} \mid |x - x_0| < \delta\} \notin I$, then it does not make sense to evaluate the function at all of these points because some of them (at least) are not in the domain of the function!

Fact 4.17: Equivalent Ways to Define Continuity at a Point

The following statements are equivalent for a function $f : I \rightarrow \mathbb{R}$, where I is an interval:

(a) f is continuous at x_0 by Def. 4.16.

(b) $\lim_{x \rightarrow x_0^-} f(x) = f(x_0)$ and $\lim_{x \rightarrow x_0^+} f(x) = f(x_0)$, meaning that Def. 4.11 of continuity at a point holds.

(c) $\lim_{x \rightarrow x_0} f(x) = f(x_0)$ (equivalent to the above because two one-sided limits are always equivalent to a single one-sided limit.)

More examples of using limits and checking continuity with the modified formulations would be repetitive. Instead, let's list some common continuous functions.

Proposition 4.18: Common Continuous Functions

The following functions are continuous for all $x \in \mathbb{R}$ unless indicated otherwise:

Polynomial Functions

- Constant function: $f(x) = c$
- Linear function: $f(x) = ax + b$
- Quadratic function: $f(x) = ax^2 + bx + c$
- Higher-degree polynomials: $f(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0$

Trigonometric Functions

- Sine: $f(x) = \sin(x)$
- Cosine: $f(x) = \cos(x)$
- Tangent (continuous except at odd multiples of $\pm\frac{\pi}{2}$): $f(x) = \tan(x)$

Power Function

- Power (continuous for $x > 0$ and $y \in \mathbb{R}$): $f(x) = x^y$

Exponential and Logarithmic Functions

- Exponential: $f(x) = a^x$ (where $a > 1$)
- Natural Exponential: $f(x) = e^x$
- Logarithm (continuous for $x > 0$): $f(x) = \log_a(x)$ (where $a > 1$)
- Natural Logarithm (continuous for $x > 0$): $f(x) = \ln(x)$

Rational Functions (subject to domain restrictions)

- $f(x) = \frac{p(x)}{q(x)}$ where $p(x)$ and $q(x)$ are polynomials and $q(x) \neq 0$ (means, rational functions are continuous for all $x \in \mathbb{R}$ such that $q(x) \neq 0$)

Root Functions

- Square Root (continuous for $x \geq 0$): $f(x) = \sqrt{x}$
- Cube Root: $f(x) = \sqrt[3]{x}$
- n -th Root (continuous for $x \geq 0$ if $n \in \mathbb{N}$ is even, all $x \in \mathbb{R}$ if $n \in \mathbb{N}$ is odd): $f(x) = \sqrt[n]{x}$

Other Special Functions

- Absolute Value: $f(x) = |x|$
- Triangle: $f(x) = \text{tri}(x)$
- Radial Basis Function: $f(x) = e^{-\frac{(x-x_c)^2}{s^2}}$
- Gaussian (Normal Distribution) with mean μ and standard deviation σ : $f(x) = \frac{1}{\sqrt{2\pi\sigma^2}}e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

Reassurance on Limits and Continuity

As a practicing engineer with a BSE, you will rarely, if ever, need to use the epsilon-delta definition of a limit or continuity at a point. You will mostly handle these concepts by “inspection”, and if you’re a good engineer, you’ll “mostly get the answers right”.

Depending on your goals, if you go to graduate school, being comfortable with theory may become more important. Hence, we’ve tried so far in the textbook to represent the perspective of both kinds of learners: those who chose engineering out of a love for math and physics per se and those who like math and physics, but really came to engineering out of a desire to build stuff and improve the lives of people. Your author has advised amazing graduate students on both ends of the spectrum and everywhere in between. Everyone’s path is different. That is why Michigan’s Robotics Graduate Program starts with **ROB 501 Mathematics for Robotics**: to level the mathematical playing field for all learners.

If, out of personal interest, you are paying close attention to the theoretical discussions so far (e.g., the proofs at the end of each Chapter), then you may wish to target Michigan’s MATH 451 course (or its equivalent). Where we pull back and hesitate on the details, MATH 451 goes full speed ahead. If you want the extra rigor now: **Rigorous limit proof ultimate study**

guide: (ϵ - δ , ϵ - N , M - δ , and M - N proofs) by **BlackPenRedPen**.

4.3.3 Key Properties of Continuous Functions or Why We Care about Continuity

For the purpose of developing and understanding the fundamentals of Calculus, the following Proposition ranks among the most important results in this course.

Proposition 4.19: Taking Limits Inside a Function

Consider two functions $g : X \rightarrow Y$ and $f : Y \rightarrow Z$, where X , Y , and Z are subsets of \mathbb{R} . Suppose that f is continuous at a point $y_0 \in Y$. Then, the following hold:

- (a) If $x_0 \in \mathbb{R}$ and $\lim_{x \rightarrow x_0} g(x) = y_0$ (i.e., the limit exists and equals y_0), then $\lim_{x \rightarrow x_0} f(g(x)) = f(\lim_{x \rightarrow x_0} g(x)) = f(y_0)$. **We say the limit can be taken “inside” the function f .**
- (b) Similarly, if $\lim_{x \rightarrow \infty} g(x) = y_0$, then $\lim_{x \rightarrow \infty} f(g(x)) = f(\lim_{x \rightarrow \infty} g(x)) = f(y_0)$.
- (c) If $\lim_{x \rightarrow -\infty} g(x) = y_0$, then $\lim_{x \rightarrow -\infty} f(g(x)) = f(\lim_{x \rightarrow -\infty} g(x)) = f(y_0)$.

In short, if $g(x)$ has a limit (as x approaches a finite point or as x approaches $\pm\infty$), and f is continuous at the limit of g (i.e., y_0), then, when evaluating a limit of $f(g(x)) = f \circ g(x)$, the limit can be taken inside the function f and applied directly to g .

You may enjoy the video [Mu Prime Math](#) “When can we switch the limit and function?”

Example 4.20. Use Prop. 4.19 to evaluate the following limits. In each case, as part of your solution, identify the functions $g : X \rightarrow Y$ and $f : Y \rightarrow Z$, as well as the limit point of g , which we have been calling y_0 .

$$(a) \lim_{x \rightarrow 0} \sin\left(\frac{x}{1+x^2}\right)$$

$$(b) \lim_{x \rightarrow \frac{\pi}{2}} \frac{3 + \sin^3(x)}{\sin^2(x) + 3 \sin(x) + 4}$$

$$(c) \lim_{x \rightarrow \infty} \ln(e^{-2x} + e^{-x} + 1)$$

Solutions:

$$(a) \text{Ans. } \lim_{x \rightarrow 0} \sin\left(\frac{1+x}{1+x^2}\right) = \sin(1) \approx 0.8415.$$

We define $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(y) = \sin(y)$, $g : \mathbb{R} \rightarrow \mathbb{R}$ by $g(x) = \frac{1+x}{1+x^2}$, and note that $\lim_{x \rightarrow 0} g(x) = \lim_{x \rightarrow 0} \frac{1+x}{1+x^2} = 1 =: y_0$. The limit exists because g is continuous everywhere on its domain of definition per Prop. 4.18 Common Continuous Functions. Finally, also by Prop. 4.18, $f(y)$ is continuous everywhere and hence at the point $y_0 := 1$.

$$(b) \text{Ans. } \lim_{x \rightarrow \frac{\pi}{2}} \frac{3 + \sin^3(x)}{\sin^2(x) + 3 \sin(x) + 4} = \frac{1}{2} \text{ (because } \sin(\pi/2) = 1).$$

We define $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(y) = \frac{3+y^3}{y^2+3y+4}$, $g : \mathbb{R} \rightarrow \mathbb{R}$ by $g(x) = \sin(x)$, and note that $\lim_{x \rightarrow \frac{\pi}{2}} g(x) = \lim_{x \rightarrow \frac{\pi}{2}} \sin(x) = 1.0 =: y_0$.

The limit exists because g is continuous everywhere on its domain of definition per Prop. 4.18 Common Continuous Functions. Finally, also by Prop. 4.18, $f(y)$ is continuous everywhere and hence at the point $y_0 := 0$.

(c) **Ans.** $\lim_{x \rightarrow \infty} \ln(e^{-2x} + e^{-x} + 1) = 0$.

This problem requires, given what we know at this point in the Chapter, two applications of Prop. 4.19. We will write $\ln(e^{-2x} + e^{-x} + 1)$ as the composition of three functions, $\ln(e^{-2x} + e^{-x} + 1) =: f(g(h(x)))$, where

$$\begin{aligned} f : (0, \infty) &\rightarrow \mathbb{R} \text{ by } f(z) = \ln(z) \\ g : (0, \infty) &\rightarrow (1, \infty) \text{ by } g(y) = y^2 + y + 1 \\ h : \mathbb{R} &\rightarrow (0, \infty) \text{ by } h(x) = e^{-x}, \end{aligned}$$

where we note that $g(h(x)) = (e^{-x})^2 + e^{-x} + 1 = e^{-2x} + e^{-x} + 1$ because $(a^\alpha)^\beta = a^{\alpha \cdot \beta}$ for all $a > 0$, and real numbers α and β .

All three functions are continuous everywhere on their domains of definition. Moreover,

$$\begin{aligned} \lim_{x \rightarrow \infty} h(x) &= \lim_{x \rightarrow \infty} e^{-x} = 0 =: y_0 \\ g(y_0) &= g(0) = 1 =: z_0 \\ f(z_0) &= f(1) = \ln(1) = 0. \end{aligned}$$

Therefore,

$$\begin{aligned} \lim_{x \rightarrow \infty} \ln(e^{-2x} + e^{-x} + 1) &= \lim_{x \rightarrow \infty} f(g(h(x))) \\ &= f\left(\lim_{x \rightarrow \infty} g(h(x))\right) \\ &= f\left(g\left(\underbrace{\lim_{x \rightarrow \infty} h(x)}_{y_0=0}\right)\right) \\ &= f\left(g(0)\right) \\ &= \underbrace{f(1)}_{z_0=1} \\ &= \ln(1) \\ &= 0. \end{aligned}$$

Note: All of these limits can be approximated numerically using the same techniques we have already demonstrated. ■

The above problem was mainly an exercise in writing a “complicated” function as a composition of simpler functions. The next two problems are more representative of how Prop. 4.19 is typically used. Their solutions will illustrate how a change of variable can make an apparently intractable problem become tractable.

Example 4.21. Using the fact that $\lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n =: e$, compute the following limits:

(a) $\lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right)^{-n}$;

(b) $\lim_{n \rightarrow \infty} \left(1 + \frac{x}{n}\right)^n$, where $x \in \mathbb{R}$ is arbitrary; and

(c) $\lim_{h \rightarrow 0} \frac{\ln(x+h) - \ln(x)}{h}$, where $x > 0$.

Hints: For (a), note that $\left(1 - \frac{1}{n}\right)^{-1} = \left(\frac{n-1}{n}\right)^{-1} = \left(\frac{n}{n-1}\right) = \left(\frac{n-1+1}{n-1}\right) = 1 + \frac{1}{n-1}$. Also, Prop. 2.35 is helpful. Finally, remember your power rules!

For (b), when $x = 0$, the limit is clearly one, leaving two cases to investigate, $x > 0$ and $x < 0$. First, assume $x > 0$ and use the substitution $m := \frac{n}{x}$ and look for the power function $f : (0, \infty) \rightarrow (0, \infty)$ by $f(y) = y^x$. Note that here, the variable is y , and x is a fixed positive constant. Follow a similar pattern for $x < 0$.

For (c), recall your logarithm rules for $\ln(x) \pm \ln(y)$ and $\ln(x^y)$. Try the substitution, $n := \frac{1}{h}$.

Solutions:

(a) **Ans.** $\lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right)^{-n} = e$.

Using the Hint, we have

$$\begin{aligned} \lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right)^{-n} &= \lim_{n \rightarrow \infty} \left(\left(1 - \frac{1}{n}\right)^{-1}\right)^n \\ &= \lim_{n \rightarrow \infty} \left(1 + \frac{1}{n-1}\right)^n. \end{aligned}$$

This is very close to the limit for e , but we're not quite there yet. We do a change of variable $m := n - 1$ so that

$$n \rightarrow \infty \iff m \rightarrow \infty.$$

Hence,

$$\begin{aligned} \lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right)^{-n} &= \lim_{n \rightarrow \infty} \left(1 + \frac{1}{n-1}\right)^n \\ &= \lim_{m \rightarrow \infty} \left(1 + \frac{1}{m}\right)^{m+1} \\ &= \lim_{m \rightarrow \infty} \underbrace{\left(1 + \frac{1}{m}\right)^m}_{f(m)} \cdot \underbrace{\left(1 + \frac{1}{m}\right)}_{g(m)}, \end{aligned}$$

setting us up to apply Prop. 2.35, because

$$\lim_{m \rightarrow \infty} g(m) = 1 \text{ and } \lim_{m \rightarrow \infty} f(m) = e.$$

We conclude therefore that $\lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right)^{-n} = e$.

Note: Can you show that for all integers $k \geq 1$, $\lim_{n \rightarrow \infty} \left(1 + \frac{1}{n \pm k}\right)^n = \lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^{n \pm k} = e$? The key point is that k is a fixed (finite) value when applying Prop. 2.35.

(b) **Ans.** $\lim_{n \rightarrow \infty} \left(1 + \frac{x}{n}\right)^n = e^x$.

Using the hint, we assume $x > 0$ and set $m := \frac{n}{x}$ so that

$$n = m \cdot x \quad \text{and} \quad m \rightarrow \infty \iff n \rightarrow \infty.$$

Making the substitution, we have

$$\begin{aligned} \lim_{n \rightarrow \infty} \left(1 + \frac{x}{n}\right)^n &= \lim_{m \rightarrow \infty} \left(1 + \frac{1}{m}\right)^{m \cdot x} \\ &= \lim_{m \rightarrow \infty} \left(\left(1 + \frac{1}{m}\right)^m\right)^x. \end{aligned}$$

We define $f(y) := y^x$ and note that it is continuous for all $y > 0$. Hence,

$$\begin{aligned}\lim_{n \rightarrow \infty} \left(1 + \frac{x}{n}\right)^n &= \lim_{m \rightarrow \infty} \left(\left(1 + \frac{1}{m}\right)^m\right)^x \\ &= \left(\lim_{m \rightarrow \infty} \left(1 + \frac{1}{m}\right)^m\right)^x \quad (\text{take limit inside } f) \\ &= (e)^x \quad (\text{recognizing the limit}) \\ &= e^x.\end{aligned}$$

We still have $x < 0$ to go, but the pattern is now pretty clear. Note that because $x < 0$, $x = -|x|$ and $|x| \neq 0$. Hence, we set $m := \frac{n}{|x|}$ so that

$$n = m \cdot |x| \quad \text{and } m \rightarrow \infty \iff n \rightarrow \infty.$$

Making the indicated substitution, we have

$$\begin{aligned}\lim_{n \rightarrow \infty} \left(1 + \frac{x}{n}\right)^n &= \lim_{n \rightarrow \infty} \left(1 - \frac{|x|}{n}\right)^n \\ &= \lim_{m \rightarrow \infty} \left(1 - \frac{1}{m}\right)^{m \cdot |x|} \\ &= \lim_{m \rightarrow \infty} \left(\left(1 - \frac{1}{m}\right)^{-m}\right)^{-|x|},\end{aligned}$$

where on the last line, we used our power rules. We define $f(y) := y^x = y^{-|x|}$ and note that it is a continuous function of y for all $y > 0$. Hence,

$$\begin{aligned}\lim_{n \rightarrow \infty} \left(1 + \frac{x}{n}\right)^n &= \lim_{m \rightarrow \infty} \left(\left(1 - \frac{1}{m}\right)^{-m}\right)^{-|x|} \\ &= \left(\lim_{m \rightarrow \infty} \left(1 - \frac{1}{m}\right)^{-m}\right)^{-|x|} \quad (\text{take limit inside } f) \\ &= (e)^{-|x|} \quad (\text{recognizing the limit}) \\ &= e^{-|x|} \\ &= e^x \quad (\text{because } x = -|x|).\end{aligned}$$

- (c) **Ans.** For all $x > 0$, $\lim_{h \rightarrow 0} \frac{\ln(x+h) - \ln(x)}{h} = \frac{1}{x}$.

Using our logarithm rules,

$$\begin{aligned}\frac{\ln(x+h) - \ln(x)}{h} &= \frac{1}{h} \ln\left(\frac{x+h}{x}\right) \\ &= \ln\left(\left(\frac{x+h}{x}\right)^{\frac{1}{h}}\right) \\ &= \ln\left(\left(1 + \frac{h}{x}\right)^{\frac{1}{h}}\right).\end{aligned}$$

We make the substitution $n := \frac{1}{h}$ so that

$$n \rightarrow \infty \iff h \rightarrow 0^+.$$

Using this substitution, we have

$$\begin{aligned}
\lim_{h \rightarrow 0^+} \frac{\ln(x+h) - \ln(x)}{h} &= \lim_{h \rightarrow 0^+} \ln \left(\left(1 + \frac{h}{x}\right)^{\frac{1}{h}} \right) \\
&= \lim_{n \rightarrow \infty} \ln \left(\left(1 + \frac{1}{nx}\right)^n \right) \\
&= \ln \left(\lim_{n \rightarrow \infty} \left(1 + \frac{1}{nx}\right)^n \right) \quad (\text{taking the limit inside because } \ln(x) \text{ is continuous}) \\
&= \ln \left(e^{\frac{1}{x}} \right) \quad (\text{recognizing the limit by part (b) of this problem}) \\
&= \frac{1}{x} \quad (\text{using logarithm rules.})
\end{aligned}$$

In theory, we should now consider $\lim_{h \rightarrow 0^-} \frac{\ln(x+h) - \ln(x)}{h}$, but we leave it to the learner. ■

This following (challenging) limit is the key to understanding the definite integral, $\int_a^b e^x dx$.

Example 4.22. Evaluate $\lim_{h \rightarrow 0} \frac{e^h - 1}{h}$.

Solution: **Ans.** $\lim_{h \rightarrow 0} \frac{e^h - 1}{h} = 1$.

At first blush, none of the methods we have covered so far seems applicable. This is more typical of how real engineering problems go. Initially, you are totally convinced that nothing you learned in School is useful for the problem at hand, but then, all of a sudden, you see how the problem can be tackled if you look at it the right way. Here, we need to make a change of variable and apply logarithm rules, even though there is no logarithm in sight!

The power move is to define $u := e^h - 1$ so that $e^h = u + 1 \iff h = \ln(1 + u)$. We note $\ln(x)$ is continuous everywhere on $(0, \infty)$, and that $h \rightarrow 0 \iff u \rightarrow 0$. With this substitution, we have

$$\lim_{h \rightarrow 0} \frac{e^h - 1}{h} = \lim_{u \rightarrow 0} \frac{u}{\ln(1 + u)} = \lim_{u \rightarrow 0} \frac{1}{\frac{1}{u} \cdot \ln(1 + u)}.$$

From our logarithm properties, recall that

$$y \cdot \ln(x) = \ln(x^y),$$

and hence, $\frac{1}{u} \cdot \ln(1 + u) = \ln((1 + u)^{\frac{1}{u}})$, yielding

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

Now we have to stare at the problem a bit and realize that $(1 + u)^{\frac{1}{u}}$ looks vaguely familiar to us. In fact, if we defined $n := \frac{1}{u}$, we'd have

$$(1 + u)^{\frac{1}{u}} = \left(1 + \frac{1}{n}\right)^n \xrightarrow{n \rightarrow \infty} e,$$

Euler's number, and we know that $\ln(e) = 1$. How can these observations be useful to us?

We write $\frac{1}{\ln((1 + \frac{1}{n})^n)} = f(g(k(n)))$, where

$$\begin{aligned}
f(\beta) &:= \frac{1}{\beta}, \quad \text{continuous for all } \beta > 0 \\
g(\alpha) &:= \ln(\alpha), \quad \text{continuous for all } \alpha > 0 \\
k(n) &:= \left(1 + \frac{1}{n}\right)^n \xrightarrow{n \rightarrow \infty} e.
\end{aligned}$$

Then, as in Example 4.20, we have that

$$\begin{aligned}
\lim_{h \rightarrow 0^+} \frac{e^h - 1}{h} &= \lim_{n \rightarrow \infty} \frac{1}{\ln((1 + \frac{1}{n})^n)} \quad (\text{after two changes of variable}) \\
&= \lim_{n \rightarrow \infty} f(g(k(n))) \quad (\text{using our defined functions, two being continuous and one having a known limit}) \\
&= f\left(\lim_{n \rightarrow \infty} g(k(n))\right) \quad (\text{take limit inside } f) \\
&= f\left(g\left(\lim_{n \rightarrow \infty} k(n)\right)\right) \quad (\text{take limit inside } g) \\
&= f\left(g\left(\lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n\right)\right) \quad (\text{identify } k(n)) \\
&= f(g(e)) \quad (\text{use known limit}) \\
&= f(\ln(e)) \quad (\text{now we are on autopilot}) \\
&= f(1) \\
&= \frac{1}{1} \\
&= 1.
\end{aligned}$$

Note: Because we showed every step, the solution looks impossibly arduous. With practice, you will skip most of the steps because “they become obvious with sufficient practice”. The key ideas were to make a useful substitution and subsequently recognize when limits can be taken inside functions. ■

Example 4.23. Determine the limit $\lim_{n \rightarrow \infty} \left(\frac{n+3}{n}\right)^n$.

Solution: We first note that $\frac{n+3}{n} = 1 + \frac{3}{n}$, and hence, $\left(\frac{n+3}{n}\right)^n = \left(1 + \frac{3}{n}\right)^n$. We recognize immediately that this has something to do with Euler’s number, e . Indeed, from Example 4.21-(b), we have

$$\lim_{n \rightarrow \infty} \left(\frac{n+3}{n}\right)^n = \lim_{n \rightarrow \infty} \left(1 + \frac{3}{n}\right)^n = e^3.$$

How is your limit intuition coming along? **BlackPenRedPen** has a [video](#) to find out! No epsilons and deltas. He uses continuity and qualitative reasoning. The answer pivots on the difference between 0^- and 0^+ . Intrigued?

4.3.4 Application to Integrating the Natural Exponential, e^x

Proposition 4.24: Integrating e^x

For $x > 0$ finite,

$$\int_0^x e^y dy = e^x - 1. \quad (4.10)$$

Proof: In Chapter 3.1.3, we developed closed-form expressions for integrating monomials, namely, for $k \in \mathbb{N}$

$$\int_0^x y^k dy = \frac{y^{k+1}}{k+1}.$$

Here, we follow the same pattern to discover the closed-form expression for exponentials, namely,

- Form the Riemann Lower and Upper Sums as indicated in Fig. 4.5.
- By taking $\Delta y := \frac{x}{n}$, we will compute lower and upper bounds of (4.10), $\text{Area}_n^{\text{Low}}$ and $\text{Area}_n^{\text{Up}}$, that depend on n .

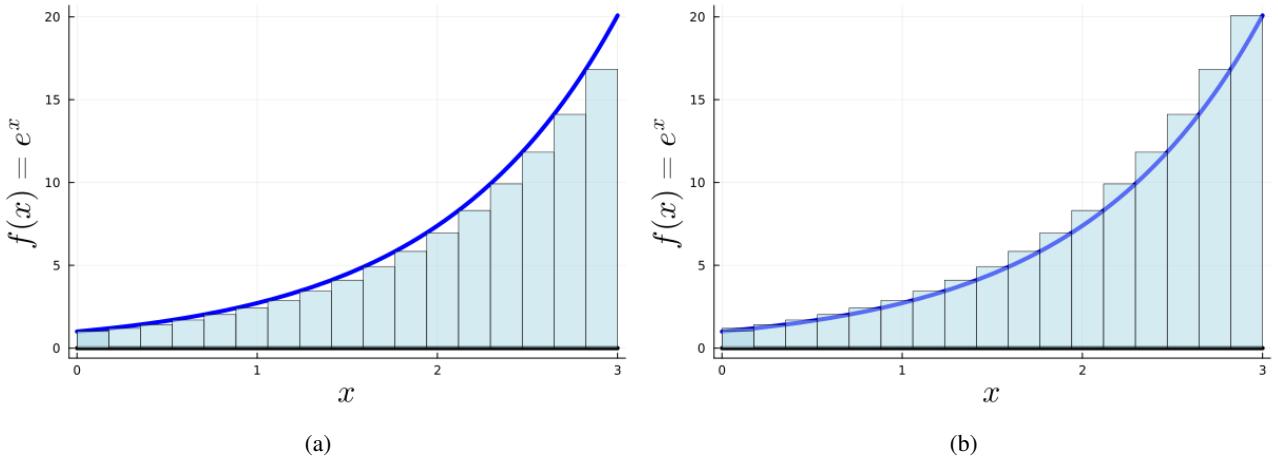


Figure 4.5: Under- and overapproximations for the area under e^x . (a) An underapproximation for the area. (b) An overapproximation for the area.

- Seek an exact value for (4.10) by evaluating $\lim_{n \rightarrow \infty} \text{Area}_n^{\text{Low}}$ and $\lim_{n \rightarrow \infty} \text{Area}_n^{\text{Up}}$.
- If the two limits exist and are equal, then we will have determined the definite integral.

Let's do it!

We work this out in full detail so as to review the **Approximation Principle** that undergirds most of Calculus. We let $n > 1$ be an integer, and because we are working on the interval $[0, x]$ for $x > 0$, we define $\Delta y := \frac{x-0}{n}$ along with $y_i := a + (i-1)\Delta y$, yielding

$$0 =: a = y_1 < y_2 < \cdots < y_n < y_{n+1} = b := x.$$

Because $f(y) = e^y$ is increasing, $f(y_i) \leq f(y) \leq f(y_{i+1})$ for $y \in [y_i, y_{i+1}]$, it is easy to determine the min and max values of the function over subintervals.

Important Observation: $e^{k\Delta y} = (e^{\Delta y})^k$ by Prop. 1.55-(g).

Riemann Lower Sum: We underestimate the total area between $[0, x]$ by summing up the underapproximations, $f(y_i) \cdot \Delta y$, for each interval $[y_i, y_{i+1}]$,

$$\begin{aligned} \text{Area}_n^{\text{Low}} &:= \sum_{i=1}^n f(y_i) \cdot \Delta y \\ &= \sum_{i=1}^n e^{(i-1)\Delta y} \cdot \Delta y \quad (\text{substituting in terms}) \\ &= \sum_{i=1}^n (e^{\Delta y})^{(i-1)} \cdot \Delta y \quad (\text{by the observation}) \\ &= \Delta y \cdot \sum_{i=1}^n (e^{\Delta y})^{(i-1)} \quad (\text{take terms that do not depend on } i \text{ outside the sum}) \\ &= \frac{x}{n} \cdot \sum_{i=0}^{n-1} (e^{\Delta y})^i \quad (\text{change of index}) \\ &= \frac{x}{n} \cdot \left(\frac{1 - r^n}{1 - r} \right) \quad (\text{Prop. 2.37}), \end{aligned} \tag{4.11}$$

where

$$r := e^{\Delta y}. \tag{4.12}$$

Hence, after noting that $r^n = e^{n\Delta y} = e^x$, we have

$$\text{Area}_n^{\text{Low}} = \frac{x}{n} \cdot \frac{1 - e^x}{1 - e^{\frac{x}{n}}} = \frac{x}{n} \cdot \frac{e^x - 1}{e^{\frac{x}{n}} - 1}. \tag{4.13}$$

Therefore,

$$\begin{aligned}
\text{Area}^{\text{Low}} &:= \lim_{n \rightarrow \infty} \text{Area}_n^{\text{Low}} \\
&= \lim_{n \rightarrow \infty} \frac{x}{n} \cdot \frac{e^x - 1}{e^{\frac{x}{n}} - 1} \\
&= (e^x - 1) \cdot \lim_{n \rightarrow \infty} \frac{x}{n} \cdot \frac{1}{e^{\frac{x}{n}} - 1} \quad (e^x - 1 \text{ does not depend on } n, \text{ so factor out of the limit}) \\
&= (e^x - 1) \cdot \lim_{h \rightarrow 0^+} h \cdot \frac{1}{e^h - 1} \quad (\text{change of variable } h := \frac{x}{n} \iff n = \frac{x}{h}) \\
&= (e^x - 1) \cdot \lim_{h \rightarrow 0^+} \frac{1}{\frac{e^h - 1}{h}} \\
&= (e^x - 1) \cdot \lim_{h \rightarrow 0^+} \frac{1}{e^h - 1} \quad (\text{by Prop. 4.19}) \\
&\quad \lim_{h \rightarrow 0^+} \frac{1}{h} \\
&= (e^x - 1) \cdot \frac{1}{1} \\
&= e^x - 1,
\end{aligned} \tag{4.14}$$

because, by Example 4.22, $\lim_{h \rightarrow 0^+} \frac{e^h - 1}{h} = 1$. Next, we repeat this process for the Riemann Upper Sum².

Riemann Upper Sum: We overestimate the total area between $[0, x]$ by summing up the overapproximations, $f(y_{i+1}) \cdot \Delta y$, for each interval $[y_i, y_{i+1}]$,

$$\begin{aligned}
\text{Area}_n^{\text{UP}} &:= \sum_{i=1}^n f(y_{i+1}) \cdot \Delta y \\
&= \sum_{i=1}^n e^{i\Delta y} \cdot \Delta y \quad (\text{substituting in terms}) \\
&= \sum_{i=1}^n (e^{\Delta y})^i \cdot \Delta y \quad (\text{by the observation}) \\
&= \Delta y \cdot \sum_{i=1}^n (e^{\Delta y})^i \quad (\text{take terms that do not depend on } i \text{ outside the sum}) \\
&= \frac{x}{n} \cdot \left(\sum_{i=0}^n (e^{\Delta y})^i - 1 \right) \quad (\text{adding and subtracting one for the geometric sum}) \\
&= \frac{x}{n} \cdot \left(\frac{1 - r^{n+1}}{1 - r} - 1 \right) \quad (\text{Prop. 2.37}) \\
&= \frac{x}{n} \cdot r \cdot \left(\frac{r^n - 1}{r - 1} \right) \quad (\text{algebra})
\end{aligned} \tag{4.15}$$

where, as before, $r = e^{\Delta y} = e^{\frac{x}{n}}$.

Hence, after noting that $r^n = (e^{\Delta y})^n = e^{n\Delta y} = e^x$, we have

$$\text{Area}_n^{\text{UP}} = \frac{x}{n} \cdot e^{\frac{x}{n}} \cdot \left(\frac{e^x - 1}{e^{\frac{x}{n}} - 1} \right). \tag{4.16}$$

²Because e^x is a continuous function, we know the Riemann integral exists. Hence, the two Riemann sums converge to a common finite value, meaning that we really only have to do one of them to compute the value of the integral. Here, we do both because it is our last opportunity to practice the definition of the Riemann integral!

Therefore,

$$\begin{aligned}
\text{Area}^{\text{Up}} &:= \lim_{n \rightarrow \infty} \text{Area}_n^{\text{Up}} \\
&= \lim_{n \rightarrow \infty} \frac{x}{n} \cdot e^{\frac{x}{n}} \cdot \left(\frac{e^x - 1}{e^{\frac{x}{n}} - 1} \right) \\
&= (e^x - 1) \cdot \lim_{n \rightarrow \infty} \frac{x}{n} \cdot e^{\frac{x}{n}} \cdot \left(\frac{1}{e^{\frac{x}{n}} - 1} \right) \\
&= (e^x - 1) \cdot \lim_{h \rightarrow 0^+} h \cdot e^h \cdot \frac{1}{e^h - 1} \quad (\text{change of variable } h := \frac{x}{n}) \\
&= (e^x - 1) \cdot \lim_{h \rightarrow 0^+} \frac{1}{\frac{e^h - 1}{h}} \cdot \lim_{h \rightarrow 0^+} e^h \quad (\text{by Prop. 4.19}) \\
&= (e^x - 1) \cdot 1 \cdot 1 \\
&= e^x - 1,
\end{aligned} \tag{4.17}$$

because,

- by Example 4.22, $\lim_{h \rightarrow 0^+} \frac{e^h - 1}{h} = 1$, and
- because e^x is continuous, $\lim_{h \rightarrow 0^+} e^h = e^{\left(\lim_{h \rightarrow 0^+} h\right)} = e^0 = 1$.

The limits of the Riemann Lower and Upper Sums both exist and are equal. Hence, our proof is complete. ■

Proposition 4.25: Integrating $e^{\alpha x}$

For $x > 0$ finite and $\alpha \in \mathbb{R}$ non-zero,

$$\int_0^x e^{\alpha y} dy = \frac{1}{\alpha} \cdot (e^{\alpha x} - 1). \tag{4.18}$$

Proof: By Prop. 3.12, Integration and Scaling,

$$\begin{aligned}
\int_0^x e^{\alpha y} dy &= \frac{1}{\alpha} \int_{\alpha \cdot 0}^{\alpha \cdot x} e^y dy \\
&= \frac{1}{\alpha} e^y \Big|_{\alpha \cdot 0}^{\alpha \cdot x} \\
&= \frac{1}{\alpha} (e^{\alpha x} - 1)
\end{aligned}$$

Example 4.26. Evaluate the limit as $n \rightarrow \infty$ of the finite sum,

$$\frac{1}{n} \sum_{k=1}^n e^{\frac{k}{n}} = \frac{\sqrt[n]{e^1} + \sqrt[n]{e^2} + \dots + \sqrt[n]{e^n}}{n}.$$

Solution: See the video, “the trick you need to know for limits”, by Michael Penn ■

Fact 4.27: Euler's Formula for Integrating Sine and Cosine

Recall that by Euler's formula,

$$\cos(\omega_0 x) = \frac{e^{i\omega_0 x} + e^{-i\omega_0 x}}{2} \quad (4.19)$$

$$\sin(\omega_0 x) = \frac{e^{i\omega_0 x} - e^{-i\omega_0 x}}{2i}.$$

While we have not proved it, Prop. 4.25 also holds for $\alpha \in \mathbb{C}$ non-zero. Hence,

$$\int_0^x e^{i\omega_0 y} dy = \frac{1}{i\omega_0} (e^{i\omega_0 x} - 1) \quad (4.20)$$

$$\int_0^x e^{-i\omega_0 y} dy = \frac{-1}{i\omega_0} (e^{-i\omega_0 x} - 1).$$

Hence, after some algebra, we obtain

$$\int_0^x \cos(\omega_0 y) dy = \frac{1}{\omega_0} \sin(\omega_0 x) \quad (4.21)$$

$$\int_0^x \sin(\omega_0 y) dy = \frac{1}{\omega_0} (1 - \cos(\omega_0 x)).$$

We will not belabor this point here because, in Chapter 7, we'll find another means to arrive at the same result. Nevertheless, it's pretty cool how things tie together.

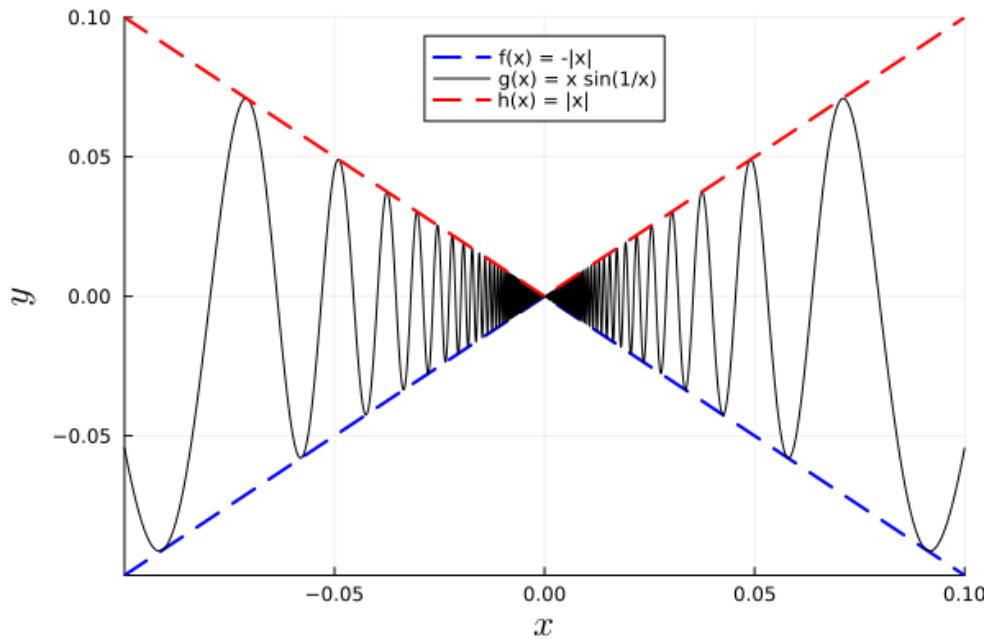


Figure 4.6: (**Squeeze or Sandwich Theorem:**) If a function is sandwiched between two functions that have a common limit, then the function itself also has the same limit.

4.4 The Squeeze Theorem

Proposition 4.28: Squeeze or Sandwich Theorem

Suppose that $a < x_0 < b$ and that f , g , and h are functions from the interval (a, b) to \mathbb{R} , except possibly being undefined at x_0 itself (the usual situation we face when computing limits). If for every $x \in (a, b)$, $x \neq x_0$, it holds that,

$$f(x) \leq g(x) \leq h(x)$$

and if

$$\lim_{x \rightarrow x_0} f(x) = \lim_{x \rightarrow x_0} h(x) = L,$$

then

$$\lim_{x \rightarrow x_0} g(x) = L.$$

Note: The result also holds for limits from the left or right at x_0 . To see this, simply redefine all three functions to the limit from the left for $x > x_0$ (or the limit from the right, for $x < x_0$). A proof is given in [Wikipedia](#), which also has some good examples.

Example 4.29. Use the Squeeze Theorem to evaluate the following limits.

$$(a) \lim_{x \rightarrow 0} x \sin\left(\frac{1}{x}\right).$$

$$(b) \lim_{x \rightarrow 0} \frac{\sin(x)}{x}.$$

$$(c) \lim_{x \rightarrow 0} \frac{1 - \cos(x)}{x}.$$

Solution:

$$(a) \text{ Ans. } \lim_{x \rightarrow 0} x \sin\left(\frac{1}{x}\right) = 0.$$

The intuition is illustrated in Fig. 4.6. For an analytical proof, because $-1 \leq \sin\left(\frac{1}{x}\right) \leq 1$ for all $x \neq 0$, we have that

$$\begin{aligned} -x &\leq x \sin\left(\frac{1}{x}\right) \leq x \quad (\text{multiply both sides by } x > 0) \\ -x &\geq x \sin\left(\frac{1}{x}\right) \geq x \quad (\text{multiply both sides by } x < 0 \text{ and note the flipped inequalities}) \\ x &\leq x \sin\left(\frac{1}{x}\right) \leq -x \quad (\text{rewrite the above line from right to left, still for } x < 0). \end{aligned}$$

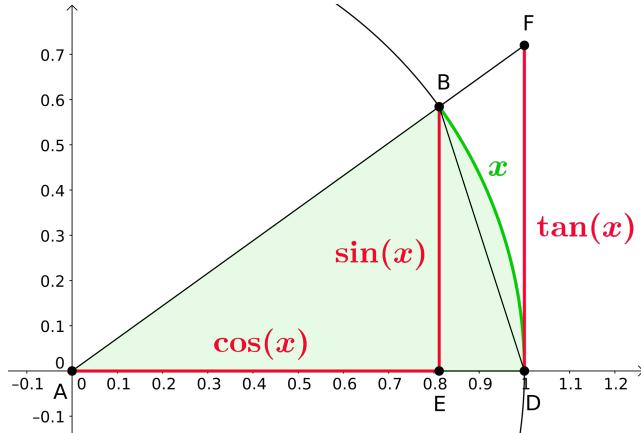
From the first and last lines, we have

$$-|x| \leq x \sin\left(\frac{1}{x}\right) \leq |x|,$$

as illustrated in Fig. 4.6. Hence, we can take $f(x) := -|x|$, $g(x) := x \sin\left(\frac{1}{x}\right)$, and $h(x) := |x|$. Because

$$\lim_{x \rightarrow 0} f(x) = \lim_{x \rightarrow 0} h(x) = 0,$$

the Squeeze Theorem gives the stated conclusion.



(a)

	$\text{Area}(\triangle ADB)$	\leq	$\text{Area}(\text{sector } ADB)$	\leq	$\text{Area}(\triangle ADF)$
\implies	$\frac{1}{2} \cdot \sin(x) \cdot 1$	\leq	$\frac{x}{2\pi} \cdot \pi$	\leq	$\frac{1}{2} \cdot \tan(x) \cdot 1$
\implies	$\sin(x)$	\leq	x	\leq	$\frac{\sin(x)}{\cos(x)}$
\implies	$\frac{\cos(x)}{\sin(x)}$	\leq	$\frac{1}{x}$	\leq	$\frac{1}{\sin(x)}$
\implies	$\cos(x)$	\leq	$\frac{\sin(x)}{x}$	\leq	1

(b)

Figure 4.7: Image and Table from [Wikipedia](#). The image shows a right triangle $\triangle ADB$, a sector of the unit circle, sector ADB having perimeter x radians, and a right triangle $\triangle ADF$. The table notes their respective areas and relative sizes. From the formulas for area, useful implications (\implies) are obtained for bounding $\sin(x)$ and $\cos(x)$.

(b) **Ans.** $\lim_{x \rightarrow 0} \frac{\sin(x)}{x} = 1$.

The hard work is done in Fig. 4.7, which shows that we can take $f(x) := \cos(x)$, $g(x) := \frac{\sin(x)}{x}$, and $h(x) := 1$. Because

$$\lim_{x \rightarrow 0} f(x) = \lim_{x \rightarrow 0} h(x) = 1,$$

the Squeeze Theorem gives the stated conclusion.

(c) **Ans.** $\lim_{x \rightarrow 0} \frac{1 - \cos(x)}{x} = 0$.

This is the most challenging problem. For $x > 0$,

$$\begin{aligned}
 \cos(x) &\leq \frac{\sin(x)}{x} \leq 1 \\
 &\Downarrow \\
 \cos(x) &\leq \frac{\sqrt{1 - \cos^2(x)}}{x} \leq 1 \quad (\text{Basic trig identity}) \\
 &\Downarrow \\
 \cos^2(x) &\leq \frac{1 - \cos^2(x)}{x^2} \leq 1 \quad (\text{Square "all" sides}) \\
 &\Downarrow \\
 x \cos^2(x) &\leq \frac{1 - \cos^2(x)}{x} \leq x \quad (\text{Multiply "all" sides by } x > 0) \\
 &\Downarrow \\
 x \cos^2(x) &\leq \frac{(1 - \cos(x))(1 + \cos(x))}{x} \leq x \quad (\text{Factor the middle term}).
 \end{aligned} \tag{4.22}$$

Applying the Squeeze Theorem to the last line of (4.22), we obtain

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

And because $\lim_{x \rightarrow 0^+} (1 + \cos(x)) = 1$, it follows that $\lim_{x \rightarrow 0^+} \frac{(1 - \cos(x))}{x} = 0$.

To complete the problem, we need to compute the limit from the left, which we leave for the learner. The key issue is that we used $x > 0$ in the fourth line of (4.22). When multiplying by $x < 0$, the inequalities in that line need to be flipped. ■

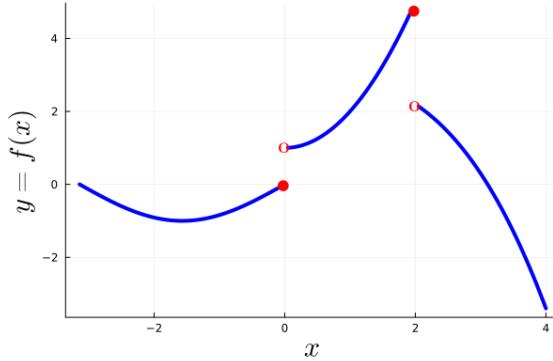


Figure 4.8: The function $f : [-\pi, 4] \rightarrow \mathbb{R}$ is an example of what we will call a **piecewise continuous function**. If we look at the function on the open interval $(-\pi, 0)$, the function is continuous and has finite one-sided limits at $-\pi$ and 0 (i.e., the two ends of the interval); the function is also continuous on the interval $(0, 2)$ and has finite one-sided limits at 0 and 2 ; finally, the function is continuous on the open interval $(2, 4)$ and has finite one-sided limits at 2 and 4 .

4.5 Piecewise Continuity

When a function is discontinuous, it may be continuous on subintervals of its domain of definition, as illustrated in Fig. 4.8. Such functions occur in Engineering when we are controlling the speed of a robot, for example. When we update a planned speed profile, it is not always convenient to make the command continuous. Allowing for jumps gives us more freedom in directing the robot.

When we studied definite integrals, we noted that a continuous function on a closed interval $[a, b]$ could always be integrated. We will define piecewise continuous functions in such a way that they can be integrated as well.

Piecewise Continuity

Definition 4.30. A real-valued function $f : [a, b] \rightarrow \mathbb{R}$, $a < b$ both finite, is **piecewise continuous** if there exists a partition of $[a, b]$,

$$a =: a_1 < a_2 < \dots < a_{n+1} := b, \quad (4.23)$$

such that for all $1 \leq i \leq n$,

- (a) $f(a_i, a_{i+1}) \rightarrow \mathbb{R}$ is continuous (i.e., f is continuous on the open interval (a_i, a_{i+1})),
- (b) $\lim_{x \rightarrow a_i^+} f(x)$ and $\lim_{x \rightarrow a_{i+1}^-} f(x)$ both exist and are finite, and
- (c) $\lim_{x \rightarrow a^+} f(x) = f(a)$ and $\lim_{x \rightarrow b^-} f(x) = f(b)$ (matches the function definition at the two endpoints).

A function $f : I \rightarrow \mathbb{R}$, for I an **unbounded interval**, is **piecewise continuous** if for all (bounded) intervals $[a, b] \subset I$, $a < b$ finite,

$$\bar{f} : [a, b] \rightarrow \mathbb{R} \text{ by } \bar{f}(x) := f(x)$$

satisfies the above conditions for piecewise continuity.

Notes:

- We've broken the interval $[a, b]$ into pieces, hence the adjective “**piecewise**”, and we impose properties on the function for each of the pieces. Since we are defining “piecewise continuity,” the condition in (a) is evidently important. What is the role of the two limits in (b)? They prevent the function from “exploding at the points a_i , $2 \leq i \leq n$ ” as we will illustrate shortly. You can take a look back at Example 4.10 to see what this means!
- What about functions defined on an open interval, $f : (a, b) \rightarrow \mathbb{R}$? If the two limits $\lim_{x \rightarrow a^+} f(x)$ and $\lim_{x \rightarrow b^-} f(x)$ exist and are finite, then simply define $f_e : [a, b] \rightarrow \mathbb{R}$ by

$$f_e(x) := \begin{cases} \lim_{x \rightarrow a^+} f(x) & x = a \\ f(x) & a < x < b \\ \lim_{x \rightarrow b^-} f(x) & x = b. \end{cases}$$

You can then apply the definition to the extended function, f_e . If one of the two one-sided limits does not exist or is infinite, then the function is automatically not piecewise continuous.

- Some books do not include the requirement that the function has finite limits at the “boundary points”, a_i . We will always require that the two limits exist and are finite, but in other contexts, you have to be careful.
- If $f : [a, b] \rightarrow \mathbb{R}$ is continuous, then it is also piecewise continuous. The converse is evidently false, as exemplified in Fig. 4.8.
- There is no such thing as “piecewise continuous at a point.”

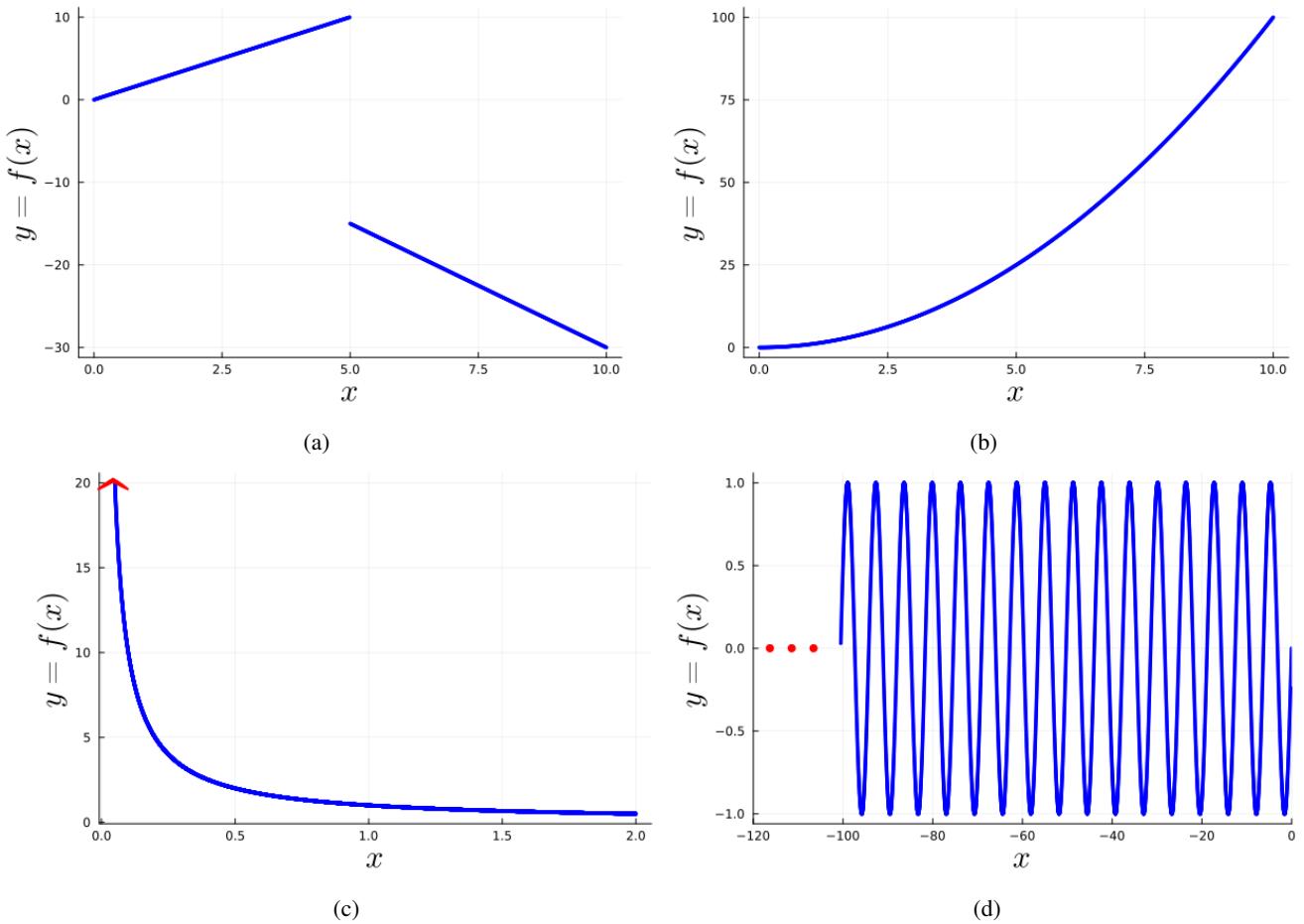


Figure 4.9: Graphs for Example 4.31. In (c), the graph heads to ∞ as $x \rightarrow 0^+$. In (d), the graph extends infinitely to the left, oscillating between -1 and $+1$ as sine waves like to do! The graphs in (a) and (b) show all the values of the functions.

Example 4.31. Determine which, if any, of the following functions are piecewise continuous.

$$(a) f : [0, 10] \rightarrow \mathbb{R} \text{ by } f(x) = \begin{cases} 2x & 0 \leq x < 5 \\ -3x & 5 \leq x \leq 10 \end{cases}.$$

$$(b) f : [0, 10] \rightarrow \mathbb{R} \text{ by } f(x) = x^2.$$

$$(c) f : (0, 2] \rightarrow \mathbb{R} \text{ by } f(x) = \frac{1}{x}.$$

$$(d) f : (-\infty, 0) \rightarrow \mathbb{R} \text{ by } f(x) = \sin(x).$$

Solutions: Plots are provided in Fig. 4.9; analytical reasoning is provided below.

(a) **Ans.** Piecewise continuous.

We break the interval $[0, 10]$ into two pieces, $A_1 := (0, 5)$ and $A_2 := (5, 10)$ with

$$f_1 : (0, 5) \rightarrow \mathbb{R}, \text{ by } f_1(x) = 2x \text{ and } f_2 : (5, 10) \rightarrow \mathbb{R}, \text{ by } f_2(x) = -3x.$$

Each of these functions is continuous and

$$\lim_{x \rightarrow 0^+} f_1(x) = 0 = f(0), \lim_{x \rightarrow 5^-} f_1(x) = 10, \lim_{x \rightarrow 5^+} f_2(x) = -15, \text{ and } \lim_{x \rightarrow 10^-} f_2(x) = -30 = f(10).$$

Hence, all of the conditions are met.

(b) **Ans.** Piecewise continuous.

$f(x) = x^2$ has a finite left limit at $x = 10$ and it is equal to $(10)^2$. Therefore, we have $f_e : [0, 10] = x^2$, which we know is continuous everywhere, and hence it is piecewise continuous.

(c) **Ans.** Not piecewise continuous because $\lim_{x \rightarrow 0^+} \frac{1}{x} = \infty$. Our definition of piecewise continuity requires finite limits.

(d) **Ans.** Piecewise continuous, because

- the function has a finite left limit at zero, namely, $\lim_{x \rightarrow 0^-} \sin(x) = 0$, and
- the function is continuous everywhere for all bounded intervals $[a, 0)$, for $a < 0$ finite.

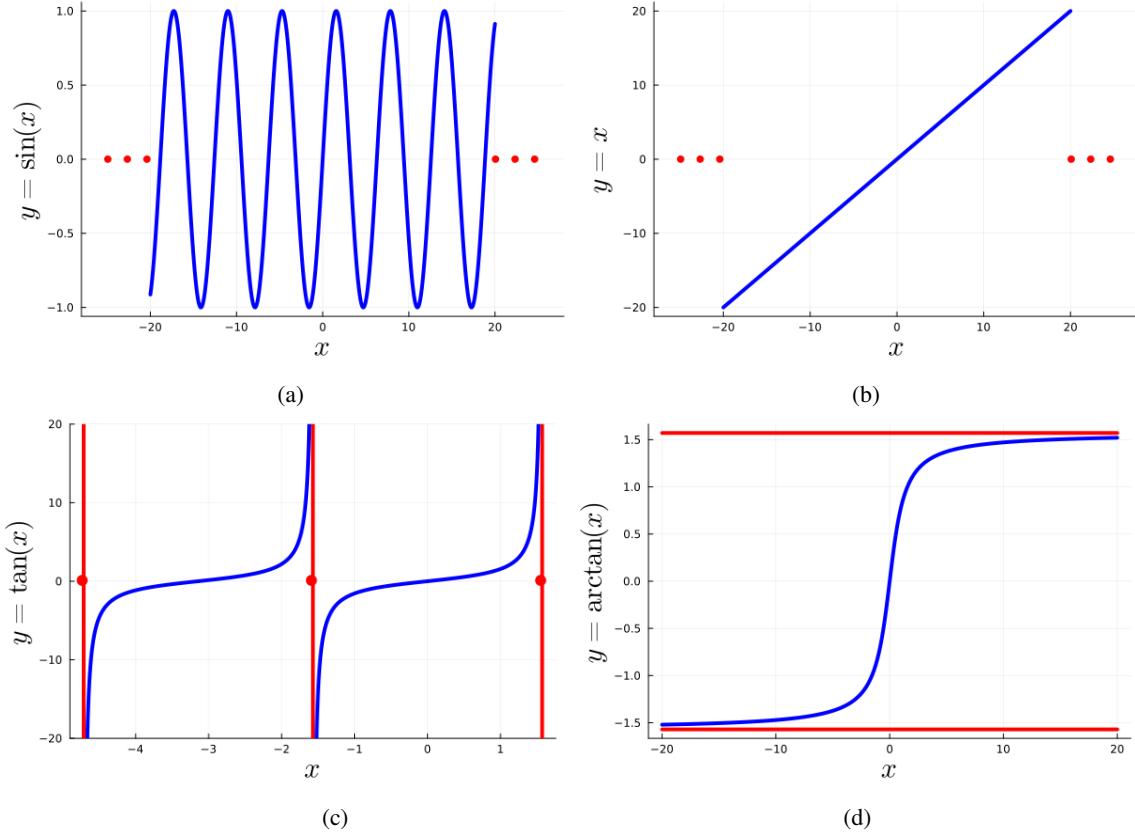


Figure 4.10: Functions can be bounded or unbounded, have horizontal or vertical asymptotes, or not! (a) Bounded, with no asymptotes of any kind. (b) Represents $f(x) = x$ and is unbounded with no asymptotes of any kind (c) Unbounded with vertical asymptotes at odd multiples of $\frac{\pi}{2}$ indicated in red. (d) Bounded with horizontal asymptotes at $\pm\frac{\pi}{2}$ indicated in red

4.6 Bounded vs Unbounded Functions and Asymptotes

Figure 4.10 illustrates further characteristics of functions that, when present, often aid in their analysis or, in the case of functions of time (aka, trajectories), aid in the analysis of the system that produces them. Bounded functions, such as the sine function, are confined within a fixed range, offering predictability and stability, qualities essential in fields like engineering and physics where maintaining variables within limits is crucial. On the other hand, asymptotes describe the behavior of functions as they approach specific values or infinity, as seen in the tangent and arc-tangent functions. These concepts are not just abstract mathematical ideas

but are pivotal in understanding growth rates, stability, and constraints in real-world applications ranging from economic models to biological processes.

Bounded Functions

Definition 4.32. Let $A \subset \mathbb{R}$ be a subset of \mathbb{R} or \mathbb{Z} , and $f : A \rightarrow \mathbb{R}$ a function.

- (a) f is **bounded from above** if there exists $M < \infty$ such that $f(x) \leq M$ for all $x \in A$.
- (b) f is **bounded from below** if there exists $M > -\infty$ such that $f(x) \geq M$ for all $x \in A$.
- (c) f is **bounded** if it is both bounded from above and from below. This is equivalent to there exists $M < \infty$ such that $|f(x)| \leq M$ for all $x \in A$.
- (d) f is **unbounded** if it is **not bounded** (that is, either not bounded from above or not bounded from below). This is equivalent to for all $M < \infty$, there $x \in A$, such that $|f(x)| \geq M$. In particular, a function that is bounded from above but not from below is unbounded, and similarly, a function that is bounded from below but not from above is unbounded.
- (e) Recall that when $A \subset \mathbb{N}$, then $a_n := f(n)$ defines a sequence $\{a_n\}_{n \in A}$.

Illustration: See Fig. 4.10-(a) and (b).

Notes: The notion of “boundedness” appears in:

- **Numerical Methods:** Bounded functions are easier to approximate numerically. For example, numerical integration techniques often rely on the function being bounded to guarantee a certain level of accuracy.
- **Stability:** In control theory and engineering, bounded functions are often desirable because they indicate a system’s stability. A bounded output in response to a bounded input usually signifies a stable system.
- **Optimization:** Bounded functions are often easier to optimize because you can be certain that the optimum lies within a specific range. This is particularly useful in machine learning algorithms.
- **Convergence:** In mathematics and computer science, bounded functions are important for analyzing the convergence of sequences and series.

Before we dive into examples, it is worthwhile to understand the related notion of **asymptotes**.

Asymptotes come in two Flavors: Horizontal and Vertical

Definition 4.33. Let $A \subset \mathbb{R}$ be a subset of \mathbb{R} and $f : A \rightarrow \mathbb{R}$ a function. A **horizontal asymptote** is a horizontal line $y = b$ where the function f approaches b as x approaches positive infinity (∞) or negative infinity ($-\infty$). In other symbols, a horizontal line $y = b$ is a **horizontal asymptote** if at least one of the following two conditions holds:

1. $\lim_{x \rightarrow \infty} f(x) = b$;
2. $\lim_{x \rightarrow -\infty} f(x) = b$.

Illustration: See Fig. 4.10-(c) and (d).

Definition 4.34. Let $A \subset \mathbb{R}$ be a subset of \mathbb{R} and $f : A \rightarrow \mathbb{R}$ a function. A **vertical asymptote** is a vertical line at $x = x_0$ where the function f approaches $\pm\infty$ as x approaches x_0 from the left or the right. In other symbols, a vertical line $\{(x_0, y) \mid y \in \mathbb{R}\}$ is a **vertical asymptote** if at least one of the following two conditions holds:

1. $\lim_{x \rightarrow x_0^+} f(x) = \pm\infty$ (the limit exists and is unbounded);

2. $\lim_{x \rightarrow x_0} f(x) = \pm\infty$ (the limit exists and is unbounded).

Note: At least one of the limits needs to exist for x_0 to be a vertical asymptote. As an example, the function $f(x) := \frac{\sin(\frac{1}{x})}{x^2}$ “blows up” as x approaches the origin from the left or from the right, but neither limit exists. Hence, $x_0 = 0$ is not a vertical asymptote.

Notes: The notion of “asymptotes” is useful because:

- **Behavior at Infinity:** Asymptotes help us understand the behavior of functions as points in their domain approach infinity or a finite point. This is crucial in calculus for evaluating limits and integrals.
- **Graphical Interpretation:** Asymptotes provide a way to graphically represent the behavior of functions, making it easier to visualize and understand their properties.
- **Modeling:** In physics, engineering, and economics, asymptotic behavior can model saturation effects. For example, the law of diminishing returns in economics and saturation current in electrical engineering can be modeled using functions with horizontal asymptotes.
- **Simplification:** Asymptotic approximations can often simplify complex functions into more manageable forms, making it easier to perform calculations or gain qualitative insights.

As illustrated in Fig. 4.10, a function can have at most two horizontal asymptotes, while it can have multiple (even infinite in number) vertical asymptotes or none at all.

Example 4.35. Determine which of the following functions are bounded versus unbounded. Determine also if they have any asymptotes or not.

(a) $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = e^x$.

(b) $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = x^2$.

(c) $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = \frac{.001x^2+x}{1+x^2}$.

(d) $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = \begin{cases} \frac{1}{(x-1)(x+1)} & |x| \neq 1 \\ 0 & \text{otherwise} \end{cases}$.

(e) $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = \begin{cases} \frac{x-1}{x^2-1} & |x| \neq 1 \\ 0 & \text{otherwise} \end{cases}$. **Hint:** Compare to the previous problem.

(f) $f : \mathbb{N} \rightarrow \mathbb{R}$ by $f(n) = (1 - 1/n)^{-n}$. **Hint:** See Example 4.21-(a).

(g) $f : \mathbb{N} \rightarrow \mathbb{R}$ by $f(n) = \sum_{k=0}^n \frac{1}{k!}$. **Hint:** See Chapter 1.18.1.

Solutions:

(a) **Ans.** $f(x) = e^x$ is unbounded. It has a horizontal asymptote at $y = 0$, corresponding to $\lim_{x \rightarrow -\infty} e^x = 0$.

(b) **Ans.** $f(x) = x^2$ is unbounded and has no asymptotes.

(c) **Ans.** $f(x) = \frac{-0.001x^2+x}{1+x^2}$ is bounded and has a horizontal asymptote at $y = -10^{-3}$, corresponding to $\lim_{x \rightarrow \pm\infty} f(x) = -10^{-3}$.

- (d) **Ans.** $f(x) = \begin{cases} \frac{1}{(x-1)(x+1)} & |x| \neq 1 \\ 0 & \text{otherwise} \end{cases}$ is unbounded with two vertical asymptotes. One is at $x_0 = -1$, corresponding to $\lim_{x \rightarrow -1^-} f(x) = +\infty$ and $\lim_{x \rightarrow -1^+} f(x) = -\infty$; the other is at $x_0 = +1$, corresponding to $\lim_{x \rightarrow 1^-} f(x) = -\infty$ and $\lim_{x \rightarrow 1^+} f(x) = +\infty$. If you plot the function, the answer jumps out at you. Here, we numerically approximate the limits.

```

1 f (x) = 1 / (x2-1)
2 del= 1e-5
3 X0 = [-1, 1]
4 Data = Array{Float64} (undef, 0, 2)
5 for k = 1:2
6   x=X0[k]-del
7   Data = [Data; x f(x) ]
8   x=X0[k]+del
9   Data = [Data; x f(x) ]
10 end
11 Data
```

Output The first two rows are for $x_0 = -1$ and the last two are for $x_0 = 1$.

4×2 Matrix{Float64} :

-1.00001	49999.8
-0.99999	-50000.3
0.99999	-50000.3
1.00001	49999.8

- (e) **Ans.** $f(x) = \begin{cases} \frac{x-1}{x^2-1} & |x| \neq 1 \\ 0 & \text{otherwise} \end{cases}$ has a vertical asymptote at $x_0 = -1$, but does not have an additional vertical asymptote at $x_0 = 1$ because the “apparent singularity” at +1 is canceled by the term in the numerator. You can look at the problem in two ways:

First: Because $x^2 - 1 = (x - 1)(x + 1)$, $\frac{x-1}{x^2-1} = \frac{x-1}{(x-1)(x+1)} = \frac{1}{x+1}$.

Second: Numerically approximate the limits as in the previous problem and see what you get!

```

1 f (x) = (x-1) / (x2-1)
2 del= 1e-5
3 X0 = [-1, 1]
4 Data = Array{Float64} (undef, 0, 2)
5 for k = 1:2
6   x=X0[k]-del
7   Data = [Data; x f(x) ]
8   x=X0[k]+del
9   Data = [Data; x f(x) ]
10 end
11 Data
```

Output The first two rows are for $x_0 = -1$ and show that the right and left limits diverge to $\pm\infty$; the last two are for $x_0 = 1$, and indicate that the left and right limits converge to a half.

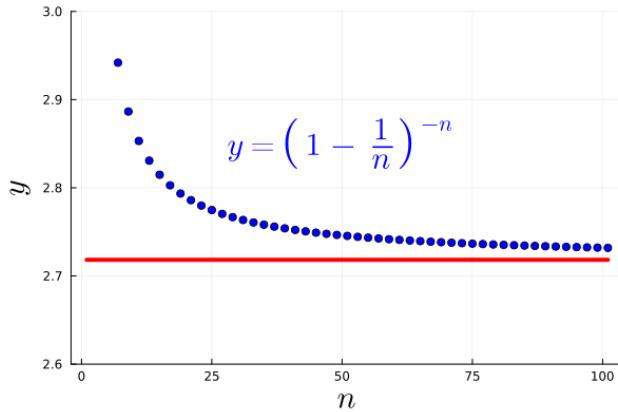
4×2 Matrix{Float64} :

-1.00001	-100000.0
-0.99999	100000.0
0.99999	0.500003
1.00001	0.499998

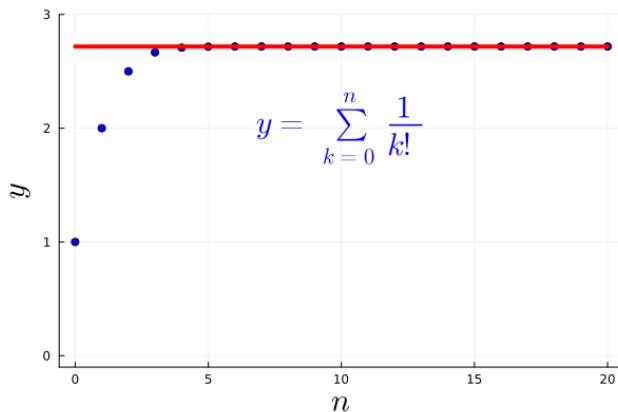
Hence, $\lim_{x \rightarrow 1^-} f(x) = 0.5$ and $\lim_{x \rightarrow 1^+} f(x) = 0.5$.

Both methods are valid.

(f) **Ans.** $f : \mathbb{N} \rightarrow \mathbb{R}$ by $f(n) = (1 - 1/n)^{-n}$ has a horizontal asymptote at e because $\lim_{n \rightarrow \infty} (1 - 1/n)^{-n} = e$, Euler's number. We computed the limit in Example 4.21-(a). We illustrate the limit in the following plot.



(g) **Ans.** $f : \mathbb{N} \rightarrow \mathbb{R}$ by $f(n) = \sum_{k=0}^n \frac{1}{k!}$ has a horizontal asymptote at e . We illustrate this via a plot. Note how fast the summation approaches e . Euler really was smart!



4.7 When does a Function have a Definite Riemann Integral?

We developed the Riemann Integral in Chapter 3 as the signed area under a curve. Our original development was for continuous functions $f : [a, b] \rightarrow \mathbb{R}$, where $[a, b]$ is a closed bounded interval. We now, finally, understand what is a continuous function. We also now understand other properties of functions, such as piecewise continuous functions, bounded functions, and functions that have vertical asymptotes. How do these properties relate to the Riemann Integral?

Proposition 4.36: Handy List of Functions Having a Riemann Integral

Functions Commonly Encountered by Engineers:

- Continuous Functions:** All continuous functions on a closed and bounded interval $[a, b]$ are Riemann integrable. This is the most straightforward class of functions that are Riemann integrable. Riemann integrability may fail for continuous functions on an open and bounded interval (a, b) or half-open intervals. It is very important that $a < b$ are finite real numbers and the function is continuous on $[a, b]$.

- **Piecewise Continuous Functions:** A function that is continuous on $[a, b]$ except for a finite number of discontinuities, while having finite left and right limits at the discontinuities, is also Riemann integrable. If the one-sided limits at points of discontinuity either do not exist or are unbounded, then the result can be a function that is not Riemann integrable. Hence, the presence of vertical asymptotes, while not a deal breaker, signals that caution is needed.
- **Bounded Functions with a Finite Number of Discontinuities (not so common, but interesting):** More generally, any function that is bounded on $[a, b]$ and has a finite number of discontinuities is Riemann integrable. The function doesn't have to be continuous or even piecewise continuous. It does not have to possess one-sided limits at the discontinuities. This result requires graduate-level mathematics to prove. *Mostly, when an engineer encounters such functions, they are piecewise continuous. An example that is not piecewise continuous would be $f(0) := 0$, and for $0 < x < 1$, $f(x) := \sin(1/x)$.* The function $f : [0, 1] \rightarrow [-1, 1]$ is bounded between ± 1 and has a single discontinuity, namely at $x_0 = 0$. Amazingly, it is Riemann integrable.
- **Composite Functions:** If $f : [a, b] \rightarrow \mathbb{R}$ and $g : [a, b] \rightarrow \mathbb{R}$ are both Riemann integrable, then so is $f(g(x))$, provided $f(x)$ is continuous. The function g does not have to be continuous or piecewise continuous. It could be one of the more exotic functions mentioned elsewhere in this list.

Functions $f : (a, b) \rightarrow \mathbb{R}$ with Vertical Asymptotes or Functions Defined on Unbounded Intervals:

- Such integrals are called **Improper Integrals**, and we treat them in Chapter 8.
- Sometimes, we can make sense of them, and sometimes not, as indicated above.
- Interesting examples include $\int_0^\infty e^{-x} dx$ (finite), $\int_0^\infty x^k \cdot e^{-x} dx$ (finite), and $\int_0^1 \frac{1}{x} dx$ (infinite), all of which are encountered frequently.

More Exotic Functions, Beyond the Scope of this Course:

- **Monotonic Functions:** Functions that are monotonic (either entirely non-increasing or entirely non-decreasing) on a closed and bounded interval $[a, b]$ are Riemann integrable. **Such functions could have a countably infinite number of discontinuities and still be Riemann integrable.** For open bounded intervals, (a, b) , half-open bounded intervals $[a, b)$ and $(a, b]$, or unbounded intervals, it is NOT true that all monotonic functions are guaranteed to be Riemann integrable.
- **Functions Equivalent to Continuous Functions:** If a function $f : [a, b] \rightarrow \mathbb{R}$ is (i) bounded and (ii) its points of discontinuity have **Lebesgue measure** zero, then f is Riemann integrable. Moreover, if (iii), f differs from a continuous function $g : [a, b] \rightarrow \mathbb{R}$ on a set of **Lebesgue measure** zero, then $f(x)$ is Riemann integrable, and its integral is the same as that of $g(x)$. The Lebesgue measure is treated in graduate-level courses on Functional Analysis or Measure Theory.
- **Functions with Bounded Variation:** Functions of **bounded variation** on a closed and bounded interval are **Riemann-Stieltjes** integrable with respect to any other function of bounded variation. It is noted that the **Lebesgue integral** generalizes the Riemann integral and can integrate a broader class of functions. However, every Riemann integrable function is also Lebesgue integrable, and the two integrals will agree on the value for such functions.

4.8 Maximum and Minimum Values of Sets and Their Generalizations

This is the first time we come close to the “nitty gritty” of the real numbers!

Max and Min Values of a Set

Definition 4.37. Let $A \subset \mathbb{R}$ be a nonempty subset of the reals.

- (a) a^* is a **maximum element** of the set A if (i) $a^* \in A$ and (ii) $a^* \geq a$ for all $a \in A$. The notation is

$$a^* = \max\{A\}.$$

(b) a_* is a **minimum element** of the set A if (i) $a_* \in A$ and (ii) $a_* \leq a$ for all $a \in A$. The notation is

$$a_* = \min\{A\}.$$

Notes:

- A maximum or minimum element must be in the set.
- A finite set always has both maximum and minimum elements. Here, “finite set” means a finite number of elements.
- A set with an infinite number of elements may not have a maximum or minimum element.

Example 4.38. Indicate the maximum and minimum elements of the following sets, if they exist.

- (a) $A := [3, 4]$.
- (b) $B := (3, 4]$.
- (c) $C := [3, 4)$.
- (d) $D := (3, 4)$.
- (e) $E := (0, 1) \cap \mathbb{Q}$, the set of rational numbers in the open interval $(0, 1)$.
- (f) $F := \{x \in [0, 1] \mid x \notin \mathbb{Q}\}$, the set of irrational numbers in the closed interval $[0, 1]$.

Solutions:

- (a) **Ans.** $\min\{A\} = 3$ and $\max\{A\} = 4$.

While we take this as obvious, you might try writing out a proof by appealing to the definition. It will be good practice. Many proofs in 100- and 200-level math courses boil down to applying the definition to the data at hand.

- (b) **Ans.** $\min\{B\}$ does not exist while $\max\{B\} = 4$.

To prove that the set B does not have a minimum element, let $x \in B$ be arbitrary and define $\delta := x - 3$. Then $\delta > 0$ because $3 < x \leq 4$. Hence, we have

$$3 < 3 + \delta/2 < x$$

and thus $3 + \delta/2$ is both smaller than x and an element of B ; in symbols, $3 + \delta/2 < x$ and $(3 + \delta/2) \in B$. Hence, x cannot be a minimum, and because x was arbitrary, no element of B can be a minimum.

- (c) **Ans.** $\max\{C\}$ does not exist while $\min\{C\} = 3$.

To prove that the set C does not have a maximum element, let $x \in C$ be arbitrary and define $\delta := 4 - x$. Then $\delta > 0$ because $4 \notin C$. Hence, we have

$$x < x + \delta/2 < 4$$

and thus $x + \delta/2 > x$ and $(x + \delta/2) \in C$. Hence, x cannot be a maximum, and because x was arbitrary, no element of C can be a maximum.

- (d) **Ans.** Neither $\min\{D\}$ nor $\max\{D\}$ exist. The reasoning is the same as in the previous two parts of this Example.

- (e) **Ans.** Neither $\min\{E\}$ nor $\max\{E\}$ exist.

Let $r \in E$ be an arbitrary rational number in the open interval $(0, 1)$ and define $\delta := 1 - r$. Then $\delta > 0$ because $1 \notin E$ and δ is rational because it is the difference of two rational numbers. Then we have

$$r < r + \delta/2 < 1$$

and thus $r + \delta/2 > r$ and $(r + \delta/2) \in E$ because $r + \delta/2$ is also rational. Hence, r cannot be a maximum. Similar reasoning shows that r cannot be a minimum.

- (f) **Ans.** Neither $\min\{F\}$ nor $\max\{F\}$ exist. The reasoning is the same as in the previous part of this Example, because, if $x \in F$ is irrational, then $1 - x$, $1 - x/2$, and $x/2$ are all irrational numbers belonging to F . ■

Sup and Inf Values of a Set

Definition 4.39. Let $A \subset \mathbb{R}$ be a nonempty subset of the reals.

- (a) a^* is the **least upper bound** (aka, **supremum**) of the set A if (i) $a^* \in \mathbb{R}$, (ii) $a^* \geq a$ for all $a \in A$, and (iii) if $b \in \mathbb{R}$ satisfies $b \geq a$ for all $a \in A$, then $a^* \leq b$. The notation is

$$a^* = \sup\{A\},$$

where “sup” is short for “supremum”. Note that (ii) says that a^* is an upper bound and (iii) says that, among all upper bounds, it is the least upper bound.

- (b) a_* is the **greatest lower bound** (aka, **infimum**) of the set A if (i) $a_* \in \mathbb{R}$, (ii) $a_* \leq a$ for all $a \in A$, and (iii) if $b \in \mathbb{R}$ satisfies $b \leq a$ for all $a \in A$, then $a_* \geq b$. The notation is

$$a_* = \inf\{A\},$$

where “inf” is short for “infimum”. Note that (ii) says that a_* is a lower bound and (iii) says that, among all lower bounds, it is the greatest lower bound.

Remark 4.40. Properties of the infimum and supremum.

- The key difference between a maximum or minimum element and the corresponding supremum or infimum element is where they live! $\max\{A\}$ and $\min\{A\}$ must belong to A , whereas $\sup\{A\}$ and $\inf\{A\}$ belong to \mathbb{R} . In fact,

$$\sup\{\mathbf{A}\} = \max\{\mathbf{A}\} \iff \sup\{\mathbf{A}\} \in \mathbf{A} \text{ and } \inf\{\mathbf{A}\} = \min\{\mathbf{A}\} \iff \inf\{\mathbf{A}\} \in \mathbf{A}.$$

- Every bounded set of real numbers has a supremum and an infimum, whereas we noticed in the previous example many bounded sets of real numbers do not have a maximum or a minimum.
- The existence of the supremum and the infimum for any bounded set of the reals is actually part of the definition of the real numbers. Michigan’s Math 451 digs more deeply into this topic.
- If a set is unbounded from above, such as $[0, \infty)$, it is common to define the supremum as ∞ . Similarly, if a set is unbounded from below, it is common to define the infimum as $-\infty$. In this way, all subsets of the reals have an infimum and a supremum.

Example 4.41. Find the supremum and infimum for the following sets, if they exist. If a supremum or infimum is also a maximum or minimum, indicate that as well.

- $A := [3, 4]$.
- $B := (3, 4]$.
- $C := [3, 4)$.
- $D := (3, 4)$.
- $E := (0, 1) \cap \mathbb{Q}$, the set of rational numbers in the open interval $(0, 1)$.
- $F := \{x \in \mathbb{R} \mid x \notin \mathbb{Q}\}$, the set of irrational numbers.

Solutions:

- (a) **Ans.** $\inf\{A\} = 3$ and $\sup\{A\} = 4$. Because $3 = \inf\{A\} \in A$, it is also a minimum. Similarly, because $4 = \sup\{A\} \in A$, it is also a maximum.
- (b) **Ans.** $\inf\{B\} = 3$ and $\sup\{B\} = 4$. Because $3 = \inf\{B\} \notin B$, it is NOT a minimum. However, because $4 = \sup\{B\} \in B$, it is also a maximum.
- (c) **Ans.** $\inf\{C\} = 3$ and $\sup\{C\} = 4$. Because $3 = \inf\{C\} \in C$, it is also a minimum. However, because $4 = \sup\{C\} \notin C$, it is NOT a maximum.
- (d) **Ans.** $\inf\{D\} = 3$ and $\sup\{D\} = 4$. Because $3 = \inf\{D\} \notin D$, it is NOT a minimum. Similarly, because $4 = \sup\{D\} \notin D$, it is NOT a maximum.
- (e) **Ans.** $\inf\{E\} = 0$ and $\sup\{E\} = 1$. Because $0 = \inf\{E\} \notin E$, it is NOT a minimum. Similarly, because $1 = \sup\{E\} \notin E$, it is NOT a maximum.
- (f) **Ans.** $\inf\{F\} = -\infty$ and $\sup\{F\} = \infty$. Because $-\infty = \inf\{F\} \notin F$, it is NOT a minimum. Similarly, because $\infty = \sup\{F\} \notin F$, it is NOT a maximum. Recall that $\pm\infty$ are concepts, not real numbers.

■

4.9 Maximum and Minimum Values of Functions and Their Generalizations

In some sense, this material should be super easy: it is a repeat of the previous definitions of max, min, sup, and inf applied to the range of a function! In practice, it seems harder to grasp for most learners. Perhaps this is because the range of a function can be very complicated? In any case, we'll change things up and first present the sup and inf, before moving on to max and min.

Extremal Values of a Function (or a Function's Range)

Definition 4.42. Let $f : A \rightarrow \mathbb{R}$ be a function, with its domain A being nonempty. Define $Y := f(A) := \{y \in \mathbb{R} \mid y = f(x) \text{ for some } x \in A\}$ to be the **range** of f .

- (a) $f^* := \sup\{Y\}$ is the **supremum** of the range of f , though it is most commonly called the **supremum** of f . The notation is

$$f^* = \sup_{x \in A} f(x), \text{ or } f^* = \sup \{ \{y = f(x) \mid x \in A\} \},$$

which provides good insight: f^* is the least upper bound of $f(x)$ as x varies over the entire domain, A .

- (b) $f_* := \inf\{Y\}$ is the **infimum** of the range of f , though it is most commonly called the **infimum** of f . The notation is

$$f_* = \inf_{x \in A} f(x), \text{ or } f^* = \inf \{ \{y = f(x) \mid x \in A\} \},$$

which provides good insight: f_* is the greatest lower bound of $f(x)$ as x varies over the entire domain, A .

Recalling that the key difference between sup and max is “where they live”, we define the max and min of a function as follows.

- (c) $f^* := \sup\{Y\}$ is the **maximum of f** if $f^* \in Y$, that is, there exists $\bar{a} \in A$ such that $f(\bar{a}) = f^*$. The notation is

$$f^* = \max_{x \in A} f(x),$$

which provides good insight: f^* is the maximum element of $f(x)$ as x varies over the entire domain, A .

- (d) $f_* := \inf\{Y\}$ is the **minimum of f** if $f_* \in Y$, that is, there exists $\underline{a} \in A$ such that $f(\underline{a}) = f_*$. The notation is

$$f_* = \min_{x \in A} f(x),$$

which provides good insight: f_* is the minimum element of $f(x)$ as x varies over the entire domain, A .

Notes:

- For the maximum (resp., minimum) of a function to exist, there must be an element in the domain that achieves the maximum (resp., minimum) value of the function. This is not the case of the supremum (resp., infimum).
- $\sup_{x \in A} f(x)$ and $\inf_{x \in A} f(x)$ always exist, while $\max_{x \in A} f(x)$ and $\min_{x \in A} f(x)$ sometimes exist, and sometimes do not.

Proposition 4.43: Extreme Value Theorem due to Bolzano-1830

Let $f : [a, b] \rightarrow \mathbb{R}$ be a **continuous function on a closed and bounded interval** $[a, b]$. Then f always achieves its extreme values (aka, maximum and minimum values). In other symbols, there always exist $x^* \in [a, b]$ and $x_* \in [a, b]$ such that

$$f(x^*) = \max_{a \in A} f(x) \text{ and } f(x_*) = \min_{a \in A} f(x).$$

Notes: The proof is given in Math 451 and ROB 501; it uses really cool techniques.

Example 4.44. For the following functions, determine their sup (resp., inf) values, and if they exist, also give the max (resp. min) values.

(a) $f : (-\infty, \infty) \rightarrow \mathbb{R}$ by $f(x) = e^{-x}$.

(b) $f : [0, 10] \rightarrow \mathbb{R}$ by $f(x) = \frac{e^x}{x^2 + 1}$.

(c) $f : (-\infty, \infty) \rightarrow \mathbb{R}$ by $f(x) = e^{-x^2} \sin(x)$

(d) $f : [0, \infty) \rightarrow \mathbb{R}$ by $f(x) = \frac{2x+1}{x^2+1}$.

(e) $f : [0, 1] \rightarrow \mathbb{R}$ by $f(x) = \frac{\tan(x)-x^3}{x^2+2x \sin(x)+10}$.

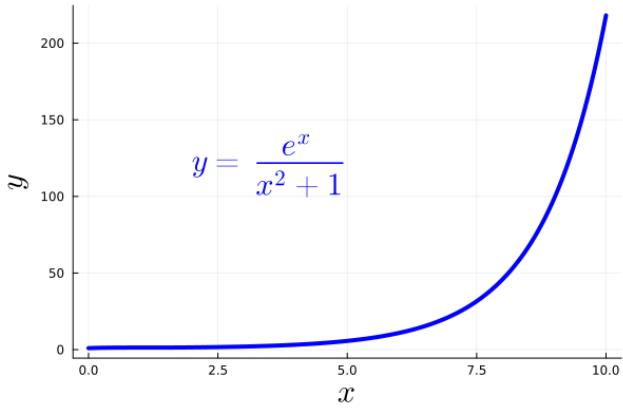
Solutions:

(a) **Ans.** $\inf_{-\infty < x < \infty} e^{-x} = 0$ and $\sup_{-\infty < x < \infty} e^{-x} = \infty$. The min and max of f do not exist.

e^{-x} is strictly monotonically decreasing. Hence, $\inf_{-\infty < x < \infty} e^{-x} = \lim_{x \rightarrow \infty} e^{-x} = 0$ and $\sup_{-\infty < x < \infty} e^{-x} = \lim_{x \rightarrow -\infty} e^{-x} = \infty$. You can verify this with a plot, if you wish!

(b) **Ans.** $\min_{x \in [0, 10]} \frac{e^x}{x^2 + 1} = 1$ and $\max_{x \in [0, 10]} \frac{e^x}{x^2 + 1} = 218.1$.

Because the denominator never vanishes, $f : [0, 10] \rightarrow \mathbb{R}$ is a continuous function on the closed and bounded interval, $[0, 10]$. Hence, the max and min values do exist. If the function is strictly monotonically increasing, then the extreme values will be at the endpoints of the interval. We do not know enough Calculus at the moment to check whether or not the function is monotonic, so we'll determine the extreme values numerically. The code below also provides values of x that achieve the max and min of f .



```

1 xmin = 0
2 xmax = 10
3 xVals=xmin:.01:xmax
4
5 f(x) = exp(x) / (1 + x^2)
6
7 y = f.(xVals)
8
9 p1 = plot(xVals, y, linewidth=4, color=:blue, label=false)
10 plot!( xlabel=L"$x$",
11          ylabel=L"$y=\frac{e^x}{x^2+1}$")
12
13 @show MaxF = maximum(y)
14 @show MinF = minimum(y)
15
16 max_value, index = findmax(y)
17 println ("Maximum value: ", max_value)
18 println ("x*: ", xVals[index])
19
20 min_value, index = findmin(y)
21 println ("Minimum value: ", min_value)
22 println ("x_*: ", xVals[index])
23
24 display(p1)
25
26 png(p1, "MaxMinB")

```

Output

```

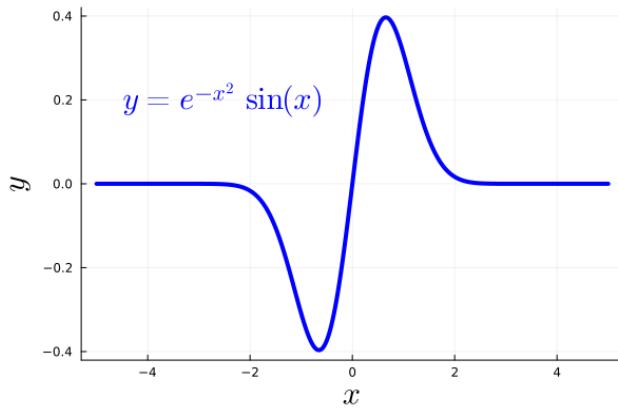
MaxF = maximum(y) = 218.08381975056156
MinF = minimum(y) = 1.0
Maximum value: 218.08381975056156
x*: 10.0
Minimum value: 1.0
x_*: 0.0

```

(c) **Ans.** $\min_{x \in \mathbb{R}} e^{-x^2} \sin(x) = -0.397$ and $\max_{x \in \mathbb{R}} e^{-x^2} \sin(x) = 0.397$.

Because the domain of the function is unbounded, Prop. 4.43 does not apply. **However, the function is continuous everywhere on its domain, and it is easy to see that the extreme values do not occur at $\pm\infty$; in fact, they occur within the closed bounded set $[-2, 2]$, where the Proposition does apply.** How did we see all of that? Well, $\lim_{x \rightarrow \pm\infty} e^{-x^2} \sin(x) = 0$. Yet,

we know that e^{-x^2} is always positive, while $\sin(x)$ takes on positive and negative values. Hence, $f(x)$ takes on positive and negative values, and, therefore, the infimum will be negative, and the supremum will be positive. We then need to graph the function to find a range where the extreme values of the function will occur.



```

1 xmin = -5
2 xmax = 5
3 xVals=xmin:.01:xmax
4
5 f(x) = exp(-x^2)*sin(x)
6
7 y = f.(xVals)
8
9 p1 = plot(xVals,y, linewidth=4, color=:blue, label=false)
10 plot!( xlabel=L"$x$",
11 ylabel=L"$y=e^{-x^2} \sim \sin(x)$")
12 @show MaxF = maximum(y)
13 @show MinF = minimum(y)
14
15 max_value, index = findmax(y)
16 println("Maximum value: ", max_value)
17 println("x^*: ", xVals[index])
18
19 min_value, index = findmin(y)
20 println("Minimum value: ", min_value)
21 println("x_*: ", xVals[index])
22
23 display(p1)
24
25 png(p1, "MaxMinC")

```

Output

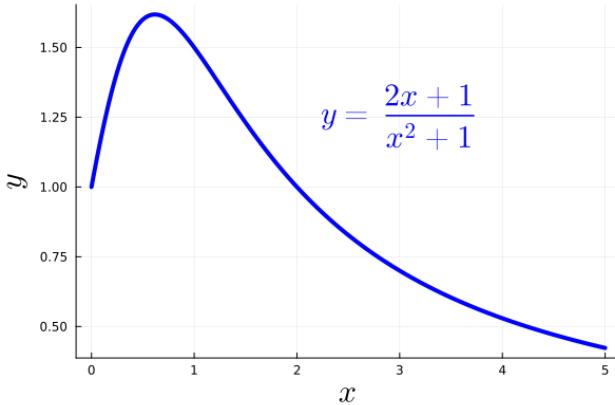
```

MaxF = maximum(y) = 0.3966429553529812
MinF = minimum(y) = -0.3966429553529812
Maximum value: 0.3966429553529812
x^*: 0.65
Minimum value: -0.3966429553529812
x_*: -0.65

```

$$(d) \text{ Ans. } \min_{x \in \mathbb{R}} \frac{2x+1}{x^2+1} = 0.423 \text{ and } \max_{x \in \mathbb{R}} \frac{2x+1}{x^2+1} = 1.618.$$

Because the domain of the function is unbounded, Prop. 4.43 does not apply. **However, the function is continuous everywhere on its domain, and it is easy to see that the extreme values do not occur at ∞ ; in fact, they occur within the closed bounded set $[0, 5]$, where the Proposition does apply.**



```

1 xmin = 0
2 xmax = 5
3 xVals=xmin:.01:xmax
4
5 f(x) = (2x + 1) / (x^2 + 1)
6
7 y = f.(xVals)
8
9 p1 = plot(xVals, y, linewidth=4, color=:blue, label=false)
10 plot!(xlabel=L"$x$",
11 ylabel=L"$y=\frac{2x+1}{x^2+1}$")
12 @show MaxF = maximum(y)
13 @show MinF = minimum(y)
14
15 max_value, index = findmax(y)
16
17 println("Maximum value: ", max_value)
18 println("x^*: ", xVals[index])
19
20 display(p1)
21
22 png(p1, "MaxMinD")

```

Output

```

MaxF = maximum(y) = 1.6180294712510837
MinF = minimum(y) = 0.4230769230769231
Maximum value: 1.6180294712510837
x^*: 0.62

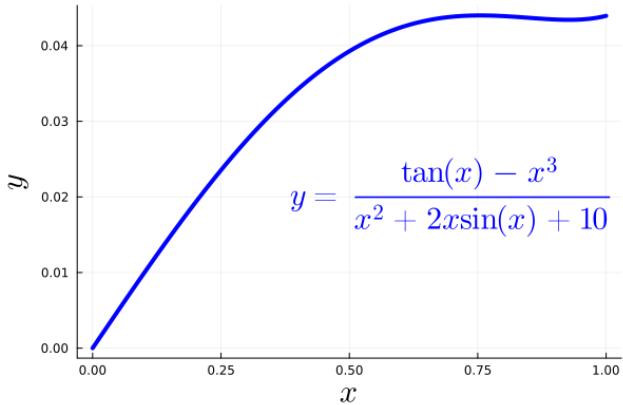
```

$$(e) \text{ Ans. } \min_{x \in [0,1]} \frac{\tan(x) - x^3}{x^2 + 2x \sin(x) + 10} = 0 \text{ and } \max_{x \in [0,1]} \frac{\tan(x) - x^3}{x^2 + 2x \sin(x) + 10} = 0.044.$$

Because the denominator never vanishes and $1 < \pi/2$, $f : [0, 1] \rightarrow \mathbb{R}$ is a continuous function on the closed and bounded interval, $[0, 1]$. Hence, the max and min values do exist. How do we know the denominator does not vanish? We complete the square,

$$x^2 + 2x \sin(x) + 10 = x^2 + 2x \sin(x) + 10 + \sin^2(x) - \sin^2(x) = (x + \sin(x))^2 + 10 - \sin^2(x),$$

and we note that, for all x , $10 - \sin^2(x) \geq 9$ and $(x + \sin(x))^2 \geq 0$.



```

1 xmin = 0
2 xmax = 1
3 xVals=xmin:.001:xmax
4
5 f(x) = (tan(x) - x^3) / (x^2 + 2x * sin(x) + 10)
6
7 y = f.(xVals)
8
9 p1 = plot(xVals,y, linewidth=4, color=:blue, label=false)
10 plot!( xlabel=L"$x$",
11 ylabel=L"$y=\frac{\tan(x) - x^3}{x^2 + 2x\sin(x) + 10}$")
12 @show MaxF = maximum(y)
13 @show MinF = minimum(y)
14
15 max_value, index = findmax(y)
16
17 println("Maximum value: ", max_value)
18 println("x^*: ", xVals[index])
19
20 display(p1)
21
22 png(p1, "MaxMinE")

```

Output

```

MaxF = maximum(y) = 0.043999918166140344
MinF = minimum(y) = 0.0
Maximum value: 0.043999918166140344
x^*: 0.755

```

■

4.10 (Optional Read:) Intermediate Value Theorem and the Mean Value Theorem

The Intermediate Value Theorem is a fundamental result in calculus and real analysis. It asserts that for any continuous function defined on a closed interval, the function takes on every value between its minimum and maximum on that interval. This theorem has important implications in various areas of mathematics and its applications. We used it to establish the soundness of the Bisection Algorithm.

Proposition 4.45: Intermediate Value Theorem

Let $f : [a, b] \rightarrow \mathbb{R}$ be a continuous function on the closed interval $[a, b]$, where $a < b$. If d is any number between $f(a)$ and $f(b)$, then there exists at least one $c \in [a, b]$ such that $f(c) = d$.

Proof: Without loss of generality, assume $f(a) < d < f(b)$. Define the set $S = \{x \in [a, b] \mid f(x) \leq d\}$.

Because $f(a) \leq d$, we have that $a \in S$, and hence the set S is non-empty. Moreover, the set S is bounded above by b because $S \subset [a, b]$. Finally, by the completeness property of the real numbers, S has a supremum. Define

$$c := \sup S.$$

We now want to show that $f(c) = d$. Because $c = \sup S$ and $c \in [a, b]$, c is the least upper bound of S . Therefore, for any $\epsilon > 0$, there exists an $x \in S$ such that $c - \epsilon < x \leq c$.

- Since $x \in S$, $f(x) \leq d$. Given f is continuous, we have $\lim_{x \rightarrow c^-} f(x) = f(c) \leq d$.
- By definition of c being the supremum, for any $\epsilon > 0$, there exists a point $y \in [c, c + \epsilon]$ such that $y \notin S$. Thus, $f(y) > d$.
- Since f is continuous at c , we have $\lim_{y \rightarrow c^+} f(y) = f(c) \geq d$.

Combining these limits, we obtain $f(c) = d$, and hence, there exists at least one $c \in [a, b]$ such that $f(c) = d$. ■

The Mean Value Theorem is a central result in calculus and real analysis. It states that for any function continuous on a closed interval and differentiable on the open interval, there exists at least one point in the open interval where the derivative of the function equals the average rate of change over the closed interval.

Proposition 4.46: Mean Value Theorem

Let $f : [a, b] \rightarrow \mathbb{R}$ be a continuous function on the closed interval $[a, b]$ and differentiable on the open interval (a, b) , where $a < b$. Then, there exists at least one $c \in (a, b)$ such that

$$f'(c) = \frac{f(b) - f(a)}{b - a}.$$

Proof: Define the function $g(x) = f(x) - \frac{f(b) - f(a)}{b - a}(x - a)$. Note that g is continuous on $[a, b]$ and differentiable on (a, b) .

We first observe that:

$$g(a) = f(a) - \frac{f(b) - f(a)}{b - a}(a - a) = f(a),$$

and

$$g(b) = f(b) - \frac{f(b) - f(a)}{b - a}(b - a) = f(b) - (f(b) - f(a)) = f(a).$$

Thus, $g(a) = g(b)$. By [Rolle's Theorem](#), since g is continuous on $[a, b]$, differentiable on (a, b) , and $g(a) = g(b)$, there exists at least one $c \in (a, b)$ such that $g'(c) = 0$. Next, we compute $g'(x)$:

$$g'(x) = f'(x) - \frac{f(b) - f(a)}{b - a}.$$

At $x = c$, we have $g'(c) = 0$, which gives us:

$$f'(c) - \frac{f(b) - f(a)}{b - a} = 0,$$

or equivalently,

$$f'(c) = \frac{f(b) - f(a)}{b - a}.$$

Thus, there exists at least one $c \in (a, b)$ such that

$$f'(c) = \frac{f(b) - f(a)}{b - a}.$$

Related Videos:

- [Intermediate Value Theorem](#) by The Organic Chemistry Tutor
- [Mean Value Theorem](#) by @MindSphereYT.

4.11 (Optional Read:) Proofs Associated with the Chapter

We first state and prove a few lemmas, and then we get into the proofs.

Lemma (Product of Two Continuous Functions is Continuous): Let $f : [a, b] \rightarrow \mathbb{R}$ and $g : [a, b] \rightarrow \mathbb{R}$ be continuous functions at a point $x_0 \in [a, b]$. Then, the function $h(x) := f(x)g(x)$ is also continuous at x_0 .

Proof: To prove that $h(x)$ is continuous at x_0 , we need to show that for every $\epsilon > 0$, there exists a $\delta > 0$ such that if $|x - x_0| < \delta$, then $|h(x) - h(x_0)| < \epsilon$.

Because f and g are continuous at x_0 , for any $\epsilon_1 > 0$, there exists $\delta_1 > 0$ such that if $|x - x_0| < \delta_1$, then $|f(x) - f(x_0)| < \epsilon_1$. Similarly, for any $\epsilon_2 > 0$, there exists $\delta_2 > 0$ such that if $|x - x_0| < \delta_2$, then $|g(x) - g(x_0)| < \epsilon_2$.

Let $\epsilon > 0$. Choose $\epsilon_1 = \frac{\epsilon}{2(|g(x_0)| + 1)}$ and $\epsilon_2 = \frac{\epsilon}{2(|f(x_0)| + 1)}$. By the continuity of f and g , there exist δ_1 and δ_2 corresponding to ϵ_1 and ϵ_2 , respectively. Let $\delta = \min\{\delta_1, \delta_2\}$. Then, for $|x - x_0| < \delta$, we have,

$$\begin{aligned}|f(x)g(x) - f(x_0)g(x_0)| &= |f(x)g(x) - f(x)g(x_0) + f(x)g(x_0) - f(x_0)g(x_0)| \\&\leq |f(x)||g(x) - g(x_0)| + |g(x_0)||f(x) - f(x_0)| \\&< |f(x)|\epsilon_2 + |g(x_0)|\epsilon_1 \\&\leq (|f(x_0)| + 1)\epsilon_2 + (|g(x_0)| + 1)\epsilon_1 \\&= \frac{\epsilon}{2} + \frac{\epsilon}{2} \\&= \epsilon.\end{aligned}$$

Therefore, $h(x) = f(x)g(x)$ is continuous at x_0 . ■

Lemma (Sum of Two Continuous Functions is Continuous): Let $f : [a, b] \rightarrow \mathbb{R}$ and $g : [a, b] \rightarrow \mathbb{R}$ be continuous functions on the interval $[a, b]$. Then, the function $h(x) := f(x) + g(x)$ is also continuous on $[a, b]$.

Proof: To prove that $h(x)$ is continuous on $[a, b]$, we need to show that for every point $x_0 \in [a, b]$ and for every $\epsilon > 0$, there exists a $\delta > 0$ such that if $|x - x_0| < \delta$, then $|h(x) - h(x_0)| < \epsilon$.

Since f and g are continuous at x_0 , for any $\epsilon > 0$, there exist $\delta_f > 0$ and $\delta_g > 0$ such that:

- If $|x - x_0| < \delta_f$, then $|f(x) - f(x_0)| < \frac{\epsilon}{2}$.
- If $|x - x_0| < \delta_g$, then $|g(x) - g(x_0)| < \frac{\epsilon}{2}$.

Let $\delta = \min\{\delta_f, \delta_g\}$. Now, for $|x - x_0| < \delta$, we have by the triangle inequality,

$$\begin{aligned} |h(x) - h(x_0)| &= |(f(x) + g(x)) - (f(x_0) + g(x_0))| \\ &= |(f(x) - f(x_0)) + (g(x) - g(x_0))| \\ &\leq |f(x) - f(x_0)| + |g(x) - g(x_0)| \\ &< \frac{\epsilon}{2} + \frac{\epsilon}{2} \\ &= \epsilon. \end{aligned}$$

Thus, $h(x) = f(x) + g(x)$ is continuous at x_0 . Since x_0 was arbitrary, $h(x)$ is continuous on the entire interval $[a, b]$. ■.

Lemma (Composition of Two Continuous Functions is Continuous): Let $f : X \rightarrow Y$ and $g : Y \rightarrow Z$ be continuous functions, where X, Y , and Z are subsets of \mathbb{R} . Then, the composition $g \circ f : X \rightarrow Z$, defined by $(g \circ f)(x) := g(f(x))$, is continuous.

Proof:

To show that $g \circ f$ is continuous at any point $x_0 \in X$, we need to verify that for every $\epsilon > 0$, there exists a $\delta > 0$ such that if $|x - x_0| < \delta$, then $|(g \circ f)(x) - (g \circ f)(x_0)| < \epsilon$.

Because g is continuous at $f(x_0) \in Y$, for every $\epsilon > 0$, there exists an $\eta > 0$ such that if $|y - f(x_0)| < \eta$, then $|g(y) - g(f(x_0))| < \epsilon$, where $y = f(x)$. Furthermore, since f is continuous at x_0 , for the $\eta > 0$ found above, there exists a $\delta > 0$ such that if $|x - x_0| < \delta$, then $|f(x) - f(x_0)| < \eta$.

Combining these two observations, if $|x - x_0| < \delta$, then $|f(x) - f(x_0)| < \eta$, which implies $|g(f(x)) - g(f(x_0))| < \epsilon$. This means that $|(g \circ f)(x) - (g \circ f)(x_0)| < \epsilon$. Therefore, by the ϵ - δ definition of continuity, $g \circ f$ is continuous at x_0 . Since x_0 was arbitrary, $g \circ f$ is continuous on X . ■

We leave unproven that all three of these results can be extended to $n > 2$ functions. You may wish to prove at least one of them for yourself. Each of the extensions can be proven by induction.

Proof of Prop. 4.18 (Common Continuous Functions)

We'll prove a few of the dozen or more functions listed in Prop. 4.18: the following functions are continuous for all $x \in \mathbb{R}$ unless indicated otherwise:

- (a) Monomials: $f(x) = x^k$
- (b) Polynomials: $f(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0$
- (c) Sine: $f(x) = \sin(x)$
- (d) Natural Exponential: $f(x) = e^x$
- (e) Natural Logarithm (continuous for $x > 0$): $f(x) = \ln(x)$
- (f) Power (continuous for $x > 0$ and $y \in \mathbb{R}$): $f(x) = x^y$
- (g) Square Root (continuous for $x \geq 0$): $f(x) = \sqrt{x}$

- (a) Monomials: $f(x) = x^k$.

We prove by induction that the function $f(x) = x^k$ is continuous for all positive integers k .

Base Cases:

- For $k = 0$, $f(x) = x^0 = 1$, a constant function, which is continuous everywhere.
- For $k = 1$, $f(x) = x$, the identity function, which is also continuous everywhere.

Inductive Step: Assume $f(x) = x^k$ is continuous. We need to show that $f(x) = x^{k+1}$ is also continuous. Given $f(x) = x^k$ is continuous by our inductive hypothesis, and knowing that $g(x) = x$ is continuous, consider $h(x) := x^{k+1} = x^k \cdot x$. Since the product of two continuous functions is continuous, it follows that x^{k+1} is continuous.

Therefore, by mathematical induction, $f(x) = x^k$ is continuous for all positive integers k . ■

- (b) Polynomials: $f(x) = a_n x^n + a_{n-1} x^{n-1} + \cdots + a_1 x + a_0$

This follows immediately by the above result for monomials and the lemma on the sum of continuous functions.

- (c) Sine: $f(x) = \sin(x)$

To prove continuity at $x_0 \in \mathbb{R}$, we need to show that for every $\epsilon > 0$, there exists a $\delta > 0$ such that if $|x - x_0| < \delta$, then $|\sin(x) - \sin(x_0)| < \epsilon$.

Using the trigonometric identity for the difference of two sines,

$$|\sin(x) - \sin(x_0)| = |2 \cos\left(\frac{x+x_0}{2}\right) \sin\left(\frac{x-x_0}{2}\right)|.$$

Because $|\cos(\theta)| \leq 1$ for all θ , and $|\sin(\theta)| \leq |\theta|$ for all θ , we can further simplify the inequality to:

$$|\sin(x) - \sin(x_0)| \leq |2 \sin\left(\frac{x-x_0}{2}\right)| \leq |x - x_0|.$$

To ensure $|\sin(x) - \sin(x_0)| < \epsilon$, it suffices to choose $\delta = \epsilon$. Thus, for $|x - x_0| < \delta$, we indeed have $|\sin(x) - \sin(x_0)| < \epsilon$.

Therefore, by the ϵ - δ criterion of continuity, $f(x) = \sin(x)$ is continuous at x_0 . ■

- (d) Natural Exponential: $f(x) = e^x$

Let x_0 be an arbitrary point in \mathbb{R} . To prove that $f(x) = e^x$ is continuous at x_0 , consider any $\epsilon > 0$. We aim to show that there exists a $\delta > 0$ such that if $|x - x_0| < \delta$, then $|e^x - e^{x_0}| < \epsilon$.

Given $\epsilon > 0$, define $\theta > 1$ such that

$$e^{x_0} - \epsilon < \frac{e^{x_0}}{\theta} < e^{x_0} < \theta e^{x_0} < e^{x_0} + \epsilon.$$

It follows from the definition of θ that

$$\frac{e^{x_0}}{\theta} < e^{x_0} < \theta e^{x_0},$$

which simplifies to

$$-\ln(\theta) < x - x_0 < \ln(\theta).$$

Let $\delta = \ln(\theta)$. Then, for $x \in (x_0 - \delta, x_0 + \delta)$, we have

$$e^{x_0 - \delta} < e^x < e^{x_0 + \delta}.$$

This implies

$$\frac{e^{x_0}}{\theta} < e^x < \theta e^{x_0},$$

or equivalently,

$$e^{x_0} - \epsilon < e^x < e^{x_0} + \epsilon.$$

Therefore, $|e^x - e^{x_0}| < \epsilon$ whenever $|x - x_0| < \delta$, proving that $f(x) = e^x$ is continuous at x_0 . Since x_0 was arbitrary, $f(x) = e^x$ is continuous for all $x \in \mathbb{R}$. ■

(e) Natural Logarithm (continuous for $x > 0$): $f(x) = \ln(x)$

Let $x_0 > 0$ be an arbitrary point. To prove that $f(x) = \ln(x)$ is continuous at x_0 , consider any $\varepsilon > 0$. We must show that there exists a $\delta > 0$ such that if $|x - x_0| < \delta$, then $|\ln(x) - \ln(x_0)| < \varepsilon$.

Given $\varepsilon > 0$, consider the intervals $(e^{-\varepsilon}x_0, e^{\varepsilon}x_0)$. For x in this interval, we have $e^{-\varepsilon} < \frac{x}{x_0} < e^{\varepsilon}$, implying that $-\varepsilon < \ln\left(\frac{x}{x_0}\right) < \varepsilon$, which simplifies to $-\varepsilon < \ln(x) - \ln(x_0) < \varepsilon$, or equivalently, $|\ln(x) - \ln(x_0)| < \varepsilon$.

Let $\delta = \min\{x_0 - e^{-\varepsilon}x_0, e^{\varepsilon}x_0 - x_0\}$. Then, for $|x - x_0| < \delta$, x lies within the interval $(e^{-\varepsilon}x_0, e^{\varepsilon}x_0)$, ensuring that $|\ln(x) - \ln(x_0)| < \varepsilon$. Therefore, by the ε - δ definition of continuity, $f(x) = \ln(x)$ is continuous at x_0 . Since x_0 was arbitrary and $x_0 > 0$, $f(x) = \ln(x)$ is continuous for all $x_0 > 0$. ■

(f) Power (continuous for $x > 0$ and $y \in \mathbb{R}$): $f(x) = x^y$

For $x > 0$ and logarithm properties, we have that $x^y = e^{y \ln x}$, which is the composition of the continuous functions, e^x and $y \ln x$. ■

(g) Square Root (continuous for $x \geq 0$): $f(x) = \sqrt{x}$

This is a special case of the power function, with $y = \frac{1}{2}$. ■

Proof of Prop. 4.19 (Taking Limits Inside a Function)

Consider two functions $g : X \rightarrow Y$ and $f : Y \rightarrow Z$, where X , Y , and Z are subsets of \mathbb{R} . Suppose that f is continuous at a point $y_0 \in Y$. Then, the following hold:

- (a) If $x_0 \in \mathbb{R}$ and $\lim_{x \rightarrow x_0} g(x) = y_0$ (i.e., the limit exists and equals y_0), then $\lim_{x \rightarrow x_0} f(g(x)) = f(\lim_{x \rightarrow x_0} g(x)) = f(y_0)$. **We say the limit can be taken “inside” the function f .**
- (b) Similarly, if $\lim_{x \rightarrow \infty} g(x) = y_0$, then $\lim_{x \rightarrow \infty} f(g(x)) = f(\lim_{x \rightarrow \infty} g(x)) = f(y_0)$.
- (c) If $\lim_{x \rightarrow -\infty} g(x) = y_0$, then $\lim_{x \rightarrow -\infty} f(g(x)) = f(\lim_{x \rightarrow -\infty} g(x)) = f(y_0)$.

In short, if $g(x)$ has a limit (as x approaches a finite point or as x approaches $\pm\infty$), and f is continuous at the limit of g (i.e., y_0), then, when evaluating a limit of $f(g(x)) = f \circ g(x)$, the limit can be taken inside the function f and applied directly to g .

Proof:

- (a) If $x_0 \in \mathbb{R}$ and $\lim_{x \rightarrow x_0} g(x) = y_0$ (i.e., the limit exists and equals y_0), then $\lim_{x \rightarrow x_0} f(g(x)) = f(\lim_{x \rightarrow x_0} g(x)) = f(y_0)$. **We say the limit can be taken “inside” the function f .**
- (b) Similarly, if $\lim_{x \rightarrow \infty} g(x) = y_0$, then $\lim_{x \rightarrow \infty} f(g(x)) = f(\lim_{x \rightarrow \infty} g(x)) = f(y_0)$.
- (c) If $\lim_{x \rightarrow -\infty} g(x) = y_0$, then $\lim_{x \rightarrow -\infty} f(g(x)) = f(\lim_{x \rightarrow -\infty} g(x)) = f(y_0)$.

Proof:

- (a) Given that $\lim_{x \rightarrow x_0} g(x) = y_0$, we know that for every $\varepsilon > 0$ there exists a $\delta > 0$ such that if $0 < |x - x_0| < \delta$, then $|g(x) - y_0| < \varepsilon$. Since f is continuous at y_0 , for this ε , there exists a $\delta_f > 0$ such that $|y - y_0| < \delta_f$ implies $|f(y) - f(y_0)| < \varepsilon$ for all y in the domain of f .

Since $g(x)$ approaches y_0 as x approaches x_0 , and f maps values near y_0 to values near $f(y_0)$, it follows that as x approaches x_0 , $f(g(x))$ approaches $f(y_0)$. Hence, $\lim_{x \rightarrow x_0} f(g(x)) = f(y_0)$.

- (b) The proofs for parts (b) and (c) follow a similar reasoning to part (a), but considering the limits as x approaches ∞ and $-\infty$, respectively. The key is the continuous behavior of f at y_0 and the fact that $g(x)$ approaches y_0 in each scenario. By the continuity of f , the image under f of the values of $g(x)$ that approach y_0 will approach $f(y_0)$, thus allowing the limit to be taken inside the function f .

- (c) See above.

In conclusion, if $g(x)$ converges to a point y_0 as x approaches a finite value or $\pm\infty$, and f is continuous at y_0 , the composition $f \circ g(x)$ behaves such that the limit operation can be directly applied after the composition, yielding $f(\lim_{x \rightarrow x_0} g(x)) = f(y_0)$, where x_0 represents a finite value, ∞ , or $-\infty$. ■

Proof of Prop. 4.36 (Handy List of Functions Having a Riemann Integral)

Functions Commonly Encountered by Engineers:

- (a) **Continuous Functions:** All continuous functions on a closed and bounded interval $[a, b]$ are Riemann integrable. This is the most straightforward class of functions that are Riemann integrable. Riemann integrability may fail for continuous functions on an open and bounded interval (a, b) or half-open intervals. It is very important that $a < b$ are finite real numbers and the function is continuous on $[a, b]$.
- (b) **Piecewise Continuous Functions:** A function that is continuous on $[a, b]$ except for a finite number of discontinuities, while having finite left and right limits at the discontinuities, is also Riemann integrable. If the one-sided limits at points of discontinuity either do not exist or are unbounded, then the result can be a function that is not Riemann integrable.

Proofs: We realize this is far from the entire list given in Prop. 4.36, but even these two are on the edge of what we can handle in the course.

- (a) **Continuous Functions:** Because f is continuous on $[a, b]$, it is also **uniformly continuous** on this interval due to the interval being closed and bounded. This means that for any $\epsilon > 0$, there exists a $\delta > 0$ such that for all $x, y \in [a, b]$, if $|x - y| < \delta$, then $|f(x) - f(y)| < \frac{\epsilon}{b-a}$. We have not covered this “flavor” of continuity in the textbook.

Uniformly partition the interval $[a, b]$ into n subintervals $\{[x_i, x_{i+1}]\}_{i=1}^n$, where $x_1 = a$ and $x_{n+1} = b$, and each subinterval has a width equal to $\Delta x := \frac{b-a}{n} \leq \delta$. This ensures that the supremum and infimum of f in each subinterval,

$$h_i^{\text{Up}} := \sup_{x \in [x_i, x_{i+1}]} f(x)$$

$$h_i^{\text{Low}} := \inf_{x \in [x_i, x_{i+1}]} f(x)$$

satisfy

$$\omega_i := h_i^{\text{up}} - h_i^{\text{low}} \leq \frac{\epsilon}{b-a}.$$

The difference between the Riemann upper sum and lower sum for the partition can be expressed as the sum of the “oscillations”, ω_i , over all subintervals, weighted by the width of each subinterval:

$$\sum_{i=1}^n \omega_i \Delta x \leq \sum_{i=1}^n \frac{\epsilon}{b-a} \cdot \frac{b-a}{n} = \sum_{i=1}^n \frac{\epsilon}{n} = \epsilon.$$

This shows that the difference between the upper and lower Riemann sums for f over $[a, b]$ can be made less than any given $\epsilon > 0$, which proves that f is Riemann integrable on $[a, b]$.

- (b) **Piecewise Continuous Functions:** By the way we have defined piecewise continuous functions, we can write them as a finite sum of functions that are continuous on disjoint sub-intervals of the form (a_i, b_i) , have bounded limits at the endpoints, and are zero everywhere else. Next, the proof of Prop. 3.4 (First Additivity Property) can be adapted to cover any jumps at the endpoints of the intervals. ■

Chapter 5

Differentiation

Learning Objectives

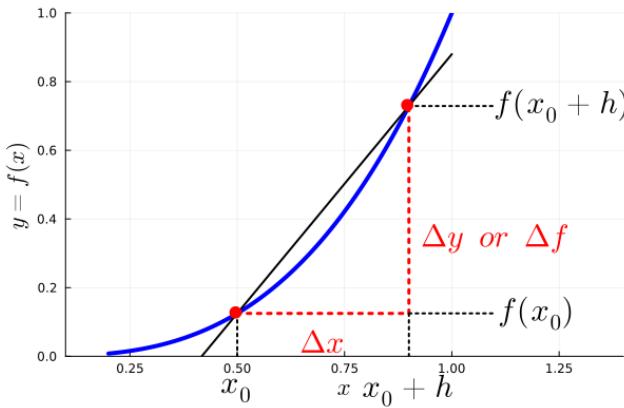
By the end of this chapter, the student should be able to:

- Understand two conceptual views of a single-variable derivative
- Appreciate the effectiveness of various software tools for computing derivatives
- See real problems where single-variable derivatives are important in engineering.
- Make the leap to partial derivatives, which are single-variable derivatives applied to multivariable functions.

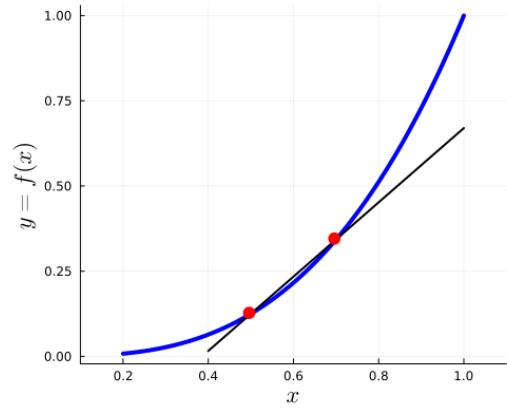
Outcomes

Upon successful completion of this chapter, students will be able to:

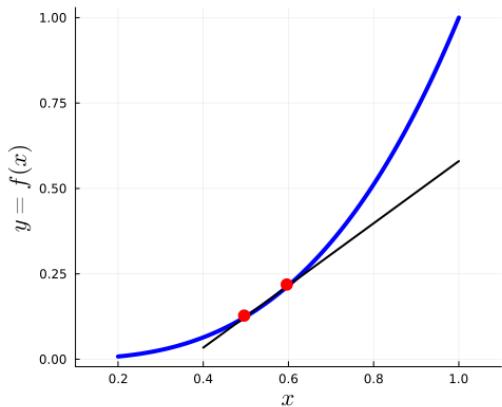
- Apply the definition of the derivative as rise over run in the limit.
- Understand the centrality of the derivative for linear approximation of a function.
- Compute a few derivatives via the rise over run definition.
- Obtain a sense of when common functions are differentiable and when they are not.
- Learn how to compute derivatives with various software packages.
- Master the Rules of Differentiation and Understand their Origin.
- Apply single-variable derivatives to determine speed from position.
- Apply the fact that a strictly positive derivative implies the function is strictly monotonically increasing.
- Use L'Hôpital's Rule for limits of indeterminate form.
- Compute Jacobians, gradients, and Hessians, with examples using software.
- Understand and use the total derivative.



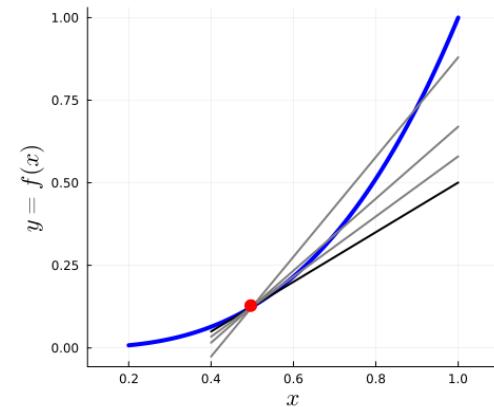
(a)



(b)



(c)



(d)

Figure 5.1: The derivative of $f : (a, b) \rightarrow \mathbb{R}$ at $x_0 \in (a, b)$ is the limiting value of the “rise over run,” as $h \rightarrow 0$. (a) Graphical definition of all the terms in (5.1). The “rise” is $\Delta f := f(x_0 + h) - f(x_0)$, and the “run” is $\Delta x := (x_0 + h) - x_0 = h$; the slope, of course, is their ratio, $\frac{\Delta f}{\Delta x}$. A line passing through the pair of points $(x_0, f(x_0))$ and $(x_0 + h, f(x_0 + h))$ is called a **secant line**, where the term “secant” is Latin for “cutting” (it cuts the function at the two indicated points, as opposed to a “tangent line”, which touches the function at only one point). (b) and (c) show that as the run, $\Delta x := (x_0 + h) - x_0 = h$, gets smaller and smaller, so does the rise, $\Delta f := f(x_0 + h) - f(x_0)$. (d) Shows the “limiting case” where the secant line becomes a tangent line (indicated by the dark line), the slope of which is the value of the derivative at x_0 . The lighter lines are the previous secant lines so that we can see the convergence process playing out as $h \rightarrow 0$.

5.1 The Derivative as the Local Slope of a Function

Figure 5.1 illustrates the notion of the **derivative of a function at a point** as the **limiting slope**, or **rise over run**, as the run approaches zero (from both sides).

Derivative at a Point

Definition 5.1. *The derivative of $f : (a, b) \rightarrow \mathbb{R}$ at a point $x_0 \in (a, b)$ is defined as the limiting value of the rise over run,*

$$\frac{df(x_0)}{dx} := \lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{(x_0 + h) - x_0} = \lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h}. \quad (5.1)$$

Notes: Recall that the double-sided limit $\lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h}$ is equivalent to two single-sided limits,

$$\lim_{h \rightarrow 0^-} \frac{f(x_0 + h) - f(x_0)}{h} \text{ and } \lim_{h \rightarrow 0^+} \frac{f(x_0 + h) - f(x_0)}{h}.$$

Both have to exist, be finite, and equal for the single-sided limit to exist. To show that a function is **not differentiable at x_0** , it is enough to show any one of the following:

- The two one-sided limits both exist and are finite, but they have different values.
- One of the limits is infinite.
- One of the limits does not exist (aka, is undefined).

If $f : [a, b] \rightarrow \mathbb{R}$ is defined on a closed interval, the derivatives at the two endpoints can be defined via single-sided limits, namely

$$\begin{aligned}\frac{df(a)}{dx} &:= \lim_{h \rightarrow 0^+} \frac{f(x_0 + h) - f(x_0)}{h} \\ \frac{df(b)}{dx} &:= \lim_{h \rightarrow 0^-} \frac{f(x_0 + h) - f(x_0)}{h}.\end{aligned}\tag{5.2}$$

The derivative at points $a < x_0 < b$ is given by (5.1).

Example 5.2. Explain why the function $f(x) := |x|$ is not differentiable at $x_0 = 0$.

Solution: The key is to compute both single-sided limits at the origin. To make this particularly easy, we note that

$$f(x) := |x| = \begin{cases} x & x > 0 \\ -x & x < 0 \\ 0 & x = 0. \end{cases}$$

Hence,

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

Both limits exist and are finite, but they are not equal. Hence, $|x|$ is not differentiable at the origin. You can check that it is differentiable at all other points. ■

Example 5.3. Is the function

$$f(x) := \begin{cases} 1 + x & x > 0 \\ x & x \leq 0. \end{cases}$$

differentiable at $x_0 = 0$?

Solution: Ans. No. The key, once again, is to compute a single-sided limit at the origin. If any one of them does not exist or is unbounded, then the derivative is undefined at that point.

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

This shows that f is not differentiable at the origin. For completeness, we compute the other derivative, namely

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

Note that the “jump” discontinuity at the origin made the function not differentiable at that point. ■

The following two Propositions are actually equivalent (aka, they state the same fact). At this point in your math education, it’s OK not to recognize this immediately.

Proposition 5.4: Never Differentiable at a Point of Discontinuity

Suppose $f : [a, b] \rightarrow \mathbb{R}$ is discontinuous at $x_0 \in [a, b]$. Then, f is nondifferentiable^a at x_0 .

Note: x_0 could be a point of discontinuity that arises from a “jump” in the function’s value at x_0 , as in Example 5.3, or it could be due to something wilder, such as

$$f(x) := \begin{cases} \sin(1/x) & x \neq 0 \\ 0 & x = 0 \end{cases}$$

not being continuous at $x_0 = 0$ due to “infinitely rapid variation” in the function’s values near the origin.

^aAKA, not differentiable.

Proposition 5.5: Differentiable Implies Continuous

Suppose $f : [a, b] \rightarrow \mathbb{R}$ is differentiable at $x_0 \in [a, b]$. Then, f is continuous at x_0 .

Why are they equivalent? Well, it’s the classic rule in Logic: $(A \implies B) \iff (\neg B \implies \neg A)$. In our case, let’s select A to be “ f is differentiable at x_0 ” and B to be “ f is continuous at x_0 ”. Then $(A \implies B)$ is Prop. 5.5. Continuing, $\neg B$ is “ f is discontinuous at x_0 ” and $\neg A$ is “ f is nondifferentiable at x_0 ”. The content of Prop. 5.4 is $(\neg B \implies \neg A)$.

Importantly, the two propositions are NOT telling us that a function is differentiable if, and only if, it is continuous. We’ve already seen that $|x|$ is not differentiable at the origin even though it is continuous everywhere.

Example 5.6. Compute the derivatives of the functions below at an arbitrary point $x \in \mathbb{R}$.

(a) $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = x^2$.

(b) $g : \mathbb{R} \rightarrow \mathbb{R}$ by $g(x) = e^x$.

(c) $k : \mathbb{R} \rightarrow \mathbb{R}$ by $k(x) = \sin(x)$.

(d) $\ell : (0, \infty) \rightarrow \mathbb{R}$ by $\ell(x) = \ln(x)$.

Solutions:

(a) $f(x) = x^2$. **Ans.** $\frac{df(x)}{dx} = 2x$. You can also write the answer so that $\frac{d}{dx}$ looks like an operation being performed on the function: $\frac{d}{dx}(x^2) = 2x$. The operation is one of taking the derivative (when it exists).

In this problem, as requested, we let the arbitrary point be x instead of x_0 . This will help us to understand that the derivative

of a function is often another function.

$$\begin{aligned}
 \frac{df(x)}{dx} &= \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \quad (\text{apply definition}) \\
 &= \lim_{h \rightarrow 0} \frac{(x+h)^2 - x^2}{h} \quad (\text{plug in the function}) \\
 &= \lim_{h \rightarrow 0} \frac{(x^2 + 2hx + h^2) - x^2}{h} \quad (\text{expand the square}) \\
 &= \lim_{h \rightarrow 0} \frac{2hx + h^2}{h} \quad (\text{algebra}) \\
 &= \lim_{h \rightarrow 0} 2x + h \quad (\text{more algebra}) \\
 &= 2x.
 \end{aligned}$$

(b) $g(x) = e^x$. **Ans.** $\frac{dg(x)}{dx} = e^x$. The solution to this problem uses Example 4.22, which showed that $\lim_{h \rightarrow 0} \frac{e^h - 1}{h} = 1$.

$$\begin{aligned}
 \frac{dg(x)}{dx} &= \lim_{h \rightarrow 0} \frac{g(x+h) - g(x)}{h} \quad (\text{apply definition}) \\
 &= \lim_{h \rightarrow 0} \frac{e^{(x+h)} - e^x}{h} \quad (\text{plug in the function}) \\
 &= \lim_{h \rightarrow 0} e^x \cdot \frac{e^h - 1}{h} \quad (\text{factor out } e^x) \\
 &= e^x \cdot \lim_{h \rightarrow 0} \frac{e^h - 1}{h} \quad (\text{move } e^x \text{ outside the limit because it is independent of } h) \\
 &= e^x \cdot 1 \quad (\text{Example 4.22}) \\
 &= e^x.
 \end{aligned}$$

(c) $k(x) = \sin(x)$. **Ans.** $\frac{dk(x)}{dx} = \cos(x)$. The solution requires two trig identities, namely,

$$\begin{aligned}
 \sin(a+b) &= \cos(a) \cdot \sin(b) + \sin(a) \cdot \cos(b) \\
 2\sin^2\left(\frac{c}{2}\right) &= 1 - \cos(c).
 \end{aligned} \tag{5.3}$$

and $\lim_{h \rightarrow 0} \frac{\sin(h)}{h} = 1$ from Example 4.29. Related video [The Limit \(do not use L'Hôpital's Rule\)](#) by **BlackPenRedPen**.

$$\begin{aligned}
 \frac{dk(x)}{dx} &= \lim_{h \rightarrow 0} \frac{k(x+h) - k(x)}{h} \quad (\text{apply definition}) \\
 &= \lim_{h \rightarrow 0} \frac{\sin(x+h) - \sin(x)}{h} \quad (\text{plug in the function}) \\
 &= \lim_{h \rightarrow 0} \frac{(\cos(x) \cdot \sin(h) + \sin(x) \cdot \cos(h)) - \sin(x)}{h} \quad (\text{use the first identity}) \\
 &= \cos(x) \cdot \lim_{h \rightarrow 0} \frac{\sin(h)}{h} + \sin(x) \cdot \lim_{h \rightarrow 0} \frac{(\cos(h) - 1)}{h} \quad (\text{algebra}) \\
 &= \cos(x) \cdot \lim_{h \rightarrow 0} \frac{\sin(h)}{h} - \sin(x) \cdot \lim_{h \rightarrow 0} 2 \cdot \frac{\sin^2\left(\frac{h}{2}\right)}{h} \quad (\text{use the second identity}) \\
 &= \cos(x) \cdot 1 - \sin(x) \cdot \lim_{h \rightarrow 0} \sin\left(\frac{h}{2}\right) \cdot \frac{\sin\left(\frac{h}{2}\right)}{\frac{h}{2}} \quad (\text{more algebra}) \\
 &= \cos(x),
 \end{aligned}$$

because from Prop. 2.35, when we have a product of functions where at least one has a finite limit, then the limit of the product is the product of the limits. Moreover,

$$\lim_{h \rightarrow 0} \sin\left(\frac{h}{2}\right) = 0 \text{ and } \lim_{h \rightarrow 0} \frac{\sin\left(\frac{h}{2}\right)}{\frac{h}{2}} = 1.$$

Therefore,

$$\begin{aligned} \lim_{h \rightarrow 0} \left(\sin\left(\frac{h}{2}\right) \right) \cdot \left(\frac{\sin\left(\frac{h}{2}\right)}{\frac{h}{2}} \right) &= \underbrace{\left(\lim_{h \rightarrow 0} \sin\left(\frac{h}{2}\right) \right)}_0 \cdot \underbrace{\left(\lim_{h \rightarrow 0} \frac{\sin\left(\frac{h}{2}\right)}{\frac{h}{2}} \right)}_1 \quad (\text{Applying Prop. 2.35 and the above limits}) \\ &= 0 \cdot 1 \\ &= 0. \end{aligned}$$

(d) $\ell(x) = \ln(x)$. **Ans.** $\frac{d\ell(x)}{dx} = \frac{1}{x}$.

We worked this one in Example 4.21-(c). ■

That's how derivatives are computed via the definition. As a practicing engineer, you will do that like, never in your career, and in your engineering courses, you are unlikely to be asked to compute a derivative by the definition. Hence, we'll not belabor it. Instead, as with continuous functions, we'll close this Section with a list of common functions and their derivatives...once we introduce a second common notation for the derivative of a function.

Notation for Differentiation

Definition 5.7. In Newton's notation for the derivative, the primary objects are functions, such as $f(x) = x^2$, and derivatives are written with a prime, as in $f'(x) = 2x$; some claim that Lagrange introduced this notation on the basis of Newton's work. In Leibniz's notation, the primary objects are relationships, such as $y = x^2$, and derivatives are written as a ratio, as in $\frac{dy}{dx} = 2x$. Both notations are quite common, and we will even see a third notation (based on dots) when we focus on derivatives with respect to time of a trajectory, such as $x(t)$.

$$\begin{aligned} f'(x) &:= \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \quad \text{Newton Style} \\ \frac{df(x)}{dx} &:= \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \quad \text{Leibniz Style}. \end{aligned} \tag{5.4}$$

Note: Once again, we could create a perfect bubble in the course where one and only one notation is used, but then you would be at a disadvantage in other courses and life.

Proposition 5.8: Common Functions, Their Derivatives, and Domains

The following functions are differentiable on the indicated domains. Here, we use Newton's notation for the derivative, namely

$$f'(x) := \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}.$$

Function	Derivative	Domain
$f(x) = c$	$f'(x) = 0$	$x \in \mathbb{R}$
$f(x) = x^n$	$f'(x) = nx^{n-1}$	$x \in \mathbb{R}$
$f(x) = e^x$	$f'(x) = e^x$	$x \in \mathbb{R}$
$f(x) = \ln(x)$	$f'(x) = \frac{1}{x}$	$x > 0$
$f(x) = \sin(x)$	$f'(x) = \cos(x)$	$x \in \mathbb{R}$
$f(x) = \cos(x)$	$f'(x) = -\sin(x)$	$x \in \mathbb{R}$
$f(x) = \tan(x)$	$f'(x) = \sec^2(x)$	$x \neq \frac{\pi}{2} + k\pi, k \in \mathbb{Z}$
$f(x) = \tan(x)$	$f'(x) = 1 + \tan^2(x)$	$x \neq \frac{\pi}{2} + k\pi, k \in \mathbb{Z}$
$f(x) = \cot(x)$	$f'(x) = -\csc^2(x)$	$x \neq k\pi, k \in \mathbb{Z}$
$f(x) = \sec(x)$	$f'(x) = \sec(x) \tan(x)$	$x \neq \frac{\pi}{2} + k\pi, k \in \mathbb{Z}$
$f(x) = \csc(x)$	$f'(x) = -\csc(x) \cot(x)$	$x \neq k\pi, k \in \mathbb{Z}$
$f(x) = \sinh(x)$	$f'(x) = \cosh(x)$	$x \in \mathbb{R}$
$f(x) = \cosh(x)$	$f'(x) = \sinh(x)$	$x \in \mathbb{R}$
$f(x) = \tanh(x)$	$f'(x) = \operatorname{sech}^2(x)$	$x \in \mathbb{R}$

Note: $1 + \tan^2(x) = \left(\frac{\cos(x)}{\cos(x)}\right)^2 + \left(\frac{\sin(x)}{\cos(x)}\right)^2 = \frac{\sin^2(x) + \cos^2(x)}{\cos^2(x)} = \left(\frac{1}{\cos(x)}\right)^2 = \sec^2(x)$. Both forms are common for the derivative of $\tan(x)$. In a similar manner, $-\csc^2(x) = -\frac{1}{\sin^2(x)} = -\frac{\sin^2(x) + \cos^2(x)}{\sin^2(x)} = -\left(1 + \cot^2(x)\right)$.

5.2 The Derivative as a Local Linear Approximation of a Function

Let $f : (a, b) \rightarrow \mathbb{R}$ and let $x_0 \in (a, b)$. Suppose we wanted to approximate the function by a line, $y = y_0 + m \cdot (x - x_0)$, near the point x_0 . Fig. 5.1 suggests that the derivative is doing something like that. Let's explore it a bit more deeply.

We pose

$$f(x) \approx y_0 + m \cdot (x - x_0) \quad (5.5)$$

for $|x - x_0| < \delta$ and some $\delta > 0$. What values should we assign to y_0 and m ? Well, if we plug $x = x_0$ into (5.5), we obtain

$$f(x_0) = y_0.$$

Hence,

$$f(x) \approx f(x_0) + m \cdot (x - x_0). \quad (5.6)$$

Next, because we are used to limits, we plug in $x = x_0 + h$ into (5.6), yielding

$$f(x_0 + h) \approx f(x_0) + \underbrace{m \cdot h}_{h=(x_0+h)-x_0}.$$

A little algebra gives us

$$(f(x_0 + h) - f(x_0) \approx m \cdot h) \iff \left(m \approx \frac{f(x_0 + h) - f(x_0)}{h}\right). \quad (5.7)$$

The approximation should be better and better with smaller and smaller h . In the limit, we have

$$m = \lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h}, \quad (5.8)$$

as long as the limit exists. In other symbols, if f is differentiable at x_0 , then

$$f(x) \approx f(x_0) + \frac{df(x_0)}{dx} \cdot (x - x_0) \quad (5.9)$$

is a **local linear approximation** to the function near the point x_0 .

(Optional Read:) More Formal Definition of a Linear Approximation of a Function Near a Point

Definition 5.9. Let $f(x)$ be a differentiable function at x_0 . The function $y(x) = y_0 + m(x - x_0)$ is said to be a **linear approximation** of $f(x)$ near x_0 if

$$\lim_{x \rightarrow x_0} \frac{f(x) - (y_0 + m(x - x_0))}{x - x_0} = 0.$$

Roughly speaking, this definition implies that the error between $f(x)$ and its linear approximation $y(x)$ “quadratic” in $x - x_0$.

Claim: The only possible values for y_0 and m such that $y(x) = y_0 + m(x - x_0)$ is a linear approximation of $f(x)$ near x_0 are $y_0 = f(x_0)$ and $m = f'(x_0)$.

Proof: Assume that $y(x) = y_0 + m(x - x_0)$ is a linear approximation of $f(x)$ near x_0 . By the definition of linear approximation, we have the condition,

$$\lim_{x \rightarrow x_0} \frac{f(x) - (y_0 + m(x - x_0))}{x - x_0} = 0.$$

Expanding the numerator gives,

$$f(x) - (y_0 + m(x - x_0)) = f(x) - y_0 - m(x - x_0).$$

thus, the condition becomes,

$$\lim_{x \rightarrow x_0} \frac{f(x) - y_0 - m(x - x_0)}{x - x_0} = 0.$$

We can rewrite this as,

$$\lim_{x \rightarrow x_0} \left(\frac{f(x) - f(x_0)}{x - x_0} + \frac{f(x_0) - y_0}{x - x_0} - m \right) = 0.$$

Using the fact that $\lim_{x \rightarrow x_0} \frac{f(x) - f(x_0)}{x - x_0} = f'(x_0)$, we obtain,

$$f'(x_0) + \lim_{x \rightarrow x_0} \frac{f(x_0) - y_0}{x - x_0} - m = 0.$$

Now, observe that $\lim_{x \rightarrow x_0} \frac{f(x_0) - y_0}{x - x_0}$ is only well-defined if $f(x_0) = y_0$, otherwise the limit would not exist (it would go to infinity). Therefore, we must have,

$$y_0 = f(x_0).$$

Substituting this back into the equation, we obtain,

$$f'(x_0) - m = 0,$$

which implied,

$$m = f'(x_0).$$

Thus, the only possible values for y_0 and m are $y_0 = f(x_0)$ and $m = f'(x_0)$. ■

Later in the Chapter, we will see how derivatives help us to build local quadratic, cubic, ..., n -th degree polynomial approximations to functions. These have two names: Maclaurin series when the approximation is performed at the origin (i.e., $x_0 = 0$) and **Taylor Series** when $x_0 \neq 0$.

Fact 5.10: (Optional Read:) Quality of the Approximation of a Linear Approximation

To approximate a function $f(x)$ at a point x using its behavior near x , we can expand it in terms of a small increment h . The first-order approximation, often referred to as the linear approximation, assumes that $f(x + h)$ can be expressed as:

$$f(x + h) \approx f(x) + f'(x) \cdot h,$$

which implies that the change in f is primarily governed by its derivative at x . This approximation is valid for small values of h , where higher-order terms become negligible.

To rigorously understand the accuracy of this approximation, we consider the error term,

$$\text{Error}(h) := f(x + h) - (f(x) + f'(x) \cdot h).$$

The first-order approximation assumes that this error is of order h^2 , meaning it diminishes much faster than h as $h \rightarrow 0$. Mathematically, we can write:

$$\lim_{h \rightarrow 0} \frac{f(x + h) - f(x) - f'(x) \cdot h}{h^2} = 0.$$

This statement ensures that the remainder of the approximation is proportional to h^2 and becomes insignificant for small h . Thus, the first-order approximation provides a simple yet powerful tool to estimate $f(x + h)$, while the error term analysis assures its accuracy.

It's not worth doing much more with derivatives until you believe they are very easy to compute with modern software tools. Once you see that you can throw a “derivative party” without breaking a sweat, you’ll be more inclined to understand how useful derivatives can be.

5.3 Software Tools for Computing Derivatives

This section is based on ROB 101 *Computational Linear Algebra*: Lab 10, developed by Dr. Alphonsus Antwi Adu-Bredu when he was a PhD student at the University of Michigan, Department of Robotics. He presents three methods to compute derivatives. Here is a link to the [ROB 101 Programming Guide](#).

5.3.1 Symbolic Derivatives

Symbolic means computing derivatives and manipulating functions the way a human would on a sheet of paper. You define variables that are “symbolic” (aka, do not take on numerical values, such as x , y , or z), and you manipulate them as variables. It’s pretty incredible that a machine designed to crunch numbers can do that.

```
1 using Symbolics # Only need to do this once per session
2
3 # Create a custom function to make differentiation convenient
4 function deriv(f, x)
5     return expand_derivatives(Differential(x)(f))
6 end
7
8 # Sample Usage
9 @variables x # Makes x a symbolic variable
10 f = x^3*cos(x) + x^2*sin(x)^2 + log(3x) # Note we did not say f(x) =, just f =
11 dfdx = deriv(f, x)
```

Output

$$\frac{1}{x} - x^3 \sin(x) + 2 \sin^2(x)x + 3x^2 \cos(x) + 2x^2 \cos(x) \sin(x)$$

Recall that in Julia, $\log(x) = \ln(x)$ and all trig functions are in radians.

The above gives us the derivative as a function of x. If we want to evaluate it at a point, say $x_0 = 0.5$, we can do so as follows, with a convenient, custom-made function

```
1 function evaluate(f, x0)
2     return Symbolics.value(substitute(f, (Dict(x => x0) ) ))
3 end
4
5 # Example use
6 @show dfdx
7 evaluate(dfdx, .5)
```

Output

```
dfdx = 1 / x + 2x*(sin(x)^2) + 3(x^2)*cos(x) + 2(x^2)*cos(x)*sin(x) - (x^3)*sin(x)
3.0384753223601586
```

Note: The output of `substitute(f, (Dict(x => x0)))` has Type Num as you can check with the command `typeof`

```
1 val_numeric = substitute(dfdx, (Dict(x => 0.5))) # Will be Type Num
2
3 @show typeof(val_numeric)
4
5 val = Symbolics.value(val_numeric) # converts to Float64
```

Output

```
typeof(val_numeric) = Num
3.0384753223601586
```

It is recommended that you use¹ the “evaluate” command created above. It is, for some reason, not a part of the Symbolic.jl package.

Continuing with our demonstration of symbolic differentiation, one can even compute the derivative of the derivative! It goes like so:

```
1 d2fdx2 = deriv(dfdx, x)
```

Output

$$\frac{-1}{x^2} + 2\sin^2(x) + 2\cos^2(x)x^2 - x^3\cos(x) + 6x\cos(x) - 2\sin^2(x)x^2 - 6x^2\sin(x) + 8x\cos(x)\sin(x)$$

And, we can go completely bonkers and put the derivative operation in a loop, like so

```
1 @variables x
2 f = x^3*cos(x) + x^2*sin(x)^2 + log(3x)
3 N=5
4 temp = f
5 for k = 1:N
6     temp= deriv(temp, x)
7 end
8 dNfdxN = temp
```

¹Yes, you have to write your version of the function!

Output

$$\begin{aligned} & \frac{-32}{x^5} - 60 \cos(x) - 80 \cos^2(x)x - 4x^3 \frac{2}{x^8} - x^3 \sin(x) + 60x \sin(x) \\ & + 80 \sin^2(x)x - 64x^{11} \frac{-1}{x^{16}} + 15x^2 \cos(x) \\ & - 160 \cos(x) \sin(x) + 32x^2 \cos(x) \sin(x). \end{aligned}$$

Note that the awkward expression,

$$-4x^3 \frac{2}{x^8} - 64x^{11} \frac{-1}{x^{16}},$$

was not simplified. We can fix that as follows,

```
1 # Added a simplify operation
2 function deriv(f, x)
3     df = Differential(x)(f)
4     return simplify(expand_derivatives(df)) # Note the simplify command
5 end
6
7 # Example use
8 @variables x
9 f = x^3*cos(x) + x^2*sin(x)^2 + log(3x)
10 N=5
11 temp = f
12 for k = 1:N
13     temp= deriv(temp, x)
14 end
15 dNfdxN = temp
```

Output

$$\frac{24}{x^5} + 6(-\cos(x) - 4\sin(2x)) - 54\cos(x) - 56\sin(2x) - x^3 \sin(x) + 15x^2 \cos(x) + 60x \sin(x) + 16x^2 \sin(2x) - 80x \cos(2x)$$

The “cost” is that the function runs slower. Hence, it’s your call: to `simplify` or not to `simplify`, that is the question!

Advantages of Symbolic Differentiation

- Symbolic Differentiation allows you to compute exact derivatives up to machine precision.
- Symbolic Differentiation is also a great tool for analysis. Being able to generate closed-form derivatives of a model facilitates the understanding of the dynamics of the model.
- Symbolic Differentiation is also significantly useful for generating closed-form derivatives of functions as code. Closed-loop controllers running on real-time systems often have to run at frequencies in the 1kHz range. This constrains the control loop to a time budget of 1 millisecond. Given this time budget, operations in the controller (such as computing derivatives and Jacobians) have to be as fast as possible. Being able to generate closed-form parameterized functions for derivatives and Jacobians offline using Symbolic differentiation that can then be evaluated online in microseconds can be very useful in speeding up your controller and meeting the controller time budget.

Disadvantages of Symbolic Differentiation

- The main disadvantage of symbolic derivatives is that they tend to get gnarly and exponentially long very quickly. This phenomenon is called **Expression swell**. Imagine taking the derivative of a function

$$h(x) = f(x)g(x)$$

where $f(x)$ is itself the product $f(x) = u(x)v(x)$, the derivative of $h(x)$ ends up as this gnarly function

$$\frac{d}{dx} h(x) = \left(\frac{d}{dx} u(x)v(x) + u(x) \frac{d}{dx} v(x) \right) g(x) + u(x)v(x) \frac{d}{dx} g(x)$$

leading to inefficient code. We will learn the **Product Rule** in Chapter 5.4.

5.3.2 Numerical Differentiation

Numerical Differentiation, also called **finite differences**, was the primary method we used in ROB 101 *Computational Linear Algebra*.

Numerical Approximations of a Derivative

Common approximations of a derivative are based on “not letting the limit go all the way to zero” but stopping at h small.

$$\frac{df(x_0)}{dx} := \lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h}. \quad (5.10)$$

In practice, we will use “small values” for h and never compute the exact limit. Hence, we have an **approximation for the derivative** at a point, namely

$$\frac{df(x_0)}{dx} \approx \frac{f(x_0 + h) - f(x_0)}{h}, \quad (5.11)$$

which is called a **forward difference approximation to the derivative**.

We can also do a **backward difference approximation to the derivative**,

$$\frac{df(x_0)}{dx} \approx \frac{f(x_0) - f(x_0 - h)}{h}, \quad (5.12)$$

and a **symmetric difference approximation**, where we go both forward and backward from the point x_0 ,

$$\frac{df(x_0)}{dx} \approx \frac{f(x_0 + h) - f(x_0 - h)}{2h}. \quad (5.13)$$

The forward and backward difference approximations to the derivative are, in fact, exact for linear functions, while the symmetric difference approximation is *exact* for quadratic polynomials. The symmetric difference is also sometimes called a **central difference**.

If the derivative of $f(x)$ at a point x_0 exists, then for h sufficiently small, the forward difference, backward difference, and symmetric difference approximations to the derivative will always agree. If they provide different answers, then the limit in (5.10) does not exist, and the function is said to be nondifferentiable.

In Julia, you do numerical derivatives as follows.

```
1 using FiniteDiff # Julia package for finite difference approximations of derivatives.  
2  
3 # Example function  
4 g(x) = x^3*cos(x) + x^2*sin(x)^2 + log(3x)  
5  
6 val = FiniteDiff.finite_difference_derivative(g, 0.5)
```

Output

3.038475322499568

The “true” (symbolic) value was 3.0384753223601586 ; their difference is 1.394×10^{-10} . That seems pretty good! Note that the FiniteDiff package chooses the value of “ $h > 0$ ” for you, which is pretty sweet.

Advantages of Numerical Differentiation

- Numerical differentiation is the simplest differentiation method to implement. It only depends on the function values and not on any set of differential calculus rules.
- In view of this, numerical differentiation can be applied to any differentiable function to get reasonably accurate results.

Disadvantages of Numerical Differentiation

- Finite Difference suffers from *Truncation error* resulting in inaccurate derivatives when the step size is large.
- Finite Difference can also be expensive to compute for functions with large input dimensions. This is more important when we get to Jacobians and gradients.
- It is not easy to “iterate” numerical differentiation methods. We had to develop special formulas for that in ROB 101.

5.3.3 Automatic Differentiation

Similar to Symbolic Differentiation, **Automatic Differentiation (aka, AD)** uses differential calculus rules to compute derivatives. **However, its goal is to return a numerical solution instead of a single closed-form analytical expression.** As such, Automatic Differentiation keeps track of intermediate variables and their derivatives to speed up computation. It decomposes the function into primitive operations using the “chain rule”, evaluates the operations and their derivatives, and stores them in intermediate variables. We cover the **chain rule** in Chapter 5.4. Your author thinks of Automatic Differentiation as “the chain rule meets databases”! It is very efficient.

There are two primary forms of automatic differentiation, **Forward Mode AD**

- In forward mode AD, the derivative information is propagated from the input variables through to the output of the function.
- This is done by augmenting each intermediate variable in the computation with its derivative (or tangent) with respect to the input variables.
- Forward mode is particularly efficient when the number of input variables is small compared to the number of output variables.

There is also **Reverse Mode AD**

- In reverse mode AD, the derivative information is propagated from the output back to the input variables.
- This is done by computing the adjoint or cotangent of each intermediate variable, which represents the rate of change of the output with respect to that intermediate variable.
- The adjoints are then used to compute the gradient of the output with respect to the input variables.
- Reverse mode is particularly efficient when the number of output variables is small compared to the number of input variables. This property makes it especially popular in deep learning, where the objective is typically scalar (e.g., loss) but the input (e.g., parameters of a neural network) is high-dimensional.

Both forward and reverse modes of AD are exact methods (up to machine precision) and are more accurate than numerical differentiation, which can suffer from truncation and rounding errors. In addition to these primary forms, there are also hybrid methods that combine aspects of both forward and reverse modes to optimize the computation of derivatives for specific types of functions or computational graphs (e.g., Neural Networks).

At this point in your Calculus education, keeping all these methods straight is not important. What you should take away is that differentiation is an operation to extract information from functions. Until recently, differentiation was a painful manual process. Today, you have tons of options for computing derivatives. Modern AI would not exist without these methods; it relies especially on Automatic Differentiation for the training process of Neural Networks. Each tool has its place. **In other science and engineering courses, you may be terrified by the complicated functions you are asked to differentiate because you have to do it by hand, maybe even on an exam. As a practicing engineer, it's cake. Crazy world, huh?**

To appreciate how automatic differentiation works, let's take a look at an example. Consider the function,

$$f(x) = \sin(x^2 + 2x).$$

Using forward mode automatic differentiation, we can compute its derivative by breaking down the function into elementary operations.

Intermediate Variables

$$\begin{aligned}v_1 &= x^2 \\v_2 &= 2x \\v_3 &= v_1 + v_2 \\f(x) &= \sin(v_3)\end{aligned}$$

Derivatives

Here, we are using the “prime” notation for the first derivative of a quantity.

$$\begin{aligned}v'_1 &= 2x \quad (\text{derivative of } x^2) \\v'_2 &= 2 \quad (\text{derivative of } 2x) \\v'_3 &= v'_1 + v'_2 \quad (\text{derivative of a sum equals sum of the derivatives}) \\f'(x) &= \cos(v_3) \cdot v'_3 \quad (\text{chain rule, which we have not yet covered})\end{aligned}$$

Evaluation

For $x = 1$:

$$\begin{aligned}v_1 &= 1 \\v_2 &= 2 \\v_3 &= 3 \\f(1) &= \sin(3) \\v'_1 &= 2 \\v'_2 &= 2 \\v'_3 &= 4 \\f'(1) &= \cos(3) \cdot 4\end{aligned}$$

Thus, the value of the derivative of the function at $x = 1$ is $4 \cos(3) = -3.9599699864017817$.

To keep things coherent, we return to our previous function and compute its derivative using Forward-mode Automatic Differentiation.

```
1 using ForwardDiff # load the ForwardDiff.jl package as usual
2                         # only need to do once per session
3
4 f2(x) = x^3*cos(x) + x^2*sin(x)^2 + log(3x) # define the function
5
6 x0 = 0.5 # specific point where we want to evaluate the derivative
7
8 df2dx = ForwardDiff.derivative(f2, x0)
```

Output

3.0384753223601586

5.3.4 Guidelines on Choosing Differentiation Methods

Symbolic Differentiation

- **When to Use:** For functions that have a known algebraic form and when you need an explicit formula for the derivative.
- **Advantages:**
 - Provides exact, closed-form expressions for derivatives.
 - Useful for understanding the underlying behavior and properties of a function.
- **Limitations:**
 - Can lead to complex and lengthy expressions for even moderately complicated functions.
 - Not suitable for functions without an explicit mathematical expression.

Numerical Differentiation

- **When to Use:** When you have discrete data points or when an analytical solution is difficult to obtain.
- **Advantages:**
 - Simple to implement and understand.
 - Doesn't require knowledge of the function's underlying algebraic form.
- **Limitations:**
 - Can introduce errors due to truncation and rounding.
 - Accuracy depends on the choice of step size.

Automatic Differentiation (AD)

- **When to Use:** For complex functions where symbolic differentiation is impractical, and when high accuracy is needed. Common in computational applications and machine learning.
- **Advantages:**
 - Provides exact derivatives up to machine precision.
 - Can handle a wide range of functions, including those defined by algorithms or programs.
- **Limitations:**
 - Might be overkill for simple functions.
 - Requires specialized software or libraries.

General Advice

For learners just starting with calculus, it's beneficial to begin with symbolic differentiation to grasp the fundamental concepts. As you progress and encounter more complex functions or real-world data, exploring numerical and automatic differentiation methods will be valuable.

5.4 Differentiation Rules that all Engineers are Expected to Understand

Proposition 5.8 provided the derivatives of many common or elementary functions, but failed to consider more complicated functions, such as

- $e^{\sin(x)}$
- $\cos(\omega_0 x + \theta_0)$
- $\sin(x^3) + \tan(x)$
- $x \cdot \ln(x),$

which are compositions of elementary functions. Determining their derivatives via the limit definition would be very slow and painful! Instead, as hinted at in our discussion of Automatic Differentiation, there are **Differentiation Rules** that are used to compute the derivatives.

As a practicing engineer, in discussions with colleagues, you will be expected to know the following means of computing derivatives. **More immediately, in most of your engineering courses, you will need competency in the following derivative rules in order to understand lectures, complete HW, and pass exams. All of the differentiation rules can be derived from our result in Chapter 5.2 relating differentiation to linear approximations of functions.** An animated, quite beautiful demonstration of these rules using this idea is provided by **Mathologer: Why is calculus so ... EASY ?**; here, we've jumped ahead in the video to the relevant location. While you are welcome to start it from the beginning, we are not ready for that just yet; you may wish to come back to this video when we do **antiderivatives** in Chapter 7.

5.4.1 Differentiating by the Rules

The Rules all in one place; only the Rules.

Proposition 5.11: Differentiation Rules

Here are the rules of differentiation:

- (a) **Sum/Difference Rule:** If $f(x) = g(x) + h(x)$ or $f(x) = g(x) - h(x)$, then

$$f'(x) = g'(x) + h'(x) \text{ or } f'(x) = g'(x) - h'(x).$$

- (b) **Product Rule:** If $f(x) = g(x) \cdot h(x)$, then

$$f'(x) = g'(x) \cdot h(x) + g(x) \cdot h'(x).$$

In particular, if $g(x) = c$, a constant, then

$$f'(x) = c \cdot h'(x).$$

- (c) **Quotient Rule or Ratio Rule:** If $f(x) = \frac{g(x)}{h(x)}$ and $h(x) \neq 0$, then

$$f'(x) = \frac{g'(x) \cdot h(x) - g(x) \cdot h'(x)}{(h(x))^2}.$$

In particular, if $g(x) = c$, a constant, then

$$f'(x) = -c \cdot \frac{h'(x)}{(h(x))^2}.$$

- (d) **Chain Rule:** If $f(x) = g(h(x))$, then

$$f'(x) = g'(h(x)) \cdot h'(x).$$

- (e) **Chain Rule** (Another way to state it): Mentally, decompose $f(x) = g(h(x))$ as $g(y)$ evaluated at $y = h(x)$. Then,

$$\frac{d}{dx}(g(h(x))) = \frac{dg(y)}{dy} \Big|_{y=h(x)} \cdot \frac{dh(x)}{dx},$$

where $|_{y=h(x)}$ means to evaluate y at $h(x)$, or equivalently, substitute in $h(x)$ for y .

Find one of these forms of the Chain Rule that you can master and stick with it! The Chain Rule shows up everywhere in Calculus.

- (f) **Exponential Rule:** If $f(x) = a^x$, where a is a positive constant, then

$$f'(x) = \ln(a) \cdot a^x.$$

In particular, if $a = e$, Euler's constant, then

$$f'(x) = e^x.$$

- (g) **Logarithm Rule:** If $f(x) = \log_a(x)$, where a is a positive constant, then

$$f'(x) = \frac{1}{x \ln(a)}.$$

In particular, if $a = e$, Euler's constant, then

$$f'(x) = \frac{1}{x}.$$

- (h) **Power Rule:** If $f(x) = x^n$, where $n \in \mathbb{N}$, the counting numbers, then

$$f'(x) = nx^{n-1}$$

and there is no restriction on x .

- (i) **Generalized Power Rule:** If $f(x) = x^\alpha$, where $\alpha \in \mathbb{R}$ and $x > 0$, then

$$f'(x) = \alpha x^{\alpha-1}.$$

Here, when the exponent (power) is an arbitrary real number, we need $x > 0$. When the exponent is a counting number, $f'(x) = nx^{n-1}$, and there is no restriction on x .

It's pretty incredible what the `Symbolics` and `ForwardDiff` packages are doing for us! There are many online sources for the proofs of the differentiation rules. Here are a few.

- Paul's Online Notes
- openstax.org Differentiation Rules
- **Mathologer: Why is calculus so ... EASY ?** (Recall that Leibniz's notation for the derivative is $\frac{df}{dx}$, an increment in f over an increment in x . The prime notation is due to Newton. **Mathologer shows how the $\frac{df}{dx}$ notation eases the derivation of the rules of differentiation.**)

As we show below, it's really all about linear approximations of functions about a point.

5.4.2 The How and Why of the Rules

In this section, some of the examples explain why the rules are true, while other examples illustrate the application of the rules.

Example 5.12. Derive the Chain Rule using our result in Chapter 5.2 relating differentiation to linear approximations of functions.

Solution: Using the prime notation for derivatives, we write the linear approximations for f and g about points x and y as

$$g(x+h) \approx g(x) + g'(x) \cdot h \text{ for } h \text{ small} \iff \lim_{h \rightarrow 0} \frac{g(x+h) - g(x)}{h} = g'(x)$$

$$f(y+\delta) \approx f(y) + f'(y) \cdot \delta \text{ for } \delta \text{ small} \iff \lim_{\delta \rightarrow 0} \frac{f(y+\delta) - f(y)}{\delta} = f'(y);$$

shortly, we'll be motivated to take $y = g(x)$ and $\delta = g'(x) \cdot h$.

By the definition of the derivative of $f \circ g$ at x ,

$$(f \circ g(x))' = (f(g(x)))' = \lim_{h \rightarrow 0} \frac{f(g(x+h)) - f(g(x))}{h}.$$

As Step 1, in the above formula, we replace $g(x+h)$ by its linear approximation about x ,

$$(f(g(x)))' = \lim_{h \rightarrow 0} \frac{f(g(x) + g'(x) \cdot h) - f(g(x))}{h}.$$

As Step 2, in the above formula, we identify $f(g(x) + g'(x) \cdot h)$ as $f(y+\delta)$. We can then replace f with its linear approximation about $y = g(x)$ and $\delta = \delta(h) := g'(x) \cdot h$, where we note that $h \rightarrow 0$ implies $\delta(h) \rightarrow 0$. This gives us,

$$\begin{aligned} (f(g(x)))' &= \lim_{h \rightarrow 0} \frac{f(g(x)) + f'(g(x)) \cdot g'(x) \cdot h - f(g(x))}{h} \\ &= \lim_{h \rightarrow 0} \frac{f'(g(x)) \cdot g'(x) \cdot h}{h} \\ &= f'(g(x)) \cdot g'(x). \end{aligned}$$

Again, the Mathologer video animates the above proof. The animation even includes the important point that all of this works because of the equivalence between differentiation and linear approximation of a function near a point. ■

Example 5.13. Derive the Product Rule using our result in Chapter 5.2 relating differentiation to linear approximations of functions.

Solution: Using the prime notation for derivatives, we write the linear approximations for f and g about a point x as

$$g(x+h) \approx g(x) + g'(x) \cdot h \text{ for } h \text{ small} \iff \lim_{h \rightarrow 0} \frac{g(x+h) - g(x)}{h} = g'(x)$$

$$f(x+h) \approx f(x) + f'(x) \cdot h \text{ for } h \text{ small} \iff \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} = f'(x).$$

By the definition of the derivative of $f(x) \cdot g(x)$ at x ,

$$= (f(x) \cdot g(x))' = \lim_{h \rightarrow 0} \frac{f(x+h) \cdot g(x+h) - f(x) \cdot g(x)}{h}.$$

We replace both $f(x+h)$ and $g(x+h)$ by their linear approximations about x , and expand out the multiplication

$$\begin{aligned}
(f(x) \cdot g(x))' &= \lim_{h \rightarrow 0} \frac{(f(x) + f'(x) \cdot h) \cdot (g(x) + g'(x) \cdot h) - f(x) \cdot g(x)}{h} \\
&= \lim_{h \rightarrow 0} \frac{(f(x) \cdot g(x) + f(x) \cdot g'(x) \cdot h + f'(x) \cdot g(x) \cdot h + f'(x) \cdot g'(x) \cdot h^2) - f(x) \cdot g(x)}{h} \\
&= \lim_{h \rightarrow 0} \frac{f(x) \cdot g'(x) \cdot h + f'(x) \cdot g(x) \cdot h + f'(x) \cdot g'(x) \cdot h^2}{h} \\
&= f(x) \cdot g'(x) + f'(x) \cdot g(x) + \lim_{h \rightarrow 0} \frac{f'(x) \cdot g'(x) \cdot h^2}{h} \\
&= f(x) \cdot g'(x) + f'(x) \cdot g(x)
\end{aligned}$$

because

$$\lim_{h \rightarrow 0} \frac{f'(x) \cdot g'(x) \cdot h^2}{h} = \lim_{h \rightarrow 0} f'(x) \cdot g'(x) \cdot h = f'(x) \cdot g'(x) \cdot \lim_{h \rightarrow 0} h = 0.$$

Again, the Mathologer [video](#) animates the above proof. ■

Fact 5.14: Practice Videos

A great source of practice examples: [A Nancy Pi Playlist on Differentiation](#). Below are some of the key videos.

- Derivatives... How? (basic intro, product rule, quotient rule)
- The Chain Rule... How? When? **Excellent Video! Be sure to watch to the end.** We will borrow some of Nancy Pi's ideas in the examples that follow.
- More Chain Rule means more practice,

Other sources are [Calculus I Tutorial: how to use the product rule, quotient rule, and chain rule](#) by **BlackPenRedPen**, and [Differentiation Formulas](#) by The Organic Chemistry Tutor, is basically a list that is written out on the screen, which does not sound enticing. However, it has a lot of derivatives!

In ROB 101, you had a HW problem similar to this video [Derivative without Equations](#) by **BlackPenRedPen**.

Finally, your author finds this video's topic spot on: [HOW TO READ CALCULUS OUT LOUD! LIMITS, DERIVATIVES & INTEGRAL SYMBOLS](#) by **I Need to Pass Calculus...** (what a cool channel name). Share your thoughts in class.

Function	c	x^n	e^x	$\ln(x)$	$\sin(x)$	$\cos(x)$	$\tan(x)$	$\text{atan}(x)$
Derivative	0	nx^{n-1}	e^x	$\frac{1}{x}$	$\cos(x)$	$-\sin(x)$	$1 + \tan^2(x)$	$\frac{1}{1+x^2}$

Table 5.15: Common Elementary Functions and their Derivatives (Worth committing to memory to showcase your expertise).

Example 5.16. Using Table 5.15 and the Rules of Differentiation, compute by hand the derivatives of the following functions.

$$(a) f(x) = 3x^2 + 2x + 4$$

$$(b) f(x) = \frac{1}{x}$$

$$(c) f(x) = \frac{1}{x^2}$$

$$(d) f(x) = \frac{\sin(x)}{x}$$

$$(e) f(x) = \sin(x^3 + 4x)$$

$$(f) f(x) = e^{\sin(x^3 + 4x)}$$

$$(g) f(x) = e^{-(x - x_c)^2/s^2} = \exp\left(-\frac{(x - x_c)^2}{s^2}\right) \text{ (Radial Basis Function)}$$

It is suggested to check your work using the `Symbolic` differentiation tools from Chapter 5.3.1.

Solutions:

(a) $f(x) = 3x^2 + 2x + 4$ **Ans.** $f'(x) = 6x + 2$ uses the Sum Rule, Constant Rule $(c)' = 0$, and the Power Rule $(x^k)' = kx^{k-1}$. The rest you can handle.

(b) $f(x) = \frac{1}{x}$ **Ans.** $f'(x) = -\frac{1}{x^2}$ uses the Quotient Rule and noting else. Here, $g(x) = 1$ and $h(x) = x$. Hence,

$$f'(x) = \frac{g'(x) \cdot h(x) - g(x) \cdot h'(x)}{(h(x))^2} = \frac{0 \cdot x - 1 \cdot 1}{x^2} = -\frac{1}{x^2}.$$

(c) $f(x) = \frac{1}{x^2}$ **Ans.** $f'(x) = -\frac{2x}{x^4} = -2\frac{1}{x^3}$ uses the Quotient Rule. Here, $g(x) = 1$ and $h(x) = x^2$. Hence,

$$f'(x) = \frac{g'(x) \cdot h(x) - g(x) \cdot h'(x)}{(h(x))^2} = \frac{0 \cdot x^2 - 1 \cdot 2x}{(x^2)^2} = -2\frac{1}{x^3}.$$

(d) $f(x) = \frac{\sin(x)}{x}$ **Ans.** $f'(x) = \frac{x \cdot \cos(x) - \sin(x)}{x^2}$ uses the Quotient Rule. Here, $g(x) = \sin(x)$ and $h(x) = x$. Hence,

$$f'(x) = \frac{g'(x) \cdot h(x) - g(x) \cdot h'(x)}{(h(x))^2} = \frac{\cos(x) \cdot x - 1 \cdot \sin(x)}{(x)^2} = \frac{x \cdot \cos(x) - \sin(x)}{x^2}.$$

(e) $f(x) = \sin(x^3 + 4x)$ **Ans.** $f'(x) = (3x^2 + 4) \cdot \cos(x^3 + 4x)$ uses the Chain Rule and the Power Rule. Here, $g(x) = \sin(x)$ and $h(x) = x^3 + 4x$. Hence,

$$f'(x) = g'(h(x)) \cdot h'(x) = \underbrace{\cos(x^3 + 4x)}_{g'(h(x))} \cdot \underbrace{(3x^2 + 4)}_{h'(x)} = (3x^2 + 4) \cdot \cos(x^3 + 4x),$$

where it is traditional to move the polynomial terms in front of the trigonometric terms, but it is not mandatory to do so.

(f) $f(x) = e^{\sin(x^3 + 4x)}$ **Ans.** $f'(x) = (3x^2 + 4) \cdot \cos(x^3 + 4x) \cdot e^{\sin(x^3 + 4x)}$ uses the Chain Rule **TWICE** and the Power Rule. Here, $g(x) = e^x$ and $h(x) = \sin(x^3 + 4x)$. Fortunately for us, we already computed $h'(x) = (3x^2 + 4) \cdot \cos(x^3 + 4x)$ in the previous part via the Chain Rule, so we only have one application of the Chain Rule to do here! Hence,

$$f'(x) = g'(h(x)) \cdot h'(x) = \underbrace{e^{\sin(x^3 + 4x)}}_{g'(h(x))} \cdot \underbrace{(3x^2 + 4) \cdot \cos(x^3 + 4x)}_{h'(x)} = (3x^2 + 4) \cdot \cos(x^3 + 4x) \cdot e^{\sin(x^3 + 4x)}.$$

You can write the terms in any order you wish!

(g) $f(x) = e^{-(x - x_c)^2/s^2}$ **Ans.** $f'(x) = -2\frac{(x - x_c)}{s^2} \cdot e^{-(x - x_c)^2/s^2}$ uses the Chain Rule and the Power Rule. Here, $g(x) = e^x$ and $h(x) = -\frac{(x - x_c)^2}{s^2}$. Hence,

$$f'(x) = g'(h(x)) \cdot h'(x) = \underbrace{\left(e^{-\frac{(x - x_c)^2}{s^2}}\right)}_{g'(h(x))} \cdot \underbrace{\left(-2\frac{(x - x_c)}{s^2}\right)}_{h'(x)} = -2\frac{(x - x_c)}{s^2} \cdot e^{-(x - x_c)^2/s^2}.$$

You can write the terms in any order.

```

1 @variables x
2
3 fa = 3x^2 + 2x + 4
4
5 fb = 1/x
6
7 fc = 1/(x^2)
8
9 fd = sin(x)/x
10
11 fe = sin(x^3 + 4x)
12
13 ff = exp( ( sin(x^3 + 4x) ) )
14
15 @variables xc s # additional variables in the radial basis function
16           # which are treated as constants because we are differentiating
17           # with respect to x
18
19 fg = exp(-(x-xc)^2 / s^2)
20
21 # With the simplify operation
22 function deriv(f, x)
23     df = Differential(x)(f)
24     return simplify(expand_derivatives(df)) # Note the simplify command
25 end

```

Output

```
deriv (generic function with 1 method)
```

```

1 @show deriv(fa, x)
2
3 @show deriv(fb, x)
4
5 @show deriv(fc, x)
6
7 @show deriv(fd, x)
8
9 @show deriv(fe, x)
10
11 @show deriv(ff, x)
12
13 @show deriv(fg, x);

```

Output

```
deriv(fa, x) = 2 + 6x
deriv(fb, x) = -1 / (x^2)
deriv(fc, x) = -2 / (x^3)
deriv(fd, x) = (x*cos(x) - sin(x)) / (x^2)
deriv(fe, x) = (4 + 3(x^2))*cos(x^3 + 4x)
deriv(ff, x) = (4 + 3(x^2))*cos(x^3 + 4x)*exp(sin(x^3 + 4x))
deriv(fg, x) = (-2(x - xc)*exp((xc - x)*(x - xc)) / (s^2)) / (s^2)
```

We next work examples that are “similar” to the power rule.

Example 5.17. The Generalized Power Rule, Exponential Rule, and Logarithm Rules of Differentiation in Prop. 5.11 can all be derived from the Chain Rule and the Logarithm Rules in Prop. 1.61. Use this idea to derive the derivatives of the following functions.

(a) $f(x) = x^y$, where $x > 0$ and $y \in \mathbb{R}$ is a constant.

(b) $f(x) = \sqrt[a]{x}$, where $x > 0$ and $a \neq 0$ is a constant.

(c) $f(x) = a^x$, where $x \in \mathbb{R}$ and $a > 0$ is a constant.

Hint: For $u > 0$ and $v \in \mathbb{R}$, the Logarithm Rules from Prop. 1.61 yield,

$$u^v = e^{\ln(u^v)} = e^{v \cdot \ln(u)}.$$

This “observation” is key to computing the derivatives of the above functions.

Solutions:

(a) $f(x) = x^y$, where $x > 0$ and $y \in \mathbb{R}$ is a constant. **Ans.** $f'(x) = y \cdot x^{y-1}$ uses the Logarithm Rules and the Chain Rule. As in the Hint, we write

$$f(x) = e^{\ln(x^y)} = e^{y \cdot \ln(x)}.$$

We choose $g(x) = e^x$ and $h(x) = y \cdot \ln(x)$. We note that

$$h'(x) = y \cdot (\ln(x))' = \frac{y}{x}$$

because $(\ln(x))' = \frac{d}{dx} \ln(x) = \frac{1}{x}$ and y is a constant. Hence,

$$f'(x) = g'(h(x)) \cdot h'(x) = \underbrace{(e^{y \cdot \ln(x)})}_{g'(h(x))} \cdot \underbrace{\left(\frac{y}{x}\right)}_{h'(x)} = \underbrace{(x^y)}_{g'(h(x))} \cdot \underbrace{\left(\frac{y}{x}\right)}_{h'(x)} = y \cdot x^{y-1},$$

where we used $e^{y \cdot \ln(x)} = x^y$.

Note: As long as you are careful about the domain of the function, namely, $x > 0$, this works like the Power Rule.

(b) $f(x) = \sqrt[a]{x}$ where $x > 0$ and $a \neq 0$ is a constant. **Ans.** $f'(x) = \frac{1}{a} \cdot x^{\frac{1-a}{a}}$ follows directly from part (a) because

$$= \sqrt[a]{x} = x^{\frac{1}{a}},$$

and we identify $y = \frac{1}{a}$. The rest is algebra; indeed,

$$y - 1 = \frac{1}{a} - 1 = \frac{1-a}{a}.$$

Note: If we take $n = 2$, we have

$$(\sqrt{x})' = \frac{d}{dx} \sqrt{x} = \frac{d}{dx} \sqrt{x} = \frac{1}{2} x^{-\frac{1}{2}} = \frac{1}{2\sqrt{x}},$$

where we have emphasized the multiple ways in which the derivative operation can be written. You are allowed to choose a favorite, but you need to be able to read all of them.

(c) $f(x) = a^x$, where $x \in \mathbb{R}$ and $a > 0$ is a constant. **Ans.** $f'(x) = \ln(a) \cdot a^x$ uses the Logarithm Rules and the Chain Rule. As in the Hint, we write

$$f(x) = e^{\ln(a^x)} = e^{x \cdot \ln(a)}.$$

We choose $g(x) = e^x$ and $h(x) = x \cdot \ln(a)$. We note that $h(x)$ is linear in x because a is a constant and, therefore,

$$h'(x) = \ln(a).$$

Hence,

$$f'(x) = g'(h(x)) \cdot h'(x) = \underbrace{(e^{x \cdot \ln(a)})}_{g'(h(x))} \cdot \underbrace{\ln(a)}_{h'(x)} = \underbrace{(a^x)}_{g'(h(x))} \cdot \underbrace{\ln(a)}_{h'(x)} = \ln(a) \cdot a^x,$$

where we used $e^{x \cdot \ln(a)} = a^x$.

Note: If $a = e$, Euler's constant, then $\ln(a) = \ln(e) = 1$ and $f'(x) = (e^x)' = e^x$.

Second Note: Almost all of you will, at least once in your academic career, say that $(a^x)' = \frac{d}{dx} a^x = a^x$, forgetting the term $\ln(a)$. So, be careful with this one! The good news is, in most engineering courses, you will only see the natural exponential, e^x .

```

1 @variables x y n a # y, n, and a are all treated as constants
2                                # because we are taking derivatives with respect
3                                # to the variable x
4
5 fa = x^y
6
7 fb = x^(1/n)
8
9 fc = a^x
10
11 @show deriv(fa, x)
12
13 @show deriv(fb, x)
14
15 @show deriv(fc, x);

```

Output

```

deriv(fa, x) = y*(x^(y - 1))
deriv(fb, x) = (x^((1 - n) / n)) / n
deriv(fc, x) = (a^x)*log(a)

```

Remark 5.18. We have derived all of the Rules of Differentiation with the exception of the Quotient Rule. It can be derived following the same pattern as in Examples 5.12 and 5.13. **Look for it in HW!**

5.4.3 Bonus Examples

Example 5.19. (Looks Can be Deceiving) From the Generalized Power Rule,

$$\frac{d}{dx} \sqrt{x} = \frac{d}{dx} x^{\frac{1}{2}} = \frac{1}{2} \cdot x^{\frac{1}{2}-1} = \frac{1}{2} \cdot x^{-\frac{1}{2}} = \frac{1}{2} \cdot \frac{1}{\sqrt{x}}.$$

But, what would be

$$\frac{d}{dx} (\sqrt{x}^{\sqrt{x}})?$$

Solution: We give two solutions based on a common technique developed in Example 5.17.

Option 1: The problem looks terrifying until we recall our power rules: $(x^a)^b = x^{a \cdot b}$. For us, $a = \frac{1}{2}$ and $b = \sqrt{x} = x^{\frac{1}{2}}$. Hence, we seek to differentiate

$$\sqrt{x}^{\sqrt{x}} = x^{\frac{1}{2}\sqrt{x}} = e^{\ln(x^{\frac{1}{2}\sqrt{x}})} = e^{\frac{1}{2}\sqrt{x} \cdot \ln(x)}.$$

Key Differentiation Fact: For $h(x)$ differentiable, $\frac{d}{dx} e^{h(x)} = e^{h(x)} \cdot h'(x) = e^{h(x)} \cdot \frac{d}{dx} h(x)$.

Hence,

$$\begin{aligned}
\frac{d}{dx} (\sqrt{x}^{\sqrt{x}}) &= \frac{d}{dx} (e^{\frac{1}{2}\sqrt{x} \cdot \ln(x)}) \\
&= e^{\frac{1}{2}\sqrt{x} \cdot \ln(x)} \cdot \frac{d}{dx} \left(\frac{1}{2}\sqrt{x} \cdot \ln(x) \right) \\
&= e^{\frac{1}{2}\sqrt{x} \cdot \ln(x)} \cdot \frac{1}{2} \cdot \frac{d}{dx} (\sqrt{x} \cdot \ln(x)) \\
&= \sqrt{x}^{\sqrt{x}} \cdot \frac{1}{2} \left(\frac{1}{2} \cdot \frac{1}{\sqrt{x}} \cdot \ln(x) + \sqrt{x} \cdot \frac{1}{x} \right) \\
&= \sqrt{x}^{\sqrt{x}} \cdot \frac{1}{4\sqrt{x}} (\ln(x) + 2) \\
&= \frac{1}{4} (\ln(x) + 2) \cdot \sqrt{x}^{\sqrt{x} - \frac{1}{2}}.
\end{aligned}$$

Option 2: We solve a more general case first! That's right, sometimes generality can oblige us to focus. What is $\frac{d}{dx} (f(x)^{g(x)})$, assuming $f : \mathbb{R} \rightarrow (0, \infty)$ and $g : \mathbb{R} \rightarrow \mathbb{R}$ are both differentiable? It's not so bad, because

$$f(x)^{g(x)} = e^{\ln(f(x)^{g(x)})} = e^{g(x) \cdot \ln(f(x))}.$$

Key Differentiation Facts: For $h(x)$ differentiable,

$$\begin{aligned}
\frac{d}{dx} e^{h(x)} &= e^{h(x)} \cdot h'(x) = e^{h(x)} \cdot \frac{d}{dx} h(x) \\
\frac{d}{dx} \ln(h(x)) &= \frac{1}{h(x)} \cdot h'(x) \quad (\text{for } h(x) \text{ positive}).
\end{aligned}$$

Hence, the Chain Rule and Product Rule give

$$\begin{aligned}
\frac{d}{dx} e^{g(x) \cdot \ln(f(x))} &= e^{g(x) \cdot \ln(f(x))} \cdot \frac{d}{dx} (g(x) \cdot \ln(f(x))) \\
&= e^{g(x) \cdot \ln(f(x))} \cdot \left(g'(x) \cdot \ln(f(x)) + g(x) \cdot \frac{1}{f(x)} \cdot f'(x) \right) \\
&= f(x)^{g(x)} \cdot \left(g'(x) \cdot \ln(f(x)) + \frac{g(x)}{f(x)} \cdot f'(x) \right).
\end{aligned}$$

In our problem $f(x) = g(x) = \sqrt{x}$, and the problem statement gave us $f'(x) = g'(x) = \frac{1}{2} \cdot \frac{1}{\sqrt{x}}$. Hence,

$$\begin{aligned}
\frac{d}{dx} (\sqrt{x}^{\sqrt{x}}) &= \sqrt{x}^{\sqrt{x}} \cdot \left(\frac{1}{2} \cdot \frac{1}{\sqrt{x}} \cdot \ln(\sqrt{x}) + \frac{\sqrt{x}}{\sqrt{x}} \cdot \frac{1}{2} \cdot \frac{1}{\sqrt{x}} \right) \\
&= \frac{1}{2} \sqrt{x}^{\sqrt{x}} \cdot \frac{1}{\sqrt{x}} \cdot \left(\frac{1}{2} \ln(x) + 1 \right) \\
&= \frac{1}{4} (\ln(x) + 2) \cdot \sqrt{x}^{\sqrt{x} - \frac{1}{2}}.
\end{aligned}$$

The problem comes from [in-class derivative vs derivative on the test!](#) by **BlackPenRedPen**. The take-home message: when confronted with a complicated-looking differentiation problem that involves powers (aka, exponents), $x^y = e^{\ln(x^y)}$ can be a powerful opening gambit, as long as the quantity x is positive, thereby ensuring that $\ln(x)$ makes sense. ■.

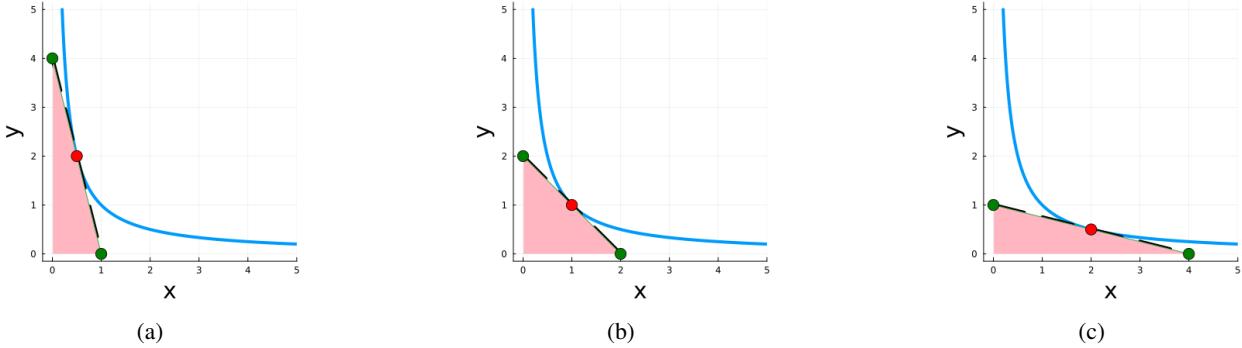


Figure 5.2: The blue line is the graph of $\frac{1}{x}$. The black dashed line is tangent to the red dot. The two green dots are the x - and y -intercepts for the tangent line. The pink triangle is the area between the tangent line and the (x, y) -axes. Its area is the same for all choices of the tangent point. Why the area under the tangent line of $1/x$ is special by BlackPenRedPen.

Example 5.20. Compute the area highlighted in Fig. 5.2.

Solution: The formula for the tangent line is precisely the linear approximation of $f(x) = \frac{1}{x}$ at a point $x_0 = a$. Hence, from (5.9),

$$y = f(a) + f'(a) \cdot (x - a).$$

Because $f'(x) = -\frac{1}{x^2}$, we arrive at

$$y = \frac{1}{a} - \frac{1}{a^2}(x - a) = \frac{2}{a} - \frac{1}{a^2}x.$$

Hence, the y -intercept of the line is the point $(0, \frac{2}{a})$ and the x -intercept is the point $(2a, 0)$. The area of the triangle is

$$\frac{1}{2}b \cdot h = \frac{1}{2} \cdot 2a \cdot \frac{2}{a} = 2,$$

independent of the choice of $a > 0$.

■

5.5 Use Cases of the Single-variable Derivative

That's really all the theory for derivatives of a function depending on a single variable. Before moving on to multivariable functions, we illustrate why derivatives are such a useful tool.

5.5.1 From Position to Velocity

The Planning Team of our company, **Algorithms in Motion**®, has given us the paths in \mathbb{R}^2 shown in Fig. 5.4. Each path gets us from

$$\text{Point A} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \text{ to Point B} = \begin{bmatrix} 10 \\ 10 \end{bmatrix},$$

while avoiding the obstacle shown in red. Good job, Planning Team! Shown by a dotted black line is a straight line connecting the two points, which goes right through the obstacle and would destroy the robot. That's why we called on the Planning Team for help! Our job is the following:

Task: Select the path that has the lowest peak speed so as to minimize the chance of the wheels slipping while maneuvering, where speed is defined to be the Euclidean norm of the velocity vector. At each point in time, velocity has x - and y -components, making it a 2-vector. The norm is

$$\|(v_x, v_y)\| = \sqrt{(v_x)^2 + (v_y)^2}.$$



Figure 5.3: The same small mobile robot that we showed in Chapter 3.4.1. The user commands wheel speed. If both wheels move at the same speed, the robot moves in a straight line. If one wheel turns faster than the other, the robot turns toward the side of the slower wheel. Yet, somehow, our “Planning Team” has given us position trajectories (aka, position curves in \mathbb{R}^2), instead of velocity. We need to convert the position trajectories to velocity commands because our “Control Team” is still working on the position trajectory tracking controller for the robot. It will not be available until Chapter 10.

Yep, we have one task, but to complete it, we have to do several things

- (a) Have the Planning Team send us the equations for the position curves, parameterized by time; once again, they just sent us a plot of p_y versus p_x ; we want the actual functions of time.
- (b) Differentiate the position curves (functions) to obtain velocity.
- (c) Compute the instantaneous speed of each curve from the instantaneous velocity.
- (d) Search over the speed curves to find the peak speed.
- (e) Choose the best curve for our robot, the one that demands the lowest peak speed (thereby inducing the least chance of wheel slippage).

Below are the formulas that were finally sent to us. They are a little awkward, but Julia will not care.

$$\begin{aligned}
 \text{path}_1(t) &= \begin{bmatrix} 0.37500000t - 0.00007813t^3 \\ 0.00937500t^2 - 0.00007813t^3 \end{bmatrix} \\
 \text{path}_2(t) &= \begin{bmatrix} -0.00562500t^2 + 0.37500000t + 0.00006250t^3 \\ -0.00187500t^2 + 0.15000000t + 0.00010938t^3 \end{bmatrix} \\
 \text{path}_3(t) &= \begin{bmatrix} 0.00187500t^2 + 0.22500000t - 0.00003125t^3 \\ -0.03187500t^2 + 0.75000000t + 0.00048438t^3 \end{bmatrix}.
 \end{aligned} \tag{5.14}$$

Next, we place the functions in Julia and do the rest of the work there. Why? Because differentiating the polynomials line by line using the Power Rule is now too easy for you. And because doing the work in code provides documentation that we can share with other members of [Algorithms in Motion](#)[®].

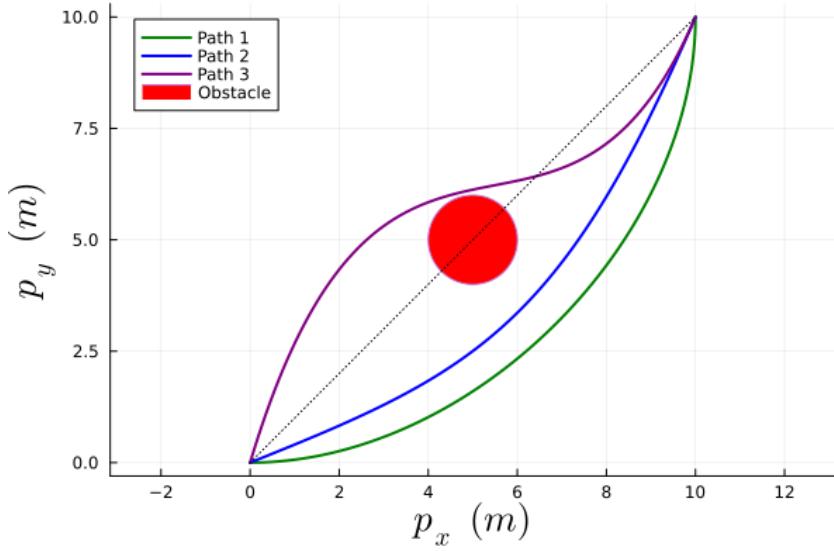


Figure 5.4: The Planning Team has proposed the indicated paths for the robot to follow, each of which avoids the obstacle indicated in red. As explained in Fig. 5.3, our robot accepts the velocity of the paths. We can use Calculus to convert the paths to velocity commands because velocity is the (time) rate of change of position. Hence, $v(t) = \frac{d}{dt}p(t)$.

```

1 # Paths from the Planning Team
2 using LinearAlgebra
3 @variables t
4
5 path1 = [0.375t - 7.812500000000002e-5(t^3),
6     0.009375000000000001(t^2) - 7.812500000000002e-5(t^3)]
7
8 path2 = [0.375t + 6.250000000000001e-5(t^3) - 0.005625000000000015(t^2),
9     0.15000000000000002t + 0.00010937500000000003(t^3) - 0.0018750000000000004(t^2)]
10
11 path3 = [0.001875000000000004(t^2) + 0.225t - 3.125000000000001e-5(t^3),
12     0.75t + 0.000484375000000001(t^3) - 0.03187500000000001(t^2)]
13
14 # Derivatives of the paths
15 v1 = deriv(path1, t)
16 v2 = deriv(path2, t)
17 v3 = deriv(path3, t)
18
19 # Range of t values
20 t_values = collect(0:0.01:40)
21
22
23 v1_traj = evaluate(v1, t_values)
24 v2_traj = evaluate(v2, t_values)
25 v3_traj = evaluate(v3, t_values)
26
27 # The rows are indexed by time. We want to apply
28 # norm to the columns. That is easy to do in Julia
29 speed_1 = norm.(eachcol(v1_traj))
30 speed_2 = norm.(eachcol(v2_traj))
31 speed_3 = norm.(eachcol(v3_traj))
32
33 # Compute max speeds

```

```

34 @show maximum(speed_1)
35 @show maximum(speed_2)
36 @show maximum(speed_3)
37
38 # Make Speed Plots vs time
39
40 # Define an array of colors
41 colors = [:green, :blue, :purple, :orange]
42
43 plot1=plot(t_values, speed_1, linewidth=2, linecolor=colors[1], label="Speed 1",
44 xlabel=L"$t \sim (s)$", ylabel=L"$Speed \sim (m/s)$",
45 yguidefontsize=18, xguidefontsize=18)
46 plot1=plot!(t_values, speed_2, linewidth=2, linecolor=colors[2], label="Speed 2")
47 plot1=plot!(t_values, speed_3, linewidth=2, linecolor=colors[3], label="Speed 3")
48
49 display(plot1)
50 png(plot1, "MobileRobotSpeeds")

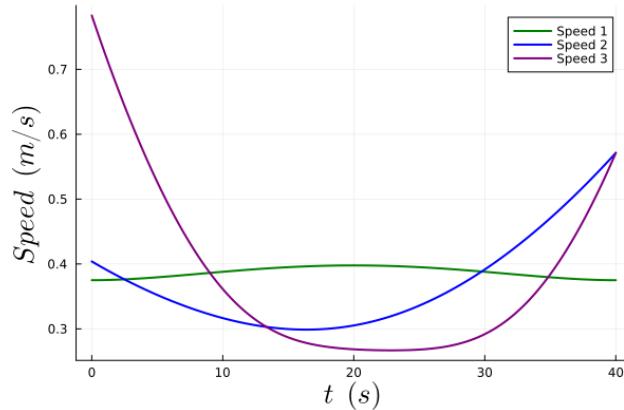
```

Output

```

maximum(speed_1) = 0.39774756441743303
maximum(speed_2) = 0.5711829829397932
maximum(speed_3) = 0.7830229881682913

```



Conclusion: Each of the paths gets us from Point A to Point B in 40 seconds. **Path 1 is the best by our criterion of lowest peak speed.** All of them will likely work. We can give props to the Planning Team!

Here is the “upgraded” evaluate function, using Julia’s multiple dispatch feature.

```

1 # evaluate has here been extended to use multiple dispatch
2
3 # Method when f is a symbolic expression and x0 is a Float64
4 function evaluate(f::Any, x0::Float64)
5     vars = Symbolics.get_variables(f)
6     return Symbolics.value(substitute(f, (Dict(first(vars) => x0) ) ))
7 end
8
9 # Method when f is a function and x0 is a Float64
10 function evaluate(f::Function, x0::Float64)
11     return f(x0)

```

```

12 end
13
14 # Method when f is a vector of symbolic expressions and x0 is a vector of Float64
15 function evaluate(f::Vector{Any}, x0::Vector{Float64})
16     vars = Symbolics.get_variables(f[1])
17     var = first(vars)
18     return hcat([float.(Symbolics.value.(substitute.(f, Ref(var => x_val)))) for x_val in x0]...)
19 end
20
21
22 # Method when f is a vector of symbolic expressions and x0 is a vector of Float64
23 function evaluate(f::Vector{Num}, x0::Vector{Float64})
24     vars = Symbolics.get_variables(f[1])
25     var = first(vars)
26     return hcat([float.(Symbolics.value.(substitute.(f, Ref(var => x_val)))) for x_val in x0]...)
27 end
28
29 # Method when f is a vector of functions and x0 is a vector of Float64
30 function evaluate(f::Vector{Function}, x0::Vector{Float64})
31     return [fi(xi) for (fi, xi) in zip(f, x0)]
32 end

```

Output

```
evaluate (generic function with 5 methods)
```

Secrets of the Arcane 5.21: Taking Stock: Relations between Differentiation and Integration

In the above, we differentiated with respect to time the (x, y) -coordinates of a position trajectory and obtained a velocity profile as a function of time. Back in Chapter 3.4.1, we did the opposite: we started with (x, y) -coordinates of a velocity profile and integrated them to obtain position trajectories. In symbols, this is what we have done,

$$v(t) = \begin{bmatrix} v_x(t) \\ v_y(t) \end{bmatrix} = \begin{bmatrix} \frac{dp_x(t)}{dt} \\ \frac{dp_y(t)}{dt} \end{bmatrix} \quad (5.15)$$

$$p(t) = \begin{bmatrix} p_x(t) \\ p_y(t) \end{bmatrix} = \begin{bmatrix} p_x(t_0) + \int_{t_0}^t v_x(\tau) d\tau \\ p_y(t_0) + \int_{t_0}^t v_y(\tau) d\tau \end{bmatrix}.$$

In other words, it seems that we can go back and forth in the sense that

- If we apply the differentiation operator to a function (or vector of functions), we can “undo it” via integration; and
- if we apply the integration operator to a function (or vector of functions), we can “undo it” via differentiation.

In plain words, one operation is the inverse of the other! This is remarkable. Differentiation takes a very local view of a function to extract its slope at a point, while definite integration focuses on an arbitrarily large (but bounded) region of the function’s domain to compute the area under the function.

One possibility is that this is simply an accident of physics ($F = ma$), and for more general functions, the inverse relationships between differentiation and integration fall apart. Another possibility is that they hold very broadly. Which is it?

Chapter 7 shows that the inverse relationships between differentiation and integration hold very broadly. Moreover, these relationships will give us a new way of computing integrals. There is lots of cool stuff to come!

Function ($f(x)$)	Derivative ($f'(x)$)	Monotonicity
x^3	$3x^2$	Strictly Increasing
$-x^3$	$-3x^2$	Strictly Decreasing
$\ln(x)$	$\frac{1}{x}$	Strictly Increasing for $x > 0$
e^x	e^x	Strictly Increasing
e^{-x}	$-e^{-x}$	Strictly Decreasing
$\tan(x)$	$1 + \tan^2(x)$	Strictly Increasing (in its domain)
$\frac{1}{1+x^2}$	$-\frac{2x}{(1+x^2)^2}$	Strictly Decreasing for $x \geq 0$
$\frac{1}{1+x^2}$	$-\frac{2x}{(1+x^2)^2}$	Strictly Increasing for $x \leq 0$
$\frac{x^2}{x^2+1}$	$\frac{2x}{(x^2+1)^2}$	Strictly Increasing for $x \geq 0$

Table 5.22: Monotonicity of functions through the lens of their derivatives.

5.5.2 If the Derivative does not Change Sign, the Function is Monotonic

In Chapter 1.5.5, we learned that a strictly monotonic function (increasing or decreasing) has an inverse as long as the codomain of the function has been chosen properly (it must equal the function's range). We address here how to use a function's derivative to check monotonicity, both strict and non-strict.

Let $I \subset \mathbb{R}$ be an open interval and let $f : I \rightarrow \mathbb{R}$ be a function on I . For convenience, we recall the various forms of monotonicity:

- f is **nondecreasing** if $y > x \implies f(y) \geq f(x)$;
- f is **strictly increasing** if $y > x \implies f(y) > f(x)$. Similarly,
- f is **nonincreasing** if $x < y \implies f(x) \leq f(y)$;
- f is **strictly decreasing** if $x < y \implies f(x) < f(y)$.
- **Note:** A constant function is both nondecreasing and nonincreasing but is neither strictly increasing nor strictly decreasing.

Proposition 5.23: Monotonic Functions

Let $I \subset \mathbb{R}$ be an open interval and let $f : I \rightarrow \mathbb{R}$ be a differentiable function on I . The following statements are true.

- If $\frac{df(x)}{dx} \geq 0$ for all $x \in I$, then f is nondecreasing.
- If $\frac{df(x)}{dx} \leq 0$ for all $x \in I$, then f is nonincreasing.
- If $\frac{df(x)}{dx} > 0$ for all $x \in I$, then f is strictly increasing.
- If $\frac{df(x)}{dx} < 0$ for all $x \in I$, then f is strictly decreasing.

The last two statements can be relaxed just a bit.

- If $\frac{df(x)}{dx} = 0$ for a finite number of $x \in I$ and $\frac{df(x)}{dx} > 0$ elsewhere, then f is strictly increasing.
- If $\frac{df(x)}{dx} = 0$ for a finite number of $x \in I$ and $\frac{df(x)}{dx} < 0$ elsewhere, then f is strictly decreasing.

These “friendly amendments” are important for a function such as x^3 , where the derivative $3x^2$ vanishes at the origin. Otherwise, you would be saying the function is strictly increasing on $(-\infty, 0)$ and $(0, \infty)$, and be perplexed about what is happening for $x = 0$. That would not be good!

Several examples of Prop. 5.23 are given in Table 5.22. Traditionally, a proof of Prop. 5.23 is based on the Fundamental Theorem of Calculus, developed later in the book. Fortunately for us, the limit definition of a derivative provides good insight into why it is true! Indeed, consider

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0^+} \frac{f(x+h) - f(x)}{h},$$

where $x + h > x$ plays the role of $y > x$ because h is greater than zero when we take the limit from the right. For the derivative to be strictly greater than zero at x , it must be true that $f(x+h) > f(x)$. And then $f(x+2h) > f(x+h)$ follows from $\frac{df(x+h)}{dx} > 0$, and so forth. While this reasoning is only qualitative, you get the idea.

Secrets of the Arcane 5.24: Once Upon A Time ...

... you could purchase mathematical results and claim them as your own! L'Hôpital's Rule is named for the French mathematician Guillaume-François-Antoine, Marquis de L'Hôpital, who purchased the formula from his teacher, the Swiss mathematician Johann Bernoulli. Recall that the number e was discovered by the Swiss mathematician Jacob Bernoulli. They were **brothers**. Can you imagine buying some math results to improve your standing in society?

Note: L'Hôpital is also spelled L'Hospital. In French, the circumflex is used to indicate that an “s” has been dropped. The circumflex is also present in the French word “hôtel”, which we use in English as “hotel”, but we also have “hostel”, as in a “Youth Hostel”. The ways we have incorporated French into English can be very confusing. One example: the French word “entrée” means “starter” or “appetizer” on a “carte (menu)” in France. We, of course, use it as “plat principal”, aka the main dish. Yikes!

5.5.3 L'Hôpital's Rule for Evaluating Limits

L'Hôpital's Rule is a powerful tool used to evaluate limits that take the form $\frac{0}{0}$ or $\frac{\pm\infty}{\pm\infty}$. Such forms are called **indeterminate**. Named after the French mathematician Guillaume de l'Hôpital, this rule simplifies the process of finding such limits by differentiating the numerator and the denominator of the function.

Proposition 5.25: L'Hôpital's Rule

Let $I := (a, b)$ be an open interval, $x_0 \in I$ a point in I , and $f : I \rightarrow \mathbb{R}$ and $g : I \rightarrow \mathbb{R}$ two functions that are differentiable on I except possibly at x_0 . If f and g satisfy the additional conditions,

- (a) For all $x \in I$, $x \neq x_0$, $g'(x) \neq 0$.
- (b) $\lim_{x \rightarrow x_0} f(x) = \lim_{x \rightarrow x_0} g(x) = 0$ or $|\lim_{x \rightarrow x_0} f(x)| = |\lim_{x \rightarrow x_0} g(x)| = \infty$ (means the limits of f and g can be $\pm\infty$).
- (c) $\lim_{x \rightarrow x_0} \frac{f'(x)}{g'(x)}$ exists (can be finite or infinite, but it must exist).

Then,

$$\lim_{x \rightarrow x_0} \frac{f(x)}{g(x)} = \lim_{x \rightarrow x_0} \frac{f'(x)}{g'(x)}. \quad (5.16)$$

Notes:

- In the above we have expressed the result using two-sided limits. It also works for one-sided limits, $x \rightarrow x_0^-$ or $x \rightarrow x_0^+$, and for $x \rightarrow \pm\infty$.
- What if $\lim_{x \rightarrow x_0} \frac{f'(x)}{g'(x)}$ is also indeterminate? Then you can apply L'Hôpital's Rule to $\frac{f'(x)}{g'(x)}$ and work with the second derivatives, $\frac{f''(x)}{g''(x)}$, etc. This can go on until you run out of steam or get an answer. Usually, one application suffices.

Example 5.26. Consider the function $\frac{\sin(x)}{x}$ as x approaches 0. Both the numerator and the denominator approach 0, placing us in the indeterminate form $\frac{0}{0}$. Can we apply L'Hôpital's Rule?

Solution: Differentiating the numerator (aka, $f(x) = \sin(x)$), we obtain $\cos(x)$, and differentiating the denominator (aka, $g(x) = x$), yields 1, which is non-zero. Therefore, the limit becomes,

$$\lim_{x \rightarrow 0} \frac{\cos(x)}{1} = \cos(0) = 1.$$

Because all the conditions for L'Hôpital's Rule are met, we conclude

$$\lim_{x \rightarrow 0} \frac{\sin(x)}{x} = \lim_{x \rightarrow 0} \frac{\cos(x)}{1} = \cos(0) = 1.$$

■

Bottom line: L'Hôpital's Rule is a convenient method for evaluating limits of indeterminate forms. It saves time and simplifies the process by allowing us to differentiate the numerator and the denominator to find the limit.

Another **indeterminate form** is $0 \cdot \infty$. Can L'Hôpital's Rule help there as well? Yes, because we can rewrite it in the standard form treated in Prop. 5.25, as we illustrate next.

Example 5.27. Determine the following limits.

$$(a) \lim_{x \rightarrow \infty} xe^{-x}.$$

$$(b) \lim_{x \rightarrow \infty} \frac{x^2 + 3x + 2}{2x^2 + 4x + 1}$$

$$(c) \lim_{x \rightarrow 0} \frac{\ln(1+x)}{x}$$

$$(d) \lim_{x \rightarrow \infty} x \ln\left(\frac{1}{x}\right)$$

$$(e) \lim_{x \rightarrow 0} \frac{e^x - \cos(x)}{x}$$

Note: The first two problems can be solved using the results of Chapter 2.6 and Chapter 2.6.3. It's good to see them through the lens of L'Hôpital. Part (d) of the problem could **easily trip you up**. Definitely take a look at it for yourself before reading the solutions.

Solutions:

$$(a) \text{Ans. } \lim_{x \rightarrow \infty} xe^{-x} = 0.$$

This problem has the indeterminate form $\infty \cdot 0$. We rewrite it in the form $\frac{\infty}{\infty}$, via

$$\begin{aligned} \lim_{x \rightarrow \infty} xe^{-x} &= \lim_{x \rightarrow \infty} \frac{x}{e^x} \quad (\text{reciprocal property of exponentials}) \\ &= \lim_{x \rightarrow \infty} \frac{1}{e^x} \quad (\text{differentiate numerator and denominator}) \\ &= 0. \end{aligned}$$

We could have also rewritten it in the form $\frac{0}{0}$, but as you will see, it leaves us with a harder problem to solve in this case. In other problems, the opposite is true. Hence, there is "art" in the application of mathematical tools! And thank goodness for that! Because if there weren't any art involved (aka, you could just turn a crank), a bot could easily take our jobs as engineers.

$$\begin{aligned} \lim_{x \rightarrow \infty} xe^{-x} &= \lim_{x \rightarrow \infty} \frac{e^{-x}}{\frac{1}{x}} \quad (x = \frac{1}{1/x}) \\ &= \lim_{x \rightarrow \infty} \frac{-e^{-x}}{\frac{-1}{x^2}} \quad (\text{differentiate numerator and denominator}) \\ &= \text{a mess!} \end{aligned}$$

(b) **Ans.** $\lim_{x \rightarrow \infty} \frac{x^2 + 3x + 2}{2x^2 + 4x + 1} = \frac{1}{2}$

This problem uses L'Hôpital's Rule twice. Its indeterminate form is $\frac{\infty}{\infty}$, so it falls directly into the L'Hôpital framework.

$$\begin{aligned}\lim_{x \rightarrow \infty} \frac{x^2 + 3x + 2}{2x^2 + 4x + 1} &= \lim_{x \rightarrow \infty} \frac{2x + 3}{4x + 4} \quad (\text{differentiate numerator and denominator}) \\ &= \lim_{x \rightarrow \infty} \frac{2}{4} \quad (\text{differentiate numerator and denominator once more}) \\ &= \frac{1}{2}.\end{aligned}$$

(c) **Ans.** $\lim_{x \rightarrow 0} \frac{\ln(1+x)}{x} = 1$.

The indeterminate form of this problem is $\frac{0}{0}$, so it falls directly into the L'Hôpital framework.

$$\begin{aligned}\lim_{x \rightarrow 0} \frac{\ln(1+x)}{x} &= \lim_{x \rightarrow 0} \frac{\frac{1}{1+x}}{1} \quad (\text{differentiate numerator and denominator}) \\ &= 1.\end{aligned}$$

(d) **Ans.** $\lim_{x \rightarrow \infty} x \ln\left(\frac{1}{x}\right) = -\infty$.

This is not an indeterminate form², and thus L'Hôpital's Rule does not apply. It is still valid to assign $f(x) := x$ and $g(x) := \ln\left(\frac{1}{x}\right)$. Then, we note,

$$\begin{aligned}\lim_{x \rightarrow \infty} f(x) &= \lim_{x \rightarrow \infty} x \\ &= \infty\end{aligned}$$

$$\begin{aligned}\lim_{x \rightarrow \infty} g(x) &= \lim_{x \rightarrow \infty} \ln\left(\frac{1}{x}\right) \\ &= \lim_{x \rightarrow \infty} -\ln(x) \quad (\text{logarithm rules}) \\ &= -\infty.\end{aligned}$$

Hence,

$$\lim_{x \rightarrow \infty} f(x) \cdot g(x) = -\infty,$$

because something large and positive times something large and negative is something (super) large and negative.

Now, suppose we skipped the above analysis and applied L'Hôpital's Rule **blindly!** We might do something like this

$$\begin{aligned}\lim_{x \rightarrow \infty} x \ln\left(\frac{1}{x}\right) &= \lim_{x \rightarrow \infty} \frac{\ln\left(\frac{1}{x}\right)}{\frac{1}{x}} \quad (x = \frac{1}{1/x}) \\ &= \lim_{x \rightarrow \infty} \frac{x \cdot \frac{-1}{x^2}}{\frac{-1}{x^2}} \quad (\text{differentiate numerator and denominator with the Chain Rule}) \\ &= \lim_{x \rightarrow \infty} x \quad (\text{cancel like terms in the numerator and denominator}) \\ &= \infty.\end{aligned}$$

That means we blew the sign on the limit! Or, we might even say, our error is $-2 \cdot \infty$ ($-\infty - \infty$, actual limit minus computed limit). Crazy! But, that is what happens when you misapply L'Hopital's Rule, as ChatGPT4 did for your author! You gotta

²Plus infinity times minus infinity is not indeterminate. To be indeterminate, we need $\frac{0}{0}$ or $\frac{\pm\infty}{\pm\infty}$.

be careful.

Here is another VALID solution using logarithm rules.

$$\begin{aligned}\lim_{x \rightarrow \infty} x \ln\left(\frac{1}{x}\right) &= \lim_{x \rightarrow \infty} x \cdot (\ln(1) - \ln(x)) \quad (\ln(a/b) = \ln(a) - \ln(b)) \\ &= \lim_{x \rightarrow \infty} -x \cdot \ln(x) \quad (\ln(1) = 0) \\ &= -\infty,\end{aligned}$$

because, $\lim_{x \rightarrow \infty} x = \infty$ and $\lim_{x \rightarrow \infty} \ln(x) = \infty$.

(e) **Ans.** $\lim_{x \rightarrow 0} \frac{e^x - \cos(x)}{x} = 1$

This problem's indeterminate form is $\frac{0}{0}$, so it falls directly into the L'Hôpital framework.

$$\begin{aligned}\lim_{x \rightarrow 0} \frac{e^x - \cos(x)}{x} &= \lim_{x \rightarrow 0} \frac{e^x + \sin(x)}{1} \quad (\text{differentiate numerator and denominator}) \\ &= 1.\end{aligned}$$

Fact 5.28: Need to Apply L'Hôpital's Rule with Care

All four conditions for L'Hôpital's rule are **necessary** (from Wikipedia):

- **Indeterminacy of form:** $\lim_{x \rightarrow x_0} f(x) = \lim_{x \rightarrow x_0} g(x) = 0$ or $\pm\infty$;
- **Differentiability of functions:** $f(x)$ and $g(x)$ are differentiable on an open interval I except possibly at a point x_0 contained in I (the same point used in the limit);
- **Non-zero derivative of denominator:** $g'(x) \neq 0$ for all x in I with $x \neq x_0$; and
- **Existence of limit of the quotient of the derivatives:** $\lim_{x \rightarrow x_0} \frac{f'(x)}{g'(x)}$ exists.

If even one of the above conditions is not satisfied, L'Hôpital's rule is not valid (i.e., it **CANNOT** be applied).

Example 5.29. Compute $\lim_{x \rightarrow \infty} \sqrt[3]{x}$.

Solution: 1 in 1 Minute by BlackPenRedPen (do not hesitate to adjust the playback speed). ■

Want more examples?

- **L'Hôpital's Rule: Ultimate Study guide** by BlackPenRedPen
- **L'Hopital's rule** by The Organic Chemistry Tutor.
- **What is L'Hospital's rule?** by Krista King. (She also talks about Johann Bernoulli.)

1 # Reminder about our function

2

3 using Symbolics

4

5 # For scalar functions

```

6 function deriv(f::Symbolics.Num, x::Symbolics.Num)
7     df = Differential(x)(f)
8     return simplify(expand_derivatives(df)) # Note the simplify command
9 end
10
11 # For vector-valued functions
12 function deriv(f::Vector{Symbolics.Num}, x::Symbolics.Num)
13     df = []
14     for component in f
15         df_comp = Differential(x)(component)
16         push!(df, simplify(expand_derivatives(df_comp))) # Note the simplify command
17     end
18     return df
19 end

```

Output

```
deriv (generic function with 2 methods)
```

```

1 @variables x1 x2 x3
2
3 fa = [x1 * (x2)^2; exp(3*x1) + cos(4*x2)]
4 fb = exp(3*x1) + cos(4*x2) # not a vector, hence no brackets!
5 fc = [x1 * (x2)^2; x2*(sin(x1))^3 + 11; 4*(x1)^2 + (x2 + 3)^2]
6 fd = [x1*x2*(x3)^3; x2*exp(x1*x3); x1]
7
8
9
10 @show dfadx1 = deriv.(fa, x1)
11 @show dfadx2 = deriv.(fa, x2)
12 println(" ")
13 @show dfbdx1 = deriv.(fb, x1)
14 @show dfbdx2 = deriv.(fb, x2)
15 println(" ")
16 @show dfcdx1 = deriv.(fc, x1)
17 @show dfccx2 = deriv.(fc, x2)
18 println(" ")
19 @show dfddx1 = deriv.(fd, x1) # the dfddxi is really dfd/dxi
20 @show dfddx2 = deriv.(fd, x2)
21 @show dfddx3 = deriv.(fd, x3);

```

Output

```
dfadx1 = deriv.(fa, x1) = Num[x2^2, 3exp(3x1)]
dfadx2 = deriv.(fa, x2) = Num[2x1*x2, -4sin(4x2)]

dfbdx1 = deriv.(fb, x1) = 3exp(3x1)
dfbdx2 = deriv.(fb, x2) = -4sin(4x2)

dfcdx1 = deriv.(fc, x1) = Num[x2^2, 3x2*(sin(x1)^2)*cos(x1), 8x1]
dfccx2 = deriv.(fc, x2) = Num[2x1*x2, sin(x1)^3, 2(3 + x2)]

dfddx1 = deriv.(fd, x1) = Num[x2*(x3^3), x2*x3*exp(x1*x3), 1]
dfddx2 = deriv.(fd, x2) = Num[x1*(x3^3), exp(x1*x3), 0]
dfddx3 = deriv.(fd, x3) = Num[3x1*x2*(x3^2), x1*x2*exp(x1*x3), 0]
```

Explaining the Num Type in Julia

I asked ChatGpt4 to explain the Num Type in the expression `Num[x22, 3 · x2 · (sin(x1)2) · cos(x1), 8x1]`. Here's what it said. “**Num:** In the context of Julia and symbolic computation: Num is a type that holds symbolic numbers or expressions. These are not evaluated until you explicitly ask Julia to do so. They can be variables, standard mathematical functions, or more complex expressions.” ChatGPT4 did not say this, but it should have: **The key thing is, you can continue taking derivatives like so,**

```
1 @show d2fadx1dx1 = deriv.(dfadx1, x1)
2 @show d2fadx1dx2 = deriv.(dfadx1, x2)
3 println(" ")
4 @show d2fadx2dx1 = deriv.(dfadx2, x1)
5 @show d2fadx2dx2 = deriv.(dfadx2, x2)
6 println(" ")
```

Output

```
d2fadx1dx1 = deriv.(dfadx1, x1) = Num[0, 9exp(3x1)]
d2fadx1dx2 = deriv.(dfadx1, x2) = Num[2x2, 0]

d2fadx2dx1 = deriv.(dfadx2, x1) = Num[2x2, 0]
d2fadx2dx2 = deriv.(dfadx2, x2) = Num[2x1, -16cos(4x2)]
```

where the correct mathematical notation would be

$$\frac{\partial^2 f_a(x_1, x_2)}{\partial x_1^2} = \frac{\partial^2 f_a(x_1, x_2)}{\partial x_1 \partial x_1} = \begin{bmatrix} 0 \\ 9e^{3x_1} \end{bmatrix} \text{ and } \frac{\partial^2 f_a(x_1, x_2)}{\partial x_1 \partial x_2} = \begin{bmatrix} 2x_2 \\ 0 \end{bmatrix}$$
$$\frac{\partial^2 f_a(x_1, x_2)}{\partial x_2 \partial x_1} = \begin{bmatrix} 2x_2 \\ 0 \end{bmatrix} \text{ and } \frac{\partial^2 f_a(x_1, x_2)}{\partial x_2^2} = \frac{\partial^2 f_a(x_1, x_2)}{\partial x_2 \partial x_2} = \begin{bmatrix} 2x_1 \\ -16\cos(4x_2) \end{bmatrix}.$$

And then you take derivatives of those expressions as well; see [Introduction to Higher Order Partial Derivatives Notation and Example](#) by the Math Sorcerer. **Our final remark here is the following:**

Is it OK to use $\frac{\partial}{\partial x_1}$ for computing the derivative of a function that only depends on x_1 ? Technically, it's not wrong, but it'll make you stand out in an uncomfortable sort of way! It's better to show your finesse and adopt the “customary” notation, $\frac{d}{dx}$. And besides, other Profs and GSIs **might mark you wrong**. Just don't do it!

5.5.4 Higher-order Derivatives Lead to Taylor's and Maclaurin's Expansions

In Chapter 5.2, we showed how the notion of a linear approximation to a function near a point and the derivative of the function at that same point go hand in hand; they are really two sides of the same coin. As hinted at in Fig. 5.5, why stop at a linear approximation of a function? Why not use higher-order derivatives to propose a **polynomial approximation** that may be able to provide better accuracy? Progress on this question began 700 years ago by the Indian mathematician, Madhava of Sangamagrama, and then accelerated under “the relatively more recent” investigation by Issac Newton in the mid to late 1600’s. The general theory came together under Brook Taylor in 1715. The importance of Taylor’s work became more clear when Colin Maclaurin used it to characterize maxima, minima, and inflection points of functions in 1725. Today, the uses of Taylor approximations range from the creation of faster and more accurate numerical integration routines to highly efficient **ray tracing** algorithms that render photo-realistic images in video games.

Fundamental Idea

We start with a very simple proposition that lays bare the essence of Taylor series.

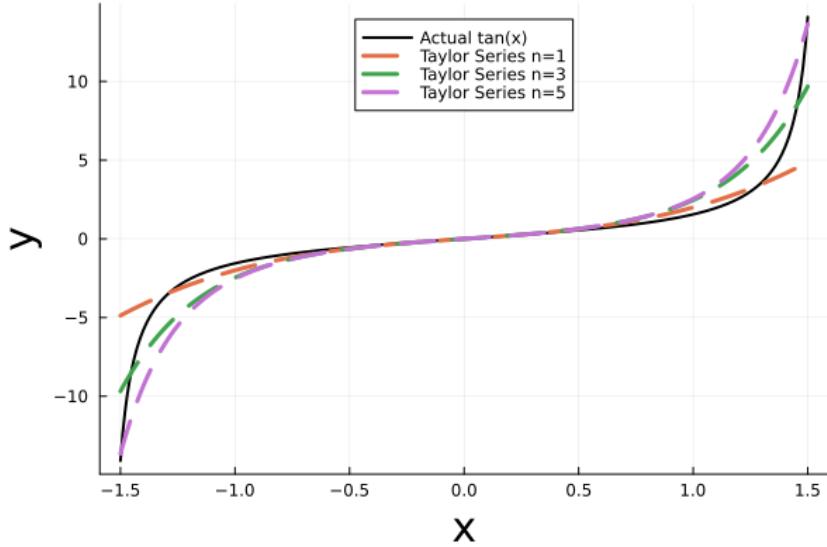


Figure 5.5: (**Taylor-Maclaurin Expansion for Tangent**) Taylor and Maclaurin studied polynomial approximations of functions based on higher-order derivatives. As such, their work extends the notion of a function's linearization, which only uses the first derivative. Because of Taylor series, people often say “**first-order Taylor approximation**” when they really mean the “**linearization of a function**”. Hey, it is what it is!

Proposition 5.30: Local Polynomial Approximation

Consider a function $f : (a, b) \rightarrow \mathbb{R}$ that is n -times differentiable at $0 \in (a, b)$, the origin. Then the polynomial

$$p(x) := a_n \frac{x^n}{n!} + a_{n-1} \frac{x^{(n-1)}}{(n-1)!} + \cdots + a_2 \frac{x^2}{2!} + a_1 x + a_0 \quad (5.17)$$

matches the values of f and its first n derivatives at the origin if, and only if,

$$a_k = f^{(k)}(0) = \left. \frac{d^k f(x)}{dx^k} \right|_{x=0}. \quad (5.18)$$

Note: We have added factorial terms to the monomials. This is because

$$\frac{d}{dx} \left(\frac{x^k}{k!} \right) = \frac{x^{k-1}}{(k-1)!}, \quad \frac{d^2}{dx^2} \left(\frac{x^k}{k!} \right) = \frac{x^{k-2}}{(k-2)!}, \quad \frac{d^3}{dx^3} \left(\frac{x^k}{k!} \right) = \frac{x^{k-3}}{(k-3)!}, \quad \text{etc.}$$

Note 2: A similar result holds for

$$q(x) := \sum_{k=0}^n a_k \frac{(x-x_0)^k}{k!} = a_n \frac{(x-x_0)^n}{n!} + \cdots + a_2 \frac{(x-x_0)^2}{2!} + a_1 (x-x_0) + a_0,$$

that is, $p(x)$ shifted by x_0 . Then we need $a_k = f^{(k)}(x_0)$.

Proof: It is super informative to look at the second-order case first. We are interested in constructing a polynomial

$$p(x) = a_2 \frac{x^2}{2!} + a_1 x + a_0$$

that matches the function f and its first two derivatives at the origin.

Function Value at the Origin: For the polynomial $p(x)$ to match the value of $f(x)$ at $x = 0$, we must have

$$p(0) = a_0 = f(0).$$

This establishes that the constant term a_0 is equal to $f(0)$, the value of the function at the origin.

First Derivative: The first derivative of $p(x)$ is

$$p'(x) = 2a_2 \frac{x}{2!} + a_1 = a_2 x + a_1.$$

Matching the first derivative of $p(x)$ with that of $f(x)$ at $x = 0$ yields

$$p'(0) = a_1 = f'(0).$$

This shows that the coefficient a_1 is determined by the first derivative of f at the origin.

Second Derivative: The second derivative of $p(x)$ is

$$p''(x) = a_2.$$

Matching the second derivative of $p(x)$ with that of $f(x)$ at $x = 0$ gives

$$p''(0) = a_2 = f''(0).$$

This confirms that the coefficient a_2 is equal to the second derivative of f at the origin.

(Optional Read) Proof by Induction: **Base Case:** For $k = 0$, we have $p(0) = f(0)$, which implies $a_0 = f(0)$. This is because all other terms in $p(x)$ contain x and therefore vanish at $x = 0$.

Induction Hypothesis: Assume that for some $k \geq 0$, the k th derivative of $p(x)$ equals

$$p^{(k)}(x) = a_k + a_{k+1} \frac{x}{1!} + a_{k+2} \frac{x^2}{2!} + \cdots + a_n \frac{x^{n-k}}{(n-k)!},$$

and that it matches the k th derivative of $f(x)$ at 0, i.e.,

$$p^{(k)}(0) = f^{(k)}(0) = a_k.$$

Inductive Step: We must show that if the hypothesis holds for k , it also holds for $k + 1$. Consider the $(k + 1)$ st derivative of $p(x)$:

$$p^{(k+1)}(x) = a_{k+1} + a_{k+2} \frac{x}{1!} + \cdots + a_n \frac{x^{n-(k+1)}}{(n-(k+1))!}.$$

Evaluating this at $x = 0$ gives

$$p^{(k+1)}(0) = a_{k+1},$$

which, by the definition of $p(x)$, must equal $f^{(k+1)}(0)$. Therefore, for the polynomial $p(x)$ to match f and its first n derivatives at 0, we have

$$a_{k+1} = f^{(k+1)}(0) := \left. \frac{d^{k+1} f(x)}{dx^{k+1}} \right|_{x=0}.$$

By the principle of mathematical induction, this completes the proof, showing that each coefficient a_k of the polynomial $p(x)$ must be chosen as $a_k = f^{(k)}(0)$ for the polynomial to match f and its first n derivatives at the origin. ■

The Actual Taylor and Maclaurin Polynomials

Taylor and Maclaurin Polynomials

Let $f : (a, b) \rightarrow \mathbb{R}$ be a function defined on an open interval (a, b) and let $x_0 \in (a, b)$. Recall that if f is once differentiable at x_0 , then

$$f(x) \approx f(x_0) + \frac{df(x_0)}{dx} \cdot (x - x_0) \text{ or with the prime notation } f(x) \approx f(x_0) + f'(x_0) \cdot (x - x_0)$$

is the **linear approximation** of the function about x_0 .

Definition 5.31. Suppose that $f : (a, b) \rightarrow \mathbb{R}$ is n -times differentiable at x_0 . Then

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0) + f''(x_0) \cdot \frac{(x - x_0)^2}{2!} + \dots + f^{(n)}(x_0) \cdot \frac{(x - x_0)^n}{n!}, \quad (5.19)$$

is the **n -th degree Taylor expansion of f about the point x_0** . For the special case $x_0 = 0$ (the origin),

$$f(x) \approx f(0) + f'(0) \cdot x + f''(0) \cdot \frac{x^2}{2!} + \dots + f^{(n)}(0) \cdot \frac{x^n}{n!} \quad (5.20)$$

is called an **n -th degree Maclaurin polynomial or Maclaurin expansion**.

Note: The Taylor expansion is the unique polynomial of degree n that exactly matches the values of f and its first n derivatives at the point x_0 .

Note 2: There is a nice theory for quantifying the approximation error when using a Taylor expansion. Here is a very clear [written document](#) by the University of South Carolina and here is a [video](#) by the Organic Chemistry Tutor.

Taking it to the Limit

Many functions can be differentiated an infinite number of times. For example, e^x is **infinitely differentiable** because $\frac{d}{dx} e^x = e^x$ implies that $\frac{d^k}{dx^k} e^x = e^x$ for all $k \geq 1$. Sine and cosine are also infinitely differentiable because of the simple chain of derivatives,

$$\frac{d}{dx} \sin(x) = \cos(x) \quad \text{and} \quad \frac{d}{dx} \cos(x) = -\sin(x).$$

In fact, the k -th derivative of $\sin(x)$ is

$$\frac{d^k}{dx^k} \sin(x) = \begin{cases} \sin(x) & \text{if } k \equiv 0 \pmod{4} \\ \cos(x) & \text{if } k \equiv 1 \pmod{4} \\ -\sin(x) & \text{if } k \equiv 2 \pmod{4} \\ -\cos(x) & \text{if } k \equiv 3 \pmod{4} \end{cases}$$

and the k -th derivative of $\cos(x)$ is

$$\frac{d^k}{dx^k} \cos(x) = \begin{cases} \cos(x) & \text{if } k \equiv 0 \pmod{4} \\ -\sin(x) & \text{if } k \equiv 1 \pmod{4} \\ -\cos(x) & \text{if } k \equiv 2 \pmod{4} \\ \sin(x) & \text{if } k \equiv 3 \pmod{4} \end{cases}$$

For functions that are infinitely differentiable, one can propose an **infinite** Taylor or Maclaurin series representation of the function.

Fact 5.32: (Infinite) Taylor Series Expansion

The Taylor series expansion of an infinitely differentiable function $f(x)$ around a point x_0 is given by,

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \frac{f'''(x_0)}{3!}(x - x_0)^3 + \dots = \lim_{n \rightarrow \infty} \sum_{k=0}^n \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k, \quad (5.21)$$

where $f^{(k)}(x_0)$ denotes the k -th derivative of f evaluated at the point x_0 .

Whether or not the limit in an infinite Taylor series converges is a highly technical question that we will not explore. Nevertheless, there are three CONVERGENT infinite Maclaurin series that show up frequently in Engineering and STEM. We highlight them below:

Expansion for e^x

$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots = \lim_{n \rightarrow \infty} \sum_{k=0}^n \frac{x^k}{k!}, \quad (5.22)$$

Expansion for $\sin(x)$

$$\sin(x) = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots = \lim_{n \rightarrow \infty} \sum_{k=0}^n \frac{(-1)^k x^{2k+1}}{(2k+1)!}, \quad (5.23)$$

Expansion for $\cos(x)$

$$\cos(x) = 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \frac{x^6}{6!} + \dots = \lim_{n \rightarrow \infty} \sum_{k=0}^n \frac{(-1)^k x^{2k}}{(2k)!}. \quad (5.24)$$

If we were to add a fourth one, it would be the expansion for $\ln(1+x)$.

Deriving Euler's Formula from Fact 5.32

Euler's formula, $e^{ix} = \cos(x) + i\sin(x)$, follows from the three Taylor series in Fact 5.32. Here is how it goes. We first DEFINE e^{ix} by substituting ix for x into (5.22). Yes, otherwise, we do not know how to evaluate the function with a complex argument. However, we know how to perform algebra on terms like $(ix)^k$.

This approach to defining e^{ix} through a power series is in and of itself a powerful use of Taylor series; continuing, we have

$$e^{ix} := 1 + ix + \frac{(ix)^2}{2!} + \frac{(ix)^3}{3!} + \frac{(ix)^4}{4!} + \dots$$

Then, because $i^k = \begin{cases} 1 & k = 0 \pmod 4 \\ i & k = 1 \pmod 4 \\ -1 & k = 2 \pmod 4 \\ -i & k = 3 \pmod 4, \end{cases}$ the above simplifies to,

$$e^{ix} = 1 + ix - \frac{x^2}{2!} - i\frac{x^3}{3!} + \frac{x^4}{4!} + i\frac{x^5}{5!} \dots$$

Grouping the real and imaginary parts, yields

$$e^{ix} = \left(1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \dots\right) + i\left(x - \frac{x^3}{3!} + \frac{x^5}{5!} - \dots\right)$$

The series within the parentheses are recognized as the Taylor series expansions of $\cos(x)$ and $\sin(x)$, respectively,

$$\cos(x) = 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \dots$$

$$\sin(x) = x - \frac{x^3}{3!} + \frac{x^5}{5!} - \dots$$

Therefore, by direct comparison, we obtain

$$e^{ix} = \cos(x) + i\sin(x).$$

■

Don't Fall for the Hype

Thanks to ROB 101 *Computational Linear Algebra*, you know how to fit polynomials to functions. You also know how to build approximations with radial basis functions. For many applications, these sorts of approximations are superior to Taylor expansions. Taylor expansions are one tool among many in your algorithmic toolbox.

Secrets of the Arcane 5.33: Taylor Series are often Over Hyped

When the pinnacle of calculus applications you have to offer a student is the Taylor series, consider yourself in dire straits. This blunt truth cuts through the academic grandeur that has elevated the Taylor series to a pedestal for over seven decades. Despite its theoretical elegance, the harsh reality is that, beyond the realm of first-order linearizations, Taylor series rarely grace the practical world of engineering. This stark contrast exposes a glaring discrepancy between traditional calculus education and the actual demands of engineering practice. The hyping of Taylor series as a major application of Calculus not only misleads but also underprepares students for the challenges they will face in applying Calculus creatively outside the classroom. In essence, if the Taylor series is your go-to case study for the wonders of Calculus, you're not just propagating an untruth; you're supporting a profound educational failure that has persisted for far too long.

For Further Learning

- Higher order derivatives by [3Blue1Brown](#).
- Taylor series by [3Blue1Brown](#).
- Taylor & Maclaurin Polynomials: Intro Khan Academy.
- Taylor Series and Power Series Made Easy (with Pictures): Review of Calculus by Steve Brunton.
- Taylor Series of the Exponential Function and Euler's Formula! by Steve Brunton.

m	Expression for $\text{Im}\{(1+i)^m\}$	Answer
0	$\text{Im}\{(1+i)^0\}$	0
1	$\text{Im}\{(1+i)^1\}$	1
2	$\text{Im}\{(1+i)^2\}$	2
3	$\text{Im}\{(1+i)^3\}$	2
4	$\text{Im}\{(1+i)^4\}$	0
5	$\text{Im}\{(1+i)^5\}$	-2^2
6	$\text{Im}\{(1+i)^6\}$	-2^3
7	$\text{Im}\{(1+i)^7\}$	-2^3

Table 5.34: A Y2K generalization can be found in this table of values for m , expressions for $\text{Im}\{(1+i)^m\}$, and computed numbers.

5.5.5 (Optional Read:) Y2K or BlackPenRedPen Strikes Again

While the original [Y2K Problem](#) was of potential major significance for our engineered world, the following problem has zero engineering relevance. Nevertheless, many of us got into engineering because we liked math. Hence, scratching that itch is a worthwhile endeavor.

Problem: Compute $\left. \frac{d^{2000}}{dx^{2000}} (e^x \cdot \sin(x)) \right|_{x=0}$. Yes, the Y2K derivative of $e^x \cdot \sin(x)$ evaluated at the origin.

Ans. 1 BCE³.

Why: Recall Euler's formula: $e^{ix} = \cos(x) + i \sin(x)$. Hence, we can write

$$\sin(x) = \text{Im}\{e^{ix}\},$$

³The year immediately before 1 CE (Common Era) is 1 BCE (Before Common Era). There is no year 0 in the Gregorian calendar system, so it jumps directly from 1 BCE to 1 CE. Does this mean that 1 BCE, being the year before 1 CE, functions as zero?

where Im stands for the imaginary part. Our problem can, therefore, be written as

$$\frac{d^{2000}}{dx^{2000}} (\text{Im}\{e^x \cdot e^{ix}\}) \Big|_{x=0} = \text{Im} \left\{ \frac{d^{2000}}{dx^{2000}} (e^x \cdot e^{ix}) \Big|_{x=0} \right\} = \text{Im} \left\{ \frac{d^{2000}}{dx^{2000}} (e^{(1+i)x}) \Big|_{x=0} \right\}.$$

This greatly simplifies things because

$$\frac{d^k}{dx^k} (e^{\alpha x}) \Big|_{x=0} = \alpha^k \cdot e^{\alpha x} \Big|_{x=0} = \alpha^k,$$

even for $\alpha \in \mathbb{C}$. It follows that,

$$\frac{d^{2000}}{dx^{2000}} (e^x \cdot \sin(x)) \Big|_{x=0} = \text{Im}\{(1+i)^{2000}\},$$

and while this may look painful, it is not! We first note that

$$(1+i)^2 = (1+2i+(i^2)) = (1+2i-1) = 2i.$$

Then, we note that

$$(1+i)^8 = (2i)^4 = 2^4,$$

because $i^4 = (-1)^2 = 1$. Do you see where this is going? $2000 = 8 \cdot 250$. Hence, using our exponent rules,

$$(1+i)^{2000} = ((1+i)^8)^{250} = (2^4)^{250} = 2^{1000},$$

a real number. We conclude that

$$\frac{d^{Y2K}}{dx^{Y2K}} (e^x \cdot \sin(x)) \Big|_{x=0} = \frac{d^{2000}}{dx^{2000}} (e^x \cdot \sin(x)) \Big|_{x=0} = \text{Im}\{(1+i)^{2000}\} = \text{Im}\{2^{1000}\} = 0.$$

Here's the link to **BlackPenRedPen**'s video solving the tiebreaker problems from the 2022 Berkeley Math Tournament. The Berkeley version of the problem is ever so slightly different. After watching the video, can you generalize the problem to

$$\text{Find: } \frac{d^{\text{Year}}}{dx^{\text{Year}}} (e^x \cdot \sin(x)) \Big|_{x=0} ?$$

Hint: Given in the footnote⁴ so that you can play with the problem first. Scratch that itch!

5.6 Partial Derivatives

Because you have taken Linear Algebra already, where you solved problems involving thousands of variables, working with functions of a single variable must seem somewhat “tame”. The training of modern AI Neural Networks involves literally billions of variables where their derivatives are computed via Automatic Differentiation (AD); see Chapter 5.3.3. While we will NOT attempt any examples of this magnitude (look for a course on Deep Learning), we will treat **multivariable** (aka, **multivariate**) functions in $n > 1$ variables. It is suggested that you review Chapter 11.5 in the ROB 101 *Computational Linear Algebra* textbook.

5.6.1 Intuition

Imagine you're standing somewhere on the side of a mountain, and you want to know how fast you're ascending or descending. You can use partial derivatives to figure this out. Partial derivatives are a way of measuring how a function changes when one of its inputs changes while keeping all of the other inputs the same. Since you're standing on a mountain, let's take the function as your altitude (aka, height) and the inputs as your (x, y) position on the mountain.

If you compute the local slope of your function as x varies while holding y constant, it is called a **partial derivative with respect to x** . **Why partial?** Because you are only considering a “small part of the variables” (one at a time) and not “all of them” (together). The partial derivative with respect to x would tell you the rate of change of your altitude as you move a small amount in an eastward direction (i.e., move to the right), **while keeping y constant**. A partial derivative with respect to y would tell you the rate of change of your altitude as you move a small amount northward (i.e., straight ahead), **while keeping x constant**. It is the “slope” or rate

⁴Write $\text{Year} = 8 \cdot N + m$, where $m \in \{0, 1, 2, \dots, 7\}$. Then the answer is $2^{4N} \cdot \text{Im}\{(1+i)^m\}$. You are almost there.

of change taken one variable at a time, a kind of divide and conquer approach. Don't worry; we can still obtain information on the "slope" when all variables are considered.

A partial derivative is a little like how a basis for vectors allows you to understand how much each individual basis vector contributes to a given vector, and by adding up all of the contributions, you can reconstruct the entire vector. Recall, if

$$\{u_1, u_2, \dots, u_n\}$$

is a basis for \mathbb{R}^n , then we can write any vector $v \in \mathbb{R}^n$ as a linear combination of the basis vectors, that is,

$$v = \alpha_1 u_1 + \alpha_2 u_2 + \dots + \alpha_n u_n.$$

The quantity $\alpha_i u_i$ is the **component of v lying in the direction u_i** .

We will see that partial derivatives are EXACTLY LIKE THIS, and in the best possible way. **How's that? Best possible way!** We use the natural basis vectors, $\{e_1, e_2, \dots, e_n\}$, to define partial derivatives. For definiteness, we recall that in \mathbb{R}^n , the **natural basis vectors (aka, canonical basis vectors)** are the columns of the $n \times n$ identity matrix, that is,

$$e_1 := \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, e_2 := \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, e_3 := \begin{bmatrix} 0 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}, \dots, e_n := \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

5.6.2 Partial Derivatives: Single-variable Derivatives meet Multivariable Functions

To set the stage, let's note that if $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a vector-valued function depending on $x = (x_1, x_2, \dots, x_n)$, $x_0 \in \mathbb{R}^n$, and $h \in \mathbb{R}$, then $g(h) := f(x_0 + h \cdot e_i)$ is only a function of h ; it is a single-variable function. As an example, consider

$$f(x_1, x_2) := \begin{bmatrix} x_1 \cdot x_2 \\ e^{x_1} + \sin(x_2) \end{bmatrix}$$

and let $x_0 = (3, \frac{\pi}{4})$. Then

$$g(h) := f(x_0 + h \cdot e_1) = f\left(\begin{bmatrix} 3 \\ \frac{\pi}{4} \end{bmatrix} + h \cdot \begin{bmatrix} 1 \\ 0 \end{bmatrix}\right) = f(3 + h, \frac{\pi}{4}) = \begin{bmatrix} (3 + h) \cdot \frac{\pi}{4} \\ e^{3+h} + \sin\left(\frac{\pi}{4}\right) \end{bmatrix},$$

and

$$\gamma(h) := f(x_0 + h \cdot e_2) = f\left(\begin{bmatrix} 3 \\ \frac{\pi}{4} \end{bmatrix} + h \cdot \begin{bmatrix} 0 \\ 1 \end{bmatrix}\right) = f(3, \frac{\pi}{4} + h) = \begin{bmatrix} 3 \cdot \left(\frac{\pi}{4} + h\right) \\ e^3 + \sin\left(\frac{\pi}{4} + h\right) \end{bmatrix}.$$

Each of these functions depends on h alone. And where h shows up in the function is dictated by which basis vector e_i is used in forming $f(x_0 + h \cdot e_i)$. The first function depends on h everywhere that f depends on x_1 , and the second function depends on h everywhere that f depends on x_2 . Once again, this is a divide-and-conquer strategy that works very well. Moreover, because these functions depend on a single variable, we can use our limit concept from Def. 5.1 to define a partial derivative! **Yes, partial derivatives are “rise over run” computed one variable at a time.**

Partial Derivative at a Point

Definition 5.35. Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a vector-valued function depending on $x = (x_1, x_2, \dots, x_n)$ and that $x_0 \in \mathbb{R}^n$ is a

point (aka, a fixed vector). The **partial derivative of f with respect to x_i at the point x_0** is

$$\begin{aligned}\frac{\partial f(x_0)}{\partial x_i} &:= \lim_{h \rightarrow 0} \frac{f(x_0 + h \cdot e_i) - f(x_0)}{h} \\ \frac{\partial f(x_0)}{\partial \mathbf{x}_i} &:= \lim_{h \rightarrow 0} \frac{f(x_0 + h \cdot \mathbf{e}_i) - f(x_0)}{h} \quad (\text{same as above, highlighting } \mathbf{x}_i \text{ and } \mathbf{e}_i),\end{aligned}\tag{5.25}$$

assuming the limit exists and is finite.

Notes:

- To distinguish a “partial” derivative from an “ordinary” derivative, we replace d with the symbol ∂ , which is pronounced, “partial”.
- The symbol $\frac{\partial f(x_0)}{\partial x_i}$ is called the “partial derivative of f with respect to x_i , evaluated at x_0 . Yeah, that’s a mouthful. That’s why math symbols are useful. A few of them can stand for a ton of words as we saw in Chapter 1.2, when we discussed Descartes vs. the Babylonians.
- Because the function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ has m -components, that is,

$$f(x_1, x_2, \dots, x_n) = \begin{bmatrix} f_1(x_1, x_2, \dots, x_n) \\ f_2(x_1, x_2, \dots, x_n) \\ \vdots \\ f_m(x_1, x_2, \dots, x_n) \end{bmatrix},$$

the limit in (5.25) must exist and be finite for each component of the function, that is, for each $f_i(x_1, x_2, \dots, x_n)$, or the partial derivative does not exist.

Didn’t we just stop using limits to compute the derivative with respect to a single variable? Yes! We compute partial derivatives the same way as ordinary derivatives because we are only differentiating with respect to one variable at a time, as in

$$\frac{\partial f(x_0)}{\partial \mathbf{x}_i} = \frac{d}{dh} f(x_0 + h \cdot \mathbf{e}_i) \Big|_{h=0},$$

and hence, all of the rules in Prop. 5.11 for ordinary derivatives apply to partial derivatives as well. Moreover, we have Software tools for turning gnarly problems into manageable problems. We illustrate all of this next.

Example 5.36. Compute—by hand—the partial derivatives of

$$f(x_1, x_2) := \begin{bmatrix} x_1 \cdot x_2 \\ e^{x_1} + \sin(x_2) \end{bmatrix}$$

at a general point, $x \in \mathbb{R}^2$. In other words, you do not have to evaluate them at a specific point.

Solution:

$$\frac{\partial}{\partial x_1} f(x_1, x_2) = \frac{\partial}{\partial x_1} \begin{bmatrix} x_1 \cdot x_2 \\ e^{x_1} + \sin(x_2) \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial x_1}(x_1 \cdot x_2) \\ \frac{\partial}{\partial x_1}(e^{x_1} + \sin(x_2)) \end{bmatrix} = \begin{bmatrix} x_2 \\ e^{x_1} \end{bmatrix}.$$

Indeed, when taking the derivative with respect to x_1 , we treat everything else as a constant,

- $\frac{\partial}{\partial x_1}(x_1 \cdot x_2) = x_2$ (we think of $x_2 = c$, a constant),
- $\frac{\partial}{\partial x_1}(\sin(x_2)) = 0$ (we think of $x_2 = c$, a constant), and
- $\frac{\partial}{\partial x_1}(e^{x_1}) = e^{x_1}$ (because the derivative of e^{x_1} is itself).

Next,

$$\frac{\partial}{\partial x_2} f(x_1, x_2) = \frac{\partial}{\partial x_2} \begin{bmatrix} x_1 \cdot x_2 \\ e^{x_1} + \sin(x_2) \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial x_2}(x_1 \cdot x_2) \\ \frac{\partial}{\partial x_2}(e^{x_1} + \sin(x_2)) \end{bmatrix} = \begin{bmatrix} x_1 \\ \cos(x_2) \end{bmatrix}.$$

Indeed, when taking the derivative with respect to x_2 , we treat everything else as a constant,

- $\frac{\partial}{\partial x_2}(x_1 \cdot x_2) = x_1$ (we think of $x_1 = c$, a constant),
- $\frac{\partial}{\partial x_2}(e^{x_1}) = 0$ (we think of $x_1 = c$, a constant), and
- $\frac{\partial}{\partial x_2}(\sin(x_2)) = \cos(x_2)$ (because the derivative of $\sin(x_2)$ is $\cos(x_2)$).

Remark 5.37. Why do we sometimes write $\frac{\partial f(x_1, x_2)}{\partial x_1}$ and sometimes $\frac{\partial}{\partial x_1} f(x_1, x_2)$? Because

$$\frac{\partial}{\partial x_1} \begin{bmatrix} x_1 \cdot x_2 \\ e^{x_1} + \sin(x_2) \end{bmatrix} \text{ looks much better than } \frac{\partial \begin{bmatrix} x_1 \cdot x_2 \\ e^{x_1} + \sin(x_2) \end{bmatrix}}{\partial x_1}.$$

Don't you agree? Also, it's good to recognize differentiation as an operator that is applied to functions. ■

5.7 Packaging Partial Derivatives to form Jacobians, Gradients, and Hessians

Jacobians and gradients are matrix and vector forms, respectively, of partial derivatives. They collect all of the “local slope information” into a single object, just as we learned to do for linear equations, $Ax = b$. The Hessian is the Jacobian of the gradient. Much of this material is based on Chapter 11.5 of the ROB 101 *Computational Linear Algebra* textbook.

5.7.1 The Jacobian

Consider a multivariable vector-valued function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$. The notation introduced in Chapter 5.6.2 for partial derivatives allows us to make the following very compact definition.

Jacobian of a Function

Definition 5.38. The **Jacobian** of $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is

$$\frac{\partial f(x)}{\partial x} := \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} & \frac{\partial f(x)}{\partial x_2} & \dots & \frac{\partial f(x)}{\partial x_n} \end{bmatrix}, \quad (5.26)$$

assuming the indicated partial derivatives exist. The Jacobian of a function is formed by packaging the column vectors $\frac{\partial f(x)}{\partial x_j}$ into a matrix. We need to keep in mind that, for each value of $x \in \mathbb{R}^n$, the compact and innocent-looking object

$$\frac{\partial f(x)}{\partial x}$$

is really an $m \times n$ matrix: because there are n columns of the form $\frac{\partial f(x)}{\partial x_j}$, and each column is an m -vector. There are software tools that compute the Jacobian in one fell swoop.

Note: The Jacobian of a function has two common notations (i) $\frac{\partial f(x)}{\partial x}$ as we used above, and (ii) $J_f(x)$.

Just for the record, we write out $\frac{\partial f(x)}{\partial x}$ as an $m \times n$ matrix. Suppose

$$f(x) = \begin{bmatrix} f_1(x) \\ f_2(x) \\ \vdots \\ f_m(x) \end{bmatrix} = \begin{bmatrix} f_1(x_1, x_2, \dots, x_n) \\ f_2(x_1, x_2, \dots, x_n) \\ \vdots \\ f_m(x_1, x_2, \dots, x_n) \end{bmatrix}.$$

Writing out all of the entries in the $m \times n$ Jacobian matrix gives

$$\frac{\partial f(x)}{\partial x} := \left[\begin{array}{cccc} \frac{\partial f_1(x)}{\partial x_1} & \frac{\partial f_1(x)}{\partial x_2} & \dots & \frac{\partial f_1(x)}{\partial x_n} \\ \frac{\partial f_2(x)}{\partial x_1} & \frac{\partial f_2(x)}{\partial x_2} & \dots & \frac{\partial f_2(x)}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m(x)}{\partial x_1} & \frac{\partial f_m(x)}{\partial x_2} & \dots & \frac{\partial f_m(x)}{\partial x_n} \end{array} \right], \quad (5.27)$$

which is much more intimidating than (5.26). In other symbols, the ij component of $J_f(x) = \frac{\partial f(x)}{\partial x}$ is

$$\begin{aligned} \left[\frac{\partial f(x)}{\partial x} \right]_{ij} &= \frac{\partial f_i(x)}{\partial x_j} \\ J_{ij}(x) &= \frac{\partial f_i(x)}{\partial x_j}, \end{aligned} \quad (5.28)$$

If you want to toss out all of the benefits of vector-matrix notation, you can compute each entry of the Jacobian matrix one by one, as in (5.28). When doing hand computations for simple examples, such as $m = n = 2$, the above very explicit 1×1 scalar (i.e., non-vectorized) calculations can be informative.

Proposition 5.39: Linear Approximation at a Point for Multivariable Functions

The linear approximation of a (nonlinear) function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ at a point x_0 is defined to be

$$f(x) \approx f(x_0) + A(x - x_0) := f(x_0) + \frac{\partial f(x_0)}{\partial x}(x - x_0), \quad (5.29)$$

where the $m \times n$ matrix $A := J_f(x_0) = \frac{\partial f(x_0)}{\partial x}$ is the Jacobian of f at the point x_0 . We note that the dimensions make sense because $f(x_0)$ is $m \times 1$, $(x - x_0)$ is $n \times 1$, and therefore $J_f(x_0)(x - x_0) = \frac{\partial f(x_0)}{\partial x}(x - x_0)$ is $(m \times n) \cdot (n \times 1) = m \times 1$.

Note: The linear approximation is also called the **first-order Taylor expansion of f about the point x_0** . **Probably, at least 50% of engineering uses this local linear approximation of a function.**

Proof: We pose

$$f(x) \approx y_0 + A \cdot (x - x_0) \quad (5.30)$$

for $\|x - x_0\| < \delta$ and some $\delta > 0$. What values should we assign to y_0 and A ? Just as in the case of a scalar-valued function of a single variable, if we plug $x = x_0$ into (5.30), we obtain

$$f(x_0) = y_0.$$

Hence,

$$f(x) \approx f(x_0) + A \cdot (x - x_0). \quad (5.31)$$

Next, we plug in $x = x_0 + h \cdot e_i$ into (5.31), yielding

$$\begin{aligned} f(x_0 + h \cdot e_i) &\approx f(x_0) + \underbrace{A \cdot h \cdot e_i}_{h \cdot e_i = (x_0 + h \cdot e_i) - x_0} \\ &\approx f(x_0) + h \cdot A \cdot e_i \\ &\approx f(x_0) + h \cdot A_i, \end{aligned}$$

where we note that $A \cdot e_i = A_i$, the i -th column of A . A little algebra gives us

$$(f(x_0 + h \cdot e_i) - f(x_0) \approx h \cdot A_i) \iff \left(A_i \approx \frac{f(x_0 + h \cdot e_i) - f(x_0)}{h} \right). \quad (5.32)$$

The approximation should be better and better with smaller and smaller h . In the limit, we have

$$A_i = \lim_{h \rightarrow 0} \frac{f(x_0 + h \cdot e_i) - f(x_0)}{h} =: \frac{\partial f(x_0)}{\partial x_i}, \quad (5.33)$$

as long as the limit exists. In other symbols, if f is differentiable at x_0 , once we assemble the columns of

$$A = \begin{bmatrix} \frac{\partial f(x_0)}{\partial x_1} & \frac{\partial f(x_0)}{\partial x_2} & \dots & \frac{\partial f(x_0)}{\partial x_n} \end{bmatrix},$$

we have

$$f(x) \approx f(x_0) + \frac{\partial f(x_0)}{\partial x} \cdot (x - x_0) \quad (5.34)$$

is a **local linear approximation** to the function near the point x_0 . ■

Secrets of the Arcane 5.40: Use of First-Order Taylor Expansions in Engineering

The use of **first-order Taylor expansions**, often referred to as **linear approximations**, is pervasive in various fields of engineering. While it's difficult to quantify "how much" of engineering is based on this concept, it's safe to say that it's a fundamental tool used frequently for analysis, design, and optimization. Here are some areas where it's commonly used:

Feedback Control Systems: In control theory, nonlinear systems are often linearized around an operating point to make them easier to analyze and control. The linearized models are then used to design controllers. **We will make**

use of linearization when we design feedback controllers for the Ball Bot. A cool thing about the use of linearizations in feedback control is that **the controller is trying to keep an error term small**, and this error is often $(x - x_0)$. In other words, the controller is helping to validate the assumptions required for using a linear approximation in the first place.

Structural Engineering: In the analysis of structures like beams, bridges, and buildings, linear approximations are often used to simplify complex equations describing deformation and stress.

Electrical Engineering: In circuit analysis, devices like transistors are often modeled using linear approximations around an operating point, especially for small-signal analysis.

Fluid Mechanics: In fluid dynamics, linear approximations are used to simplify the Navier-Stokes equations under certain conditions, leading to more tractable models like potential flow.

Thermodynamics: In thermodynamics, properties like enthalpy and entropy are often linearized to analyze processes like isothermal and adiabatic expansions.

Signal Processing: In signal processing, systems are often linearized to make them easier to analyze using Fourier and Laplace transforms.

Optimization: In optimization problems, the objective function and constraints are often linearized to find a local minimum or maximum.

Economics and Operations Research: These fields also make extensive use of linear approximations for problems like supply chain optimization, market equilibrium analysis, etc.

Limitations: It's important to note that linear approximations are only valid "near the point of expansion"^a and may not capture the behavior of a nonlinear function (or system model) accurately over a large domain. Engineers must exercise caution when using these approximations, especially for systems that exhibit "strong"^b nonlinearities.

^aIn practice, it is an art to understand what "near" means in any given application. Once again, this is great news as it keeps domain experts, like you, employed!

^bNonlinear terms, such as x^2 or e^x that "blow up quickly". Functions that are bounded, such as $\sin(x)$, are considered "mild nonlinearities".

Example 5.41. For the following functions $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$, compute their Jacobians at a general point $x \in \mathbb{R}^n$ and compute their linearization about the given point.

$$(a) f(x_1, x_2) := \begin{bmatrix} x_1 \cdot (x_2)^2 \\ e^{3x_1} + \cos(4x_2) \end{bmatrix} \text{ and } x_0 = \begin{bmatrix} \sqrt{2} \\ \pi \end{bmatrix}.$$

$$(b) f(x_1, x_2) := e^{3x_1} + \cos(4x_2) \text{ and } x_0 = \begin{bmatrix} 0.77 \\ 11.6 \end{bmatrix}.$$

$$(c) f(x_1, x_2) := \begin{bmatrix} x_1 \cdot (x_2)^2 \\ x_2 \cdot \sin^3(x_1) + 11 \\ 4(x_1)^2 + (x_2 + 3)^2 \end{bmatrix} \text{ and } x_0 = \begin{bmatrix} \sqrt{7} \\ 1.0 \end{bmatrix}.$$

$$(d) f(x_1, x_2, x_3) := \begin{bmatrix} x_1 \cdot x_2 \cdot (x_3)^3 \\ x_2 \cdot e^{x_1 \cdot x_3} \\ x_1 \end{bmatrix} \text{ and } x_0 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}.$$

Solutions: These are the same functions from Example 5.41. We could assemble the various partial derivatives we have already computed by hand and verified with the help of the `Symbolics` package. Instead, we'll do everything in software and use the functions `Symbolics.jacobian` and `ForwardDiff.jacobian`. Would you accept anything else?

```
1 using Symbolics
2
3 # Define variables
4 @variables x1 x2 x3
5
6 fa = [x1 * (x2)^2; exp(3*x1) + cos(4*x2)]
7 fb = [exp(3*x1) + cos(4*x2)]
8 fc = [x1 * (x2)^2; x2*(sin(x1))^3 + 11; 4*(x1)^2 + (x2 + 3)^2]
9 fd = [x1*x2*(x3)^3; x2*exp(x1*x3); x1]
10
11 # Symbolic Jacobians
12
13 #a)
14 Ja = Symbolics.jacobian(fa, [x1, x2])
15 println("(a) ")
16 display(Ja)
17 println(Ja)
18 println(" ")
19
20 #b)
21
22 Jb = Symbolics.jacobian(fb, [x1, x2]) # The Jacobian function want a vector
23 # hence, the brackets
24 println("(b) ")
25 display(Jb)
26 println(Jb)
27 println(" ")
28
29 #c)
30 Jc = Symbolics.jacobian(fc, [x1, x2])
31 println("(c) ")
32 display(Jc)
33 println(Jc)
34 println(" ")
35
36 #d)
37 Jd = Symbolics.jacobian(fd, [x1, x2, x3])
38 println("(d) ")
39 display(Jd)
40 println(Jd)
41 println(" ")
```

Output

(a)

$$\begin{bmatrix} x2^{2.0000} & 2.0000x1x2 \\ 3.0000e^{3.0000x1} & -4.0000 \sin(4.0000x2) \end{bmatrix}$$

Num[x2^2 2x1*x2; 3exp(3x1) -4sin(4x2)]

(b)

$$\begin{bmatrix} 3.0000e^{3.0000x1} & -4.0000 \sin(4.0000x2) \end{bmatrix}$$

Num[3exp(3x1) -4sin(4x2)]

(c)

$$\begin{bmatrix} x2^{2.0000} & 2.0000x1x2 \\ 3.0000 \sin^{2.0000}(x1)x2 \cos(x1) & \sin^{3.0000}(x1) \\ 8.0000x1 & 2.0000(3.0000 + x2) \end{bmatrix}$$

Num[x2^2 2x1*x2; 3x2*(sin(x1)^2)*cos(x1) sin(x1)^3; 8x1 2(3 + x2)]

(d)

$$\begin{bmatrix} x3^{3.0000}x2 & x3^{3.0000}x1 & 3.0000x3^{2.0000}x1x2 \\ x2x3e^{x1x3} & e^{x1x3} & x1x2e^{x1x3} \\ 1.0000 & 0.0000 & 0.0000 \end{bmatrix}$$

Num[x2*(x3^3) x1*(x3^3) 3x1*x2*(x3^2); x2*x3*exp(x1*x3) exp(x1*x3) x1*x2*exp(x1*x3); 1 0 0]

Next, we evaluate the Jacobians at a point.

```
1 using ForwardDiff
2 # Example of converting from Symbolic to Numerical Functions
3 using Symbolics, ForwardDiff
4
5 #a)
6 x0 = [sqrt(2); pi]
7 # Build a numerical function from the symbolic expression
8 fa_num = build_function(fa, [x1, x2])
9 fa_num = eval(fa_num[1])
10 # Compute the Jacobian of the new function at the test point
11 Jaca = ForwardDiff.jacobian(fa_num, x0)
12 println("(a) ")
13 display(Jaca)
14 println(" ")
15
16 #b)
17 x0 = [0.77; 11.6]
18 fb_num = build_function(fb, [x1, x2])
19 fb_num = eval(fb_num[1])
20 # Compute the Jacobian of the new function at the test point
21 Jacb = ForwardDiff.jacobian(fb_num, x0)
22 println("(b) ")
23 display(Jacb)
24 println(" ")
25
26 #c)
```

```

27 x0 = [sqrt(7); 1.0]
28 fc_num = build_function(fc, [x1, x2])
29 fc_num = eval(fc_num[1])
30 # Compute the Jacobian of the new function at the test point
31 Jacc = ForwardDiff.jacobian(fc_num, x0)
32 println("(c) ")
33 display(Jacc)
34 println(" ")
35
36 #d)
37 x0 = [0.0; 0; 0]
38 # Build a numerical function from the symbolic expression
39 fd_num = build_function(fd, [x1, x2, x3])
40 fd_num = eval(fd_num[1])
41 # Compute the Jacobian of the new function at the test point
42 Jacd = ForwardDiff.jacobian(fd_num, x0)
43 println("(d) ")
44 display(Jacd)
45 println(" ")

```

Output

(a)

$$2 \times 2 \text{ Matrix}\{\text{Float64}\}:$$

$$\begin{matrix} 9.8696 & 8.88577 \\ 208.774 & 1.95943e-15 \end{matrix}$$

(b)

$$1 \times 2 \text{ Matrix}\{\text{Float64}\}:$$

$$\begin{matrix} 30.2233 & -2.64922 \end{matrix}$$

(c)

$$3 \times 2 \text{ Matrix}\{\text{Float64}\}:$$

$$\begin{matrix} 1.0 & 5.2915 \\ -0.597294 & 0.107695 \\ 21.166 & 8.0 \end{matrix}$$

(d)

$$3 \times 3 \text{ Matrix}\{\text{Float64}\}:$$

$$\begin{matrix} 0.0 & 0.0 & 0.0 \\ 0.0 & 1.0 & 0.0 \\ 1.0 & 0.0 & 0.0 \end{matrix}$$

In each case, the first-order linear approximation is $f(x_0) + \frac{\partial f(x_0)}{\partial x} \cdot (x - x_0)$. All of the Jacobians have been evaluated.

Converting Symbolic Functions to Numerical Functions in Julia

The various Julia packages for computing derivatives do not yet play so nicely with one another. They will someday. The following method was found by applying numerous prompts to ChatGPT4 and Bard (25 Sept 2023). After many false starts, it's actually pretty clean and is what we used above. In the following, the code is explained.

```

1 # Example of converting from Symbolic to Numerical Functions
2 using Symbolics, ForwardDiff
3
4 @variables x1 x2 x3
5 fd = [x1*x2*(x3)^3; x2*exp(x1*x3); x1]

```

```

6
7 # Build a numerical function from the symbolic expression
8 fd_num = build_function(fd, [x1, x2, x3])
9 fd_num = eval(fd_num[1])
10
11 # Test the new function
12 x0 = [1.0, 2.0, 3.0]
13 result = fd_num(x0)
14 println("Result of fd_num(x0) : ", result)
15
16 # Compute the Jacobian of the new function at the test point
17 Jac = ForwardDiff.jacobian(fd_num, x0)
18 println("Jacobian: ", Jac)

```

The above code snippet demonstrates how to convert a symbolic function to a numerical function in Julia using the Symbolics and ForwardDiff packages. The symbolic function is defined with variables x_1, x_2, x_3 and consists of three components. The numerical function is then tested at a point $x_0 = [1.0, 2.0, 3.0]$, and its Jacobian is computed at this point.

Importing Packages

```
1 using Symbolics, ForwardDiff
```

The Symbolics package is used for symbolic computations, and ForwardDiff is used for automatic differentiation.

Defining Symbolic Variables and Function

```

1 @variables x1 x2 x3
2 fd = [x1*x2*(x3)^3; x2*exp(x1*x3); x1]

```

Symbolic variables x_1, x_2, x_3 are defined, and a symbolic function fd is created as a vector of three components.

Building the Numerical Function

```

1 fd_num = build_function(fd, [x1, x2, x3])
2 fd_num = eval(fd_num[1])

```

The `build_function` method from the Symbolics package is used to convert the symbolic function into a numerical function. The `eval` function is then used to make it callable.

Testing the Numerical Function

```

1 x0 = [1.0, 2.0, 3.0]
2 result = fd_num(x0)
3 println("Result of fd_num(x0) : ", result)

```

The numerical function is tested at a point $x_0 = [1.0, 2.0, 3.0]$.

Computing the Jacobian

```

1 Jac = ForwardDiff.jacobian(fd_num, x0)
2 println("Jacobian: ", Jac)

```

The Jacobian of the numerical function at the point x_0 is computed using ForwardDiff's `jacobian` function. **Output**

```
Result of fd_num(x0) : [54.0, 40.171073846375336, 1.0]
Jacobian: [54.0 27.0 54.0; 120.51322153912601 20.085536923187668 40.171073846375336;
1.0 0.0 0.0]
```

You may prefer the following alternative method. It is simpler but has some manual intervention. The purest among you will not like that. Those who just want to get their HW completed, however, may love it!

```

1 # Exercise in textbook
2 using Symbolics
3 # define variables
4 @variables x[1:3]
5
6 fa = [x[1] * (x[2])^2; exp(3*x[1]) + cos(4*x[2])]
7 fb = [exp(3*x[1]) + cos(4*x[2])]
8 fc = [x[1] * (x[2])^2; x[2]*sin(x[1]) )^3 + 11; 4*(x[1])^2 + ( x[2] +3 )^2]
9 fd = [x[1]*x[2]*(x[3])^3; x[2]*exp(x[1]*x[3]); x[1]]
10
11 # Compute Symbolic Jacobians
12 #a)
13 Ja = Symbolics.jacobian(fa, x[1:2])
14 println("(a) ")
15 display(Ja)
16 println(Ja)
17 println(" ")
18
19 #b)
20
21 Jb = Symbolics.jacobian(fb, x[1:2]) # The Jacobian function want a vector
22 # hence, the brackets
23 println("(b) ")
24 display(Jb)
25 println(Jb)
26 println(" ")
27
28 #c)
29 Jc = Symbolics.jacobian(fc, x[1:2])
30 println("(c) ")
31 display(Jc)
32 println(Jc)
33 println(" ")
34
35 #d)
36 Jd = Symbolics.jacobian(fd, x)
37 println("(d) ")
38 display(Jd)
39 println(Jd)
40 println(" ")

```

Output

(a)

$$\begin{bmatrix} x_2^{2.0000} & 2.0000x_1x_2 \\ 3.0000e^{3.0000x_1} & -4.0000 \sin(4.0000x_2) \end{bmatrix}$$

Num[x[2]^2 2x[1]*x[2]; 3exp(3x[1]) -4sin(4x[2])]

(b)

$$\begin{bmatrix} 3.0000e^{3.0000x_1} & -4.0000 \sin(4.0000x_2) \end{bmatrix}$$

Num[3exp(3x[1]) -4sin(4x[2])]

(c)

$$\begin{bmatrix} x_2^{2.0000} & 2.0000x_1x_2 \\ 3.0000 \sin^{2.0000}(x_1) \cos(x_1)x_2 & \sin^{3.0000}(x_1) \\ 8.0000x_1 & 2.0000(3.0000 + x_2) \end{bmatrix}$$

Num[x[2]^2 2x[1]*x[2]; 3(sin(x[1])^2)*cos(x[1])*x[2] sin(x[1])^3; 8x[1] 2(3 + x[2])]

(d)

$$\begin{bmatrix} x_3^{3.0000}x_2 & x_3^{3.0000}x_1 & 3.0000x_3^{2.0000}x_1x_2 \\ e^{x_1x_3}x_2x_3 & e^{x_1x_3} & e^{x_1x_3}x_1x_2 \\ 1.0000 & 0.0000 & 0.0000 \end{bmatrix}$$

Num[(x[3]^3)*x[2] (x[3]^3)*x[1] 3(x[3]^2)*x[1]*x[2]; exp(x[1]*x[3])*x[2]*x[3] exp(x[1]*x[3]) exp(x[1]*x[3])*x[1]*x[2]; 1 0 0]

using ForwardDiff

```

1
2
3 # Copy and paste to create regular Julia functions. The names are slightly altered so as
4 # not to have a conflict with the Symbolic functions fa, fb, fc, fd
5
6 fA(x) = [x[1] * (x[2])^2; exp(3*x[1]) + cos(4*x[2])]
7 fB(x) = [exp(3*x[1]) + cos(4*x[2])]
8 fC(x) = [x[1] * (x[2])^2; x[2]*(sin(x[1]))^3 + 11; 4*(x[1])^2 + (x[2] + 3)^2]
9 fD(x) = [x[1]*x[2]*(x[3])^3; x[2]*exp(x[1]*x[3]); x[1]];
10
11 #a)
12 x0 = [sqrt(2); pi]
13 Ja = ForwardDiff.jacobian(fA, x0)
14 println ("(a) ")
15 display(Ja)
16 println (" ")
17
18 #b)
19 x0 = [0.77; 11.6]
20 Jb = ForwardDiff.jacobian(fB, x0)
21 println ("(b) ")
22 display(Jb)
23 println (" ")
24
25 #c)
26 x0 = [sqrt(7); 1.0]
27 Jc = ForwardDiff.jacobian(fC, x0)
28 println ("(c) ")
29 display(Jc)
30 println (" ")
31
32 #d)

```

```

33 x0 = [0.0; 0; 0]
34 Jd = ForwardDiff.jacobian(fD, x0)
35 println("(d) ")
36 display(Jd)
37 println(" ")

```

Output

```

(a)
2×2 Matrix{Float64}:
 9.8696  8.88577
 208.774 1.95943e-15

(b)
1×2 Matrix{Float64}:
 30.2233 -2.64922

(c)
3×2 Matrix{Float64}:
 1.0      5.2915
 -0.597294 0.107695
 21.166    8.0

(d)
3×3 Matrix{Float64}:
 0.0  0.0  0.0
 0.0  1.0  0.0
 1.0  0.0  0.0

```

One more example, just for fun.

```

1 @variables t x y z a b
2
3 # Define a gnarly function with variables t, x, y, z, a, b and 3 components
4 f_gnarly = [
5     sin(a*x^3) + cos(b*exp(y)) + t*z^2 + exp(a - b),
6     a*x^2*y - b*t*z + log(x + y + z + 1),
7     a*b*t*x*y*z + sin(t) + cos(a) + (atan(exp(b))+11)^z]
8
9 # Compute the Jacobian symbolically
10 Jacobian_gnarly = Symbolics.jacobian(f_gnarly, [t, x, y, z, a, b])
11
12 @show f_gnarly
13 Jacobian_gnarly

```

Output To save space, we replace $f_{\text{gnarly}}(t, x, y, z, a, b)$ with $f(\dots)$, and note that

$$\text{Jac } f(\dots) = \begin{bmatrix} \frac{\partial f(\dots)}{\partial t} & \frac{\partial f(\dots)}{\partial x} & \frac{\partial f(\dots)}{\partial y} & \frac{\partial f(\dots)}{\partial z} & \frac{\partial f(\dots)}{\partial a} & \frac{\partial f(\dots)}{\partial b} \end{bmatrix},$$

where

$$\frac{\partial f(\dots)}{\partial t} = \begin{bmatrix} z^2 \\ -bz \\ abxyz + \cos(t) \end{bmatrix}, \quad \frac{\partial f(\dots)}{\partial x} = \begin{bmatrix} 3x^2a \cos(x^3a) \\ \frac{1}{1+x+y+z} + 2axy \\ abtyz \end{bmatrix}, \quad \frac{\partial f(\dots)}{\partial y} = \begin{bmatrix} -be^y \sin(be^y) \\ x^2a + \frac{1}{1+x+y+z} \\ abtxz \end{bmatrix}$$

$$\frac{\partial f(\dots)}{\partial z} = \begin{bmatrix} 2tz \\ \frac{1}{1+x+y+z} - bt \\ (11 + \arctan(e^b))^z \log(11 + \arctan(e^b)) + abtxy \end{bmatrix}$$

$$\frac{\partial f(\dots)}{\partial a} = \begin{bmatrix} x^3 \cos(x^3a) + e^{a-b} \\ x^2y \\ -\sin(a) + btxyz \end{bmatrix}, \quad \frac{\partial f(\dots)}{\partial b} = \begin{bmatrix} -e^{a-b} - e^y \sin(be^y) \\ -tz \\ \frac{(11+\arctan(e^b))^{-1+z}ze^b}{1+(e^b)^2} + atxyz \end{bmatrix}.$$

That would be hard to do by hand. ■

5.7.2 The Gradient

The **gradient** is the special name given to the **transpose of the Jacobian** of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, (i.e., for each $x \in \mathbb{R}^n$, $f(x) \in \mathbb{R}$ is a scalar). In Chapter 6.2.2, where we deal with optimization, we will show that the gradient of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ **points in the direction of maximum increase of the function**, and hence the negative of the gradient points in the **direction of maximum decrease**. All of the marvelous AI tools at our disposal today rely on the gradient for *training* the neural networks on which the tools are based. **That makes the gradient one of the most important objects in all of Calculus**...and yet, it is just a collection of partial derivatives stacked into a column vector. Who could have imagined that?

The Gradient

Definition 5.42. The **gradient** of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is the partial derivatives of f with respect to x_i arranged to form a **column vector**^a, instead of a row vector,

$$\nabla f(x_0) := \begin{bmatrix} \frac{\partial f(x_0)}{\partial x_1} \\ \frac{\partial f(x_0)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x_0)}{\partial x_n} \end{bmatrix}. \quad (5.35)$$

The cool symbol ∇ is usually pronounced as “**grad**” and one says “**grad f**” when speaking of ∇f .

Remark 5.43. The **Jacobian** of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is the row vector, $\frac{\partial f(x)}{\partial x} := \left[\frac{\partial f(x)}{\partial x_1} \quad \frac{\partial f(x)}{\partial x_2} \quad \dots \quad \frac{\partial f(x)}{\partial x_n} \right]$. It is a $1 \times n$ row vector because each entry, $\frac{\partial f(x)}{\partial x_i}$, is 1×1 due to f having codomain \mathbb{R} , and it has n columns because $x \in \mathbb{R}^n$.

The gradient is the transpose of the Jacobian. Hence the linearization of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ can be written as

$$\begin{aligned} f(x) &\approx f(x_0) + (\nabla f(x_0))^\top \cdot (x - x_0) \\ &\approx f(x_0) + \nabla f(x_0) \bullet (x - x_0), \end{aligned} \quad (5.36)$$

where $\nabla f(x_0) \bullet (x - x_0)$ is the dot product of the $n \times 1$ gradient vector with the $n \times 1$ vector $(x - x_0)$.

^aIn ROB 101 *Computational Linear Algebra*, we did not assume Calculus. When we introduced the gradient, for simplicity, we left it as a row vector.

Example 5.44. Compute the gradient of $f : \mathbb{R}^2 \rightarrow \mathbb{R}$, for $f(x_1, x_2) = x_1 \cos(x_2)$ and $x_0 = [2 \ \pi/4]^\top$.

Solution: We'll compute $\nabla f(x)$ first by hand and then compute both $\nabla f(x)$ and $\nabla f(x_0)$ in software.

¹ using Symbolics

²

³ @variables x[1:2]

⁴ f = x[1] * cos(x[2])

```

5 grad_f = Symbolics.gradient(f, [x[1], x[2]])
6
7 @show f
8 @show grad_f
9
10 x0 = [2; pi/4]
11 f_num = build_function(f, x)
12 f_num = eval(f_num)
13 # Compute the Jacobian of the new function at the test point
14 Grad_f_num= ForwardDiff.gradient(f_num, x0)
15 println(" ")
16 display(Grad_f_num)
17 println(" ")

```

Output

$$\frac{\partial}{\partial x} (x_1 \cos(x_2)) = \begin{bmatrix} \frac{\partial}{\partial x_1} (x_1 \cos(x_2)) \\ \frac{\partial}{\partial x_2} (x_1 \cos(x_2)) \end{bmatrix} = \begin{bmatrix} \cos(x_2) \\ -x_1 \sin(x_2) \end{bmatrix}.$$

$$\left. \frac{\partial}{\partial x} (x_1 \cos(x_2)) \right|_{x_0} = \begin{bmatrix} 0.70710678 \\ -1.41421356 \end{bmatrix}$$

One could easily do a “gnarly” example with the gradient as well. We leave that to you! ■

5.7.3 A Handy Product Rule for the Gradient and the Jacobian

On one hand, because partial derivatives are ordinary derivatives taken one variable at a time while holding the others constant, all of the rules for differentiation presented in Chapter 5.4 apply to individual partial derivatives, $\frac{\partial}{\partial x_i}$. On the other hand, applying these rules in a matrix-vector form typically requires tensor notation, which we do not want to touch in this textbook. You are likely very happy about that! There is, however, one useful case where things work out nicely, and we treat that here.

Proposition 5.45: Gradient and Jacobian Meet the Dot Product

Suppose $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ are vector-valued functions and define $h : \mathbb{R}^n \rightarrow \mathbb{R}$ by

$$h(x) := (f(x))^\top \cdot g(x) = f(x) \bullet g(x), \quad (5.37)$$

the dot product of f and g . If f and g are both differentiable, then so is h and

$$J_h(x) := \begin{bmatrix} \frac{\partial h(x)}{\partial x_1} & \frac{\partial h(x)}{\partial x_2} & \dots & \frac{\partial h(x)}{\partial x_n} \end{bmatrix} = [g(x)]^\top \cdot \frac{\partial f(x)}{\partial x} + [f(x)]^\top \cdot \frac{\partial g(x)}{\partial x} \quad (5.38)$$

and because the gradient is the transpose of the Jacobian, and the transpose of a product equals the product of the transposes in reverse order, we have

$$\nabla h(x) := \begin{bmatrix} \frac{\partial h(x)}{\partial x_1} \\ \frac{\partial h(x)}{\partial x_2} \\ \vdots \\ \frac{\partial h(x)}{\partial x_n} \end{bmatrix} = \left[\frac{\partial g(x)}{\partial x} \right]^\top \cdot f(x) + \left[\frac{\partial f(x)}{\partial x} \right]^\top \cdot g(x). \quad (5.39)$$

Proof: The function $h(x)$ is equal to

$$h(x) = \sum_{j=1}^n f_j(x)g_j(x).$$

Taking the partial derivative of h with respect to a specific component x_i , we can apply the product rule for single-variable derivatives to obtain,

$$\begin{aligned}\frac{\partial h(x)}{\partial x_i} &= \sum_{j=1}^m \left(\frac{\partial f_j(x)}{\partial x_i} g_j(x) + f_j(x) \frac{\partial g_j(x)}{\partial x_i} \right) \quad (\text{next, rewrite each sum as a row vector times a column vector}) \\ &= \begin{bmatrix} g_1(x) & g_2(x) & \cdots & g_m(x) \end{bmatrix} \cdot \begin{bmatrix} \frac{\partial f_1(x)}{\partial x_i} \\ \frac{\partial f_2(x)}{\partial x_i} \\ \vdots \\ \frac{\partial f_m(x)}{\partial x_i} \end{bmatrix} + \begin{bmatrix} f_1(x) & f_2(x) & \cdots & f_m(x) \end{bmatrix} \cdot \begin{bmatrix} \frac{\partial g_1(x)}{\partial x_i} \\ \frac{\partial g_2(x)}{\partial x_i} \\ \vdots \\ \frac{\partial g_m(x)}{\partial x_i} \end{bmatrix} \\ &= [g(x)]^\top \cdot \begin{bmatrix} \frac{\partial f(x)}{\partial x} \end{bmatrix}_i + [f(x)]^\top \cdot \begin{bmatrix} \frac{\partial g(x)}{\partial x} \end{bmatrix}_i,\end{aligned}$$

where $\begin{bmatrix} \frac{\partial f(x)}{\partial x} \end{bmatrix}_i$ and $\begin{bmatrix} \frac{\partial g(x)}{\partial x} \end{bmatrix}_i$ denote the i -th columns of the Jacobian matrices. Hence, putting it all together,

$$J_h(x) := \begin{bmatrix} \frac{\partial h(x)}{\partial x_1} & \frac{\partial h(x)}{\partial x_2} & \cdots & \frac{\partial h(x)}{\partial x_n} \end{bmatrix} = [g(x)]^\top \cdot \frac{\partial f(x)}{\partial x} + [f(x)]^\top \cdot \frac{\partial g(x)}{\partial x}.$$

■

5.7.4 (Optional Read:) The Hessian

Assume that the gradient of the scalar-valued function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ exists. Then $\nabla f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is an $(n \times 1)$ vector-valued function. If its partial derivatives with respect to $x = (x_1, x_2, \dots, x_n)$ exist, then we can compute its Jacobian, which will be an $n \times n$ matrix. The name **Hessian** for this object is in honor of its inventor, Ludwig Otto Hesse.

Hessian of a Function

Definition 5.46. The **Hessian** of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is

$$H_f(x) := \frac{\partial \nabla f(x)}{\partial x} := \frac{\partial}{\partial x} (\nabla f(x)) = \begin{bmatrix} \frac{\partial \nabla f(x)}{\partial x_1} & \frac{\partial \nabla f(x)}{\partial x_2} & \cdots & \frac{\partial \nabla f(x)}{\partial x_n} \end{bmatrix}, \quad (5.40)$$

assuming the indicated partial derivatives exist. The Hessian of a function is formed by packaging the column vectors $\frac{\partial \nabla f(x)}{\partial x_j}$ into a matrix. We need to keep in mind that, for each value of $x \in \mathbb{R}^n$, the compact and awesome-looking object

$$H_f(x) := \frac{\partial \nabla f(x)}{\partial x}$$

is really an $n \times n$ matrix: because there are n columns of the form $\frac{\partial \nabla f(x)}{\partial x_j}$, and each column is an n -vector. There are software tools that compute the Hessian in one fell swoop. You do not need to first compute the gradient and then the Jacobian, as in the definition.

Secrets of the Arcane 5.47: Point Where Many Calc III Students Throw Up Their Hands and Run!

- (a) Yes, we forgot to mention that Jacobians, gradients, and Hessians are typically taught in Calc III, Multivariable Calculus, aka Vector Calculus^a. Because we introduced all of these topics in ROB 101 *Computational Linear Algebra* in the context of root finding for vector-valued functions (Newton-Raphson Algorithm) and minimization of a scalar-valued function depending on m -variables (Gradient Descent), we knew you could handle it.
- (b) The Julia commands for the Hessian are `Symbolics.hessian(f, x)` and `ForwardDiff.hessian(f, x0)`, one line of code.

- (c) Below, we write out the Hessian explicitly. It's a mess! Just the thought of having to compute it by hand is terrifying. In math courses, the teachers often take pity on the students and only assign trivial examples for computing the Hessian or have students compute just a few of the terms in the matrix. Consequently, the whole exercise seems pointless to the students. Whence, this textbook.

- (d) **A math fact on mixed partial derivatives.** If

$$\frac{\partial^2 f(x)}{\partial x_i \partial x_j} \text{ and } \frac{\partial^2 f(x)}{\partial x_j \partial x_i}$$

are both continuous functions, then

$$\frac{\partial^2 f(x)}{\partial x_i \partial x_j} = \frac{\partial^2 f(x)}{\partial x_j \partial x_i}.$$

Hence, for most engineering applications, the Hessian matrix is symmetric. In math courses, the Hessian is defined to be the transpose of (5.40). For examples of interest to us, the difference is unimportant because our Hessians are square symmetric matrices.

^aMath 215 Multiv & Vector Calc, at Michigan. In addition to differentiation, it covers multivariable integration as well.

Just for the record, we write out $H_f(x) := \frac{\partial \nabla f(x)}{\partial x}$ as an $n \times n$ matrix. Suppose $f(x) = f(x_1, x_2, \dots, x_n)$. Then,

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix}.$$

Computing the Jacobian of the gradient, and writing out all of the entries in their $m \times m$ matrix glory, gives

$$H_f(x) = \frac{\partial \nabla f(x)}{\partial x} := \left[\begin{array}{cccc} \frac{\partial \nabla f(x)}{\partial x_1} & \frac{\partial \nabla f(x)}{\partial x_2} & \cdots & \frac{\partial \nabla f(x)}{\partial x_n} \end{array} \right] = \left[\begin{array}{cccc} \frac{\partial^2 f(x)}{\partial x_1 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \cdots & \frac{\partial^2 f(x)}{\partial x_n \partial x_1} \\ \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \frac{\partial^2 f(x)}{\partial x_2 \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_n \partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_1 \partial x_n} & \frac{\partial^2 f(x)}{\partial x_2 \partial x_n} & \cdots & \frac{\partial^2 f(x)}{\partial x_n \partial x_n} \end{array} \right]. \quad (5.41)$$

In other symbols, the ij component of $H_f(x) := \frac{\partial \nabla f(x)}{\partial x}$ is

$$\left[\frac{\partial \nabla f(x)}{\partial x} \right]_{ij} = \frac{\partial^2 f(x)}{\partial x_j \partial x_i} \left(= \frac{\partial^2 f(x)}{\partial x_i \partial x_j} \text{ when the partial derivatives are continuous} \right). \quad (5.42)$$

5.7.5 (Optional Read:) More Matrix-Vector Calculus

The first thing we want to do is highlight this very useful site: [Matrix Calculus](#), an online tool that computes vector and matrix derivatives (matrix calculus), so that you do not have to memorize all the formulas. Now, this site brings up an issue we mentioned previously: there are different conventions for whether $\frac{\partial}{\partial x}$ applied to a scalar-valued function $f : \mathbb{R}^m \rightarrow \mathbb{R}$ is interpreted as a row vector, which would be a Jacobian, or a column vector, which would be a gradient. Sheesh, you are thinking, after 200 years of Calculus, can't they at least have this settled? Probably not until one of you mounts a worldwide campaign to settle this "atrocious state of affairs". We discuss the "why" of this below, in case you are interested. Let's note, though, that in spoken and written language, you deal with contradictory meanings and conventions all the time and accept that as part of being fluent. It's no different in math. **Context usually settles any possible confusion.**

With these comments out of the way, we'll list a few key formulas and derive them so that you have confidence when using matrix-vector calculus.

Proposition 5.48: Key Matrix-vector Calculus Formulas

Most useful formulas for engineers:

- (a) For C an $m \times n$ matrix and $x \in \mathbb{R}^n$, $\frac{\partial}{\partial x}(Cx) = C$.
- (b) For $x \in \mathbb{R}^n$, $\frac{\partial}{\partial x}(x^\top x) = 2x^\top$.
- (c) For $x \in \mathbb{R}^n$, $\nabla(x^\top x) = 2x$.
- (d) For Q an $n \times n$ symmetric matrix and $x \in \mathbb{R}^n$, $\frac{\partial}{\partial x}(x^\top Qx) = 2x^\top \cdot Q$.
- (e) For Q an $n \times n$ symmetric matrix and $x \in \mathbb{R}^n$, $\nabla(x^\top Qx) = 2Qx$.
- (f) For Q an $n \times n$ symmetric matrix and $x \in \mathbb{R}^m$, $\frac{\partial \nabla(x^\top Qx)}{\partial x} = 2Q$ (yes, this is the Hessian of $x^\top Qx$). You will also see it written as $\frac{\partial}{\partial x}\left(\frac{\partial(x^\top Qx)}{\partial x}\right) = 2Q$, especially when using the `Symbolics` package or `SymPy`.

Note: You will often see a $\frac{1}{2}$ inserted into $x^\top x$ and $x^\top QX$ so that $\nabla\left(\frac{1}{2}x^\top x\right) = x$ and $\nabla\left(\frac{1}{2}x^\top Qx\right) = Qx$. There has to be a two or a one-half somewhere! It's your call.

Row vs Column? How to Know: The Jacobian matrix provides a linear approximation of a multivariable function. In this context, if $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $m \geq 2$, then the Jacobian of f has to be an $m \times n$ matrix. However, when $m = 1$, all bets are off! **Because the gradient is just the transpose of the Jacobian, in many contexts, it seems redundant to use both operations.** The choice to represent either the gradient or Jacobian of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ as a row vector or a column vector, therefore, depends on conventions, contexts, and the operations you're planning to perform with it. Here's how it usually plays out:

- (i) When $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$, $m > 1$, then the Jacobian $\frac{\partial f(x)}{\partial x}$ is an $m \times n$ matrix.
- (ii) When $f : \mathbb{R}^n \rightarrow \mathbb{R}$, (i.e., $m = 1$), then it's 25-75% whether the Jacobian $\frac{\partial f(x)}{\partial x}$ is a $1 \times n$ row vector (25%) or $n \times 1$ column vector (75%), even in this textbook. Whenever possible, we will write the Jacobian as a row vector and use the gradient when we want a column vector. However, we cannot fly in the face of tradition and rewrite famous equations in a form that would be unrecognizable to most users. When treating Lagrange Multipliers in Chapter 6.3, it is very natural to use the gradient for a column of partial derivatives, and we will. When treating Lagrangian Dynamics in Chapter 6.4, indicating a column of partial derivatives with a gradient instead of a Jacobian would make the equations unrecognizable, to the point that if you showed them to someone who studied the topic in a different course, they'd think you were from Mars. Alas, sometimes tradition prevails.
- (iii) Hence, for us, the Jacobian of $x^\top x$ is a row vector. The extremely useful site, **Matrix Calculus**, the online tool that computes vector and matrix derivatives so that you do not have to, makes the Jacobian of $x^\top x$ a column vector. Outside of this course, when writing out the Jacobian or gradient, always be clear about which convention you're using, especially if your colleagues might be used to a different convention.
- (iv) When the gradient of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is represented as a column vector, it is the transpose of the Jacobian matrix of f ,

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix} = [J_f(x)]^\top.$$

This is the more common representation, especially in optimization contexts.

- (v) The gradient of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ can also be represented as a row vector, in which case it's exactly the Jacobian matrix for f ,

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} & \frac{\partial f(x)}{\partial x_2} & \dots & \frac{\partial f(x)}{\partial x_n} \end{bmatrix} = J_f(x).$$

We used this convention in ROB 101 *Computational Linear Algebra* to “keep things simple” because Calculus was not a prerequisite.

In conclusion, while mathematical conventions can guide the choice of representation, the ultimate decision often depends on the specific context and operations in which the gradient or Jacobian is used. Always prioritize clarity and correctness in your work. **Always ensure that the dimensionality matches the operations you intend to perform. That’s what we mean by context!**

Proofs of Prop. 5.48

$$\text{In the following, } x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}.$$

- (a) For C an $m \times n$ matrix and $x \in \mathbb{R}^n$, $\frac{\partial}{\partial x}(Cx) = C$. Express $C = \begin{bmatrix} C_1 & C_2 & \cdots & C_n \end{bmatrix}$ in terms of its columns. Then $Cx = x_1C_1 + x_2C_2 + \cdots + x_nC_n$. Hence, $\frac{\partial(Cx)}{\partial x_i} = C_i$, the i -th column of C . Therefore, $\frac{\partial}{\partial x}(Cx) = \begin{bmatrix} \frac{\partial Cx}{\partial x_1} & \frac{\partial Cx}{\partial x_2} & \cdots & \frac{\partial Cx}{\partial x_n} \end{bmatrix} = C$.

Notice what we did: we broke it down to first computing $\frac{\partial}{\partial x_i}$ and then stacked things into rows or columns as required. Keep this in mind when you get stuck making vector computations with derivatives.

- (b) For $x \in \mathbb{R}^n$, $\frac{\partial}{\partial x}(x^\top x) = 2x^\top$.

$x^\top x = x_1^2 + x_2^2 + \cdots + x_n^2$. Hence, $\frac{\partial}{\partial x_i}(x^\top x) = 2x_i$. Therefore, $\frac{\partial(x^\top x)}{\partial x} = \begin{bmatrix} \frac{\partial(x^\top x)}{\partial x_1} & \frac{\partial(x^\top x)}{\partial x_2} & \cdots & \frac{\partial(x^\top x)}{\partial x_n} \end{bmatrix} = \begin{bmatrix} 2x_1 & 2x_2 & \cdots & 2x_n \end{bmatrix} = 2x^\top$.

Alternative proof based on the Product Rule in Prop. 5.45. We identify $f(x) = x$ and $g(x) = x$. Hence, $\frac{\partial f(x)}{\partial x} = \frac{\partial g(x)}{\partial x} = \frac{\partial x}{\partial x} = I_n$, where I_n is the $n \times n$ identity matrix. This gives the result immediately.

- (c) For $x \in \mathbb{R}^n$, $\nabla(x^\top x) = 2x$.

The gradient of a scalar-valued function, being in this textbook, is always the transpose of the Jacobian. The result is immediate from the previous part or from Prop. 5.45.

- (d) For Q an $n \times n$ symmetric matrix and $x \in \mathbb{R}^n$, $\frac{\partial}{\partial x}(x^\top Qx) = 2x^\top \cdot Q$.

This part follows from the next part by taking transposes and noting that that Q symmetric implies that $Q^\top = Q$.

- (e) For Q an $n \times n$ symmetric matrix and $x \in \mathbb{R}^n$, $\nabla(x^\top Qx) = 2Qx$.

We apply Prop. 5.45 by identifying $f(x) = x$ and $g(x) = Qx$. Hence, $\frac{\partial f(x)}{\partial x} = I_n$, and by part (a) of this problem, $\frac{\partial g(x)}{\partial x} = Q$. Hence,

$$\nabla(x^\top Qx) = I_n^\top \cdot Qx + Q^\top \cdot x = I_n \cdot Qx + Q^\top x = 2Qx,$$

because both I_n and Q are symmetric.

- (f) For Q an $n \times n$ symmetric matrix and $x \in \mathbb{R}^n$, $\frac{\partial \nabla(x^\top Qx)}{\partial x} = 2Q$ (yes, this is the Hessian of $x^\top Qx^\top$).

Apply part (a) to part (e).

5.8 The Total Derivative or the Chain Rule on Steroids

Our journey into the realm of differentiation began with the study of functions of a single variable. It was a straightforward world where changes took place in a one-dimensional landscape. But as with most explorations, our innate curiosity led us to broaden our

horizons. We ventured into the more intricate domain of functions of several variables, acquainting ourselves with the powerful tool of Jacobians.

Now, imagine combining the simplicity of our initial one-dimensional world with the complexity of the multi-dimensional one. Enter the **total derivative**. This handy tool allows us to efficiently compute the derivative of a function of a single scalar variable, t , that emerges from the composition of two functions: $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $g : \mathbb{R} \rightarrow \mathbb{R}^n$. In essence, the total derivative marries our initial knowledge with our further explorations, showcasing the true power of differentiation and meriting its place as the last section in this Chapter.

Here's the intuition: Suppose we are given a scalar-valued function of three variables, $f(x, y, z)$, where each of the variables in turn depends on time, denoted as $x(t)$, $y(t)$, and $z(t)$. Then the **total derivative of f with respect to t** is

$$\frac{df}{dt} = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt} + \frac{\partial f}{\partial z} \frac{dz}{dt}, \quad (5.43)$$

which we recognize as being (the rate of change of f with respect to x) **times** (the rate of change of x with respect to t), plus dot dot dot, **plus** (the rate of change of f with respect to z) **times** (the rate of change of z with respect to t). Hence, (5.44) is a vector version of the Chain Rule given in Prop. 5.11 Differentiation Rules, parts (d) and (e). Because (5.43) **adds up all of the individual rates of change**, it is given the name **total derivative**.

You may be thinking, “Couldn’t we have simply plugged in $x(t)$, $y(t)$, and $z(t)$ into f and then differentiated with respect to t ?” and you would be absolutely correct. It is equivalent, and in some instances, that would be the easier thing to do. But just as it is easier to compute

$$\frac{d}{dx} (\sin(x^2)) = (\cos(x^2)) \cdot (2x) = 2x \cdot \cos(x^2)$$

with the Chain Rule, than it is to memorize a ginormous table of derivatives that includes the above function, it can be easier to compute the total derivative using a formula such as (5.43). That is all that is going on here.

This discussion is formalized in the following definition, which also treats the case of vector valued multivariable functions.

Total Derivative of a Multivariable Function

Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a differentiable function and $x : \mathbb{R} \rightarrow \mathbb{R}^n$ by $x(t) = g(t)$ is also differentiable.

Definition 5.49. The **total derivative** of $f(g(t))$ is $\frac{d}{dt}f(g(t))$, which can be computed as

$$\frac{d}{dt}f(g(t)) = \left. \frac{\partial f(x)}{\partial x} \right|_{x=g(t)} \cdot \frac{dg(t)}{dt}. \quad (5.44)$$

Notes:

- $\left. \frac{\partial f(x)}{\partial x} \right|_{x=g(t)}$ is the Jacobian of $f(x)$ evaluated at $x = g(t)$.
- $\frac{dg(t)}{dt}$ is the ordinary derivative of $g(t)$ with respect to t . When g is a vector of functions of t , you simply differentiate each of its components.
- When $f(x)$ is scalar-valued, that is $f : \mathbb{R}^n \rightarrow \mathbb{R}$, then expanding out the terms in $\frac{\partial f(x)}{\partial x}$ yields a formula like (5.43).
- Indeed, to obtain a formula similar to (5.43) for vector-valued functions, we rewrite (5.44) as follows:
 - let $J_f(x) := \frac{\partial f(x)}{\partial x}$ denote the Jacobian of $f(x)$ and
 - let $g'(t) := \frac{dg(t)}{dt}$ be the ordinary derivative of g with respect to t ;
 - then, the total derivative becomes

$$\frac{d}{dt}f(g(t)) = J_f(g(t)) \cdot g'(t).$$

We'll illustrate the total derivative on (i) two basic examples to show the mechanics of its computation and (ii) on an example that is much closer to its real use for us: an application to an equation with derivatives, aka, a differential equation. In Chapter 6.4 we begin

laying the groundwork for differential equations through the study of Lagrange's powerful reformulation of Newtonian Mechanics (i.e., $F = ma$), allowing us to derive dynamical models of multi-link robots and other things that move and go bump in the night.

Example 5.50. *Mechanics of the Total Derivative:* Suppose that $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ by

$$f(x_1, x_2) = x_1 \cdot \sin(x_2)$$

and

$$g(t) = \begin{bmatrix} g_1(t) \\ g_2(t) \end{bmatrix} = \begin{bmatrix} t^2 \\ t^3 \end{bmatrix}.$$

Compute $\frac{df(t)}{dt} := \frac{d}{dt} f(g(t))$.

Solution: $\frac{df(t)}{dt} = 2t \cdot \sin(t^3) + 3t^4 \cdot \cos(t^3)$.

Using the Chain Rule,

$$\frac{df(g(t))}{dt} = \frac{\partial f(x)}{\partial x} \Big|_{x=g(t)} \cdot \begin{bmatrix} \frac{dg_1(t)}{dt} \\ \frac{dg_2(t)}{dt} \end{bmatrix}.$$

We compute that

$$\begin{aligned} \frac{\partial f(x)}{\partial x} &= \begin{bmatrix} \sin(x_2) & x_1 \cos(x_2) \end{bmatrix} \\ \frac{\partial f(x)}{\partial x} \Big|_{x=g(t)} &= \begin{bmatrix} \sin(t^2) & t^2 \cdot \cos(t^3) \end{bmatrix}. \end{aligned}$$

and

$$\begin{bmatrix} \frac{dg_1(t)}{dt} \\ \frac{dg_2(t)}{dt} \end{bmatrix} = \begin{bmatrix} 2t \\ 3t^2 \end{bmatrix}.$$

Putting it all together, we obtain

$$\frac{df(g(t))}{dt} = \frac{\partial f(x)}{\partial x} \Big|_{x=g(t)} \cdot \begin{bmatrix} \frac{dg_1(t)}{dt} \\ \frac{dg_2(t)}{dt} \end{bmatrix} = \begin{bmatrix} \sin(t^2) & t^2 \cdot \cos(t^3) \end{bmatrix} \cdot \begin{bmatrix} 2t \\ 3t^2 \end{bmatrix} = 2t \sin(t^3) + 3t^4 \cos(t^3).$$

You will need to do examples like this by hand in various courses. In code, to check your work, you could do the following:

```

1 using Symbolics
2
3 # Create a custom function to make differentiation convenient
4 function deriv(f, x)
5     return expand_derivatives.(Symbolics.Differential(x)(f))
6 end
7
8 @variables x1, x2, t
9
10 f = x1*sin(x2)
11
12 # Compute the Jacobian of f with respect to x1, x2
13 Jac = Symbolics.jacobian([f], [x1, x2])
14
15 # Define the functions x1(t), x2(t)

```

```

16 x1_def = t^2
17 x2_def = t^3
18
19
20 # Compute the derivatives of x1(t) and x2(t) with respect to t
21 derivatives = [deriv(x1_def, t), deriv(x2_def, t)]
22
23 # Compute the total derivative of f using the chain rule (matrix multiplication)
24 df_dt_chain = Jac * derivatives
25
26 # Substitute in for x1(t) and x2(t)
27 df_dt = substitute(df_dt_chain, Dict(x1 => x1_def, x2 => x2_def))
28
29 # The result of df_dt_chain = Jac * derivatives is a 1 x 1 column vector.
30 # In Julia, we need to extract the first (and only) element of this
31 # vector using df_dt_chain[1] to produce a scalar-valued function. Why?
32 # The simplify command in Symbolics expects a scalar and not a vector.
33 # You can also use broadcasting to make it work
34 #
35 df_dt = Symbolics.simplify(df_dt[1], expand=true) # extract scalar
36
37
38 println("\nThe total derivative of f with respect to t is:\n\n", df_dt)

```

Output

The total derivative of f with respect to t is:

$$2t \sin(t^3) + 3(t^4) \cos(t^3)$$

Example 5.51. *Mechanics of the Total Derivative:* Suppose that $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ by

$$f(x, y, z) = x^3 + y^3 + z^3$$

and

$$\begin{bmatrix} x(t) \\ y(t) \\ z(t) \end{bmatrix} = \begin{bmatrix} e^t \\ e^t \sin(t) \\ e^t \cos(t) \end{bmatrix}.$$

Compute $\frac{df(t)}{dt} := \frac{d}{dt} f(x(t), y(t), z(t))$.

Solution: $\frac{df(t)}{dt} = 3e^{3t} - 3\cos^2(t) \cdot \sin(t) + 3\cos(t) \cdot \sin^2(t)$.

We solve this one using non-vector notation so that you can find a preference: Using the Chain Rule and leaving out the arguments on the functions, we have

$$\frac{df}{dt} = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt} + \frac{\partial f}{\partial z} \frac{dz}{dt}.$$

We compute that

$$\frac{\partial f}{\partial x} = 3x^2, \quad \frac{\partial f}{\partial y} = 3y^2, \quad \text{and} \quad \frac{\partial f}{\partial z} = 3z^2,$$

and

$$\frac{dx}{dt} = e^t, \quad \frac{dy}{dt} = e^t \sin(t) + e^t \cos(t), \quad \text{and} \quad \frac{dz}{dt} = e^t \cos(t) - e^t \sin(t).$$

Combining these values and doing the required substitutions, we obtain

$$\begin{aligned} \frac{df}{dt} &= \underbrace{3(e^t)^2}_{\frac{\partial f}{\partial x}} \cdot \underbrace{\left(\frac{dx}{dt}\right)}_{e^t} + \underbrace{3(e^t \sin(t))^2}_{\frac{\partial f}{\partial y}} \cdot \underbrace{\left(e^t \sin(t) + e^t \cos(t)\right)}_{\frac{dy}{dt}} + \underbrace{3(e^t \cos(t))^2}_{\frac{\partial f}{\partial z}} \cdot \underbrace{\left(e^t \cos(t) - e^t \sin(t)\right)}_{\frac{dz}{dt}} \\ &= (\text{lol, where's Julia?}) \\ &= 3e^{3t} - 3\cos^2(t) \cdot \sin(t) + 3\cos(t) \cdot \sin^2(t). \end{aligned}$$

Some Perspective: The above illustrates the “grind of Calculus”. Everyone dislikes it. You will have to do some messy calculations in your Science and Engineering courses. Once again, this is just part of the learning process, akin to running wind sprints in sports—they’re no fun, but they make you better in the game! As a practicing STEM professional, you will use “productivity tools”, such as Julia or MATLAB. As a student, you can always check your HW with Julia! **On an exam, you are unlikely to face a math calculation like this one because the instructor is evaluating your mastery of some newly acquired STEM skills (we hope) and not your ability to do fussy algebra.**

```

1 using Symbolics
2
3 # Create a custom function to make differentiation convenient
4 function deriv(f, x)
5     return expand_derivatives.(Symbolics.Differential(x)(f))
6 end
7
8 @variables x, y, z, t
9
10 f = [x^3 + y^3 + z^3]
11
12 # Compute the Jacobian of f with respect to x, y, z
13 Jac = Symbolics.jacobian(f, [x, y, z])
14
15 # Define the functions x(t), y(t), and z(t)
16 x_def = exp(t)
17 y_def = sin(t)
18 z_def = cos(t)
19
20 # Compute the derivatives of x(t), y(t), and z(t) with respect to t
21 derivatives = [deriv(x_def, t), deriv(y_def, t), deriv(z_def, t)]
22
23 # Compute the total derivative of f using the chain rule (matrix multiplication)
24 df_dt_chain = Jac * derivatives
25
26 # Substitute in for x(t), y(t), and z(t)
27 df_dt = substitute(df_dt_chain, Dict(x => x_def, y => y_def, z => z_def))
28
29 # The result of df_dt_chain = Jac * derivatives is a 1 x 1 column vector.
30 # In Julia, we need to extract the first (and only) element of this
31 # vector using df_dt_chain[1] to produce a scalar-valued function. Why?
32 # The simplify command in Symbolics expects a scalar and not a vector.
33 # You can also use broadcasting to make it work
34 #
35 df_dt = Symbolics.simplify.(df_dt, expand=true) # broadcasting, note the dot
36
37 println("\nThe total derivative of f with respect to t using the chain rule is:\n\n", df_dt)

```

Output

The total derivative of f with respect to t using the chain rule is:

$$3\exp(3t) - 3(\cos(t)^2) * \sin(t) + 3\cos(t) * (\sin(t)^2)$$

Setting the stage for the next example: The following is called a linear system of differential equations,

$$\begin{bmatrix} \frac{dx_1(t)}{dt} \\ \frac{dx_2(t)}{dt} \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -1 & -1 \end{bmatrix} \begin{bmatrix} x_1(t) \\ x_2(t) \end{bmatrix} = \begin{bmatrix} x_2(t) \\ -x_1(t) - x_2(t) \end{bmatrix}. \quad (5.45)$$

Physically, this system of equations, with derivatives on the left side and unknown (aka, to be found) functions of time on the right side, represents a linear approximation of a pendulum with friction at the pivot: $x_1(t)$ is the angle of the pendulum and $x_2(t)$ is the angular velocity of the pendulum.

We do not yet know how to find the unknown functions, $x_1(t)$ and $x_2(t)$. When we come to this later in the book, we'll understand that a solution to this system of equations is not so different than what is required for solving any system of equations: we need to find functions that, when substituted into the left side (which involves taking their derivatives), equals the right side (which involves only substituting in the functions). That's why they are called equations: an equality holds.

The next example asks (and answers) an intriguing question: without knowing the solutions to the above differential equations, can we prove that the solutions do not explode? That is,

$$\lim_{t \rightarrow \infty} (x_1(t))^2 + (x_2(t))^2 \quad (5.46)$$

is finite? And just maybe, that it even equals zero? Huh? We do not know the solution, but we seek to determine some of its properties anyway? That seems like a power mindset! **To answer these questions, we will use the total derivative.**

Example 5.52. For the differential equations in (5.45):

- (a) Show that the limit in (5.46) does not blow up.
- (b) Bonus: Give a plausible argument that the limit is, in fact, equal to zero.

Solution: Define a function $V(x_1, x_2) := (x_1)^2 + (x_2)^2$, which is the norm squared of the vector (x_1, x_2) . We compute its total derivative with respect to time, $\frac{dV(x_1(t), x_2(t))}{dt}$, using the Chain Rule,

$$\begin{aligned} \frac{d}{dt} V(x_1(t), x_2(t)) &= \frac{\partial V(x_1(t), x_2(t))}{\partial x_1} \cdot \frac{dx_1(t)}{dt} + \frac{\partial V(x_1(t), x_2(t))}{\partial x_2} \cdot \frac{dx_2(t)}{dt} \\ &= \underbrace{2x_1(t) \cdot \frac{dx_1(t)}{dt}}_{\frac{\partial V}{\partial x_1}} + \underbrace{2x_2(t) \cdot \frac{dx_2(t)}{dt}}_{\frac{\partial V}{\partial x_2}} \\ &= 2x_1(t) \cdot x_2(t) + 2x_2(t) \cdot (-x_1(t) - x_2(t)) \quad (\text{plugging in from (5.45)}) \\ &= -2(x_2(t))^2. \end{aligned}$$

Because

$$\frac{d}{dt} V(x_1(t), x_2(t)) = -2(x_2(t))^2 \leq 0 \quad (5.47)$$

for all t , Prop. 5.23 implies that the function $V(x_1(t), x_2(t))$ is nonincreasing! We do not know the actual time-function $V(t) = V(x_1(t), x_2(t))$, but we know that its initial value at time $t = 0$ is its largest possible value (otherwise, it could not be nonincreasing). Therefore, the limit in (5.46), if it exists, cannot go to infinity. **That settles the first part of the question, which is already an awesome demonstration of the total derivative in action.** What about the second part of the question?

Bonus: The limit likely goes to zero, and the reasoning is as follows. If a pendulum is not at rest at the downward position, it oscillates. That's what pendulums are designed to do. Now, as long as the angular velocity $x_2(t)$ of the pendulum is non-zero, the time rate of change of $V(x_1(t), x_2(t))$ is negative, so $V(x_1(t), x_2(t))$ continues to decrease. The angular velocity of the pendulum

goes to zero when the pendulum reaches one of its two “extreme swing angles” and reverses course. So, as long as the pendulum is not at rest in its downward position, $V(x_1(t), x_2(t))$ is decreasing. It cannot decrease below zero because $(x_1(t))^2 + (x_2(t))^2$ represents the norm squared of the vector $(x_1(t), x_2(t))$. Hence $V(x_1(t), x_2(t))$ likely converges to zero as $t \rightarrow \infty$. At Michigan, the course EECS 562: Nonlinear Systems and Control makes this argument fully rigorous.

1 using Symbolics

```

2
3 # Create a custom function to make differentiation convenient
4 function deriv(f, x)
5     return expand_derivatives.(Differential(x)(f))
6 end
7
8 @variables x1 x2
9 x=[x1 x2]
10
11 # Define the differential equations as a vector
12 dx_dt = [x2; -x1 - x2]
13
14 # Define a function V(x)
15 V = x1^2 + x2^2
16
17 # Compute the Jacobian of V with respect to x
18 Jac_V = Symbolics.jacobian([V], x)
19
20 # Compute the total derivative of V using the chain rule
21 dV_dt = Jac_V * dx_dt
22
23 # The result of the matrix multiplication Jac_V * dx_dt
24 # is a 1 x 1column vector.
25 # Since we are interested in the scalar value of the derivative,
26 # we extract the first (and only) element of this vector using dV_dt[1].
27 #
28 # Simplify the result and expand polynomial expressions
29 #
30 dV_dt = Symbolics.simplify(dV_dt[1], expand=true)
31
32 println("The total derivative of V with respect to t using the chain rule is:\n\n", dV_dt)

```

Output

The total derivative of V with respect to t using the chain rule is:

$-2(x_2^2)$

Bottom Line:

- **Multidimensional Chain Rule:** The total derivative can be thought of as the chain rule’s multidimensional counterpart. While the partial derivative captures the rate of change of a multivariable function with respect to a single variable, the total derivative encapsulates the compounded effect of each variable’s rate of change on the function. It’s the chain rule, but on a grander scale.
- **Versatility:** While plugging in $x(t)$, $y(t)$, and $z(t)$ directly into f and differentiating is a straightforward approach, the total derivative offers a more structured method. This becomes especially valuable when dealing with functions where the relationship between variables isn’t immediately clear (such as our differential equation) or when the function’s form makes direct differentiation less intuitive.
- **Efficiency:** Just as the chain rule can simplify differentiation for certain composite functions, the total derivative can be more efficient for functions of multiple variables. Instead of working out derivatives for every possible combination of functions, the

total derivative provides a systematic approach that breaks the computation down into simpler pieces. Recall that Automatic Differentiation, heavily used in Machine Learning, also uses the Chain Rule to decompose a complicated derivative into simpler parts.

Secrets of the Arcane 5.53: Let's ask ChatGPT4

Grizzle to ChatGPT4: What do you think is the most important mathematical tool used for training your neural networks?

Answer: The training of neural networks, like mine, relies on a combination of mathematical concepts, algorithms, and techniques. If we had to pinpoint the most fundamental mathematical tool, it would be:

1. **Calculus and the Chain Rule:** The backpropagation algorithm, which is central to training deep neural networks, is essentially an application of the chain rule from calculus. It computes the **gradient of the loss function with respect to each weight** by propagating the gradient backward through the network.

Other important mathematical tools and concepts include:

2. **Linear Algebra:** Neural networks involve a lot of matrix multiplications and transformations, especially in architectures like convolutional neural networks (CNNs) and recurrent neural networks (RNNs).
3. **Probability and Statistics:** These are crucial for understanding concepts like regularization, dropout, and batch normalization. They also play a role in the design of loss functions and the interpretation of neural network outputs, especially in classification tasks.
4. **Optimization:** Gradient descent and its variants (like stochastic gradient descent, Adam, RMSprop, etc.) are optimization algorithms used to adjust the weights of the network to minimize the loss function.
5. **Differential Equations:** Especially relevant in the context of understanding and designing recurrent neural networks and more recent architectures like Neural ODEs.
6. **Information Theory:** Concepts like entropy and cross-entropy, which are foundational to certain loss functions (e.g., the cross-entropy loss for classification tasks).

While the chain rule and calculus form the backbone of the training process, the successful training and deployment of neural networks require a holistic understanding and application of various mathematical tools and concepts.

Grizzle: Just like Robotics!

13 October 2023, Saint-Rémy-de-Provence, France

5.9 (Optional Read): Proofs Associated with the Chapter

Proof of Prop. 5.5 (Differentiable at x_0 Implies Continuous at x_0)

Suppose $f : [a, b] \rightarrow \mathbb{R}$ is differentiable at $x_0 \in [a, b]$. Then, f is continuous at x_0 .

Proof: If we assume f is differentiable at x_0 , then

$$L := \lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h}$$

exists and is finite. We can rewrite the above as

$$\lim_{h \rightarrow 0} \left[\frac{f(x_0 + h) - f(x_0)}{h} - L \right] = 0,$$

or in other symbols, if we define

$$\text{err}(h) := \frac{f(x_0 + h) - f(x_0)}{h} - L$$

to be the error in the limit for $h \neq 0$, then $\lim_{h \rightarrow 0} \text{err}(h) = 0$. But that, in turn by Prop. 2.35, implies that $\lim_{h \rightarrow 0} h \cdot \text{err}(h) = 0$, which gives us

$$\lim_{h \rightarrow 0} [f(x_0 + h) - f(x_0) - h \cdot L] = 0.$$

From this, we deduce that

$$\lim_{h \rightarrow 0} f(x_0 + h) = \lim_{h \rightarrow 0} [f(x_0) - h \cdot L] = f(x_0) - \lim_{h \rightarrow 0} h \cdot L = f(x_0).$$

In other words, f is continuous at x_0 by Definition 4.11. ■

Chapter 6

Engineering Applications of the Derivative

Learning Objectives

By the end of this chapter, the student should be able to:

- Analyze and model engineering problems using the principles of calculus, specifically through the application of derivatives to understand system behaviors.
- Develop strategies for solving optimization problems, both with and without constraints.
- Critically assess the role of derivatives in determining system dynamics and understand the significance of dynamic equations in engineering contexts.
- Apply the concept of energy to the analysis of mechanical systems using the framework of Lagrangian mechanics.

Outcomes

Upon successful completion of this chapter, students will be able to:

- Calculate path length and arc length for given paths.
- Solve engineering problems involving root finding and minimization.
- Use gradient descent to find local minima of functions and understand its limitations.
- Apply second derivative tests to determine the nature of critical points in functions.
- Solve optimization problems involving equality and inequality constraints using Lagrange multipliers.
- Derive and apply Lagrange's equations to solve problems in dynamics.
- Compute kinetic and potential energy for mechanical systems.
- Model and analyze the motion of complex systems like multi-link manipulators using Lagrange's formalism.
- Determine moments of inertia for various bodies and understand their effects on rotational motion.

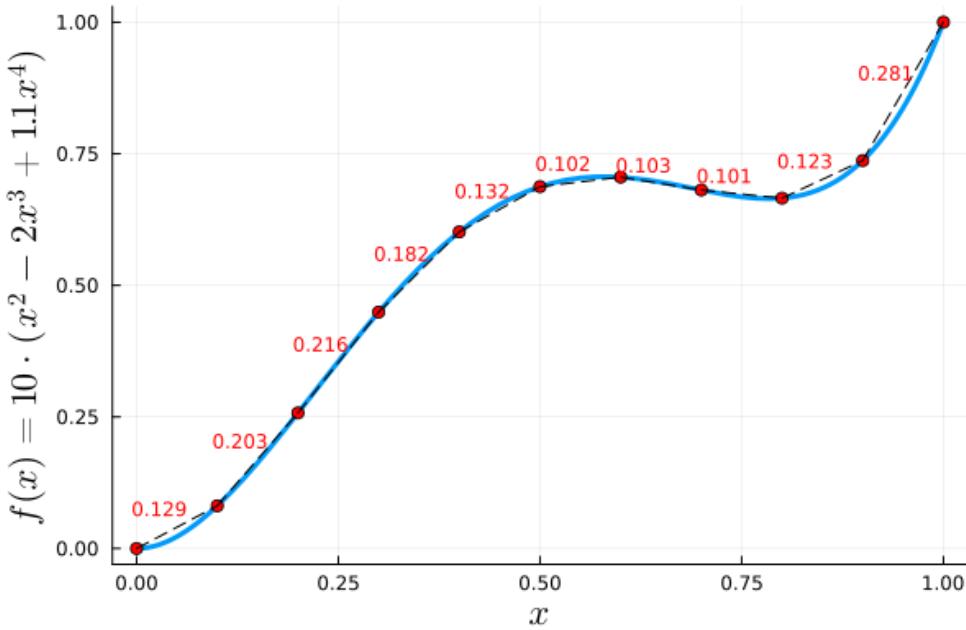


Figure 6.1: Path length is defined by chopping up the curve into short segments, connecting the endpoints of the segments by straight lines, and then adding up the lengths of the lines. As the segment size tends to zero, the true path length is obtained. Simple? Yes. Useful? Very! It can be used to formulate an important problem: determine a path of shortest length for a robot to travel without hitting obstacles; we'll explore this in homework. Table 6.1 shows that the ten segment lengths in the plot add up to 1.572; the true answer, obtained in the limit as the segment lengths shrink to zero and become infinite in number, is $1.5852977678599067 \pm 6.654127165939272e-9$. The approximation with ten segments is already quite good!

6.1 Path Length or Arc Length

You understand the length of a line segment in the plane and the length of the circumference of a circle. But for more complicated curves, what is the length from a point a to a point b ? Figure 6.1 and Table 6.1 illustrate the concept of **arc length (aka, path length)**, which is fundamental in engineering because it allows us to quantify the “length of a curve” in the plane \mathbb{R}^2 , which can be defined by a function $f : (a, b) \rightarrow \mathbb{R}$ or parametric equations $x(t), y(t)$, $t_1 \leq t \leq t_2$, the arc length is the distance one would travel along the curve between two points.

The integral formulas for arc length arise from the Pythagorean theorem. By dividing the curve into infinitesimally small segments, each segment can be approximated by a straight line whose length can be calculated using the Pythagorean theorem as the hypotenuse of a right triangle with sides dx and dy . As we sum up the lengths of all these segments (a process achieved by integrating), we obtain the total arc length of the curve.

Arc Length = Path Length

Definition 6.2. Given a differentiable function $f : [a, b] \rightarrow \mathbb{R}^2$, the arc length S of the curve described by $y = f(x)$ from $x = a$ to $x = b$ is given by the integral, $S := \int_a^b \sqrt{dx^2 + dy^2}$. After factoring out dx , we obtain the standard formula,

$$S = \int_a^b \sqrt{1 + \left(\frac{dy}{dx}\right)^2} dx. \quad (6.1)$$

For a curve defined parametrically by $x(t)$ and $y(t)$ where t ranges from t_1 to t_2 , the arc length is given by,

$$S = \int_{t_1}^{t_2} \sqrt{\left(\frac{dx}{dt}\right)^2 + \left(\frac{dy}{dt}\right)^2} dt. \quad (6.2)$$

Table 6.1: Path Length is Easy to Understand but Tedious to Compute by Hand,

$$S := \int_a^b \sqrt{dx^2 + dy^2} = \int_a^b \sqrt{1 + \left(\frac{dy}{dx}\right)^2} dx \text{ (after factoring out } dx)$$

x_i	y_i	$\sqrt{(x_i - x_{i-1})^2 + (y_i - y_{i-1})^2}$
0.0	0.0	0.0
0.1	0.0811	0.129
0.2	0.2576	0.203
0.3	0.4491	0.216
0.4	0.6016	0.182
0.5	0.6875	0.132
0.6	0.7056	0.102
0.7	0.6811	0.103
0.8	0.6656	0.101
0.9	0.7371	0.123
1.0	1.0	0.281
$\int_a^b \sqrt{dx^2 + dy^2} \approx \sum_{i=2}^{11} \sqrt{(x_i - x_{i-1})^2 + (y_i - y_{i-1})^2} =$		1.572

Notes: Equation 6.2 is the more general formula for computing path length. When the path is given by a function, as in $(x, y(x))$, (6.1) follows from (6.2) as shown below:

$$\begin{aligned} S &= \int_{a=x(t_1)}^{b=x(t_2)} \sqrt{\left(\frac{dx}{dt}\right)^2 + \left(\frac{dy}{dt}\right)^2} dt \text{ (factor out } \left(\frac{dx}{dt}\right)^2) \\ &= \int_{a=x(t_1)}^{b=x(t_2)} \sqrt{1 + \left(\frac{\frac{dy}{dt}}{\frac{dx}{dt}}\right)^2} \left(\frac{dx}{dt}\right) dt \text{ (cancel the } dt's) \\ &= \int_a^b \sqrt{1 + \left(\frac{dy}{dx}\right)^2} dx. \end{aligned}$$

Source of the formula for arc length: Figure 6.1 and Table 6.1 provide the basic idea. Let Δx and Δy be incremental changes in x and y along a curve. Then the distance covered is given by the Pythagorean Theorem, $\Delta S := \sqrt{(\Delta x)^2 + (\Delta y)^2}$. Taking the limit as Δx goes to zero, we obtain $dS = \sqrt{(dx)^2 + (dy)^2}$, and hence

$$S = \int_a^b dS = \int_a^b \sqrt{(dx)^2 + (dy)^2}.$$

But hmm, we have never seen an integral like this! Suppose that the (x, y) coordinates of the curve are parameterized by time. Then, using the linearization (aka, first-order Taylor expansion), we have

$$\Delta x = x(t + \Delta t) - x(t) \xrightarrow{\Delta t \rightarrow 0} \left(\frac{dx(t)}{dt}\right) \cdot dt$$

$$\Delta y = y(t + \Delta t) - y(t) \xrightarrow{\Delta t \rightarrow 0} \left(\frac{dy(t)}{dt}\right) \cdot dt$$

Substituting into the previous equation and factoring out dt gives the result.

Example 6.3. (*Classical Example because the Integral is Easy to Evaluate*). A circle of radius $r > 0$ about the origin can be parameterized by $(x(t) = r \sin(t), y(t) = r \cos(t))$, for $0 \leq t \leq 2\pi$. Compute the arc length from $t_1 = \pi/3$ to $t_2 = \pi$.

Solution: **Ans.** $S = 2r\frac{\pi}{3}$.

We apply (6.2).

- $x(t) = r \sin(t) \implies \frac{dx}{dt} = r \cos(t)$.
- $y(t) = r \cos(t) \implies \frac{dy}{dt} = -r \sin(t)$.
- $\sqrt{\left(\frac{dx}{dt}\right)^2 + \left(\frac{dy}{dt}\right)^2} = r$.

Hence,

$$S = \int_{t_1}^{t_2} \sqrt{\left(\frac{dx}{dt}\right)^2 + \left(\frac{dy}{dt}\right)^2} dt = \int_{\pi/3}^{\pi} r dt = r \cdot t \Big|_{\pi/3}^{\pi} = r (\pi - \pi/3) = 2r\pi/3.$$

■

Example 6.4. Compute the path length of the curve defined by $y = f(x) = x^2$ from $x = 0$ to $x = 1$.

Solution: **Ans.** $S = 1.479$.

The derivative $\frac{dy}{dx}$ is $2x$. Hence, the arc length is calculated as:

$$S = \int_0^1 \sqrt{1 + (2x)^2} dx. \quad (6.3)$$

Setting up the integral is the main skill to be mastered here. In Chapter 7, we'll learn how to compute closed-form solutions¹ for such integrals when the function defining the path is simple enough, such as $f(x) = x^2$. For now, we seek a numerical answer, which is what real engineers must do in 99.9% of cases.

using QuadGK

```

1
2
3 # Define the function for the curve
4 f(x) = x^2
5
6 # Define the derivative of the function
7 df(x) = 2x
8
9 # Define the integrand for arc length
10 function arc_length_integrand(x)
11     return sqrt(1 + df(x)^2)
12 end
13
14 # Compute the arc length using QuadGK
15 integral, error = quadgk(arc_length_integrand, 0, 1)
16
17 println("The arc length of the curve from x = 0 to x = 1 is approximately: ", integral)
18 println("The estimated error of the computation is: ", error)

```

Output

$$^1 S = \int_0^1 \sqrt{1 + (2x)^2} dx = \int_0^2 \sqrt{1 + u^2} \frac{du}{2} = \frac{1}{2} \int_0^2 \sqrt{1 + u^2} du = \operatorname{asinh}(u) \Big|_0^2 = \frac{1}{2} \operatorname{asinh}(2),$$

The arc length of the curve from $x = 0$ to $x = 1$ is approximately: 1.4789428575446002

The estimated error of the computation is: 1.0424397123287577e-8

Example 6.5. Compute the arc length of the curve defined by $y = f(x) = 10(x^2 - 2x^3 + 1.1x^4)$ from $x = 0$ to $x = 1$.

Solution: **Ans.** $S = 1.585$.

The derivative $\frac{dy}{dx}$ is $10(2x - 6x^2 + 4.4x^3)$. Hence, the arc length is calculated as:

$$S = \int_0^1 \sqrt{1 + (10(2x - 6x^2 + 4.4x^3))^2} dx. \quad (6.4)$$

Try as one may with techniques to be developed in Chapter 7, computing a closed-form solution to the integral is likely impossible. In any case, neither ChatGPT nor your author is able to do it! For QuadGK, it is cake:

```
1 using QuadGK
2
3 # Define the function for the curve
4 f(x) = 10*(x^2 - 2*x^3 + 1.1*x^4)
5
6 # Define the derivative of the function
7 df(x) = 10*(2x - 6x^2 + 4.4x^3)
8
9 # Define the integrand for arc length
10 function arc_length_integrand(x)
11     return sqrt(1 + df(x)^2)
12 end
13
14 # Compute the arc length using QuadGK
15 integral, error = quadgk(arc_length_integrand, 0, 1)
16
17 println("The arc length of the curve from x = 0 to x = 1 is approximately: ", integral)
18
19 println("The estimated error of the computation is: ", error)
```

Output

The arc length of the curve from $x = 0$ to $x = 1$ is approximately: 1.5852977678599067

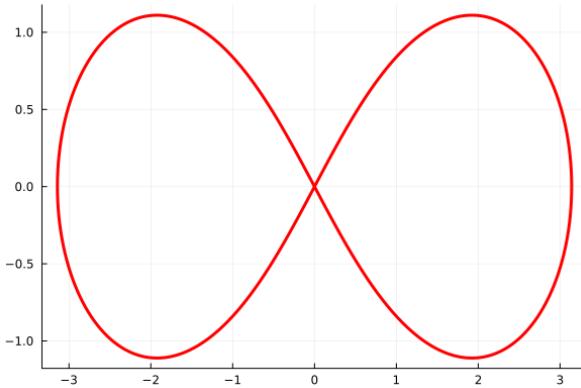
The estimated error of the computation is: 6.654127165939272e-9

Example 6.6. Figure 6.2 depicts a figure eight parameterized by

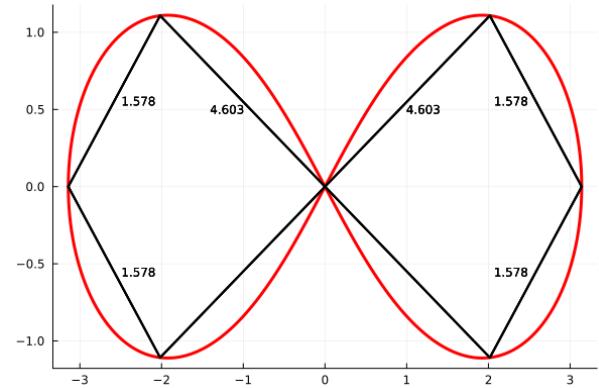
$$x(t) = \pi \cdot \frac{\cos\left(\frac{t \cdot 2\pi}{4}\right)}{1 + \sin^2\left(\frac{t \cdot 2\pi}{4}\right)}$$

$$y(t) = \pi \cdot \frac{\sin\left(\frac{t \cdot 2\pi}{4}\right) \cdot \cos\left(\frac{t \cdot 2\pi}{4}\right)}{1 + \sin^2\left(\frac{t \cdot 2\pi}{4}\right)}$$

for $0 \leq t \leq 4$. Compute its arc length.



(a)



(b)

Figure 6.2: A figure eight is clearly not a function from x to y . (a) Rather, it is a parameterized curve in the plane. Example 6.6 explains what is going on. (b) Shows a very crude straight line underapproximation of the arc length yielding $S \approx 4 \cdot (1.578) + 2 \cdot (4.603) = 15.518$. How far do you think it is off? The two lines of length 4.603 cross through the origin to hit the other side of the figure.

Solution: **Ans.** $S = 16.47$

This time, the path in \mathbb{R}^2 is not specified by a function, but instead by a parameterized curve $(x(t), y(t))$, $0 \leq t \leq 4$. This is more general than a function because the curve can cross itself, fail the vertical line test, and do all sorts of crazy and wonderful things, just like a drone flying in the air. To compute arc length, we need to apply the formula

$$S = \int_{t_1}^{t_2} \sqrt{\left(\frac{dx}{dt}\right)^2 + \left(\frac{dy}{dt}\right)^2} dt = \int_{t_1}^{t_2} \sqrt{(x'(t))^2 + (y'(t))^2} dt.$$

We compute the derivatives using software and obtain

$$\begin{aligned} x'(t) &= -\frac{\pi \cdot 2\pi \cdot \sin\left(\frac{2\pi t}{4}\right)}{4 \cdot (\sin^2\left(\frac{2\pi t}{4}\right) + 1)} - \frac{\pi \cdot 4\pi \cdot \sin\left(\frac{2\pi t}{4}\right) \cdot \cos^2\left(\frac{2\pi t}{4}\right)}{4 \cdot (\sin^2\left(\frac{2\pi t}{4}\right) + 1)^2} \\ y'(t) &= -\frac{\pi \cdot 2\pi \cdot \sin^2\left(\frac{2\pi t}{4}\right)}{4 \cdot (\sin^2\left(\frac{2\pi t}{4}\right) + 1)} + \frac{\pi \cdot 2\pi \cdot \cos^2\left(\frac{2\pi t}{4}\right)}{4 \cdot (\sin^2\left(\frac{2\pi t}{4}\right) + 1)} - \frac{\pi \cdot 4\pi \cdot \sin^2\left(\frac{2\pi t}{4}\right) \cdot \cos^2\left(\frac{2\pi t}{4}\right)}{4 \cdot (\sin^2\left(\frac{2\pi t}{4}\right) + 1)^2} \end{aligned}$$

We're all set now to do the integral numerically.

```

1 # Figure 8
2 using QuadGK, Plots, LaTeXStrings
3
4 function x_position(t)
5     a = pi; T = 4
6     a * cos(t*2*pi/T) / (1 + sin(t*2*pi/T)^2)
7 end
8 function y_position(t)
9     a = pi; T = 4
10    a * sin(t*2*pi/T) * cos(t*2*pi/T) / (1 + sin(t*2*pi/T)^2)
11 end
12
13 # Compute the arc length
14
15 # Derivatives
16 function dx_dt(t)

```

```

17   a = pi; T = 4
18   -2*pi*a*sin(2*pi*t/T) / (T*(sin(2*pi*t/T)^2 + 1)) - 4*pi*a*sin(2*pi*t/T)*cos(2*pi*t/T)^2 / (T*(sin(2*pi*t/T)^2 + 1)^2)
19 end
20
21 function dy_dt(t)
22   a = pi; T = 4
23   -2*pi*a*sin(2*pi*t/T)^2 / (T*(sin(2*pi*t/T)^2 + 1)) + 2*pi*a*cos(2*pi*t/T)^2 / (T*(sin(2*pi*t/T)^2 + 1)) - 4*pi*a*sin(2*pi*t/T)^2*cos(2*pi*t/T)^2 / (T*(sin(2*pi*t/T)^2 + 1)^2)
24 end
25
26 # Define the integrand for arc length
27 function arc_length_integrand(t)
28   return sqrt(dx_dt(t)^2 + dy_dt(t)^2)
29 end
30
31 # Compute the arc length using QuadGK
32 integral, error = quadgk(arc_length_integrand, 0, 4)
33
34 println("The arc length of the curve from x = 0 to x = 4 is approximately: ", integral)
35 println("The estimated error of the computation is: ", error)
36
37 # make a nice plot
38
39 dt = 0.01
40 t = collect(0.0:dt:4.0)
41 xvalues = x_position.(t)
42 yvalues = y_position.(t)
43
44 p1 = plot(xvalues, yvalues, lw=3, color=:red, aspect_ratio = 1, gudefont = 15, label=false)

```

Output

The arc length of the curve from $x = 0$ to $x = 4$ is approximately: 16.474873499707456

The estimated error of the computation is: 1.5850782530968388e-7

Hence, the very crude approximation in Figure 6.2-(b) is off by about 10%. ■

Helpful videos: Keep in mind that setting up the integral is the main thing to master. Once you have the problem properly formulated, then you can employ software tools to find the answer.

- Arc Length (formula explained) by **BlackPenRedPen**.
- Arc length intro | Applications of definite integrals by the Khan Academy.
- Find the Arc Length of a Function with Calculus by Quoc Dat Phung. Positive: He has strong graphics. Negative: He uses manual integration instead of numerical integration. We cover how to manually compute similar integrals in Chapter 7. For now, just watch how he sets up the problems.
- Arc Length Calculus Problems by the Organic Chemistry Tutor.

6.2 Root Finding and Unconstrained Minimization

The methods of this section are grounded in the linear approximation of a nonlinear function developed in Prop. 5.39. Root finding uses

$$f(x) \approx f(x_0) + \frac{\partial f(x_0)}{\partial x} \cdot (x - x_0), \quad f: \mathbb{R}^n \rightarrow \mathbb{R}^n,$$

while minimization uses the equivalent form developed in (5.36),

$$f(x) \approx f(x_0) + \nabla f(x_0) \bullet (x - x_0), \quad f: \mathbb{R}^n \rightarrow \mathbb{R},$$

where the heavy dot, \bullet , denotes the dot product of two column vectors.

6.2.1 Root Finding

We consider a function $f: \mathbb{R}^n \rightarrow \mathbb{R}^n$ and seek a **root**, that is, a point $x_0 \in \mathbb{R}^n$ such that $f(x_0) = 0$. Note that the domain and range are both \mathbb{R}^n , and thus, this is the nonlinear equivalent of solving a square linear system of equations, $Ax - b = 0$. We recall that $\det(A) \neq 0$ is the magic condition for the existence and uniqueness of solutions to $Ax - b = 0$.

Suppose that x_k is the current guess for a root. Recall from ROB 101 *Computational Linear Algebra* that the **Newton-Raphson Algorithm** uses the linear approximation of f about the point x_k to choose x_{k+1} so that $f(x_{k+1}) \approx 0$. In other symbols, x_{k+1} should satisfy,

$$0 = f(x_{k+1}) \approx f(x_k) + \frac{\partial f(x_k)}{\partial x} \cdot (x_{k+1} - x_k), \quad (6.5)$$

which is equivalent to

$$0 \approx f(x_k) + \frac{\partial f(x_k)}{\partial x} \cdot (x_{k+1} - x_k). \quad (6.6)$$

If the Jacobian $\frac{\partial f(x_k)}{\partial x}$ is invertible, and the approximation sign is replaced with an equals sign, a unique value for x_{k+1} is obtained. Solving this equation iteratively yields the Newton-Raphson Algorithm.

Method 6.7: Newton-Raphson Algorithm

Based on (6.6), we define^a $x_{k+1} := x_k + \Delta x_k$, where Δx_k is our update to x_k . We can then break the algorithm into two steps,

$$\left(\frac{\partial f(x_k)}{\partial x} \right) \Delta x_k = -f(x_k) \quad (\text{solve for } \Delta x_k) \quad (6.7)$$

$$x_{k+1} = x_k + \Delta x_k \quad (\text{use } \Delta x_k \text{ to update our estimate of the root}). \quad (6.8)$$

While for toy problems, one can use the matrix inverse to solve (6.7) for Δx_k , for larger problems, one uses LU Factorization or QR Factorization. Once (6.7) has been solved, x_{k+1} is updated in (6.8) and the process repeats.

A **damped Newton-Raphson Algorithm** is obtained by replacing (6.8) with

$$x_{k+1} = x_k + \epsilon \Delta x_k, \quad (6.9)$$

for some $\epsilon > 0$. The validity of the Newton-Raphson Algorithm rests upon:

- the function f being differentiable;
- the Jacobian $\frac{\partial f(x_k)}{\partial x}$ having a non-zero determinant at points generated by (6.7) and (6.8); and
- [the linear equation $f_{\text{lin}}(x) = f(x_k) + \frac{\partial f(x_k)}{\partial x}(x - x_k)$ being a good approximation to the nonlinear function, $f(x)$.]

^aNote that $\Delta x_k = x_{k+1} - x_k$.

You can find in Chapter 11 of ROB 101 *Computational Linear Algebra* numerous examples with a hand-made implementation of the Newton-Raphson Algorithm. Hence, here, the focus is on Julia Packages. We suggest [NLsolve.jl](#), which is based on Newton-Raphson and important variations of the same.

using NLsolve

```
1 # Import the NLsolve package
2 using NLsolve
3
4 # This script aims to solve a system of nonlinear equations using the Newton-Raphson
   method provided by the NLsolve package.
5
6 # Define the system of nonlinear equations.
7 # The function `system_of_equations!` modifies the input array `F` to represent the
   system of equations.
8 # The system to be solved is:
9 # 1) x[1]^2 + x[2]^2 = 2.0
10 # 2) x[1]^2 - x[2]^2 = 1.0
11 #
12 # Parameters:
13 # - F: An array that will be modified in-place to represent the system of equations.
14 # - x: An array representing the variables of the system.
15 function system_of_equations!(F, x)
16     F[1] = x[1]^2 + x[2]^2 - 2.0    # First equation
17     F[2] = x[1]^2 - x[2]^2 - 1.0    # Second equation
18 end
19
20 # Provide an initial guess for the solution.
21 # The quality of the initial guess can influence the convergence of the Newton-Raphson
   method.
22 # In this case, both x[1] and x[2] are initially guessed to be 1.0.
23 initial_guess = [1.0, 1.0]
24
25 # Use the `nlsolve` function from the NLsolve package to find the roots of the system.
26 # The function returns a result object that contains various information about the
   solution process.
27 result = nlsolve(system_of_equations!, initial_guess)
28
29 # Print the roots (solutions) of the system.
30 println ("Roots found at: ", result.zero)
```

Output

```
Roots found at: [1.2247448713915894, 0.7071067811873448]
```

As you can see from the code, it aims to find the roots of the system of equations,

$$\begin{aligned}x_1^2 + x_2^2 &= 2 \\x_1^2 - x_2^2 &= 1.\end{aligned}$$

In Julia, the ‘bang’ notation, denoted by an exclamation mark (!) at the end of a function name, is a convention used to indicate that the function may modify its arguments in-place. This convention is not enforced by the language itself but is widely adopted in the Julia community to signal the side effects of a function.

For example, consider the function,

```
function system_of_equations!(F, x)
    F[1] = x[1]^2 + x[2]^2 - 2.0
    F[2] = x[1]^2 - x[2]^2 - 1.0
end
```

The presence of the ! in the function name `system_of_equations!` suggests that the function might modify the contents of the arguments `F` and/or `x` directly, rather than returning a new modified value. **Modifying arguments in place can be more memory efficient**, especially for large data structures, as it avoids the need to allocate new memory for the returned values. This can lead to performance benefits in many scenarios, especially in iterative algorithms where the same operation is performed multiple times. However, it's essential to use this convention judiciously. **Modifying arguments in place can make code harder to understand and debug if not used carefully.**

You are NOT obliged to use modification in place. Here is the same example using simpler programming constructs.

```

1 using NLsolve
2
3 # Define the system of nonlinear equations WITHOUT
4 # using modification in place
5 function system_of_equations(x)
6     return [
7         x[1]^2 + x[2]^2 - 2.0,
8         x[1]^2 - x[2]^2 - 1.0
9     ]
10 end
11
12 # Initial guess for the solution
13 initial_guess = [1.0, 1.0]
14
15 # Solve the system of equations
16 result = nlsolve(system_of_equations, initial_guess)
17
18 # Print the solution
19 println("Roots found at: ", result.zero)
```

Output

```
Roots found at: [1.2247448713915894, 0.7071067811873448]
```

Here is a more impressive example of root finding, where we would NOT want to solve these equations by hand.

```

1 using NLsolve
2 using Random
3 Random.seed!(1958)
4
5 function Fun!(Fun, x)
6     Fun[1] = x[1]^2 + x[2] - x[3] + x[4]^3 - 10
7     Fun[2] = x[2]^2 + x[5] - x[6] + x[7]^2 - 20
8     Fun[3] = x[3]^2 + x[8] - x[9] + x[1]^2 - 30
9     Fun[4] = x[4]^2 + x[5]^3 - x[6] + x[9] - 40
10    Fun[5] = x[7]^2 + x[8]^3 - x[9] + x[2] - 50
11    Fun[6] = x[1] + x[2]^2 - x[3]^3 + x[4] - 5
12    Fun[7] = x[5] + x[6]^2 - x[7]^3 + x[8] - 15
13    Fun[8] = x[9] + x[1]^2 - x[2]^3 + x[3] - 25
14    Fun[9] = x[4] + x[5]^2 - x[6]^3 + x[7] - 35
15 end
16
17 # Trying a different initial guess
18 x0 = randn(9) * 10
19
```

```

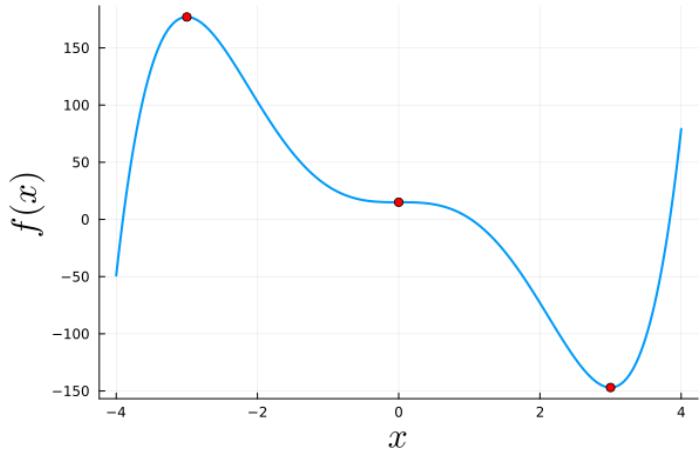
20 # Adjust solver settings
21 result = nlsolve(Fun!, x0, method = :trust_region, ftol = 1e-10,
22 xtol = 1e-10, iterations = 10000)
23
24 println("Roots found at: ", result.zero)
25 println(" ")
26
27 # After obtaining result from nlsolve
28 Fun_evaluated = Vector{Float64}(undef, length(x0))
29 Fun!(Fun_evaluated, result.zero)
30 println("Values of Fun at the root: ", Fun_evaluated)

```

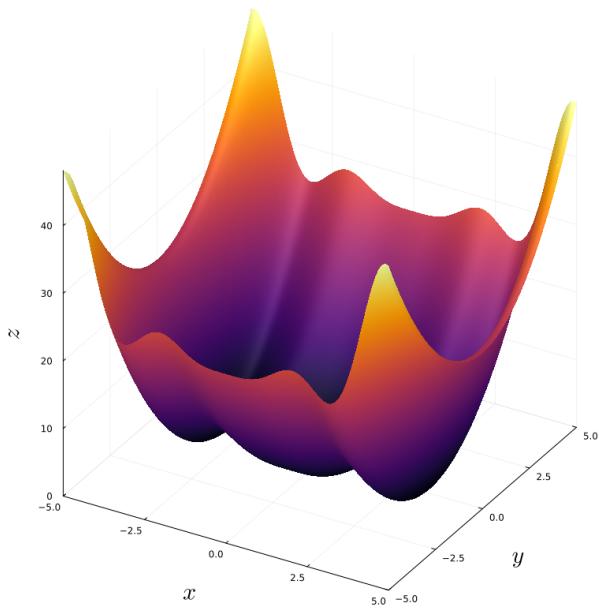
Output

Roots found at: [-7.772786278382927, 4.157370417426024, 0.9040581712865199, -3.7720366206105242]

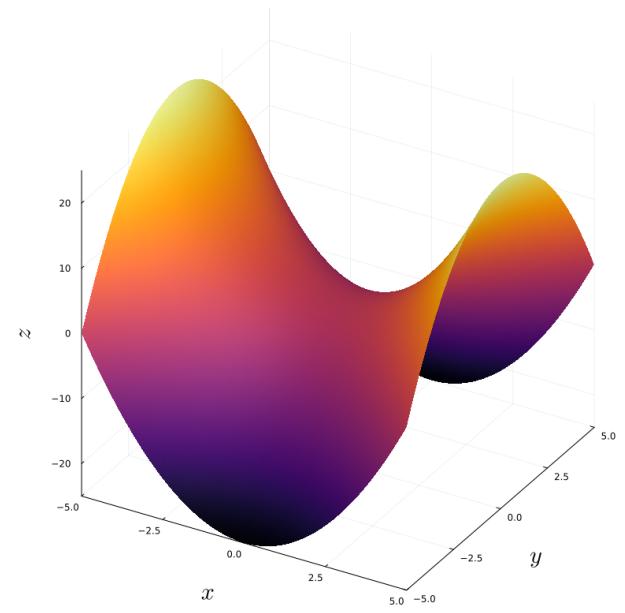
Values of Fun at the root: [7.105427357601002e-15, -3.552713678800501e-15, 3.552713678800501e-15]



(a)



(b)



(c)

Figure 6.3: Three cost functions. (a) $f(x) = x^5 - 15x^3 + 15$ has a local minimum at -3 , an inflection point at the origin, and a local maximum at 3 . (b) $f(x, y) = (x \cdot \sin(x))^2 + y^2$ is riddled with local minima and maxima. (c) $f(x, y) = x^2 - y^2$ has a saddle point at the origin. The name could not be more clear!

6.2.2 Minimization without Constraints

Optimization is the process of finding one or more values $x \in \mathbb{R}^n$ that **minimize a function** $f : \mathbb{R}^n \rightarrow \mathbb{R}$, where the scalar-valued function f is called the **cost function**. The cost function should be designed such that it has a minimum at a point x^* that is of interest to you, it should be “small” for values of x that are “near” x^* , and it should be “large” for values of x that are “far” from your preferred value, x^* . In other symbols, the function should be “bowl-shaped” like $f(x) = (x - x^*)^2$. You may recall from Chapter 12 of ROB 101 *Computational Linear Algebra*, in Machine Learning, a cost function for minimization is also called a **regret function** in the sense that you regret being far from your preferred value, x^* , and your regret is minimum at x^* . We will abuse notation² and write

$$x^* = \arg \min_{x \in \mathbb{R}^n} f(x) \quad (6.10)$$

to denote the value of x achieving the minimum, even when there may be more than one such x .

One can also do the opposite of minimization, which is to **maximize a cost function**. In this case, the Machine Learning community calls the cost function a **reward**, because hey, who does not like to maximize reward and minimize regret! We will stick to minimization throughout this section because

$$\arg \max_{x \in \mathbb{R}^n} f(x) = \arg \min_{x \in \mathbb{R}^n} (-f(x)),$$

that is, by a simple sign change, you convert a maximization problem into a minimization problem. Most of the numerical packages take advantage of this property.

Fundamental Idea and Vocabulary

Because the gradient is the transpose of the Jacobian, the linearization of f via the gradient can be written as

$$f(x) \approx f(x_0) + (\nabla f(x_0))^\top \cdot (x - x_0) = f(x_0) + \nabla f(x_0) \bullet (x - x_0), \quad (6.11)$$

where $u \bullet v$ denotes the **dot product or inner product** of the two n -vectors, u and v . Based on this equation, we make the following observations:

- (a) If the gradient evaluated at a point x_0 is the zero vector, that is, $\nabla f(x_0) = 0_{n \times 1}$, then the linear approximation of f is equal to the constant (vector) $f(x_0)$. This means that in the vicinity of x_0 , irrespective of the direction³ in which we move from x_0 , the function satisfies $f(x) \approx f(x_0)$. Such points x_0 where $\nabla f(x_0) = 0$ are referred to as **extrema** or **stationary points**. The term “extrema” is more common in Calculus courses, whereas “stationary points” seems to be the go-to term for modern optimization packages!

The essential vocabulary for discussing points where $\nabla f(x_0) = 0_{n \times 1}$ is the following:

- $x_0 \in \mathbb{R}^n$ is a **stationary point** if $\nabla f(x_0) = 0_{n \times 1}$.
- a stationary point $x_0 \in \mathbb{R}^n$ is a **local minimum** if $f(x_0) \leq f(x)$ for all x “near x_0 ” (i.e., there exists $\delta > 0$ such that for all $\|x - x_0\| < \delta$, $f(x_0) \leq f(x)$).
- a stationary point $x_0 \in \mathbb{R}^n$ is a **local maximum** if $f(x_0) \geq f(x)$ for all x “near x_0 ” (i.e., there exists $\delta > 0$ such that for all $\|x - x_0\| < \delta$, $f(x_0) \geq f(x)$).
- A stationary point x_0 that is **neither** a local minimum **nor** a local maximum is called a **saddle point**. In the special case of functions $f : \mathbb{R} \rightarrow \mathbb{R}$, saddle points are called **stationary points of inflection**.
- The crowning accomplishment in optimization is to find a **global minimum**, an x_0 that is a local minimum and $f(x_0) \leq f(x)$ for all $x \in \mathbb{R}^n$, i.e., you can remove the “near x_0 .” The same idea applies to a **global maximum**. We will content ourselves with local minima and maxima.

Figure 6.3 illustrates the various types of stationary points.

²Abusing notation means that one is being a bit sloppy with one’s use of symbols, which is what notation is! One often abuses notation when doing the right thing is painful and would cause more of a distraction to explain the good notation than to caution about using poor notation!

³In \mathbb{R}^n , there are n possible directions from x_0 . Each direction can have a positive or negative coefficient.

The subject of **inflection points** can be quite subtle as explained on [Wikipedia](#). We will stick with stationary points of inflection.

(b) If $\nabla f(x_0) \neq 0_{n \times 1}$, then we have the possibility of choosing x near x_0 that causes f to change.

- If $\nabla f(x_0) \bullet (x - x_0) > 0$, then (6.11) implies $f(x) > f(x_0)$, so we are making f increase.
- If $\nabla f(x_0) \bullet (x - x_0) < 0$, then (6.11) implies $f(x) < f(x_0)$, so we are making f decrease.
- If $\nabla f(x_0) \bullet (x - x_0) = 0$, then implies $f(x) \approx f(x_0)$, so f is remaining approximately constant. Therefore **these directions are not fruitful to explore**.

To understand the non-fruitful directions, let's see how well you recall null spaces! When $\nabla f(x_0) \neq 0_{n \times 1}$, it has rank one, and its null space has dimension⁴ $(n - 1)$. Moreover, all vectors $v \in \text{null}(\nabla f(x_0))$ satisfy

$$\nabla f(x_0) \bullet v = 0 \quad (\text{same as writing } \nabla f(x_0) \perp v).$$

That leaves us one dimension that is orthogonal to $\text{null}(\nabla f(x_0))$ and it is spanned by $\nabla f(x_0)$, the gradient of f !

(c) Moving in the direction given by $-\nabla f(x_0)$ gives rise to **gradient descent**, because

$$(x - x_0) = -\nabla f(x_0) \implies \nabla f(x_0) \bullet (x - x_0) = -\nabla f(x_0) \bullet \nabla f(x_0) = -\|\nabla f(x_0)\|^2,$$

and hence, by (6.11),

$$f(x) \approx f(x_0) - \|\nabla f(x_0)\|^2 < f(x_0).$$

When we move in the opposite direction,

$$(x - x_0) = \nabla f(x_0) \implies \nabla f(x_0) \bullet (x - x_0) = \nabla f(x_0) \bullet \nabla f(x_0) = \|\nabla f(x_0)\|^2,$$

and hence

$$f(x) \approx f(x_0) + \|\nabla f(x_0)\|^2 > f(x_0).$$

This is called **gradient ascent**.

6.2.3 Gradient Descent Algorithm

Method 6.8: Gradient Descent

We assume that we are given an initial vector $x_0 \in \mathbb{R}^n$. Then the update law,

$$x_{k+1} = x_k - s(\nabla f(x_k)), \tag{6.12}$$

where $s > 0$ is called the **step size**, moves x in the direction of a local minimum.

Warning: If $s > 0$ is selected too large, then (6.12) may not converge. Hence, the sophisticated implementations found in today's software packages adaptively adjust the step size. In fact, they often perform a **line search**, meaning, they evaluate (6.12) for multiple values of $s \in \{s_1, s_2, \dots, s_n\}$, and choose the value of s_i resulting in the smallest value of

$$f(x_{k+1}) = f(x_k - s_i [\nabla f(x_k)]).$$

Here is a short video, [Backtracking Line Search in Gradient Descent](#) by Deep Shallownet, discussing this process (the computer voice employed is so appropriate).

You can find in Chapter 12 of ROB 101 *Computational Linear Algebra* numerous examples with a hand-made implementation of the Gradient Descent Algorithm. Instead, here, the focus is on the Julia Package [JuMP](#).

⁴By the Rank-Nullity Theorem, for the $n \times 1$ matrix $\nabla f(x_0)$, we have $\text{rank}(\nabla f(x_0)) + \text{nullity}(\nabla f(x_0)) = n$.

Below, we give a detailed example of how to use JuMP. If the explanations prove inadequate, you can use the following prompt¹ for ChatGPT: Hello ChatGPT! I've recently started learning about optimization in Julia, and I've come across the JuMP package and the Ipopt solver. I'm a bit overwhelmed with how to set up and solve optimization problems using these tools. Could you provide a step-by-step explanation on how to use JuMP with Ipopt, including code examples along the way? At the end, I'd appreciate a complete sample program in Julia that demonstrates the entire process. Thank you!

¹ This was generated by ChatGPT itself in response to: "Can you give me a good prompt that a student can use for having you explain to them how to use JuMP and Ipopt, include code examples along the way, and then output a complete sample program in Julia?" **Author: This is getting very meta!**

Example 6.9. (Minimization of the Rosenbrock Function:) The Rosenbrock function, often used as a performance test for optimization algorithms, is defined as:

$$f(x, y) = (1 - x)^2 + 100(y - x^2)^2;$$

note that the term $(y - x^2)$ makes $f(x, y)$ not a simple, bowl-shaped quadratic that is familiar to us from minimum norm problems in Linear Algebra. The function appears as a narrow, curved valley, and its global minimum is inside a long, narrow, parabolic-shaped flat valley. Finding the valley is trivial, but converging to the global minimum within the valley is difficult due to the flatness of the valley, making it a challenging test for optimization algorithms.

The term $(1 - x)^2$ vanishes at $x = 1$ and the term $(y - x^2)$ vanishes at $y = x^2 = 1$. The function, therefore, has a global minimum of zero at:

$$(x, y) = (1, 1).$$

Solution:

```

1 using JuMP, Ipopt
2
3 # Initialize a new optimization model
4 # We're using the Ipopt solver, which is suitable for nonlinear optimization problems
5 model = Model(Ipopt.Optimizer)
6
7 # Define decision variables for the optimization problem
8 # We're defining two variables, x and y, both bounded between -20 and 20
9 # We also provide an initial guess for the solver to start from. This can influence the
   solution found.
10 @variable(model, -20 <= x <= 20, start = -20.0)    # Initial guess set to -20 for x
11 @variable(model, -20 <= y <= 20, start = 20.0)    # Initial guess set to 20 for y
12
13 # Define the objective function we want to minimize
14 # Here, we're using the Rosenbrock function, a common test problem for optimization
   algorithms
15 @NLobjective(model, Min, (1 - x)^2 + 100 * (y - x^2)^2)
16
17 # Adjust solver settings to improve performance and accuracy
18 # "max_iter" sets the maximum number of iterations the solver will take
19 # "tol" sets the convergence tolerance. The solver stops when improvements are smaller
   than this value.
20 set_optimizer_attribute(model, "max_iter", 1000)    # Allow up to 1000 iterations
21 set_optimizer_attribute(model, "tol", 1e-9)        # Require a very accurate solution
22
23 # Start the optimization process
24 # The solver will try to find values of x and y that minimize the Rosenbrock function
25 optimize!(model)
26
27 # Once optimization is complete, display the results
28 # "value()" retrieves the optimal value of a decision variable

```

```

29 # "objective_value()" retrieves the minimum value of the objective function
30 println("Optimal parameters (x, y): (" , value(x) , " , " , value(y) , " )")
31 println("Minimum objective value: " , objective_value(model))

```

Summary of the Output

EXIT: Optimal Solution Found.
Optimal parameters (x, y): (0.9999999792197696, 0.9999999583699614)
Minimum objective value: 4.3230208956579026e-16

Complete Output

This is Ipopt version 3.14.13, running with linear solver MUMPS 5.6.1.

Number of nonzeros in equality constraint Jacobian....:	0
Number of nonzeros in inequality constraint Jacobian.:	0
Number of nonzeros in Lagrangian Hessian.....:	3
Total number of variables.....:	2
variables with only lower bounds:	0
variables with lower and upper bounds:	2
variables with only upper bounds:	0
Total number of equality constraints.....:	0
Total number of inequality constraints.....:	0
inequality constraints with only lower bounds:	0
inequality constraints with lower and upper bounds:	0
inequality constraints with only upper bounds:	0
iter objective inf_pr inf_du lg(mu) d lg(rg) alpha_du alpha_pr ls	
0 1.3856695e+07 0.00e+00 9.70e+01 -1.0 0.00e+00 - 0.00e+00 0.00e+00 0	
1 4.1877282e+06 0.00e+00 3.94e+01 -1.0 4.81e+00 - 4.03e-02 1.00e+00f 1	
2 7.4626085e+05 0.00e+00 1.16e+01 -1.0 4.68e+00 - 9.94e-01 1.00e+00f 1	
3 1.1541790e+05 0.00e+00 3.26e+00 -1.0 2.98e+00 - 1.00e+00 1.00e+00f 1	
4 1.3175100e+04 0.00e+00 8.37e-01 -1.0 1.79e+00 - 1.00e+00 1.00e+00f 1	
5 7.1409431e+02 0.00e+00 1.60e-01 -1.7 8.86e-01 - 1.00e+00 1.00e+00f 1	
6 3.4096513e+01 0.00e+00 1.36e-02 -2.5 2.80e-01 - 1.00e+00 1.00e+00f 1	
7 2.8367005e+01 0.00e+00 2.56e-04 -3.8 1.63e-01 - 9.81e-01 1.00e+00f 1	
8 2.6050947e+01 0.00e+00 4.41e-03 -5.7 4.88e+01 - 3.33e-02 4.96e-02f 5	
9 2.3356648e+01 0.00e+00 5.03e-03 -5.7 2.43e+00 - 9.91e-01 1.00e+00f 1	
iter objective inf_pr inf_du lg(mu) d lg(rg) alpha_du alpha_pr ls	
10 2.0568742e+01 0.00e+00 2.84e-03 -5.7 1.66e+00 - 1.00e+00 1.00e+00f 1	
11 1.9006996e+01 0.00e+00 6.06e-03 -5.7 2.57e+00 - 1.00e+00 1.00e+00f 1	
12 1.5917375e+01 0.00e+00 1.05e-03 -5.7 7.47e-01 - 1.00e+00 1.00e+00f 1	
13 1.5661315e+01 0.00e+00 5.92e-03 -5.7 4.71e+00 - 1.00e+00 5.00e-01f 2	
14 1.2150777e+01 0.00e+00 5.87e-04 -5.7 3.75e-01 - 1.00e+00 1.00e+00f 1	
15 1.0995766e+01 0.00e+00 2.62e-03 -5.7 5.42e+00 - 1.00e+00 2.50e-01f 3	
16 9.2807587e+00 0.00e+00 1.09e-03 -5.7 7.27e-01 - 1.00e+00 1.00e+00f 1	
17 8.3235430e+00 0.00e+00 1.52e-03 -5.7 1.56e+00 - 1.00e+00 5.00e-01f 2	
18 7.0460250e+00 0.00e+00 1.33e-03 -5.7 8.03e-01 - 1.00e+00 1.00e+00f 1	
19 5.8858703e+00 0.00e+00 1.01e-03 -5.7 6.37e-01 - 1.00e+00 1.00e+00f 1	
iter objective inf_pr inf_du lg(mu) d lg(rg) alpha_du alpha_pr ls	
20 4.9033128e+00 0.00e+00 9.38e-04 -5.7 5.86e-01 - 1.00e+00 1.00e+00f 1	
21 3.9797479e+00 0.00e+00 5.72e-04 -5.7 3.77e-01 - 1.00e+00 1.00e+00f 1	
22 3.3193263e+00 0.00e+00 6.86e-04 -5.7 4.40e-01 - 1.00e+00 1.00e+00f 1	
23 2.5481245e+00 0.00e+00 2.37e-04 -5.7 1.31e-01 - 1.00e+00 1.00e+00f 1	
24 2.1475610e+00 0.00e+00 3.11e-04 -5.7 4.04e-01 - 1.00e+00 5.00e-01f 2	

```

25 1.6279972e+00 0.00e+00 1.61e-04 -5.7 1.54e-01 - 1.00e+00 1.00e+00f 1
26 1.3005814e+00 0.00e+00 3.15e-04 -5.7 2.19e-01 - 1.00e+00 1.00e+00f 1
27 8.8740291e-01 0.00e+00 6.30e-05 -5.7 9.78e-02 - 1.00e+00 1.00e+00f 1
28 6.9660011e-01 0.00e+00 2.02e-04 -5.7 3.22e-01 - 1.00e+00 5.00e-01f 2
29 4.5955435e-01 0.00e+00 7.80e-05 -5.7 1.09e-01 - 1.00e+00 1.00e+00f 1
iter objective inf_pr inf_du lg(mu) ||d|| lg(rg) alpha_du alpha_pr ls
30 3.7378503e-01 0.00e+00 2.58e-04 -5.7 1.98e-01 - 1.00e+00 1.00e+00f 1
31 1.7419935e-01 0.00e+00 1.86e-05 -5.7 9.55e-02 - 1.00e+00 1.00e+00f 1
32 1.1678685e-01 0.00e+00 1.61e-04 -5.7 3.13e-01 - 1.00e+00 5.00e-01f 2
33 5.1675933e-02 0.00e+00 2.27e-05 -5.7 1.03e-01 - 1.00e+00 1.00e+00f 1
34 3.9425897e-02 0.00e+00 2.04e-04 -5.7 2.09e-01 - 1.00e+00 1.00e+00f 1
35 5.1649798e-03 0.00e+00 2.71e-06 -5.7 5.44e-02 - 1.00e+00 1.00e+00f 1
36 1.9553865e-03 0.00e+00 5.70e-05 -8.6 1.23e-01 - 9.95e-01 1.00e+00f 1
37 6.6464788e-06 0.00e+00 5.67e-08 -8.6 1.02e-02 - 1.00e+00 1.00e+00f 1
38 4.4011255e-09 0.00e+00 8.67e-08 -8.6 5.14e-03 - 1.00e+00 1.00e+00f 1
39 3.3519648e-13 0.00e+00 1.29e-13 -8.6 1.55e-05 - 1.00e+00 1.00e+00f 1
iter objective inf_pr inf_du lg(mu) ||d|| lg(rg) alpha_du alpha_pr ls
40 4.3230209e-16 0.00e+00 4.08e-15 -10.0 1.12e-06 - 1.00e+00 1.00e+00f 1

```

Number of Iterations....: 40

	(scaled)	(unscaled)
Objective.....:	1.4220267008343644e-20	4.3230208956579026e-16
Dual infeasibility.....:	4.0802855406541762e-15	1.2404239415581401e-10
Constraint violation....:	0.0000000000000000e+00	0.0000000000000000e+00
Variable bound violation:	0.0000000000000000e+00	0.0000000000000000e+00
Complementarity.....:	9.0909232958894218e-11	2.7636788638282267e-06
Overall NLP error.....:	9.0909232958894218e-11	2.7636788638282267e-06

Number of objective function evaluations	= 80
Number of objective gradient evaluations	= 41
Number of equality constraint evaluations	= 0
Number of inequality constraint evaluations	= 0
Number of equality constraint Jacobian evaluations	= 0
Number of inequality constraint Jacobian evaluations	= 0
Number of Lagrangian Hessian evaluations	= 40
Total seconds in IPOPT	= 0.032

EXIT: Optimal Solution Found.

Optimal parameters (x, y): (0.9999999792197696, 0.9999999583699614)

Minimum objective value: 4.3230208956579026e-16

JuMP's native output is quite verbose. You can suppress it entirely with the following parameter setting:

```

1 # Suppress all of the output
2 set_optimizer_attribute(model, "print_level", 0)
3
4 # Solve the optimization problem
5 optimize!(model)

```

Output

Yes, nothing comes out. You can play with the `print_level` setting to find the level of output that you want. We like

```
set_optimizer_attribute(model, "print_level", 4).
```

Method 6.10: Additional Explanation on Using JuMP and Ipopt for Optimization

JuMP is a modeling language in Julia for mathematical optimization. It provides a user-friendly interface to define optimization problems and connect with various solvers. Among the powerful solvers that can be used with JuMP is **Ipopt**, which stands for Interior Point OPTimizer. It's particularly suited for nonlinear optimization problems.

(a) **Setting up the Problem:** To start, you need to define your optimization problem. This involves:

- Specifying the decision variables (i.e., the variables that can be varied in the optimization, such as $x = (x_1, x_2, \dots, x_n)$).
- Setting the objective function (i.e., the cost function).
- Adding any constraints.

(b) **Decision Variables:** In JuMP, you can define decision variables using the `@variable` macro. For example:

```
@variable(model, -2 <= x <= 2)
```

defines a variable x that is bounded between -2 and 2. If you want no bounds on x , that is, $-\infty < x < \infty$, then you use

```
@variable(model, x)
```

and if you want $x \geq 0$, you use

```
@variable(model, x >= 0)
```

That's fairly intuitive.

(c) **Objective Function:** The objective function is what you aim to minimize (or maximize). In JuMP, you can set the objective function using the `@objective` macro. For nonlinear objectives, use `@NLobjective`. For instance:

```
@NLobjective(model, Min, (1 - x)^2 + 100 * (y - x^2)^2)
```

This sets the Rosenbrock function as the objective (aka, cost function) to be minimized.

(d) **Constraints:** In Chapter 6.3, we will deal with minimization subject to constraints, such as $g(x) \leq c$ or $g(x) = c$. Constraints restrict the feasible region of the optimization problem. In JuMP, you can add constraints using the `@constraint` macro. For nonlinear constraints, use `@NLconstraint`. For example,

```
@NLconstraint(model, x^2 + y^2 <= 1)
```

ensures that the solution (x, y) lies inside a unit circle. **We are not using this option yet, but, later, we'll not repeat the long explanation of how JuMP works. Ergo, you get it here and now!**

(e) **Solving the Problem:** Once the problem is defined, you can solve it using an optimization solver. Here's how to set up and solve the problem using the Ipopt solver:

```

# Create a model with Ipopt as the solver
model = Model(Ipopt.Optimizer)

# Define variables, objective, and constraints...

# Solve the problem
optimize!(model)

```

(f) **Retrieving the Solution:** After solving, you can retrieve the optimal solution and objective value using:

```

optimal_x = value(x)
optimal_y = value(y)
min_value = objective_value(model)

```

Tips for Success:

- **Initial Guess:** The starting point can influence the solution. Use the `start` attribute to set an initial guess, e.g., `@variable(model, -2 <= x <= 2, start = -1.2)`.
- **Solver Options:** Adjust solver options for better convergence. For instance, increase the maximum number of iterations or set a tighter convergence tolerance.
- **Scaling:** Ensure that the **problem's scale is appropriate**. If variables or constraints have vastly different magnitudes, consider rescaling.
- If your code has a bug in it, paste the full code into ChatGPT and ask it to find and fix the error.

Bottom Line: JuMP, combined with the Ipopt solver, offers a powerful and flexible environment for solving optimization problems in Julia. With its intuitive syntax and robust capabilities, engineers can tackle a wide range of optimization challenges.

Warning: There is no guarantee that a GLOBAL MINIMUM will be found with any numerical optimizer. Before adjusting the parameters in JuMP, via

```

1 set_optimizer_attribute(model, "max_iter", 1000)    # Allow up to 1000 iterations
2 set_optimizer_attribute(model, "tol", 1e-9)          # Require a very accurate solution

```

it reported the “optimal solution” to the Rosenbrock Problem as $(0.7864151510570021, 0.6176983075999913)$, even though the true answer is $(1.0, 1.0)$.

```

1 using JuMP, Ipopt
2
3 # Create a new model with Ipopt as the solver
4 model = Model(Ipopt.Optimizer)
5
6 # Define variables
7 @variable(model, -2 <= x <= 2, start = -1.2)    # Set initial guess
8 @variable(model, -2 <= y <= 2, start = 1.0)        # Set initial guess
9
10 # Define the Rosenbrock function
11 @NLobjective(model, Min, (1 - x)^2 + 100 * (y - x^2)^2)
12
13 # Solve the optimization problem
14 optimize!(model)
15

```

```

16 # Display results
17 println("Optimal parameters (x, y): (" + value(x) + ", " + value(y) + ")")
18 println("Minimum objective value: " + objective_value(model))

```

Abridged Output

EXIT: Optimal Solution Found.
Optimal parameters (x, y): (0.7864151510570021, 0.6176983075999913)
Minimum objective value: 0.04567481005305849

A recommended example: Gradient descent, how neural networks learn by **3Blue1Brown**.

6.2.4 Second Derivative Tests for Local Min and Max

We cover both single-variable and multivariable problems.

Second-Order Derivative Test for a Scalar Variable

Proposition 6.11: Scalar Variable Second-order Test for Local Max and Min

Given a function $f : \mathbb{R} \rightarrow \mathbb{R}$, with a stationary point at $x = x^*$, that is, its first derivative $\frac{df(x^*)}{dx} = 0$, the nature of the stationary point can be determined using the second derivative of f at x^* .

- If $\frac{d^2f(x^*)}{dx^2} > 0$, then f has a local minimum at $x = x^*$.
- If $\frac{d^2f(x^*)}{dx^2} < 0$, then f has a local maximum at $x = x^*$.
- If $\frac{d^2f(x^*)}{dx^2} = 0$, the test is inconclusive, and f may have a stationary point of inflection at $x = x^*$.

Note: To remember these, take the second derivatives of $f(x) = x^2$, which has a minimum at the origin, $f(x) = -x^2$, which has a maximum at the origin, $f(x) = x^3$, which has an inflection point at the origin, and $f(x) = x^4$, which has a minimum at the origin, but the derivative test is inconclusive because $\frac{d^2f(0)}{dx^2} = 12x^2|_{x^*=0} = 0$. Alternatively, for some of you, Example 6.14, which shows how Taylor expansions can be used to establish the result, is what you need to recall the second-derivative test.

Example 6.12. $f(x) = x^2 - x^3 + x^4$ has a stationary point at the origin. Does the second derivative test confirm it to be a local min, a local max, or is it inconclusive?

Solution: Ans. Local minimum because $\frac{d^2f(x)}{dx^2} = 2 - 6x + 12x^2$ and thus, $\frac{d^2f(0)}{dx^2} = 2 > 0$. ■

Example 6.13. Use Jump to find a local minimum of each of the following functions,

$$(a) f(x) = x^4 \sin(1 + x^2)$$

$$(b) f(x) = e^x \sin(1 + x^2),$$

and then verify or refute the results using the second-derivative test.

Solutions:

- (a) Jump reports a local minimum for $f(x) = x^4 \sin(1 + x^2)$ at $x^* = 0.0$ but the second-derivative test is inconclusive

- (b) Jump reports a local minimum for $f(x) = e^x \sin(1 + x^2)$ at $x^* = 0.0$, which is confirmed by the second-derivative test.

The second-derivative test is applied using FiniteDiff.

```

1 using JuMP, Ipopt, FiniteDiff
2
3 # Define the function
4 if false
5     f(x) = x^4 * sin(1 + x^2)
6 else
7     f(x) = exp(x) * sin(1 + x^2)
8 end
9
10 # Create a new model with Ipopt as the solver
11 model = Model(Ipopt.Optimizer)
12
13 # Manually register the function with the model to use automatic differentiation
14 register(model, :f, 1, f; autodiff = true)
15
16 # Define variables
17 @variable(model, x)
18
19 # Define the cost function using the registered function
20 @NLobjective(model, Min, f(x))
21
22 # Fine-tune a few parameters
23 set_optimizer_attribute(model, "max_iter", 1000) # Allow up to 1000 iterations
24 set_optimizer_attribute(model, "tol", 1e-9) # Require a very accurate solution
25
26 # Suppress much of the output
27 set_optimizer_attribute(model, "print_level", 0)
28
29 # Solve the optimization problem
30 optimize!(model)
31
32 # Minimizing value
33 xStar = value(x)
34
35 # Display results
36 println("\n Optimal parameters x: ", xStar)
37 println("Minimum objective value: ", objective_value(model))
38
39 # Define a method to compute the first derivative
40 function df(x)
41     return FiniteDiff.finite_difference_derivative(f, x)
42 end
43
44 # Compute the second derivative using finite_difference_derivative on the first
# derivative function
45 second_derivative = FiniteDiff.finite_difference_derivative(df, xStar)
46
47 println("The second derivative of f(x) at x = $xStar is: ", second_derivative[1, 1])

```

Output for $f(x) = x^4 \sin(1 + x^2)$

Optimal parameters x: 0.0

Minimum objective value: 0.0

The second derivative of $f(x)$ at $x = 0.0$ is: 2.4684402317159913e-10

Output $f(x) = e^x \sin(1 + x^2)$

Optimal parameters x : -1.8572587077723928

Minimum objective value: -0.15073321671888845

The second derivative of $f(x)$ at $x = -1.8572587077723928$ is: 2.149336845346529

Example 6.14. Use a second-order Taylor expansion to prove (or at least, plausibly argue the veracity of) Prop. 6.11.

Solution: In the prime notation for derivatives, we are given that $f'(x^*) = 0$. Hence, the second order Taylor expansion about x^* is

$$f(x) = f(x^*) + f''(x^*) \cdot (x - x^*)^2.$$

If $f''(x^*) > 0$, then f is locally an upward pointing parabola about x^* , and thus x^* is a local minimum. If $f''(x^*) < 0$, then f is locally a downward pointing parabola about x^* , and thus x^* is a local maximum. If $f''(x^*) = 0$, then higher-order terms in the Taylor expansion would be needed to reach any conclusion. ■

(Optional Read:) Second-Order Derivative Test for a Multivariable Function

For multivariable functions, stationary points are the roots of the gradient. To check if a stationary point is a local max, min, or neither, we must use the Hessian from Def. 5.46.

Proposition 6.15: Multivariable Second-order Test for Local Max and Min

Consider a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ that is twice continuously differentiable. Let $\mathbf{x}^* = (x_1^*, x_2^*, \dots, x_n^*)$ be a stationary point of f , i.e., the gradient $\nabla f(\mathbf{x}^*) = \mathbf{0}$.

To determine whether \mathbf{x}^* is a local minimum, local maximum, or neither, we examine the Hessian matrix of f at \mathbf{x}^* , denoted by $H_f(\mathbf{x}^*)$.

- If $H_f(\mathbf{x}^*)$ is **positive definite**, then f has a local minimum at \mathbf{x}^* .
- If $H_f(\mathbf{x}^*)$ is **negative definite**, then f has a local maximum at \mathbf{x}^* .
- If $H_f(\mathbf{x}^*)$ is **indefinite** then \mathbf{x}^* is a saddle point, and f does not have a local minimum or maximum at \mathbf{x}^* .

To determine if the Hessian matrix is positive definite, negative definite, or indefinite, one can examine its eigenvalues:

- All eigenvalues of $H_f(\mathbf{x}^*)$ are positive $\implies H_f(\mathbf{x}^*)$ is positive definite.
- All eigenvalues of $H_f(\mathbf{x}^*)$ are negative $\implies H_f(\mathbf{x}^*)$ is negative definite.
- Eigenvalues of $H_f(\mathbf{x}^*)$ are of mixed signs or one is zero $\implies H_f(\mathbf{x}^*)$ is indefinite.

We give a fully numerical example for you to follow.

Instructions for ChatGPT4: This prompt guides ChatGPT in generating a Julia solution that not only addresses an interesting optimization problem but also includes educational documentation throughout the code. It ensures that the solution will be accessible and informative for students, providing a valuable learning resource on numerical optimization, the use of specific Julia

packages, and the mathematical analysis of optimization problems.

“Create a comprehensive Julia solution to analyze the Himmelblau function. The solution should include the following components, each thoroughly documented with comments in the code:

- **Function Definition:** Implement the Himmelblau function $f(x, y) = (x^2 + y - 11)^2 + (x + y^2 - 7)^2$ in Julia. Provide a brief comment on what the Himmelblau function represents in optimization theory.
- **Optimization Setup with JuMP and Ipopt:** Use the JuMP package to define and solve the optimization problem for finding critical points of the Himmelblau function. Employ Ipopt as the optimizer. Document how to set up the optimization model, define variables x and y , and specify the objective function within JuMP.
- **Moreover:** Initialize the optimization from multiple starting points: (3, 2), (-2.805118, 3.283186), (-3.779310, -3.283186), and (3.584428, -1.848126). Explain the choice of these points and their significance in exploring the function’s landscape.
- **Hessian Computation with FiniteDiff:** After identifying a critical point, use the FiniteDiff.jl package to compute the Hessian matrix at that point. Include comments explaining the purpose of computing the Hessian and how it’s done using FiniteDiff.
- **Eigenvalue Analysis with LinearAlgebra:** Calculate the eigenvalues of the Hessian matrix using the LinearAlgebra package. Document the process of eigenvalue calculation and its importance in classifying critical points.
- **Critical Point Classification:** Based on the eigenvalues of the Hessian, classify each critical point as a local minimum, maximum, or saddle point. Include a detailed explanation of how the classification is determined from the eigenvalues.

Ensure the Julia code is clean, well-organized, and fully documented, with comments explaining each step of the process and the significance of each operation. The documentation should be detailed enough for someone unfamiliar with the process to understand the methodology, purpose, and implementation details.”

You will obtain code that is more thorough than the example below!

```
1 using JuMP, Ipopt, FiniteDiff, LinearAlgebra
2
3 # Himmelblau function definition remains the same
4 function himmelblau(x, y)
5     return (x^2 + y - 11)^2 + (x + y^2 - 7)^2
6 end
7
8 # Extended function to also compute Hessian and its eigenvalues
9 function optimize_and_analyze_himmelblau(x_init, y_init)
10    model = Model(Ipopt.Optimizer)
11
12    @variable(model, x, start = x_init)
13    @variable(model, y, start = y_init)
14
15    @NLobjective(model, Min, (x^2 + y - 11)^2 + (x + y^2 - 7)^2)
16
17    set_optimizer_attribute(model, "max_iter", 1000)
18    set_optimizer_attribute(model, "tol", 1e-9)
19    set_optimizer_attribute(model, "print_level", 0)
20
21    optimize!(model)
22
23    xStar, yStar = value(x), value(y)
24    println("\nOptimal parameters (x, y) : (", xStar, ", ", yStar, ")")
25    println("Minimum objective value : ", objective_value(model))
```

```

26
27 # Define a wrapper function for FiniteDiff
28 function himmelblau_wrapper(v)
29     return himmelblau(v[1], v[2])
30 end
31
32 # Compute the Hessian at the optimal point
33 hessian = FiniteDiff.finite_difference_hessian(himmelblau_wrapper, [xStar, yStar])
34
35 # Compute the eigenvalues of the Hessian
36 eigenvalues = eigvals(hessian)
37
38 println("Hessian at the optimal point: \n", hessian)
39 println("Eigenvalues of the Hessian: ", eigenvalues)
40
41 # Determine the nature of the critical point based on the eigenvalues
42 if all(eig .-> eig > 0, eigenvalues)
43     println("The point is a local minimum.")
44 elseif all(eig .-> eig < 0, eigenvalues)
45     println("The point is a local maximum.")
46 else
47     println("The point is a saddle point or the test is inconclusive.")
48 end
49 end
50
51 # Starting points near the four local minima
52 starting_points = [(3, 2), (-2.805118, 3.283186), (-3.779310, -3.283186), (3.584428,
53 -1.848126)]
54
55 # Optimize Himmelblau's function and analyze each result from the starting points
56 for (x_init, y_init) in starting_points
57     println("Starting from point: ($x_init, $y_init)")
58     optimize_and_analyze_himmelblau(x_init, y_init)
59 end

```

Output

```

Starting from point: (3, 2)

Optimal parameters (x, y): (3.0, 2.0)
Minimum objective value: 0.0
Hessian at the optimal point:
[74.0000002682209 20.0; 20.0 34.00000011920929]
Eigenvalues of the Hessian: [25.715728893569633, 82.28427149386056]
The point is a local minimum.

Starting from point: (-2.805118, 3.283186)

Optimal parameters (x, y): (-2.805118086952745, 3.131312518250573)
Minimum objective value: 7.888609052210118e-31
Hessian at the optimal point:
[64.949500088456 1.3047777251870691; 1.3047777251870691 80.44094498775873]
Eigenvalues of the Hessian: [64.84037300554613, 80.5500720706686]
The point is a local minimum.

Starting from point: (-3.77931, -3.283186)

```

```

Optimal parameters (x, y): (-3.779310253377774, -3.2831859912861767)
Minimum objective value: 3.996369345849646e-26
Hessian at the optimal point:
[116.26548835580718 -28.24998497866954; -28.24998497866954 88.23448234835634]
Eigenvalues of the Hessian: [70.71435509450659, 133.78561560965693]
The point is a local minimum.

```

Starting from point: (3.584428, -1.848126)

```

Optimal parameters (x, y): (3.584428340330554, -1.8481265269646447)
Minimum objective value: 9.516628633969033e-25
Hessian at the optimal point:
[104.78501259863707 6.945207253476179; 6.945207253476179 29.32457337913531]
Eigenvalues of the Hessian: [28.690677260855026, 105.41890871691734]
The point is a local minimum.

```

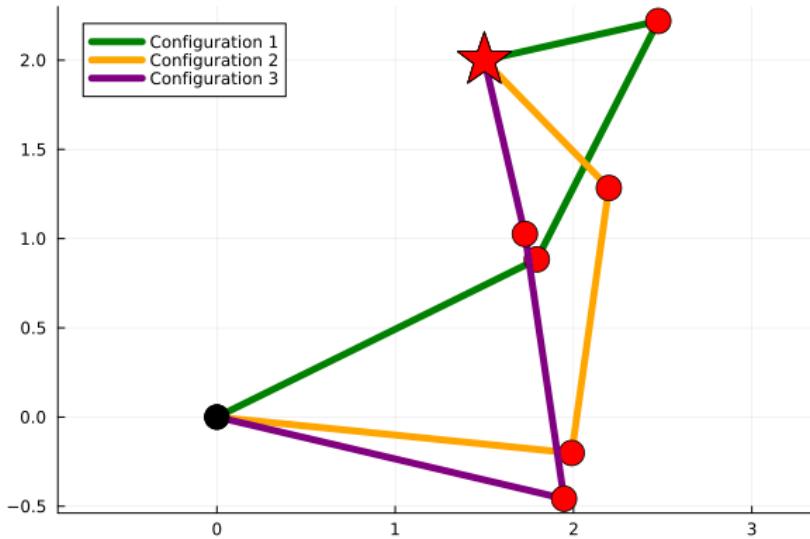


Figure 6.4: This is the same 3-link manipulator from Fig. 1.6-(b). Here, the red star (\star) indicates the desired position of the end of the last link, meaning it is where we are trying to position the robot's **end effector**. The precise point in Cartesian space is $p_3 = (1.5, 2)$. We see that there are numerous configurations (i.e., different values for θ_i) of the robot that achieve the same position of the end effector (we only show three of an infinite number of possible configurations). Hence the problem is **underdetermined**. What is an aspiring Roboticist to do with this messy situation? **The solution is constrained minimization.**

Example 6.16. Sketch out a plausible argument for Proposition 6.15.

Solution: We use facts about symmetric matrices that are stated in the Appendix of the ROB 101 textbook, *Computational Linear Algebra*, and proved in the ROB 501 textbook, *Mathematics for Robotics*; the textbooks can be found [here](#).

Let $x^* \in \mathbb{R}^n$ be a critical point of $f(x)$, meaning that $\nabla f(x^*) = 0$. To determine the nature of stationary point, we will analyze the behavior of the second-order Taylor expansion of $f(x)$ around x^* .

The Hessian matrix of f at x^* , denoted $H_f(x^*)$, is symmetric. Therefore, the eigenvalues of $H_f(x^*)$, denoted $\{\lambda_1, \lambda_2, \dots, \lambda_n\}$, are real, and $H_f(x^*)$ has a full set of linearly independent eigenvectors, $\{v_1, \dots, v_n\}$. We can select these eigenvectors such that the matrix

$$Q := [v_1 \ v_2 \ \dots \ v_n]$$

is orthogonal, i.e., $Q^\top \cdot Q = I_n$.

Step 1: Change of Variables

We perform a change of variables by introducing the new coordinate system $y = Q^\top(x - x^*)$. In this coordinate system, we will examine the Hessian and show that it is diagonal. Since Q is orthogonal, this transformation preserves distances, so the nature of the critical point does not change.

Step 2: Taylor Expansion

The second-order Taylor expansion of $f(x)$ around x^* is:

$$f(x) \approx f(x^*) + \frac{1}{2}(x - x^*)^\top H_f(x^*)(x - x^*).$$

Because $\nabla f(x^*) = 0$, there is no first-order term. Now, using the change of variables $y = Q^\top(x - x^*)$, we express the quadratic form in terms of y :

$$f(x) \approx f(x^*) + \frac{1}{2}y^\top Q^\top H_f(x^*) Q y$$

Step 3: Diagonalizing the Hessian

Because Q is the matrix of eigenvectors of $H_f(x^*)$, we know that:

$$Q^\top H_f(x^*) Q = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$$

This is a diagonal matrix with the eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$ of $H_f(x^*)$ on the diagonal. Therefore, in the new coordinates y , the second-order term becomes:

$$f(x) \approx f(x^*) + \frac{1}{2} \sum_{i=1}^n \lambda_i y_i^2$$

Step 4: Analyzing the Eigenvalues

- If all the eigenvalues $\{\lambda_1, \dots, \lambda_n\}$ of the Hessian matrix are positive, then the quadratic form $\frac{1}{2} \sum_{i=1}^n \lambda_i y_i^2$ is strictly positive for all non-zero y . This implies that the function $f(x)$ is locally bowl in a neighborhood around x^* , meaning x^* is a local minimum of $f(x)$. Similarly, if all of the eigenvalues are negative, we have the analogy of a downward pointing parabola, and hence x^* is a local maximum.
- When the eigenvalues have mixed signs, such as $y_1^2 - y_2^2$ we have the classic saddle point shown in Fig. 6.3.
- When one of the eigenvalues is zero, then the test is inconclusive, just as in the scalar case.

Step 5: Chain Rule Explanation

In this proof, we used the chain rule implicitly during the change of variables. When we transformed the variable from x to $y = Q^\top(x - x^*)$, the Hessian transforms according to the rule:

$$H_y = Q^\top H_f(x^*) Q$$

This is a standard result in Math 215 Multivariable Calculus, showing how the Hessian matrix transforms under a change of coordinates, which diagonalizes it when expressed in the eigenvector basis.

■

6.3 Constrained Optimization

In this section, we'll first motivate the problem being addressed via a physical example. Then we cover the typical material that you would learn in Calc III for use in an engineering course. You cannot do much with it in real life, except get good grades! Hence, we'll also show you how to set up problems for that next internship and the eventual soul-satisfying, technically demanding, well-paying career.

6.3.1 Motivating Problems

As motivation, we pose two example problems based on Fig. 6.4. Recall that the end effector position of the robot as a function of $\theta = (\theta_1, \theta_2, \theta_3)$ was given in (1.17). We'll first take just the x -component of the end effector and define the function

$$g_1(\theta) := L1 \cdot \cos(\theta_1) + L2 \cdot \cos(\theta_1 + \theta_2) + L3 \cdot \cos(\theta_1 + \theta_2 + \theta_3) - 1.5, \quad (6.13)$$

so that $g_1(\theta)$ vanishes when the end effector is at the desired x -position, namely 1.5. Given this function, we pose,

Motivating Problem 1:

$$\begin{aligned} & \text{Minimize } f(\theta) := (\theta_1)^2 + (\theta_2)^2 + (\theta_3)^2 \\ & \text{subject to } g_1(\theta) = 0. \end{aligned} \quad (6.14)$$

The “subject to” means that a solution to this problem is **constrained** to find values for θ_1 , θ_2 , and θ_3 such that $g_1(\theta_1, \theta_2, \theta_3) = 0$. Clever us designed g_1 so that it only vanishes for values of θ where the x -position of the end effector is at the desired value of 1.5. **This is why Engineers use constraints in optimization problems: we often have goals to achieve that are distinct from minimizing some quantity.**

What is the role of the cost function? Well, we already saw in Fig. 6.4 that there are an infinite number of values of θ_1 , θ_2 , and θ_3 that achieve the full constraint

$$p_3 = \begin{bmatrix} 1.5 \\ 2.0 \end{bmatrix},$$

so there must be even more solutions that only constrain the x -component of p_3 and allows the y -component to be **unconstrained**. The cost function sorts through this infinite set of values and finds the one that minimizes the cost while satisfying the constraint. **This is a mainstay of modern engineering.**

Next, we define

$$g(\theta) := \begin{bmatrix} g_1(\theta) \\ g_2(\theta) \end{bmatrix} = \begin{bmatrix} L1 \cdot \cos(\theta_1) + L2 \cdot \cos(\theta_1 + \theta_2) + L3 \cdot \cos(\theta_1 + \theta_2 + \theta_3) - 1.5 \\ L1 \cdot \sin(\theta_1) + L2 \cdot \sin(\theta_1 + \theta_2) + L3 \cdot \sin(\theta_1 + \theta_2 + \theta_3) - 2.0 \end{bmatrix}, \quad (6.15)$$

so that $g(\theta)$ vanishes when the end effector is at the desired x - and y -positions. Given this function, we pose,

Motivating Problem 2:

$$\begin{aligned} & \text{Minimize } f(\theta) := (\theta_1)^2 + (\theta_2)^2 + (\theta_3)^2 \\ & \text{subject to } g(\theta) = 0. \end{aligned} \quad (6.16)$$

You could also write the problem as

$$\begin{aligned} & \text{Minimize } f(\theta) := (\theta_1)^2 + (\theta_2)^2 + (\theta_3)^2 \\ & \text{subject to } g_1(\theta) = 0 \text{ and } g_2(\theta) = 0. \end{aligned} \quad (6.17)$$

Once again, the “subject to” means that a solution to this problem is **constrained** to find values for θ_1 , θ_2 , and θ_3 such that $g(\theta_1, \theta_2, \theta_3) = 0$. And once again, we designed g so that it only vanishes when the (x, y) -position of the end effector is at the desired point in Cartesian space. **We repeat: this is why Engineers use constraints in optimization problems: we have goals to achieve that are distinct from minimizing some quantity.** The cost function takes care of the ambiguity in trying to place the end effector (that depends on three angles) at the desired location (specified by a 2-vector in the plane): we have three variables and only two constraints. Hence, we have an underdetermined problem, and constrained optimization is a perfect method for resolving the ambiguity associated with undetermined problems.

6.3.2 Vocabulary of Constrained Optimization

In the context of constrained optimization, several key terms and concepts are fundamental to understanding and solving optimization problems. Below is an overview of the primary vocabulary:

- **Cost Function:** Also known as the objective function, this is the function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ that we aim to minimize or maximize within the feasible set. The cost function quantifies the “cost” associated with a particular choice of $x = (x_1, x_2, \dots, x_n)$.

- **Decision Variables:** For a cost function, $f(x) = f(x_1, x_2, \dots, x_n)$, the components $x_i, 1 \leq i \leq n$ are the decision variables. They are the variables we can search over to find values that extremize the cost function.
- **Constraints:** Conditions that the solution x must satisfy. Constraints can be classified into two types:
 - **Equality Constraints:** These require that the constraint functions equal zero, $g(x) = 0$.
 - **Inequality Constraints:** These specify that the constraint functions must be greater than or equal to zero, $h(x) \geq 0$, or less than or equal to zero, depending on the problem formulation.
- **Feasible Set:** The set of all points x that satisfy the constraints of the optimization problem. For equality constraints, this set includes all x such that $g(x) = 0$, where $g(x)$ represents the constraint functions.
- **Constraint Qualification:** A set of conditions that ensure the mathematical program has a well-defined solution. One common qualification is that the rows of the Jacobian matrix of the constraint functions $g(x)$ are linearly independent at all points in the feasible set. Equivalently, the gradients of the constraints are linearly independent as column vectors.
- **Lagrange Multipliers:** Numbers λ associated with each constraint that measure the sensitivity of the objective function's optimal value to changes in the constraint boundary. The method of Lagrange multipliers is a strategy for finding the local maxima and minima of a function subject to equality constraints.
- **KKT Conditions (Karush-Kuhn-Tucker):** Necessary conditions for a solution in nonlinear programming to be optimal, given certain regularity conditions. These conditions generalize the method of Lagrange multipliers to handle inequality constraints. (Typically covered in 400-level courses; is not treated here.)
- **Slack Variables:** Additional variables introduced to transform inequality constraints into equality constraints, making it easier to apply certain optimization methods. (Typically covered in 400-level courses; is not treated here.)

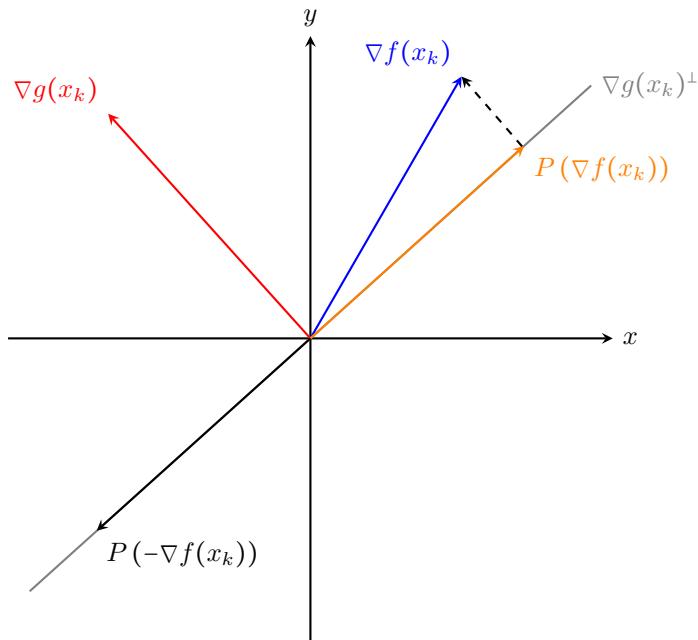


Figure 6.5: Geometry of gradient descent with an equality constraint. In the image, the line labeled $\nabla g(x_k)^\perp$ is the set of points (actually, the subspace) orthogonal to $\nabla g(x_k)$, the gradient of the cost. The angle between $\nabla f(x_k)$ and $\nabla g(x_k)$ is 72° . The dotted line shows $\alpha \cdot \nabla g(x_k)$, the proportion of $\nabla g(x_k)$ that must be removed (aka, subtracted) from $\nabla f(x_k)$ so that $\nabla f(x_k) - \alpha \nabla g(x_k)$ is orthogonal to the constraint gradient. The new vector, denoted $P(\nabla f(x_k))$, is an ascent direction that respects the constraint, whereas $P(-\nabla f(x_k))$, is a descent direction that respects the constraint.

6.3.3 Gradient Descent with a Single Equality Constraint

Suppose that the functions $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g : \mathbb{R}^n \rightarrow \mathbb{R}$ define the cost and the constraint in the constrained minimization problem

$$\begin{aligned} & \text{Minimize} && f(x) \\ & \text{subject to} && g(x) = 0. \end{aligned} \tag{6.18}$$

We seek to modify the gradient descent algorithm to accommodate the equality constraint. Let $x_k \in \mathbb{R}^n$ be the current value of the vector of decision variables, which we assume to be feasible, that is $g(x_k) = 0$. We seek Δx such that

$$\begin{aligned} f(x_k + \Delta x) &< f(x_k) \\ g(x_k + \Delta x) &= 0. \end{aligned} \tag{6.19}$$

We use our gradient notation to linearize these equations about x_k ,

$$\begin{aligned} f(x_k + \Delta x) &\approx f(x_k) + \nabla f(x_k) \bullet \Delta x \\ g(x_k + \Delta x) &\approx \nabla g(x_k) \bullet \Delta x. \end{aligned} \tag{6.20}$$

From (6.19) and (6.20), we deduce that Δx must satisfy

$$\begin{aligned} \nabla f(x_k) \bullet \Delta x &< 0 \\ \nabla g(x_k) \bullet \Delta x &= 0. \end{aligned} \tag{6.21}$$

We pose $\Delta x := -\nabla f(x_k) + \alpha \nabla g(x_k)$ and seek to choose α such that $\nabla g(x_k) \bullet \Delta x = 0$. Using results from ROB 101 Computational Linear Algebra, we compute

$$\begin{aligned} 0 &= \nabla g(x_k) \bullet \Delta x \\ &\Downarrow \\ 0 &= \nabla g(x_k) \bullet (-\nabla f(x_k) + \alpha \nabla g(x_k)) \\ &\Downarrow \\ 0 &= -\nabla g(x_k) \bullet \nabla f(x_k) + \alpha \nabla g(x_k) \bullet \nabla g(x_k) \end{aligned} \tag{6.22}$$

Hence, if $\nabla g(x_k) \neq 0_{n \times 1}$, then

$$\alpha = \frac{\nabla g(x_k) \bullet \nabla f(x_k)}{\nabla g(x_k) \bullet \nabla g(x_k)} \tag{6.23}$$

and we have

$$\Delta x = -\nabla f(x_k) + \frac{\nabla g(x_k) \bullet \nabla f(x_k)}{\nabla g(x_k) \bullet \nabla g(x_k)} \cdot \nabla g(x_k). \tag{6.24}$$

Fig 6.5 shows that this is the orthogonal projection of $-\nabla f(x_k)$ onto $[\nabla g(x_k)]^\perp := \{v \in \mathbb{R}^n \mid v \perp \nabla g(x_k)\}$. In simpler terms, we have removed from $-\nabla f(x_k)$ its component lying in the direction of $\nabla g(x_k)$. This gives the basic setup for gradient descent with a single equality constraint.

Method 6.17: Gradient Descent with a Single Equality Constraint

Consider the constrained optimization problem (6.18) for functions $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g : \mathbb{R}^n \rightarrow \mathbb{R}$ with optimal value x^* . If the gradients vector $\nabla g(x^*) \neq 0_{n \times 1}$, then a locally optimal x^* can be found as follows:

Initialization: Find a feasible point $x_0 \in \mathbb{R}^n$, that is, a point such that $g(x_0) = 0$; this can be accomplished with `NLsolve`.

Loop: Given a point $x_k \in \mathbb{R}^n$, repeat the following steps to either find x_{k+1} or terminate.

(a) **Compute:** $\nabla f(x_k)$, the gradient of f at x_k and $\nabla g(x_k)$, the gradient of the constraint vector.

(b) **Compute a Descent Direction that Respects the Constraint:** $\Delta x = -\nabla f(x_k) + \frac{\nabla g(x_k) \bullet \nabla f(x_k)}{\nabla g(x_k) \bullet \nabla g(x_k)} \cdot \nabla g(x_k)$

(c) **Update Step:** $x_{k+1} = x_k + s \cdot \Delta x$, where $s > 0$ is the step size.

- (d) **Termination Condition:** $\|\Delta x\| < a_{tol}$; in other words, no more progress can be made, and thus x_k is (to numerical precision) a stationary point of the optimization problem.

Notes:

- In practice, the updates x_k will “slowly” exit the feasible set. Hence, when $|g(x_k)|$ exceeds a user-defined tolerance, a new call to `NLsolve` should be made. This is illustrated in the following code block.
- We’ll see shortly that **JuMP does all of this and a whole lot more**. It is your best option for use in professional projects. Algorithms, like the one presented here, are essential for building fundamental understanding of gradient descent, the workhorse of Machine Learning (ML) and Inverse Kinematics!

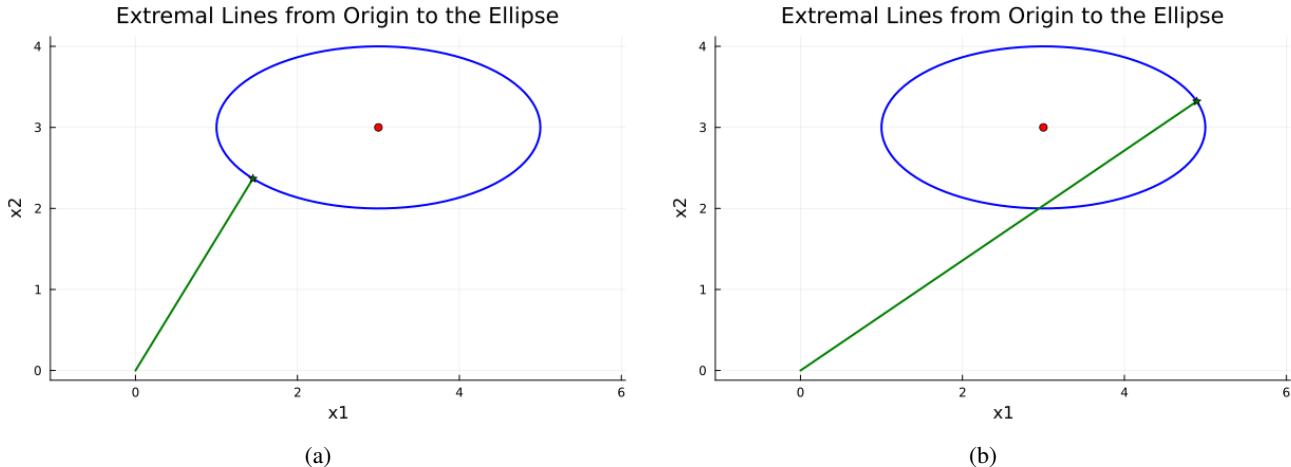


Figure 6.6: Optimization with equality constraints. In these examples, the cost function is $f(x, y) = (x)^2 + (y)^2$, distance squared from the origin to (x, y) . The constraint is $g(x, y) = \frac{(x-x_c)^2}{a^2} + \frac{(y-y_c)^2}{b^2} - 1 = 0$, an ellipse. (a) Shows the minimum distance solution subject to (x, y) lies on the ellipse, while (b) shows the maximum distance solution. The first problem uses gradient descent with one equality constraint, while the second uses gradient ascent with one equality constraint; the sign on the update step has been swapped.

Example 6.18. Solve the optimization problem (6.18) where the cost is $f(x, y) = x^2 + y^2$, the distance squared from the origin, and the constraint is that (x, y) lies on the ellipse defined by $g(x, y) = 0$, with $g(x, y) = \frac{(x-x_c)^2}{a^2} + \frac{(y-y_c)^2}{b^2} - 1 = 0$. Also find the maximizing solution.

Solution: The solution is worked out in the code block below. The results are plotted in Fig. 6.6. ■

1 `using ForwardDiff, LinearAlgebra, NLsolve`

```

2
3 # Define the cost function f(x) and the equality constraints G(x)
4 xc = 3.0; yc = 3.0 # centers
5 a = 2 # x semi-major axis
6 b = 1.0 #y semi-major axis
7
8 # note: a = b = 1/r for a circle
9
10 # Define the cost (distance squared)
11 function cost(x)
12     return x[1]^2 + x[2]^2
13 end

```

```

14
15 # Define the equality constraint: point must lie on the ellipse
16 function geq(x)
17     # note: a = b = 1/r for a circle
18     return (1/a)^2 * (x[1] - xc)^2 + (1/b)^2 * (x[2] - yc)^2 - 1
19 end
20
21 # Compute the gradient of the cost and constraint using ForwardDiff
22 grad_cost(x) = ForwardDiff.gradient(cost, x)
23 grad_geq(x) = ForwardDiff.gradient(geq, x)
24
25 # Remove from the gradient of the cost any component lying along the gradient(s) of the
26 # constraint(s)
26 function modGradientCost(x)
27     gradCost = grad_cost(x)
28     gradConstraint = grad_geq(x)
29
30     if norm(gradConstraint) < 1e-6
31         println("Gradient of the constraint is essentially zero and was ignored")
32         alpha = 0.0
33     else
34         alpha = dot(gradCost, gradConstraint) / dot(gradConstraint, gradConstraint)
35     end
36
37     modGradCost = gradCost - alpha * gradConstraint
38     return modGradCost
39 end
40
41 # Function to find a feasible starting point using NLsolve
42 function find_feasible_point(geq, x_guess)
43     # Build function for NLsolve
44     function constraint!(F, x)
45         F[1] = geq(x)
46     end
47     xFeasible = nlsolve(constraint!, x_guess)
48     return xFeasible.zero
49 end
50
51 # Starting guess (not feasible initially)
52 x_guess = [0.1; 0.1]
53 x0 = find_feasible_point(geq, x_guess) # Make the initial guess feasible
54
55 # Set step size, max iterations, and tolerances
56 s = 0.1
57 maxIter = 1000
58 gradTol = 1e-6
59 constraintTol = 1e-4
60
61 xk = x0
62 for k = 1:maxIter
63     modGradCost = modGradientCost(xk) # Use the modified gradient cost function
64     if norm(modGradCost) < gradTol # Convergence check
65         println("---- Iteration $k ----")
66         println("Converged.")
67         println("Cost: ", round.(cost(xk), digits=5))
68         println("Constraints (should be near 0): ", round.(geq(xk), digits=5))
69         println("-----")

```

```

70     break
71
72 else
73     xk = xk - s * modGradCost    # Update using the modified gradient
74 end
75 if norm(geq(xk)) > constraintTol
76     println("Updated x to maintain g(xk)=0 at interation = $k")
77     xk = find_feasible_point(geq, xk)    # Ensure feasibility
78 end
79
80 xStar = xk
81 costStar = cost(xStar)
82 geqStar = geq(xStar)
83
84 println("\nFinal Results:")
85 println("-----")
86 println("Optimal x: ", round.(xStar, digits=5))
87 println("Optimal cost(xStar) : ", round(costStar, digits=5))
88 println("Constraint(xStar) : ", round(geqStar, digits=5))

```

Output

```

Updated x to maintain g(xk)=0 at interation = 1
Updated x to maintain g(xk)=0 at interation = 2
Updated x to maintain g(xk)=0 at interation = 3
Updated x to maintain g(xk)=0 at interation = 4
Updated x to maintain g(xk)=0 at interation = 5
Updated x to maintain g(xk)=0 at interation = 6
Updated x to maintain g(xk)=0 at interation = 7
Updated x to maintain g(xk)=0 at interation = 8
---- Iteration 25 ----
Converged.
Cost: 7.70976
Constraints (should be near 0): 9.0e-5
-----
Final Results:
-----
Optimal x: [1.45048, 2.36767]
Optimal cost(xStar): 7.70976
Constraint(xStar): 9.0e-5
-----
```

6.3.4 Gradient Descent with a Vector of Equality Constraints

Suppose now that in addition to the cost function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, there are m equality constraints, $g_i : \mathbb{R}^n \rightarrow \mathbb{R}$ in the problem

$$\begin{aligned} & \text{minimize} && f(x) \\ & \text{subject to} && g_i(x) = 0, 1 \leq i \leq m \end{aligned} \tag{6.25}$$

We seek to modify the gradient descent algorithm to accommodate a vector of equality constraints. Let $x_k \in \mathbb{R}^n$ be the current value of the vector of decision variables, which we assume to be feasible, that is $g_i(x_k) = 0, 1 \leq i \leq m$. We seek Δx such that

$$\begin{aligned} & f(x_k + \Delta x) < f(x_k) \\ & g_i(x_k + \Delta x) = 0, 1 \leq i \leq m. \end{aligned} \tag{6.26}$$

We use our gradient notation to linearize these equations about x_k ,

$$\begin{aligned} f(x_k + \Delta x) &\approx f(x_k) + \nabla f(x_k) \bullet \Delta x \\ g_i(x_k + \Delta x) &\approx \nabla g_i(x_k) \bullet \Delta x, \quad 1 \leq i \leq m. \end{aligned} \tag{6.27}$$

From (6.19) and (6.27), we deduce that Δx must satisfy

$$\begin{aligned} \nabla f(x_k) \bullet \Delta x &< 0 \\ \nabla g_i(x_k) \bullet \Delta x &= 0, \quad 1 \leq i \leq m. \end{aligned} \tag{6.28}$$

Analogously to the single-constraint case, we pose

$$\Delta x := -\nabla f(x_k) + \sum_{i=1}^m \alpha_i \nabla g_i(x_k) = -\nabla f(x_k) + \underbrace{\begin{bmatrix} \nabla g_1(x_k) & \nabla g_2(x_k) & \cdots & \nabla g_m(x_k) \end{bmatrix}}_{\nabla G(x_k)} \cdot \underbrace{\begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_m \end{bmatrix}}_{\alpha} \tag{6.29}$$

and seek to solve for $\{\alpha_1, \alpha_2, \dots, \alpha_m\}$ such that

$$\nabla g_i(x_k) \bullet \Delta x = 0, \quad 1 \leq i \leq m. \tag{6.30}$$

Using results from ROB 101 Computational Linear Algebra, the above is equivalent to $\nabla G(x_k)^\top \cdot \Delta x = 0_{m \times 1}$, yielding

$$-\nabla G(x_k)^\top \cdot \nabla f(x_k) + \nabla G(x_k)^\top \cdot \nabla G(x_k) \cdot \alpha = 0_{m \times 1}. \tag{6.31}$$

If the columns of $\nabla G(x_k)$ are linearly independent, then the above equation has a unique solution. For small problems, you can invert the matrix $\nabla G(x_k)^\top \cdot \nabla G(x_k)$ to solve for the vector of coefficients α , but for larger problems, it is better to use matrix factorization, such as LU or QR. In Julia, the backslash command will choose the most appropriate method for the given problem. Also relying on ROB101, an alternative method to compute Δx in (6.29), is to apply Gram-Schmidt to the vectors

$$\{\nabla g_1(x_k), \nabla g_2(x_k), \dots, \nabla g_m(x_k), -\nabla f(x_k)\}, \tag{6.32}$$

where it is very important to place the gradient of the cost **after** the gradients of the constraints.

Example 6.19. Define the cost by $f : \mathbb{R}^4 \rightarrow \mathbb{R}$ be $f(x_1, x_2, x_3, x_4) = \sum_{i=1}^4 (x_i)^2$, distance squared from the origin to a point $x \in \mathbb{R}^4$, with the constraint that x must satisfy

$$0_{2 \times 1} = g(x) = \begin{bmatrix} g_1(x) \\ g_2(x) \end{bmatrix} = \begin{bmatrix} \frac{(x_1-3)^2}{2^2} + \frac{(x_2-3)^2}{1^2} - 1 \\ x_1 + x_2 + x_3 + x_4 \end{bmatrix}.$$

This is a minimum distance problem from the origin to the feasible set, that is, $\{x \in \mathbb{R}^4 \mid g(x) = 0_{2 \times 1}\}$. Use gradient descent with equality constraints to solve (6.25).

Solution: Using Gram-Schmidt to remove the components of $\nabla g_1(x_k)$ and $\nabla g_2(x_k)$ from $\nabla f(x_k)$.

```

1 """
2     gram_schmidt (U::Matrix)
3
4     Performs the Gram-Schmidt process on the columns of the matrix `U` (assumed to be
5     linearly independent)
6     to produce an orthogonal set of vectors stored in matrix `V` and a set of orthonormal
7     vectors stored
8     in `Vnormalized`.
9
# Arguments
- `U::Matrix`: A matrix whose columns are linearly independent vectors.

```

```

10
11 # Returns
12 - `V::Matrix`: A matrix containing the orthogonal vectors corresponding to the columns
13 - `Vnormalized::Matrix`: A matrix containing the orthonormal vectors obtained by
14 - normalizing the columns of `V`.
15 " "
16
17 using LinearAlgebra
18
19 function gram_schmidt (U::Matrix)
20     m, n = size(U)    # m = number of rows, n = number of columns (vectors)
21     V = Matrix{Float64}(undef, m, n)    # Matrix to store orthogonal vectors
22     Vnormalized = Matrix{Float64}(undef, m, n)    # Matrix to store normalized orthogonal
23     vectors
24
25     for i in 1:n
26         # Start with the original vector from U
27         vi = U[:, i]
28
29         # Subtract projections onto previous vectors
30         for j in 1:i-1
31             vj = V[:, j]
32             vi -= (dot(vi, vj) / dot(vj, vj)) * vj
33         end
34
35         # Store the orthogonal vector
36         V[:, i] = vi
37
38         # Normalize the vector and store in Vnormalized
39         Vnormalized[:, i] = vi / norm(vi)
40     end
41
42     return (V=V, Vnormalized=Vnormalized)
43 end

```

Output

```
gram_schmidt (generic function with 1 method)
```

```

1 # Example gradient descent with several equality constraints
2 # Not in the textbook
3
4 using ForwardDiff, LinearAlgebra, NLsolve, Random
5
6 # Set the seed to a specific value, for example, 1234
7 Random.seed!(1234)
8
9
10
11 # Define the cost function f(x) and the equality constraint G(x)
12 xc = 3.0; yc = 3.0 # centers
13 a = 2.0 # x semi-major axis
14 b = 1.0 # y semi-major axis
15
16 # Define the cost (distance squared)

```

```

17 function cost (x)
18     return x[1]^2 + x[2]^2 + x[3]^2 + x[4]^2
19 end
20
21 # Define the equality constraint: point must lie on the ellipse
22 function Geq (x)
23     # Equality constraints evaluated at x
24     g1 = (1/a)^2 * (x[1] - xc)^2 + (1/b)^2 * (x[2] - yc)^2 - 1 # ellipse in R4
25     g2 = x[1] + x[2] + x[3] + x[4] # 3 dimensional plane through the origin
26     return [g1, g2] # Return evaluated constraints as an array
27 end
28
29 # Compute the gradient of the cost and constraint using ForwardDiff
30 grad_cost(x) = ForwardDiff.gradient(cost, x)
31 Jac_constraints(x) = ForwardDiff.jacobian(Geq, x)
32
33 # Remove the component of the gradient of the cost lying along the gradient of the
34 # constraints
35 function modNegGradientCost (x)
36     gradCost = grad_cost (x)
37     matGradientsConstraints = Jac_constraints (x)' # transpose of Jacobian for constraint
38     # gradients
39
40     # Perform Gram-Schmidt on constraint gradients and cost gradient
41     U = [matGradientsConstraints gradCost]
42     F = gram_schmidt (U)
43     modNegGradCost = -F.V[:,end] #it's the negative of the last vector in the matrix when
44     # doing Gradient Descent
45     return modNegGradCost
46 end
47
48 # Function to find a feasible starting point using NLsolve
49 function find_feasible_point (Geq, x_guess)
50     function constraint! (F, x)
51         F[1] = Geq (x) [1]
52         F[2] = Geq (x) [2]
53     end
54     result = nlsolve (constraint!, x_guess)
55     return result.zero
56 end
57
58 # Set step size, max iterations, and tolerances
59 s = 0.1
60 maxIter = 1000
61 gradTol = 1e-6
62 constraintTol = 1e-5
63
64 # Starting guess (not feasible initially)
65 x_guess = rand(4,1)
66 xk = find_feasible_point (Geq, x_guess) # Make the initial guess feasible
67
68 # Gradient Descent with equality constraints loop
69 for k = 1:maxIter
70     modNegGradCost = modNegGradientCost (xk) # Use the modified gradient cost function
71     if norm(modNegGradCost) < gradTol # Convergence check
72         println ("\n---- Iteration $k ----")
73         println ("Converged .")

```

```

71    println("Cost: ", round.(cost(xk), digits=5))
72    println("Constraints (should be near 0): ", round.(Geq(xk), digits=6))
73    println("-----")
74    break
75  else
76    xk = xk + s * modNegGradCost # Update using the modified gradient
77  end
78
79  if norm(Geq(xk)) > constraintTol # Feasibility check
80    println("Updating for feasibility: G(xk) = 0_{m x 1} at iteration k = $k")
81    xk = find_feasible_point(Geq, xk) # Ensure feasibility
82  end
83 end
84
85 # Final results
86 xStar = xk
87 costStar = cost(xStar)
88
89 println("\nFinal Results:")
90 println("-----")
91 println("Optimal x: ", round.(xStar, digits=5))
92 println("Optimal cost(xStar): ", round(costStar, digits=5))
93 println("Minimal Distance is sqrt(cost(xStar)): ", round(sqrt(costStar), digits=5))

```

Output

Updating for feasibility: $G(x_k) = 0_{m \times 1}$ at iteration $k = 1$
 Updating for feasibility: $G(x_k) = 0_{m \times 1}$ at iteration $k = 2$

---- Iteration 63 ----
 Converged.
 Cost: 14.91277
 Constraints (should be near 0): [0.0, -0.0]

Final Results:

 Optimal x: [1.33912; 2.44289; -1.89101; -1.89101; ;]
 Optimal cost(xStar): 14.91277
 Minimal Distance is sqrt(cost(xStar)): 3.86171

Solution: This time using (6.31) to update Δx .

```

1 # Example gradient descent with several equality constraints
2 # Not in the textbook
3
4 using ForwardDiff, LinearAlgebra, NLsolve, Random
5
6 # Set the seed to a specific value, for example, 1234
7 Random.seed!(1234)
8
9
10
11 # Define the cost function f(x) and the equality constraint G(x)
12 xc = 3.0; yc = 3.0 # centers
13 a = 2.0 # x semi-major axis

```

```

14 b = 1.0 # y semi-major axis
15
16 # Define the cost (distance squared)
17 function cost (x)
18     return x[1]^2 + x[2]^2 + x[3]^2 + x[4]^2
19 end
20
21 # Define the equality constraint: point must lie on the ellipse
22 function Geq (x)
23     # Equality constraints evaluated at x
24     g1 = (1/a)^2 * (x[1] - xc)^2 + (1/b)^2 * (x[2] - yc)^2 - 1 # ellipse in R4
25     g2 = x[1] + x[2] + x[3] + x[4] # 3 dimensional plane through the origin
26     return [g1, g2] # Return evaluated constraints as an array
27 end
28
29 # Compute the gradient of the cost and constraint using ForwardDiff
30 grad_cost (x) = ForwardDiff.gradient(cost, x)
31 Jac_constraints (x) = ForwardDiff.jacobian(Geq, x)
32
33 # Remove the component of the gradient of the cost lying along the gradient of the
34 # constraints
35 function modNegGradientCost (x)
36     gradCost = grad_cost (x)
37     matGradientsConstraints = Jac_constraints (x)' # transpose of Jacobian for constraint
38     gradients
39
40     alpha = (matGradientsConstraints)' * matGradientsConstraints \ matGradientsConstraints' * gradCost
41
42     modNegGradCost = -gradCost + matGradientsConstraints * alpha #it's the negative of the last
43     vector in the matrix when doing Gradient Descent
44     return modNegGradCost
45 end
46
47 # Function to find a feasible starting point using NLsolve
48 function find_feasible_point (Geq, x_guess)
49     function constraint! (F, x)
50         F[1] = Geq (x) [1]
51         F[2] = Geq (x) [2]
52     end
53     result = nlsolve (constraint!, x_guess)
54     return result.zero
55 end
56
57 # Set step size, max iterations, and tolerances
58 s = 0.1
59 maxIter = 1000
60 gradTol = 1e-6
61 constraintTol = 1e-5
62
63 # Starting guess (not feasible initially)
64 x_guess = rand (4,1)
65 xk = find_feasible_point (Geq, x_guess) # Make the initial guess feasible
66
67 # Gradient Descent with equality constraints loop
68 for k = 1:maxIter
69     modNegGradCost = modNegGradientCost (xk) # Use the modified gradient cost function
70     if norm (modNegGradCost) < gradTol # Convergence check

```

```

68    println (" \n---- Iteration $k ----")
69    println ("Converged . ")
70    println ("Cost: ", round.(cost(xk), digits=5) )
71    println ("Constraints (should be near 0): ", round.(Geq(xk), digits=6) )
72    println ("----- ")
73    break
74 else
75     xk = xk + s * modNegGradCost # Update using the modified gradient
76 end
77
78 if norm(Geq(xk)) > constraintTol # Feasibility check
79     println ("Updating for feasibility: G(xk) = 0_{m x 1} at iteration k = $k")
80     xk = find_feasible_point(Geq, xk) # Ensure feasibility
81 end
82 end
83
84 # Final results
85 xStar = xk
86 costStar = cost(xStar)
87
88 println ("\nFinal Results:")
89 println ("-----")
90 println ("Optimal x: ", round.(xStar, digits=5) ) )
91 println ("Optimal cost(xStar): ", round(costStar, digits=5) )
92 println ("Minimal Distance is sqrt(cost(xStar)): ", round(sqrt(costStar), digits=5) )

```

Output

Updating for feasibility: $G(x_k) = 0_{m \times 1}$ at iteration $k = 1$
 Updating for feasibility: $G(x_k) = 0_{m \times 1}$ at iteration $k = 2$

---- Iteration 63 ----
 Converged.
 Cost: 14.91277
 Constraints (should be near 0): [0.0, -0.0]

Final Results:

 Optimal x: [1.33912; 2.44289; -1.89101; -1.89101;;]
 Optimal cost(xStar): 14.91277
 Minimal Distance is sqrt(cost(xStar)): 3.86171

6.3.5 Lagrange Multiplier for a Problem with a Single Equality Constraint

The method of Lagrange multipliers is a tool that reduces optimization problems subject to equality constraints to what looks like an unconstrained problem. The new problem has dimension $(n+1)$ instead of n , but at least you can apply the methods of Chapter 6.2.2 to solve it. Mathematicians are very good at taking what looks like a new, super-hard problem and reducing its solution to a problem they already know how to solve. It's a power mindset! Engineers have learned to do this, too, by the way.

Proposition 6.20: Lagrange Multiplier Theorem for Constrained Optimization

Suppose that the functions $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g : \mathbb{R}^n \rightarrow \mathbb{R}$ define the cost and the constraint in the constrained minimization problem

$$\begin{aligned} &\text{Minimize} && f(x) \\ &\text{subject to} && g(x) = 0. \end{aligned} \tag{6.33}$$

If x^* is a solution to (6.33) and $\frac{\partial g(x^*)}{\partial x} \neq 0$, then there exists a unique scalar $\lambda^* \in \mathbb{R}$ such that

$$\nabla f(x^*) + \lambda^* \nabla g(x^*) = 0, \quad (6.34)$$

in other words, x^* is a **stationary point of the unconstrained (cost) function**

$$L(x, \lambda^*) := f(x) = \lambda^* g(x).$$

Remark 6.21. For reasons we will see shortly,

- (a) the constant λ^* is called the **Lagrange Multiplier** associated with the constrained minimization problem, (6.33), and
- (b) $L(x, \lambda) : f(x) + \lambda g(x)$ (with a general λ instead of a specific value) is the **Lagrangian** associated with (6.33).
- (c) So far, we have not told you how to find the Lagrange multiplier; that is done in Method 6.22. But at least you understand why it is important: it reduces a problem we do not know how to solve to one we do know how to solve.
- (d) Because maximizing a function is the same as minimizing its negative, all of the above holds for constrained maximization problems as well. In fact, if you succeed in computing all of the stationary points of L , you will have found all of the minima of f , all of the maxima of f , and its saddle points. The stationary conditions are agnostic to min versus max.
- (e) A similar result can be given when there are several constraints, as in our Motivating Problem 2 in (6.16). The multi-constraint result will be placed in an “Optional Read” section along with an algebraic proof of the Method of Lagrange Multipliers, based on knowledge from ROB 101 Computational Linear Algebra.

Method 6.22: Method of Lagrange Multipliers

Consider the constrained optimization problem (6.33) for functions $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g : \mathbb{R}^n \rightarrow \mathbb{R}$ with optimal value x^* . If $\frac{\partial g(x^*)}{\partial x} \neq 0$, then the optimal x^* and λ^* can be found as follows:

Step 1: Form the **Lagrangian** function

$$L(x, \lambda) := f(x) + \lambda g(x). \quad (6.35)$$

Step 2: Solve for the **stationary points**^a of L , that is, find all solutions to

$$0_{(n+1) \times 1} = \nabla L(x, \lambda). \quad (6.36)$$

When the compact notation in (6.36) is expanded out, the terms needed in a calculation become more evident, namely

$$0_{n \times 1} = \begin{bmatrix} \frac{\partial L(x, \lambda)}{\partial x_1} \\ \frac{\partial L(x, \lambda)}{\partial x_2} \\ \vdots \\ \frac{\partial L(x, \lambda)}{\partial x_n} \end{bmatrix} = \nabla f(x) + \lambda \nabla g(x) \quad (6.37)$$

$$0_{1 \times 1} = \frac{\partial L(x, \lambda)}{\partial \lambda} = g(x).$$

For small problems, the calculations are easy to do by hand. For “smallish to medium-sized problems”, the `Symbolics` package makes short work of computing the gradient of L , and `NLSolve` can be “coaxed” into finding the roots of the gradient, that is, the stationary points.

The notation, $0_{1 \times 1}$, instead of 0, emphasizes that this line is scalar-valued; in fact, it imposes the scalar-valued (feasibility) constraint, $g(x) = 0$. We note that (6.37) has $(n+1)$ equations in $(n+1)$ unknowns, (x, λ) , which is why we **can now solve**

for the Lagrange multiplier λ^* , in addition to x^* .

Step 3: Evaluate the cost function $f(x)$ at each of the stationary points and choose the one (or ones) yielding a minimum. At this stage, we no longer have to worry about the constraint, $g(x)$, because it was taken care of in the computation of the stationary points of the Lagrangian, L .

^aA nice discussion on this aspect of Lagrange's Method is given in [Example 5: Numerical optimization](#) of Wikipedia.

We next solve a few examples that are on the level of difficulty that you will find in your 200- and 300-level Engineering courses. **With the exception of Robotics and IOE, the problems must be solvable algebraically because you will not be assumed to know anything about numerical methods.** You will have to solve them by hand, with paper and pencil.

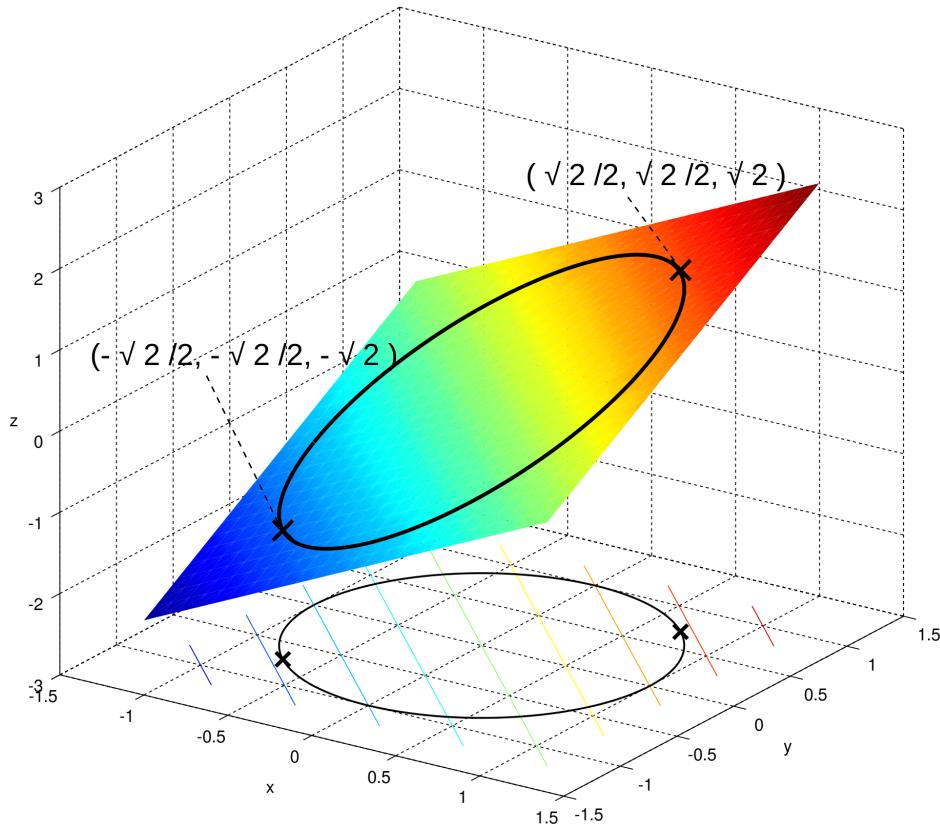


Figure 6.7: Illustration of a simple constrained optimization problem. For those who prefer geometric reasoning, the lines shown in the (x, y) -plane are the contour lines of the cost function, that is, lines along which $z = f(x, y) = x + y$ is constant. The unit circle is the constraint set, $g(x, y) = x^2 + y^2 - 1 = 0$. Hence, the set of feasible points is the unit circle. **Image Credit:** Wikipedia: Lagrange Multipliers: Example 1.

Example 6.23. Use the Method of Lagrange Multipliers to minimize $x + y$ subject to the constraint $x^2 + y^2 = 1$. Note that the problem is employing standard “non-vector” notation where one uses (x, y) to denote a point in \mathbb{R}^2 instead of $x = (x_1, x_2)^\top$.

The Example is based on [Wikipedia, Example 1](#). You may also want to look at [Wikipedia, Example 2](#), which we will not work.

Solution: Though it is not required, we write down the constrained optimization problem,

$$\begin{aligned} \text{Minimize } & f(x, y) = x + y \\ \text{subject to } & 0 = g(x, y) = x^2 + y^2 - 1, \end{aligned}$$

so as to make the identification of $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ and $g : \mathbb{R}^2 \rightarrow \mathbb{R}$ crystal clear. Following the steps of Method 6.22, we define

$$L(x, y, \lambda) := (x + y) + \lambda(x^2 + y^2 - 1).$$

We then compute the required partial derivatives

$$\begin{aligned} 0 &= \frac{\partial L(x, y, \lambda)}{\partial x} = 1 + 2\lambda \cdot x \\ 0 &= \frac{\partial L(x, y, \lambda)}{\partial y} = 1 + 2\lambda \cdot y \\ 0 &= \frac{\partial L(x, y, \lambda)}{\partial \lambda} = x^2 + y^2 - 1. \end{aligned} \tag{6.38}$$

Note that the first two lines of the above equation correspond to the first 2×1 block of (6.37), and the last line of the above equation corresponds to the last line of (6.37).

Fortunately for us –**and it has to be like this for the problem to be doable by hand**– we can easily solve the stationary conditions in (6.38) for x and y in terms of λ , giving us

$$\begin{aligned} x &= -\frac{1}{2\lambda} \\ y &= -\frac{1}{2\lambda}. \end{aligned}$$

Plugging these values into the constraint yields

$$0 = -1 + x^2 + y^2 = -1 + \left(-\frac{1}{2\lambda}\right)^2 + \left(-\frac{1}{2\lambda}\right)^2 = -1 + \frac{1}{2} \frac{1}{\lambda^2},$$

or

$$\lambda^2 = \frac{1}{2},$$

and therefore, $\lambda = \pm \frac{1}{\sqrt{2}} = \pm \frac{\sqrt{2}}{2}$. Putting all of this together and doing the messy, mistake-prone algebra, yields the hand-computed stationary points of the Lagrangian are

$$\begin{bmatrix} x^* \\ y^* \\ \lambda^* \end{bmatrix} = \begin{bmatrix} \frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} \\ -\frac{\sqrt{2}}{2} \end{bmatrix} \text{ and } \begin{bmatrix} x^* \\ y^* \\ \lambda^* \end{bmatrix} = \begin{bmatrix} -\frac{\sqrt{2}}{2} \\ -\frac{\sqrt{2}}{2} \\ \frac{\sqrt{2}}{2} \end{bmatrix}.$$

As promised, we have found λ^* ! (Exclamation point, and not factorial.)

Evaluating the cost function $f(x, y) = x + y$ at these points yields

$$f\left(\frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2}\right) = \sqrt{2}, \text{ and } f\left(-\frac{\sqrt{2}}{2}, -\frac{\sqrt{2}}{2}\right) = -\sqrt{2}.$$

Thus the constrained minimum is $-\sqrt{2}$, and if we had wanted it, the constrained maximum is $\sqrt{2}$. ■

Example 6.24. Let's check that Prop. 6.20 really does work by re-solving Example 6.23 with the side information that

$$\lambda^* = \frac{\sqrt{2}}{2}.$$

The goal is to check that we can find the minimum to the constrained problem by computing the stationary points of the unconstrained problem

$$L^*(x, y) := L(x, y, \lambda^*) = x + y + \underbrace{\frac{\sqrt{2}}{2}}_{\lambda^*} (x^2 + y^2 - 1).$$

Solution: We compute the gradient of L^* and set it equal to zero, namely

$$0_{2 \times 1} = \nabla L^*(x, y) = \begin{bmatrix} \frac{\partial L^*(x, y)}{\partial x} \\ \frac{\partial L^*(x, y)}{\partial y} \end{bmatrix} = \begin{bmatrix} 1 + \sqrt{2}x \\ 1 + \sqrt{2}y \end{bmatrix}.$$

Solving for x and y gives the same answer we computed in Example 6.23, once we simplify via $\frac{1}{\sqrt{2}} = \frac{\sqrt{2}}{2}$.

Notes:

- If we had wanted to find the maximum via the unconstrained problem, we would have used $\lambda^* = -\frac{\sqrt{2}}{2}$.
- All of this is just a theoretical exercise to underline that Lagrange's Method is intimately tied to the gradient of an appropriately defined cost function. If you don't get it, just realize that no 200- or 300-level Engineering course will test you on the theoretical underpinnings of Lagrange's Method. And if it is broached in a 400-level course, it will be taught (again) in the course. Second time is charm!

■

This next problem is a Calc III Classic: Minimize the surface area of a cylindrical can, subject to a constraint on its volume.

Why minimize surface area? Because that accounts for the amount of material it takes to manufacture the can, that is, its cost. Why constrain volume? To sell you a fixed amount of stuff. To formulate this problem, we need to recall two facts from High School Geometry:

- Volume of a **right circular cylinder** of height $h > 0$ and radius $r > 0$ is $V(r, h) = \pi \cdot r^2 \cdot h$.
- Surface area of the same cylinder is $S(r, h) = \underbrace{2\pi \cdot r \cdot h}_{\text{side of cylinder}} + \underbrace{2\pi \cdot r^2}_{\text{area of ends}} = 2\pi \cdot (r \cdot h + r^2)$.

Yes, we are writing them as functions. And yes, a can has two ends! If you want to know more about the cost of manufacturing tin cans, see [Tin Can Production Cost Reports](#).

Example 6.25. Minimize the surface area of a cylindrical can of radius r and height h subject to its volume being equal to V_0 .

Solution: Though it is not required, we write down the constrained optimization problem,

$$\begin{aligned} \text{Minimize } f(r, h) &= S(r, h) = 2\pi \cdot (r \cdot h + r^2) \\ \text{subject to } 0 &= g(r, h) = V(r, h) - V_0 = 2\pi \cdot r^2 \cdot h - V_0, \end{aligned}$$

so as to make the identification of $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ and $g : \mathbb{R}^2 \rightarrow \mathbb{R}$ crystal clear. Following the steps of Method 6.22, we define

$$L(r, h, \lambda) := 2\pi \cdot (r \cdot h + r^2) + \lambda (2\pi \cdot r^2 \cdot h - V_0).$$

We then compute the required partial derivatives

$$\begin{aligned} 0 &= \frac{\partial L(r, h, \lambda)}{\partial r} = 2\pi \cdot (h + 2r) + 4\pi \cdot \lambda \cdot r \cdot h = 2\pi \cdot (h + 2r + 2\lambda \cdot r \cdot h) \\ 0 &= \frac{\partial L(r, h, \lambda)}{\partial h} = 2\pi \cdot r + 2\pi \lambda \cdot r^2 = 2\pi \cdot r \cdot (1 + \lambda \cdot r) \\ 0 &= \frac{\partial L(r, h, \lambda)}{\partial \lambda} = 2\pi \cdot r^2 \cdot h - V_0. \end{aligned} \tag{6.39}$$

Again, note that the first two lines of the above equation correspond to the first 2×1 block of (6.37), and the last line of the above equation corresponds to the last line of (6.37).

Fortunately for us –**and it has to be like this for the problem to be doable by hand**– we can solve for r , h and λ using elementary algebraic manipulations. Indeed, from the second stationary condition, if $r \neq 0$, which is required for a physical can to have a non-zero volume, we obtain

$$\lambda \cdot r = -1.$$

Plugging this value into the first stationary condition, we obtain

$$(0 = h + 2r - 2h) \iff (h = 2r).$$

Plugging this into the last stationary condition, that is, the constraint, yields

$$(0 = 2\pi \cdot r^2 \cdot (2r) - V_0) \iff V_0 = 4\pi \cdot r^3.$$

Hence,

$$r = \sqrt[3]{\frac{V_0}{4\pi}}$$

$$h = 2 \cdot \sqrt[3]{\frac{V_0}{4\pi}}.$$

We could stop here. However, just to be super complete, putting all of this together after using

$$(\lambda \cdot r = -1) \iff \left(\lambda = -\frac{1}{r} \right),$$

we obtain

$$\begin{bmatrix} r^* \\ h^* \\ \lambda^* \end{bmatrix} = \begin{bmatrix} \sqrt[3]{\frac{V_0}{4\pi}} \\ 2 \cdot \sqrt[3]{\frac{V_0}{4\pi}} \\ -\sqrt[3]{\frac{4\pi}{V_0}} \end{bmatrix}.$$

As promised, by using Lagrange's Method, we have found λ^* as well as the proper dimensions of the can. Does someone have a patent on this? We mean the can's dimensions. Who's going to buy a Lagrange Multiplier? Maybe L'Hôpital?

In this case, we only found one extremal solution, the minimizing solution. You might ask yourself, why? **Why can't you maximize the surface area of a tin can?** Well, theoretically, you can do it. If we go back to our second stationary condition,

$$0 = \frac{\partial L(r, h, \lambda)}{\partial h} = 2\pi \cdot r \cdot (1 + \lambda \cdot r),$$

there was a second solution, namely, $r = 0$. Well, replacing that with $r = \epsilon > 0$ and solving for h such that $V_0 = 2\pi \cdot r^2 \cdot h$ gives

$$h = \frac{V_0}{2\pi\epsilon^2}.$$

Now, take the limit as $\epsilon \rightarrow 0^+$ and you get a “can” that is infinitely tall and narrow with a fixed volume. What is its surface area?

$$S = 2\pi \cdot (r \cdot h + r^2) = 2\pi \cdot (\epsilon \cdot \frac{V_0}{2\pi\epsilon^2} + \epsilon^2) = 2\pi \cdot (\frac{V_0}{2\pi\epsilon} + \epsilon^2) \rightarrow \infty \text{ as } \epsilon \rightarrow 0^+.$$

Hence, if you pose the nonsensical question of maximizing the surface area subject to a fixed volume, you can compute the nonsensical answer, as long as you believe in limits! To be clear, if you give this sort of problem to a numerical solver, it will go nuts. ■.

If you enjoy the idea of objects that have strange properties, such as finite volume and infinite surface area, you may enjoy the video [Gabriel's Horn Paradox](#) by Numberphile (aka, The [Painter's Paradox](#)) or [Torricelli's 2nd Paradox](#) and its [14th-century genius monk resolution](#) by Mathologer. The volume of the horn is finite, it is built from infinitely thin material so that the surface areas inside and outside the horn are the same, has a finite “bell” or “flare” at one end (instead of our uniformly thin cylinder), and yet it cannot be painted with a finite volume of paint. What's the paradox? Well, you can fill the horn with a finite volume of paint, and that volume is touching the inner surface of the horn, isn't it? Don't let this hurt your brain too much: this is a souped-up version of maximizing the surface area of a cylinder subject to a fixed volume.

Example 6.26. For $x \in \mathbb{R}^n$, minimize the norm squared of x subject to $a \bullet x = a^\top x = b_0$, where $0_{n \times 1} \neq a \in \mathbb{R}^n$ and $b \in \mathbb{R}$.

Solution: Though it is not required, we write down the constrained optimization problem,

$$\begin{aligned} & \text{Minimize} && f(x) = x \bullet x \\ & \text{subject to} && 0 = g(x) = a \bullet x - b_0, \end{aligned}$$

so as to make the identification of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g : \mathbb{R}^n \rightarrow \mathbb{R}$ crystal clear. Following the steps of Method 6.22, we define

$$L(x, \lambda) := x \bullet x + \lambda(a \bullet x - b_0).$$

We then compute the required partial derivatives in (6.37) using the vector product rule in Prop. 5.45:

$$0_{n \times 1} = \nabla f(x) + \lambda \nabla g(x) = 2x + \lambda a$$

$$0_{1 \times 1} = \frac{\partial L(x, \lambda)}{\partial \lambda} = a \bullet x - b_0 = a^\top x - b_0.$$

Hence, the first equation gives $x = -\frac{1}{2}\lambda a$, and substituting this into the second equation gives

$$a^\top \left(-\frac{1}{2}\lambda a \right) - b_0 = 0,$$

which in turn gives

$$\lambda = -2 \frac{1}{a^\top \cdot a} \cdot b_0.$$

Therefore, substituting this back into $x = -\frac{1}{2}\lambda a$ gives

$$x = \frac{b_0}{a^\top a} \cdot a.$$

■

6.3.6 (Optional Read:) Lagrange Multipliers for a Problem with a Vector of Equality Constraints

The method of Lagrange multipliers has a natural extension to problems with a vector of equality constraints. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be the real-valued cost function as before and assume we have a vector of equality constraints $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$, where

$$g(x) = \begin{bmatrix} g_1(x) \\ g_2(x) \\ \vdots \\ g_m(x) \end{bmatrix} = \begin{bmatrix} g_1(x_1, x_2, \dots, x_n) \\ g_2(x_1, x_2, \dots, x_n) \\ \vdots \\ g_m(x_1, x_2, \dots, x_n) \end{bmatrix}.$$

Optimization Problem with Several Constraints

$$\begin{aligned} & \text{Minimize} && f(x) \\ & \text{subject to} && g_i(x) = 0, \quad 1 \leq i \leq m, \end{aligned} \tag{6.40}$$

We'll skip the vector version of Prop. 6.20 and go straight to the heart of the matter, Lagrange's Method, now that you see how it works.

Method 6.27: Method of Lagrange Multipliers with a Vector of Equality Constraints

Consider the constrained optimization problem (6.40) for functions $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ with optimal value x^* . If the columns vectors from the gradients $\{\nabla g_1(x^*), \nabla g_2(x^*), \dots, \nabla g_m(x^*)\}$ are linearly independent (equivalently, the rows

of the Jacobian $\frac{\partial g(x^*)}{\partial x}$ are linearly independent), then a locally optimal x^* **and** vector of Lagrange multipliers

$$\lambda^* = \begin{bmatrix} \lambda_1^* \\ \lambda_2^* \\ \vdots \\ \lambda_m^* \end{bmatrix}$$

can be found as follows:

Step 1: Form the **Lagrangian**

$$L(x, \lambda) := f(x) + \sum_{i=1}^m \lambda_i g_i(x) = f(x) + \lambda_1 g_1(x) + \lambda_2 g_2(x) + \cdots + \lambda_m g_m(x). \quad (6.41)$$

Step 2: Solve for the **stationary points** of L , that is, find all solutions to

$$0_{(n+m) \times 1} = \nabla L(x, \lambda);$$

when expanded out, we have

$$\begin{aligned} 0_{n \times 1} &= \begin{bmatrix} \frac{\partial L(x, \lambda)}{\partial x_1} \\ \frac{\partial L(x, \lambda)}{\partial x_2} \\ \vdots \\ \frac{\partial L(x, \lambda)}{\partial x_n} \end{bmatrix} = \nabla f(x) + \sum_{i=1}^m \lambda_i \nabla g_i(x) \\ 0_{m \times 1} &= \begin{bmatrix} \frac{\partial L(x, \lambda)}{\partial \lambda_1} \\ \frac{\partial L(x, \lambda)}{\partial \lambda_2} \\ \vdots \\ \frac{\partial L(x, \lambda)}{\partial \lambda_m} \end{bmatrix} = \begin{bmatrix} g_1(x) \\ g_2(x) \\ \vdots \\ g_m(x) \end{bmatrix}. \end{aligned} \quad (6.42)$$

We note that the bottom $m \times 1$ block of (6.42), imposes the vector of (feasibility) constraints, $g(x) = 0$. Moreover, (6.42) has $(n+m)$ equations in $(n+m)$ unknowns, (x, λ) , giving us enough information to solve for x^* and the Lagrange multipliers, λ^* . When the equations are nonlinear, there can be multiple solutions or no solutions.

Step 3: Evaluate the cost function $f(x)$ at each of the stationary points and choose the one (or ones) yielding a minimum. Good luck!

Note: The condition that the columns of the gradient $\nabla g(x^*)$ be linearly independent (equivalently, the rows of the Jacobian, $\frac{\partial g(x^*)}{\partial x}$ are linearly independent), ensures that the Lagrange multipliers are unique.

We next solve an example by “pen and paper” to illustrate the method..

Example 6.28. Maximize $f(x, y, z) = (x + y + z)$ subject to the constraints $x^2 + y^2 = 1$ and $x + z = 2$. The maximization is not a typo. We wish to illustrate our “trick” of minimizing the negative of the objective function.

The way we are approaching the problem, namely through Lagrange Multipliers, we can find all of the extreme points, whether we include the minus sign on the cost or not. With typical numerical solvers, this is not the case.

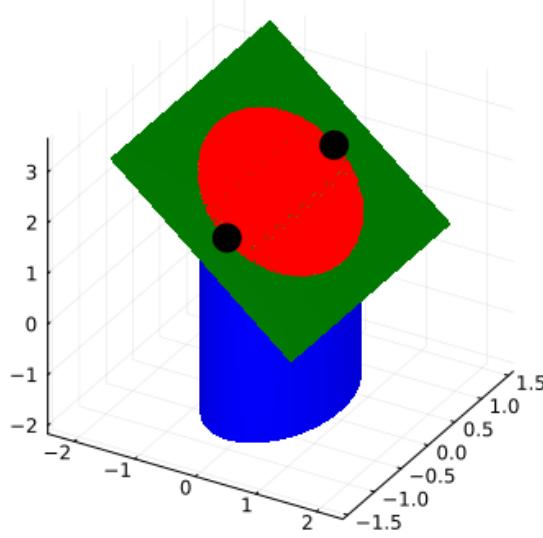


Figure 6.8: The first constraint in Example 6.28 gives a cylinder about the z -axis having radius one, while the second constraint defines a plane. The red elliptical/circular shape in the plane, excluding its interior, is the intersection of the two constraints, thereby defining the feasible set over which the cost is to be optimized. The two black points are the extrema, one corresponding to a maximum of $f(x, y, z) = x + y + z$ and the other to a minimum.

Solution: Though it is not required, we write down the constrained optimization problem,

$$\begin{aligned} \text{Minimize } & f(x, y, z) = -(x + y + z) \\ \text{subject to } & \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} g_1(x, y, z) \\ g_2(x, y, z) \end{bmatrix} = \begin{bmatrix} x^2 + y^2 - 1 \\ x + z - 2 \end{bmatrix}, \end{aligned}$$

so as to make the identification of $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ and $g : \mathbb{R}^3 \rightarrow \mathbb{R}^2$ crystal clear. Following the steps of Method 6.27, we define

$$L(x, y, z, \lambda_1, \lambda_2) = -(x + y + z) + \lambda_1(x^2 + y^2 - 1) + \lambda_2(x + z - 2).$$

Taking the required partial derivatives, stacking them in a column, and setting them to zero, we obtain the stationary conditions,

$$\begin{aligned} 0 &= \frac{\partial L(x, y, z, \lambda_1, \lambda_2)}{\partial x} = -1 + 2\lambda_1 \cdot x + \lambda_2 \\ 0 &= \frac{\partial L(x, y, z, \lambda_1, \lambda_2)}{\partial y} = -1 + 2\lambda_1 \cdot y \\ 0 &= \frac{\partial L(x, y, z, \lambda_1, \lambda_2)}{\partial z} = -1 + \lambda_2 \\ 0 &= \frac{\partial L(x, y, z, \lambda_1, \lambda_2)}{\partial \lambda_1} = x^2 + y^2 - 1 \\ 0 &= \frac{\partial L(x, y, z, \lambda_1, \lambda_2)}{\partial \lambda_2} = x + z - 2. \end{aligned} \tag{6.43}$$

Fortunately for us –and it has to be like this for the problem to be doable by hand– we can solve for x, y, z, λ_1 and λ_2 using

elementary algebraic manipulations. Indeed,

$$\begin{aligned} \text{eqn # 3: } 0 &= -1 + \lambda_2 \implies \lambda_2 = 1 \\ \text{eqn # 1: } 0 &= 2\lambda_1 \cdot x \implies x = 0 \text{ (because eqn # 2: implies neither } y \text{ nor } \lambda_1 \text{ can equal zero)} \\ \text{eqn # 5: } 0 &= x + z = -2 \implies 0 = \frac{x}{y} \implies z = 2 \\ \text{eqn # 4: } 0 &= x^2 + y^2 - 1 \implies y = \pm 1 \\ \text{eqn # 2: } 0 &= -1 + 2\lambda_1 \cdot y \implies \lambda_1 = \text{sign}(y) \cdot \frac{1}{2}, \end{aligned}$$

where eqn # i: corresponds to row i of (6.43). If there is a confusing point, it is the conclusion that eqn # 1: implies $x = 0$. This is true because if $\lambda_1 = 0$, then eqn # 2: becomes $0 = -1$, which is impossible. Hence, λ_1 cannot be equal to zero, and hence, eqn # 1: implies $x = 0$. This is basic algebra. It's not fun math, perhaps, but still, just basic algebra. **Calculus gave us the equations; algebra is being used to solve them!**

Putting all of this together for the “real” cost function $f(x, y, z) = x + y + z$ we wanted to maximize, **the cost without the minus sign**, yields

$$\begin{bmatrix} x^* \\ y^* \\ z^* \\ \lambda_1^* \\ \lambda_2^* \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ 2 \\ \frac{1}{2} \\ 1 \end{bmatrix} \text{ giving a maximum of } f^* = 3, \text{ and} \quad \begin{bmatrix} x^* \\ y^* \\ z^* \\ \lambda_1^* \\ \lambda_2^* \end{bmatrix} = \begin{bmatrix} 0 \\ -1 \\ 2 \\ -\frac{1}{2} \\ 1 \end{bmatrix} \text{ giving a minimum of } f^* = 1.$$

■

Let’s now see how to do Lagrange Multipliers in code! The solution combines our ability to compute gradients symbolically, understand that stationary points correspond to roots of the gradient, and numerically solve root-finding problems. Not bad! This puts it all together.

```

1 using Symbolics, NLsolve, Random
2
3 # -----
4 # PROBLEM DEFINITION
5 # -----
6
7 # Define symbolic variables for the optimization
8 @variables var[1:3] # Actual problem variables
9 @variables lam[1:2] # Lagrange multipliers
10
11 # Objective function to be minimized
12 f = -(var[1] + var[2] + var[3])
13
14 # Constraint equations
15 g1 = var[1]^2 + var[2]^2 - 1
16 g2 = var[1] + var[3] - 2
17
18 # Formulate the Lagrangian: objective function + lagrange multipliers for constraints
19 L = f + lam[1]*g1 + lam[2]*g2
20
21 # Compute the gradient of the Lagrangian with respect to both variables and Lagrange
     multipliers

```

```

22 gradL = Symbolics.gradient(L, [var; lam])
23
24 # For debugging: display the gradient or convert it to LaTeX for better visualization
25 if false
26     display(gradL)
27 else
28     latexify(gradL) |> println
29 end
30
31 # -----
32 # NUMERICAL SOLUTION APPROACH
33 # -----
34
35 # Convert the symbolic gradient expression to a numerical function for computational
36 # efficiency
36 gradL_num = build_function(gradL, [var; lam])
37 gradL_num = eval(gradL_num[1])
38
39 # Initialize an array to store solutions
40 f_solutions = Float64[]
41
42 # Solve the optimization problem for multiple initial guesses
43 # This is done to explore different solutions since NLsolve can get trapped in local
44 # minima
44 for k = 1:100
45     # Generate a random initial guess for the optimization variables and Lagrange
46     # multipliers
47     # randn is the normal or gaussian distribution, which includes positive and negative
48     # values
47     initial_guess = randn(5)
49
50     # Use nlsolve to find where the gradient is zero, which corresponds to an extremum
50     result = nlsolve(gradL_num, initial_guess, xtol=1e-6, ftol=1e-7)
51
52
53 # Process the solution only if nlsolve successfully converged
54 if converged(result)
55     # Extract the values of the problem variables from the solution
56     var_val = result.zero[1:3]
57     lam_val = result.zero[4:5]
58
59     # Compute the value of the objective function at the found extremum
60     f_extremePoint = -sum(var_val)
61
62     # Add the found solution only if it is sufficiently different from the previously
63     # found solutions
63     if all(abs(f_extremePoint - sol) > 1e-2 for sol in f_solutions)
64         push!(f_solutions, f_extremePoint)
65         @show [var_val; lam_val]
66     end
67 end
68 end
69
70 # Sort solutions for easier analysis
71 sorted_solutions = sort(f_solutions)

```

Output

$$\begin{bmatrix} -1.0000 + lam2 + 2.0000lam1var1 \\ -1.0000 + 2.0000lam1var2 \\ -1.0000 + lam2 \\ -1.0000 + var1^2 + var2^2 \\ -2.0000 + var1 + var3 \end{bmatrix} \quad (6.44)$$

```
[var_val; lam_val] = [2.5557001655682698e-9, -1.000000021988782, 1.999999974442995,
-0.499999970071819, 1.0]
[var_val; lam_val] = [2.672589761676145e-12, 1.000000000001114, 1.9999999999973277,
0.4999999999892936, 1.0]

2-element Vector{Float64}:
-3.000000000001114
-0.999999978011215
```

Remark 6.29. The above problem has two distinct roots. To which one will `NLsolve` converge? It depends on the initial condition. That is why the solver is run 100 times with random initial conditions, giving us a chance to find both of them.

Here is a solution using `JuMP`. As a professional, this is preferred over the above solution using Lagrange Multipliers. Why? **JuMP's code has been optimized for speed and tested by enough users that it is pretty bug-free**. For doing HW in Engineering courses, being able to run the symbolic version of Lagrange Multipliers will help you tremendously. Moreover, you now know several ways to check your HW and learn from your mistakes.

```

1 using JuMP, Ipopt
2
3 # Create a new model with Ipopt as the solver
4 model = Model(Ipopt.Optimizer)
5
6 # Define the variables
7 @variable(model, var[1:3])
8
9 # Set the objective function
10 @NLobjective(model, Min, -(var[1] + var[2] + var[3]))
11
12 # Add the constraint equations
13 @NLconstraint(model, var[1]^2 + var[2]^2 == 1)
14 @NLconstraint(model, var[1] + var[3] == 2)
15
16 # Fine tune a few params
17 set_optimizer_attribute(model, "max_iter", 1000)    # Allow up to 1000 iterations
18 set_optimizer_attribute(model, "tol", 1e-9)          # Require a very accurate solution
19
20 # Suppress much of the output
21 set_optimizer_attribute(model, "print_level", 4)
22
23 # Solve the problem
24 optimize!(model)
25
26 # Display the results
27 println("Objective value: ", objective_value(model))
```

```

28 println ("Solution: ")
29 println ("var[1] = ", value(var[1]))
30 println ("var[2] = ", value(var[2]))
31 println ("var[3] = ", value(var[3]))

```

Output

```

Total number of variables.....: 3
    variables with only lower bounds: 0
    variables with lower and upper bounds: 0
    variables with only upper bounds: 0
Total number of equality constraints.....: 2
Total number of inequality constraints.....:
    inequality constraints with only lower bounds: 0
    inequality constraints with lower and upper bounds: 0
    inequality constraints with only upper bounds: 0

Number of Iterations....: 15

                                (scaled)          (unscaled)
Objective.....: -3.000000000000000e+00 -3.000000000000000e+00
Dual infeasibility.....: 0.000000000000000e+00 0.000000000000000e+00
Constraint violation....: 0.000000000000000e+00 0.000000000000000e+00
Variable bound violation: 0.000000000000000e+00 0.000000000000000e+00
Complementarity.....: 0.000000000000000e+00 0.000000000000000e+00
Overall NLP error.....: 0.000000000000000e+00 0.000000000000000e+00

Number of objective function evaluations      = 16
Number of objective gradient evaluations     = 16
Number of equality constraint evaluations    = 16
Number of inequality constraint evaluations = 0
Number of equality constraint Jacobian evals = 16
Number of inequality constraint Jacobian evals = 0
Number of Lagrangian Hessian evaluations     = 15
Total seconds in IPOPT                      = 0.012

EXIT: Optimal Solution Found.
Objective value: -3.0
Solution:
var[1] = 4.579141502905379e-17
var[2] = 1.0
var[3] = 2.0

```

For our final example, we recover the solution to the underdetermined least squares problem treated in Chapter 9.9 of the ROB 101 *Computational Linear Algebra* textbook.

Example 6.30. Consider an underdetermined system of linear equations $Ax = b$, where A is $m \times n$. When the rows of A are linearly independent (equivalently, the columns of A^\top are linearly independent), determine a minimum norm squared solution to $Ax = b$.

Solutions: You know the routine: even though it is not required, we write down the constrained optimization problem

$$\begin{aligned} \text{Minimize } & f(x) = x \bullet x = x^\top \cdot x \\ \text{subject to } & 0_{m \times 1} = g(x) = Ax - b, \end{aligned}$$

so as to make the identification of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ crystal clear. Following the steps of Method 6.22, we define

$$L(x, \lambda) := x \bullet x + \lambda_1(a_1 x - b_1) + \lambda_2(a_2 x - b_2) + \cdots + \lambda_m(a_m x - b_m),$$

where a_i are the rows of A and b_i are the entries of b . We then compute the required derivatives

$$\begin{aligned} 0_{n \times 1} &= \nabla(x \bullet x) + \lambda_1 \nabla(a_1 x - b_1) + \lambda_2 \nabla(a_2 x - b_2) + \dots + \lambda_m \nabla(a_m x - b_m) \\ &= 2x + \lambda_1 \cdot (a_1)^\top + \lambda_2 \cdot (a_2)^\top + \dots + \lambda_m \cdot (a_m)^\top \\ &= 2x + A^\top \cdot \lambda \\ 0_{m \times 1} &= Ax - b. \end{aligned}$$

From the first n equations, we obtain $x = -\frac{1}{2}A^\top \cdot \lambda$, and from the last m equations we obtain the constraint, $Ax = b$. Substituting the first into the second and doing a bit of algebra yields

$$(A \cdot A^\top) \cdot \lambda = -2b.$$

Solving for λ and substituting into $x = -\frac{1}{2}A^\top \cdot \lambda$ yields

$$x = A^\top \cdot (A \cdot A^\top)^{-1} \cdot b,$$

the same answer we had in ROB 101, namely,

$$x^* = \arg \min_{Ax=b} \|x\|^2 \iff x^* = A^\top \cdot (A \cdot A^\top)^{-1} b \iff x^* = A^\top \alpha \text{ and } A \cdot A^\top \alpha = b. \quad (6.45)$$

Once we identify $\alpha = -\frac{1}{2}\lambda$, the correspondence is spot on. As in ROB 101, we recommend that the minimum norm squared solution be computed with the right-hand side of (6.45) so that the matrix inverse is avoided, but for small problems, the middle answer is fine. ■

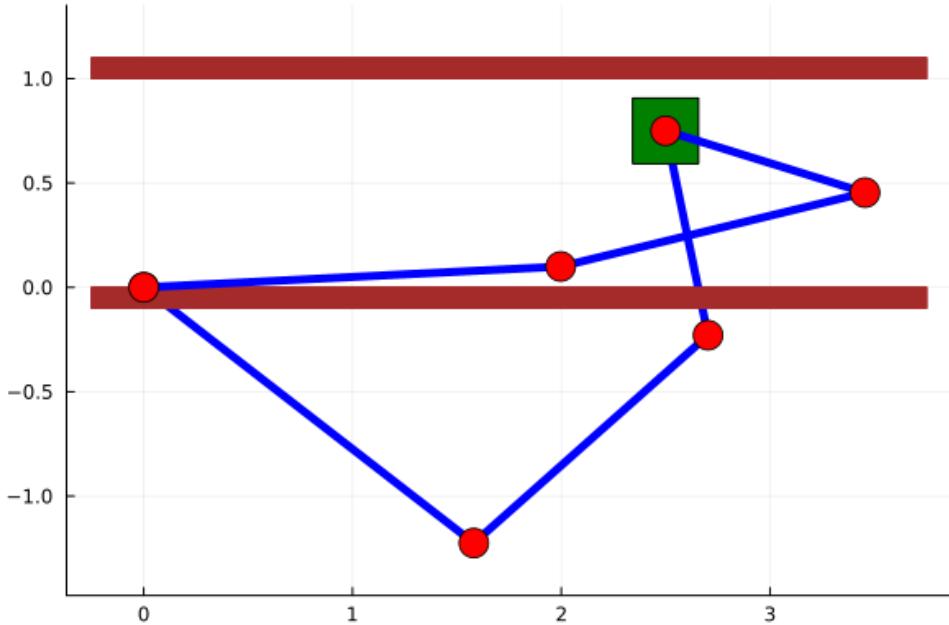


Figure 6.9: A pseudo-realistic manipulation task that requires inequality constraints in addition to equality constraints. The brown objects represent obstacles we must avoid while reaching into the green box with the robot's end effector. A solution with only equality constraints does reach into the green box, but “pass through” the obstacles. By adding inequality constraints, Jump finds the desired solution. The cost of the tightly constrained problem is approximately double of the other.

6.3.7 3-Link Manipulator Meets Inequality Constraints

We return to the three-link manipulator. Our goal is to place the end effector in a box, that will use **equality constraints**, and we want to avoid contacting an obstacle, which will use **inequality constraints**. We first set the problem up and then solve it using JuMP.

A standard form for a constrained optimization problem includes both **equality and inequality constraints**. Many solvers require you to write the constraints in the form

$$g_j^{\text{in}}(x) \leq 0,$$

as we show in (6.46). Some solvers, such as JuMP do not make you jump through so many hoops (pun intended) and accept constraints in the form

- $g_j^{\text{in}}(x) \leq c_j$

- $g_j^{\text{in}}(x) < c_j$

- $g_j^{\text{in}}(x) > c_j$

- $g_j^{\text{in}}(x) \geq c_j$.

Standard form for many Optimization Problems with inequality and equality constraints:

$$\begin{aligned} & \text{Minimize} && f(x) \\ & \text{subject to} && g_i^{\text{eq}}(x) = 0, \quad 1 \leq i \leq m, \\ & && g_j^{\text{in}}(x) \leq 0, \quad 1 \leq j \leq p, \end{aligned} \tag{6.46}$$

For the 3-link manipulator, we'll take the cost function as

$$f(\theta) = (\theta_1)^2 + (\theta_2)^2 + (\theta_3)^2.$$

We need two equality constraints to place the end effector in the box located at $(3.5, 0)$, namely

$$g^{\text{eq}}(\theta) = \begin{bmatrix} g_1^{\text{eq}}(\theta) \\ g_2^{\text{eq}}(\theta) \end{bmatrix} = \begin{bmatrix} L1 \cdot \cos(\theta_1) + L2 \cdot \cos(\theta_1 + \theta_2) + L3 \cdot \cos(\theta_1 + \theta_2 + \theta_3) - 3.5 \\ L1 \cdot \sin(\theta_1) + L2 \cdot \sin(\theta_1 + \theta_2) + L3 \cdot \sin(\theta_1 + \theta_2 + \theta_3) - 0.0 \end{bmatrix},$$

and we'll impose four inequality constraints on the vertical components of p_1 and p_2 of the robot, namely,

$$g^{\text{in}}(\theta) = \begin{bmatrix} g_1^{\text{in}}(\theta) \\ g_2^{\text{in}}(\theta) \\ g_3^{\text{in}}(\theta) \\ g_4^{\text{in}}(\theta) \end{bmatrix} = \begin{bmatrix} L1 \cdot \sin(\theta_1) \geq 0.1 \\ L1 \cdot \sin(\theta_1) + L2 \cdot \sin(\theta_1 + \theta_2) \geq 0.1 \\ L1 \cdot \sin(\theta_1) \leq 0.95 \\ L1 \cdot \sin(\theta_1) + L2 \cdot \sin(\theta_1 + \theta_2) \leq 0.95 \end{bmatrix}.$$

We do not impose constraints on their horizontal positions because we don't care where the joints end up. We simply want the links to be above the obstacle.

Here is what it looks like in JuMP

```

1 # JuMP solves the problem
2 using Plots, JuMP, Ipopt
3
4 # Define the model parameters for the robot manipulator
5 function modelParameters()
6     # Lengths of the three links
7     L1, L2, L3 = [2, 1.5, 1]

```

```

8     return (L1=L1, L2=L2, L3=L3)
9 end
10
11 # Calculate the positions of the robot's links based on the joint angles
12 function linkPositions(th1, th2, th3)
13     params = modelParameters()
14     # Base position
15     p0 = [0.0; 0.0]
16     # Calculate positions of the other joints
17     p1 = p0 + [params.L1 * cos(th1); params.L1 * sin(th1)]
18     p2 = p1 + [params.L2 * cos(th1+th2); params.L2 * sin(th1+th2)]
19     p3 = p2 + [params.L3 * cos(th1+th2+th3); params.L3 * sin(th1+th2+th3)]
20     return (p0=p0, p1=p1, p2=p2, p3=p3)
21 end
22
23 # Plot a rectangle on the current plot
24 function plot_rectangle(xc, yc, w, h; color="brown")
25     # Calculate the corners of the rectangle
26     x1 = xc - w/2
27     x2 = xc + w/2
28     y1 = yc
29     y2 = yc + h
30     # Plot the rectangle
31     plot!([x1, x2, x2, x1], [y1, y1, y2, y1], fill=true, fillalpha=1.0,
32           fillcolor=color, linecolor=color, legend=false)
33 end
34
35 # Plot the robot's configuration
36 function plot_points(positions, label; line_thickness=5, ball_size=10)
37     # Extract joint positions
38     p0 = positions.p0
39     p1 = positions.p1
40     p2 = positions.p2
41     p3 = positions.p3
42     # Plot obstacles
43     plot_rectangle(1.75, 1, 4, .1)
44     plot_rectangle(1.75, -.1, 4, .1)
45     # Plot the robot's configuration
46     scatter!([p3[1]], [p3[2]], color=:green, markersize=2*ball_size, marker=:square, label=nothing)
47     plot!([p0[1], p1[1], p2[1], p3[1]], [p0[2], p1[2], p2[2], p3[2]],
48           linewidth=line_thickness, color=:blue, label=label)
49     scatter!([p0[1], p1[1], p2[1], p3[1]], [p0[2], p1[2], p2[2], p3[2]], color=:red, markersize=
50         ball_size, label=nothing)
51 end
52
53 # Solve the optimization problem and plot the solution
54 function solve_and_plot(use_inequality_constraints)
55     params = modelParameters()
56     model = Model(Ipopt.Optimizer)
57     # Solver settings
58     set_optimizer_attribute(model, "max_iter", 1000)
59     set_optimizer_attribute(model, "tol", 1e-9)
60     set_optimizer_attribute(model, "print_level", 1)
61     # Define optimization variables and constraints
62     @variable(model, -π <= th1 <= π)
63     @variable(model, -π <= th2 <= π)
64     @variable(model, -π <= th3 <= π)

```

```

64 @constraint(model, g1, params.L1*cos(th1) + params.L2*cos(th1 + th2) + params.L3*cos(th1 + th2
+ th3) - 2.5 == 0)
65 @constraint(model, g2, params.L1*sin(th1) + params.L2*sin(th1 + th2) + params.L3*sin(th1 + th2
+ th3) - 0.75 == 0)
66 if use_inequality_constraints
67     @constraint(model, g3, params.L1*sin(th1) <= 0.95)
68     @constraint(model, g4, params.L1*sin(th1) + params.L2*sin(th1 + th2) <= 0.95)
69     @constraint(model, g5, params.L1*sin(th1) >= 0.1)
70     @constraint(model, g6, params.L1*sin(th1) + params.L2*sin(th1 + th2) >= 0.1)
71 end
72 # Define the objective function
73 @objective(model, Min, th1^2 + th2^2 + th3^2)
74 # Solve the optimization problem
75 optimize!(model)
76 # Display the results
77 label = use_inequality_constraints ? "With Inequality Constraints" : "No Inequality Constraints"
78 println(label)
79 println("Objective value: ", objective_value(model))
80 println("Solution: ")
81 println("th1Star = ", value(th1))
82 println("th2Star = ", value(th2))
83 println("th3Star = ", value(th3))
84 println(" ")
85 positions = linkPostions(value(th1), value(th2), value(th3))
86 plot_points(positions, label)
87 end
88
89 # Initialize the plot
90 fig = plot(; aspect_ratio=:equal, legend=false)
91 # Solve and plot without and with inequality constraints
92 solve_and_plot(false)
93 solve_and_plot(true)
94 # Display and save the plot
95 display(fig)
96 png(fig, "constrainedInequalityOptimizationFor3LinkRobotKinematicChain")

```

Output (Angles in radians)

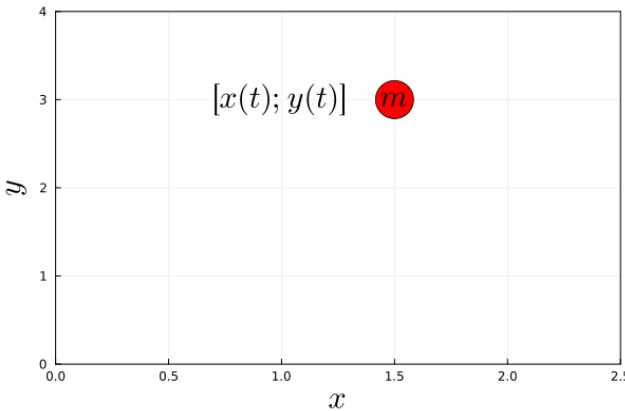
```

No Inequality Constraints
Objective value: 3.4536798144131877
Solution:
th1Star = -0.6590530918483652
th2Star = 1.3848362097307525
th3Star = 1.0495510986878442

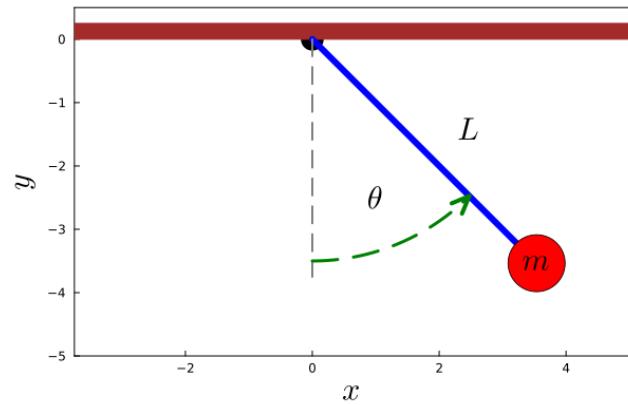
With Inequality Constraints
Objective value: 6.812490744951165
Solution:
th1Star = 0.05002085181111906
th2Star = 0.18815229215108328
th3Star = 2.6028037525509964

```

The configuration of the robot is shown Fig. 6.9.



(a)



(b)

Figure 6.10: (a) A point mass in the plane with its Cartesian coordinates. The origin is in the bottom left corner.(b) The classic pendulum formed by a massless rod with a point mass attached at its end. The angle θ is measured CCW from the vertical. The origin is placed at the pivot point attached to the “ceiling”.

6.4 Dynamics à la Lagrange

You know for sure that $F = m \cdot a$ for a point mass, as shown in Fig. 6.10-(a), and you’ve likely seen the equation for the pendulum of Fig. 6.10-(b),

$$\ddot{\theta}(t) = -\frac{g}{L} \cdot \sin(\theta). \quad (6.47)$$

But could you derive the equations for a person riding a Hoverboard? Or the 3-link manipulator in Fig. 1.6-(b)? At the end of this section, you will be able to do all of this and more. You’ll be able to derive by hand the equations for simple models such as the pendulum, and, via the `Symbolics` package, handle more complicated models.

Secrets of the Arcane 6.31: Fun Facts about Lagrange

Joseph-Louis Lagrange (1736–1813) did not have the advantage of computers or modern software to assist with computations, as we do today. In his time, mathematics was done by hand using pen and paper. Though some mathematicians had colleagues or students with whom they might discuss ideas or collaborate on problems, they didn’t have “assistants” in the sense of someone performing computational tasks for them.

Oh yeah, the Fun Facts: Lagrange was largely self-taught and did not obtain a university degree. Fascinated by maxima and minima of functions, he was the principal founder of the calculus of variations. In a far-reaching reformulation of Isaac Newton’s laws, Lagrange created a brilliant new vision of mechanics. We are working toward that now in this Section!

ChatGPT’s Remarks on the Fun Facts: The statement is mostly accurate, and here’s a deeper analysis of its points:

1. **Self-taught:** Lagrange was largely self-taught in his early mathematical education. He exhibited an innate interest in mathematics from a tender age and initiated his study autonomously.
2. **University Degree:** While it’s true that the nature of formal education in the 18th century was distinct from today, Lagrange did not follow the prototypical academic trajectory that we associate with eminent mathematicians today. Nevertheless, by the age of 19, he was already delivering significant papers to the mathematical domain. His pioneering contributions were acknowledged early, leading him to assume several esteemed academic roles throughout his life.
3. **Calculus of Variations:** Indeed, Lagrange made seminal contributions to the calculus of variations. He held a specific intrigue for the problem of the brachistochrone, which pertains to the quickest descent path under the influence of gravity. This inquiry propelled him to devise what is contemporarily referred to as the Euler-Lagrange equation, pivotal in the calculus of variations.
4. **Reformulation of Newton’s Laws:** In his magnum opus “Mécanique Analytique”, Lagrange made profound contributions to mechanics. Within, he presented a novel interpretation of classical mechanics rooted entirely in energy

principles, laying the groundwork for the Lagrangian mechanics we recognize today. This methodology deploys a singular scalar entity, termed the Lagrangian, to elucidate system dynamics, as opposed to the vectorial forces characteristic of Newtonian mechanics. Lagrange's conceptualization offers a more holistic and adaptable perspective, especially pertinent for intricate systems and scenarios.

In summation, while the provided statement offers a concise encapsulation of Lagrange's achievements and background, the life and contributions of such a distinguished individual warrant a more detailed exploration. Thank you, ChatGPT! 10 Oct. 2023.

6.4.1 The Dot Notation for Derivatives

In (6.47), $\ddot{\theta}(t)$ denotes the **angular acceleration** of the pendulum. The dot notation for derivatives is a very popular quasi-standard notation for denoting derivatives with respect to time. For example, if the vector

$$p(t) := \begin{bmatrix} x(t) \\ y(t) \end{bmatrix} \quad (6.48)$$

is the Cartesian position of the particle in the plane, then its Cartesian **velocity** is the first derivative with respect to time of the position vector, namely,

$$v(t) := \begin{bmatrix} \dot{x}(t) \\ \dot{y}(t) \end{bmatrix} := \frac{d}{dt} \begin{bmatrix} x(t) \\ y(t) \end{bmatrix} = \begin{bmatrix} \frac{dx(t)}{dt} \\ \frac{dy(t)}{dt} \end{bmatrix}, \quad (6.49)$$

read as x and y dot. In a similar manner, the derivative of the velocity vector with respect to time is the Cartesian **acceleration vector**,

$$a(t) := \begin{bmatrix} \ddot{x}(t) \\ \ddot{y}(t) \end{bmatrix} := \frac{d}{dt} \begin{bmatrix} \dot{x}(t) \\ \dot{y}(t) \end{bmatrix} = \begin{bmatrix} \frac{d^2x(t)}{dt^2} \\ \frac{d^2y(t)}{dt^2} \end{bmatrix} \quad (6.50)$$

read as x and y double dot.

For a point mass in the plane, as in Fig. 6.10-(a), with its Cartesian position, velocity, and acceleration given in (6.48) to (6.50), you may recall from a Physics course (and if not, please accept this for a fact), that when the only force acting on the mass is gravity, $g = 9.81 \frac{\text{m}}{\text{s}^2}$, Newton's laws give

$$\begin{bmatrix} m \cdot \ddot{x}(t) \\ m \cdot \ddot{y}(t) \end{bmatrix} = \begin{bmatrix} 0.0 \\ -m \cdot g \end{bmatrix}. \quad (6.51)$$

Equations like this are also called the **equations of motion** or EoM for short. The equations of motion for the **Cassie bipedal robot** are “somewhat” similar to the above, except they would fill 100 pages!

We will not need the third derivative with respect to time in this Section, but you might note that it is called the **jerk**. Car and elevator manufacturers are very careful about the “smoothness of acceleration” (keeping its rate of change continuous and not too large), so they use bounds on **jerk**, the third derivative of position, in the design of their products. When you see clothing⁵ with

“Don’t be a $\frac{d^3x(t)}{dt^3}$ ”,

you are now in the know.

6.4.2 Kinetic and Potential Energy of a Point Mass in the Plane

We define:

- **Potential Energy of a Point Mass** is equal to “mass times gravity times height”, or in symbols

$$V := m \cdot g \cdot y,$$

where we have deliberately dropped the explicit dependence on time; keep in mind, however, that the particle’s height does vary with time.

⁵Easily available on the web. Here, three dots would be too subtle.

- **Kinetic Energy of a Point Mass** by “one half mass times the magnitude (or norm) of the velocity squared”, or in symbols,

$$K := \frac{1}{2}m \cdot (\dot{x}^2 + \dot{y}^2), \quad (6.52)$$

where, again, we have deliberately dropped the explicit dependence on time; keep in mind, however, that the particle’s velocity does vary with time.

For later use, we note that both the potential and kinetic energies are scalar quantities. Their units are **Joules**. We will not make a big deal about units, though they are very important. **We do want to observe that the kinetic energy is quadratic in the velocity components of the particle.** In fact, the Cartesian velocity, $v(t)$, is a 2-vector (because we are working in the plane) and its kinetic energy is

$$K := \frac{1}{2}m \cdot \|v(t)\|^2.$$

We note this, so that in 3-space, you know how to compute the kinetic energy of a point mass! It is always the same formula: one-half the mass times the norm of the velocity squared.

This all looks very innocent and perhaps a touch boring. For an unconstrained particle in the plane, computing potential and kinetic energy is absolutely cake! Let’s consider, instead, the pendulum, which has a mass attached to a massless rod. We emphasize that the rod is massless because it then has no potential or kinetic energy. The potential and kinetic energies of the pendulum reside solely in the bob of mass $m > 0$.

Just as we did with the 3-link robot in Chapter 1.6.3, we use basic trigonometry to write down the Cartesian position of the bob, namely,

$$\begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} L \cdot \sin(\theta) \\ -L \cdot \cos(\theta) \end{bmatrix}. \quad (6.53)$$

The sines and cosines show up in different places because the angle θ for the pendulum is now referenced with respect to the vertical, with $\theta = 0$ corresponding to straight down, i.e., the pendulum at rest. We can now compute the potential energy for the pendulum,

$$V_{\text{pend}} = m \cdot g \cdot y = -m \cdot g \cdot L \cdot \cos(\theta). \quad (6.54)$$

The potential energy is at its lowest when the pendulum is hanging downwards and at its peak when the pendulum has rotated to $\theta = \pi$, straight up. Moreover, the potential energy is a nonlinear function of θ .

To determine the kinetic energy of the pendulum’s bob, we must differentiate the position, namely,

$$\begin{bmatrix} \dot{x} \\ \dot{y} \end{bmatrix} = \begin{bmatrix} \frac{d}{dt}(L \cdot \sin(\theta)) \\ \frac{d}{dt}(-L \cdot \cos(\theta)) \end{bmatrix} = \begin{bmatrix} L \cdot \cos(\theta) \cdot \dot{\theta} \\ L \cdot \sin(\theta) \cdot \dot{\theta} \end{bmatrix}, \quad (6.55)$$

where we used the chain rule. For the moment, let’s accept the velocity calculation as being correct, and compute the kinetic energy,

$$\begin{aligned} K_{\text{pend}} &= \frac{1}{2}m \cdot ((\dot{x})^2 + (\dot{y})^2) \\ &= \frac{1}{2}m \cdot ((L \cdot \cos(\theta) \cdot \dot{\theta})^2 + (L \cdot \sin(\theta) \cdot \dot{\theta})^2) \quad (\text{substituting in}) \\ &= \frac{1}{2}m \cdot L^2 \cdot (\cos^2(\theta) + \sin^2(\theta)) \cdot \dot{\theta}^2 \quad (\text{algebra}) \\ &= \frac{1}{2}m \cdot L^2 \cdot \dot{\theta}^2 \quad (\text{using the one trig identity everyone knows}). \end{aligned} \quad (6.56)$$

The kinetic energy is now quadratic function of the **angular velocity**, $\dot{\theta}$. The kinematics of the pendulum are specified by a single quantity, its angle θ , and the time rate of change of that quantity shows up in the kinetic energy.

The Chain Rule: We come back to computing \dot{x} and \dot{y} in (6.55).

- We need to differentiate $x(t) = L \cdot \sin(t)$ with respect to time, t .

Define $f : \mathbb{R} \rightarrow \mathbb{R}$, by $f(z) = \sin(z)$ and $g : \mathbb{R} \rightarrow \mathbb{R}$ by $g(t) = \theta(t)$. Then $x(t) = f(g(t))$; yes, z in $f(z)$ is substituted in by $g(t) = \theta(t)$. By the Chain Rule, using the prime notation for differentiation,

$$\dot{x}(t) = x'(t) = L \cdot (f(g(t)))' = L \cdot f'(g(t)) \cdot g'(t).$$

We compute the various terms,

$$\begin{aligned} f'(z) &= (\sin(z))' = \cos(z) \\ g'(t) &= (\theta(t))' = \frac{d}{dt}\theta(t) = \dot{\theta}(t). \end{aligned}$$

Putting the pieces together yields,

$$\begin{aligned} \dot{x}(t) &= L \cdot f'(g(t)) \cdot g'(t) \\ &= L \cdot \sin(g(t)) \cdot g'(t) \\ &= L \cdot \sin(\theta(t)) \cdot \dot{\theta}(t) \\ &= L \cdot \sin(\theta) \cdot \dot{\theta} \text{ after dropping the t's.} \end{aligned}$$

Professionals drop the t's because they know they are there and carrying them along for the ride in the derivations is “heavy”. You may not need the next one, but a minus sign does disappear, and it’s good to see how that happens.

- We need to differentiate $y(t) = -L \cdot \cos(t)$ with respect to time, t .

Define $f : \mathbb{R} \rightarrow \mathbb{R}$, by $f(z) = \cos(z)$ and $g : \mathbb{R} \rightarrow \mathbb{R}$ by $g(t) = \theta(t)$. Then $y(t) = f(g(t))$; again, z in $f(z)$ is substituted in by $g(t) = \theta(t)$. By the Chain Rule, using the prime notation for differentiation,

$$\dot{y}(t) = y'(t) = -L \cdot (f(g(t)))' = -L \cdot f'(g(t)) \cdot g'(t).$$

We compute the various terms,

$$\begin{aligned} f'(z) &= (\cos(z))' = -\sin(z) \\ g'(t) &= (\theta(t))' = \frac{d}{dt}\theta(t) = \dot{\theta}(t). \end{aligned}$$

Putting the pieces together yields,

$$\begin{aligned} \dot{y}(t) &= -L \cdot f'(g(t)) \cdot g'(t) \\ &= -L \cdot (-\sin(g(t))) \cdot g'(t) \\ &= L \cdot \sin(\theta(t)) \cdot \dot{\theta}(t) \text{ (double minus sign cancelled)} \\ &= L \cdot \sin(\theta) \cdot \dot{\theta} \text{ after dropping the t's.} \end{aligned}$$

In the next part, we’ll compute the potential and kinetic energies of the 3-link robot. While it can be done by hand, it’s so painful that we will approach it symbolically. This will put you one step closer to how real dynamical models of engineering systems are determined.

6.4.3 Kinetic and Potential Energy of the 3-Link Manipulator

Just as with the pendulum, we assume the links of the robot are massless⁶ except for point masses $m_i > 0$ at the ends of each link; the red balls in Fig. 1.6-(b) and Fig. 6.4.

The kinetic energy functions become quite complex quite quickly! Hence, as a warm-up exercise, we leave off the third link in the manipulator, giving us a 2-link robot. To set us up for **Lagrange’s Method** in the next section (yes, him again! what an amazing mathematician.), we’re going to use what are called **generalized coordinates**, where we replace

$$\begin{aligned} \theta_1 &\rightarrow q_1 \\ \theta_2 &\rightarrow q_2 \\ \theta_3 &\rightarrow q_3. \end{aligned}$$

Even for the 2-link version of the robot, the kinetic energy will not fit on one line! Fortunately for us Engineers, kinetic energy has some beautiful structure, which we highlight in the following result.

⁶A bad assumption in practice, but a good one for learning.

Proposition 6.32: Kinetic Energy is a Quadratic Form

The kinetic energy of any robot comprised of rigid links (e.g., our 3-link manipulator) can be expressed in the following quadratic form

$$K(q, \dot{q}) = \frac{1}{2} \dot{q}^\top D(q) \dot{q}, \quad (6.57)$$

where

$$\begin{bmatrix} q_1 \\ q_2 \\ \vdots \\ q_n \end{bmatrix} \text{ and } \begin{bmatrix} \dot{q}_1 \\ \dot{q}_2 \\ \vdots \\ \dot{q}_n \end{bmatrix}$$

are the generalized position and velocity coordinates, respectively, and

$$D(q) = \begin{bmatrix} d_{11}(q) & d_{12}(q) & \cdots & d_{1n}(q) \\ d_{21}(q) & d_{22}(q) & \cdots & d_{2n}(q) \\ \vdots & \vdots & \ddots & \vdots \\ d_{n1}(q) & d_{n2}(q) & \cdots & d_{nn}(q) \end{bmatrix} \quad (6.58)$$

is a **positive-definite^a** symmetric matrix.

$D(q)$ is called the **mass-inertia matrix**. Its determinant is always positive (meaning greater than zero), and hence $D(q)$ is always invertible.

^aAll eigenvalues are positive real numbers is the characterization we used in ROB 101 *Computational Linear Algebra*. The Appendix of the ROB 101 textbook has a deeper treatment of positive definite matrices, in case you are curious.

Warm-up Problems

Example 6.33. Determine the mass-inertia matrices for the systems in Fig. 6.10, that is, a point mass in the plane moving under the influence of gravity and a pendulum consisting of a point mass at the end of a massless bar.

Solutions:

For Fig. 6.10-(a), we defined the kinetic energy in (6.52). All we have to do is write it in the form (6.57). Doing so yields

$$\begin{aligned} K &= \frac{1}{2} m \cdot (\dot{x}^2 + \dot{y}^2) \\ &= \frac{1}{2} \begin{bmatrix} \dot{x} \\ \dot{y} \end{bmatrix}^\top \begin{bmatrix} m & 0 \\ 0 & m \end{bmatrix} \begin{bmatrix} \dot{x} \\ \dot{y} \end{bmatrix}. \end{aligned}$$

Hence,

$$D(x, y) = \begin{bmatrix} m & 0 \\ 0 & m \end{bmatrix},$$

a diagonal matrix with the mass on the diagonal. Not so interesting, except it is clearly symmetric and positive definite as long as $m > 0$!

For Fig. 6.10-(b), we computed in (6.56) that

$$KE_{\text{pend}} = \frac{1}{2} m \cdot L^2 \cdot \dot{\theta}^2.$$

Hence,

$$D(\theta) = m \cdot L^2,$$

a 1×1 matrix depending on mass and length. Not so interesting, except it is clearly symmetric and positive definite as long as $m > 0$ and $L > 0$.

Will the multi-link robots be more interesting? In fact, they will be so interesting that writing KE by hand in the form of (6.57) will be highly error-prone.

Method 6.34: Recovering the Mass-Inertia Matrix from the Kinetic Energy

Given an arbitrarily complicated (correct) expression $KE(q, \dot{q})$ for the Kinetic Energy of a rigid n -link mechanism, the mass-intertia matrix can be extracted via a double Jacobian, that is,

$$D(q) = \frac{\partial}{\partial \dot{q}} \left[\left(\frac{\partial}{\partial q} KE(q, \dot{q}) \right)^T \right].$$

Note: This follows from part (f) of Prop. 5.48 in Chapter 5.7.5, which is optional reading. Hence, **you are not expected to go check it out, unless you are really into the math.** We'll give you the required symbolic code to compute $D(q)$. When you see the double Jacobian, we wanted you to know what's up. Enough said.

Final Warm-Up Problem, a 2-Link Manipulator

Here is the Symbolic code for computing the energies of a two-link robot.

```

1 using Symbolics, SymPy, Latexify
2 set_default(fmt = "% .1f", convert_unicode = false)
3 # latexify(b) |> println
4
5 # This is for a two-link manipulator
6
7 # For the angles, see the schematic of the 3-link manipulator
8 # and remove the last link
9
10 @Symbolics.variables q1 q2 dq1 dq2 g m1 m2 L1 L2
11
12 # We have substituted theta1 and theta2 by q1 and q2
13 # Not only is this more compact, but it's also
14 # standard notation for Lagrange's Method for computing
15 # dynamical equations of motion
16
17 # Form vectors for use in Jacobians
18 q = [q1, q2]
19 dq = [dq1, dq2]
20
21 # Kinematics
22 p1 = L1 * [cos(q1), sin(q1)]
23 p2 = p1 + L2 * [cos(q1+q2), sin(q1+q2)]
24
25 # Potential Energy
26 PE = g*m1*p1[2] + g*m2*p2[2]
27 PE = simple(PE)
28 println("PE = ", PE)
29
30 # Velocity vectors via the Jacobian!
31 v1 = Symbolics.jacobian(p1, q) * dq
32 v2 = Symbolics.jacobian(p2, q) * dq

```

```

33
34 # Kinetic Energy of each mass on the Robot
35 KE1 = 0.5 * m1 * sum(v1.*v1)
36 KE2 = 0.5 * m2 * sum(v2.*v2)
37
38 # Sum the individual KEs to get the total Kinetic Energy
39 KE = KE1 + KE2
40
41 # Extract the mass-inertia matrix from KE
42 if false
43     D = Symbolics.jacobian(Symbolics.jacobian([KE], dq), dq)
44 else
45     D = Symbolics.hessian(KE, dq)
46 end
47
48 D = simple(D)
49 # Choose display method for the mass-inertia matrix
50 if false
51     print("D = ")
52     display(D)
53 else
54     n= size(D, 1)
55     for i = 1:n
56         for j = i:n
57             print("\n \n D[$i,$j] = ", D[i,j])
58         end
59     end
60 end

```

Output After Simplification Using a Custom Function Developed in Collaboration with ChatGPT4, 09 to 14 Oct 2023. More on this later.

PE = g*(L1*m1*sin(q1) + m2*(L1*sin(q1) + L2*sin(q1 + q2)))

D[1,1] = L1^2*m1 + m2*(L1^2 + 2*L1*L2*cos(q2) + L2^2)

D[1,2] = L2*m2*(L1*cos(q2) + L2)

D[2,1] = L2^2*m2

D[2,2] = L2^2*m2

Typeset Mass-intertia Matrix

$$D(q) = \begin{bmatrix} L_1^2 \cdot m_1 + m_2 \cdot (L_1^2 + 2 \cdot L_1 \cdot L_2 \cdot \cos(q_2) + L_2^2) & L_2 \cdot m_2 \cdot (L_1 \cdot \cos(q_2) + L_2) \\ L_2 \cdot m_2 \cdot (L_1 \cdot \cos(q_2) + L_2) & L_2^2 \cdot m_2 \end{bmatrix}; \quad (6.59)$$

this is a more interesting mass-inertia matrix! We are done with the warmup.

3-Link Manipulator: The Real Deal

```

1 using Symbolics, SymPy, Latexify
2 set_default(fmt = "% .1f", convert_unicode = false)
3 # latexify(b) |> println
4
5 # This is for a three-link manipulator

```

```

6
7 # For the angles, see the schematic of the 3-link manipulator
8 # and remove the last link
9
10 @Symbolics.variables q1 q2 q3 dq1 dq2 dq3 g m1 m2 m3 L1 L2 L3
11
12 # We have substituted theta_i by q_i
13 # Not only is this more compact, but it's also
14 # standard notation for Lagrange's Method for computing
15 # dynamical equations of motion
16
17 # Form vectors for use in Jacobians
18 q = [q1, q2, q3]
19 dq = [dq1, dq2, dq3]
20
21 # Kinematics
22 p1 = L1*[cos(q1), sin(q1)]
23 p2 = p1 + L2*[cos(q1+q2), sin(q1+q2)]
24 p3 = p2 + L3*[cos(q1+q2+q3), sin(q1+q2+q3)]
25
26 # Potential Energy
27 PE = g*m1*p1[2] + g*m2*p2[2] + g*m3*p3[2]
28 PE = simple(PE)
29 println("PE = ", PE)
30
31 # Velocity vectors via the Jacobian!
32 v1 = Symbolics.jacobian(p1, q) * dq
33 v2 = Symbolics.jacobian(p2, q) * dq
34 v3 = Symbolics.jacobian(p3, q) * dq
35
36 # Kinetic Energy of each mass on the Robot
37 KE1 = 0.5 * m1 * sum(v1.*v1)
38 KE2 = 0.5 * m2 * sum(v2.*v2)
39 KE3 = 0.5 * m3 * sum(v3.*v3)
40
41 # Sum the individual KEs to get the total Kinetic Energy
42 KE = KE1 + KE2 + KE3
43 KE = Symbolics.simplify(KE)
44
45 # Extract the mass-inertia matrix from KE
46 if false
47     D = Symbolics.jacobian(Symbolics.jacobian([KE], dq), dq)
48 else
49     D = Symbolics.hessian(KE, dq)
50 end
51
52 # Do simplification on D with a custom function
53 D = simple(D)
54
55 # Choose display method for the mass-inertia matrix
56 if false
57     print("D = ")
58     display(D)
59 else
60     n= size(D, 1)
61     for i = 1:n
62         for j = i:n

```

```

63     println("\n D[$i,$j] = ", D[i,j])
64     #display(D[1,j])
65 end
66 end
67 end

```

Output

```

PE = g*(L1*m1*sin(q1) + m2*(L1*sin(q1) + L2*sin(q1 + q2)) + m3*(L1*sin(q1) +
    L2*sin(q1 + q2) + L3*sin(q1 + q2 + q3)))

D[1,1] = L1^2*m1 + m2*(L1^2 + 2*L1*L2*cos(q2) + L2^2) + m3*(L1^2 + 2*L1*L2*cos(q2) +
    2*L1*L3*cos(q2 + q3) + L2^2 + 2*L2*L3*cos(q3) + L3^2)

D[1,2] = L2*m2*(L1*cos(q2) + L2) + m3*(L1*L2*cos(q2) + L1*L3*cos(q2 + q3) + L2^2 +
    2*L2*L3*cos(q3) + L3^2)

D[1,3] = L3*m3*(L1*cos(q2 + q3) + L2*cos(q3) + L3)

D[2,2] = L2^2*m2 + m3*(L2^2 + 2*L2*L3*cos(q3) + L3^2)

D[2,3] = L3*m3*(L2*cos(q3) + L3)

D[3,3] = L3^2*m3

```

The typeset terms are not that much more edifying, so we skip that exercise. You see enough to know that you do not want to derive the kinetic energy by hand. Back in the day, your author needed similar equations for a 3-link bipedal walker (diagram on page 33) (see also Fig. 6.12), and he derived everything by hand. It was painful.

6.4.4 Remarks on Simplifying Symbolic Equations

We've been using the `Symbolics` package in Julia. It works quite well and has native Julia syntax. Its main weakness lies in its simplification of trigonometric expressions: if you really want to combine all the trig terms, it comes up short. The competing package, `Sympy`, is based in Python and has a "Julia wrapper" so that it can be used in Julia. Its syntax is Python-centric, which is good for some folks and awkward for others. A strength is that its simplification command is more powerful: if there is a way to combine some gnarly trig expressions to produce a shorter output string, it will likely find it.

Currently, it is hard to go back and forth between the two packages. ChatGPT4 helped your author build a simplification tool that translates a `Symbolic` expression to a `Sympy` expression for simplification. We give it below.

Warning: Once the function below has been applied to a `Symbolic` expression, you cannot perform additional operations on it. Hence, only use it at the very end, when you are "printing out" your functions. For intermediate simplification of expressions, use `Symbolics.simplify`.

```

1 # Updated 14 October
2 # Working with ChatGPT4
3 # 9 through 14 Oct 2023
4
5 using SymPy, Symbolics
6
7 function simple(expr)
8     # If the expression is a matrix, handle each element individually
9     if expr isa Matrix
10        rows, cols = size(expr)
11        new_mat = Matrix{Any}(undef, rows, cols)
12        for i in 1:rows

```

```

13     for j in 1:cols
14         new_mat[i, j] = simple(expr[i, j])
15     end
16 end
17 return new_mat
18
19
20 # Convert Symbolics.jl expression to string
21 expr_str = string(expr)
22
23 # Replace implicit multiplication (like "0.5(") with explicit multiplication
24 ("0.5*(")
25 # expr_str = replace(expr_str, r"([\d\.]+)\(" => s"\1*\(")
26 expr_str = replace(expr_str, r"([\d\.]+)([a-zA-Z\()]" => s"\1*\2")
27
28 # Remove unwanted artifacts (e.g., "Num")
29 expr_str = replace(expr_str, "Num[" => "[")
30 expr_str = replace(expr_str, "Symbol (" => "")
31 expr_str = replace(expr_str, "Integer (" => "")
32 # Replace double divide
33 expr_str = replace(expr_str, "//" => "/")
34
35 # Convert the modified string to SymPy.jl expression
36 sympy_expr = SymPy.sympify(expr_str)
37
38 # Use SymPy to simplify the expression
39 simplified_expr = SymPy.simplify(sympy_expr)
40
41 # Convert simplified expression back to string
42 simplified_str = string(simplified_expr)
43
44 # Remove 1.0* from the simplified expression string
45 final = replace(simplified_str, r"(?<!\\d)1\\.0\\*(?!\\d\\.)" => "")
46
47 return final
end

```

Output

```
simple (generic function with 1 method)
```

6.4.5 Lagrange's Equations

Lagrange's equation is a profound and elegant reformulation of classical mechanics. Instead of navigating the intricacies of vector forces and torques in Newtonian mechanics, Lagrange's equation offers a scalar approach, focusing on the potential and kinetic energy of moving bodies. Working with the concept of generalized coordinates, it provides a systematic method to derive equations of motion for complex mechanical systems, even those with constraints. This approach not only simplifies the mathematical treatment of complicated systems but also offers deeper insights into the conservation laws and symmetries inherent in physical systems. In essence, Lagrange's equation bridges the gap between the microscopic laws of forces as encapsulated in Newton's “ $F = ma$ ” and the macroscopic behavior of systems, making it an indispensable tool for physicists and engineers alike.

Consider a mechanical system with generalized position coordinates, $q = (q_1, q_2, \dots, q_n)$, and generalized velocities, $\dot{q} = (\dot{q}_1, \dot{q}_2, \dots, \dot{q}_n)$; both q and \dot{q} are column vectors. For the problems we will study, there is one generalized coordinate for each independent link or moving part in our robots.

Lagrange's Equations of Motion

Let $V(q)$ and $K(q, \dot{q})$ be the potential and kinetic energies, respectively, of a mechanical system.

Definition 6.35. *The Lagrangian is defined as*

$$\mathcal{L}(q, \dot{q}) := K(q, \dot{q}) - V(q). \quad (6.60)$$

Lagrange's Equations of Motion (EoM) are

$$\frac{d}{dt} \frac{\partial \mathcal{L}(q, \dot{q})}{\partial \dot{q}} - \frac{\partial \mathcal{L}(q, \dot{q})}{\partial q} = \Gamma, \quad (6.61)$$

where the partial derivatives are column vectors and Γ is a column vector of external forces acting on the system. For us, they will primarily be motor torques.

Note that we are using the traditional convention (meaning nearly universal) in Lagrange's equations of motion that the partial derivatives are organized into a column vector.

Note: Written in the compact form (6.61), Lagrange's equations of motion are at once awe-inspiring due to their elegance and terrifying due to the many complicated-looking symbols. We will unpack (6.61) little by little to make it digestible, and then we'll turn loose the `Symbolics` package on it and compute some real equations of motion! After that, we'll deal with Γ . One thing at a time!

As a first step, we unpack the partial derivatives, just as we did with the Method of Lagrange Multipliers, giving us

$$\begin{bmatrix} \frac{d}{dt} \frac{\partial \mathcal{L}(q, \dot{q})}{\partial \dot{q}_1} \\ \frac{d}{dt} \frac{\partial \mathcal{L}(q, \dot{q})}{\partial \dot{q}_2} \\ \vdots \\ \frac{d}{dt} \frac{\partial \mathcal{L}(q, \dot{q})}{\partial \dot{q}_n} \end{bmatrix} - \begin{bmatrix} \frac{\partial \mathcal{L}(q, \dot{q})}{\partial q_1} \\ \frac{\partial \mathcal{L}(q, \dot{q})}{\partial q_2} \\ \vdots \\ \frac{\partial \mathcal{L}(q, \dot{q})}{\partial q_n} \end{bmatrix} = \begin{bmatrix} \Gamma_1 \\ \Gamma_2 \\ \vdots \\ \Gamma_n \end{bmatrix}. \quad (6.62)$$

While that $\frac{d}{dt}$ looks like it could cause complications, it is simply the total derivative covered in Chapter 5.8. You may recall that it means we need to use a vector version of the Chain Rule.

Exercise 6.36. Apply Lagrange's equations of motion to the unconstrained point mass in the plane, shown in Fig. 6.10-(a).

Solution: From Chapter 6.4.2, we have

$$\begin{aligned} V(x, y) &= m \cdot g \cdot y \\ K(x, y, \dot{x}, \dot{y}) &= \frac{1}{2}m(\dot{x}^2 + \dot{y}^2). \end{aligned}$$

Therefore, the Lagrangian is

$$\mathcal{L}(x, y, \dot{x}, \dot{y}) = \frac{1}{2}m(\dot{x}^2 + \dot{y}^2) - m \cdot g \cdot y,$$

which is not so intimidating.

In Lagrange's equations, we have the two column vectors $q = (x, y)$ and $\dot{q} = (\dot{x}, \dot{y})$. Let's compute some partial derivatives:

$$\frac{\partial \mathcal{L}(q, \dot{q})}{\partial q_1} = \frac{\partial}{\partial x} \left(\frac{1}{2}m(\dot{x}^2 + \dot{y}^2) - m \cdot g \cdot y \right) = 0 \quad (x \text{ does not appear in the Lagrangian})$$

$$\frac{\partial \mathcal{L}(q, \dot{q})}{\partial q_2} = \frac{\partial}{\partial y} \left(\frac{1}{2}m(\dot{x}^2 + \dot{y}^2) - m \cdot g \cdot y \right) = -m \cdot g \quad (y \text{ appears linearly in the Lagrangian}).$$

That was not so hard. Continuing with the partial derivatives,

$$\frac{\partial \mathcal{L}(q, \dot{q})}{\partial \dot{q}_1} = \frac{\partial \mathcal{L}(x, y, \dot{x}, \dot{y})}{\partial \dot{x}} = m \cdot \dot{x} \quad (\dot{x} \text{ appears quadratically in the Lagrangian, and there's a one half})$$

$$\frac{\partial \mathcal{L}(q, \dot{q})}{\partial \dot{q}_2} = \frac{\partial \mathcal{L}(x, y, \dot{x}, \dot{y})}{\partial \dot{y}} = m \cdot \dot{y} \quad (\dot{y} \text{ appears quadratically in the Lagrangian, and there's a one half}).$$

It remains to handle $\frac{d}{dt}$; continuing

$$\frac{d}{dt} \frac{\partial \mathcal{L}(q, \dot{q})}{\partial \dot{q}_1} = \frac{d}{dt} (m \cdot \dot{x}) = m \cdot \ddot{x}$$

$$\frac{d}{dt} \frac{\partial \mathcal{L}(q, \dot{q})}{\partial \dot{q}_2} = \frac{d}{dt} (m \cdot \dot{y}) = m \cdot \ddot{y}.$$

Putting everything together, Lagrange's equations give

$$\underbrace{\begin{bmatrix} m \cdot \ddot{x} \\ m \cdot \ddot{y} \end{bmatrix}}_{\frac{d}{dt} \frac{\partial \mathcal{L}(q, \dot{q})}{\partial \dot{q}}} - \underbrace{\begin{bmatrix} 0 \\ -mg \end{bmatrix}}_{\frac{\partial \mathcal{L}(q, \dot{q})}{\partial q}} = \underbrace{\begin{bmatrix} 0 \\ 0 \end{bmatrix}}_{\Gamma}, \quad (6.63)$$

where $\Gamma = 0_{2 \times 1}$ because there are no external forces acting on the particle. How's that? Gravity is in the picture. Yes, but gravity and springs are (almost) always included in the potential energy term. External forces, in our case, would mean a jet engine attached to the particle because it is free-floating in the plane.

Bottom Line: For this simple system, Lagrange's Equations of Motion, (EoM) for short, are overkill. However, we certainly recovered the same model we had previously. ■.

Exercise 6.37. Apply Lagrange's equations of motion to the constrained point mass in the plane, shown in Fig. 6.10-(b), otherwise known as a pendulum.

Solution: From Chapter 6.4.2, we have

$$V(\theta) = -m \cdot g \cdot L \cdot \cos(\theta)$$

$$K(\theta, \dot{\theta}) = \frac{1}{2} m \cdot L^2 \cdot \dot{\theta}^2.$$

Therefore, the Lagrangian is

$$\mathcal{L}(\theta, \dot{\theta}) = \frac{1}{2} m \cdot L^2 \cdot \dot{\theta}^2 + m \cdot g \cdot L \cdot \cos(\theta),$$

where we replaced the double minus signs with a single plus sign.

In Lagrange's equations, we have scalar generalized coordinates $q = \theta$ and $\dot{q} = \dot{\theta}$. Let's compute some partial derivatives:

$$\frac{\partial \mathcal{L}(q, \dot{q})}{\partial q} = \frac{\partial}{\partial \theta} \left(\frac{1}{2} m \cdot L^2 \cdot \dot{\theta}^2 + m \cdot g \cdot L \cdot \cos(\theta) \right) = -m \cdot g \cdot L \cdot \sin(\theta).$$

Continuing with the partial derivatives,

$$\frac{\partial \mathcal{L}(q, \dot{q})}{\partial \dot{q}} = \frac{\partial}{\partial \dot{\theta}} \left(\frac{1}{2} m \cdot L^2 \cdot \dot{\theta}^2 + m \cdot g \cdot L \cdot \cos(\theta) \right) = m \cdot L^2 \cdot \dot{\theta}.$$

It remains to handle $\frac{d}{dt}$; continuing

$$\frac{d}{dt} \frac{\partial \mathcal{L}(q, \dot{q})}{\partial \dot{q}_1} = \frac{d}{dt} (m \cdot L^2 \cdot \dot{\theta}) = m \cdot L^2 \cdot \ddot{\theta}$$

because the mass and length are constant. If the bob were a bucket of water, with water splashing over the sides as it oscillated, we would have to compute \dot{m} ; but not in our case. When you "pump your legs and torso on a swing", you are moving the center of

mass of your body up and down, so we'd then have to account for $L(t)$ and $\dot{L}(t)$; but not in our case. Do you see how general and powerful Lagrange's method is?

Putting everything together, Lagrange's equations give

$$\underbrace{\frac{d}{dt} \frac{\partial \mathcal{L}(q, \dot{q})}{\partial \dot{q}}}_{m \cdot L^2 \cdot \ddot{\theta}} - \underbrace{(-m \cdot g \cdot L \cdot \sin(\theta))}_{\frac{\partial \mathcal{L}(q, \dot{q})}{\partial q}} = \underbrace{\Gamma}_{0}, \quad (6.64)$$

where $\Gamma = 0_{1 \times 1}$ because we have not yet installed a motor at the pivot. With a tiny bit of algebra, we recover the equations we had earlier in (6.47).

Bottom Line: Already, for this simple system, Lagrange's EoM are useful. How do you think Lagrange's EoM will be for the 3-link manipulator? Awesome, of course. ■.

The following result will help us to interpret the output of Lagrange's EoM for the 3-link manipulator and, in fact, for all models of connected rigid bodies.

Proposition 6.38: The Robot Equations

When Lagrange's EoM in (6.61) are expanded out for a standard mechanical system (i.e., most robots), they take the form

$$(\text{"the Robot Equations"}) \quad D(q) \cdot \ddot{q} + C(q, \dot{q}) \cdot \dot{q} + G(q) = \Gamma, \quad (6.65)$$

where,

- $D(q)$ is the mass-inertia matrix from the kinetic energy and is therefore invertible,
- $G(q) := \nabla V(q)$, is the gradient of the potential energy,
- and $C(q, \dot{q}) \cdot \dot{q}$ contains all the remaining terms in (6.61).

Note: If you take advanced courses in dynamics^a, you will learn that the $n \times n$ matrix $C(q, \dot{q})$ has very important properties dealing with **centrifugal and centripetal forces**, a topic that is way beyond our needs and our abilities in this course. For now, we suggest you **Don't Worry Bout It!**^b.

^aFor example, at Michigan, AERO 343 Spacecraft Dynamics and ME 440 Intermediate Dynamics and Vibrations.

^bCredit: Kings and ChatGPT.

Exercise 6.39. Identify the terms in the robot equations for the unconstrained point mass in the plane and the pendulum.

Solution:

The unconstrained point mass in the plane has equations of motion,

$$\begin{bmatrix} \ddot{x} \\ \ddot{y} \end{bmatrix} - \begin{bmatrix} 0 \\ -mg \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}.$$

Hence,

$$q = \begin{bmatrix} x \\ y \end{bmatrix}, \quad D(q) = I_2, \quad C(q, \dot{q}) = 0_{2 \times 2}, \quad G(q) = \begin{bmatrix} 0 \\ -mg \end{bmatrix}, \quad \text{and } \Gamma = \begin{bmatrix} 0 \\ 0 \end{bmatrix}.$$

The pendulum has equations of motion,

$$m \cdot L^2 \cdot \ddot{\theta} + (m \cdot g \cdot L \cdot \sin(\theta)) = 0.$$

Hence,

$$q = \theta, \quad D(q) = m L^2, \quad C(q, \dot{q}) = 0_{1 \times 1}, \quad G(q) = m \cdot g \cdot L \cdot \sin(\theta), \quad \text{and } \Gamma = 0.$$

Remark 6.40. If you find the Symbolic computations of the next section impressive, please read [Secrets of the Arcane 6.41](#). Today, robot models can be extracted directly from the CAD files used to design the robot in a **Unified Robot Description Format**, or (URDF) for short. You can think of a URDF as a written blueprint for a robot: “Link-1 is connected to Link-2 via a revolute joint. Link-2 is connected to Link-3 via a prismatic joint.” Moreover, the URDF contains all of the mechanical parameters of each link, such as length, mass, center of mass, moments of inertia about various points on the link, etc. It is all coming together so that you can spend your time on creative electromechanical design, perception systems and signal processing, planning and mapping, feedback control design, [human-robot interaction \(HRI\)](#), etc. It’s a great time to be a Roboticist!

Videos on the Calculus of Variations and Lagrange’s Equation: While this is definitely advanced material, it may inspire future course selections in mathematics and physics.

- [Introduction to Variational Calculus - Deriving the Euler-Lagrange Equation](#) by Good Vibrations with Freeball.
- [Derivation of the Euler-Lagrange Equation | Calculus of Variations](#) by Faculty of Khan.
- [Derivation of the Euler-Lagrange Equation](#) by The Caribbean Bookworm.

6.4.6 Lagrange’s Equations of Motion for the 3-Link Manipulator

```

1 using Symbolics, SymPy, Latexify
2 set_default(fmt = "% .1f", convert_unicode = false)
3 # latexify(b) |> println
4
5 # This is for a three link manipulator
6
7 # For the angles, see the schematic of the 3-link manipulator
8 # and remove the last link
9
10 @Symbolics.variables q1 q2 q3 dq1 dq2 dq2 dq3 g m1 m2 m3 L1 L2 L3
11
12 # We have substituted theta_i by q_i
13 # Not only is this more compact, but it's also
14 # standard notation for Lagrange's Method for computing
15 # dynamical equations of motion
16
17 # Form vectors for use in Jacobians
18 q = [q1, q2, q3]
19 dq = [dq1, dq2, dq3]
20
21 # Kinematics
22 p1 = L1*[cos(q1), sin(q1)]
23 p2 = p1 + L2*[cos(q1+q2), sin(q1+q2)]
24 p3 = p2 + L3*[cos(q1+q2+q3), sin(q1+q2+q3)]
25
26 # Potential Energy
27 PE = g*m1*p1[2] + g*m2*p2[2] + g*m3*p3[2]
28
29 # Use native Symbolics simplification
30 PE = Symbolics.simplify(PE)
31
32 # Velocity vectors via the Jacobian!
33 v1 = Symbolics.jacobian(p1, q) * dq
34 v2 = Symbolics.jacobian(p2, q) * dq
35 v3 = Symbolics.jacobian(p3, q) * dq
36
37 # Kinetic Energy of each mass on the Robot

```

```

38 KE1 = 0.5 * m1 * sum(v1.*v1)
39 KE2 = 0.5 * m2 * sum(v2.*v2)
40 KE3 = 0.5 * m3 * sum(v3.*v3)
41
42 # Sum the individual KEs to get the total Kinetic Energy
43 KE = KE1 + KE2 + KE3
44
45 # Use native Symbolics simplification
46 KE = Symbolics.simplify(KE)
47
48 # Form the Lagrangian
49 L = KE - PE
50
51 # Compute terms in the Robot Equations
52 G = Symbolics.gradient(PE, q)
53 D = Symbolics.hessian(KE, dq) # Strange name explained below
54
55 # Compute C (advanced, uses Christoffel symbols)
56 n = length(q)
57 C = Array{Num}(undef, n, n) # Num is Symbolic Type
58 # Strange name explained below
59 for k in 1:n
60     for j in 1:n
61         C[k, j] = 0
62         for i in 1:n
63             C[k, j] += 0.5 * (Symbolics.gradient(D[k, j], q)[i] +
64             Symbolics.gradient(D[k, i], q)[j] - Symbolics.gradient(D[i, j], q)[k]) * dq[i]
65         end
66     end
67 end
68
69 # Place an actuator at the hip on
70 # the relative angle between the torso and
71 # each leg
72 # Use D'Alembert's Principle
73 E = [q1; q2; q3]
74 B = Symbolics.jacobian(E, q)
75 B = B'
76 JacG = Symbolics.jacobian(G, q) # Used when we discuss lienearization of NL models
77
78
79 # Only run simple after all Symbolic manipulation is done and you
80 # are ready to create a function
81 if true
82     # Simplify as much as possible at the very end
83     D = simple(D)
84     C = simple(C)
85 end
86
87 # Create a function that computes all terms
88 # in the Robot Equations
89 fcn_name = "dyn_mod_3LinkManipulator"
90 modelParamString = "g, L1, L2, L3, m1, m2, m3 = modelParameters()"
91 # Need to include line break command \n
92 variableNamesString = "q1, q2, q3 = q \ndq1, dq2, dq3 = dq"
93 writeEOM(fcn_name, D, C, G, B, JacG, modelParamString, variableNamesString)
94

```

```

95 # bring the model into the workspace for later use
96 include("dyn_mod_3LinkManipulator.jl")
97
98 # If the line above fails, it is because Julia is trying to read the file before it has
99 # been made available to it by the file system. Simply copy the command into a new cell
100 # and
101 # run it. The time it takes you to do that will allow the file system to catch up! :-

```

Output

```
[creating DYN_MOD_3LINKMANIPULATOR.jl]
File dyn_mod_3LinkManipulator.jl created successfully!
```

```
dyn_mod_3LinkManipulator (generic function with 1 method)
```

When you open the file, you have a function that computes all the terms in Lagrange's EOM, formatted in a more-or-less human-readable form.

```

1 function dyn_mod_3LinkManipulator(q, dq)
2 # DYN_MOD_3LINKMANIPULATOR
3 # 2023-12-03 08:48:06
4 #
5 # Author: Grizzle
6 #
7 # Model NOTATION: D(q)ddq + C(q,dq)*dq + G(q) = B*tau
8 # The Robot Equations: From Lagrange's Equations of Motion
9 #
10 g, L1, L2, L3, m1, m2, m3 = modelParameters()
11 #
12 # Variable names for the model
13 q1, q2, q3 = q
14 dq1, dq2, dq3 = dq
15 #
16 D = zeros(3, 3)
17 D[1, 1] = L1^2*m1 + m2*(L1^2 + 2*L1*L2*cos(q2) + L2^2) + m3*(L1^2 + 2*L1*
18 L2*cos(q2) + 2*L1*L3*cos(q2 + q3) + L2^2 + 2*L2*L3*cos(q3) +
19 L3^2)
20 D[1, 2] = L2*m2*(L1*cos(q2) + L2) + m3*(L1*L2*cos(q2) + L1*L3*cos(q2 +
21 q3) + L2^2 + 2*L2*L3*cos(q3) + L3^2)
22 D[1, 3] = L3*m3*(L1*cos(q2 + q3) + L2*cos(q3) + L3)
23 D[2, 1] = L2*m2*(L1*cos(q2) + L2) + m3*(L1*L2*cos(q2) + L1*L3*cos(q2 +
24 q3) + L2^2 + 2*L2*L3*cos(q3) + L3^2)
25 D[2, 2] = L2^2*m2 + m3*(L2^2 + 2*L2*L3*cos(q3) + L3^2)
26 D[2, 3] = L3*m3*(L2*cos(q3) + L3)
27 D[3, 1] = L3*m3*(L1*cos(q2 + q3) + L2*cos(q3) + L3)
28 D[3, 2] = L3*m3*(L2*cos(q3) + L3)
29 D[3, 3] = L3^2*m3
30 #
31 C = zeros(3, 3)
32 C[1, 1] = -L1*L2*dq2*m2*sin(q2) - L1*L2*dq2*m3*sin(q2) - L1*L3*dq2*m3*
33 sin(q2 + q3) - L1*L3*dq3*m3*sin(q2 + q3) - L2*L3*dq3*m3*sin(q3)
34 C[1, 2] = -L1*L2*dq1*m2*sin(q2) - L1*L2*dq1*m3*sin(q2) - L1*L2*dq2*m2*
35 sin(q2) - L1*L2*dq2*m3*sin(q2) - L1*L3*dq1*m3*sin(q2 + q3) - L1*
36 L3*dq2*m3*sin(q2 + q3) - L1*L3*dq3*m3*sin(q2 + q3) - L2*L3*dq3*
37 m3*sin(q3)
38 C[1, 3] = -L3*m3*(L1*sin(q2 + q3) + L2*sin(q3))*(dq1 + dq2 + dq3)

```

```

39 C[2, 1] = L1*L2*dq1*m2*sin(q2) + L1*L2*dq1*m3*sin(q2) + L1*L3*dq1*m3*
40           sin(q2 + q3) - L2*L3*dq3*m3*sin(q3)
41 C[2, 2] = -L2*L3*dq3*m3*sin(q3)
42 C[2, 3] = -L2*L3*m3*(dq1 + dq2 + dq3)*sin(q3)
43 C[3, 1] = L3*m3*(L1*dq1*sin(q2 + q3) + L2*dq1*sin(q3) + L2*dq2*sin(q3))
44 C[3, 2] = L2*L3*m3*(dq1 + dq2)*sin(q3)
45 C[3, 3] = 0.0
46 #
47 G = zeros(3)
48 G[1] = g*m2*(L1*cos(q1) + L2*cos(q1 + q2)) + g*m3*(L1*cos(q1) + L2*
49           cos(q1 + q2) + L3*cos(q1 + q2 + q3)) + L1*g*m1*cos(q1)
50 G[2] = g*m3*(L2*cos(q1 + q2) + L3*cos(q1 + q2 + q3)) + L2*g*m2*cos(q1 +
51           q2)
52 G[3] = L3*g*m3*cos(q1 + q2 + q3)
53 #
54 B = zeros(3, 3)
55 B[1, 1] = 1
56 B[2, 2] = 1
57 B[3, 3] = 1
58 #
59 JacG = zeros(3, 3)
60 JacG[1, 1] = g*m2*(-L1*sin(q1) - L2*sin(q1 + q2)) + g*m3*(-L1*sin(q1) - L2*
61           sin(q1 + q2) - L3*sin(q1 + q2 + q3)) - L1*g*m1*sin(q1)
62 JacG[1, 2] = g*m3*(-L2*sin(q1 + q2) - L3*sin(q1 + q2 + q3)) - L2*g*m2*sin(q1 +
63           q2)
64 JacG[1, 3] = -L3*g*m3*sin(q1 + q2 + q3)
65 JacG[2, 1] = g*m3*(-L2*sin(q1 + q2) - L3*sin(q1 + q2 + q3)) - L2*g*m2*sin(q1 +
66           q2)
67 JacG[2, 2] = g*m3*(-L2*sin(q1 + q2) - L3*sin(q1 + q2 + q3)) - L2*g*m2*sin(q1 +
68           q2)
69 JacG[2, 3] = -L3*g*m3*sin(q1 + q2 + q3)
70 JacG[3, 1] = -L3*g*m3*sin(q1 + q2 + q3)
71 JacG[3, 2] = -L3*g*m3*sin(q1 + q2 + q3)
72 JacG[3, 3] = -L3*g*m3*sin(q1 + q2 + q3)
73 #
74 return (D=D, C=C, G=G, B=B, JacG=JacG)
75 end

```

This brings the model into the workspace.

```
include("dyn_mod_3LinkManipulator.jl") # bring the model into the workspace
```

Output

```
dyn_mod_3LinkManipulator (generic function with 1 method)
```

6.5 (Optional Read:) More on Lagrangian Dynamics and Mechanical Properties of Planar Bodies

6.5.1 A Crude Method to Numerically Integrate the Robot Equations

This section describes a crude but straightforward approach to “integrating” the **robot dynamics equations**. Since ordinary differential equations (ODEs) have not yet been introduced in the course, we use a **basic integration method** combining the **Riemann sum** and the **trapezoidal rule**. The included code stub assumes a 2-link manipulator.

Problem Setup

We consider the robot equations without input, namely

$$D(q)\ddot{q} + C(q, \dot{q})\dot{q} + G(q) = 0_{n \times 1},$$

where, as before,

- q typically represents joint angles,
- \dot{q} typically represents angular velocities, and
- \ddot{q} represents the angular accelerations.

Moreover,

- $D(q)$ is the $n \times n$ **inertia matrix**,
- $C(q, \dot{q})$ represents the $n \times n$ **Coriolis and centrifugal effects**,
- $G(q)$ is the $n \times 1$ **gravity vector**.

To be extra clear, we are not including motor torques at the joints.

Angular Acceleration

For later use, the angular acceleration can be solved for and used to update the velocity term,

$$\ddot{q} = D(q) \backslash (-C(q, \dot{q})\dot{q} - G(q)).$$

Note that we always choose to avoid matrix inverses whenever possible.

Numerical Integration Approach

We solve this system over a finite time interval T using a **discrete time step**:

$$t = \{0, \Delta t, 2\Delta t, \dots, T\}$$

where $\Delta t = \frac{T}{N}$ is the time step, and N is the total number of steps.

The numerical integration proceeds as follows:

1. Notation for the code

- $q[:, i]$ is q (angular position) at time $t[i]$.
- $dq[:, i]$ is \dot{q} (angular velocity) at time $t[i]$.
- $ddq[:, i]$ is \ddot{q} (angular acceleration) at time $t[i]$.

2. Compute angular acceleration \ddot{q} at the current time step using:

$$ddq[:, i] = D(q[:, i]) \backslash (-C(q[:, i], dq[:, i]) \cdot dq[:, i] - G(q[:, i])).$$

3. Update velocity \dot{q} using a (left) **Riemann Sum**:

$$dq[:, i+1] = dq[:, i] + \Delta t \cdot ddq[:, i]$$

4. Update position q using the **Trapezoidal Rule**:

$$q[:, i+1] = q[:, i] + \Delta t \cdot \left(\frac{dq[:, i] + dq[:, i+1]}{2} \right)$$

```

1
2 # bring the model into the workspace. We assume a 2-link manipulator
3 include ("dyn_mod_2LinkManipulator.jl")
4
5 # Define parameters for the model, all in one place
6 # 2-Link Manipulator Model
7 function modelParameters()
8     g = 9.81
9     L1 = 1.0
10    L2 = 0.5
11    m1 = 3.0
12    m2 = 1.0
13    return g, L1, L2, m1, m2
14 end
15
16 # Initial values for the angles and the angular velocities
17 q0 = [-pi/4; 0.0] # initial angles of the 2-links
18 dq0 = [0.0; 0.0]   # initial angular velocities of the 2-links
19
20
21 # set up a time vector
22 T = 5.0 #seconds
23 N = 200
24 t = range(0.0, T, N+1)
25 dt = T/N
26
27
28 # Create vectors to store the angles q, and the angular velocities, dq
29 q = zeros(2,N+1)
30 dq = zeros(2,N+1)
31
32 # Initialize everything
33 q[:,1] = q0
34 dq[:,1] = dq0
35
36 # Set up a for loop to run a combination of a Riemann sum and TrapZ
37 println ("\nThe computed solution is only accurate over a short time interval due to our crude integratin
            method\n")
38 println ("Numerical errors will excite the pendulum and cause an increase in amplitude. This is not physical
            behavior.\n \n")
39
40 for i = 1:N
41     # Call the model and solve for the angular acceleration of the links at time t[i]
42     D, C, G, ~, ~ = dyn_mod_2LinkManipulator(q[:,i], dq[:,i]) # model terms at time t[i]
43     ddq_i = D\(-C*dq[:,i] - G) # acceleration at time t[i]
44     dq[:, i+1] = dq[:, i] + dt * ddq_i # integrate acceleration to compute velocity using
        Riemann Sum
45     q[:, i+1] = q[:, i] + dt * (dq[:, i] + dq[:, i+1]) / 2 # integrate velocity to compute
        position
46                                         # using TrapZ
47 end

```

Observations and Limitations

- This method is only accurate over a short time interval due to its numerical simplicity.

- The **Riemann sum** for velocity and the **trapezoidal rule** for position introduce small integration errors that accumulate over time.
- These numerical errors **excite the pendulum**, leading to unrealistic, unstable behavior over long simulation times.
- More sophisticated integration methods (e.g., **Runge-Kutta**, implicit solvers) are needed for **long-term accuracy**.
- Despite its simplicity, this approach provides a **basic understanding of numerical integration** without requiring any significant knowledge of **ODE solvers**. Later in the course, we will introduce **professional solvers** to improve stability and accuracy.

6.5.2 How the Robot Equations Arise from Lagrange's Equations of Motion

The Robt Equations in (6.65) arise from a careful evaluation of the terms in

$$\frac{d}{dt} \frac{\partial \mathcal{L}(q, \dot{q})}{\partial \dot{q}} - \frac{\partial \mathcal{L}(q, \dot{q})}{\partial q} = \Gamma$$

for

$$\mathcal{L}(q, \dot{q}) := K(q, \dot{q}) - V(q).$$

From Prop. 6.32, and in particular, (6.57), the kinetic energy has the quadratic form

$$K(q, \dot{q}) = \frac{1}{2} \dot{q}_1^\top D(q) \dot{q}.$$

Moreover, the potential energy $V(q)$ does not depend on the velocity terms. These observations yield

$$\begin{aligned} \frac{\partial}{\partial \dot{q}} \mathcal{L}(q, \dot{q}) &= \frac{\partial}{\partial \dot{q}} \left(\frac{1}{2} \dot{q}_1^\top D(q) \dot{q} - V(q) \right) \quad (\text{derivative of a sum}) \\ &= \frac{\partial}{\partial \dot{q}} \left(\frac{1}{2} \dot{q}_1^\top D(q) \dot{q} \right) - \frac{\partial}{\partial \dot{q}} (V(q)) \quad (\text{equals sum of derivatives}) \\ &= \frac{\partial}{\partial \dot{q}} \left(\frac{1}{2} \dot{q}_1^\top D(q) \dot{q} \right) \quad (V(q) \text{ independent of } \dot{q}) \\ &= D(q) \cdot \dot{q} \quad (\text{column convention for the partial derivative}). \end{aligned}$$

The term $\frac{d}{dt} \frac{\partial \mathcal{L}(q, \dot{q})}{\partial \dot{q}} = \frac{d}{dt} (D(q) \dot{q})$ is the total derivative from Chapter 5.8. Inspired by the Product Rule of single-variable differentiation, one is tempted to write

$$\frac{d}{dt} \frac{\partial \mathcal{L}(q, \dot{q})}{\partial \dot{q}} = \frac{d}{dt} (D(q) \cdot \dot{q}) = \dot{D}(q, \dot{q}) \cdot \dot{q} + D(q) \cdot \ddot{q},$$

and, in fact, it is possible to make sense of

$$\dot{D}(q, \dot{q}) := \frac{d}{dt} D(q).$$

Indeed, as in (6.58), let $d_{ij}(q)$ denote the ij -term of the mass-inertia matrix, $D(q)$. Then, by Def. 5.49, the Chain Rule applied to the ij -term of $\frac{d}{dt} D(q)$ is

$$\frac{d}{dt} d_{ij}(q) = \underbrace{\frac{\partial d_{ij}(q)}{\partial q}}_{1 \times n \text{ Jacobian}} \cdot \underbrace{\dot{q}}_{n \times 1},$$

a scalar. In other symbols,

$$\left[\frac{d}{dt} D(q) \right]_{ij} = \frac{\partial d_{ij}(q)}{\partial q} \cdot \dot{q} =: [\dot{D}(q, \dot{q})]_{ij}.$$

Next,

$$\begin{aligned}\frac{\partial \mathcal{L}(q, \dot{q})}{\partial q} &= \frac{\partial}{\partial q} (K(q, \dot{q}) - V(q)) \\ &= \frac{\partial}{\partial q} (K(q, \dot{q})) - \frac{\partial}{\partial q} V(q) \\ &= \frac{\partial}{\partial q} (K(q, \dot{q})) - \nabla V(q).\end{aligned}$$

Hence, doing all of the required bookkeeping yields

$$\frac{d}{dt} \frac{\partial \mathcal{L}(q, \dot{q})}{\partial \dot{q}} - \frac{\partial \mathcal{L}(q, \dot{q})}{\partial q} = D(q) \cdot \ddot{q} + \left(\dot{D}(q, \dot{q}) \cdot \dot{q} - \frac{\partial}{\partial q} (K(q, \dot{q})) \right) + \nabla V(q).$$

In this course, we'll not expand more upon the middle term in parentheses, except to say that it provides the term traditionally denoted as $C(q, \dot{q}) \cdot \dot{q}$, that is,

$$C(q, \dot{q}) \cdot \dot{q} := \dot{D}(q, \dot{q}) \cdot \dot{q} - \frac{\partial K(q, \dot{q})}{\partial q}.$$

In case you are thinking, “There is no way I’m doing this calculation!” we agree: for all but the simplest of mechanical systems, computing Lagrange’s EoM by hand is extremely tedious and error-prone. Today, however, Symbolic packages/toolboxes allow the calculations to be done accurately and quickly. Moreover, as we saw with the 3-link manipulator, we can easily turn the computed quantities into executable code. **This underlines, once again, why a Calculus course like this one is needed: without seeing with one’s own eyes the software tools of today in action, it is all too easy to dismiss an amazing body of work as completely beyond your abilities.**

Things have been taken even one step further, where the robot model is extracted directly from the CAD files used to design the robot! Roboticists have found it beneficial to describe robot geometries in a manner that’s both human-friendly and not tied to any specific programming language. Imagine having a textual blueprint: “part-one is situated 1 meter to the left of part-two and showcases this particular triangular mesh for visualization,” etc, but only much more compact and precise. The **Unified Robot Description Format (URDF)** stands out as the leading programming-language-agnostic medium for this purpose at the present time.

Secrets of the Arcane 6.41: URDF Files

Here are some resources that could be beneficial for undergraduates looking to learn URDF (Unified Robot Description Format) files for robots:

1. ROS Wiki Tutorials:

- The ROS (Robot Operating System) Wiki provides a series of tutorials on URDF. It covers building a visual robot model with URDF from scratch, defining movable joints in URDF, and adding physical and collision properties to a URDF model. These tutorials provide a step-by-step approach which could be very beneficial for learning URDF. ([source](#))

2. Medium Article on Creating URDF files:

- A Medium article discusses creating URDF files for robots. It seems to cover forward kinematics and mentions that it will later cover inverse kinematics as well. The author has uploaded some files on GitHub for reference. This article could provide practical insights and hands-on examples for learning URDF. ([source](#))

3. Learning ROS for Robotics Programming:

- Another book that provides a section on creating the first URDF file, which is part of building a mobile robot with four wheels and an arm with a gripper. This resource might provide a good hands-on approach to learning URDF through a practical project. ([source](#))

6.5.3 Including Actuators in Lagrange’s Equations

In the 3-link manipulator, we included the actuation term without explaining how it is done. This is not a course on Lagrangian Dynamics, so it’s necessary to take some shortcuts. Just in case you are curious, this is how to compute Γ in the robot equations for planar robots with n -revolute joints when the motors, if any, are mounted either

- on the ground (aka, to the world frame), with link-1 attached to the output shaft of the motor, or for $i > 1$,
- at the end of link- i with link- $i + 1$ attached to the output shaft of the motor.

Not all links have to be actuated; humanoids are famously **underactuated**, which is the term for robots that do not have a motor on every joint (Why not include a lot of motors? They are heavy and use a lot of energy!) We also assume that the generalized coordinates have been assigned as we've done with the 3-link manipulator, with the base link associated to an absolute angle q_1 and subsequent links assigned relative angles q_i .

Let q_a be an $m \times 1$ column vector containing the list of actuated joint angles. For the 3-link manipulator, we assumed all of the joints were actuated, so we took $m = 3$ and

$$q_a := \begin{bmatrix} q_1 \\ q_2 \\ q_3 \end{bmatrix}.$$

The recipe for computing Γ is then

$$\Gamma := B \cdot \tau$$

where,

$$B := \left(\frac{\partial q_a}{\partial q} \right)^\top,$$

an $n \times m$ matrix, and τ is the $m \times 1$ vector of motor torques

$$\tau := \begin{bmatrix} \tau_1 \\ \vdots \\ \tau_m \end{bmatrix}.$$

Usually, there is gearing between a motor and a joint. Moreover, the spinning of the motor's rotor contributes significantly to the kinetic energy of the overall system. All of these effects can be taken into account when building the Lagrangian model of your robot. To learn from a simplified presentation, see the bottom of page 433 of [Feedback Control of Dynamic Bipedal Locomotion](#) by Westervelt et al., **Accounting for motors and rigid gear trains**. You might also watch the video [Actuation, Gearing, and Friction](#) by Prof. Kevin Lynch of Northwestern University.

6.5.4 Moment of Inertia

This is a repeat of Section 3.4.5, with the idea that you may be more interested in it now.

The **moment of inertia** (I) of an object about a specific axis is a measure of how difficult it is to change the rotational motion of that object. It depends on the mass distribution of the object relative to the axis of rotation. The moment of inertia is given by the integral

$$I = \int r^2 dm,$$

where r is the distance of a small mass element (dm) from the axis of rotation. This formula highlights that the further the mass is distributed from the axis, the higher the moment of inertia, thus making it harder to rotate the object.

For practical applications, common shapes have standard formulas for their moments of inertia. For example, for a solid cylinder of radius R rotating about its central axis, the moment of inertia is

$$I = \frac{1}{2}mR^2,$$

where m is the mass of the cylinder. This value helps in understanding the rotational characteristics of the cylinder and plays a crucial role in dynamics and engineering.

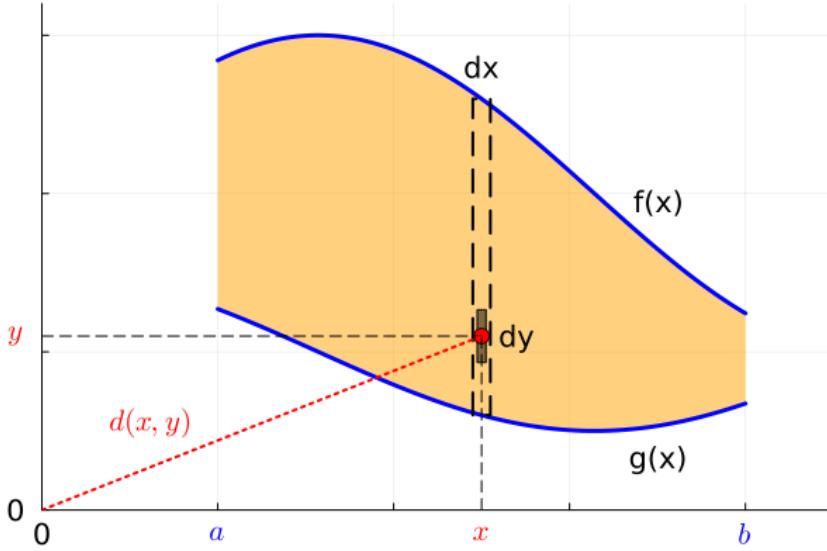


Figure 6.11: (**Moment of Inertia about the z -Axis**): This figure illustrates how to use a rectangle of width dx , upper bounded by $f(x)$ and lower bounded by $g(x)$, to compute the moment of inertia about the z -axis, that is, the axis perpendicular to the (x, y) -plane and passing through the origin. The body is assumed to be of constant density. The computation requires two integrals. The first one, along the y -direction, is a very simple computation of the moment of inertia of an infinitesimally thin rectangle of width dx . The second one, along the x -direction, adds up the contribution of each thin rectangle. The final result is given in (6.66).

Proposition 6.42: Moment of Inertia about the z -Axis for Planar Bodies of Constant Density

Given the planar body in Fig. 6.11 bounded by $y = f(x)$ and $y = g(x)$, with constant density ρ , the moment of inertia about the z -axis can be computed by a single integral from $x = a$ to $x = b$,

$$I_z = \int_a^b \rho \left(x^2(f(x) - g(x)) + \frac{1}{3}(f^3(x) - g^3(x)) \right) dx. \quad (6.66)$$

Proof: The proof is another example of the power of the rectangle in solving Calculus problems. Please refer to Fig. 6.11 throughout the derivation.

Consider the vertical rectangle of width dx , and then within the rectangle, a small rectangle of height dy . The differential area of the infinitesimal rectangle is the product of height times width, that is, $dA := dy dx$ and the differential mass is the product of the density ρ times the differential area, giving us

$$dm = \rho dA = \rho dy dx. \quad (6.67)$$

The differential moment of inertia about the z -axis is defined to be the distance squared times the differential mass, where

$$r^2 = d^2(x, y) = x^2 + y^2.$$

We need to add up the contribution of each differential mass along the y -axis of the thin rectangle, which results in the integral

$$\begin{aligned} dI_z(x) &= \left[\int_{g(x)}^{f(x)} d^2(x, y) \rho dy \right] dx \quad (\text{integral w.r.t. } y) \\ &= \left[\int_{g(x)}^{f(x)} (x^2 + y^2) \rho dy \right] dx \quad (x \text{ is treated like a constant when integrating w.r.t. } y) \\ &= \left[x^2 y + \frac{y^3}{3} \right]_{g(x)}^{f(x)} \rho dx \\ &= \left(x^2(f(x) - g(x)) + \frac{1}{3}(f^3(x) - g^3(x)) \right) \rho dx. \end{aligned} \quad (6.68)$$

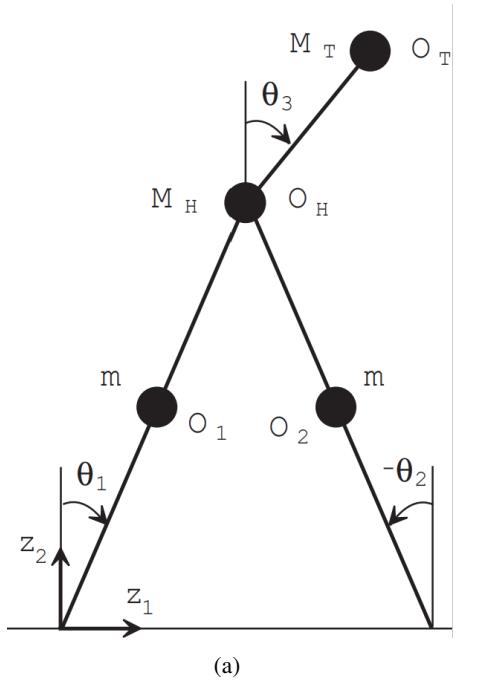
Next, we need to integrate with respect to x so as to add up the contribution of each thin rectangle because that is how Calculus works!

$$\begin{aligned} I_z &= \int_a^b dI_z(x) \\ &= \int_a^b \left(x^2(f(x) - g(x)) + \frac{1}{3}(f^3(x) - g^3(x)) \right) \rho dx, \end{aligned} \quad (6.69)$$

which establishes (6.66). ■

Suggested Videos:

- [More on moment of inertia | Moments, torque, and angular momentum](#) by Khan Academy.
- [Demonstrating Rotational Inertia \(or Moment of Inertia\)](#) by Flipping Physics. It's a bit goofy but full of information if you can stick through it. It's also a bit advanced.
- [Moment of Inertia](#) by DoctorPhys. This is Lesson 1 of many! Next, we call out Lesson 9 because it pertains to the BallBot.
- [Moment of Inertia for a Spherical Shell](#) by DoctorPhys. This is the moment of inertia of a basketball rolling on a floor.



(a)



(b)



(c)

Figure 6.12: Michigan's foray into bipedal locomotion began with the study of the 3-link walker shown in (a), and over the years, progressed to Cassie in (b), and collaborating with [Wandercraft](#) on the exoskeleton, Atalante, shown in (c). Atalante allows patients with paraplegia to walk again. In addition to powering the patient's legs, it provides the overall balance for the patient and itself. It is truly a wearable robot.

6.5.5 A Simple Planar Biped

Most of the terms in the model are defined in Fig. 6.12-(a). Missing is the length of the torso (line from O_H to O_T) is L and the four line segments in the legs each have length r .

¹ # Julia code
² using Symbolics

```

3
4 # 3-link Walker Model is from 1999.
5 # All angles are ABSOLUTE and are
6 # Positive in the CW Direction (not CCW)
7
8 # Variable declarations
9 @variables th1 th2 th3 dth1 dth2 dth3 m Mh Mt g L r
10
11 # Generalized coordinates
12 q = [th1, th2, th3]
13
14 # First derivative of generalized coordinates
15 dq = [dth1, dth2, dth3]
16
17
18 # Position of masses in the system
19 p_m1 = [r/2*sin(th1), r/2*cos(th1)]
20 p_Mh = [r*sin(th1), r*cos(th1)]
21 p_Mt = p_Mh + [L*sin(th3), L*cos(th3)]
22 p_m2 = p_Mh - [r/2*sin(th2), r/2*cos(th2)]
23
24 # Potential energy of masses in the system
25 PE_m1 = p_m1[2] * m * g
26 PE_Mh = p_Mh[2] * Mh * g
27 PE_Mt = p_Mt[2] * Mt * g
28 PE_m2 = p_m2[2] * m * g
29
30 # Total potential energy of the system
31 PE = PE_m1 + PE_Mh + PE_Mt + PE_m2
32 PE = Symbolics.simplify(PE)
33
34 # Velocities of masses in the system
35 v_m1 = Symbolics.jacobian(p_m1, q) * dq
36 v_Mh = Symbolics.jacobian(p_Mh, q) * dq
37 v_Mt = Symbolics.jacobian(p_Mt, q) * dq
38 v_m2 = Symbolics.jacobian(p_m2, q) * dq
39
40 # Kinetic Energy of each mass on the Robot
41 KE_m1 = 0.5 * m * sum(v_m1.*v_m1)
42 KE_Mh = 0.5 * Mh * sum(v_Mh.*v_Mh)
43 KE_Mt = 0.5 * Mt * sum(v_Mt.*v_Mt)
44 KE_m2 = 0.5 * m * sum(v_m2.*v_m2)
45
46 # Total kinetic energy of the system
47 KE = KE_m1 + KE_Mh + KE_Mt + KE_m2
48 KE = Symbolics.simplify(KE)
49
50 # Form the Lagrangian
51 L = KE - PE
52
53 # Compute terms in the Robot Equations
54 G = Symbolics.gradient(PE, q)
55 D = Symbolics.hessian(KE, dq) # Strange name explained below
56
57 # Compute C (advanced, uses Christoffel symbols)
58 n = length(q)
59 C = Array{Num}(undef, n, n) # Num is Symbolic Type

```

```

60 # Strange name explained below
61 for k in 1:n
62     for j in 1:n
63         C[k, j] = 0
64     for i in 1:n
65         C[k, j] += 0.5 * (Symbolics.gradient(D[k, j], q)[i] +
66             Symbolics.gradient(D[k, i], q)[j] - Symbolics.gradient(D[i, j], q)[k]) * dq[i]
67     end
68 end
69 end
70
71 # Place an actuator at the hip on
72 # the relative angle between the torso and
73 # each leg
74 # Use D'Alembert's Principle
75 E = [th3-th1; th3-th2]
76 B = Symbolics.jacobian(E, q)
77 B = B'
78 JacG = Symbolics.jacobian(G, q) # Used when we discuss linearization of NL models
79
80
81
82 # Only run simple after all Symbolic manipulation is done and you
83 # are ready to create a function
84 if true
85     # Simplify as much as possible at the very end
86     D = simple(D)
87     C = simple(C)
88 end
89
90
91 # Create a function that computes all terms
92 # in the Robot Equations
93 fcn_name = "dyn_mod_3LinkWalker"
94 modelParamString = "g, r, L, m, Mh, Mt = modelParameters3LinkWalker()"
95 # Need to include line break command \n
96 variableNamesString = "th1, th2, th3 = q \ndth1, dth2, dth3 = dq"
97 writeEOM(fcn_name, D, C, G, B, JacG, modelParamString, variableNamesString)
98
99 # bring the model into the workspace for later use
100 include("dyn_mod_3LinkWalker.jl")
101
102
103 # If the line above fails, it is because Julia is trying to read the file before it has
104 # been made available to it by the file system. Simply copy the command into a new cell
105 # and
# run it. The time it takes you to do that will allow the file system to catch up! :-)

```

Output

```
[creating DYN_MOD_3LINKWALKER.jl]
File dyn_mod_3LinkWalker.jl created successfully!
dyn_mod_3LinkWalker (generic function with 1 method)
```

Here is the Lagrangian Model.

```

1 function modelParameters3LinkWalker ()
2 #
3     r=1;      # length of each leg
4     m=5;      # mass of a leg
5     Mh=15;    # mass of hips
6     Mt=10;    #mass of torso
7     L=0.5;    # distance between hips and torso mass
8     g=9.8;    # acceleration due to gravity
9 #
10    return g, r, L, m, Mh, Mt
11 end
12
13
14 function dyn_mod_3LinkWalker (q, dq)
15 # DYN_MOD_3LINKWALKER
16 # 2024-02-21 17:48:08
17 #
18 # Author: Grizzle
19 #
20 # Model NOTATION: D(q)ddq + C(q,dq)*dq + G(q) = B*tau
21 # The Robot Equations: From Lagrange's Equations of Motion
22 #
23 g, r, L, m, Mh, Mt = modelParameters3LinkWalker ()
24 #
25 # Variable names for the model
26 th1, th2, th3 = q
27 dth1, dth2, dth3 = dq
28 #
29 D = zeros (3, 3)
30 D[1, 1] = r^2*(Mh + Mt + 5*m/4)
31 D[1, 2] = -0.5*m*r^2*cos(th1 - th2)
32 D[1, 3] = L*Mt*r*cos(th1 - th3)
33 D[2, 1] = -0.5*m*r^2*cos(th1 - th2)
34 D[2, 2] = 0.25*m*r^2
35 D[2, 3] = 0
36 D[3, 1] = L*Mt*r*cos(th1 - th3)
37 D[3, 2] = 0
38 D[3, 3] = L^2*Mt
39 #
40 C = zeros (3, 3)
41 C[1, 1] = 0.0
42 C[1, 2] = -0.5*dth2*m*r^2*sin(th1 - th2)
43 C[1, 3] = L*Mt*dth3*r*sin(th1 - th3)
44 C[2, 1] = 0.5*dth1*m*r^2*sin(th1 - th2)
45 C[2, 2] = 0.0
46 C[2, 3] = 0.0
47 C[3, 1] = -L*Mt*dth1*r*sin(th1 - th3)
48 C[3, 2] = 0.0
49 C[3, 3] = 0.0
50 #
51 G = zeros (3)
52 G[1] = -Mh*g*r*sin(th1) - Mt*g*r*sin(th1) - (3/2)*g*m*r*sin(th1)
53 G[2] = (1/2)*g*m*r*sin(th2)
54 G[3] = -L*Mt*g*sin(th3)
55 #
56 B = zeros (3, 2)

```

```

57  B[1, 1] = -1
58  B[2, 2] = -1
59  B[3, 1] = 1
60  B[3, 2] = 1
61 #
62 JacG = zeros(3, 3)
63 JacG[1, 1] = -Mh*g*r*cos(th1) - Mt*g*r*cos(th1) - (3/2)*g*m*r*cos(th1)
64 JacG[2, 2] = (1/2)*g*m*r*cos(th2)
65 JacG[3, 3] = -L*Mt*g*cos(th3)
66 #
67 return (D=D, C=C, G=G, B=B, JacG=JacG)
68 end

```

6.6 (Optional Read:) Proof behind Lagrange Multipliers for Equality Constraints

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a differentiable cost function, $g_i : \mathbb{R}^n \rightarrow \mathbb{R}$ be differentiable equality constraints for $1 \leq i \leq m$, and let $x_0 \in \mathbb{R}^n$ be an arbitrary point. We seek necessary conditions for x_0 to be a local extremum of the constrained optimization problem (6.40), which we repeat here for convenience,

$$\begin{aligned} & \text{Minimize} && f(x) \\ & \text{subject to} && g_i(x) = 0, \quad 1 \leq i \leq m, \end{aligned} \tag{6.70}$$

We use the Jacobian to define $(n+1)$ constant row vectors

$$\begin{aligned} J_f(x_0) &:= \frac{\partial f(x_0)}{\partial x} = \left[\begin{array}{cccc} \frac{\partial f(x_0)}{\partial x_1} & \frac{\partial f(x_0)}{\partial x_2} & \dots & \frac{\partial f(x_0)}{\partial x_n} \end{array} \right], \\ J_{g_i}(x_0) &:= \frac{\partial g_i(x_0)}{\partial x} = \left[\begin{array}{cccc} \frac{\partial g_i(x_0)}{\partial x_1} & \frac{\partial g_i(x_0)}{\partial x_2} & \dots & \frac{\partial g_i(x_0)}{\partial x_n} \end{array} \right], \quad 1 \leq i \leq m. \end{aligned} \tag{6.71}$$

By definition⁷, we note that $J_f(x_0)$ is a linear combination of the set of row vectors $\{J_{g_1}(x_0), J_{g_2}(x_0), \dots, J_{g_m}(x_0)\}$, if, and only if, there exist scalars $\lambda_1, \lambda_2, \dots, \lambda_m \in \mathbb{R}$ such that

$$J_f(x_0) + \lambda_1 J_{g_1}(x_0) + \lambda_2 J_{g_2}(x_0) + \dots + \lambda_m J_{g_m}(x_0) = 0. \tag{6.72}$$

In other words,

- $J_f(x_0)$ being a linear combination of $\{J_{g_1}(x_0), J_{g_2}(x_0), \dots, J_{g_m}(x_0)\}$ will be intimately related to the existence of Lagrange multipliers transforming the constrained optimization problem (6.70) into the unconstrained problem of finding stationary points for the Lagrangian

$$L(x, \lambda) := f(x) + \lambda_1 g_1(x) + \lambda_2 g_2(x) + \dots + \lambda_m g_m(x). \tag{6.73}$$

- We also note that the set of row vectors $\{J_{g_1}(x_0), J_{g_2}(x_0), \dots, J_{g_m}(x_0)\}$ is **linearly independent** if, and only if, the constraint qualification condition holds at x_0 ,

$$\text{rank} \left(\frac{\partial g(x_0)}{\partial x} \right) = \text{rank} \left(\begin{bmatrix} J_{g_1}(x_0) \\ J_{g_2}(x_0) \\ \vdots \\ J_{g_m}(x_0) \end{bmatrix} \right) = m. \tag{6.74}$$

- The coefficients in (6.72) are unique if, and only if, the set of row vectors $\{J_{g_1}(x_0), J_{g_2}(x_0), \dots, J_{g_m}(x_0)\}$ is linearly independent.

We next state a result that makes it easy to relate conditions for x_0 being a local extremum of (6.70) to x_0 being a stationary point of the Lagrangian (6.73).

⁷Write $J_f(x_0) = \alpha_1 J_{g_1}(x_0) + \alpha_2 J_{g_2}(x_0) + \dots + \alpha_m J_{g_m}(x_0)$, and then define $\lambda_i := -\alpha_i$. Moving all the terms to the left side of the equation gives (6.72).

Proposition 6.43: When can a Row Vector be Expressed as a Linear Combination of other Row Vectors

The row vector $J_f(x_0)$ is a linear combination of the set of row vectors $\{J_{g_1}(x_0), J_{g_2}(x_0), \dots, J_{g_m}(x_0)\}$ if, and only if,

$$\bigcap_{i=1}^n \text{null}(J_{g_i}(x_0)) \subset \text{null}(J_f(x_0)). \quad (6.75)$$

We delay the proof of Prop. 6.43 until later.

Assumptions:

- x_0 is feasible; that is $g_i(x_0) = 0, 1 \leq i \leq m$.
- x_0 satisfies the constraint qualification condition, that is, (6.74) holds.

Given x_0 as above, we call $v \in \mathbb{R}^n$ a **feasible direction** if $v \in \bigcap_{i=1}^n \text{null}(J_{g_i}(x_0))$. Why this name? One way to understand it is that by the linear approximation of $g_i : \mathbb{R}^n \rightarrow \mathbb{R}$ at a feasible point x_0 ,

$$\begin{aligned} g_i(x_0 + v) &\approx g_i(x_0) + J_{g_i}(x_0) \cdot v \\ &\approx J_{g_i}(x_0) \cdot v \\ &\approx 0. \end{aligned}$$

Hence, if x_0 is feasible, then so is $x_0 + v$ to first order⁸.

Let $v \in \bigcap_{i=1}^n \text{null}(J_{g_i}(x_0))$ be arbitrary, and consider $f(x_0 + v)$ through the lens of its linear approximation,

$$f(x_0 + v) \approx f(x_0) + J_f(x_0) \cdot v.$$

From the linearization of f at x_0 , we make the following conclusions:

- If $v \notin \text{null}(J_f(x_0))$, meaning $J_f(x_0) \cdot v \neq 0$, then we can move along v to make $f(x_0 + v)$ smaller or larger than $f(x_0)$. In this case, x_0 is not an extreme point of (6.70).
- Therefore, the only way for x_0 to be an extreme point is for $J_f(x_0) \cdot v = 0$, that is $v \in \text{null}(J_f(x_0))$.

Because $v \in \bigcap_{i=1}^n \text{null}(J_{g_i}(x_0))$ was arbitrary, we conclude that the only way for x_0 to be an extreme point of (6.70) is for (6.75) to hold, and thus $J_f(x_0)$ can be expressed as the linear combination

$$J_f(x_0) + \lambda_1 J_{g_1}(x_0) + \lambda_2 J_{g_2}(x_0) + \dots + \lambda_m J_{g_m}(x_0) = 0. \quad (6.76)$$

We are now essentially done, with only some bookkeeping to complete the proof. Using the coefficients in (6.76), define the Lagrangian as in (6.73). Then, by (6.76), x_0 is a stationary point of L .

That seems almost too easy! But it was easy only because we knew about Prop. 6.43 for checking that a given row vector could be expressed as a linear combination of other row vectors. Knowledge is power. ■

Proof of Prop. 6.43 Without loss of generality, we can assume that the set of row vectors $\{J_{g_1}(x_0), J_{g_2}(x_0), \dots, J_{g_m}(x_0)\}$ is linearly independent because any dependent vectors do not change the intersection of their null spaces.

⁸A more formal formulation uses the concept of a **directional derivative**, defined as, for v a unit vector,

$$\frac{d}{dt} g_i(x_0 + t \cdot v) \Big|_{t=0} = \frac{\partial g_i(x_0)}{\partial x} \cdot v = J_{g_i}(x_0) \cdot v.$$

Hence, the directional derivative vanishes for $v \in \text{null}(J_{g_i}(x_0))$.

If there exist $\lambda_1, \dots, \lambda_m$ such that (6.72) holds, then v satisfying $J_{g_i}(x_0) \cdot v = 0$, $1 \leq i \leq m$, implies $J_f(x_0) \cdot v = 0$. For the other direction, we assume that $\{J_{g_1}(x_0), J_{g_2}(x_0), \dots, J_{g_m}(x_0)\}$ is linearly independent and that $J_f(x_0)$ is not a linear combination of these vectors. Hence, the $(m+1) \times n$ matrix

$$A := \begin{bmatrix} J_f(x_0) \\ J_{g_1}(x_0) \\ J_{g_2}(x_0) \\ \vdots \\ J_{g_m}(x_0) \end{bmatrix}$$

has full row rank. Therefore, there exist $v \in \mathbb{R}^n$ such that

$$A \cdot v = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}.$$

Hence, there exists $v \in \bigcap_{i=1}^n \text{null}(J_{g_i}(x_0))$ such that $v \notin \text{null}(J_f(x_0))$. This completes the proof. ■

6.7 (Optional Read:) Proof Underlying Gradient Descent with Equality Constraints

It is a bit ironic, but the above proof of the Method of Lagrange Multipliers shows how to accomplish constrained optimization without explicitly computing Lagrange multipliers! In fact, it shows how to update gradient descent to work in the presence of equality constraints. We retain all of the assumptions made for the proof of the Method of Lagrange Multipliers with Equality Constraints.

Method 6.44: Gradient Descent with Equality Constraints

Consider the constrained optimization problem (6.40) for functions $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ with optimal value x^* . If the columns vectors from the gradients $\{\nabla g_1(x^*), \nabla g_2(x^*), \dots, \nabla g_m(x^*)\}$ are linearly independent (equivalently, the rows of the Jacobian $\frac{\partial g(x^*)}{\partial x}$ are linearly independent), then a locally optimal x^* can be found as follows:

Initialization: Find a feasible point $x_0 \in \mathbb{R}^n$, that is, a point such that $g_i(x_0) = 0$ for $1 \leq i \leq m$. This can be accomplished with `NLsolve`.

Loop: Given a point $x_k \in \mathbb{R}^n$, repeat the following steps to either find x_{k+1} or terminate.

- (a) **Compute:** $\nabla f(x_k)$, the gradient of f at x_k and $\frac{\partial g(x_k)}{\partial x}$, the Jacobian of the constraint vector.

From the proof of the Method of Lagrange Multipliers, the set of feasible directions is

$$V := \bigcap_{i=1}^n \text{null}(J_{g_i}(x_k)) = \text{null}\left(\frac{\partial g(x_k)}{\partial x}\right),$$

where $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is the vector of equality constraints.

Wow, it's good to know about null spaces. They are not just abstract nonsense after all!

- (b) **Compute the orthogonal projection of $\nabla f(x_k)$ onto V , the set of feasible directions:** This can be accomplished by

applying Gram-Schmidt. Indeed, let $\{b_1, b_2, \dots, b_m\}$ be an orthonormal basis for the columns of $\left(\frac{\partial g(x_k)}{\partial x}\right)^\top$. Then

$$P_{\nabla f(x_k)}^V := \nabla f(x_k) - \sum_{i=1}^m (\nabla f(x_k) \bullet b_i) \cdot b_i$$

is the desired projection.

Why: Each term $(\nabla f(x_k) \bullet b_i) \cdot b_i$ is a component of $\nabla f(x_k)$ that is orthogonal to V , the feasible directions. From ROB 101 *Computational Linear Algebra*, Chapter 9, once these components are subtracted off, only the part of $\nabla f(x_k)$ lying in V is left.

Wow, it's good to know about Gram-Schmidt and orthogonal vectors. Like null spaces, these are very useful tools.

(c) **Update or Terminate, That is the Question:** If the termination condition is not met, define

$$x_{k+1} = x_k - s \cdot P_{\nabla f(x_k)}^V,$$

where $s > 0$ is the step size. Loop back to (a).

(d) **Termination Condition:** $\|P_{\nabla f(x_k)}^V\| < a_{\text{tol}}$; in other words, $J_f(x_k)$ is to numerical precision a linear combination of $\{J_{g_1}(x_k), J_{g_2}(x_k), \dots, J_{g_m}(x_k)\}$, and thus x_k is (to numerical precision) a stationary point of the Lagrangian.

Notes:

- In practice, the updates x_k will “slowly” exit the feasible set. Hence, when $\|g(x_k)\|$ exceeds a user-defined tolerance, a new call to `NLsolve` should be made.
- Lagrange likely understood that an algorithm like this was possible, but given that his “computer” consisted of pen and paper, it was not feasible to pursue this version of his method.
- If there are no constraints, then the null space of $\left(\frac{\partial g(x_k)}{\partial x}\right)$ is all of \mathbb{R}^n . The projected gradient is then the regular gradient.
- An algorithm like this one cannot be presented in most Calc III courses **because Linear Algebra traditionally comes after Calc I through Calc IV**. Michigan Robotics is pushing to reverse this sequence and put Linear Algebra before Calculus.
- **JuMP does all of this and a whole lot more.** It is your best option for use in professional projects. For learning about gradient descent, the workhorse of Machine Learning (ML), algorithms like the one presented here are helpful.

Here is an instantiation of the above version of Lagrange’s Method for Optimization with Equality Constraints.

```

1 using ForwardDiff, LinearAlgebra, NLsolve
2
3 # Define the cost function f(x) and the equality constraints G(x)
4 function f(x)
5     return x[1]^2 + x[2]^2 + x[3]^2 + x[4]^2
6 end
7
8 function G(x)
9     return [x[1] + x[2] - 3, sin(pi*x[1]) - x[3]^2 - 1 + x[4]^3]
10 end
11
12 # Compute the gradient of f and the Jacobian of G using ForwardDiff
13 grad_f(x) = ForwardDiff.gradient(f, x)
14 JacG(x) = ForwardDiff.jacobian(G, x)

```

```

15
16 # Define functions used for computing the projection of grad_f onto
17 # the nullspace of JacG, that is, onto the feasible directions
18
19 function gram_schmidt (jacobian_G)
20     n, m = size (jacobian_G)
21     B = zeros (n, m) # Basis vectors
22
23     for j = 1:m
24         b = jacobian_G [:, j]
25         for i = 1:j-1
26             b -= dot (B [:, i], jacobian_G [:, j]) * B [:, i]
27         end
28         B [:, j] = b / norm (b)
29     end
30
31     return B
32 end
33
34 # Project the gradient onto the feasible directions,
35 # that is, nullspace of JacG
36 function project_grad_f (grad_f, B)
37     # Iterate over columns of B, subtracting their projection from grad_f
38     projected_grad_f = grad_f
39     for j = 1:size (B, 2)
40         projected_grad_f = projected_grad_f - dot (grad_f, B [:, j]) * B [:, j]
41     end
42     return projected_grad_f
43 end
44
45
46 # Function to find a feasible starting point using NLsolve
47 # or to push xk back onto the feasible set when it
48 # drifts off
49 function find_feasible_start (G, x_guess)
50     # build function for NLsolve
51     function constraint! (F, x)
52         F [1:2] = G (x)
53     end
54     xFeasible = nlsolve (constraint!, x_guess)
55     return xFeasible .zero
56 end
57 #
58 ## End define functions for projected gradient descent
59
60 ## Actual Gradient Descent Algorithm Starts Here
61
62 x_guess = [0.0, 0.0, 0.0, 0.0] # Not feasible
63 x0 = find_feasible_start (G, x_guess) # Make feasible
64
65
66 # set step size, maxIter, tolerances
67 s=0.1
68 maxIter=10000
69 gradTol = 1e-6
70 constraintTol = 1e-4
71

```

```

72 xk = x0
73 for k=1 :maxIter
74     jacobian_G_transposed = transpose (JacG (xk) )
75     B = gram_schmidt (jacobian_G_transposed)
76     projected_grad_f = project_grad_f (grad_f (xk) , B)
77     if norm (projected_grad_f) < gradTol
78         println ("---- Iteration $k ----")
79         println ("Met tolerance . ")
80         println ("Jacf applied to null (JacG) (should be near 0) : ",
81             round . ( (grad_f (xk) ') *nullspace (JacG (xk) ) , digits=5) )
82         println ("Constraints (should be near 0) : " , round . (G (xk) , digits=5) )
83         println ("-----")
84         break
85     else
86         # Update
87         xk = xk - s*projected_grad_f
88     end
89     if norm (G (xk) ) > constraintTol
90         xk = find_feasible_start (G, xk) # Again Feasible
91     end
92 end
93 xStar = xk
94 fStar = f (xStar)
95
96 println ("\nFinal Results : ")
97 println ("-----")
98 println ("Optimal x: " , round . (xStar, digits=5) )
99 #
100 println ("Optimal Cost f(xStar) : " , round (fStar, digits=5) )

```

Output

```

---- Iteration 184 ----
Met tolerances.
Jacf applied to null (JacG) (should be near 0) : [0.0 -0.0]
Constraints (should be near 0) : [-0.0, 2.0e-5]
-----

Final Results:
-----
Optimal x: [0.95451, 2.04549, 0.0, 0.95008]
Optimal Cost f(xStar) : 5.99778

```

6.8 Words of Wisdom from one of the Planet's Greatest Living Mathematicians

Math isn't actually Sorcery by Prof. Terrence (Terry) Tao, an Australian-American mathematician who is widely regarded as one of the greatest living mathematicians (if you hear an accent in the video, it's Aussie).

Summarized from [Wikipedia](#): Born in 1975 to Chinese immigrant parents in Adelaide, Tao was a child prodigy who exhibited extraordinary mathematical abilities from a very young age. **He skipped six grades and attended university-level mathematics courses at age 9.** At just eight years old, he scored 760 on the SAT math section, a feat achieved by only three children in the history of the Johns Hopkins Study of Exceptional Talent program.

Tao's prodigious talent was further demonstrated through his participation in the International Mathematical Olympiad. He became the youngest participant in the competition's history, winning bronze, silver, and gold medals at ages 10, 11, and 13, respectively. He

remains the youngest winner of each medal type in the Olympiad's history.

Currently, Tao is a professor of mathematics at the University of California, Los Angeles (UCLA), where he holds the James and Carol Collins Chair in the College of Letters and Sciences. His research spans various areas of mathematics, including harmonic analysis, partial differential equations, combinatorics, probability theory, and analytic number theory.

Tao's contributions to mathematics have earned him numerous prestigious awards, including the Fields Medal in 2006, the Royal Medal and Breakthrough Prize in Mathematics in 2014, and a MacArthur Fellowship in 2006. He has authored or co-authored over three hundred research papers, solidifying his reputation as one of the most influential mathematicians of his generation.

Relevance to Calculus for the Modern Engineer: At each step of our adventure, we've broken a new topic into smaller building blocks from which we could assemble an approach to solving a problem. Often, the building blocks, put end-to-end, formed an algorithm. Other times, they provided an Ah Ha moment of insight into how a complicated looking formula was a direct consequence of understanding rectangles (e.g., theory of integration and its applications) or linear approximations of functions (e.g., derivation of differentiation rules and the applications of derivatives, scalar and vector). Limits of one kind or another have been the cement holding the blocks together.

Chapter 7

Antiderivatives and the Fundamental Theorems of Calculus

Learning Objectives

By the end of this chapter, students should be able to:

- Understand that antidifferentiation is essentially the inverse operation of differentiation.
- Comprehend the two Fundamental Theorems of Calculus and explain the precise sense in which differentiation and integration are inverse operations on functions.
- Understand that while antidifferentiation and definite integration are technically distinct operations, they are often perceived as identical in the context of Calculus education. This perception is partially justified by their closely related concepts. However, it's crucial to acknowledge the potential drawbacks of this viewpoint, particularly the confusion and discouragement that can arise from the complex rules associated with manually computing antiderivatives.
- Appreciate the historical and practical context of computing antiderivatives by hand, and recognize the modern approaches that render hand computations almost obsolete for all but the simplest of cases.

Outcomes

Upon successful completion of this chapter, students will:

- Be introduced to the concept of antiderivatives and familiarize themselves with elementary techniques for finding them.
- Gain a thorough understanding of the Fundamental Theorems of Calculus, including their implications in both geometric and analytic contexts.
- Develop an insight into the interplay between differentiation, antiderivatives, and definite integration, enhancing their comprehension of calculus as a whole.
- Cultivate an appreciation for the traditional methods of finding antiderivatives while also understanding the value and efficiency of using computational tools for real-world applications.

Antiderivative $\int f(x) dx$	Function $f(x)$	Derivative $f'(x)$
C	0	0
$kx + C$	k	0
$\frac{1}{3}x^3 + C$	x^2	$2x$
$e^x + C$	e^x	e^x
$-\cos(x) + C$	$\sin(x)$	$\cos(x)$
$x \ln(x) - x + C$	$\ln(x)$	$\frac{1}{x}$
$\ln x + C$	$\frac{1}{x}$	$-\frac{1}{x^2}$
$-\ln \cos(x) + C$	$\tan(x)$	$1 + \tan^2(x)$

Table 7.1: C is an arbitrary constant, later called a **constant of integration**. $F(x)$ is an antiderivative of $f(x)$ if $F'(x) = f(x)$. This table illustrates the relationships between an antiderivative, a function, and a derivative. The functions in the middle are the derivatives of the functions in the left column, and the functions in the right column are the derivatives of the functions in the middle column. Going the other direction, the functions in the middle column are antiderivatives of the functions in the right column. As labeled, the functions in the left column are antiderivatives of the functions in the middle column. Bottom line: $F(x) + C$, for C a constant, is always an antiderivative for $f(x) := \frac{dF(x)}{dx}$.

7.1 Introduction

This chapter marks a major waypoint in our journey to discover and understand the essence of Calculus. The **concept of an antiderivative** is a bridge connecting the micro (differential calculus, focused on the infinitesimal, df , dx , and their ratio) to the macro (integral calculus, focused on adding up an infinite number of vanishingly small areas, $dA = f(x) dx$). As illustrated in Table 7.1, an antiderivative of a function $f(x)$ is another function $F(x)$ such that $F'(x) = f(x)$. The Fundamental Theorems of Calculus elegantly tie together differentiation and integration, showing that these two operations are, in fact, inverse processes. As we unveil the theorems, we will discover the profound connection between the definite integral and antiderivatives, which will further lead us to study indefinite integration and, eventually, improper integrals and differential equations in subsequent chapters.

Because of the Fundamental Theorems of Calculus, the operation of determining antiderivatives shares the same notation as integration, namely

$$\int f(x) dx = F(x) + C,$$

where C is a **constant of integration** and $F'(x) = f(x)$. **Later, we may regard this elegant formula as the proverbial “wolf in sheep’s clothing”. Intrigued?**

Antiderivative of a Function

Definition 7.2. Let $f : I \rightarrow \mathbb{R}$ be a function defined on an interval I . An **antiderivative** of f on the interval I is a function $F : I \rightarrow \mathbb{R}$ such that

$$F'(x) = f(x)$$

for all $x \in I$.

Note: Because $\frac{dC}{dx} = 0$ for all constants C , if F is an antiderivative of f , then so is $F(x) + C$. In other words, if a function has one antiderivative, it has an infinite number of antiderivatives, and moreover, they all differ by a constant.

Example 7.3. In the list of functions (a) - (e) below, identify those functions that are antiderivatives of other functions in the same list.

$$(a) f_a(x) = x \cdot \ln(x) - x$$

$$(b) f_b(x) = \ln(x)$$

- (c) $f_c(x) = \ln(x^2)$
 (d) $f_d(x) = \cos(x^2)$
 (e) $f_e(x) = -2x \cdot \sin(x^2)$

Solution: This problem is acquainting us with how to check if a given function $F(x)$ is an antiderivative for another function $f(x)$ or not. We apply the definition of what it means to be an antiderivative and check whether $F'(x) = f(x)$ or $F'(x) \neq f(x)$. It is that simple.

- (a) $\frac{df_a(x)}{dx} = \ln(x) = f_b(x)$, and hence $f_a(x)$ is an antiderivative of $f_b(x)$. Why? We verify $(f_a(x))' = f_b(x)$ by applying the product rule to $x \cdot \ln(x)$ to obtain $\ln(x) + 1$. We then differentiate $-x$ to obtain -1 . Adding up the two derivatives gives us the answer.
- (b) $\frac{df_b(x)}{dx} = \frac{1}{x}$, and hence $f_b(x)$ is NOT an antiderivative of any function in the list. Why? We computed $(f_b(x))'$, and it is not equal to any of the functions in the list.
- (c) $\frac{df_c(x)}{dx} = \frac{2}{x}$, and hence $f_c(x)$ is NOT an antiderivative of any function in the list.
- (d) $\frac{df_d(x)}{dx} = -2x \sin(x^2) = f_e(x)$, and hence $f_d(x)$ is an antiderivative of $f_e(x)$.
- (e) $\frac{df_e(x)}{dx} = -2 \sin(x^2) - 4x^2 \cos(x^2)$, and hence $f_e(x)$ is NOT an antiderivative of any function in the list.

1 using Symbolics

```

2
3 # Create a custom function to make differentiation convenient
4 function deriv(f, x)
5     return expand_derivatives.(Differential(x)(f))
6 end
7
8 # Sample Usage
9 @variables x # Makes x a symbolic variable
10 f = x^3*cos(x) + x^2*sin(x)^2 + log(3x)
11 dfdx = deriv(f, x)
12
13 # Note we do not say f(x) =, just f =
14
15 fa = x * log(x) - x
16 fb = log(x)
17 fc = log(x^2)
18 fd = cos(x^2)
19 fe = -2x * sin(x^2);
```

Output

(nothing)

1 using Symbolics

```

2
3 @show dfadx = deriv(fa, x)
4
5 @show dfbdx = deriv(fb, x)
6
7 @show dfcdx = deriv(fc, x)
8
```

```

9 @show dfddx = deriv(fd, x)
10
11 @show dfedx = deriv(fe, x);

```

Output

```

dfadx = deriv(fa, x) = log(x)
dfbdx = deriv(fb, x) = 1 / x
dfcdx = deriv(fc, x) = 2 / x
dfddx = deriv(fd, x) = -2x * sin(x^2)
dfedx = deriv(fe, x) = -2sin(x^2) - 4(x^2) * cos(x^2)

```

How to verify that $F(x)$ is an antiderivative of $f(x)$? Simply apply the definition: differentiate $F(x)$ and compare the result to $f(x)$. Hence, given a candidate antiderivative for $f(x)$, verifying or refuting the candidate antiderivative is straightforward. We'll see that finding candidate antiderivatives is, for the most part, anything but straightforward.

Example 7.4. ChatGPT4 (on its own, without an important plugin) has proposed the following antiderivatives for the functions in Example 7.3. Determine if they are correct or not.

- (a) For $f_a(x) = x \log(x) - x$: the antiderivative is $\int f_a(x) dx = \frac{1}{2}x^2 \log(x) - \frac{1}{4}x^2 + C$
- (b) For $f_b(x) = \log(x)$: the antiderivative is $\int f_b(x) dx = x \log(x) - x + C$
- (c) For $f_c(x) = \log(x^2)$: the antiderivative is $\int f_c(x) dx = x \log(x^2) - x + C$
- (d) For $f_d(x) = \cos(x^2)$: the antiderivative is $\int f_d(x) dx = \frac{1}{2}\sqrt{\pi}FresnelS\left(\sqrt{\frac{2}{\pi}}x\right) + C$
- (e) For $f_e(x) = -2x \sin(x^2)$: the antiderivative is $\int f_e(x) dx = \cos(x^2) + C$

Solution: The antiderivative for $f_d(x)$, as discussed in [Wikipedia](#), is a special function that we'll talk about later.

```

1 using Symbolics
2
3 # Define the symbolic variable
4 @variables x
5
6 # Define the functions as per their original definitions
7 fa(x) = x * log(x) - x
8 fb(x) = log(x)
9 fc(x) = log(x^2) # Which is equivalent to 2*log(x)
10 fd(x) = cos(x^2)
11 fe(x) = -2x * sin(x^2)
12
13
14 # Antiderivatives proposed by ChatGPT on 30 March 2025.
15 # The constants of integration have been removed.
16
17 AntiDeriv_fa = (x^2 * log(x)) / 2 - (3x^2) / 4 # Updated to match fa(x)
18 AntiDeriv_fb = x * log(x) - x
19 AntiDeriv_fc = 2x * log(x) - 2x # Correct antiderivative for fc(x)
20 AntiDeriv_fe = cos(x^2)
21
22 # Perform differentiation and compare with the original functions

```

```

23 result_fa = simplify(deriv(AntiDeriv_fa, x))
24 result_fb = simplify(deriv(AntiDeriv_fb, x))
25 result_fc = simplify(deriv(AntiDeriv_fc, x))
26 result_fe = simplify(deriv(AntiDeriv_fe, x))

27
28 println("Derivative of AntiDeriv_fa compared to fa(x) : ", result_fa, " == ", fa(x))
29 println("Derivative of AntiDeriv_fb compared to fb(x) : ", result_fb, " == ", fb(x))
30 println("Derivative of AntiDeriv_fc compared to fc(x) : ", result_fc, " == ", fc(x))
31 println("Derivative of AntiDeriv_fe compared to fe(x) : ", result_fe, " == ", fe(x))

```

Output

Derivative of AntiDeriv_fa compared to fa(x) : $x \log(x) - x == x \log(x) - x$
 Derivative of AntiDeriv_fb compared to fb(x) : $\log(x) == \log(x)$
 Derivative of AntiDeriv_fc compared to fc(x) : $2\log(x) == \log(x^2)$
 Derivative of AntiDeriv_fe compared to fe(x) : $-2x \sin(x^2) == -2x \sin(x^2)$

Using logarithm rules, $\ln(x^2) = 2\ln(x)$, and hence the third one is correct. In October 2023, ChatGPT got two of them wrong. Much progress has been made since then.

Secrets of the Arcane 7.5: Pro Tip for Determining Antiderivatives

The Symbolics.jl package can compute simple antiderivatives. The competing package, SymPy.jl, is somewhat better. However, the hands-down winners are MapleSoft and the tools from Wolfram, namely Wolfram Alpha Pro and Mathematica. Because of its ease of use, ChatGPT is your author's first recourse, but he always checks the result!.

Example 7.6. Compute the antiderivative for

$$f(x) := \frac{2x^3 \cos(x) \sin(x) e^{2+\sin^2(x)}}{1 + (e^{2+\sin^2(x)})^2} + 3x^2 \arctan(e^{2+\sin^2(x)}).$$

Solution: ChaptGPT-4 failed. Both Wolfram Alpha Pro and ChatGPT-o3mini-high computed

$$x^3 \arctan(e^{\frac{1}{2}(5-\cos(2x))}) + C.$$

Moreover, the answer is correct because the problem was created by differentiating $F(x) := x^3 \cdot \text{atan}(\exp(\sin(x)^2 + 2)) \cdot \log(x^2 + 1)$ with respect to x . The trig identity, $\sin^2(x) = \frac{1-\cos(2x)}{2}$, allows one to show that $\sin(x)^2 + 2 = \frac{(5-\cos(2x))}{2}$. Few humans would want to take on this “integration problem” using the tools taught in Calc II. **What ChatGPT-o3mini-high accomplished is impressive.**

Example 7.7. Compute the antiderivative for

$$\begin{aligned}
 f(x) &= 3x^2 \sqrt{2 + (\log(1 + x^2))^2} \arcsin(e^{2+\sin^2(x)}) \\
 &\quad + \frac{2x^4 \log(1 + x^2) \arcsin(e^{2+\sin^2(x)})}{\sqrt{2 + (\log(1 + x^2))^2} (1 + x^2)} \\
 &\quad + \frac{2x^3 \cos(x) \sin(x) e^{2+\sin^2(x)} \sqrt{2 + (\log(1 + x^2))^2}}{\sqrt{1 - (e^{2+\sin^2(x)})^2}}.
 \end{aligned} \tag{7.1}$$

Solution: ChatGPT4 failed. ChatGPT-o3mini-high computed

$$f(x) = \frac{d}{dx} \underbrace{\left(x^3 \cdot \arcsin \left(e^{2.0+\sin^2(x)} \right) \cdot \sqrt{2 + (\log(1+x^2))^2} \right)}_{F(x)}, \quad (7.2)$$

and hence, a “relatively simple” antiderivative $F(x)$ exists. ■

7.2 Uniting Integration and Differentiation through the Fundamental Theorems of Calculus

The **First Fundamental Theorem of Calculus** states that every continuous function has an antiderivative. The **Second Fundamental Theorem of Calculus** states that all definite integrals can be evaluated through antiderivatives. Together, these two results seamlessly tie together integration and differentiation as “inverse operations”.

7.2.1 Fundamental Theorems

Proposition 7.8: Fundamental Theorems of Calculus

First Fundamental Theorem of Calculus: Let $f : [a, b] \rightarrow \mathbb{R}$ be a continuous real-valued function defined on a closed interval $[a, b]$. Then, the function $F : [a, b] \rightarrow \mathbb{R}$ defined by

$$F(x) = \int_a^x f(t) dt \quad (7.3)$$

for all $x \in [a, b]$, is continuous on $[a, b]$, differentiable on the open interval (a, b) , and $F'(x) = f(x)$ for all $x \in (a, b)$. You will also see it written as

$$\frac{d}{dx} F(x) = \frac{d}{dx} \left[\int_a^x f(t) dt \right] = f(x), \quad (7.4)$$

to emphasize that differentiation undoes, or inverses, integration.

Second Fundamental Theorem of Calculus: Let $f : [a, b] \rightarrow \mathbb{R}$ be a real-valued function defined on a closed interval $[a, b]$ that admits an antiderivative F on $[a, b]$, that is, $F : [a, b] \rightarrow \mathbb{R}$ is differentiable on (a, b) and for all $x \in (a, b)$, $F'(x) = f(x)$. Then, f is integrable on $[a, b]$ and

$$\int_a^b f(x) dx = F(x) \Big|_a^b := F(b) - F(a). \quad (7.5)$$

You will also see it written as

$$\int_a^b \frac{dF(x)}{dx} dx = F(x) \Big|_a^b := F(b) - F(a), \quad (7.6)$$

to underline that integration undoes, or inverses, differentiation.

Notes: For a nuanced presentation of the First Fundamental Theorem of Calculus, we recommend [3 Levels of Proving the Fundamental Theorem of Calculus](#) by Trivial. The first part of the video [The Essence of \(Multivariable\) Calculus](#) by Foolish Chemist provides an intuitive treatment of the Second Fundamental Theorem of Calculus. If you exceed the 6:20 mark, you’ll be into multivariable versions of the Fundamental Theorems. Yikes!

7.2.2 Using the Second Fundamental Theorem of Calculus for Definite Integration

Let’s roll up our sleeves and get to work using the Second Fundamental Theorem of Calculus to evaluate some definite integrals. It looks super easy, right? In Table 5.15 we had a list of functions and their derivatives. Based on the First Fundamental Theorem of Calculus, we can repurpose our table to give us a list of antiderivatives, like so,

Function	C	$\frac{x^{k+1}}{k+1} + C$	$e^x + C$	$\ln x + C$	$\sin(x) + C$	$\cos(x) + C$	$\tan(x) + C$	$\text{atan}(x) + C$	Antiderivative
Derivative	0	x^k	e^x	$\frac{1}{x}$	$\cos(x)$	$-\sin(x)$	$1 + \tan^2(x)$	$\frac{1}{1+x^2}$	Function

Table 7.9: Common Elementary Functions and their Antiderivatives (Worth committing to memory to showcase your expertise). When doing antiderivatives, it is more convenient to have the monomials presented as given here than how we did the monomial entry in Table 5.15. In addition, we've extended the natural logarithm to negative numbers via an absolute value. Otherwise, the two tables are the same.

Why is there an absolute value on the natural logarithm? Note that $\frac{1}{x}$ is defined for all $x \neq 0$. Hence, its domain can be $(0, \infty)$, the positive reals, or $(-\infty, 0)$, the negative reals. Logarithms, on the other hand, are only defined for the positive reals. When $x > 0$, $|x| = x$, and hence, $\frac{d}{dx} \ln|x| = \frac{d}{dx} \ln(x) = \frac{1}{x}$. When $x < 0$, $|x| = -x$, and then, by the Chain Rule,

$$\frac{d}{dx} \ln|x| = \frac{d}{dx} \ln(-x) = \frac{1}{-x} \cdot (-1) = \frac{1}{x},$$

where the (-1) comes from differentiating $-x$ with respect to x . While it is never wrong to include the absolute value sign in the antiderivative of $\frac{1}{x}$, when the context makes it clear that $x > 0$, it is standard practice to drop it because, then, $|x| = x$.

You can also think about the derivative of the natural logarithm as follows. From the Chain Rule, $\frac{d}{dx} \ln(g(x)) = \frac{1}{g(x)} \cdot \frac{dg(x)}{dx}$, for $g(x) > 0$. Substituting in $g(x) = |x|$ and $x \neq 0$, yields, $\frac{d}{dx} \ln(|x|) = \frac{1}{|x|} \cdot \frac{d}{dx} |x|$. When $x > 0$, $\frac{d}{dx} |x| = 1$ and when $x < 0$, $\frac{d}{dx} |x| = -1$. Hence,

$$\frac{1}{|x|} \cdot \frac{d}{dx} |x| = \begin{cases} \frac{1}{x} \cdot 1 & x > 0 \\ \frac{1}{-x} \cdot (-1) & x < 0 \end{cases} = \frac{1}{x} \quad \text{for } x \neq 0.$$

Example 7.10. On the basis of Table 7.9 and the Fundamental Theorems of Calculus, compute the following definite integrals.

$$(a) \int_1^5 (2x^3 + e^x) dx$$

$$(b) \int_1^e \left(\frac{1}{x} + \frac{1}{1+x^2} \right) dx$$

$$(c) \int_0^\pi (\sin(x) + 3\cos(x)) dx$$

$$(d) \int_0^1 \left(\frac{x}{1+x^2} \right) dx$$

Solution:

(a) The antiderivative of x^3 is $\frac{x^4}{4}$ and the antiderivative of e^x is itself, e^x . Hence,

$$\int_1^5 (2x^3 + e^x) dx = \underbrace{\left(2 \cdot \frac{x^4}{4} + e^x \right)}_{\substack{\text{antideriv. } F(x) \\ \text{limits}}} \Big|_{a=1}^{b=5} = \underbrace{\left(\frac{5^4}{2} + e^5 \right)}_{F(b)} - \underbrace{\left(\frac{1^4}{2} + e^1 \right)}_{F(a)} = 312 + e^5 - e \approx 457.69.$$

(b) The antiderivative of $\frac{1}{x}$ is $\ln(x)$ and the antiderivative of $\frac{1}{1+x^2}$ is $\text{atan}(x)$. Hence,

$$\begin{aligned} \int_1^e \left(\frac{1}{x} + \frac{1}{1+x^2} \right) dx &= \underbrace{(\ln(x) + \text{atan}(x))}_{\substack{\text{antideriv. } F(x) \\ \text{limits}}} \Big|_{a=1}^{b=e} \\ &= \underbrace{(\ln(e) + \text{atan}(e))}_{F(b)} - \underbrace{(\ln(1) + \text{atan}(1))}_{F(a)} = 1 + \text{atan}(e) - \text{atan}(1) \approx 1.43. \end{aligned}$$

(c) The antiderivative of $\sin(x)$ is $-\cos(x)$ and the antiderivative of $\cos(x)$ is $\sin(x)$. Hence,

$$\int_0^\pi (\sin(x) + 3\cos(x)) \, dx = (-\cos(x) + 3\sin(x)) \Big|_0^\pi = \underbrace{(-\cos(x) + 3\sin(x))}_{\text{antideriv. } F(x)} \Big|_{a=0}^{b=\pi} = \underbrace{(-\cos(\pi) + 3\sin(\pi))}_{F(b)} - \underbrace{(-\cos(0) + 3\sin(0))}_{F(a)} = 2.$$

(d) The antiderivative of $\frac{x}{1+x^2}$ is NOT in Table 7.9. Hence, we either have to compute it via ChatGPT+Wolfram or learn a more powerful means for generating antiderivatives. From ChatGPT+Wolfram, we have the antiderivative of $\frac{x}{1+x^2}$ is $0.5 \log(x^2+1)$. Hence,

$$\int_0^1 \left(\frac{x}{1+x^2} \right) dx = 0.5 \log(x^2+1) \Big|_0^1 = 0.5 \log(2) - 0.5 \log(1) = 0.5 \log(2) \approx 0.34657$$

```

1 # Our prior means of computing definite integrals as a double check
2 using QuadGK
3 ansA,   = quadgk(x -> 2x^3 + exp(x), 1, 5)
4 ansB,   = quadgk(x ->(1/x) + (1/(1+x^2)), 1, exp(1))
5 ansC,   = quadgk(x -> sin(x) + 3cos(x), 0, pi)
6 ansD,   = quadgk(x ->x/(1+x^2), 0, 1);
7 [ansA, ansB, ansC, ansD]
```

Output

```
4-element Vector{Float64}:
457.69487727411763
1.4328847416198296
2.0000000000000001
0.3465735902799727
```

Taking Stock: For simple, “everyday” functions, it is relatively easy to use the Second Fundamental Theorem of Calculus for evaluating definite integrals. It can be a bit tedious because we have to find the antiderivative for each term in the integral, and then we have to evaluate the antiderivatives at the limits of integration, but conceptually, it is straightforward. Obtaining the antiderivatives is clearly the most difficult part of this process. For some super common functions, we have the antiderivatives available in a standard table, or even memorized. For less common functions, we have to determine the antiderivatives ourselves, one way or another. In Julia, the process of determining the antiderivatives and plugging in values is certainly harder than just running `quadgk` on the original functions. **So, when should we compute antiderivatives to evaluate integrals?**

- When you seek to establish a general result about an integral, having the closed-form solution provided by an antiderivative is highly useful. What is the volume of a sphere of radius $r > 0$? `Quadgk` can only compute the result for a specific value of r , and not for $r > 0$ a variable.
- In real-time coding for robotics, your available computing resources may be too slow to run `quadgk` over and over. Being able to compute the answer to an integral in closed form may provide for ultra-fast function evaluation. In this case, knowledge (i.e., Fundamental Theorems of Calculus) is power!
- **When you do not have Julia handy, such as in an exam.**
- When a HW problem in an engineering course asks you to integrate a function and show your work, 99.9% of the time, you are expected to use antiderivatives and the Fundamental Theorems of Calculus because that is what the instructors and

teaching assistants learned. It will take time for the Education Community to accept numerical computation in the classroom. In engineering practice, don't worry. Industry is all over this.

- When you want to enter an **Integration Bee**. Just like a **Spelling Bee**, it takes a lot of practice and rote memorization. Just as being a Spelling Bee Champ does not put you on the road to becoming a creative writer, being an Antiderivative Champ does not put you on the road to becoming a creative engineer or mathematician.

7.2.3 Conflating Integration and Antiderivatives is a Pedagogical Pitfall in Calculus

In this book, we have endeavored to maintain a clear distinction between definite integrals and integration via antiderivatives. For starters, we treat them in separate chapters. However, in broader discourse (aka, the big-bad world outside the course), the distinction might not always be as clear-cut, and the sense of the action verb “to integrate” must be extracted from the surrounding context. The notation of calculus, while often elegant and succinct, is not helpful in this regard.

In the context of definite integrals, the integral sign \int is accompanied by limits of integration, as in $\int_a^b f(x) dx$, indicating the interval over which the function f is to be integrated. However, when we venture into the domain of antiderivatives or indefinite integrals, the integral sign stands alone,

$$\int f(x) dx,$$

and represents a family of functions whose derivative is f . An alternative notation for denoting antiderivatives, albeit less common, is to use the differentiation symbol with a negative exponent, as in

$$D^{-1}f \quad \text{or} \quad \frac{d^{-1}}{dx^{-1}} f(x).$$

This notation explicitly mirrors the inverse operation being performed relative to differentiation. Yet, it's not widely adopted, perhaps due to its unwieldiness or the potential confusion with reciprocal functions. Because the notation has yet to “go viral” —three hundred years and waiting—we do not use it.

The transition from $\int_a^b f(x) dx$ to $\int f(x) dx$ marks the evolution from a concrete computation to a more abstract concept, mirroring the increase in abstraction from arithmetic to algebra. On one hand, antidifferentiation does amplify the power of Calculus. On the other hand, all too often, the numerous methods for computing antiderivatives become the “be all and end all” of a standard Calc II course, leading many students to throw up their hands and walk away.

How could that be? As we'll explain in the next Section, the methods for computing antiderivatives are based on inverting the differentiation rules in Prop. 5.11. Becoming proficient in them takes a lot of effort and is typically accompanied by deep frustration on the part of the learner. The frustration arises because (a) knowing when to use one method over another requires significant experience; (b) each method typically involves multiple steps where a simple minus sign error or confusing a multiply by a constant with a divide by the same constant destroys the entire calculation, and finally, (c), once you have computed the antiderivative, beyond the satisfaction of having solved a cool puzzle, you rarely get to use it for **something meaningful**. In real engineering, you are unlikely to use the hand methods of computing antiderivatives with sufficient frequency to maintain fluency, supposing you master them while in College.

You were once taught how to extract a square root by hand. How often do you do that now? **This textbook takes the attitude that while there was a time when it made sense to learn how to execute by hand all of the methods for computing antiderivatives, that time has passed.** Yes, it's almost always true that “knowledge is power”. Understanding the concepts underlying the “hand methods of integration” is sometimes useful. In Linear algebra, however, once we understood the basic workings of a determinant or the LU Factorization, we were happy to turn the actual computations over to an algorithm; the same is true of the rules for computing antiderivatives and definite integrals: we want an algorithm.

The next section will work through enough methods and examples so that when you are required to compute antiderivatives by hand in an engineering course, you are prepared.

7.3 The Art of the Antiderivative: Inverting Differentiation Rules to Find Antiderivatives

We go through, roughly in order of importance, various strategies for finding antiderivatives of functions.

7.3.1 The Power Rule of Integration: Inverts the Power Rule of Differentiation

From the power rule for differentiation, we have that $\frac{d}{dx}x^{k+1} = (k+1)x^k$. Hence, by the Second Fundamental Theorem of Calculus,

$$\int x^k dx = \frac{x^{k+1}}{k+1} + C. \quad (7.7)$$

Example 7.11. Compute an antiderivative for $4x^2$.

Solution: By the Power Rule, $\int 4x^2 dx = 4\left(\frac{x^3}{3} + C_1\right) = \frac{4}{3}x^3 + C$, where $C := 4C_1$. ■

7.3.2 The Fundamental Rule: Integrating the Differential

When doing an integral by hand that is harder than a monomial, our first thought should be: can we recognize the integrand (aka, the function being integrated) as being the derivative of a commonly known function. The most fundamental relationship between differentiation and integration is given by the total differential of a function,

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

The Second Fundamental Theorem of Calculus says that integrating both sides with respect to x gives us back the original function (up to a constant of integration):

$$\int df(x) = \int \frac{df(x)}{dx} dx = \int f'(x) dx = f(x) + C. \quad (7.8)$$

Example 7.12. Compute an antiderivative for $3x^2$.

Solution: Here, we recognize that $x^2 = \frac{d}{dx}\left(\frac{x^3}{3}\right) = \left(\frac{x^3}{3}\right)'$. Therefore, the Second Fundamental Theorem of Calculus gives us

$$\int 3x^2 dx = \int 3\left(\frac{x^3}{3}\right)' dx = 3\frac{x^3}{3} + C = x^3 + C.$$

Example 7.13. Determine $\int 13(1 + \tan^2(x)) dx$.

Solution: This is an antiderivative problem. From Table 7.9, we recognize that $1 + \tan^2(x) = \frac{d}{dx}(\tan(x)) = (\tan(x))'$. Therefore, the Second Fundamental Theorem of Calculus gives us

$$\int 13(1 + \tan^2(x)) dx = \int 13(\tan(x))' dx = 13\tan(x) + C.$$

Admittedly, these are not impressive integrals, but as the saying goes, do not look a gift horse in the mouth.

7.3.3 Inverting the Chain Rule: Integration by Substitution, aka u-Substitution

The chain rule for differentiation states that if $f'(x) = \frac{df(x)}{dx}$ and $g'(x) = \frac{dg(x)}{dx}$ exist, then

$$\frac{d}{dx}(f(g(x))) = (f(g(x)))' = f'(g(x)) \cdot g'(x).$$

To invert this, we use **integration by substitution**. If we define $u = g(x)$, then $du = \frac{dg(x)}{dx} dx = g'(x) dx$, and

$$\int \underbrace{f'(g(x)) \cdot g'(x)}_{f'(u)} dx = \int f'(u) du = f(u) + C = f(g(x)) + C, \quad (7.9)$$

after substituting in $u = g(x)$. In other words, by a clever substitution, we turn something that did not look like a total differential into something that is a recognizable total differential so that the Fundamental Rule in (7.8) can be applied. Yeah, computing $\int f'(u) du$ is easy, while $\int f'(g(x)) \cdot g'(x) dx$ looks hard before making the substitution. This is an example of useful pattern matching: looking for instances of the chain rule appearing as the integrand (aka, under the integral sign).

Remark 7.14. (*u*-Substitutions) It is common practice to use $u(x)$ instead of $g(x)$ AND to call this the Method of *u*-Substitution. Why u ? When you find a satisfying answer, let us all know.

Example 7.15. Use integration to find $\int 3x \cos(x^2 + 1) dx$.

Solution: Integrate, here, means to find an antiderivative. Table 7.9 does not contain an antiderivative for $3x \cos(x^2 + 1)$. However, we see that the cosine function has been composed with $x^2 + 1$, and lucky for us, its derivative, $2x$, is ALMOST sitting in front of the function. Inspired by this, we let $f'(x) = \cos(x)$ and $u(x) = x^2 + 1$ so that $du = 2x dx$. Then,

$$\int 3x \cos(x^2 + 1) dx = \frac{3}{2} \cdot \int 2x \cos(x^2 + 1) dx = \frac{3}{2} \cdot \int \cos(u) du = \frac{3}{2} \cdot \sin(u) + C = \frac{3}{2} \cdot \sin(x^2 + 1) + C,$$

after substituting in $u = x^2 + 1$. ■

Example 7.16. Compute the definite integral $\int_0^\pi 2x \cos(x^2 + 1) dx$.

Solution: We give two solutions, one that obliges us to complete the last step of substituting $u = u(x)$ back into $f(u)$ and one that allows us to avoid that step.

Approach One: From example 7.15, we know the antiderivative of $2x \cos(x^2 + 1)$ is $F(x) := f(u(x)) = \sin(x^2 + 1) + C$. By the Second Fundamental Theorem of Calculus,

$$\int_0^\pi 2x \cos(x^2 + 1) dx = F(\pi) - F(0) = (\sin(\pi^2 + 1) + C) - (\sin(0^2 + 1) + C) = \sin(\pi^2 + 1) - \sin(1).$$

Because the constant of integration C always cancels out, it is common not even to write it down when using antiderivatives for definite integrals.

Approach Two: The limits of integration $a = 0$ and $b = \pi$ really mean the x -lower limit of integration is 0 and the x -upper limit of integration is π . When we perform a change of variables $u(x)$, we can also change the limits of integration to present the u -lower limit of integration and the u -upper limit of integration, which are $u(a) = u(0) = \sin(1)$ and $u(b) = u(\pi) = \sin(\pi^2 + 1)$. Hence, we'd compute the integral as

$$\int_0^\pi 2x \cos(x^2 + 1) dx = \int_{u(0)}^{u(\pi)} \cos(u) du = \sin(u) \Big|_1^{\pi^2+1} = \sin(\pi^2 + 1) - \sin(1).$$

When the antiderivative is more complicated, the second approach becomes more valuable. In addition, when the first substitution only gets you halfway there, and a second substitution is required, this second approach really shines. It's also common to add extra notation to help one keep track of the limits carefully when doing a change of variable, such as

$$\int_{x=0}^{x=\pi} 2x \cos(x^2 + 1) dx = \int_{u=u(0)}^{u=u(\pi)} \cos(u) du = \sin(u) \Big|_{u=1}^{u=\pi^2+1} = \sin(\pi^2 + 1) - \sin(1);$$

this helps avoid a really bad error, where the original limits are carried forward as in

$$\int_0^\pi 2x \cos(x^2 + 1) dx = \int_0^\pi \cos(u) du = \sin(u) \Big|_0^\pi = \sin(\pi) - \sin(0) = 0,$$

and everything in red is wrong! ■

Example 7.17. Suppose that $f(x)$ is the antiderivative of $f'(x)$. Find the antiderivative of $f'(\alpha x - x_0)$, where $\alpha \neq 0$ and $x_0 \in \mathbb{R}$ are constants.

Solution: We define $u(x) := \alpha x - x_0$ and note that $f'(u(x)) = f'(\alpha x - x_0)$. Moreover, $du(x) = \alpha dx \implies dx = \frac{1}{\alpha} du$. Hence,

$$\int f'(\alpha x - x_0) dx = \int f'(u) \frac{1}{\alpha} du = \frac{1}{\alpha} f(u) + C = \frac{1}{\alpha} f(\alpha x - x_0) + C. \quad (7.10)$$

Using the above result, we can immediately expand our table of known antiderivatives to include scaling and shifting.

Antiderivative	$\frac{1}{\alpha} \frac{(\alpha x - x_0)^{k+1}}{k+1}$	$\frac{1}{\alpha} e^{(\alpha x - x_0)}$	$\frac{1}{\alpha} \ln \alpha x - x_0 $	$\frac{1}{\alpha} \sin(\alpha x - x_0)$	$\frac{1}{\alpha} \cos(\alpha x - x_0)$	$\frac{1}{\alpha} \tan(\alpha x - x_0)$	$\frac{1}{\alpha} \operatorname{atan}(\alpha x - x_0)$
Function	$(\alpha x - x_0)^k$	$e^{(\alpha x - x_0)}$	$\frac{1}{(\alpha x - x_0)}$	$\cos(\alpha x - x_0)$	$-\sin(\alpha x - x_0)$	$1 + \tan^2(\alpha x - x_0)$	$\frac{1}{1+(\alpha x - x_0)^2}$

Table 7.18: Common Elementary Functions and their Antiderivatives. The constants of integration, C , have been dropped to save space.

Example 7.19. Compute the indefinite integral $\int \frac{(\ln|x|)^2}{x} dx$.

Solution: Ans. $\int \frac{(\ln|x|)^2}{x} dx = \frac{(\ln(x))^3}{3} + C$ for any domain $[a, b]$ that excludes $x = 0$.

The problem is using different vocabulary for the antiderivative so that we get used to the idea that for most students of Calculus, integration and antiderivation are one and the same, two peas in the same pod, identical twins; you get the picture. Once again, we cannot find the integrand in Table 7.9. If we take $f'(x) = x^2$ and $u(x) = \ln|x|$, then $u'(x) = \frac{1}{x}$ and $du = \frac{1}{x} dx$. Hence, we have that

$$\frac{(\ln|x|)^2}{x} dx = u^2 du = f'(u) du$$

Because we know that the antiderivative of $f'(u) = u^2$ is $f(u) = \frac{u^3}{3}$, we have

$$\int \frac{(\ln|x|)^2}{x} dx = \int u^2 du = \frac{u^3}{3} + C = \frac{(\ln|x|)^3}{3} + C.$$

Secrets of the Arcane 7.20: “Danger, Will Robinson!” is 1960’s Slang for “Red Flag!”

Because **Danger, Will Robinson!** is tied to a robot, we use it in this textbook.

The phrase "Danger, Will Robinson" is a classic line from the 1960s TV show "**Lost in Space**" used by the Robot to alert young Will Robinson to imminent danger. This iconic phrase has become synonymous with signaling a warning or drawing attention to a potential problem. It's a part of pop culture lexicon that transcends generations, often used humorously or dramatically to indicate caution in various contexts.

In contemporary terms, expressions like "Red flag!" have gained popularity, especially on social media and in daily communication, to signal caution or alert someone to potential issues, dangers, or problematic situations. The term "red flag" is widely used to describe warning signs in relationships, behaviors, or situations that may be considered suspicious, dangerous, or indicative of deeper issues.

While "Danger, Will Robinson" harks back to a specific context of science fiction and the protective instincts of a robotic character, "Red Flag" is more universally applied across various scenarios, including personal interactions, professional environments, and online discourse. Both serve a similar purpose in signaling caution but reflect different cultural contexts and eras.

Fact 7.21: Danger, Danger Will Robinson!

When computing antiderivatives, it is easy to forget about the domains on which the integrand and the antiderivative are defined. The integrand $\frac{(\ln|x|)^2}{x}$ is defined and continuous on domains $[a, b]$ that do not include $x = 0$. Hence, its Riemann integral exists and will be finite for $[a, b] \subset (-\infty, 0)$ or $[a, b] \subset (0, \infty)$. What about the antiderivative for $\frac{(\ln|x|)^2}{x}$? It should exist and be finite in this same domain, and in fact, we observe that $\frac{(\ln|x|)^3}{3}$ exists and is finite for $[a, b] \subset (-\infty, 0)$ and $[a, b] \subset (0, \infty)$.

Where you can get into trouble: Suppose you want to evaluate the definite integral $\int_{-e}^e \frac{(\ln|x|)^2}{x} dx$. From our analysis above, we know the problem does not make sense because $0 \in [-e, e]$. However, it is ever so tempting to write

$$\int_{-e}^e \frac{(\ln|x|)^2}{x} dx = \frac{(\ln|x|)^3}{3} \Big|_{-e}^e = \frac{(\ln|e|)^3}{3} - \frac{(\ln|-e|)^3}{3} = 0 - 0 = 0,$$

which is wrong. The correct answer is: the definite integral is not defined. When computing definite integrals, it is easy to make “innocent” mistakes like this one. **It is especially easy to make errors on made-up problems where you are not personally invested in the answer.**

In case everything has seemed totally straightforward so far, this next example will recalibrate your expectations when it comes to finding antiderivatives by hand.

Example 7.22. (A Tricky Antiderivative) Find an antiderivative for $\sec(x) := \frac{1}{\cos(x)}$.

Solution: Ans. : $\int \sec(x) dx = \ln|\sec(x)| + \tan(x)| + C$, with its domain analyzed at the end of the solution.

To find the antiderivative of $\sec(x)$, we multiply and divide by¹ $\sec(x) + \tan(x)$ to create a convenient form for u -substitution,

$$\begin{aligned} \int \sec(x) dx &= \int \sec(x) \frac{\sec(x) + \tan(x)}{\sec(x) + \tan(x)} dx \\ &= \int \frac{\sec^2(x) + \sec(x)\tan(x)}{\sec(x) + \tan(x)} dx. \end{aligned}$$

Now, let $u = \sec(x) + \tan(x)$. Then, $du = (\sec(x)\tan(x) + \sec^2(x)) dx$, which is the numerator of the integrand. Thus, we can rewrite the integral in terms of u ,

$$\begin{aligned} \int \sec(x) dx &= \int \frac{1}{u} du \\ &= \ln|u| + C \\ &= \ln|\sec(x) + \tan(x)| + C, \end{aligned}$$

where C is the constant of integration.

Where did the idea to multiply and divide by $\sec(x) + \tan(x)$ come from? Probably from knowing way too many trig identities! Finding antiderivatives by hand is definitely an art form. As with any form of art, some of you will love it, and others hate it.

After all of this, we need to ask ourselves, in what domain does the integration problem make sense, and the same for the proposed answer.

- We cannot compute an integral over an interval $[a, b]$ where $\sec(x) := \frac{1}{\cos(x)}$ blows up. Hence, we need to avoid intervals $[a, b]$ that contain $\frac{\pi}{2} + k\pi$, for some $k \in \mathbb{Z}$.
- Does the antiderivative blow up anywhere? The natural logarithm blows up at $x = 0$. Now, $\sec(x) + \tan(x) = \frac{1}{\cos(x)} + \frac{\sin(x)}{\cos(x)} = \frac{1+\sin(x)}{\cos(x)}$, and we note that $1 + \sin(x) = 0$ for $x = -\frac{\pi}{2} + 2k\pi$, $k \in \mathbb{Z}$, which is a subset of where $\cos(x)$ vanishes. What about the

¹Where did that come from? And are we multiplying and dividing by zero at any point?

ratio? Applying L'Hôpital's Rule, we have

$$\lim_{x \rightarrow -\frac{\pi}{2}} \frac{1 + \sin(x)}{\cos(x)} = \lim_{x \rightarrow -\frac{\pi}{2}} \frac{\cos(x)}{-\sin(x)} = 0.$$

Hence, the antiderivative makes sense everywhere that the integrand makes sense, which is how it should be. ■

This and the next example come from [Advanced Integration by U-Substitution - Calculus II](#) by the Calculus God.

Example 7.23. Determine $\int \cos^3(x) dx$.

Solution: We'll give two solutions.

Option 1: Sine and cosine raised to a power can always be expressed as a sum of sine and cosine without any powers; see Prop. 1.75. By the identity for the cosine of a triple angle,

$$\cos(3x) = 4\cos^3(x) - 3\cos(x).$$

Solving for $\cos^3(x)$ gives

$$\cos^3(x) = \frac{1}{4}\cos(3x) + \frac{3}{4}\cos(x).$$

Hence, $\int \cos^3(x) dx = \frac{1}{12}\sin(3x) + \frac{3}{4}\sin(x)$.

Option 2 (by the Calculus God): $\int \cos^3(x) dx = \int \cos(x) \cdot \cos^2(x) dx = \int \cos(x) \cdot (1 - \sin^2(x)) dx = \int \cos(x) dx - \int \cos(x) \sin^2(x) dx$. The first antiderivative is easy; the second one looks harder.

Let $u = \sin(x)$ so that $du = \cos(x) dx$. Then,

$$\begin{aligned} \int \cos^3(x) dx &= \sin(x) - \int u^2 du \\ &= \sin(x) - \frac{u^3}{3} \\ &= \sin(x) - \frac{1}{3} \cdot \sin^3(x). \end{aligned}$$

The two proposed antiderivatives appear to be different, but by a trig substitution, $\sin^3(x) = \frac{3}{4}\sin(x) - \frac{1}{4}\sin(3x)$, which shows that $\sin(x) - \frac{1}{3} \cdot \sin^3(x) = \frac{1}{12}\sin(3x) + \frac{3}{4}\sin(x)$. ■

Example 7.24. Determine $\int \frac{x}{\sqrt{x+1}} dx$.

Solution: We'll give two solutions. We note that the integrand, $\frac{x}{\sqrt{x+1}}$, is a well defined function for $x+1 > 0$, which is equivalent to $x > -1$.

Option 1: For the u-substitution, we'll let $u = x+1$, which we note is strictly greater than zero. In addition, we have $du = dx$ and $x = u-1$. Doing the substitutions gives us

$$\begin{aligned} \int \frac{x}{\sqrt{x+1}} dx &= \int \frac{u-1}{\sqrt{u}} du \\ &= \int \frac{u-1}{u^{1/2}} du \\ &= \int u^{1/2} - u^{-1/2} du \end{aligned}$$

Applying the generalized power rule gives

$$\begin{aligned} \int \frac{x}{\sqrt{x+1}} dx &= \int u^{1/2} - u^{-1/2} du \\ &= \frac{u^{3/2}}{3/2} - \frac{u^{1/2}}{1/2} \\ &= \frac{2}{3}u^{3/2} - 2u^{1/2}. \end{aligned}$$

Substituting in $u = x + 1$ gives

$$\int \frac{x}{\sqrt{x+1}} dx = \frac{2}{3} \sqrt{(x+1)^3} - 2\sqrt{x+1}.$$

Option 2 (by the Calculus God): For the u-substitution, we'll let $u = \sqrt{x+1}$, and for later use, we note that $u > 0$. From $u = \sqrt{x+1}$, we obtain $u^2 = x+1$ and $2u du = dx$. Doing the substitutions gives us

$$\begin{aligned}\int \frac{x}{\sqrt{x+1}} dx &= \int \frac{u^2 - 1}{u} \cdot 2u du \quad (\text{note that the } u \text{ in the denominator}) \\ &= \int 2(u^2 - 1) du \quad (\text{cancels with the } u \text{ in } 2u du) \\ &= 2\left(\frac{u^3}{3} - u\right).\end{aligned}$$

Plugging back in for x , we have

$$\int \frac{x}{\sqrt{x+1}} dx = \frac{2}{3} \left(\sqrt{(x+1)}\right)^3 - 2\sqrt{x+1},$$

which agrees with the answer in Option 1, because $\left(\sqrt{(x+1)}\right)^3 = \sqrt{(x+1)^3}$ by our power rules. ■.

Hopefully, it is now abundantly clear that the construction of antiderivatives can be done in many ways, and even for a given problem, there is not a single best way to work it. ChatGPT+Wolfram is still a good approach.

Recommended Videos for Additional Practice

- How to Integrate Using U-Substitution by Nancy Pi.
- Understand U-substitution, the Idea! by BlackPenRedPen.
- How To Integrate Using U-Substitution by the Organic Chemistry Tutor.
- U-substitution With Definite Integrals by the Organic Chemistry Tutor.

Danger, Will Robinson! When using any antiderivative, you must be careful to avoid making the “innocent” error highlighted in Fact. 7.21. It does not matter if you used ChatGPT+Wolfram or did the computation yourself by hand. You should always ask yourself: “What is an allowed interval of integration for the problem (aka, what is an allowed domain for the integrand)?” and “Does the answer make sense in that domain?” On YouTube, this is rarely, if ever, discussed. And when you have landed that nice internship or steady employment, it’s always a good idea to either differentiate the antiderivative to make sure you recover the integrand or compare to quadgk for a set of (definite) integration bounds.

7.3.4 Inverting the Product Rule: Integration by Parts

The product rule states that $\frac{d}{dx}(u \cdot v) = u \cdot \frac{dv}{dx} + v \cdot \frac{du}{dx}$. Or, in terms of the total differential, we have

$$d(u \cdot v) = u \cdot dv + v \cdot du.$$

Integration by parts is the inverse of the product rule, based on rearranging the above equation

$$u \cdot dv = d(u \cdot v) - v \cdot du.$$

Integrating both sides gives

$$\int u \, dv = u \cdot v - \int v \, du, \tag{7.11}$$

where $\int d(u \cdot v) = u \cdot v$, because $d(u \cdot v)$ is a total differential. This “trick” is useful when the integral you need to evaluate, $\int u \, dv$, is hard, but its other half coming from the product rule of differentiation, $\int v \, du$, is easier. The method is called **Integration by Parts** because you have broken the integral into two parts, one that is trivial $\int d(u \cdot v) = u \cdot v$, and one, $\int v \, du$, that you “hope” is easier than the original integral, $\int u \, dv$. We will return to this method when we study Laplace Transforms in Chapter 10.

Example 7.25. Solve for $\int x \sin(x) dx$.

Solution: For this first example, we do not have a lot of choices to make when defining $u(x)$ and $dv(x)$. One of them has to be associated with x and one has to be associated with $\sin(x)$. We explore both options.

Option 1: Let $u = x$ and $dv = \sin(x) dx$. Then $du = dx$ and $v = -\cos(x)$. It was very important to choose dv such that it has a known antiderivative. Armed with this information, we have

$$\int \underbrace{x \sin(x)}_{u \, dv} dx = \underbrace{-x \cos(x)}_{u \cdot v} + \int \underbrace{\cos(x) dx}_{v \, du} = -x \cos(x) + \sin(x) + C.$$

Option 2: Let $u = \sin(x)$ and $dv = x dx$. Then $du = \cos(x) dx$ and $v = \frac{x^2}{2}$. For this problem, it looks like both choices are yielding easy antiderivatives. Armed with this information, we proceed and find that

$$\int \underbrace{x \sin(x) dx}_{u \, dv} = \underbrace{\frac{x^2}{2} \sin(x)}_{u \cdot v} - \int \underbrace{\frac{x^2}{2} \cos(x) dx}_{v \, du}.$$

At this point, we go Ooops! We've created an even more challenging antiderivative to compute, namely, the antiderivative for $\frac{x^2}{2} \cos(x)$ has a higher power of x than we started with.

The second option illustrates another aspect of the **Art of the Antiderivative**; you have to make good choices for $u(x)$ and $v(x)$. In this problem, the x term is causing the problem. We know how to integrate $\sin(x)$. It is smart to choose $u(x) = x$ because when we differentiate $u(x)$, we get a constant, which puts us back to dealing with an integral with a simple antiderivative. ■.

Example 7.26. Find an antiderivative for $\ln(x)$.

Solution: Without an integral sign in sight, it's hard to even think about applying integration by parts, so we rewrite the problem as evaluate $\int \ln(x) dx$. We do not know how to integrate $\ln(x)$, but we sure know how to differentiate it. Hence, we set $u(x) = \ln(x)$ and $dv = dx$. The latter looks so ridiculously trivial that we hesitate even to pose it, right? But with these choices, $du = \frac{1}{x} dx$ and $v(x) = x$. Hence, our integration by parts formula becomes

$$\int \underbrace{\ln(x) dx}_{u \, dv} = \underbrace{\ln(x) \cdot x}_{u \cdot v} - \int \underbrace{x \cdot \frac{1}{x} dx}_{v \, du} = \ln(x) \cdot x - \int 1 dx = x \cdot \ln(x) - x + C.$$

■.

Example 7.27. Evaluate $\int x^3 \cdot e^x dx$.

Solution: This problem requires three successive applications of the integration by parts formula,

$$\int u \, dv = uv - \int v \, du$$

For the first application, we choose:

$$\begin{aligned} u_1 &= x^3 &\implies du_1 &= 3x^2 dx \\ dv_1 &= e^x dx &\implies v_1 &= e^x \end{aligned}$$

Applying the integration by parts formula:

$$\begin{aligned} \int x^3 e^x dx &= u_1 v_1 - \int v_1 du_1 \\ &= x^3 e^x - \int 3x^2 e^x dx \end{aligned}$$

For the second application, we choose:

$$\begin{aligned} u_2 &= x^2 &\implies du_2 &= 2x dx \\ dv_2 &= e^x dx &\implies v_2 &= e^x \end{aligned}$$

Applying the integration by parts formula again:

$$\begin{aligned}\int 3x^2 e^x dx &= 3(u_2 v_2 - \int v_2 du_2) \\ &= 3(x^2 e^x - \int 2x e^x dx)\end{aligned}$$

For the third application, we choose:

$$\begin{aligned}u_3 &= x &\implies du_3 &= dx \\ dv_3 &= e^x dx &\implies v_3 &= e^x\end{aligned}$$

Applying the integration by parts formula a third time:

$$\begin{aligned}\int 2x e^x dx &= 2(u_3 v_3 - \int v_3 du_3) \\ &= 2(xe^x - \int e^x dx) \\ &= 2xe^x - 2e^x\end{aligned}$$

Combining all the applications, we get:

$$\begin{aligned}\int x^3 e^x dx &= x^3 e^x - 3(x^2 e^x - 2(xe^x - e^x)) \\ &= x^3 e^x - 3x^2 e^x + 6x e^x - 6e^x + C\end{aligned}$$

Thus, the antiderivative of $x^3 e^x$ is

$$x^3 e^x - 3x^2 e^x + 6x e^x - 6e^x + C,$$

where C is the constant of integration. ■

Method 7.28: (Optional Read: DI Method or Integration by Parts on Steroids)

For those who wish to master problems requiring multiple applications of integration by parts, please see [Integration by Parts, DI Method, VERY EASY](#) by **BlackPenRedPen**. The DI Method is a tabular form of integration by parts that makes it easy to keep track of the signs and understand how the initial choices of u and dv propagate through successive steps of integration by parts. The D-column contains u , the function that is differentiated when applying integration by parts, while the I-column contains dv , the function that is to be integrated. After watching the video by **BlackPenRedPen**, you will be able to build a table like the one below and extract from it the antiderivative of $x^3 \cdot e^x$. Not only is the table very compact, but the alternating signs in the left column make it easy to avoid sign errors.

	D	I
+	x^3	e^x
-	$3x^2$	e^x
+	$6x$	e^x
-	6	e^x
+	0	e^x

$$\int x^3 e^x dx = x^3 e^x - 3x^2 e^x + 6x e^x - 6e^x + C$$

You can perfect your technique by watching a few more [DI Method Videos](#) by **BlackPenRedPen**.

Recommended Videos for Additional Practice

- [Integration by Parts... How?](#) by Nancy Pi.
- [Integration by Parts](#) by the Organic Chemistry Tutor.
- [Integration by Parts, Calculus 2, AP Calculus BC](#) by BlackPenRedPen.

7.3.5 Antiderivatives of Rational Functions by Partial Fraction Expansion (PFE)

Consider a rational function expanded as follows,

$$R(x) = \frac{k_1}{x - r_1} + \frac{k_2}{x - r_2} + \cdots + \frac{k_n}{x - r_n}, \quad (7.12)$$

for real “roots” $r_i \in \mathbb{R}$ and coefficients $k_i \in \mathbb{R}$. On the basis of Table 7.18, computing its antiderivative is straightforward,

$$\int R(x) dx = k_1 \ln|x - r_1| + k_2 \ln|x - r_2| + \cdots + k_n \ln|x - r_n|,$$

because $\frac{d}{dx} \ln|x - x_0| = \frac{1}{x - x_0}$.

Note: The antiderivative “blows up” when the argument of the natural logarithm is zero. Hence, the domain of the antiderivative is equal to $\{x \in \mathbb{R} \mid x \notin \{r_1, r_2, \dots, r_n\}\}$, which is also the domain of $R(x)$. Everything checks out.

Partial Fraction Expansion for Real Roots

If we were given (7.12) in the form $R(x) = \frac{P(x)}{Q(x)}$, then it would be far less obvious how to compute an antiderivative. As an illustration, if we start out with a few terms similar to (7.12), and put them over a common denominator,

$$\begin{aligned} R(x) &= \frac{1}{x-1} + \frac{2}{x+2} + \frac{3}{x-3} \\ &= \frac{(x+2)(x-3)}{(x-1)(x+2)(x-3)} + \frac{2(x-1)(x-3)}{(x-1)(x+2)(x-3)} + \frac{3(x-1)(x+2)}{(x-1)(x+2)(x-3)} \\ &= \frac{6x^2 - 6x - 6}{(x-1)(x+2)(x-3)} \\ &= \frac{6x^2 - 6x - 6}{x^3 - 2x^2 - 5x + 6} \\ &=: \frac{P(x)}{Q(x)}, \end{aligned} \quad (7.13)$$

it's even hard to know where to start. Finding the antiderivative of $\frac{6x^2 - 6x - 6}{x^3 - 2x^2 - 5x + 6}$ looks much more terrifying than finding the antiderivative of $\frac{1}{x-1} + \frac{2}{x+2} + \frac{3}{x-3}$.

The process called **Partial Fraction Expansion (PFE)** was invented to take a rational function like (7.13) and turn it into a sum of first-order terms, as in (7.12). We'll begin with a simplified version of partial fraction expansion that, like (7.13), works for rational functions satisfying the following two properties:

- $\deg(P(x)) < \deg(Q(x))$, the numerator polynomial has a strictly lower degree than the denominator polynomial.
- the roots of the denominator polynomial, $Q(x)$, are real and distinct (means no complex roots and no repeated roots).

For (7.13), we note that $2 = \deg(P(x)) < \deg(Q(x)) = 3$, and that the roots of $Q(x) = x^3 - 2x^2 - 5x + 6 = (x-1)(x+2)(x-3)$, being $r_1 = 1$, $r_2 = -2$, and $r_3 = 3$, are real and distinct.

Proposition 7.29: Antiderivative with a Simplified Partial Fraction Expansion (PFE)

Consider a rational function $R(x) = \frac{P(x)}{Q(x)}$ where

- the degree of $P(x)$ is strictly lower than the degree of $Q(x)$,
- $Q(x)$ is monic (aka, leading coefficient is one),
- the roots of $Q(x)$, $\{r_1, r_2, \dots, r_m\}$, are real and distinct, so that $Q(x) = \prod_{i=1}^m (x - r_i)$.

Then $R(x)$ always has a partial fraction expansion

$$R(x) = \frac{k_1}{x - r_1} + \frac{k_2}{x - r_2} + \dots + \frac{k_m}{x - r_m}, \quad (7.14)$$

where the coefficients k_1, k_2, \dots, k_m satisfy

$$k_i = \lim_{x \rightarrow r_i} \frac{(x - r_i)P(x)}{Q(x)} = \frac{P(r_i)}{(r_i - r_1) \cdots (r_i - r_{i-1})(r_i - r_{i+1}) \cdots (r_i - r_m)}. \quad (7.15)$$

Note: Because we assumed the roots to be distinct, once we canceled the term $(x - r_i)$ in its numerator and denominator,

$$\frac{(x - r_i)P(x)}{Q(x)} = \frac{(x - r_i)P(x)}{(x - r_1) \cdots (x - r_{i-1}) \cdot (x - r_i) \cdot (x - r_{i+1}) \cdots (x - r_m)},$$

the limit is simply the evaluation of the reduced polynomial at $x = r_i$.

Note: There is an alternative method for determining the coefficients k_i that involves setting up a system of linear equations. In the end, it is better to use software tools for anything more than a quadratic.

Example 7.30. Determine $\int \frac{1}{x^2 - 1} dx$.

Solution: We check to see if our simplified method of partial fraction expansion applies to the problem. The degree of the numerator is zero, which is strictly less than two, the degree of the denominator. The denominator factors as $x^2 - 1 = (x + 1) \cdot (x - 1)$, giving us distinct real roots, $r_1 = -1$ and $r_2 = 1$. Hence, we are good to go!

We can decompose this into partial fractions,

$$\frac{1}{x^2 - 1} = \frac{k_1}{(x - 1)} + \frac{k_2}{(x + 1)}.$$

To compute the coefficients in the partial fraction expansion,

- $k_1 = \lim_{x \rightarrow r_1} \frac{(x + 1)}{(x + 1)(x - 1)} = \lim_{x \rightarrow -1} \frac{1}{(x - 1)} = \frac{1}{(-1 - 1)} = -\frac{1}{2}$.
- $k_2 = \lim_{x \rightarrow r_2} \frac{(x - 1)}{(x + 1)(x - 1)} = \lim_{x \rightarrow +1} \frac{1}{x + 1} = \frac{1}{1 + 1} = \frac{1}{2}$.

Armed with the coefficients, we decompose the integrand into partial fractions,

$$\frac{1}{x^2 - 1} = \frac{1}{2} \frac{1}{x - 1} - \frac{1}{2} \frac{1}{x + 1}.$$

Thus,

$$\int \frac{1}{x^2 - 1} dx = \frac{1}{2} \ln|x - 1| - \frac{1}{2} \ln|x + 1| + C.$$

Note: The antiderivative “blows up” when the argument of the natural logarithm is zero. Hence, the domain of the antiderivative is equal to $\{x \in \mathbb{R} \mid x \notin \{-1, +1\}\}$, the same points where the original rational function blows up (technically, the points where it is not a well-defined function). ■

Example 7.31. Compute the partial fraction expansion of $R(x) = \frac{6x^2 - 6x - 6}{(x-1)(x-(-2))(x-3)}$ and then use it to evaluate the antiderivative.

Solution: The degree of the numerator is one less than the denominator, and we have $r_1 = 1$, $r_2 = -2$, and $r_3 = 3$. The roots are, therefore, real and distinct, so a partial fraction expansion exists.

$$\begin{aligned} \bullet k_1 &= \lim_{x \rightarrow r_1} \frac{(x-1)(6x^2 - 6x - 6)}{(x-1)(x-(-2))(x-3)} = \frac{(6(1)^2 - 6(1) - 6)}{(1-(-2))(1-3)} = \frac{-6}{-6} = 1 \\ \bullet k_2 &= \lim_{x \rightarrow r_2} \frac{(x-(-2))(6x^2 - 6x - 6)}{(x-1)(x-(-2))(x-3)} = \frac{(6(-2)^2 - 6(-2) - 6)}{((-2)-1)((-2)-3)} = \frac{30}{15} = 2 \\ \bullet k_3 &= \lim_{x \rightarrow r_3} \frac{(x-3)(6x^2 - 6x - 6)}{(x-1)(x-(-2))(x-3)} = \frac{(6(3)^2 - 6(3) - 6)}{(3-1)(3-(-2))} = \frac{30}{10} = 3 \end{aligned}$$

Hence, we have

$$\int \frac{6x^2 - 6x - 6}{(x-1)(x-(-2))(x-3)} dx = \ln|x-1| + 2\ln|x+2| + 3\ln|x-3|. \quad \blacksquare$$

Antiderivatives for a Pair of Complex Conjugate Roots

The denominator of the real rational function

$$R(x) = \frac{ax+b}{(x-\alpha)^2 + \omega^2}$$

has complex roots, $z = \alpha + i\omega$ and $z^* = \alpha - i\omega$, where z^* is the **complex conjugate** of z . While it is true that $R(x)$ can be expanded as

$$R(x) = \frac{\kappa}{x-z} + \frac{\kappa^*}{x-z^*},$$

trying to find the antiderivative in this form would take us into the realm of the Calculus of Complex Variables, which is beyond the scope of the course. Instead, we work directly with $R(x)$ as a ratio of a linear numerator and a quadratic denominator.

Proposition 7.32: Antiderivative for a Pair of Complex Conjugate Roots

Consider a real rational function $R(x) = \frac{ax+b}{(x-\alpha)^2 + \omega^2}$. Its antiderivative is,

$$\int \frac{ax+b}{(x-\alpha)^2 + \omega^2} dx = \frac{a}{2} \ln((x-\alpha)^2 + \omega^2) + (b-a\alpha) \frac{1}{\omega} \arctan\left(\frac{x-\alpha}{\omega}\right) + C, \quad (7.16)$$

where C is a constant of integration.

Proof: The antiderivative of the function $\frac{ax+b}{(x-\alpha)^2 + \omega^2}$ is found by combining u-substitution and the recognition of a total derivative. The integral is broken down into two parts based on the numerator,

$$\int \frac{ax+b}{(x-\alpha)^2 + \omega^2} dx = \int \frac{a(x-\alpha) + (a\alpha+b)}{(x-\alpha)^2 + \omega^2} dx,$$

where splitting the integral gives,

$$= a \int \frac{x-\alpha}{(x-\alpha)^2 + \omega^2} dx + (a\alpha+b) \int \frac{1}{(x-\alpha)^2 + \omega^2} dx.$$

Step 1: Using u-Substitution:

For the first part,

$$a \int \frac{x-\alpha}{(x-\alpha)^2 + \omega^2} dx,$$

we let $u = (x - \alpha)^2 + \omega^2$, then $du = 2(x - \alpha)dx$, and $\frac{1}{2}du = (x - \alpha)dx$. Thus, this part of the integral becomes,

$$\frac{a}{2} \int \frac{1}{u} du = \frac{a}{2} \ln|u| + C_1 = \frac{a}{2} \ln((x - \alpha)^2 + \omega^2) + C_1$$

after substituting back in $u = (x - \alpha)^2 + \omega^2$. Because $(x - \alpha)^2 + \omega^2 > 0$, the absolute value can be removed; indeed, if $\omega = 0$, then the root is not complex.

Step 2: Recognizing the Derivative of Arctangent

For the second part,

$$(a\alpha + b) \int \frac{1}{(x - \alpha)^2 + \omega^2} dx,$$

is a total differential of the arctangent function, namely,

$$\frac{a\alpha + b}{\omega} \arctan\left(\frac{x - \alpha}{\omega}\right) + C_2,$$

where once again, we used the fact that because the root is complex, ω cannot equal zero.

Final Step:

Combining both parts, we get the full antiderivative,

$$\int \frac{ax + b}{(x - \alpha)^2 + \omega^2} dx = \frac{a}{2} \ln((x - \alpha)^2 + \omega^2) + \frac{a\alpha + b}{\omega} \arctan\left(\frac{x - \alpha}{\omega}\right) + C$$

where $C := C_1 + C_2$ is the constant of integration. ■

Example 7.33. Find an antiderivative for $R(x) = \frac{2x+1}{3x^2-6x+15}$.

Solution: Ans. $\frac{1}{3} \cdot \log(x^2 - 2 \cdot x + 5) + \frac{1}{2} \cdot \arctan\left(\frac{x}{2} - \frac{1}{2}\right)$

We identify $P(x) = 2x + 1$ and $Q(x) = 3x^2 - 6x + 15$, which clearly satisfy the degree of the numerator being less than the degree of the denominator. To see if the denominator has real or complex poles, we check the discriminant

$$b^2 - 4ac = (-6)^2 - 4(3)(15) = -144 < 0,$$

and hence, the roots are complex. Indeed, applying the quadratic formula yields $z = 1 \pm 2i$, and thus

$$\alpha = 1 \quad \text{and} \quad \omega = 2;$$

it is fine to chose $\omega = -2$ as it will not change the answer.

The final step for applying Proposition 7.32 is to make the denominator monic by factoring out the coefficient of x^2 ,

$$R(x) = \frac{1}{3} \cdot \frac{2x + 1}{x^2 - 2x + 5} = \frac{1}{3} \cdot \frac{2x + 1}{(x - 1)^2 + 2^2}.$$

Substituting into (7.16) yields,

$$\int R(x) dx = \frac{1}{3} \cdot \left(\log(x^2 - 2 \cdot x + 5) + \frac{3 \cdot \arctan\left(\frac{x}{2} - \frac{1}{2}\right)}{2} \right).$$

Antiderivatives for a Mixture of Real and Complex Conjugate Roots

As you are well aware, real polynomials often have a mixture of real roots and complex roots, with the complex roots occurring in complex conjugate pairs. Pairs of complex conjugate roots can be recombined to give a quadratic polynomial with real coefficients, which can then be integrated as in the previous subsection. We'll first show how the corresponding partial fraction expansions work in software, and then we'll state, for the record, the corresponding theoretical result.

Example 7.34. Use software to find first the partial fraction expansion and then the antiderivative of

$$R(x) = \frac{2 \cdot x^3 - x + 4}{x^6 - 4 \cdot x^4 + 10 \cdot x^3 - 11 \cdot x^2 + 10 \cdot x - 6}.$$

Clearly, doing this by hand would be a nightmare. You already need an algorithm to find the roots of the denominator, so why not use another algorithm to find the coefficients in the PFE?

using SymPy

```
1 # Define the variable
2 x = symbols("x")
3
4 P = 2*x^3 - x + 4
5 Q = x^6 - 4*x^4 + 10*x^3 - 11*x^2 + 10*x - 6
6
7 R = P/Q
8
9 # Compute the PFE
10 PFE = apart(R)
11
12
13
14 println(PFE)
```

Output

```
-(4*x + 3)/(10*(x^2 + 1)) - (5*x - 8)/(17*(x^2 - 2*x + 2)) + 47/(680*(x + 3)) + 5/(8*(x - 1))
```

or, nicely typeset,

$$\text{PFE}(x) = \frac{-(4 \cdot x + 3)}{10 \cdot (x^2 + 1)} - \frac{5 \cdot x - 8}{17 \cdot (x^2 - 2 \cdot x + 2)} + \frac{47}{680 \cdot (x + 3)} + \frac{5}{8 \cdot (x - 1)}$$

The `apart` function in SymPy has nicely expanded the rational function as a sum of first-order terms and second-order terms. We could not ask for more, because we know how antiderivatives for each term! But instead of doing them by hand, we continue with software.

Continuing with the antiderivative:

```
1 antiderivative = integrate(PFE, x)
2 println(antiderivative)
```

Output

```
5*log(x - 1)/8 + 47*log(x + 3)/680 - log(x^2 + 1)/5 - 5*log(x^2 - 2*x + 2)/34
- 3*atan(x)/10 + 3*atan(x - 1)/17
```

or nicely typeset,

$$\int R(x) dx = \frac{5 \cdot \log(x - 1)}{8} + \frac{47 \cdot \log(x + 3)}{680} - \frac{\log(x^2 + 1)}{5} - \frac{5 \cdot \log(x^2 - 2 \cdot x + 2)}{34} - \frac{3 \cdot \arctan(x)}{10} + \frac{3 \cdot \arctan(x - 1)}{17}.$$

As you might guess, the antiderivative can also be computed in one step, in other words, without first doing a PFE:

```

1 # Find the antiderivative of the rational function
2 antiderivative = integrate(R, x)
3 println(antiderivative)
```

Output

```
5*log(x - 1)/8 + 47*log(x + 3)/680 - log(x^2 + 1)/5 - 5*log(x^2 - 2*x + 2)/34
- 3*atan(x)/10 + 3*atan(x - 1)/17
```

or nicely typeset,

$$\int R(x) dx = \frac{5 \cdot \log(x-1)}{8} + \frac{47 \cdot \log(x+3)}{680} - \frac{\log(x^2+1)}{5} - \frac{5 \cdot \log(x^2-2 \cdot x+2)}{34} - \frac{3 \cdot \arctan(x)}{10} + \frac{3 \cdot \arctan(x-1)}{17}.$$

■

Below is how the subject is taught (if at all) in a traditional Calculus course.

Proposition 7.35: (Optional Read:) Antiderivative for a Mixture of Real and Complex Roots

Consider a real rational function $R(x) = \frac{P(x)}{Q(x)}$ where the degree of the numerator is strictly lower than the degree of the denominator, $Q(x)$ is monic, and the roots of the denominator polynomial are distinct. We list the roots as

- **real roots:** r_1, r_2, \dots, r_m
- **complex roots:** z_1, z_2, \dots, z_p followed by their **complex conjugates**, $z_1^*, z_2^*, \dots, z_p^*$, where
- $\text{Real}\{z_i\} =: \alpha_i$ and $\text{Imag}\{z_i\} =: \omega_i$, so that $z_i = \alpha_i + i\omega_i$.

Under these assumptions,

$$Q(x) = \prod_{i=1}^p (x - z_i)(x - z_i^*) = \prod_{i=1}^p ((x - \alpha_i)^2 + \omega_i^2),$$

and the rational function $R(x)$ can be expressed as a partial fraction expansion (PFE)

$$R(x) = \sum_{i=1}^m \frac{k_i}{x - r_i} + \sum_{i=1}^p \frac{a_i x + b_i}{(x - \alpha_i)^2 + \omega_i^2},$$

where the real coefficients k_i, a_i, b_i satisfy

$$k_i = \lim_{x \rightarrow r_i} \frac{(x - r_i)P(x)}{Q(x)} = \frac{P(r_i)}{\prod_{\substack{j=1 \\ j \neq i}}^n (r_i - r_j)}, \quad (7.17)$$

for the real roots, and for the complex roots,

$$a_i = \text{Real} \left\{ \lim_{x \rightarrow z_i} \frac{(x - z_i)P(x)}{Q(x)} \right\} = \text{Real} \left\{ \frac{P(z_i)}{(z_i - z_i^*)} \cdot \frac{1}{\prod_{\substack{j=1 \\ j \neq i}}^p ((z_i - \alpha_j)^2 + \omega_j^2)} \right\}, \quad (7.18)$$

$$b_i = \text{Imag} \left\{ \lim_{x \rightarrow z_i} \frac{(x - z_i)P(x)}{Q(x)} \right\} = \text{Imag} \left\{ \frac{P(z_i)}{(z_i - z_i^*)} \cdot \frac{1}{\prod_{\substack{j=1 \\ j \neq i}}^p ((z_i - \alpha_j)^2 + \omega_j^2)} \right\}.$$

Antiderivative of a Generic Term: The antiderivative for a term of the form $\frac{k}{x-r}$ can be found as follows:

$$\int \frac{k}{x-r} dx = k \ln|x-r| + C, \quad (7.19)$$

where C is the integration constant.

Antiderivative of a Generic Term: The antiderivative for a term of the form $\frac{ax+b}{(x-\alpha)^2+\omega^2}$ can be found as follows:

$$\int \frac{ax+b}{(x-\alpha)^2+\omega^2} dx = \frac{a}{2} \ln((x-\alpha)^2 + \omega^2) + (b - a\alpha) \frac{1}{\omega} \arctan\left(\frac{x-\alpha}{\omega}\right) + C, \quad (7.20)$$

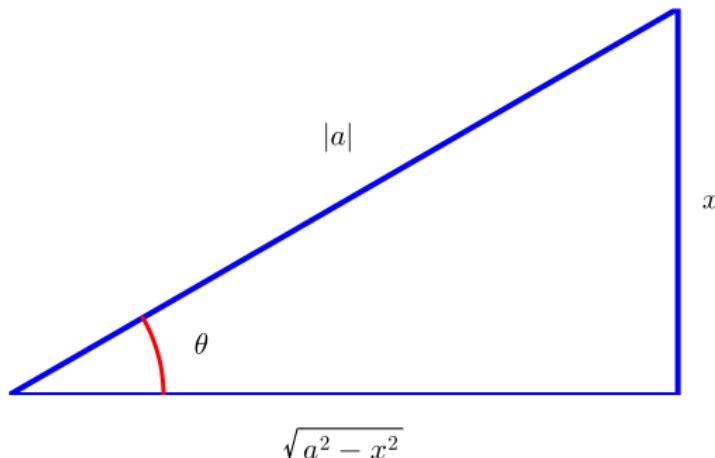
where C is the integration constant.

We do not offer practice videos because doing partial fraction expansions by hand is a fool's errand for anything but the simplest of cases. It's advisable to employ software tools for calculating PFEs for anything other than a second-order problem, and once you do that, why not go ahead and compute the antiderivative using software as well?

7.3.6 Trigonometric Substitutions for Radicals

Trigonometric substitution is a technique for computing antiderivatives (and integrals) that is particularly useful when the integrand involves expressions like $\sqrt{a^2 - x^2}$, $\sqrt{a^2 + x^2}$, or $\sqrt{x^2 - a^2}$. Here, we carefully work out examples where trigonometric substitution is commonly used, highlighting common pitfalls along the way.

Heads up: Trig substitution is the most technically and mentally taxing of all the methods for finding antiderivatives by hand. Hence, if you find it difficult, you are not alone. You might see it a few times in an engineering course, though it arises less frequently than u-substitution and integration by parts. We will include images that help you to understand how the transformations were “discovered”.



The above triangle satisfies the Pythagorean Theorem, and hence, it is a right triangle. From basic trigonometry it follows that the opposite side of the triangle is given by $x = |a| \sin(\theta)$, which motivates the change of variable we use in the next problem.

Example 7.36. Evaluate $\int \sqrt{a^2 - x^2} dx$.

Solution: Ans. For $|x| \leq |a|$, $\int \sqrt{a^2 - x^2} dx = \frac{a^2}{2} \cdot \arcsin\left(\frac{x}{|a|}\right) + \frac{x}{2} \cdot \sqrt{a^2 - x^2} + C$.

Because a is squared in the integrand, we can either replace it by $|a|$ or simply assume that $a > 0$. We'll make fewer errors if we use $|a|$, as in the diagram.

The integrand exists and is continuous for $(a^2 - x^2 \geq 0) \iff (-|a| \leq x \leq |a|)$. The standard trig substitution for this integral is $x = |a|\sin(\theta)$ and $dx = |a|\cos(\theta)d\theta$. We note that with $-\frac{\pi}{2} \leq \theta \leq \frac{\pi}{2}$, it follows that $-|a| \leq x \leq |a|$ and that $0 \leq \cos(\theta) \leq 1$. Moreover, later, when we need to inverse the transformation, we are working with a principal domain for defining $\arcsin(x)$. Hence, everything checks out!

Using this substitution and recalling that $1 - \sin^2(x) = \cos^2(x)$, the integral becomes

$$\begin{aligned} \int \left(\sqrt{a^2 - a^2 \sin^2(\theta)} \right) \cdot |a|\cos(\theta) d\theta &= \int \left(\sqrt{a^2 \cdot \cos^2(\theta)} \right) \cdot |a|\cos(\theta) d\theta \\ &= \int (|a| \cdot |\cos(\theta)|) \cdot |a|\cos(\theta) d\theta \\ &= a^2 \int \cos^2(\theta) d\theta \end{aligned}$$

where the last equality follows because we (carefully) noted earlier that $\cos(\theta) \geq 0$ for $-\frac{\pi}{2} \leq \theta \leq \frac{\pi}{2}$, the domain where the substitution was defined.

A visual derivation based on the triangle: From $x = |a|\sin(\theta)$, we deduce that $dx = |a|\cos(\theta)d\theta$. We also have the adjacent side of the triangle is $|a|\cos(\theta) = \sqrt{a^2 - x^2}$, which can be used to simplify the integrand. Substituting these into the integral, we have

$$\int \sqrt{a^2 - x^2} dx = \int \underbrace{|a|\cos(\theta)}_{\sqrt{a^2-x^2}} \cdot \underbrace{|a|\cos(\theta)d\theta}_{dx} = a^2 \int \cos^2(\theta) d\theta.$$

Using the half-angle identity, $\cos^2(\theta) = \frac{1+\cos(2\theta)}{2}$, we have

$$a^2 \int \cos^2(\theta) d\theta = \frac{a^2}{2} \int (1 + \cos(2\theta)) d\theta = \frac{a^2}{2} \left(\theta + \frac{\sin(2\theta)}{2} \right) + C.$$

To return to x , we use $\sin(\theta) = \frac{x}{|a|}$ and $\cos(\theta) = \sqrt{1 - \sin^2(\theta)} = \sqrt{1 - \left(\frac{x}{|a|}\right)^2}$. Importantly, as noted earlier, we are working with a principal domain of the function \arcsin . Thus

$$\theta = \arcsin\left(\frac{x}{|a|}\right), \text{ and } \sin(2\theta) = 2\sin(\theta)\cos(\theta) = \frac{2x}{|a|} \sqrt{1 - \left(\frac{x}{|a|}\right)^2}.$$

The final antiderivative in terms of x is

$$\int \sqrt{a^2 - x^2} dx = \frac{a^2}{2} \left(\arcsin\left(\frac{x}{|a|}\right) + \frac{x}{|a|} \sqrt{1 - \left(\frac{x}{|a|}\right)^2} \right) + C.$$

Because $-1 \leq \frac{x}{|a|} \leq 1$ for all x in the domain of the integrand, the terms $\arcsin\left(\frac{x}{|a|}\right)$ and $\sqrt{1 - \left(\frac{x}{|a|}\right)^2}$ in the antiderivative are well defined (i.e., make mathematical sense²); hence, the computed antiderivative is defined for all x in the domain of the integrand.

²As an example, if $\frac{x}{|a|}$ were greater than one, then neither the arcsine nor the square root would make sense.

A more aesthetically pleasing form of the answer can be obtained with a small bit of algebra,

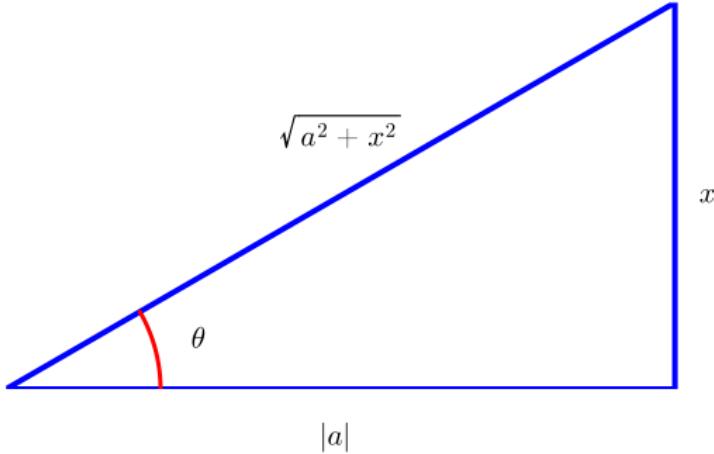
$$\boxed{\text{For } |x| \leq |a|, \int \sqrt{a^2 - x^2} dx = \frac{a^2}{2} \cdot \arcsin\left(\frac{x}{|a|}\right) + \frac{x}{2} \cdot \sqrt{a^2 - x^2} + C.}$$

Note: From the right triangle associated with this problem, we have $\theta = \arcsin\left(\frac{x}{|a|}\right)$ and $\theta = \tan^{-1}\left(\frac{x}{\sqrt{a^2 - x^2}}\right)$, and thus

$$\arcsin\left(\frac{x}{|a|}\right) = \tan^{-1}\left(\frac{x}{\sqrt{a^2 - x^2}}\right).$$

It follows that the answer can also be written as

$$\boxed{\text{For } |x| \leq |a|, \int \sqrt{a^2 - x^2} dx = \frac{a^2}{2} \cdot \tan^{-1}\left(\frac{x}{\sqrt{a^2 - x^2}}\right) + \frac{x}{2} \cdot \sqrt{a^2 - x^2} + C.}$$



Note that the above triangle satisfies the Pythagorean Theorem, and hence it is a right triangle. It follows that $\tan(\theta) = \frac{x}{|a|}$, or, $x = |a| \tan(\theta)$, which motivates the change of variable we use in the next problem.

Example 7.37. Evaluate the integral

$$\int \sqrt{a^2 + x^2} dx.$$

Solution: Ans. For $x \in \mathbb{R}$, $\int \sqrt{a^2 + x^2} dx = \frac{a^2}{2} \ln\left(x + \sqrt{a^2 + x^2}\right) + \frac{x}{2} \cdot \sqrt{a^2 + x^2} + C$.

The integrand $\sqrt{a^2 + x^2} dx$ is defined and continuous in $(-\infty, \infty)$. Because a is squared, we replace it by $|a|$ to avoid sign errors later on.

The substitution $x = |a| \tan(\theta)$ is defined for $-\frac{\pi}{2} < \theta < \frac{\pi}{2}$, and gives $dx = |a| \cdot \sec^2(\theta) d\theta$. The integral becomes

$$\int \sqrt{a^2 + a^2 \tan^2(\theta)} \cdot |a| \cdot \sec^2(\theta) d\theta = \int \sqrt{a^2 (1 + \tan^2(\theta))} \cdot |a| \cdot \sec^2(\theta) d\theta.$$

Next, we use the trig identity, $1 + \tan^2(\theta) = \sec^2(\theta)$, giving us

$$\begin{aligned} \int \sqrt{a^2(1 + \tan^2(\theta))} \cdot |a| \cdot \sec^2(\theta) d\theta &= \int \sqrt{a^2 \cdot \sec^2(\theta)} \cdot |a| \cdot \sec^2(\theta) d\theta \\ &= \int |a| |\sec(\theta)| \cdot |a| \cdot \sec^2(\theta) d\theta \\ &= a^2 \int \sec^3(\theta) d\theta \end{aligned}$$

where in the last line we used the fact that $\sec(\theta) := \frac{1}{\cos(\theta)}$ is positive for $-\frac{\pi}{2} < \theta < \frac{\pi}{2}$, the domain of the change of variable.

Computing the antiderivative of $\int \sec^3(\theta) d\theta$ requires a new technique that is useful to know when performing integration by parts. Define,

$$\begin{aligned} I &:= \int \sec^3(\theta) d\theta \\ &= \int \sec(\theta) \cdot (1 + \tan^2(x)) d\theta \quad (\text{trig identity, } 1 + \tan^2(\theta) = \sec^2(\theta)) \\ &= \int \sec(\theta) d\theta + \int \sec(\theta) \cdot \tan^2(x) d\theta \quad (\text{algebra}) \\ &= \ln |\sec(\theta) + \tan(\theta)| + \int \sec(\theta) \cdot \tan^2(x) d\theta \quad (\text{first integral is from Example 7.22}). \end{aligned} \tag{7.21}$$

The remaining integral in (7.21) requires integration by parts. Define

$$\begin{aligned} u &= \tan(\theta) \implies du = \sec^2(\theta) d\theta \\ v &= \sec(\theta) \implies dv = \sec(\theta) \cdot \tan(\theta) d\theta, \end{aligned}$$

yielding,

$$\begin{aligned} \int \underbrace{\sec(\theta) \tan^2(\theta) d\theta}_{u dv} &= \underbrace{\tan(\theta) \sec(\theta)}_{u \cdot v} - \int \underbrace{\sec(\theta) \sec^2(\theta) d\theta}_{v du} \\ &= \tan(\theta) \sec(\theta) - \int \sec^3(\theta) d\theta \\ &= \tan(\theta) \sec(\theta) - I. \end{aligned} \tag{7.22}$$

Substituting back into (7.21), we have

$$\begin{aligned} I &:= \int \sec^3(\theta) d\theta = \ln |\sec(\theta) + \tan(\theta)| + \int \sec(\theta) \cdot \tan^2(x) d\theta \\ &= \ln |\sec(\theta) + \tan(\theta)| + \tan(\theta) \sec(\theta) - I. \end{aligned} \tag{7.23}$$

Solving for I yields

$$\begin{aligned} 2I &= \ln |\sec(\theta) + \tan(\theta)| + \tan(\theta) \sec(\theta) \\ &\Downarrow \\ I &= \frac{1}{2} \ln |\sec(\theta) + \tan(\theta)| + \frac{1}{2} \tan(\theta) \sec(\theta). \end{aligned} \tag{7.24}$$

Putting the (numerous) steps together yields

$$\int \sqrt{a^2 + x^2} dx = a^2 \int \sec^3(\theta) d\theta = a^2 I = \frac{a^2}{2} \ln |\sec(\theta) + \tan(\theta)| + \frac{a^2}{2} \tan(\theta) \sec(\theta).$$

To return to an answer in x , we use the trigonometric substitutions,

$$\tan(\theta) = \frac{x}{|a|},$$

$$\sec(\theta) = \sqrt{1 + \tan^2(\theta)} = \sqrt{1 + \left(\frac{x}{|a|}\right)^2}.$$

The final antiderivative in terms of x is

$$\begin{aligned} \int \sqrt{a^2 + x^2} dx &= \frac{a^2}{2} \ln \left| \sqrt{1 + \left(\frac{x}{|a|}\right)^2} + \frac{x}{|a|} \right| + \frac{a^2}{2} \frac{x}{|a|} \cdot \sqrt{1 + \left(\frac{x}{|a|}\right)^2} + C \\ &= \frac{a^2}{2} \ln \left(\sqrt{1 + \left(\frac{x}{|a|}\right)^2} + \frac{x}{|a|} \right) + \frac{x}{2} \cdot \sqrt{a^2 + x^2} + C \end{aligned}$$

where C is the constant of integration. The absolute value signs in the natural logarithm were removed because

$$\sqrt{1 + \left(\frac{x}{|a|}\right)^2} + \frac{x}{|a|} > 0$$

for all $x \in \mathbb{R}$.

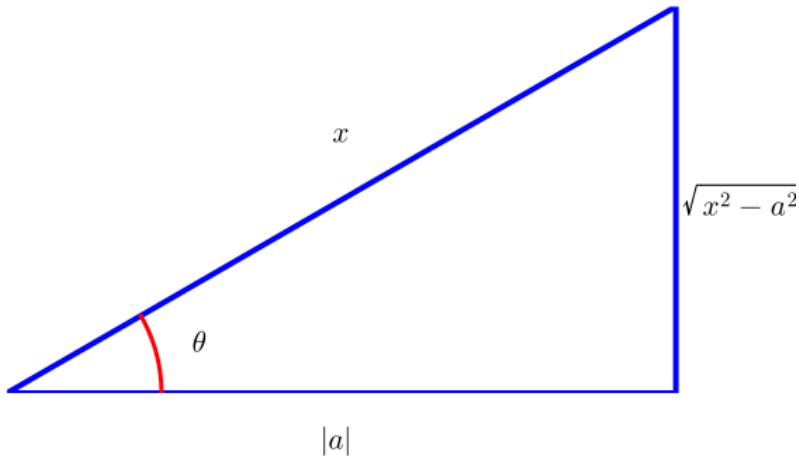
It is traditional to simplify further the answer by taking advantage of the arbitrary constant of integration, C .

$$\begin{aligned} \ln \left(\sqrt{1 + \left(\frac{x}{|a|}\right)^2} + \frac{x}{|a|} \right) &= \ln \left(\frac{1}{|a|} \cdot \left(\sqrt{a^2 + x^2} + x \right) \right) \\ &= \ln \left(\frac{1}{|a|} \right) + \ln \left(\sqrt{a^2 + x^2} + x \right) \end{aligned}$$

Absorbing the constant $\ln \left(\frac{1}{|a|} \right)$ into C gives a much nicer final answer,

$$\text{For } x \in \mathbb{R}, \int \sqrt{a^2 + x^2} dx = \frac{a^2}{2} \ln \left(x + \sqrt{a^2 + x^2} \right) + \frac{x}{2} \cdot \sqrt{a^2 + x^2} + C.$$

Note: The answer looks so innocent. Why was it so hard to compute?



Note that the above triangle satisfies the Pythagorean Theorem, and hence it is a right triangle. It follows that $\cos(\theta) = \frac{|a|}{x}$, or, $x = \frac{|a|}{\cos(\theta)} = |a|\sec(\theta)$, which motivates a common change of variables for the next problem. There is also another viable change of variable that we wish to illustrate. **Hyperbolic Trig Functions** are often used in integrals that have terms of the form $x^2 - a^2$. This is because of identities such as, $\cosh^2(x) - \sinh^2(x) = 1$, which can also be written as $\cosh^2(x) - 1 = \sinh^2(x)$. This next example follows a similar pattern to the previous examples, ensuring that the domain of the integrand is carefully delineated and each change of variables respects the domain in which the integrand is defined.

Example 7.38. Find an antiderivative for $\sqrt{x^2 - a^2}$.

Solution: Ans. For $|x| \geq |a|$, $\int \sqrt{x^2 - a^2} dx = \frac{x}{2}\sqrt{x^2 - a^2} - \text{sign}(x) \cdot \frac{a^2}{2} \ln(|x| + \sqrt{x^2 - a^2}) + C$, where

$$\text{sign}(x) = \begin{cases} +1 & x \geq 0 \\ -1 & x < 0. \end{cases}$$

First Approach: Standard Trig Substitution: The integrand of $\int \sqrt{x^2 - a^2} dx$ is defined and continuous wherever $x^2 - a^2 \geq 0$, that is $|x| \geq |a|$. We can therefore develop an antiderivative for $x \geq |a|$, and a second one for $x \leq -|a|$. Here, we will develop a solution for $x \leq -|a|$ using standard trig functions, and with the hyperbolic trig functions, we'll develop the solution for $x \geq |a|$.

For $x \leq -|a|$, we choose the change of variable $x = -|a|\sec(\theta)$, which gives $dx = -|a|\sec(\theta) \cdot \tan(\theta) d\theta$. For all $0 \leq \theta < \frac{\pi}{2}$, $1 \leq \sec(\theta) < \infty$, and hence the change of variables respects the domain of definition of the integrand. Moreover, for later use, we note that for all $0 \leq \theta < \frac{\pi}{2}$, $\tan(\theta) \geq 0$.

Using the change of variable, the integral becomes

$$\begin{aligned} \int \sqrt{x^2 - a^2} dx &= \int \sqrt{a^2 \sec^2(\theta) - a^2} \cdot (-1) \cdot |a|\sec(\theta) \cdot \tan(\theta) d\theta \\ &= -a^2 \int \sqrt{\tan^2(\theta)} \cdot \sec(\theta) \cdot \tan(\theta) d\theta \quad (\sec^2(\theta) - 1 = \tan^2(\theta)) \\ &= -a^2 \int |\tan(\theta)| \cdot \sec(\theta) \cdot \tan(\theta) d\theta \quad (\tan(\theta) \geq 0) \\ &= -a^2 \int \tan^2(\theta) \cdot \sec(\theta) d\theta \quad (\tan^2(\theta) = \sec^2(\theta) - 1) \\ &= -a^2 \int (\sec^2(\theta) - 1) \cdot \sec(\theta) d\theta \\ &= -a^2 \int \sec^3(\theta) d\theta + a^2 \int \sec(\theta) d\theta. \end{aligned}$$

Fortunately for us, we have already computed antiderivatives for both of these integrands,

$$\begin{aligned}\int \sec(\theta) d\theta &= \ln |\sec(\theta) + \tan(\theta)| + C \\ \int \sec^3(\theta) d\theta &= \frac{1}{2} \ln |\sec(\theta) + \tan(\theta)| + \frac{1}{2} \tan(\theta) \sec(\theta) + C.\end{aligned}$$

Hence,

$$\begin{aligned}\int \sqrt{x^2 - a^2} dx &= -\frac{a^2}{2} \ln |\sec(\theta) + \tan(\theta)| - \frac{a^2}{2} \tan(\theta) \sec(\theta) + a^2 \ln |\sec(\theta) + \tan(\theta)| + C \\ &= \frac{a^2}{2} \ln |\sec(\theta) + \tan(\theta)| - \frac{a^2}{2} \tan(\theta) \sec(\theta) + C\end{aligned}$$

From $x = -|a| \sec(\theta)$, we have that $\sec(\theta) = -\frac{x}{|a|}$. From the trig diagram for the problem, we have

$$\tan(\theta) = \frac{\sqrt{x^2 - a^2}}{|a|}.$$

Substituting in, we have

$$\begin{aligned}\int \sqrt{x^2 - a^2} dx &= \frac{a^2}{2} \ln \left| -\frac{x}{|a|} + \frac{\sqrt{x^2 - a^2}}{|a|} \right| - \frac{a^2}{2} \frac{\sqrt{x^2 - a^2}}{|a|} \cdot \left(\frac{-x}{|a|} \right) + C \\ &= \frac{a^2}{2} \ln |-x + \sqrt{x^2 - a^2}| - \ln(|a|) + \frac{x}{2} \sqrt{x^2 - a^2} + C \\ &= \frac{a^2}{2} \ln |-x + \sqrt{x^2 - a^2}| + \frac{x}{2} \sqrt{x^2 - a^2} + C,\end{aligned}$$

where in the last line, we absorbed $\ln(|a|)$ into the constant of integration, C .

Second Approach: Hyperbolic Trig Functions: The integrand of $\int \sqrt{x^2 - a^2} dx$ is defined and continuous wherever $x^2 - a^2 \geq 0$, that is $|x| \geq |a|$. We can therefore develop an antiderivative for $x \geq |a|$, and a second one for $x \leq -|a|$.

For $x \geq |a|$, we choose the hyperbolic substitution $x = |a| \cosh(t)$, which gives $dx = |a| \sinh(t) dt$. For all $t \geq 0$, $\cosh(t) \geq 1$, and hence the change of variables respects the domain of definition of the integrand. Moreover, for later use, we note that for all $t \geq 0$, $\sinh(t) \geq 0$.

Using the change of variable, the integral becomes

$$\begin{aligned}\int \sqrt{x^2 - a^2} dx &= \int \sqrt{a^2 \cosh^2(t) - a^2} \cdot |a| \sinh(t) dt \\ &= a^2 \int \sqrt{\cosh^2(t) - 1} \cdot \sinh(t) dt \\ &= a^2 \int \sqrt{\sinh^2(t)} \cdot \sinh(t) dt \\ &= a^2 \int |\sinh(t)| \cdot \sinh(t) dt \\ &= a^2 \int \sinh^2(t) dt,\end{aligned}$$

where we used $\sinh(t) \geq 0$.

Using the identity $\sinh^2(t) = \frac{\cosh(2t)-1}{2}$, the integral simplifies to

$$\begin{aligned} a^2 \int \sinh^2(t) dt &= \frac{a^2}{2} \int (\cosh(2t) - 1) dt \\ &= \frac{a^2}{2} \left(\frac{\sinh(2t)}{2} - t \right) + C \\ &= \frac{a^2}{4} \sinh(2t) - \frac{a^2}{2} t + C. \end{aligned}$$

To return to x , we use $\cosh(t) = \frac{x}{|a|}$ and $\sinh(t) = \sqrt{\cosh^2(t) - 1}$. Thus

$$\begin{aligned} \sinh(2t) &= 2 \cosh(t) \sinh(t) = 2 \frac{x}{|a|} \cdot \sqrt{\left(\frac{x}{|a|}\right)^2 - 1} \\ t &= \operatorname{acosh}\left(\frac{x}{|a|}\right). \end{aligned}$$

The final antiderivative in terms of x is

$$\frac{a^2}{4} \cdot 2 \frac{x}{|a|} \sqrt{\left(\frac{x}{|a|}\right)^2 - 1} - \frac{a^2}{2} \operatorname{acosh}\left(\frac{x}{|a|}\right) + C,$$

which simplifies to

$$\frac{x}{2} \sqrt{x^2 - a^2} - \frac{a^2}{2} \operatorname{acosh}\left(\frac{x}{|a|}\right) + C.$$

From Table 1.77 on the inverse hyperbolic trig functions,

$$\begin{aligned} \operatorname{acosh}\left(\frac{x}{|a|}\right) &= \ln\left(\frac{x}{|a|} + \sqrt{\frac{x^2}{a^2} - 1}\right) \\ &= \ln\left(\frac{1}{|a|} \left(x + \sqrt{x^2 - a^2}\right)\right) \\ &= \ln\left(\frac{1}{|a|}\right) + \ln\left(x + \sqrt{x^2 - a^2}\right). \end{aligned}$$

Absorbing the term $\ln\left(\frac{1}{|a|}\right)$ into the constant of integration gives, for $x \geq |a|$,

$$\int \sqrt{x^2 - a^2} dx = \frac{x}{2} \sqrt{x^2 - a^2} - \frac{a^2}{2} \ln\left(x + \sqrt{x^2 - a^2}\right) + C.$$

The answers for $x \geq |a|$ and $x \leq -|a|$ can be combined into a single formula,

$$\text{For } |x| \geq |a|, \int \sqrt{x^2 - a^2} dx = \frac{x}{2} \sqrt{x^2 - a^2} - \operatorname{sign}(x) \cdot \frac{a^2}{2} \ln(|x| + \sqrt{x^2 - a^2}) + C.$$

Recommended Videos for Additional Practice

- [Trig Substitution... How?](#) by Nancy Pi.
- [Trig substitution - How to solve?](#) by Krista King
- [How to solve EVERY trigonometric substitution problem ever!](#) by Krista King

- [Trigonometric Substitution](#) by the Organic Chemistry Tutor.
- [Calculus 2: Integration - Trig Substitution \(2 of 28\) Applicable Integrals](#) by Prof. Michael van Biezen. It's highly recommended that you take a peek to see an impressive list of integrals that can be solved with the methods we just developed.
- [Playlist for Trig Substitutions](#) by **BlackPenRedPen**. You'll have to find your favorite one.
- [Introduction to trigonometric substitution](#) by Khan Academy.

We repeat: **Trig substitutions are the “hardest to master” of all the methods for finding antiderivatives; they require great care if you want to avoid nonsense, even when using “tables”.** Moreover, when using ChatGPT+Wolfram or Wolfram Alpha, it is still on you to understand the domain on which the antiderivative is defined. You cannot just throw in a function, get an answer, snap your fingers, and say you are done! **To understand the domain of the antiderivative, analyze that of the integrand (aka, the function you are integrating).**

7.4 Coming Up for Air and Taking Stock

Now is a good time to re-read Chapter 7.2.3, **Conflating Integration and Antiderivatives is a Pedagogical Pitfall in Calculus**. It's rare in engineering practice that a function to be integrated is simple enough for you to find an antiderivative by hand, even after spending another year learning methods more advanced than those we have presented. And in the rare instances when you can find an antiderivative by hand, ChatGPT+Wolfram can almost surely do it too, and not make a sign error in the process.

In regards to integration, we recommend that you:

1. Appreciate that integration and differentiation are inverse operations (up to a constant of integration). The Fundamental Theorems of Calculus are amazing.
2. Use numerical integration tools whenever possible. At the very least, check your analytically derived integrals using computer-based methods by either symbolically differentiating the antiderivative to make sure it gives the original function, or by using numerical integration to check the integral.
3. Use ChatGPT with the Wolfram Plugin when you absolutely must determine the antiderivative of a function that is neither in a standard integral table nor among the six or eight that are worth memorizing.
4. Read [Secrets of the Arcane 7.39](#) to learn why antiderivatives limit the ability to compute definite integrals.

Secrets of the Arcane 7.39: Elementary Functions and Antiderivatives

In mathematics, when we talk about **elementary functions**, we are referring to functions that are constructed using basic arithmetic operations (addition, subtraction, multiplication, and division) and compositions of,

1. **Polynomials:** Functions like $f(x) = x^2$, $g(x) = 3x^3 - 2x + 1$.
2. **Rational Functions:** Ratios of two polynomials, e.g., $R(x) = \frac{x^2+1}{x-2}$.
3. **Root Functions:** Functions involving square roots, cube roots, etc., like $h(x) = \sqrt{x}$.
4. **Trigonometric Functions:** The standard sine, cosine, tangent, etc., and their inverses.
5. **Exponential and Logarithmic Functions:** Functions like e^x , $\ln(x)$.

These functions are termed “elementary” because they are among the first functions introduced in mathematics education and are foundational in calculus and analysis. They are well-understood and have known properties, such as **the derivative of an elementary function is also elementary**. When you combine elementary functions using operations like addition, multiplication, or composition, you get a wide variety of functions that can model many different phenomena, but they still retain the “elementary” label because they’re built from these basic components.

Most elementary functions do not have antiderivatives that are also elementary functions.^a Even the very common function, e^{-x^2} , used in Radial Basis Functions and Gaussian Probability Distributions, does not have an antiderivative that

can be expressed in terms of anything we would recognize. The same is true for the antiderivative of $\cos(x^2)$ in Example 7.4-(d). In these cases, the antiderivative is literally defined through the First Fundamental Theorem of Calculus, that is,

$$\text{the antiderivative of } f(x) \text{ is } F(x) = \int_a^x f(t) dt + C,$$

which is the integral you were trying to compute in the first place. Fortunately, for functions that arise sufficiently often, numerical analysts seek fast and accurate means to compute the new antiderivative. If you are lucky, your antiderivative may be one of the non-elementary antiderivatives that have been tabulated. Are you ready to roll the dice?

The coup de grâce in all of this: the hand methods for determining antiderivatives allow you to find antiderivatives that are already available in what used to be called **Integral Tables**. You can simply look up your antiderivative. Back in the day, engineering students purchased books with thousands of antiderivatives in them. Today, no one needs to do that. We have ChatGPT+Wolfram, which provides access to more integrals than all of the books we students of yore used to lug around.

Further Reading: [Non-elementary Integrals](#); [Proving the Non-existence of Elementary Anti-derivatives for Certain Elementary Functions](#); [Impossibility Theorems for Elementary Integration](#); [How WolframAlpha defines these nonelementary integrals](#); and [Closure of the Set of Elementary Functions](#).

^aIn plain words, seeking an antiderivative that you can express in understandable terms can be a fool's errand.

7.5 Software Tools for Finding Antiderivatives

Because some of you may arrive here without having read the previous sections, we more or less start from scratch our summary of software tools for doing “integration”.

7.5.1 SymPy in Julia

In the opinion of your author, SymPy beats Symbolics when it comes to finding antiderivatives, also known as “integration”. We illustrate several use cases here.

u-Substitution: Example 7.15

```

1 using SymPy
2
3 # Define symbolic variable
4 x = symbols("x")
5
6 # Define your integrand here
7 integrand = 3cos(x^2+1)
8 antiderivative = integrate(integrand, x)
9
10 # Print the problem and results
11 println("Integrand = ", integrand, "\n")
12 println("Antiderivative = ", antiderivative, "\n")

```

Output

Integrand = 3*cos(x^2 + 1)

Antiderivative = 3*sqrt(2)*sqrt(pi)*(cos(1)*fresnelc(sqrt(2)*x/sqrt(pi)) - sin(1)*fresnels(sqrt(2)*x/sqrt(pi)))

Integration by Parts: the “Tricky” Secant Integrand from Example 7.22

```
1 using SymPy
2 # Define symbolic variable
3 x = symbols("x")
4
5 # Define your integrand here
6 integrand = 1/cos(x)
7 antiderivative = integrate(integrand, x)
8
9 # Print the problem and results
10 println("Integrand = ", integrand, "\n")
11 println("Antiderivative = ", antiderivative, "\n")
```

Output

```
Integrand = 1/cos(x)

Antiderivative = -log(sin(x) - 1)/2 + log(sin(x) + 1)/2
```

Integration by Parts Example 7.25

```
1 using SymPy
2 # Define symbolic variable
3 x = symbols("x")
4
5 # Define your integrand here
6 integrand = x*sin(x)
7 antiderivative = integrate(integrand, x)
8
9 # Print the problem and results
10 println("Integrand = ", integrand, "\n")
11 println("Antiderivative = ", antiderivative, "\n")
```

Output

```
Integrand = x*sin(x)

Antiderivative = -x*cos(x) + sin(x)
```

Rational Functions and PFEs Example 7.34

```
1 using SymPy
2
3 # Define the variable
4 x = symbols("x")
5
6 P = 2*x^3 - x + 4
7 Q = x^6 - 4*x^4 + 10*x^3 - 11*x^2 + 10*x - 6
8
9 R = P/Q
10
```

```

11 # Define your integrand here
12 integrand = R
13 antiderivative = integrate(integrand, x)
14
15 # Print the problem and results
16 println("Integrand = ", integrand, "\n")
17 println("Antiderivative = ", antiderivative, "\n")

```

Output

Integrand = $(2*x^3 - x + 4)/(x^6 - 4*x^4 + 10*x^3 - 11*x^2 + 10*x - 6)$

Antiderivative = $5*\log(x - 1)/8 + 47*\log(x + 3)/680 - \log(x^2 + 1)/5 - 5*\log(x^2 - 2*x + 2)/34$

Trig Substitutions for Radicals Example 7.36

```

1 using SymPy
2
3 # Define the variable
4 x= symbols("x", real=true)
5
6 # Define your integrand here
7 a=pi
8 integrand =sqrt(a^2 - x^2)
9 antiderivative = integrate(integrand, x)
10
11 # Print the problem and results
12 println("Integrand = ", integrand, "\n")
13 println("Antiderivative = ", antiderivative, "\n")

```

Output

Integrand = $3.14159265358979*\sqrt{1 - 0.101321183642338*x^2}$

Antiderivative = $3.14159265358979*\text{Piecewise}((0.5*I*x*\sqrt{0.101321183642338*x^2 - 1}) - 1.570796$

Or nicely semi-nicely typeset,

$$3.14159 \cdot \begin{cases} 0.5 \cdot i \cdot x \cdot \sqrt{0.101321 \cdot x^2 - 1} - \frac{1.570796 \cdot i \cdot x \cdot \operatorname{arccosh}(0.318310 \cdot \operatorname{Abs}(x))}{\operatorname{Abs}(x)} + \frac{0.785398 \cdot \pi \cdot x}{\operatorname{Abs}(x)} & x^2 > 9.869604401089337 \\ \frac{-0.050661 \cdot x^3}{\sqrt{1 - 0.101321 \cdot x^2}} + \frac{1.570796 \cdot x \cdot \operatorname{arcsin}(0.318310 \cdot \operatorname{Abs}(x))}{\operatorname{Abs}(x)} + \frac{0.5 \cdot x}{\sqrt{1.000000 - 0.101321 \cdot x^2}} & \text{otherwise} \end{cases}$$

SymPy has given us the solution for $x^2 > \pi^2$, which makes the integrand complex-valued, as well as the answer we computed for a real integrand where $\pi^2 - x^2 \geq 0$. We'll let you dig into the intricacies of SymPy to avoid computing the complex answer in addition to the real answer you likely would want, because, as you will see in the next subsection, Wolfram Alpha kills it with supreme elegance.

7.5.2 More on Wolfram Alpha

Additional features can be accessed through the [Wolfram Alpha Website](#), such as step-by-step solutions. At the present time, the step-by-step solution is a paid feature, with student pricing available. Here is an example of the level of detail provided by Wolfram Alpha's Step-by-Step Solutions feature.



integrate sqrt(a^2 - x^2) dx steps

NATURAL LANGUAGE

MATH INPUT

EXTENDED KEYBOARD

EXAMPLES

UPLOAD

RANDOM

Indefinite integral

$$\int \sqrt{a^2 - x^2} dx = \frac{1}{2} \left(x \sqrt{a^2 - x^2} + a^2 \tan^{-1} \left(\frac{x}{\sqrt{a^2 - x^2}} \right) \right) + \text{constant}$$

$\tan^{-1}(x)$ is the inverse tangent function

Step-by-step solution

Figure 7.1: The Wolfram Alpha website provides more information than the Wolfram Plugin for ChatGPT. See the circled box: Step-by-step solution.

Indefinite Integral of $\sqrt{a^2 - x^2}$

STEP 1

Take the integral:

$$\int \sqrt{a^2 - x^2} dx$$

STEP 2

For the integrand $\sqrt{a^2 - x^2}$, (assuming all variables are positive) substitute $x = a \sin(u)$ and $dx = a \cos(u) du$. Then $\sqrt{a^2 - x^2} = \sqrt{a^2 - a^2 \sin^2(u)} = a \cos(u)$ and $u = \sin^{-1}(\frac{x}{a})$:

$$= a \int a \cos^2(u) du$$

STEP 3

Factor out constants:

$$= a^2 \int \cos^2(u) du$$

STEP 4

Write $\cos^2(u)$ as $\frac{1}{2} \cos(2u) + \frac{1}{2}$:

$$= a^2 \int \left(\frac{1}{2} \cos(2u) + \frac{1}{2} \right) du$$

STEP 5

Integrate the sum term by term and factor out constants:

$$= \frac{a^2}{2} \int \cos(2u) du + \frac{a^2}{2} \int 1 du$$

STEP 6

For the integrand $\cos(2u)$, substitute $s = 2u$ and $ds = 2 du$:

$$= \frac{a^2}{4} \int \cos(s) ds + \frac{a^2}{2} \int 1 du$$

STEP 7

The integral of $\cos(s)$ is $\sin(s)$:

$$= \frac{a^2}{4} \sin(s) + \frac{a^2}{2} \int 1 du$$

STEP 8

The integral of 1 is u :

$$= \frac{a^2}{4} \sin(s) + \frac{a^2 u}{2} + \text{constant}$$

STEP 9

Substitute back for $s = 2u$:

$$= \frac{a^2 u}{2} + \frac{a^2}{4} \sin(2u) + \text{constant}$$

STEP 10

Apply the double angle formula $\sin(2u) = 2\sin(u)\cos(u)$:

$$= \frac{a^2 u}{2} + \frac{a^2}{2} \sin(u) \cos(u) + \text{constant}$$

STEP 11

Express $\cos(u)$ in terms of $\sin(u)$ using $\cos^2(u) = 1 - \sin^2(u)$:

$$= \frac{a^2 u}{2} + \frac{a^2}{2} \sin(u) \sqrt{1 - \sin^2(u)} + \text{constant}$$

STEP 12

Substitute back for $u = \sin^{-1}\left(\frac{x}{a}\right)$:

$$= \frac{1}{2}x\sqrt{a^2 - x^2} + \frac{a^2}{2} \sin^{-1}\left(\frac{x}{a}\right) + \text{constant}$$

STEP 13

Factor the answer a different way:

$$= \frac{1}{2}\left(x\sqrt{a^2 - x^2} + a^2 \sin^{-1}\left(\frac{x}{a}\right)\right) + \text{constant}$$

STEP 14

Which is equivalent for restricted x and a values to:

$$\text{Answer: } \frac{1}{2}\left(x\sqrt{a^2 - x^2} + a^2 \tan^{-1}\left(\frac{x}{\sqrt{a^2 - x^2}}\right)\right) + \text{constant}$$

Indefinite Integral of $x^2 \sin^3(x)$ **STEP 1**

Take the integral:

$$\int x^2 \sin^3(x) dx$$

STEP 2

For the integrand $x^2 \sin^3(x)$, use the trigonometric identity $\sin^2(x) = \frac{1}{2}(1 - \cos(2x))$:

$$= \frac{1}{2} \int x^2 \sin(x)(1 - \cos(2x)) dx$$

STEP 3

Expanding the integrand $x^2 \sin(x)(1 - \cos(2x))$ gives $x^2 \sin(x) - x^2 \sin(x) \cos(2x)$:

$$= \frac{1}{2} \int (x^2 \sin(x) - x^2 \sin(x) \cos(2x)) dx$$

STEP 4

Integrate the sum term by term and factor out constants:

$$= -\frac{1}{2} \int x^2 \sin(x) \cos(2x) dx + \frac{1}{2} \int x^2 \sin(x) dx$$

STEP 5

Use the trigonometric identity $\sin(\alpha)\cos(\beta) = \frac{1}{2}(\sin(\alpha - \beta) + \sin(\alpha + \beta))$, where $\alpha = x$ and $\beta = 2x$:

$$= -\frac{1}{4} \int x^2(\sin(3x) - \sin(x)) dx + \frac{1}{2} \int x^2 \sin(x) dx$$

STEP 6

Expanding the integrand $x^2(\sin(3x) - \sin(x))$ gives $x^2 \sin(3x) - x^2 \sin(x)$:

$$= -\frac{1}{4} \int (x^2 \sin(3x) - x^2 \sin(x)) dx + \frac{1}{2} \int x^2 \sin(x) dx$$

STEP 7

Integrate the sum term by term and factor out constants:

$$= -\frac{1}{4} \int x^2 \sin(3x) dx + \frac{3}{4} \int x^2 \sin(x) dx$$

STEP 8

For the integrand $x^2 \sin(3x)$, integrate by parts, $\int f dg = fg - \int g df$, where $f = x^2$, $dg = \sin(3x) dx$, $df = 2x dx$, $g = -\frac{1}{3} \cos(3x)$:

$$= \frac{1}{12}x^2 \cos(3x) - \frac{1}{6} \int x \cos(3x) dx + \frac{3}{4} \int x^2 \sin(x) dx$$

STEP 9

For the integrand $x \cos(3x)$, integrate by parts, $\int f dg = fg - \int g df$, where $f = x$, $dg = \cos(3x) dx$, $df = dx$, $g = \frac{1}{3} \sin(3x)$:

$$= \frac{1}{12}x^2 \cos(3x) - \frac{1}{18}x \sin(3x) + \frac{1}{18} \int \sin(3x) dx + \frac{3}{4} \int x^2 \sin(x) dx$$

STEP 10

For the integrand $\sin(3x)$, substitute $u = 3x$ and $du = 3dx$:

$$= \frac{1}{12}x^2 \cos(3x) - \frac{1}{18}x \sin(3x) + \frac{1}{54} \int \sin(u) du + \frac{3}{4} \int x^2 \sin(x) dx$$

STEP 11

The integral of $\sin(u)$ is $-\cos(u)$:

$$= -\frac{\cos(u)}{54} + \frac{1}{12}x^2 \cos(3x) - \frac{1}{18}x \sin(3x) + \frac{3}{4} \int x^2 \sin(x) dx$$

STEP 12

For the integrand $x^2 \sin(x)$, integrate by parts, $\int f dg = fg - \int g df$, where $f = x^2$, $dg = \sin(x) dx$, $df = 2x dx$, $g = -\cos(x)$:

$$= -\frac{\cos(u)}{54} - \frac{3}{4}x^2 \cos(x) + \frac{1}{12}x^2 \cos(3x) - \frac{1}{18}x \sin(3x) + \frac{3}{2} \int x \cos(x) dx$$

STEP 13

For the integrand $x \cos(x)$, integrate by parts, $\int f dg = fg - \int g df$, where $f = x$, $dg = \cos(x) dx$, $df = dx$, $g = \sin(x)$:

$$= -\frac{\cos(u)}{54} - \frac{3}{4}x^2 \cos(x) + \frac{1}{12}x^2 \cos(3x) + \frac{3}{2}x \sin(x) - \frac{1}{18}x \sin(3x) - \frac{3}{2} \int \sin(x) dx$$

STEP 14

The integral of $\sin(x)$ is $-\cos(x)$:

$$= -\frac{\cos(u)}{54} - \frac{3}{4}x^2 \cos(x) + \frac{1}{12}x^2 \cos(3x) + \frac{3}{2}x \sin(x) - \frac{1}{18}x \sin(3x) + \frac{3 \cos(x)}{2} + \text{constant}$$

STEP 15

Substitute back for $u = 3x$:

$$= -\frac{3}{4}x^2 \cos(x) + \frac{1}{12}x^2 \cos(3x) + \frac{3}{2}x \sin(x) - \frac{1}{18}x \sin(3x) + \frac{3 \cos(x)}{2} - \frac{1}{54} \cos(3x) + \text{constant}$$

STEP 16

Which is equal to:

$$\text{Answer: } \frac{1}{108}(-81(x^2 - 2)\cos(x) + (9x^2 - 2)\cos(3x) - 6x(\sin(3x) - 27\sin(x))) + \text{constant}$$

7.5.3 UM-GPT: Show me how to Work Example 7.36

The integral $\int \sqrt{a^2 - x^2} dx$ can be solved using trigonometric substitution. We will consider the variable x to be $a \sin(\theta)$, where $-\frac{\pi}{2} \leq \theta \leq \frac{\pi}{2}$. This is because $\sin(\theta)$ ranges between -1 and 1, and the expression inside the square root, $a^2 - x^2$, should remain non-negative for real x .

Here is the step-by-step method to evaluate the integral:

1. **Substitute x with $a \sin(\theta)$:** Let $x = a \sin(\theta)$. Differentiate this with respect to θ to find dx :

$$dx = a \cos(\theta) d\theta$$

2. **Express the integral in terms of θ :** Substitute x and dx into the original integral:

$$\int \sqrt{a^2 - (a \sin(\theta))^2} \cdot a \cos(\theta) d\theta$$

3. **Simplify the integral:** Simplify the term under the square root and the integral becomes:

$$\int a^2 \cos^2(\theta) d\theta$$

Now, use the double-angle identity for cosine:

$$\cos^2(\theta) = \frac{1 + \cos(2\theta)}{2}$$

The integral now becomes:

$$\int \frac{a^2(1 + \cos(2\theta))}{2} d\theta$$

4. **Separate the integral:**

$$\frac{a^2}{2} \int 1 d\theta + \frac{a^2}{2} \int \cos(2\theta) d\theta$$

5. **Integrate:** The integration results in:

$$\frac{a^2}{2}\theta + \frac{a^2}{4}\sin(2\theta) + C$$

6. **Reverse the trigonometric substitution:** Recall that $x = a \sin(\theta)$, so $\sin(\theta) = \frac{x}{a}$. To find θ , use the inverse sine function (\arcsin):

$$\theta = \arcsin\left(\frac{x}{a}\right)$$

For the term with $\sin(2\theta)$, use the double-angle identity for sine:

$$\sin(2\theta) = 2\sin(\theta)\cos(\theta)$$

We know that $\sin(\theta) = \frac{x}{a}$, and since $\sin^2(\theta) + \cos^2(\theta) = 1$, $\cos(\theta)$ can be computed by:

$$\cos(\theta) = \sqrt{1 - \sin^2(\theta)} = \sqrt{1 - \left(\frac{x}{a}\right)^2}$$

Hence,

$$\sin(2\theta) = 2\left(\frac{x}{a}\right)\sqrt{1 - \left(\frac{x}{a}\right)^2}$$

7. **Plug θ and $\sin(2\theta)$ back into the expression:**

$$\frac{a^2}{2} \arcsin\left(\frac{x}{a}\right) + \frac{a^2}{4} \cdot 2\left(\frac{x}{a}\right) \sqrt{1 - \left(\frac{x}{a}\right)^2} + C$$

8. **Simplify:**

$$\frac{a^2}{2} \arcsin\left(\frac{x}{a}\right) + \frac{a}{2} x \sqrt{1 - \left(\frac{x}{a}\right)^2} + C$$

9. **Write the final answer:**

$$\int \sqrt{a^2 - x^2} dx = \frac{a^2}{2} \arcsin\left(\frac{x}{a}\right) + \frac{a}{2} x \sqrt{1 - \left(\frac{x}{a}\right)^2} + C$$

This is the final antiderivative of the initial integral in terms of the original variable x .

Are you sure the answer is correct for $a < 0$?

My apologies for the oversight. When specifying the domain with respect to the variable a , we need to consider that a represents the radius of a circle and is typically assumed to be a positive constant in the context of this integral problem. If a were negative, it would not make sense in terms of the geometry of the circle or the context of the trigonometric substitution we've done (since a radius is always non-negative).

For our integral and the substitutions made, we are working under the assumption that $a > 0$. So the answer:

$$\int \sqrt{a^2 - x^2} dx = \frac{a^2}{2} \arcsin\left(\frac{x}{a}\right) + \frac{a}{2} x \sqrt{1 - \left(\frac{x}{a}\right)^2} + C$$

is correct for the domain of x , assuming $a > 0$ and is given by $-a \leq x \leq a$. This domain corresponds to the values of x for which the integrand, $\sqrt{a^2 - x^2}$, is defined and real.

If a were to be less than zero, we would need to evaluate the integral with the absolute value $|a|$ instead, to make sure we're working with a positive radius. The final antiderivative would still be correct provided that we interpret a as $|a|$ and continue to consider x within the domain $-|a| \leq x \leq |a|$, ensuring a non-negative value under the square root. However, the problem is conventionally understood with a as positive.

Can you tell me the domain for which your answer is correct?

My apologies for the oversight. When specifying the domain with respect to the variable a , we need to consider that a represents the radius of a circle and is typically assumed to be a positive constant in the context of this integral problem. If a were negative, it would not make sense in terms of the geometry of the circle or the context of the trigonometric substitution we've done (since a radius is always non-negative).

For our integral and the substitutions made, we are working under the assumption that $a > 0$. So the answer:

$$\int \sqrt{a^2 - x^2} dx = \frac{a^2}{2} \arcsin\left(\frac{x}{a}\right) + \frac{a}{2} x \sqrt{1 - \left(\frac{x}{a}\right)^2} + C$$

is correct for the domain of x , assuming $a > 0$ and is given by $-a \leq x \leq a$. This domain corresponds to the values of x for which the integrand, $\sqrt{a^2 - x^2}$, is defined and real.

If a were to be less than zero, we would need to evaluate the integral with the absolute value $|a|$ instead, to make sure we're working with a positive radius. The final antiderivative would still be correct provided that we interpret a as $|a|$ and continue to consider x within the domain $-|a| \leq x \leq |a|$, ensuring a non-negative value under the square root. However, the problem is conventionally understood with a as positive.

7.6 (Optional Read:) Proofs Associated with the Chapter

Proof of Prop. 7.8 (First Fundamental Theorem of Calculus)

Let $f : [a, b] \rightarrow \mathbb{R}$ be a continuous real-valued function defined on a closed interval $[a, b]$. Then, the function $F : [a, b] \rightarrow \mathbb{R}$ defined by

$$F(x) = \int_a^x f(t) dt \quad (7.25)$$

for all $x \in [a, b]$, is continuous on $[a, b]$, differentiable on the open interval (a, b) , and $F'(x) = f(x)$ for all $x \in (a, b)$. You will also see it written as

$$\frac{d}{dx} F(x) = \frac{d}{dx} \left[\int_a^x f(t) dt \right] = f(x), \quad (7.26)$$

to emphasize that differentiation undoes, or inverses, integration.

Proof: We give a proof that is straightforward and fully rigorous for functions $f : [a, b] \rightarrow \mathbb{R}$ that are **non-decreasing**. The general case is proven in Michigan's Math 451.

Claim 1: Continuity of F : Because f is non-decreasing³,

$$M := \max_{x \in [a, b]} |f(x)| = \max\{|f(a)|, f(b)\} < \infty.$$

Let $\epsilon > 0$ and define $\delta := \frac{\epsilon}{M+1} > 0$. We will show for all $x, x_0 \in [a, b]$, that $|x - x_0| < \delta$ implies

$$|F(x) - F(x_0)| < \epsilon.$$

Proof: Consider the change in F as x varies from x_0 ,

$$F(x) - F(x_0) = \int_a^x f(t) dt - \int_a^{x_0} f(t) dt = \int_{x_0}^x f(t) dt.$$

This equation holds by the properties of the definite integral, specifically the additivity over adjacent intervals. The absolute value of the integral is bounded by

$$|F(x) - F(x_0)| = \left| \int_{x_0}^x f(t) dt \right| \leq \int_{x_0}^x |f(t)| dt.$$

Given that $|f(t)| \leq M$ for all t in $[a, b]$, the integral simplifies further,

$$|F(x) - F(x_0)| \leq \int_{x_0}^x M dt = M|x - x_0|.$$

The integral $\int_{x_0}^x M dt$ evaluates to $M \times |x - x_0|$ because the integrand M is constant, and the length of the interval from x_0 to x is $|x - x_0|$.

Now, by the choice of δ , if $|x - x_0| < \delta$ then:

$$|F(x) - F(x_0)| \leq M|x - x_0| < M\delta = M \left(\frac{\epsilon}{M+1} \right) < \epsilon.$$

Thus, we have shown that:

$$|F(x) - F(x_0)| < \epsilon,$$

for all x such that $|x - x_0| < \delta$. Therefore, F is continuous at x_0 , and since x_0 was arbitrary, F is continuous on $[a, b]$.

³If $f(b) \geq 0$, then no absolute value sign is needed on $f(b)$. If $f(b) < 0$, then $f(a) \leq f(b) \implies |f(a)| \geq |f(b)| \geq f(b)$. Hence, once again, no absolute value sign is needed.

Claim 2: $F'(x) = f(x)$ for all $a < x < b$:

Proof: Consider the difference quotient for F at some point x in (a, b) ,

$$\begin{aligned}\frac{F(x+h) - F(x)}{h} &= \frac{1}{h} \left(\int_a^{x+h} f(t) dt - \int_a^x f(t) dt \right) \\ &= \frac{1}{h} \int_x^{x+h} f(t) dt.\end{aligned}$$

Assume first that $h > 0$. Then, using once again that f is non-decreasing,

$$f(x) \cdot h \leq \int_x^{x+h} f(t) dt \leq f(x+h) \cdot h.$$

Plugging into the previous equation and canceling the h 's, we have

$$f(x) \leq \frac{1}{h} \int_x^{x+h} f(t) dt \leq f(x+h).$$

Because f is continuous, $\lim_{h \rightarrow 0^+} f(x+h) = f(x)$. Therefore, by the Squeeze Theorem,

$$\lim_{h \rightarrow 0^+} \frac{1}{h} \int_x^{x+h} f(t) dt = f(x).$$

The same argument can be repeated with minor variations to show that $\lim_{h \rightarrow 0^-} \frac{1}{h} \int_x^{x+h} f(t) dt = f(x)$. This completes the proof that $F'(x) = f(x)$, which is a fundamental connection between the differential and integral calculus for continuous functions. ■

Proof of Prop. 7.8 (Second Fundamental Theorem of Calculus)

Let $f : [a, b] \rightarrow \mathbb{R}$ be a real-valued function defined on a closed interval $[a, b]$ that admits an antiderivative F on $[a, b]$, that is, $F : [a, b] \rightarrow \mathbb{R}$ is differentiable on (a, b) and for all $x \in (a, b)$, $F'(x) = f(x)$. Then, f is integrable on $[a, b]$ and

$$\int_a^b f(x) dx = F(x) \Big|_a^b := F(b) - F(a). \quad (7.27)$$

We rewrite this as

$$\int_a^b F'(x) dx = F(x) \Big|_a^b := F(b) - F(a), \quad (7.28)$$

to underline that integration undoes, or inverses, differentiation.

Proof: We use the method given in [Proof of the Fundamental Theorem of Calculus](#), Math 121 Calculus II, D Joyce, Spring 2013, with one modification: we ASSUME that $f : [a, b] \rightarrow \mathbb{R}$ is continuous, thereby guaranteeing that $f(x) = F'(x)$ is integrable.

Define

$$G(x) = \int_a^x F'(t) dt, \text{ for } x \in [a, b].$$

Then, by the First Fundamental Theorem of Calculus, $G'(x) = F'(x)$. Therefore, G and F differ by a constant C , that is, $G(x) - F(x) = C$ for all $x \in [a, b]$. But

$$G(a) = \int_a^a F'(t) dt = 0,$$

and, therefore, $G(a) - F(a) = C$, so $C = -F(a)$.

It follows that $G(x) - F(x) = -F(a)$ for all $x \in [a, b]$. In particular, $G(b) - F(b) = -F(a)$, so $G(b) = F(b) - F(a)$, that is,

$$\int_a^b F'(t) dt = F(b) - F(a).$$

This completes the proof that the integral of a function over an interval is equal to the difference in the values of any of its antiderivatives evaluated at the endpoints of the interval, affirming that integration can be seen as the inverse operation to differentiation. ■

We next explore how restrictive it was to assume that a function is non-decreasing.

1) Property of max and min functions: For any $y \in \mathbb{R}$, the functions $\max(y, 0)$ and $\min(y, 0)$ have the property that:

$$\max(y, 0) + \min(y, 0) = y$$

Proof: Consider two cases based on the value of y :

Case 1: $y \geq 0$. In this case, we have

$$\max(y, 0) = y \quad \text{and} \quad \min(y, 0) = 0$$

Thus,

$$\max(y, 0) + \min(y, 0) = y + 0 = y$$

Case 2: $y < 0$. In this case, we have

$$\max(y, 0) = 0 \quad \text{and} \quad \min(y, 0) = y$$

Thus,

$$\max(y, 0) + \min(y, 0) = 0 + y = y$$

In both cases, $\max(y, 0) + \min(y, 0)$ simplifies to y , proving the statement. ■

2) Decomposition of continuously differentiable functions: If $f : [a, b] \rightarrow \mathbb{R}$ is a continuously differentiable function, then f can be expressed as the difference of two continuous, non-decreasing functions.

Construction: Define two functions g and k over $[a, b]$ by:

$$g(x) = f(a) + \int_a^x \max(f'(t), 0) dt$$

and

$$k(x) = - \int_a^x \min(f'(t), 0) dt$$

Properties: Since f is continuously differentiable, f' exists and is continuous. The function g is non-decreasing because the integrand, $\max(f'(t), 0)$, is non-negative and strictly positive whenever $f'(t) > 0$.

The function k is non-decreasing because the integrand $-\min(f'(t), 0)$ is non-negative, being zero when $f'(t) \geq 0$ and positive when $f'(t) < 0$.

Expressing f as a difference:

$$\begin{aligned} g(x) - k(x) &= f(a) + \left(\int_a^x \max(f'(t), 0) dt \right) - \left(- \int_a^x \min(f'(t), 0) dt \right) \\ &= f(a) + \int_a^x (\max(f'(t), 0) + \min(f'(t), 0)) dt \\ &= f(a) + \int_a^x f'(t) dt \\ &= f(a) + (f(x) - f(a)) \\ &= f(x). \end{aligned}$$

This construction expresses f as the difference of two continuous, non-decreasing functions over the interval $[a, b]$. Hence, the proof of Prop. 7.8 given for non-decreasing functions can be applied to functions that are continuously differentiable.

Chapter 8

Improper Integrals

Learning Objectives

By the end of this chapter, students should be able to:

- Define improper integrals and understand their significance within the context of calculus.
- Identify integrals that may pose difficulties due to their unbounded nature or the behavior of the function being integrated.
- Answer the pressing question: can an integral be doubly improper?
- Explore the applications of improper integrals, emphasizing their practical importance in Statistics and Probability.

Outcomes

Upon successful completion of this chapter, students will:

- Master the technique of integrating over unbounded domains using limits.
- Learn to manage and integrate functions with singularities or discontinuities by applying limits to circumvent infinite behavior at finite points, thus handling vertical asymptotes effectively.
- Gain a comprehensive understanding of the analytical methods for solving improper integrals, employing antiderivatives to facilitate calculation.
- Acquire skills in applying numerical methods, specifically Julia's QuadGK, for evaluating improper integrals, particularly when analytical solutions are challenging to obtain.
- Delve into two probability theory examples that utilize improper integrals, reinforcing the concept's application in real-world scenarios and theoretical studies.

8.1 Introduction

What a name, right? **Improper?** In short, the term “improper integral” serves as a descriptive label to indicate that the integral in question deviates from the standard conditions that apply to definite integrals, and hence, extra caution is advised. The word “improper” in this context doesn’t imply that there’s something wrong or less valid about the integrals we study in this Chapter. Rather, it indicates that they don’t fit within the framework of standard definite integrals and require special techniques for their evaluation. These techniques involve taking limits, as the interval of integration or the function itself is unbounded. The terminology likely originated to help mathematicians categorize and study different types of integrals, especially as the field of real analysis began to formalize concepts related to integration. **The term helps to signal that additional care or “alternative methods” may be needed to evaluate the integral.**

8.2 Type-I Improper Integrals: Unbounded Limits of Integration

We first look at cases where the integrand, f , is piecewise continuous, and hence, for any closed and bounded interval $[a, b]$, the definite integral $\int_a^b f(x) dx$ exists and is finite. We are interested in what happens as one or more of the limits of integration is unbounded, that is, $\int_a^\infty f(x) dx$, $\int_{-\infty}^b f(x) dx$, and a bit later, $\int_{-\infty}^\infty f(x) dx$. So that we can understand what the issues are, let’s look at a few examples first.

- Let $f : [0, \infty) \rightarrow \mathbb{R}$ by $f(x) = e^{-x}$. Then, $\int_0^b f(x) dx = -e^{-x} \Big|_0^b = 1 - e^{-b}$. Hence, $\lim_{b \rightarrow \infty} \int_0^b f(x) dx = \lim_{b \rightarrow \infty} (1 - e^{-b}) = 1.0$. Because the value of the integral is finite, we say it is **convergent**.
- Let $f : [0, \infty) \rightarrow \mathbb{R}$ by $f(x) = 1$. Then, $\int_0^b f(x) dx = b$. Hence, $\lim_{b \rightarrow \infty} \int_0^b f(x) dx = \infty$. Because the value of the integral is unbounded, we say it is **divergent**.
- Let $f : (-\infty, 0] \rightarrow \mathbb{R}$ by $f(x) = \cos(x)$. Then, $\int_a^0 f(x) dx = \sin(x) \Big|_a^0 = -\sin(a)$. Hence, $\lim_{a \rightarrow -\infty} \int_a^0 f(x) dx = \lim_{a \rightarrow -\infty} -\sin(a)$ is **undefined**.

The following definitions are based on the previous three examples.

Unbounded Upper or Lower Limit: Take 1

Definition 8.1. Suppose that $f : [a, \infty) \rightarrow \mathbb{R}$ is a real-valued function, and for all $b > a$ (finite), $\int_a^b f(x) dx$ exists and is finite^a. Then,

$$\int_a^\infty f(x) dx := \lim_{b \rightarrow \infty} \int_a^b f(x) dx, \quad (8.1)$$

if the limit exists. The integral is **undefined** otherwise. If the limit exists but is unbounded (i.e., equals $\pm\infty$), then we say the integral is **divergent**, while if the limit exists and is finite, we say the integral is **convergent**.

Similarly, suppose that $f : (-\infty, b] \rightarrow \mathbb{R}$ is a real-valued function, and for all $a < b$ (finite), $\int_a^b f(x) dx$ exists and is finite^b. Then,

$$\int_{-\infty}^b f(x) dx := \lim_{a \rightarrow -\infty} \int_a^b f(x) dx, \quad (8.2)$$

if the limit exists. The integral is **undefined** otherwise. If the limit exists but is unbounded, (i.e., equals $\pm\infty$), then we say the integral is **divergent**, while if the limit exists and is finite, we say the integral is **convergent**.

^aIf $f : [a, \infty) \rightarrow \mathbb{R}$ is piecewise continuous, then this condition is met.

^bIf $f : (-\infty, b] \rightarrow \mathbb{R}$ is piecewise continuous, then this condition is met.

Example 8.2. Determine values, if possible, for the following improper integrals. You may wish to use the antiderivatives in Table 7.9 and Table 7.18 because the integrands are very standard functions.

- (a) $\int_0^\infty x \, dx$.
 (b) $\int_0^\infty e^{-x} \, dx$.
 (c) $\int_{-\infty}^0 e^{-x} \, dx$.
 (d) $\int_{-\infty}^0 \sin(x) \, dx$.

Solutions:

(a) **Ans.** $\int_0^\infty x \, dx = \infty$.

From Table. 7.9, $\int_0^b x \, dx = \frac{1}{2}x^2 \Big|_0^b = \frac{b^2}{2}$. Hence, $\int_0^\infty x \, dx := \lim_{b \rightarrow \infty} \int_0^b x \, dx = \lim_{b \rightarrow \infty} \frac{1}{2}b^2 = \infty$.

(b) **Ans.** $\int_0^\infty e^{-x} \, dx = 1$.

From Table. 7.18, $\int_0^b e^{-x} \, dx = -e^{-x} \Big|_0^b = -(e^{-b} - 1) = 1 - e^{-b}$. Hence, $\int_0^\infty e^{-x} \, dx := \lim_{b \rightarrow \infty} \int_0^b e^{-x} \, dx = \lim_{b \rightarrow \infty} 1 - e^{-b} = 1$. We can also use quadGK:

```

1 using QuadGK
2
3 f(x) = exp(-x)
4 Data = Array{Float64}(undef, 0, 3)
5 #
6 for b=0:2:30
7     value, error = quadgk(f, 0, b)
8     Data = [Data; b value error]
9 end
10 Data

```

Output

```

16×3 Matrix{Float64}:
 0.0  0.0   0.0
 2.0  0.864665 7.77156e-16
 4.0  0.981684 1.00392e-11
 6.0  0.997521 1.74719e-9
 8.0  0.999665 1.02231e-11
10.0  0.999955 1.804e-10
12.0  0.999994 1.75152e-9
14.0  0.999999 1.12599e-8
16.0  1.0      2.81922e-11
18.0  1.0      7.20663e-11
20.0  1.0      2.09573e-10
22.0  1.0      6.26783e-10
24.0  1.0      1.77795e-9
26.0  1.0      4.65932e-9
28.0  1.0      1.1276e-8
30.0  1.0      2.97104e-11

```

The improper integral converges very quickly to one.

(c) **Ans.** $\int_{-\infty}^0 e^{-x} \, dx = \infty$.

From Table. 7.18, $\int_a^0 e^{-x} \, dx = -(1 - e^{-a}) = e^{-a} - 1$. Hence, $\int_{-\infty}^0 e^{-x} \, dx := \lim_{a \rightarrow -\infty} \int_a^0 e^{-x} \, dx = \lim_{a \rightarrow -\infty} e^{-a} - 1 = \infty$.

- (d) **Ans.** $\int_{-\infty}^0 \sin(x) dx$ is undefined because the limit does not exist.

Intuition: We know that $\sin(x)$ oscillates between plus and minus one. Moreover, $\sin(x + \pi) = -\sin(x)$, and hence, the integral over 2π -interval $[k\pi, (k+2)\pi]$ is equal to zero. Therefore, we do not expect the integral to converge. But we check it out anyway, just to confirm our intuition.

We can use antiderivatives for an analytical approach: $\int_a^0 \sin(x) dx = -\cos(x) \Big|_a^0 = \cos(a)$. Hence, $\int_{-\infty}^0 \sin(x) dx = \lim_{a \rightarrow -\infty} \cos(a)$, which does not exist.

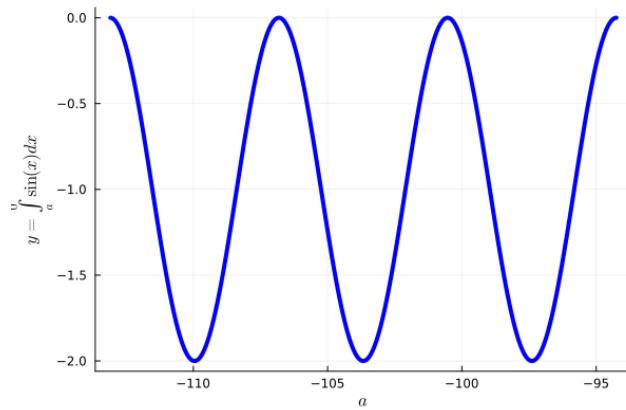
We can also see via numerical integration that something is not right!

```

1 using QuadGK, Plots, LaTeXStrings
2
3 f(x) = sin(x)
4 Data = Array{Float64}(undef, 0, 3)
5 #
6 a0=-30*pi
7 for a=a0:-.01:(a0-6pi)
8     value, error = quadgk(f, a, 0)
9     Data = [Data; a value error]
10 end
11 p1 = plot(Data[:, 1], Data[:, 2], linewidth=4, color=:blue, label=false)
12 plot!(xlabel=L"$a$",
13 ylabel=L"$y = \int_a^0 \sin(x) dx$")
14 display(p1)

```

Output



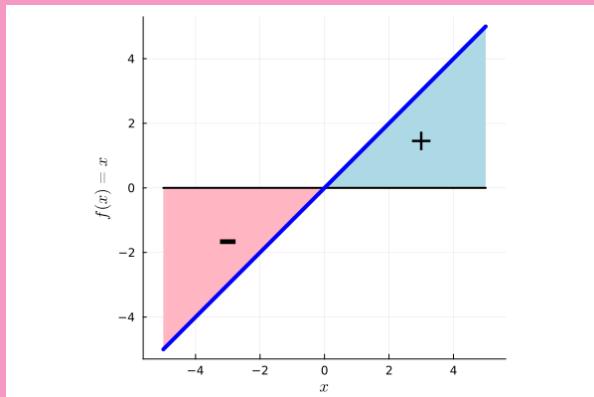
Functions defined on $(-\infty, \infty)$ pose an extra layer of difficulty. To see why, consider, $f : \mathbb{R} \rightarrow \mathbb{R}$ by $f(x) = 4x^3$. Then, $\int_a^b f(x) dx = x^4 \Big|_a^b = b^4 - a^4$. Hence, $\lim_{a \rightarrow -\infty} \int_a^b f(x) dx = \lim_{a \rightarrow -\infty} b^4 - a^4 = -\infty$. But if we evaluate $\lim_{b \rightarrow \infty} \int_a^b f(x) dx = \lim_{b \rightarrow \infty} b^4 - a^4 = \infty$. So what value should we assign to $\int_{-\infty}^{\infty} f(x) dx$?

Secrets of the Arcane 8.3: A Word of Caution for Doubly-infinite Integrals

In general, $\int_{-\infty}^{\infty} f(x) dx \neq \lim_{c \rightarrow \infty} \int_{-c}^c f(x) dx$. You must compute the two limits, $\lim_{c \rightarrow \infty} \int_{-c}^0 f(x) dx$ and $\lim_{c \rightarrow \infty} \int_0^c f(x) dx$, SEPARATELY!!!. Why? Consider

$$\int_{-c}^c x dx.$$

Then, as illustrated in the figure below, the magnitude of the negative area is exactly equal to the positive area, and hence for any interval $[-c, c]$ that is symmetric about the origin, $\int_{-c}^c x dx = 0.0$, i.e., the areas cancel out, exactly. However, $\int_{-\infty}^0 x dx = -\infty$ and $\int_0^{\infty} x dx = +\infty$, which do NOT cancel out; the their sum is undefined. Moreover, if we computed $\lim_{c \rightarrow \infty} \int_{-c}^{c+1} x dx$, we'd obtain positive infinity, and if we computed $\lim_{c \rightarrow \infty} \int_{-c^3}^{c^2} x dx$, we'd obtain negative infinity. In other words, the limit changes if we advance more quickly in one direction than the other. We avoid this problem by assessing the two limits separately and then analyzing whether their sum makes sense (or not). The sum will make sense except when both limits are unbounded and have opposite signs, giving rise to the notorious, " $\infty - \infty$ ".



Note: The symmetric limit, $\lim_{c \rightarrow \infty} \int_{-c}^c f(x) dx$, has a name: the **Cauchy Principal Value**; see Wikipedia for more information.

Unbounded Upper and Lower Limit: Take 2

Definition 8.4. Suppose that $f : (-\infty, \infty) \rightarrow \mathbb{R}$ is a real-valued function, and for all $a < b$ (both finite), $\int_a^b f(x) dx$ exists and is finite^a. Then,

$$\int_{-\infty}^{\infty} f(x) dx := \lim_{a \rightarrow -\infty} \int_a^0 f(x) dx + \lim_{b \rightarrow \infty} \int_0^b f(x) dx, \quad (8.3)$$

if both limits exist and are finite (i.e., both integrals are convergent). If one of the limits is undefined, then the integral is undefined. If both limits exist but at least one is unbounded (i.e., $\pm\infty$), more care is required to avoid $\infty - \infty$. Hence,

- if one limit is unbounded while the other is finite, then the value of the unbounded limit is the value of the integral;
- if both limits are unbounded and have the same sign, then the common value of the unbounded limits is the value of the integral; and as was said above,
- if one limit equals $+\infty$ and the other $-\infty$, then the integral is undefined.

^aIf $f : (-\infty, \infty) \rightarrow \mathbb{R}$ is piecewise continuous, then this condition is met.

Example 8.5. Determine values, if possible, for the following improper integrals. You may wish to use the antiderivatives in Table 7.9 and Table 7.18 because the integrands are very standard functions.

(a) $\int_{-\infty}^{\infty} x dx$.

(b) $\int_{-\infty}^{\infty} x^2 dx$.

$$(c) \int_{-\infty}^{\infty} e^{-x} dx.$$

Solutions:

(a) **Ans.** $\int_{-\infty}^{\infty} x dx$ is undefined because we have “ $-\infty + \infty$ ”.

From Table 7.9, $\int_a^0 x dx = -\frac{1}{2}a^2$. Hence, $\int_0^{\infty} x dx := \lim_{a \rightarrow -\infty} \int_a^0 x dx = \lim_{a \rightarrow -\infty} -\frac{1}{2}a^2 = -\infty$. From part (a) of this Example, $\int_0^{\infty} x dx = \infty$. Hence, $\int_{-\infty}^{\infty} x dx = -\infty + \infty = \text{undefined}$.

(b) **Ans.** $\int_{-\infty}^{\infty} x^2 dx = \infty$ because the two limits share the same sign and at least one is unbounded.

From Table 7.9, $\int_a^0 x^2 dx = -\frac{1}{3}a^3$. Hence, $\int_{-\infty}^0 x^2 dx := \lim_{a \rightarrow -\infty} \int_a^0 x^2 dx = \lim_{a \rightarrow -\infty} -\frac{1}{3}a^3 = +\infty$. Similarly, $\int_0^b x^2 dx = \frac{1}{3}b^3$. Hence, $\int_0^{\infty} x^2 dx := \lim_{b \rightarrow \infty} \int_0^b x^2 dx = \lim_{b \rightarrow \infty} \frac{1}{3}b^3 = \infty$. Therefore, $\int_{-\infty}^{\infty} x^2 dx = \infty + \infty = \infty$.

(c) **Ans.** $\int_{-\infty}^{\infty} e^{-x} dx = \infty$ because one limit is bounded and the other is unbounded.

In Example 8.2, we computed each limit separately. The rest is left to the learner. ■

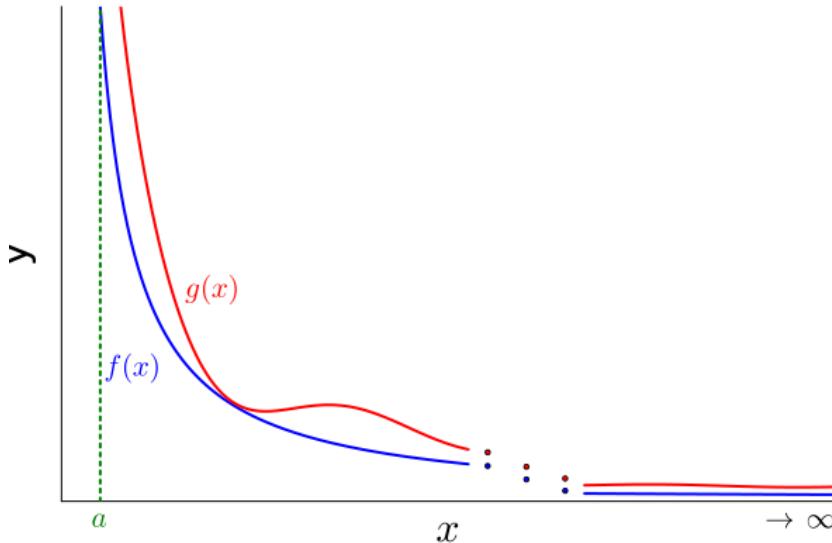


Figure 8.1: **Comparison Test:** Suppose that $f : [a, \infty) \rightarrow [0, \infty)$ and $g : [a, \infty) \rightarrow [0, \infty)$ are piecewise continuous functions so that their Riemann Integrals exist and are finite for any bounded interval $[a, b]$. When are they finite on $[a, \infty)$? In the plot, we observe that $f(x) \leq g(x)$ for all $x \in [a, \infty)$, and hence, on one hand, if the area under g is finite, so must be the area under f , and on the other hand, if the area under f is unbounded, so must be the area under g . **These simple observations are the heart of the comparison test.**

8.2.1 Comparison Test

It would be nice to have some rules of thumb for knowing when an improper integral should exist and be finite and when it is likely to be unbounded. Figure 8.1 shows that by comparing one function against another function's KNOWN behavior, we can make conclusions about an integral of the form $\int_a^{\infty} f(x) dx$ by comparing it to another integral $\int_a^{\infty} g(x) dx$. Indeed, if $0 \leq f(x) \leq g(x)$, then

- $\int_a^{\infty} g(x) dx < \infty \implies \int_a^{\infty} f(x) dx < \infty$, and

- $\int_a^\infty f(x) dx = \infty \implies \int_a^\infty g(x) dx = \infty.$

That is interesting, but what are some relevant comparison functions to use? We dig into our antiderivative bag of easy integrals and make some observations.

- (a) For all $p > 1$, $\int_1^b \frac{1}{x^p} dx = \int_1^b x^{-p} dx = \frac{x^{1-p}}{1-p} \Big|_1^b = \frac{b^{1-p}}{1-p} - \frac{1}{1-p} = \frac{1}{1-p} \cdot \left(\frac{1}{b^{p-1}} - 1 \right)$. For $p > 1$, it follows that $p-1 > 0$, and hence $\lim_{b \rightarrow \infty} \frac{1}{b^{p-1}} = 0$. Therefore, $\int_1^\infty \frac{1}{x^p} dx = \frac{-1}{1-p} = \frac{1}{p-1} < \infty$.
- (b) $\int_1^b \frac{1}{x} dx = \ln(x) \Big|_1^b = \ln(b)$ and hence $\int_1^\infty \frac{1}{x} dx = \lim_{b \rightarrow \infty} \ln(b) = \infty$. This makes sense because as $p \rightarrow 1^+$ in part (a), we see that the value of the integral, $\frac{1}{p-1}$, diverges to infinity.
- (c) For all $0 \leq p < 1$, $\int_1^b \frac{1}{x^p} dx = \int_1^b x^{-p} dx = \frac{x^{1-p}}{1-p} \Big|_1^b = \frac{b^{1-p}}{1-p} - \frac{1}{1-p} = \frac{1}{1-p} \cdot (b^{1-p} - 1)$. For $p < 1$, it follows that $1-p > 0$, and hence $\lim_{b \rightarrow \infty} b^{1-p} = \infty$. Therefore, $\int_1^\infty \frac{1}{x^p} dx = \infty$.

The above analysis says that $\frac{1}{x}$ is a “line of demarcation” between an improper integral of $\frac{1}{x^p}$ being unbounded (for $0 \leq p \leq 1$) and being finite (for $p > 1$). Hence, we can use $\frac{1}{x^p}$ in the **Comparison Test** illustrated in Fig. 8.1 and stated more formally below.

Recall that for any $a < b$ finite and any piecewise continuous function $f : [a, b] \rightarrow \mathbb{R}$, the Riemann integral, $\int_a^b f(x) dx$, exists and is finite. Hence, an improper integral, e.g., $\int_a^\infty f(x) dx$, is finite (or is infinite, or does not exist) if, and only if, for any $c \geq a$ finite, $\int_c^\infty f(x) dx$, is finite (or is infinite, or does not exist), respectively. In other words, because $\int_a^c f(x) dx$ is guaranteed to be finite, we can ignore it when trying to understand if the $\int_a^\infty f(x) dx$, is convergent or not. Why is this a helpful observation? We'll see in examples that by choosing c sufficiently large, we can sometimes use “asymptotic reasoning” to form more easily an upper bound or a lower bound on $f : [c, \infty) \rightarrow \mathbb{R}$ than we can for $f : [a, \infty) \rightarrow \mathbb{R}$.

Proposition 8.6: Comparison Test

Consider a **non-negative piecewise continuous** function $f : [0, \infty) \rightarrow [0, \infty)$ and a second non-negative piecewise continuous function, $g : [c, \infty) \rightarrow [0, \infty)$, where $c \geq 0$. Then the following are true:

(a) Test for a convergent (aka, finite) improper integral:

- Suppose that for all $x \geq c$, $f(x) \leq g(x)$ (one says that f is upper bounded by g on $[c, \infty)$, or f is eventually upper bounded by g), and
- $\int_c^\infty g(x) dx$ exists and is finite.

Then $\int_0^\infty f(x) dx$ exists and is finite. (In other words, f being eventually upper bounded by a function with a finite improper integral ensures that its improper integral is also finite).

(b) Test for a divergent (aka, unbounded) improper integral:

- Suppose that for all $x \geq c$, $g(x) \leq f(x)$ (one says that f is lower bounded by g on $[c, \infty)$, or f is eventually lower bounded by g), and
- $\int_c^\infty g(x) dx$ exists and is infinite.

Then $\int_0^\infty f(x) dx$ exists and is infinite. (In other words, f being lower bounded by a function with an unbounded improper integral implies that its improper integral is also unbounded.)

Note: We leave the cases for functions defined on $(-\infty, 0]$ to the learner.

Note: The second part of Fig. 8.1 shows that if a piecewise continuous function $f : [0, \infty) \rightarrow \mathbb{R}$ takes on both positive and negative values, then useful information is obtained by applying the comparison test to the function's absolute value.

Example 8.7. Use the Comparison Test to determine which of the following improper integrals are convergent or divergent.

- (a) $\int_0^\infty \frac{1}{x^2+1} dx$, where we note that the denominator has no real roots, and hence $\frac{1}{x^2+1}$ is finite for all $x \in \mathbb{R}$.
- (b) $\int_0^\infty \frac{\sin(x)}{x^2+1} dx$.
- (c) $\int_0^\infty \frac{x}{x^2+1} dx$.
- (d) $\int_0^\infty \frac{x^3+4}{x^5-3x+21} dx$, where the denominator has one real root at $x \approx -1.93017$ and four complex roots.
- (e) $\int_0^\infty \frac{x+e^{-x}}{x^2+1} dx$.
- (f) $\int_0^\infty e^{-x} \cdot \frac{x^4+111}{x^2+1} dx$.

Solutions:

- (a) $\int_0^\infty \frac{1}{x^2+1} dx$ is convergent. Let $f(x) := \frac{1}{x^2+1} = \frac{\frac{1}{x^2}}{1+\frac{1}{x^2}}$, where we divided the numerator and denominator by x^2 . Next, we note that $1 + \frac{1}{x^2} > 1$ for all $x \geq 0$. Hence, $f(x) \leq g(x) := \frac{1}{x^2}$ for all $x \geq 0$. We know that $\int_0^\infty g(x) dx$ is convergent, and therefore, so is $\int_0^\infty f(x) dx$.

Note: For this integrand, we know that $\text{atan}(x)$ is an antiderivative of $\frac{1}{x^2+1}$. Hence, $\int_0^\infty \frac{1}{x^2+1} dx = \text{atan}(x) \Big|_0^\infty = \frac{\pi}{2}$.

- (b) $\int_0^\infty \frac{\sin(x)}{x^2+1} dx$ is convergent. Because $|\sin(x)| \leq 1$, we have that $\left| \frac{\sin(x)}{x^2+1} \right| \leq \frac{1}{x^2+1}$. From (a), we know that $\int_0^\infty \frac{1}{x^2+1} dx$ is convergent. Hence, $\int_0^\infty \frac{\sin(x)}{x^2+1} dx$ is also convergent.

- (c) $\int_0^\infty \frac{x}{x^2+1} dx$ is divergent. Did you see that coming? Let $f(x) := \frac{x}{x^2+1} = \frac{\frac{1}{x}}{1+\frac{1}{x^2}}$, where we divided the numerator and denominator by x^2 . Next, we note that $1 + \frac{1}{x^2} \leq 2$ for all $x \geq 0$. Hence, $f(x) \geq g(x) := \frac{1}{2} \cdot \frac{1}{x}$ for all $x \geq 0$. We know that $\int_0^\infty g(x) dx$ is divergent, and therefore, so is $\int_0^\infty f(x) dx$.

- (d) $\int_0^\infty \frac{x^3+4}{x^5-3x+21} dx$ is convergent. First of all, why was it important to note that the denominator has one real root at $x \approx -1.93017$ and four complex roots? Because there are no positive real roots, the denominator does not vanish for any $x \geq 0$, and hence, we are never diving by zero. Therefore, $f(x) := \frac{x^3+4}{x^5-3x+21}$ is a continuous function on $[0, \infty)$.

Now, we turn to checking if the integral is convergent or divergent. We divide the numerator and denominator by x^5 , giving us,

$$f(x) = \frac{\frac{1}{x^2} + \frac{4}{x^5}}{1 - \frac{3}{x^4} + \frac{21}{x^5}}.$$

For all $x \geq 10$, $\text{den}(x) := 1 - \frac{3}{x^4} + \frac{21}{x^5} \geq 1 - \frac{3}{10^4} + \frac{21}{10^5} \geq 0.99$. Using this bound, we have that for all $x \geq 10$

$$f(x) = \frac{\frac{1}{x^2} + \frac{4}{x^5}}{1 - \frac{3}{x^4} + \frac{21}{x^5}} \leq \frac{\frac{1}{x^2} + \frac{4}{x^5}}{0.99}.$$

Because $\frac{1}{x^2}$ and $\frac{1}{x^5}$ are both convergent, we deduce that $\int_{10}^\infty \frac{x^3+4}{x^5-3x+21} dx$ is convergent, which then implies that $\int_0^\infty \frac{x^3+4}{x^5-3x+21} dx$ is convergent.

- (e) $\int_0^\infty \frac{x+e^{-x}}{x^2+1} dx$ is divergent. We write the integral as

$$\int_0^\infty \frac{x+e^{-x}}{x^2+1} dx = \int_0^\infty \frac{x}{x^2+1} dx + \int_0^\infty \frac{e^{-x}}{x^2+1} dx.$$

Because $e^{-x} \leq 1$ for all $x \geq 0$, $\int_0^\infty \frac{e^{-x}}{x^2+1} dx \leq \int_0^\infty \frac{1}{x^2+1} dx$, which we know to be convergent. However, we already know $\int_0^\infty \frac{x}{x^2+1} dx$ is divergent. Because the sum of something unbounded and something bounded is unbounded, we conclude that the overall integral is divergent. If both had been divergent, we would have to check that we are not committing the crime of $\infty - \infty$! Clear enough?

(f) $\int_0^\infty e^{-x} \cdot \frac{x^4+111}{x^2+1} dx$ is convergent. The intuition is that decaying exponentials crush exploding polynomials, which we established in Proposition 2.28 and Corollary 2.29. But we can do better and show this analytically. Ignoring the exponential for a moment and focusing on the rational function,

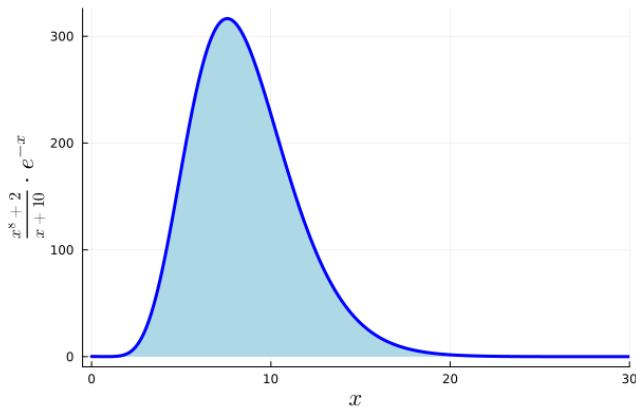
$$\frac{x^4 + 111}{x^2 + 1} = \frac{x^4}{x^2} \frac{1 + \frac{111}{x^4}}{1 + \frac{1}{x^2}} \quad (\text{factor our highest power from numerator and denominator})$$

$$\leq x^2 \frac{1 + \frac{111}{x^4}}{1} \quad (\text{lower bound on the denominator holds for all } x \geq 0)$$

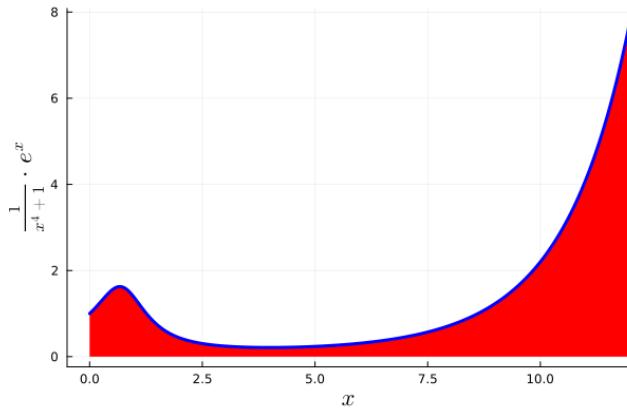
$$\leq 2x^2 \quad (\text{for all } x > \sqrt[4]{111}).$$

This yields, $\int_{10}^\infty e^{-x} \cdot \frac{x^4+111}{x^2+1} dx \leq \int_{10}^\infty 2x^2 e^{-x} dx = -2x^2 e^{-x} + 4xe^{-x} - 4e^{-x} \Big|_{10}^\infty \approx 0.003723$, and because $\int_0^{10} e^{-x} \cdot \frac{x^4+111}{x^2+1} dx$ is finite, the overall integral is then finite, in other words, convergent. ■

Way back in Proposition 2.28 and Corollary 2.29, we learned that exponentials dominate polynomials and rational functions when taking limits as $x \rightarrow \infty$. You may suspect that the presence of a decaying or exploding exponential in an integrand affects whether the integral is convergent or not, and you would be right!



(a)



(b)

Figure 8.2: When Present, Exponentials Determine Convergence of Improper Integrals: (a) Shows an exploding rational function, $\frac{x^8+2}{x+10}$, times a decaying exponential, e^{-x} . The area under the function is finite. (b) Shows a decaying rational function, $\frac{1}{x^4+1}$, times an exploding exponential, e^x . The area under the function is infinite.

Proposition 8.8: Families of Convergent and Divergent Improper Integrals: Take 1

For all real numbers $\alpha > 0$ and integers $k \geq 0$, the following hold:

(a) $\int_0^\infty x^k \cdot e^{-\alpha x} dx$ is convergent.

(b) $\int_0^\infty \frac{1}{x^k} \cdot e^{\alpha x} dx$ is divergent.

Proof: Part (a): For $k \geq 0$, define $\gamma(k) := \int_0^\infty x^k \cdot e^{-\alpha x} dx$. We prove by induction that $\gamma(k) < \infty$ for all $k \geq 0$. The base case is easy, because $\gamma(0) = \int_0^\infty e^{-\alpha x} dx = \frac{1}{\alpha}$. Next, we assume the result is proven for all some $k \geq 0$ and seek to show it also holds for $k + 1$. To

do so, we apply integration by parts to $x^{k+1} \cdot e^{-\alpha x}$ with $u := x^{k+1}$ and $dv := e^{-\alpha x} dx$. Then, for any $0 < a < \infty$,

$$\begin{aligned}\int_0^a x^{k+1} \cdot e^{-\alpha x} dx &= x^{k+1} \cdot \left(\frac{-1}{\alpha} e^{-\alpha x} \right) \Big|_0^a - \int_0^a \left(\frac{-1}{\alpha} e^{-\alpha x} \right) \cdot (kx^k) dx \\ &= x^{k+1} \cdot \left(\frac{-1}{\alpha} e^{-\alpha x} \right) \Big|_0^a + \frac{k}{\alpha} \cdot \int_0^a x^k \cdot e^{-\alpha x} dx.\end{aligned}$$

Taking limits as $a \rightarrow \infty$, we obtain,

$$\gamma(k+1) = \frac{1}{\alpha} + \frac{k}{\alpha} \cdot \gamma(k),$$

and thus, if $\gamma(k)$ is finite, so is $\gamma(k+1)$.

Part (b): From Prop. 2.28, $\lim_{x \rightarrow \infty} \frac{1}{x^k} \cdot e^{\alpha x} = \infty$. Hence, there exists $M < \infty$ such that $\frac{1}{x^k} \cdot e^{\alpha x} > 1$, for all $x > M$. It follows that the integral is unbounded. ■.

Next, we recall our “trick” from Chapter 2.6.1, namely

$$\frac{b_m x^m + b_{m-1} x^{m-1} + \dots + b_0}{a_n x^n + a_{n-1} x^{n-1} + \dots + a_0} = \boxed{\frac{x^m}{x^n}} \cdot \frac{b_m + b_{m-1} \frac{1}{x} + \dots + b_0 \frac{1}{x^m}}{a_n + a_{n-1} \frac{1}{x} + \dots + a_0 \frac{1}{x^n}}.$$

Because

$$\lim_{x \rightarrow \infty} \frac{b_m + b_{m-1} \frac{1}{x} + \dots + b_0 \frac{1}{x^m}}{a_n + a_{n-1} \frac{1}{x} + \dots + a_0 \frac{1}{x^n}} = \frac{b_m}{a_n},$$

we can immediately extend Prop. 8.8 to rational functions. The proof is left to the learner.

Proposition 8.9: Families of Convergent and Divergent Improper Integrals: Take 2

Consider a rational function $f(x) = \frac{b_m x^m + b_{m-1} x^{m-1} + \dots + b_0}{a_n x^n + a_{n-1} x^{n-1} + \dots + a_0}$ and assume it is not the zero function. If the denominator $a_n x^n + a_{n-1} x^{n-1} + \dots + a_0$ has any real roots, choose $c \geq 0$ strictly greater than the largest root; otherwise, set $c = 0$. Then the function $f : [c, \infty) \rightarrow \mathbb{R}$ is continuous and the following hold:

(a) For all $\alpha > 0$, $\int_c^\infty e^{-\alpha x} \cdot \frac{b_m x^m + b_{m-1} x^{m-1} + \dots + b_0}{a_n x^n + a_{n-1} x^{n-1} + \dots + a_0} dx$ is convergent.

(b) For all $\beta > 0$, $\int_c^\infty e^{\beta x} \cdot \frac{b_m x^m + b_{m-1} x^{m-1} + \dots + b_0}{a_n x^n + a_{n-1} x^{n-1} + \dots + a_0} dx$ is divergent.

(c) When there is no exponential term,

$$\int_c^\infty \frac{b_m x^m + b_{m-1} x^{m-1} + \dots + b_0}{a_n x^n + a_{n-1} x^{n-1} + \dots + a_0} dx \quad \begin{cases} \text{is convergent} & n \geq m+2 \\ \text{is divergent} & \text{otherwise.} \end{cases}$$

Note: When there is no exponential term present, the comparison function $\frac{1}{x^p}$ applies, as we essentially demonstrated in the solution to Example 8.7-(d). When the degree of the denominator is at least two greater than the degree of the numerator, the integrand behaves as $\frac{1}{x^p}$ for $p \geq 2$, which is convergent; otherwise, the integrand behaves as $\frac{1}{x^p}$ for $p \leq 1$, which is divergent.

When an exponential term is present, the numerator and denominator degrees can be anything; for example, $\int_1^\infty x^{1000} \cdot e^{-0.001x} dx$ is convergent and $\int_1^\infty \frac{1}{x^{1000}} \cdot e^{0.001x} dx$ is divergent. If you try these cases with a numerical integrator, your results may vary due to numerical inaccuracy. The theory is useful.

Example 8.10. Use Prop. 8.9 to determine which of the following improper integrals are convergent or divergent.

(a) $\int_0^\infty e^{-x} \cdot \frac{x^{24}+2}{x^2+1} dx.$

$$(b) \int_0^\infty e^x \cdot \frac{x^{1+1}}{x^{24}+1} dx.$$

$$(c) \int_0^\infty e^{-x} \cdot \frac{x^{24}+2}{x^2-1} dx$$

$$(d) \int_0^\infty \frac{x^{1+1}}{x^{24}+1} dx.$$

$$(e) \int_0^\infty \frac{x^{24}+2}{x^2+1} dx.$$

$$(f) \int_0^\infty \frac{x^{24}+x^7-\pi x^4+2}{x^{26}+1} dx.$$

$$(g) \int_0^\infty \frac{x^{25}+x^7-\pi x^4+2}{x^{26}+1} dx.$$

$$(h) \int_0^\infty \frac{x^{25}+x^7-\pi x^4+2}{(x^{24}+1)\cdot(x^2-1)} dx.$$

Solutions:

(a) $\int_0^\infty e^{-x} \cdot \frac{x^{24}+2}{x^2+1} dx$ is convergent due to the decaying exponential out front and $x^2 + 1$ has no real roots.

(b) $\int_0^\infty e^x \cdot \frac{x^{1+1}}{x^{24}+1} dx$ is divergent due to the exploding exponential out front.

(c) $\int_0^\infty e^{-x} \cdot \frac{x^{24}+2}{x^2-1} dx$ **cannot be addressed** with Prop. 8.9 because the denominator, $x^2 - 1 = (x+1) \cdot (x-1)$ has a real root $r_1 = 1$ within the domain of integration, $[0, \infty)$. If we changed the integral to $\int_c^\infty e^{-x} \cdot \frac{x^{24}+2}{x^2-1} dx$, for $c > 1$, then it would be convergent. **When applying a Theorem, Proposition, Corollary, etc., it is crucial to check the conditions under which it holds.**

(d) $\int_0^\infty \frac{x^{1+1}}{x^{24}+1} dx$ is convergent because (i) $x^{24} + 1$ has no real roots, and (ii) the degree of the denominator exceeds the degree of the numerator by 23, which is greater or equal to two.

(e) $\int_0^\infty \frac{x^{24}+2}{x^2+1} dx$ is divergent because (i) $x^2 + 1$ has no real roots, and (ii) the degree of the numerator exceeds the degree of the denominator.

(f) $\int_0^\infty \frac{x^{24}+x^7-\pi x^4+2}{x^{26}+1} dx$ is convergent because (i) $x^{26} + 1$ has no real roots, and (ii) the degree of the denominator exceeds the degree of the numerator by two.

(g) $\int_0^\infty \frac{x^{25}+x^7-\pi x^4+2}{x^{26}+1} dx$ is divergent because (i) $x^{26} + 1$ has no real roots, and (ii) the degree of the denominator DOES NOT EXCEED the degree of the numerator by two (it only exceeds it by one).

(h) $\int_0^\infty \frac{x^{25}+x^7-\pi x^4+2}{(x^{24}+1)\cdot(x^2-1)} dx$ **cannot be addressed** with Prop. 8.9 because the denominator, $(x^{24}+1)\cdot(x^2-1) = (x^{24}+1)\cdot(x+1)\cdot(x-1)$ has a real root $r_1 = 1$ within the domain of integration, $[0, \infty)$. If we changed the integral to $\int_c^\infty \frac{x^{25}+x^7-\pi x^4+2}{(x^{24}+1)\cdot(x^2-1)} dx$, for $c > 1$, then it would be divergent because the degree of the numerator is only one less than the degree of the denominator. **When applying a Theorem, Proposition, Corollary, etc., it is crucial to check the conditions under which it holds.**

Recommended Videos for Additional Practice

- Improper Integrals - Convergence and Divergence - Calculus 2 by the Organic Chemistry Tutor.
- Type 1 improper integrals, calculus 2 by BlackPenRedPen.
- Improper Integrals - Convergent or Divergent (Made Easy) by Quoc Dat Phung.
- Comparison test for convergence and divergence of improper integrals by BlackPenRedPen.
- Comparison Test for Improper Integrals by Dr. Trefor Bazett.

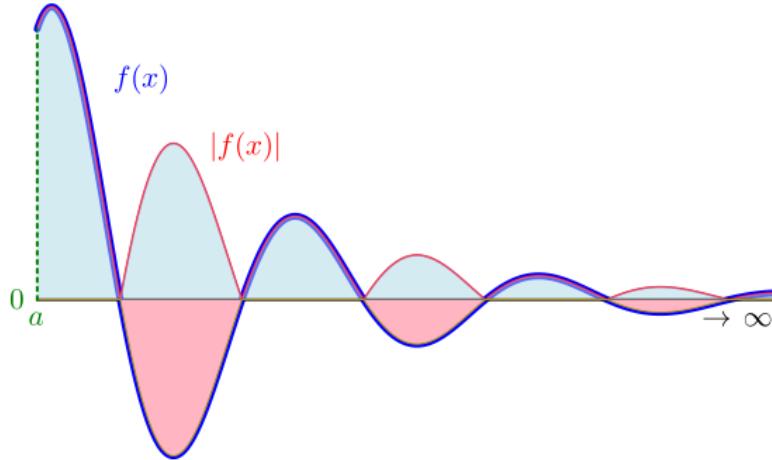


Figure 8.3: **Absolute Integrability:** Suppose that $f : [a, \infty) \rightarrow \mathbb{R}$ is a piecewise continuous function so that its Riemann Integral exists and is finite for any bounded interval $[a, b]$. We observe that the area under $|f|$ is the sum of the area where $f \geq 0$ and THE NEGATIVE OF the area where $f \leq 0$. Hence, if the area under $|f|$ is finite, then each of the positive (shaded blue) and negative (shaded pink) areas of f must be finite, and therefore, the total area under f is well-defined (aka, does not involve $\infty - \infty$) and finite.

8.2.2 Absolute Integrability

Consider a function $f : [a, \infty) \rightarrow \mathbb{R}$ that takes on both positive and negative values. As illustrated in Fig. 8.3, an easy way to make sense of its improper integral is to study the function's absolute value. Because the absolute value is non-negative, the comparison test can be applied.

Absolute Integrability

Definition 8.11. A piecewise continuous function $f : [a, \infty) \rightarrow \mathbb{R}$ is **absolutely integrable** if

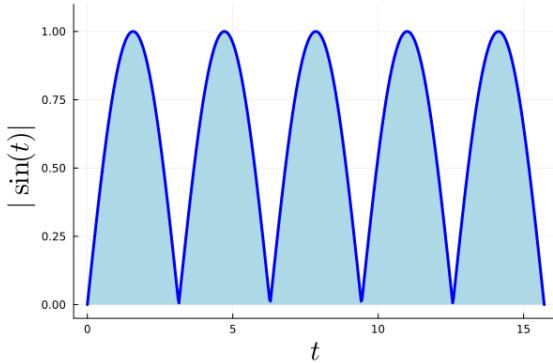
$$\int_a^\infty |f(x)| dx < \infty. \quad (8.4)$$

Proposition 8.12: Usefulness of Absolute Integrability

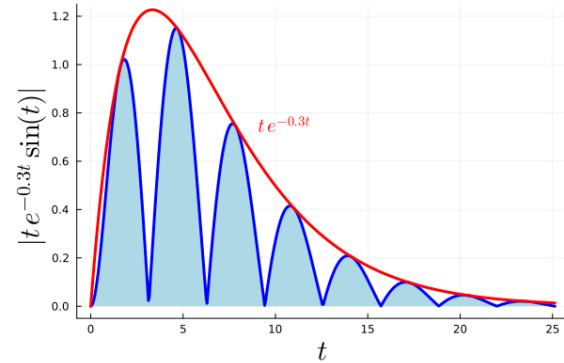
Suppose that $f : [a, \infty) \rightarrow \mathbb{R}$ is piecewise continuous. Then the following hold:

- (a) If f is absolutely integrable, then $\int_a^\infty f(t) dt$ (without absolute value signs) exists and is finite.
- (b) If $\lim_{t \rightarrow \infty} f(t)$ exists and is nonzero, then f is not absolutely integrable.
- (c) If f is differentiable on (a, ∞) and f' (its derivative) is absolutely integrable, then $\lim_{t \rightarrow \infty} f(t)$ exists and is finite.
- (d) If both f and f' (its derivative) are absolutely integrable on (a, ∞) , then $\lim_{t \rightarrow \infty} f(t) = 0$.

Notes: Each of these properties has its uses: (a) allows the comparison test to be used on functions that take on both positive and negative values; (b) provides a sufficient condition for a function NOT TO BE absolutely integrable; (c) provides a sufficient condition for a function to have a limit as t goes to infinity; and (d) provides an extra condition for that limit to be zero.



(a)



(b)

Figure 8.4: Common Functions and their Absolute Integrability: (a) shows that $\sin(t)$ is not absolutely integrable over an unbounded interval, while (b) shows that a negative exponential damps out $t \sin(t)$ so that it becomes absolutely integrable; moreover, (b) also illustrates by the red line a decaying exponential eventually overtaking an exploding monomial. The same holds for higher powers of t .

Example 8.13. Functions of the form $t^k e^{at} \sin(\omega t)$ and $t^k e^{at} \cos(\omega t)$ show up very often in engineering and STEM fields. Show that for all integers $k \geq 0$ and $\omega \in (0, \infty)$ that these functions are absolutely integrable on $[0, \infty)$ if, and only if, $a < 0$.

Solution: The proofs for sine and cosine are identical, so we only give one of them. For $\sin(\omega t)$, it is important to exclude $\omega = 0$, while for $\cos(\omega t)$, it is not. The essence of the proofs can be deduced from Fig. 8.4. If you understand the figure, you are good (the actual proof can then be an optional read); otherwise, read on.

Part 1: $a < 0 \implies \int_0^\infty |t^k e^{at} \sin(\omega t)| dt < \infty$. We break the proof into discrete steps to make it easier to follow.

- $|t^k e^{at} \sin(\omega t)| \leq t^k e^{at}$ for all $t \geq 0$ because $-1 \leq \sin(\omega t) \leq 1$ and $t^k e^{at}$ is non-negative.
- Define $b := \frac{|a|}{2} > 0$ so that $a + b = \frac{a}{2} < 0$.
- Because $b > 0$, and an exploding exponential dominates any monomial, there exists $0 < T < \infty$ such that $t^k \leq e^{bt}$ for all $t \geq T$.
- Putting these facts together, we arrive at

$$\begin{aligned}
 \int_0^\infty |t^k e^{at} \sin(\omega t)| dt &\leq \int_0^\infty t^k e^{at} dt \\
 &\leq \int_0^T t^k e^{at} dt + \int_T^\infty e^{bt} e^{at} dt \\
 &\leq \int_0^T t^k e^{at} dt + \int_T^\infty e^{\frac{a}{2}t} dt \\
 &\leq \int_0^T t^k e^{at} dt + \int_0^\infty e^{\frac{a}{2}t} dt,
 \end{aligned} \tag{8.5}$$

where the last line is true because integrating a positive function over a larger interval cannot decrease the value of the integral.

The first integral is finite because we have a continuous function over a bounded interval. The second integral evaluates to $-\frac{2}{a} = \frac{2}{|a|}$, which is finite. Hence, the sum of the two integrals is finite, completing the proof.

Part 2: $a \geq 0 \implies \int_0^\infty |t^k e^{at} \sin(\omega t)| dt = \int_0^\infty t^k e^{at} |\sin(\omega t)| dt = \infty$. Again, we break the proof into discrete steps to make it easier to follow.

- Let $T := \frac{2\pi}{\omega} > 0$ be the period of $\sin(\omega T)$.

- $\int_0^\infty t^k e^{at} |\sin(\omega t)| dt := \lim_{N \rightarrow \infty} \int_0^{NT} t^k e^{at} |\sin(\omega t)| dt.$
- $\int_0^{NT} t^k e^{at} |\sin(\omega t)| dt = \sum_{n=0}^{N-1} \int_{nT}^{(n+1)T} t^k e^{at} |\sin(\omega t)| dt.$
- Because $t^k e^{at}$ is non-decreasing¹, for all $n \geq 0$, $\int_{nT}^{(n+1)T} t^k e^{at} |\sin(\omega t)| dt \geq \int_0^T t^k e^{at} |\sin(\omega t)| dt.$
- Hence, $\sum_{n=0}^{N-1} \int_{nT}^{(n+1)T} t^k e^{at} |\sin(\omega t)| dt \geq \sum_{n=0}^{N-1} \int_0^T t^k e^{at} |\sin(\omega t)| dt = N \int_0^T t^k e^{at} |\sin(\omega t)| dt.$

Putting all of this together,

$$\int_0^\infty t^k e^{at} |\sin(\omega t)| dt \geq \lim_{N \rightarrow \infty} N \int_0^T t^k e^{at} |\sin(\omega t)| dt. \quad (8.6)$$

But, $\int_0^T t^k e^{at} |\sin(\omega t)| dt > 0$ and hence, $\int_0^\infty t^k e^{at} |\sin(\omega t)| dt = \infty$.

Part 1 and Part 2 together prove the result. ■.

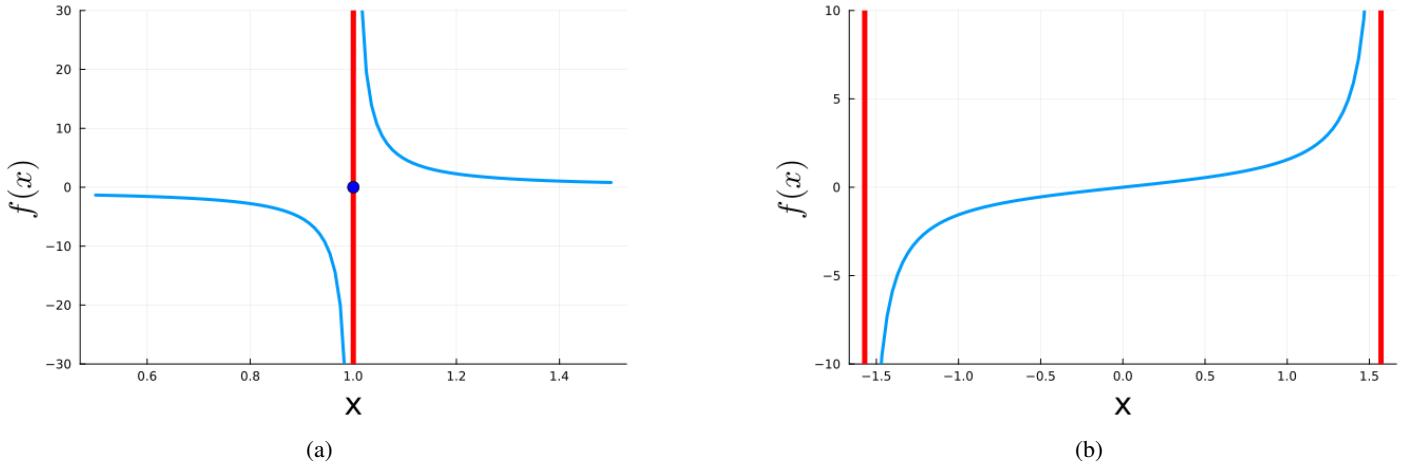


Figure 8.5: **Type-II Improper Integrals:** As shown in (a), the function $\frac{1}{x^2-1}$ has a vertical asymptote at $x = 1$. Can we make sense of the integral $\int_0^1 \frac{1}{x^2-1} dx$ or $\int_0^2 \frac{1}{x^2-1} dx$? Similarly, as shown in (b), the function $\tan(x)$ has vertical asymptotes at $x = \pm \frac{\pi}{2}$. Can we make sense of the integral $\int_0^{\frac{\pi}{2}} \tan(x) dx$?

8.3 Type-II Improper Integrals: Vertical Asymptotes

If a function has a finite vertical asymptote, then it is diverging to $\pm\infty$ at some finite point $x_0 \in \mathbb{R}$, as shown in Fig. 8.5. As you can imagine, the unbounded nature of a function near a finite point demands extra attention when evaluating its integral, akin to situations where the limits of integration are unbounded. This Section is dedicated to exploring such integrals.

Before getting into the nitty-gritty, let's work through a few basic examples.

Consider $f : (0, 2) \rightarrow \mathbb{R}$ by $f(x) = \frac{1}{x^3}$. The function is continuous on the set $(0, 2)$, and for all $0 < a < 2$, $\int_a^2 \frac{1}{x^3} dx$ exists and is finite. Indeed,

$$\int_a^2 \frac{1}{x^3} dx = \int_a^2 x^{-3} dx = \frac{1}{-2} x^{-2} \Big|_a^2 = \frac{1}{2x^2} \Big|_a^2 = \frac{1}{2a^2} - \frac{1}{8}.$$

¹For $k > 0$ or $a > 0$, it is strictly increasing, but for $0 = k = a$, it is constant.

We next take the limit as a approaches zero from the right,

$$\int_0^2 \frac{1}{x^3} dx := \lim_{a \rightarrow 0^+} \int_a^2 \frac{1}{x^3} dx = \lim_{a \rightarrow 0^+} \left(\frac{1}{2a^2} - \frac{1}{8} \right) = \infty,$$

and hence, the integral is **divergent**.

Consider $f : (0, 2) \rightarrow \mathbb{R}$ by $f(x) = \frac{1}{\sqrt{x}}$. The function is continuous on the set $(0, 2)$, and for all $0 < a < 2$, $\int_a^2 \frac{1}{\sqrt{x}} dx$ exists and is finite. Indeed,

$$\int_a^2 \frac{1}{\sqrt{x}} dx = \int_a^2 x^{-\frac{1}{2}} dx = \frac{1}{\frac{1}{2}} x^{\frac{1}{2}} \Big|_a^2 = 2(\sqrt{2} - \sqrt{a}).$$

We next take the limit as a approaches zero from the right,

$$\int_0^2 \frac{1}{\sqrt{x}} dx := \lim_{a \rightarrow 0^+} \int_a^2 \frac{1}{\sqrt{x}} dx = \lim_{a \rightarrow 0^+} 2(\sqrt{2} - \sqrt{a}) = 2 \cdot \sqrt{2},$$

and hence, the integral is **convergent**.

Consider $f : [0, 1] \rightarrow \mathbb{R}$ by $f(x) = \frac{x}{\sqrt[4]{1-x^2}}$, which has a vertical asymptote at $x_0 = 1$. The function is continuous on the set $(0, 1)$, and for all $0 < b < 1$, $\int_0^b \frac{x}{\sqrt[4]{1-x^2}} dx$ exists and is finite. Indeed, we do the change of variable $u = 1 - x^2 \implies du = -2x$, and compute

$$\int_{x=0}^{x=b} \frac{x}{\sqrt[4]{1-x^2}} dx = \int_{u=1}^{u=1-b^2} -\frac{1}{2} \frac{1}{\sqrt[4]{u}} du = \frac{1}{2} \int_{u=1-b^2}^{u=1} u^{-\frac{1}{4}} du = \frac{1}{2} \frac{1}{\frac{3}{4}} u^{\frac{3}{4}} \Big|_{u=1-b^2}^{u=1} = \frac{2}{3} \left(1 - (1-b^2)^{\frac{3}{4}} \right).$$

We next take the limit as b approaches one from the left,

$$\int_0^1 \frac{x}{\sqrt[4]{1-x^2}} dx := \lim_{b \rightarrow 1^-} \int_0^b \frac{x}{\sqrt[4]{1-x^2}} dx = \lim_{b \rightarrow 1^-} \frac{2}{3} \left(1 - (1-b^2)^{\frac{3}{4}} \right) = \frac{2}{3}.$$

and hence, the integral is **convergent**.

From the above, we see that the process is essentially identical to what we did for improper integrals of Type I: we set up the integral for a domain where the integrand is well-behaved (continuous and bounded), and then use a single-sided limit to approach the “problematic point” from the left or the right; for Type-I Integrals, the “problematic point” was $\pm\infty$, while for Type-II Integrals, it is x_0^+ or x_0^- , the right and left sides of a vertical asymptote. We’ll do one more example to illustrate what happens when a vertical asymptote is strictly between the lower and upper bounds of integration, instead of being one of the endpoints of integration.

Consider² $f : [0, 2] \rightarrow \mathbb{R}$ by

$$f(x) = \begin{cases} 0 & x = 1 \\ \frac{x}{\sqrt[4]{|1-x^2|}} & \text{otherwise} \end{cases}$$

where an absolute value was added to define the square root for $x > 1$. The function is well behaved at each of the endpoints (does not blow up), so does the integral, $\int_0^2 f(x) dx$, make sense? Especially when, as is often the case, it is written as

$$\int_0^2 \frac{x}{\sqrt[4]{|1-x^2|}} dx?$$

The answer is: to make sense of it, we need to break the integral into two pieces,

$$\underbrace{\int_0^{1^-} \frac{x}{\sqrt[4]{|1-x^2|}} dx}_{\text{left of the asymptote}} + \underbrace{\int_{1^+}^2 \frac{x}{\sqrt[4]{|1-x^2|}} dx}_{\text{right of the asymptote}},$$

and then analyze each piece separately, just as we did with an integral having unbounded lower and upper limits. Then, as long as we can avoid the dreaded “ $\infty - \infty$ ”, we can sum the two answers and arrive at a value for the overall integral. If you think of $+\infty$ as approaching infinity from the left and $-\infty$ as approaching infinity from the right, then dealing with a Type-I Improper Integral with doubly infinite limits and a Type-II Improper Integral with an asymptote between the limits of integration, are really one and the same.

²Often, the person posing the problem will not be so careful as to define the function at the vertical asymptote, but they should be that careful.

Fact 8.14: Universal Approach for all Improper Integrals

To address an improper integral, the universal approach involves initially setting up a definite integral,

$$\int_a^b f(x) dx,$$

where $f : [a, b] \rightarrow \mathbb{R}$ is piecewise continuous and bounded. Subsequently, evaluate the limit as a or b approaches either $\pm\infty$ or converges to the left or right of a vertical asymptote. This process is done one limit at a time if both bounds of integration involve limits or if the original integral involves a vertical asymptote occurring at $a < x_0 < b$ (strictly between the limits).

If you keep this process in mind, then there is no need to worry about Type-I vs Type-II Improper Integrals. In each case, taking a limit is your “get out of jail free card”.

Example 8.15. Evaluate the following integrals.

(a) $\int_0^1 \frac{1}{x} dx.$

(b) $\int_0^1 \frac{1}{\sqrt{x}} dx.$

(c) $\int_0^2 \frac{1}{\sqrt[4]{|x-1|}} dx.$

(d) $\int_0^{\frac{\pi}{2}} \tan(x) dx.$

(e) $\int_{-1}^1 \frac{x}{\sqrt{1-x^2}} dx.$

Solutions:

(a) **Ans.** $\int_0^1 \frac{1}{x} dx = \infty.$

There is a vertical asymptote at $x_0 = 0$. Hence, $\int_0^1 \frac{1}{x} dx := \lim_{a \rightarrow 0^+} \int_a^1 \frac{1}{x} dx = \lim_{a \rightarrow 0^+} \ln(x) \Big|_a^1 = \lim_{a \rightarrow 0^+} (\ln(1) - \ln(a)) = \infty.$

(b) **Ans.** $\int_0^1 \frac{1}{\sqrt{x}} dx = 1.$

There is a vertical asymptote at $x_0 = 0$. Hence, $\int_0^1 \frac{1}{\sqrt{x}} dx := \lim_{a \rightarrow 0^+} \int_a^1 x^{-\frac{1}{2}} dx = \lim_{a \rightarrow 0^+} \frac{1}{2} \sqrt{x} \Big|_a^1 = \lim_{a \rightarrow 0^+} (1 - \sqrt{a}) = 1.$

(c) **Ans.** $\int_0^2 \frac{1}{\sqrt[4]{|x-1|}} dx = \frac{8}{3}.$

There is a vertical asymptote at $x_0 = 1$. Hence, $\int_0^2 \frac{1}{\sqrt[4]{|x-1|}} dx := \int_0^{1^-} \frac{1}{\sqrt[4]{|x-1|}} dx + \int_{1^+}^2 \frac{1}{\sqrt[4]{|x-1|}} dx$. Taking them one at a time gives

$$\begin{aligned}
\int_0^{1^-} \frac{1}{\sqrt[4]{|x-1|}} dx &= \lim_{b \rightarrow 1^-} \int_0^b \frac{1}{\sqrt[4]{|x-1|}} dx \\
&= \lim_{b \rightarrow 1^-} \int_0^b \frac{1}{\sqrt[4]{1-x}} dx \quad (\text{because } x-1 \leq 0) \\
&= \lim_{b \rightarrow 1^-} \int_0^b (1-x)^{-\frac{1}{4}} dx \quad (\text{prepare to apply generalized power rule}) \\
&= \lim_{b \rightarrow 1^-} \frac{1}{\frac{3}{4}} \cdot (-1) \cdot (1-x)^{\frac{3}{4}} \Big|_0^b \quad (-1 \text{ is from the chain rule}) \\
&= \lim_{b \rightarrow 1^-} -\frac{4}{3} \left((1-b)^{\frac{3}{4}} - 1 \right) \\
&= -\frac{4}{3} \left((1-1)^{\frac{3}{4}} - 1 \right) \\
&= \frac{4}{3}
\end{aligned}$$

and, next,

$$\begin{aligned}
\int_{1^+}^2 \frac{1}{\sqrt[4]{|x-1|}} dx &= \lim_{a \rightarrow 1^+} \int_a^2 \frac{1}{\sqrt[4]{|x-1|}} dx \\
&= \lim_{a \rightarrow 1^+} \int_a^2 \frac{1}{\sqrt[4]{x-1}} dx \quad (\text{because } x-1 \geq 0) \\
&= \lim_{a \rightarrow 1^+} \int_a^2 (1-x)^{-\frac{1}{4}} dx \quad (\text{prepare to apply generalized power rule}) \\
&= \lim_{a \rightarrow 1^+} \frac{1}{\frac{3}{4}} \cdot (1-x)^{\frac{3}{4}} \Big|_a^2 \\
&= \lim_{a \rightarrow 1^+} \frac{4}{3} \left((2-1)^{\frac{3}{4}} - (a-1)^{\frac{3}{4}} \right) \\
&= \frac{4}{3} \left((2-1)^{\frac{3}{4}} - (1-1)^{\frac{3}{4}} \right) \\
&= \frac{4}{3}.
\end{aligned}$$

Hence, $\int_0^2 \frac{1}{\sqrt[4]{|x-1|}} dx = \frac{4}{3} + \frac{4}{3} = \frac{8}{3}$.

Note: In the first integral, you can do a u-substitution, via, $u = 1 - x \implies du = -dx$

$$\begin{aligned} \int_{x=0}^{x=b} (1-x)^{-\frac{1}{4}} dx &= \int_{u=1}^{u=1-b} u^{-\frac{1}{4}} (-1) du \\ &= \int_{u=1-b}^{u=1} u^{-\frac{1}{4}} du \quad (\text{swap bounds of integration to remove minus sign}) \\ &= \frac{4}{3} u^{\frac{3}{4}} \Big|_{1-b}^1. \end{aligned}$$

(d) **Ans.** $\int_0^{\frac{\pi}{2}} \tan(x) dx = \infty.$

There is a vertical asymptote at $x_0 = \frac{\pi}{2}$. Before dealing with that, let's obtain an antiderivative for $\tan(x) = \frac{\sin(x)}{\cos(x)}$ via the u-substitution, $u = \cos(x) \implies du = -\sin(x) dx$. Hence,

$$\int \frac{\sin(x)}{\cos(x)} dx = \int -\frac{1}{u} du = -\ln(u) + C = -\ln(\cos(x)) + C.$$

Knowing the antiderivative, we compute

$$\begin{aligned} \int_0^{\frac{\pi}{2}} \tan(x) dx &= \lim_{b \rightarrow \frac{\pi}{2}^-} \int_0^b \tan(x) dx \\ &= \lim_{b \rightarrow \frac{\pi}{2}^-} -\ln(\cos(x)) \Big|_0^b \\ &= \lim_{b \rightarrow \frac{\pi}{2}^-} \ln(\cos(x)) \Big|_b^0 \quad (\text{swap bounds of integration to remove minus sign}) \\ &= \lim_{b \rightarrow \frac{\pi}{2}^-} (\ln(\cos(0)) - \ln(\cos(b))) \\ &= 0 - (-\infty) \\ &= \infty. \end{aligned}$$

Note: If you were to compute $\int_{-\frac{\pi}{2}}^{\frac{\pi}{2}} \tan(x) dx$, you would have two vertical asymptotes, one at $-\frac{\pi}{2}$ and one at $\frac{\pi}{2}$, and you would arrive at the integral being undefined due to $-\infty + \infty$. If you have trouble seeing this, do a quick plot of $\tan : (-\frac{\pi}{2}, \frac{\pi}{2}) \rightarrow \mathbb{R}$.

(e) **Ans.** $\int_{-1}^1 \frac{x}{\sqrt{1-x^2}} dx = 0.$

There are vertical asymptotes at ± 1 . Before dealing with that, let's find an antiderivative for $\frac{x}{\sqrt{1-x^2}}$ via the u-substitution

$u = 1 - x^2 \implies du = -2x dx$. Hence,

$$\begin{aligned} \int \frac{x}{\sqrt{1-x^2}} dx &= \int -\frac{1}{2} \frac{1}{\sqrt{1-x^2}} (-2x) dx \\ &= \int -\frac{1}{2} \frac{1}{\sqrt{u}} du \\ &= -\frac{1}{2} \cdot \frac{1}{\frac{1}{2}} u^{-\frac{1}{2}} + C \\ &= -\sqrt{u} + C \\ &= -\sqrt{1-x^2} + C \end{aligned}$$

Now that we have an antiderivative, $\int_{-1}^1 \frac{x}{\sqrt{1-x^2}} dx = \lim_{a \rightarrow -1^+} \int_a^0 \frac{x}{\sqrt{1-x^2}} dx + \lim_{b \rightarrow 1^-} \int_0^b \frac{x}{\sqrt{1-x^2}} dx$. Taking the integrals one at a time, we have, for the first integral,

$$\begin{aligned} \lim_{a \rightarrow -1^+} \int_a^0 \frac{x}{\sqrt{1-x^2}} dx &= \lim_{a \rightarrow -1^+} -\sqrt{1-x^2} \Big|_a^0 \\ &= \lim_{a \rightarrow -1^+} \sqrt{1-x^2} \Big|_0^a \quad (\text{swap bounds of integration to remove minus sign}) \\ &= \lim_{a \rightarrow -1^+} (\sqrt{1-a^2} - \sqrt{1-0}) \\ &= \frac{1}{4} (\sqrt{1-1} - \sqrt{1}) \\ &= -1. \end{aligned}$$

Evaluating the second integral,

$$\begin{aligned} \lim_{b \rightarrow 1^-} \int_0^b \frac{x}{\sqrt{1-x^2}} dx &= \lim_{b \rightarrow 1^-} -\sqrt{1-x^2} \Big|_0^b \\ &= \lim_{b \rightarrow 1^-} \sqrt{1-x^2} \Big|_b^0 \quad (\text{swap bounds of integration to remove minus sign}) \\ &= \lim_{b \rightarrow 1^-} (\sqrt{1-0} - \sqrt{1-b^2}) \\ &= \lim_{b \rightarrow 1^-} (1 - \sqrt{1-1}) \\ &= 1. \end{aligned}$$

Hence, $\int_{-1}^1 \frac{x}{\sqrt{1-x^2}} dx = \lim_{a \rightarrow -1^+} \int_a^0 \frac{x}{\sqrt{1-x^2}} dx + \lim_{b \rightarrow 1^-} \int_0^b \frac{x}{\sqrt{1-x^2}} dx = -1 + 1 = 0$.

Example 8.16. Use, or attempt to use, quadgk to evaluate the following integrals. The primary goal is to determine whether it detects that there is a “problem” due to its integrand blowing up over the domain of integration or not. A secondary goal would be to estimate the value of the integral using quadgk.

(a) $\int_0^1 \frac{1}{x} dx.$

(b) $\int_{-2}^2 \frac{x}{\sqrt{1-x^2}} dx,$ has two vertical asymptotes within the domain of integration.

Solutions:

a) We apply numerical techniques to $\int_0^1 \frac{1}{x} dx.$

```

1 using QuadGK
2 fa(x) = 1/x
3 sol,   = quadgk(fa, 0, 1)
```

Output

```
DomainError with 3.5601181736115222e-307:
integrand produced Inf in the interval (0.0, 7.120236347223045e-307)
```

That is reassuring: quadgk detected there is an issue. So, we ask next, can we use quadgk to determine if the integral is convergent or divergent?

```

1 using QuadGK
2 k=collect(1:2:10)
3 limitRange = 10.0.^(-k)
4 for a in limitRange
5     sol,   = quadgk(fa, a, 1)
6     display([a sol])
7 end
```

Output

```
1×2 Matrix{Float64}:
 0.1  2.30259
1×2 Matrix{Float64}:
 0.001  6.90776
1×2 Matrix{Float64}:
 1.0e-5  11.5129
1×2 Matrix{Float64}:
 1.0e-7  16.1181
1×2 Matrix{Float64}:
 1.0e-9  20.7233
```

Once again, this is reassuring: as $a \rightarrow 0^+$, the integral is not converging. It would take a lot to show that it is really divergent, but we are getting a hint of it, at least.

b) We apply numerical techniques to $\int_{-2}^2 \frac{x}{\sqrt{1-x^2}} dx,$ which has vertical asymptotes at $\pm 1.$

```

1 using QuadGK
2 fe(x) = x/sqrt(abs(1 - x^2))
3 # There are two vertical asymptotes within the domain of integration
4 sol,   = quadgk(fe, -2, 2)
```

Output

```
DomainError with -1.0:
integrand produced NaN in the interval (-2.0, 0.0)
```

Sweet! It caught the first vertical asymptote at $x_0 = -1$, even though it was not on the boundary. We push our luck and continue.

```
1 using QuadGK
2 fe(x) = x/sqrt(abs(1 - x^2))
3 # There remains one vertical asymptote within the domain of integration
4 sol,   = quadgk(fe, -0.99, 2) # -0.99 approx = -1^+ = -1 + a tiny bit
```

Output

```
(1.8731177232031915, 2.1226263027166235e-8)
```

It does not complain about the vertical asymptote at +1.

```
1 using QuadGK
2 fe(x) = x/sqrt(abs(1 - x^2))
3 # There remains one vertical asymptote within the domain of integration
4 sol,   = quadgk(fe, -1.0+1e-4, 2) # -0.9999 approx = -1^+ = -1 + a tiny bit
```

Output

```
(1.746192377358709, 2.5638760207248433e-8)
```

As we will see later, that is a very good approximation of the integral. QuadGK has integrated through the singularity. You can play with the code and try `sol, = quadgk(fe, -1.0+1e-6, 2)` to see what happens.

We know where the vertical asymptotes are. Let's see how `quadgk` does if we use that information.

$$\int_{-2}^{-1^-} \frac{x}{\sqrt{1-x^2}} dx$$

```
1 # Integrate from -2 to -1
2 using QuadGK
3 k=collect(1:3:15)
4 epsilonRange = 10.0.^(-k)
5 b=-1
6 for epsilon in epsilonRange
7     sol,   = quadgk(fe, -2, b-epsilon)
8     display([b-epsilon sol])
9 end
```

Output

```
1×2 Matrix{Float64}:
-1.1  -1.27379
1×2 Matrix{Float64}:
-1.0001  -1.71791
1×2 Matrix{Float64}:
-1.0  -1.7316
1×2 Matrix{Float64}:
-1.0  -1.73204
1×2 Matrix{Float64}:
-1.0  -1.73205
```

Awesome! We get that $\int_{-2}^{-1^-} \frac{x}{\sqrt{1-x^2}} dx \approx -1.73205$. Next, we are going to cheat a bit and go for $\int_{-1^+}^{1^-} \frac{x}{\sqrt{1-x^2}} dx$. We should break it into two integrals, over $[-1^+, 0]$ and $[0, 1^-]$, but what the heck, let's roll the dice, something a real engineer would never do, right?

```

1 # Integrate from -1 to +
2 using QuadGK
3 k=collect(1:3:15)
4 epsilonRange = 10.0.^(-k)
5 a = -1
6 b = 1
7 for epsilon in epsilonRange
8     sol, = quadgk(fe, a + epsilon, b - epsilon)
9     display([a + epsilon b - epsilon sol])
10 end

```

Output

```

1×3 Matrix{Float64}:
-0.9  0.9  4.66156e-17
1×3 Matrix{Float64}:
-0.9999  0.9999  1.93258e-16
1×3 Matrix{Float64}:
-1.0  1.0  1.19294e-15
1×3 Matrix{Float64}:
-1.0  1.0  6.37856e-16
1×3 Matrix{Float64}:
-1.0  1.0  -1.35308e-16

```

Sweet, we get zero. Finally, we go for $\int_{1^+}^2 \frac{x}{\sqrt{1-x^2}} dx$

```

1 # Integrate from +1 to +2
2 using QuadGK
3 k=collect(1:3:15)
4 epsilonRange = 10.0.^(-k)
5 a = 1
6 b = 2
7 for epsilon in epsilonRange
8     sol, = quadgk(fe, a + epsilon, b - epsilon)
9     display([a + epsilon b - epsilon sol])
10 end

```

Output

```

1×3 Matrix{Float64}:
1.1  1.9  1.15729
1×3 Matrix{Float64}:
1.0001  1.9999  1.71779
1×3 Matrix{Float64}:
1.0  2.0  1.7316
1×3 Matrix{Float64}:
1.0  2.0  1.73204
1×3 Matrix{Float64}:
1.0  2.0  1.73205

```

And we obtain $\int_{1^+}^2 \frac{x}{\sqrt{1-x^2}} dx \approx 1.73205$. When we add up all the pieces, we obtain zero.

Bottom line: QuadGK is a valid option for evaluating improper integrals, though additional care is required. ■

Recommended Videos for Additional Practice

- 8 Examples of Type II Improper Integrals, Calculus 2) by BlackPenRedPen.

- [Type II Improper Integrals - Discontinuous Integrands - Convergence and Divergence](#) by K.O. Math has very nice graphics to help with intuition.
- [Type I vs Type II Improper Integrals \(10 examples\)](#) by **BlackPenRedPen**.
- [Improper Integrals \(6 of 16\) Type 2 Example](#) by Prof. Michel van Biezen.

8.4 Improper Integrals through the Lens of Probability

Improper integrals are common in probability theory, the “science of randomness”.

Secrets of the Arcane 8.17: Understanding Randomness through Exam Scores

Imagine you’re about to take an important exam. You might think that your score will depend solely on how well you’ve studied. However, even if you’re fully prepared, there are countless unforeseen factors that could influence your performance. These factors introduce an element of **randomness** into what your final score might be, where **randomness** refers to the outcome of events that are not certain and the complex ways in which these uncertain events influence outcomes. It’s a concept found everywhere around us, from the roll of dice to the weather each day.

Exam Score Factors: Consider your score on an exam as a random variable. This doesn’t mean your score is entirely left to chance; rather, it’s affected by a combination of predictable elements (like how much you studied) and unpredictable ones (randomness), such as the factors listed below:

- **Did You Sleep Well?** A good night’s sleep is crucial for memory and concentration. Poor sleep could unexpectedly lower your performance.
- **Transportation Troubles:** Being late because the bus was delayed can cause panic and rush, impacting your focus and, consequently, your score.
- **Unexpected News:** Receiving distressing news from a friend right before the exam can unsettle you, or on the flip side, hearing good news (like winning the lottery) might overly excite you, both affecting your concentration.
- **Personal Well-being:** Sudden headaches, feeling unwell, or even the room being too cold or too hot can distract you during the exam.
- **Antsy People Near You:** A nearby student who twiddles their pencil or keeps moving around in their chair can cause a loss of concentration.
- **Order the Questions are Placed on the Exam:** Nailing the first question can give you a sense of calm for proceeding to the next questions, while flailing away on the first question can undermine your confidence.
- **Adequate Time to Study:** Maybe you have a job and had to work overtime, or another course scheduled its exam on the same day with very little notice.

You would agree that most of the factors listed above are not in your direct control and that you cannot predict ahead of time with any degree of certainty whether or not they will occur for you on any given exam. That is what randomness is all about. It’s a concept that stretches far beyond the simple outcome of flipping a coin or rolling dice. Randomness infiltrates every aspect of our lives, influencing events in ways that are often subtle, sometimes profound, and always unavoidable.

The subjects of **Probability** and **Statistics** are the mathematical siblings born out of the need to understand and quantify this randomness. Through these disciplines, humanity has developed a toolkit not just for gauging the likelihood of events but also for making sense of the random variability that surrounds us. You have likely heard the terms **mean**, **variance**, and **standard deviation**, which are critical in this endeavor. Let’s delve a bit deeper into how these quantities help us grapple with randomness:

- The **mean**, often referred to as the average, gives us a central value for a set of data. When considering exam scores across many students, the mean score provides a quick snapshot of the overall performance level.

- **Variance** measures how spread out a set of data is. If we’re looking at the variance in exam scores, we’re essentially asking, “How much do students’ scores differ from each other and from the mean?” High variance means a wide dispersion of scores, indicating that factors of randomness may have had varied impacts on the students.
- The **standard deviation**, the square root of the variance, offers a digestible measure of this spread in the same units as the data itself. It tells us, on average, how far each student’s exam score deviates from the mean score. Understanding standard deviation allows students to appreciate the diversity of outcomes produced by both controllable factors (like study habits) and uncontrollable factors (like those random events we discussed).

Through Probability and Statistics, engineers learn to model, analyze, and interpret the effects of randomness. We use these mathematical methods to study random events, enabling us to predict probabilities, assess risks, and make informed decisions—even in the face of uncertainty. So, while we might not be able to control or predict every factor that affects our lives, these disciplines give us a powerful way to understand and manage the complexity of randomness.

8.4.1 What is a Probability Density?

In formal courses on Probability and Statistics, you’ll explore different ways to model (aka, mathematically represent) random quantities, including probability measures, distributions, and densities. This section highlights **probability densities**—functions $f : \mathbb{R} \rightarrow [0, \infty)$ that are non-negative, sum up (integrate) to one, and where the integral $\int_a^b f(x) dx$ signifies the probability of a random quantity falling within the range $[a, b]$, such as your exam score being between 86% and 93%. Some probability densities are defined on bounded sets, and others on unbounded sets. Those defined on unbounded sets require improper integrals.

8.4.2 The Role of Improper Integrals

Consider a web server where events correspond to users requesting information. The time between one user request and the next is a random variable. The time between these events is often modeled using the **exponential distribution**, a choice motivated by its so-called “memoryless property”: the probability of the next event occurring does not depend on how much time has elapsed since the previous event occurred.

Example 8.18. *The probability density function (PDF) for the exponential distribution is given by*

$$f(x) = \lambda e^{-\lambda x},$$

where $\lambda > 0$ is the rate parameter, and $x \geq 0$ represents the time between events in clock cycles. Given the high clock frequencies of modern processors, the time x can stretch towards infinity, making this an example of an improper integral. For the exponential distribution, do the following:

(a) Show the PDF integrates to one.

(b) Compute the mean value.

(c) Compute the variance.

Solution:

(a) We need to show $\int_0^\infty \lambda e^{-\lambda x} dx = 1$. Let’s compute this integral, considering it as an improper integral due to the infinite upper bound.

$$\begin{aligned} \int_0^\infty \lambda e^{-\lambda x} dx &= \lim_{t \rightarrow \infty} \int_0^t \lambda e^{-\lambda x} dx \\ &= \lim_{t \rightarrow \infty} \left[-e^{-\lambda x} \right]_0^t \\ &= \lim_{t \rightarrow \infty} (-e^{-\lambda t} + e^{-\lambda \cdot 0}) \\ &= \lim_{t \rightarrow \infty} (-e^{-\lambda t} + 1) \\ &= 1. \end{aligned}$$

- (b) In a course on Probability and Statistics, you will learn that the mean value, often denoted by μ , satisfies

$$\mu = \int_0^\infty xf(x) dx.$$

Specializing this formula to the exponential distribution gives

$$\mu = \int_0^\infty x\lambda e^{-\lambda x} dx.$$

To find the mean value μ , we use integration by parts, where $u = x$ and $dv = \lambda e^{-\lambda x} dx$, giving $du = dx$ and $v = -e^{-\lambda x}$. Thus,

$$\begin{aligned}\mu &= -xe^{-\lambda x}\Big|_0^\infty + \int_0^\infty e^{-\lambda x} dx \\ &= \lim_{t \rightarrow \infty} (-te^{-\lambda t}) + \left[-\frac{1}{\lambda}e^{-\lambda x}\right]_0^\infty \\ &= 0 + \left(0 - \left(-\frac{1}{\lambda}\right)\right) \\ &= \frac{1}{\lambda}.\end{aligned}$$

- (c) In a course on Probability and Statistics, you will learn that the variance, often denoted by σ^2 , satisfies

$$\sigma^2 = \int_0^\infty (x - \mu)^2 f(x) dx.$$

Specializing this formula to the exponential distribution gives

$$\sigma^2 = \int_0^\infty (x - \mu)^2 \lambda e^{-\lambda x} dx.$$

To find the variance σ^2 , we first substitute in $\mu = \frac{1}{\lambda}$ and then expand the quadratic, giving,

$$\sigma^2 = \int_0^\infty (x - \frac{1}{\lambda})^2 \lambda e^{-\lambda x} dx = \int_0^\infty x^2 \lambda e^{-\lambda x} dx - 2 \frac{1}{\lambda} \int_0^\infty x \lambda e^{-\lambda x} dx + \left(\frac{1}{\lambda}\right)^2 \int_0^\infty \lambda e^{-\lambda x} dx.$$

The second and third terms have been computed in part (b) and part (a) respectively. For the first term, using integration by parts twice (details left as an exercise), we find:

$$\int_0^\infty x^2 \lambda e^{-\lambda x} dx = \frac{2}{\lambda^2}.$$

Thus,

$$\sigma^2 = \frac{2}{\lambda^2} - 2 \left(\frac{1}{\lambda}\right) \left(\frac{1}{\lambda}\right) + \left(\frac{1}{\lambda}\right)^2 = \frac{1}{\lambda^2}.$$

We focus next on the normal distribution, which is one of the most fundamental distributions in probability and statistics due to the **Central Limit Theorem**. Unlike the exponential distribution example, the integration to show that the normal distribution density integrates to one is complex and typically beyond the scope of a first calculus course. Thus, we'll concentrate on verifying the mean.

Example 8.19. The probability density function (PDF) for the normal distribution is given by:

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2},$$

where μ is (supposed to represent) the mean, σ is (supposed to be) the standard deviation, and $x \in \mathbb{R}$ represents the value of the random variable. Verify that the mean value really is given by μ .

Solution:

The integral to find the mean μ is

$$\mu = \int_{-\infty}^{\infty} x \cdot \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} dx.$$

This is a doubly improper integral. We break it into two parts

$$\begin{aligned} \int_{\mu}^{\infty} x \cdot \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} dx &:= \lim_{b \rightarrow \infty} \int_{\mu}^b x \cdot \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} dx, \\ \int_{-\infty}^{\mu} x \cdot \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} dx &:= \lim_{a \rightarrow -\infty} \int_a^{\mu} x \cdot \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} dx. \end{aligned}$$

Focusing on the first integral, substitute $z = \frac{x-\mu}{\sigma}$, thus $dx = \sigma dz$, and $x = \mu + \sigma z$, yielding

$$\lim_{b \rightarrow \infty} \int_0^{\frac{b-\mu}{\sigma}} (\mu + \sigma z) \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2} \sigma dz.$$

Splitting the integral gives

$$\lim_{b \rightarrow \infty} \mu \int_0^{\frac{b-\mu}{\sigma}} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2} \sigma dz + \lim_{b \rightarrow \infty} \sigma \int_0^{\frac{b-\mu}{\sigma}} z \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2} \sigma dz.$$

To evaluate the first integral, you have to know that the integral of the Gaussian density is one, and half of the Gaussian density gives $\frac{1}{2}$. The second integral can be done by u-substitution and yields one. Hence,

$$\lim_{b \rightarrow \infty} \int_0^{\frac{b-\mu}{\sigma}} (\mu + \sigma z) \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2} \sigma dz = \frac{\mu}{2} + \sigma.$$

The second integral is similar and yields,

$$\lim_{a \rightarrow -\infty} \int_a^{\mu} x \cdot \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} dx = \frac{\mu}{2} - \sigma.$$

Adding them up, completes the verification. ■

Example 8.20. While proving that the Gaussian density really does integrate to one is a bit advanced for a first Calculus course, showing that the integral exists and is bounded is much easier. So that you can focus on the essential elements, prove that $\int_0^{\infty} e^{-x^2} dx$ exists and is finite.

Solution: We use the Comparison Test. For $x \geq 1$, $e^{-x^2} \leq x e^{-x^2}$. Hence,

$$\begin{aligned} \int_1^{\infty} e^{-x^2} dx &\leq \int_1^{\infty} x e^{-x^2} dx \\ &\leq \int_0^{\infty} x e^{-x^2} dx \quad (\text{because the integrand is non-negative}) \\ &\leq \frac{1}{2} \int_0^{\infty} e^{-u} du \quad (\text{u-substitution } u = x^2, du = 2xdx) \\ &\leq \frac{1}{2} e^{-u} \Big|_0^{\infty} \\ &\leq \frac{1}{2}. \end{aligned}$$

Because e^{-x^2} is continuous, $\int_0^1 e^{-x^2} dx$ exists and is finite. Hence, by combining this with the previous bound, we obtain the desired result.

Alternatively, one can apply the Comparison Test like this: for $x \geq 1$, $e^{-x^2} \leq e^{-x}$. We ask the learner to adapt the above proof to this new bound. ■.

Functions of the form $c \frac{1}{x^k}$, $x \geq x_{\min} > 0$ are called **power-law densities** or **Pareto** densities. They occur in Psychology to model forgetting, Climate Science to model energy dissipation in cyclones, in Computer Engineering to model the relationship between a CPU's cache size and the number of cache misses, and in Economics to model relationships varying from the distribution of artists by the average price of their artworks to income distribution in a market economy.

Example 8.21. Assume a power-law density $f(x) = c \frac{1}{x^3}$, $x \in [1, \infty)$. Do the following:

- (a) Find c such that f is a density.
- (b) Compute the mean value.
- (c) Compute the variance.

Solution:

- (a) To find c such that f is a density, we integrate $f(x)$ over its domain and set the integral equal to 1, namely

$$1 = \int_1^\infty c \frac{1}{x^3} dx = c \frac{x^{-2}}{-2} \Big|_1^\infty = \frac{c}{2}.$$

Solving this, we find $c = 2$.

- (b) The mean value μ is computed as,

$$\mu = \int_1^\infty x f(x) dx = \int_1^\infty x \frac{2}{x^3} dx = 2 \frac{x^{-1}}{-1} \Big|_1^\infty = 2.$$

- (c) The variance σ^2 is theoretically obtained via,

$$\sigma^2 = \int_1^\infty (x - \mu)^2 f(x) dx.$$

However, the improper integral is divergent. In a course on Probability, one would say that the infinite variance is due to the power-law's "heavy tail" (meaning, the tail of the density does not drop off to zero rapidly enough to yield a finite integral).

To see qualitatively that the integral is divergent, note that x^2 times $\frac{1}{x^3}$ behaves like $\frac{1}{x}$, which is the line of demarcation for convergent and divergent integrals, as shown in the preliminaries to Chapter 8.2.1. To show divergence analytically, note that

$$(x - \mu)^2 f(x) = (x^2 - 2\mu x + \mu^2) \frac{2}{x^3},$$

and that

$$\int_1^\infty x^2 \frac{2}{x^3} dx = \int_1^\infty \frac{2}{x} dx = 2 \ln(x) \Big|_1^\infty = \infty,$$

while

$$\int_1^\infty (-2\mu x) \frac{2}{x^3} dx = -4 \int_1^\infty x \frac{2}{x^3} dx = -8$$

and

$$\int_1^\infty \mu^2 \frac{2}{x^3} dx = 4 \int_1^\infty x \frac{2}{x^3} dx = 8.$$

We did the last two integrals just to assure you that we did not land in $\infty - \infty$, which would have landed us in the "undefined" category, instead of divergent. If that seems like splitting hairs to you, for most engineering applications, you're probably right! ■.

In Michigan's IOE 265 *Probability and Statistics for Engineers* or EECS 301 *Probabilistic Methods in Engineering*, you will encounter integration challenges similar to the above.

8.5 (Optional Read:) Proofs Associated with the Chapter

Proof of Prop. 8.12 (Usefulness of Absolute Integrability)

Suppose that $f : [a, \infty) \rightarrow \mathbb{R}$ is piecewise continuous.

- (a) If f is absolutely integrable, then $\int_a^\infty f(t) dt$ (without absolute value signs) exists and is finite.
- (b) If $\lim_{t \rightarrow \infty} f(t)$ exists and is nonzero, then f is not absolutely integrable.
- (c) If f is differentiable on (a, ∞) and f' (its derivative) is absolutely integrable, then $\lim_{t \rightarrow \infty} f(t)$ exists and is finite.
- (d) If both f and f' (its derivative) are absolutely integrable on (a, ∞) , then $\lim_{t \rightarrow \infty} f(t) = 0$.

Proof:

- (a) Decompose f into $f^+ := \max\{f, 0\}$ and $f^- := \min\{f, 0\}$. Both f^+ and f^- are piecewise continuous since f is. The absolute integrability of f means $\int_a^\infty |f(t)| dt < \infty$.

We note that by construction, $f = f^+ + f^-$, and $|f| = f^+ - f^-$. Moreover, for all $t \in [a, \infty)$, $0 \leq f^+(t) \leq |f(t)|$ and $0 \leq -f^-(t) \leq |f(t)|$. Hence, by the Comparison Test, both f^+ and $-f^-$ are integrable³ over $[a, \infty)$. Because $-f^-$ is integrable, so is its negative, $-(-f^-) = f^-$. Then, because the sum of two integrable functions over $[a, \infty)$ is also integrable, the proof is complete.

The last step can also be written as, because

$$\lim_{T \rightarrow \infty} \int_a^T f^+(t) dt \text{ exists and is finite, and because}$$

$$\lim_{T \rightarrow \infty} \int_a^T f^-(t) dt \text{ exists and is finite, it follows that there sum,}$$

$$\lim_{T \rightarrow \infty} \int_a^T f^+(t) dt + \lim_{T \rightarrow \infty} \int_a^T f^-(t) dt = \lim_{T \rightarrow \infty} \int_a^T \underbrace{(f^+(t) - f^-(t))}_{f(t)} dt \text{ exists and is finite.}$$

- (b) We assume that $L := \lim_{t \rightarrow \infty} f(t)$ exists and is nonzero. Moreover, we'll assume that $L > 0$ as the proof for $L < 0$ follows the same reasoning. Because the limit exists, for $\epsilon = \frac{L}{2} > 0$, there exists $T < \infty$ such that $|f(t) - L| \leq \frac{L}{2}$ for all $t \geq T$, which in turn implies that $f(t) \geq \frac{L}{2}$. Because a positive constant $\frac{L}{2}$ is not absolutely integrable⁴, by the Comparison Test, neither is $f(t)$.

- (c) By the Fundamental Theorem of Calculus, $f(T) = f(a) + \int_a^T f'(t) dt$. Then, by part (a), because $f'(t)$ is absolutely integrable,

$$\lim_{T \rightarrow \infty} f(T) = \lim_{T \rightarrow \infty} \left(f(a) + \int_a^T f'(t) dt \right) = f(a) + \lim_{T \rightarrow \infty} \int_a^T f'(t) dt$$

exists and is finite.

- (d) By part (b), if we assume in addition that f is absolutely integrable, then it cannot be the case that the limit exists and is nonzero. But by (c), the limit does exist, and hence, it must be zero.

■

³Means, $\lim_{b \rightarrow \infty} \int_a^b f^+(t) dt$ exists and is finite; same for $-f^-$.

⁴ $\lim_{T \rightarrow \infty} \int_a^T \frac{L}{2} dt = \lim_{T \rightarrow \infty} (T - a) \frac{L}{2} = \infty$ when $\frac{L}{2} > 0$.

Chapter 9

Ordinary Differential Equations

Learning Objectives

By completing this chapter, students will:

- Gain an introduction to differential equations, focusing on their significance and the various classifications.
- Learn the foundational theories and techniques for solving one-dimensional first-order ordinary differential equations (ODEs).
- Acquire knowledge on the fundamentals and solution methods for vector first-order ODEs.
- Explore the application of numerical methods in solving vector first-order ODEs.
- Delve into the characteristics of linear systems of first-order ODEs and their implications for mathematical modeling.

Outcomes

Upon mastering the content of this chapter, students will be equipped to:

- Understand the definition and classification of ODEs based on their order, linearity, and whether they are homogeneous or non-homogeneous.
- Engage with various strategies for tackling first-order ODEs, covering separable, linear, and select nonlinear equations.
- Dive into the theory and application of solving systems of linear differential equations, highlighting their importance across diverse applications.
- Comprehend the utilization of ODEs in modeling and solving real-world problems across various domains.
- Grasp the fundamentals of numerical methods designed for solving ODEs in scenarios where analytic solutions are challenging or unattainable.
- Deeply investigate linear systems of ODEs, focusing on:
 - The structure and solution of linear vector ODEs, symbolized as $\dot{x} = Ax, x(t_0) = x_0$.
 - The concept and significance of the matrix exponential, denoted as e^{At} .
 - Essential attributes of the matrix exponential, such as $\frac{d}{dt}e^{At} = Ae^{At} = e^{At}A$.
 - Methods for solving linear vector ODEs, exemplified by $x(t) = e^{A(t-t_0)}x_0$.
 - The role of complex scalar exponentials, $e^{(a+i\omega)t}$, in the context of ODE solutions.
 - The interplay between eigenvectors, eigenvalues, and the matrix exponential, such as $Av = \lambda v \implies e^{A(t-t_0)}v = e^{\lambda(t-t_0)}v$.
 - The criteria for stability in linear systems, $\text{real}(\lambda) < 0 \implies \lim_{t \rightarrow \infty} e^{A(t-t_0)}v = 0_{n \times 1}$, underlining how the real parts of eigenvalues influence the system's long-term behavior.
 - The application of these principles to general initial conditions when A has a complete set of eigenvectors.

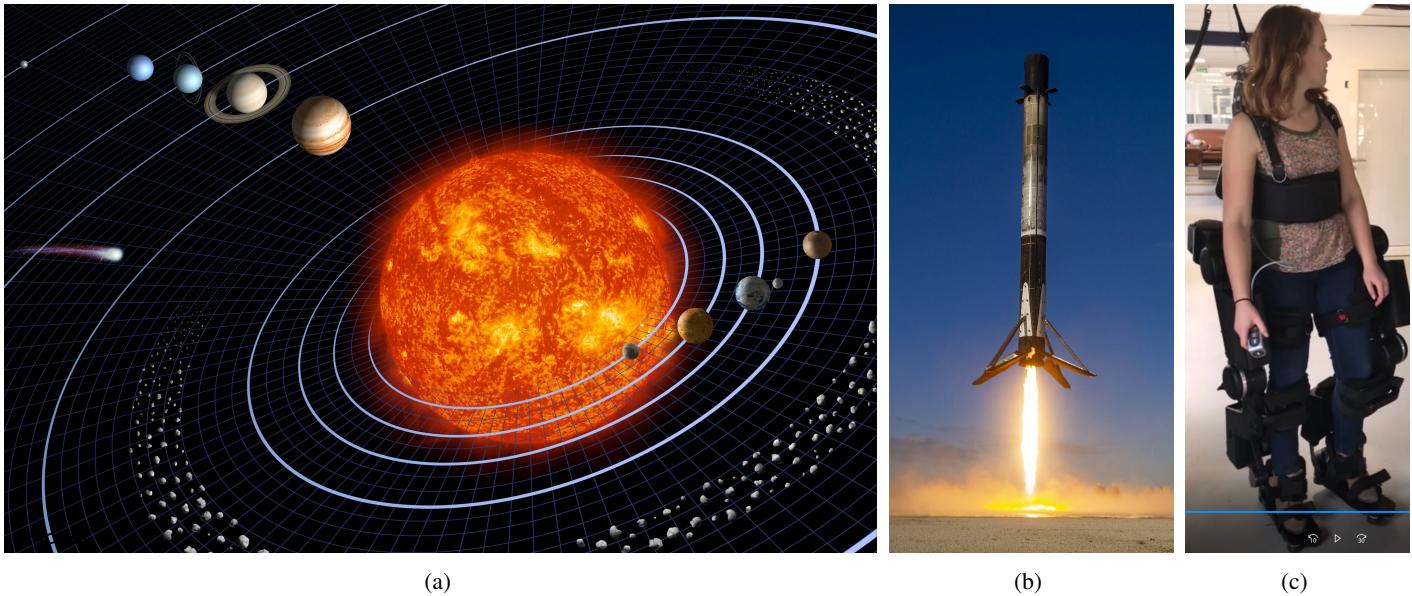


Figure 9.1: **What do we have in common?** (a) Illustration of the solar system. (b) Space-X landing a booster. (c) Wandercraft's exoskeleton allows patients with paraplegia (paralysis from the waist down) to walk again. Sources: (a) [World History Encyclopedia](#) (b) [SpaceX-Booster](#) (c) [TF1 Info: Wandercraft fait remarcher les paraplégiques grâce à un exosquelette](#)

9.1 Introduction

Figure 9.1 illustrates three “objects” where **ordinary differential equations (ODEs)** are required to understand how they function. A differential equation is an equation that includes one or more derivatives, such as

$$\begin{aligned} \dot{x}(t) + x(t) - \sin(t) &= 0, \text{ or} \\ \ddot{x}(t) + \dot{x}(t) + (x(t) + 1)^3 &= 0; \end{aligned}$$

if there were no derivatives, we could simply solve the equations for the unknown $x(t)$ and be done with it. Got it. Then what makes an equation with derivatives “ordinary”? It is an “ordinary” differential equation if the included derivatives are all with respect to the same scalar variable. If there are derivatives with respect to several variables (or a vector of variables), then they are called partial derivatives, and we would have a **partial differential equation**, or PDE for short. We’ve encountered a few ODEs already, in Example 5.52 dealing with the total derivative, and in Chapter 6.4 on Lagrangian Dynamics, where we derived the equations of motion for a point mass in \mathbb{R}^2 , a pendulum, a three-link manipulator, and a three-link bipedal walker.

In the present Chapter, we explore what it means for a function (eventually, a vector of functions) to be a solution to an ODE. Similar to antiderivatives, we’ll see that in some cases, a closed-form solution to an ODE can be found, while in other cases, the solution cannot be expressed in terms of simple functions and we rely on numerical methods for determining a solution.

Secrets of the Arcane 9.1: Understanding the Solar System has been a Major Driver of Scientific Inquiry

Watch: [We Did This](#) by Neil deGrasse Tyson. Differential equations are amazing!

The understanding of the solar system has evolved significantly over the centuries^a. In the beginning, we humans postulated paths or orbits for the planets, such as perfect circles around the earth. Eventually, data showed that model was wrong, and it was replaced by ellipses around the sun. These early models were **descriptive**, meaning they described the observational data but did not attempt to explain the “how” or the “why” of the planets’ orbits.

Newton proposed the first **mechanistic** model, with his famous law of gravitational attraction. It was eventually realized that Newton’s laws were highly accurate for all of the planets except Mercury, which presented a significant puzzle to astronomers for many years. The issue was with its perihelion, the point in its orbit closest to the Sun. Over time, the perihelion of Mercury was observed to advance or shift in a way that couldn’t be explained by Newtonian mechanics. *Mercury’s perihelion advanced by an additional 43 arcseconds per century more than what was predicted by Newton’s laws.*

Mercury's orbital discrepancy was resolved by Albert Einstein's General Theory of Relativity in 1915, which is grounded in an advanced form of Calculus on surfaces. General Relativity describes gravity not as a force, but as a curvature of spacetime caused by mass. According to this theory, the path of a planet is a geodesic (a path of shortest distance, such as a great circle route on the earth) in the curved spacetime caused by the Sun.

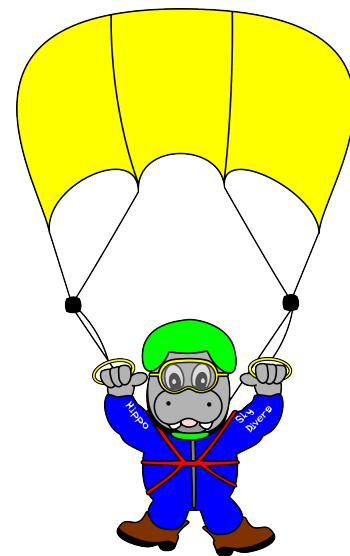
Timeline

1. **Geocentric Model (Ancient Times - 16th Century):** *Ancient Greece (circa 4th century BCE)*: The geocentric model, with Earth at the center of the universe, was developed by Greek philosophers like Aristotle and later refined by Ptolemy (Egypt) in the 2nd century ACE.
2. **Ptolemaic System (2nd century ACE):** *Claudius Ptolemy* detailed the geocentric model in his work "Almagest," introducing epicycles to explain the retrograde motion of planets.
3. **Nicolaus Copernicus (1543):** *Copernicus* proposed the heliocentric model in "De revolutionibus orbium coelestium," suggesting that Earth and other planets revolve around the Sun.
4. **Tycho Brahe (late 16th century):** *Brahe* proposed a hybrid model combining aspects of both geocentric and heliocentric theories.
5. **Kepler's Laws of Planetary Motion (Early 17th Century):** *Johannes Kepler (1609 - 1619)*: Kepler formulated his three laws of planetary motion, describing elliptical orbits of planets around the Sun.
6. **Newton's Law of Universal Gravitation (Late 17th Century):** *Isaac Newton (1687)*: In "Principia Mathematica," Newton formulated the law of universal gravitation, explaining planetary motions and laying the groundwork for Newtonian physics.
7. **Einstein's Theory of General Relativity (Early 20th Century):** *Albert Einstein (1915)*: Einstein's theory of general relativity further refined the understanding of gravitational forces and introduced the concept of spacetime curvature.

^a<https://youtu.be/9xBoC26HBIw>The day the universe changed 05 - James Burke—starts off quite slow, but gets very interesting at 33 minutes, where you see Calculus thinking being employed 100 years before Calculus was formally invented by Newton and Leibniz.



(a)



(b)

Figure 9.2: Using a chute as a brake. (a) A drag chute used for braking the Space Shuttle Endeavor; image complements of Almany and (b) a parachutist; image complements of pixabay.

9.2 Let's Start Simple: One equation, One Unknown, One Derivative

Our mission is to understand how ODEs can be used to model physical phenomena and how “to solve” some of the simplest possible ODEs. In fact, part of our mission is to understand what it even means to be a solution to an ODE.

Fact 9.2: Four Models of Chutes

We use Newton’s Law, $F = ma$, written as $m \frac{dv(t)}{dt} = F$ to model an object braking under air resistance^a while rolling on flat ground and an object falling vertically under the combined influence of gravity and air resistance. In the following, v (m/s) is the velocity of the object, m (kg) is the mass, κ (units depend on the model) is a drag coefficient, and g (m/s²) is the gravitational constant.

- (a) **Model 1: Parallel Motion, Linear Drag:** $m \frac{dv(t)}{dt} = -\kappa v(t)$, $v(t_0) = v_0$.
- (b) **Model 2: Parallel Motion, Quadratic Drag:** $m \frac{dv(t)}{dt} = -\kappa v^2(t)$, $v(t_0) = v_0$.
- (c) **Model 3: Vertical Motion, Linear Drag:** $m \frac{dv(t)}{dt} = -\kappa v(t) - mg$, $v(t_0) = v_0$.
- (d) **Model 4: Vertical Motion, Quadratic Drag:** $m \frac{dv(t)}{dt} = \kappa v^2(t) - mg$, $v(t_0) = v_0$.

The term “ $v(t_0) = v_0$ ” is called the **initial condition** of the differential equation. We understand intuitively that the velocity of a ball at some future time depends on its velocity when it leaves our hand. For example, if you drop a ball at rest versus tossing it in the air at one meter per second, you know that its evolution will be different. Yet, the moment the ball leaves your hand, the physics is the same: the ball evolves under the action of gravity and air resistance.

Regarding the signs in Models 2 and 4, the correct model for drag is $-\kappa|v(t)| \cdot v(t)$, so that the **drag force acts opposite to the direction of motion**. For the space shuttle, we can assume that it is moving left to right at touchdown, and hence $v(t) > 0$. Therefore, the model for drag simplifies to $-\kappa v^2(t)$. For a parachutist, however, $v(t) < 0$, and therefore, the model for drag simplifies to $+\kappa v^2(t)$, where the plus sign is included to emphasize that the drag force is acting “opposite to the downward force of gravity” for the parachutist.

Note: The **order of an ordinary differential equation** is the degree of the highest derivative appearing in the equation. A **homogeneous ordinary differential equation (ODE)** is a type of differential equation in which every term is a function of the dependent variable and its derivatives. In particular, there are no terms that are just constants or functions of the independent variable alone. Hence, in the above, Model 1 is a first-order linear homogeneous ODE and Model 2 is a first-order nonlinear homogeneous ODE. Models 3 and 4 are first-order ODEs, but they are not homogeneous due to the inclusion of gravity in the model. Models that are not homogeneous are **inhomogeneous**, but if you say, **non-homogeneous**, everyone will understand you just fine.

^aThe video [Physics of a Parachute](#) by EA Reach provides an explanation of drag, aka air resistance. The only significant error in the video is the claim that the opening of a parachute “initially causes a parachutist to move upward until gravity gradually takes over.”. While it surely feels that way to a parachutist, our senses often lie to us. Can you explain the feeling? The key is in Fig. 9.4-(b).

Informal Definition of a Solution: A differentiable function of time, $\varphi(t)$, is a solution to the ODE $\frac{dx(t)}{dt} = f(x(t))$, $x(t_0) = x_0$ if (i) it satisfies the initial condition, meaning, $\varphi(t_0) = x_0$, and (ii) it and its derivative together satisfy $\frac{d\varphi(t)}{dt} = f(\varphi(t))$, i.e., the left side equals the right side, as in any equation!

Note: For $\frac{dx(t)}{dt} = f(x(t))$ $x(t_0) = x_0$, it is common practice to denote the solution by $x(t)$, the variable used in writing down the ODE, instead of using a different variable name, such as φ . For the drag chute and parachute models in Fact 9.2, we’ll consequently denote their solutions by $v(t)$.

9.2.1 Analytical Solutions via Antiderivatives

Example 9.3. Compute and plot a solution to the first-order linear homogeneous ordinary differential equation,

$$m \frac{dv(t)}{dt} = -\kappa v(t), \quad v(t_0) = v_0. \quad (9.1)$$

Solution: Ans. $v(t) = v_0 \cdot e^{-\frac{\kappa}{m} \cdot (t-t_0)}$.

First, we write v in place of $v(t)$ and replace dt with $d\tau$. Why? The next steps will be more clear when we use something other than t as the dummy variable of integration. Proceeding,

$$\begin{aligned} \frac{dv}{d\tau} &= -\frac{\kappa}{m} v \\ &\Downarrow \quad (\text{multiply both sides by } \frac{d\tau}{v}) \\ \frac{dv}{d\tau} \cdot \frac{d\tau}{v} &= -\frac{\kappa}{m} v \cdot \frac{d\tau}{v} \\ &\Downarrow \\ \frac{dv}{v} &= -\frac{\kappa}{m} d\tau. \end{aligned}$$

We now have v on one side of the equation and τ on the other side. It is for this reason that the method is called **separation of variables**. We next integrate both sides to obtain

$$\int_{v_0}^{v(t)} \frac{dv}{v} = -\frac{\kappa}{m} \int_{t_0}^t d\tau \quad (\text{this is why we wanted } \tau \text{ instead of } t \text{ as the independent variable}).$$

In the above, t_0 is the time the shuttle Endeavor touches the runway, and $v_0 = v(t_0)$ is the initial velocity of Endeavor at touchdown; often, we take $t_0 = 0$, but we could also use a different time reference, such as the time since the shuttle reentered the atmosphere.

We know antiderivatives for both integrands, giving us

$$\begin{aligned} \ln(v) \Big|_{v_0}^{v(t)} &= -\frac{\kappa}{m} \cdot \tau \Big|_{t_0}^t \\ &\Downarrow \quad (\text{using } \ln(v(t)) - \ln(v_0) = \ln\left(\frac{v(t)}{v_0}\right)) \\ \ln\left(\frac{v(t)}{v_0}\right) &= -\frac{\kappa}{m} \cdot (t - t_0) \\ &\Downarrow \quad (\text{take exponential of both sides}) \\ \frac{v(t)}{v_0} &= e^{-\frac{\kappa}{m} \cdot (t - t_0)} \\ &\Downarrow \quad (\text{multiply both sides by } v_0) \\ v(t) &= v_0 \cdot e^{-\frac{\kappa}{m} \cdot (t - t_0)}. \end{aligned}$$

This gives us the general solution to the ODE as,

$$m \frac{dv(t)}{dt} = -\kappa v(t), \quad v(t_0) = v_0 \iff v(t) = v_0 \cdot e^{-\frac{\kappa}{m} \cdot (t - t_0)}, \quad (9.2)$$

which was the main point of the exercise.

Bonus: Let's do a bit more with it because we are engineers. What should the value of κ be so that in 3,000 m, the speed of the space shuttle is reduced to 10% of its landing speed of 120 m/s (over 400 km/hr), given that it weighs 78,000 kg? The assumption is that friction brakes could take over at that point.

Well, distance is the integral of speed. If we take $t_0 = 0$, we have

$$x(t) = \int_0^T v_0 \cdot e^{-\frac{\kappa}{m} \cdot \tau} d\tau = v_0 \frac{m}{\kappa} \cdot \left(1 - e^{-\frac{\kappa}{m} \cdot T}\right),$$

where T is the time it takes to reduce the velocity to $0.1v_0$. Using our expression for velocity, we compute that T satisfies,

$$0.1v_0 = v_0 \cdot e^{-\frac{\kappa}{m} \cdot T} \iff e^{-\frac{\kappa}{m} \cdot T} = 0.1$$

Substituting this into our formula for distance, we have that

$$3000 = v_0 \frac{m}{\kappa} \cdot \left(1 - e^{-\frac{\kappa}{m} \cdot T}\right) = v_0 \frac{m}{\kappa} \cdot (1 - 0.1) \implies \boxed{\kappa = \frac{0.9 \cdot v_0 \cdot m}{3000} = 2808}.$$

A plot of $v(t)$ vs t is given in Fig. 9.3-(a). The red dot marks the time at which the shuttle has traveled 3 km and reduced its speed to 12 m/s or approx 27 mph. The runway is 4km long. ■

Table 9.4: Everyday Speed Examples in mph and m/s. For a quick conversion, using a factor of two is not bad, while 2.2 and 0.45 get you even closer.

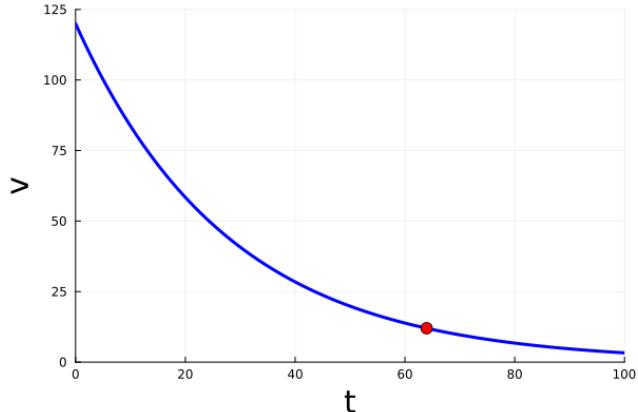
Example	Speed (mph)	Speed (m/s)
Average Human Walking	3.1	1.4
Cyclist (moderate speed)	15	6.7
Car (city driving)	30	13.4
Car (highway speed)	60	26.8
Propeller Airplane	150	67.1
Commercial Jet	550	245.9
High-Speed Train	200	89.4
Cheetah	75	33.5
Sound Speed at Sea Level	767	343

Videos on the Method of Separation of Variables

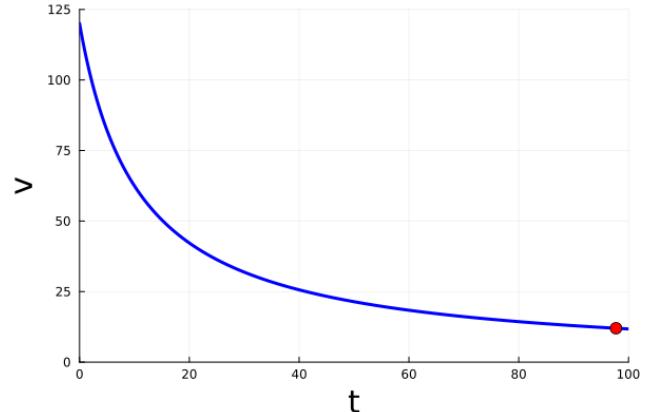
- Separable differential equations introduction: First order differential equations by Khan Academy.
- Separable First Order Differential Equations - Basic Introduction by The Organic Chemistry Tutor.
- Separation of Variables // Differential Equations by Dr. Trefor Bazett.

Example 9.5. Compute and plot a solution to the first-order nonlinear, homogeneous ordinary differential equation,

$$m \frac{dv(t)}{dt} = -\kappa v^2(t), \quad v(t_0) = v_0. \tag{9.3}$$



(a)



(b)

Figure 9.3: Velocity profiles of the Space Shuttle Endeavor braking with a chute in meters per second (m/s). (a) Linear drag model. (b) Nonlinear drag model. The linear model optimistically predicts that Endeavor's speed is reduced to 12 m/s approximately 30 seconds sooner than the (probably more accurate) nonlinear drag model.

Solution: Ans. $v(t) = \frac{v_0}{1 + v_0 \cdot \frac{\kappa}{m} \cdot (t - t_0)}$.

Things will be a bit more clear in the following steps if we write v in place of $v(t)$ and replace dt with $d\tau$. Doing so yields,

$$\begin{aligned} \frac{dv}{d\tau} &= -\frac{\kappa}{m} v^2 \\ \Updownarrow \quad (\text{multiply both sides by } \frac{d\tau}{v}) \quad & \\ \frac{dv}{d\tau} \cdot \frac{d\tau}{v} &= -\frac{\kappa}{m} v^2 \cdot \frac{d\tau}{v} \\ \Updownarrow \quad & \\ \frac{dv}{v^2} &= -\frac{\kappa}{m} d\tau. \end{aligned}$$

Separation of variables: We now have v on one side of the equation and τ on the other side. We next integrate both sides to obtain

$$\int_{v_0}^{v(t)} \frac{dv}{v^2} = -\frac{\kappa}{m} \int_{t_0}^t d\tau.$$

In the above, t_0 is the time of touchdown, and $v_0 = v(t_0)$ is the initial velocity of Endeavor at touchdown. We know antiderivatives

for these two integrands, giving us

$$\begin{aligned}
 -\frac{1}{v} \Big|_{v_0}^{v(t)} &= -\frac{\kappa}{m} \cdot \tau \Big|_{t_0}^t \\
 \Downarrow \\
 -\left(\frac{1}{v(t)} - \frac{1}{v_0}\right) &= -\frac{\kappa}{m} \cdot (t - t_0) \\
 \Downarrow \\
 -\frac{1}{v(t)} &= -\frac{1}{v_0} - \frac{\kappa}{m} \cdot (t - t_0) \\
 \Downarrow \text{(solve for } v(t)) \\
 v(t) &= \frac{1}{\frac{1}{v_0} + \frac{\kappa}{m} \cdot (t - t_0)} \\
 \Downarrow \\
 v(t) &= \frac{v_0}{1 + v_0 \cdot \frac{\kappa}{m} \cdot (t - t_0)}.
 \end{aligned}$$

This gives us the general solution to the ODE as,

$$m \frac{dv(t)}{dt} = -\kappa v^2(t), \quad v(t_0) = v_0 \iff v(t) = \frac{v_0}{1 + v_0 \cdot \frac{\kappa}{m} \cdot (t - t_0)}, \quad (9.4)$$

which was the main point of the exercise.

Bonus: Let's do a bit more with it because we are engineers. With this nonlinear model of drag, what should the value of κ be so that in 3,000 m, the speed of the space shuttle is reduced to 10% of its landing speed of 120 m/s (over 400 km/hr), given that it weighs 78,000 kg? The assumption, again, is that friction brakes could take over at that point.

The time it takes to slow down to $0.1v_0$ is

$$0.1v_0 = \frac{v_0}{1 + v_0 \cdot \frac{\kappa}{m} \cdot T} \implies 10 = 1 + v_0 \cdot \frac{\kappa}{m} \cdot T \implies T = \frac{9m}{\kappa v_0}.$$

This time, integrating velocity to obtain distance will require a u-substitution, $u = 1 + v_0 \cdot \frac{\kappa}{m} \cdot t$, $du = v_0 \cdot \frac{\kappa}{m} dt$, giving,

$$\begin{aligned}
 x(T) &= \int_{t=0}^{t=T} \frac{v_0}{1 + v_0 \cdot \frac{\kappa}{m} \cdot t} dt \\
 &= \int_{u=1}^{u=10} \frac{m}{\kappa} \frac{1}{u} du \\
 &= \frac{m}{\kappa} \cdot \ln(10).
 \end{aligned}$$

Imposing $3000 = \frac{m}{\kappa} \cdot \ln(10) \implies \boxed{\kappa = \frac{m}{3000} \cdot \ln(10) = 59.8}$. A plot of $v(t)$ vs t is given in Fig. 9.3-(b).

Note: The required drag coefficient predicted by the nonlinear model is very different from the value predicted by the linear model, $\kappa = 59.8$ versus $\kappa = 2808$. This is because of the extra factor of v in the nonlinear drag model, $\kappa \cdot v \cdot v$, versus the linear drag model $\kappa \cdot v$. Dividing the linear drag coefficient by 120 gives 23.4, which is in the same ballpark. When you understand that a linear approximation of $\kappa \cdot v \cdot v$ about the landing speed is $2\kappa \cdot v_0 \cdot (v - v_0)$, the values match up even better, $\kappa = 59.8$ versus $\kappa = 46.8$. ■

We now turn to analyzing the models for a parachutist¹, where gravity has an effect. The differential equations are now inhomogeneous. For the linear model, we need a new idea when seeking to compute a solution.

¹A “parachutist” is someone who is trained and experienced in the art of skydiving. They have likely completed many jumps and have received formal training and certification. On the other hand, a “parachuter” is someone who simply jumps out of a plane with a parachute, without necessarily having any formal training or experience.

Fact 9.6: Integrating Factors and Integration by Parts

The concept of an **integrating factor** in solving differential equations and the technique of **integration by parts** in integral calculus are both based on the **product rule of differentiation**.

1. Integration by Parts in Integral Calculus:

- Integration by parts is a technique derived from the product rule of differentiation. It is used to solve integrals involving the product of two functions and is expressed as $\int u \, dv = uv - \int v \, du$.
- It allows you to exchange the task of finding an antiderivative for $\int u \, dv$ for the “hopefully” easier task of finding an antiderivative for $\int v \, du$.
- This formula is a direct consequence of the product rule. If you take the derivative of the product of two functions $u \cdot v$ using the product rule and then integrate both sides, you arrive at the integration by parts formula.

2. Integrating Factor in Differential Equations:

- In the context of first-order linear differential equations, the integrating factor method involves finding a function (the integrating factor) that, when multiplied with the original differential equation, transforms the original ODE into a form where the left side becomes a total differential, as in $d(u \cdot v)$, making it trivial to integrate. The “hope” is that an antiderivative can be found for the right side of the equation.
- For a differential equation of the form $\frac{dy(t)}{dt} + P(t)y(t) = Q(t)$, an integrating factor is $e^{\int_a^t P(\tau) d\tau}$, for a an arbitrary constant, because:
 - Multiplying both sides of the ODE by $e^{\int_a^t P(\tau) d\tau}$ gives

$$\left(e^{\int_a^t P(\tau) d\tau} \right) \cdot \left(\frac{dy(t)}{dt} + P(t)y(t) \right) = \left(e^{\int_a^t P(\tau) d\tau} \right) \cdot Q(t).$$

- By the Chain Rule and the First Fundamental Theorem of Calculus,

$$\frac{d}{dt} \left(e^{\int_a^t P(\tau) d\tau} \right) = \left(e^{\int_a^t P(\tau) d\tau} \right) \cdot \frac{d}{dt} \int_{t_0}^t P(\tau) d\tau = e^{\int_a^t P(\tau) d\tau} \cdot P(t);$$

- By the Product Rule,

$$\begin{aligned} \frac{d}{dt} \left(y(t) \cdot e^{\int_a^t P(\tau) d\tau} \right) &= \frac{dy(t)}{dt} \cdot \left(e^{\int_a^t P(\tau) d\tau} \right) + y(t) \cdot \frac{d}{dt} \left(e^{\int_a^t P(\tau) d\tau} \right) \\ &= \frac{dy(t)}{dt} \cdot e^{\int_a^t P(\tau) d\tau} + y(t) \cdot e^{\int_a^t P(\tau) d\tau} \cdot P(t); \end{aligned}$$

- And hence, $e^{\int_a^t P(\tau) d\tau} \cdot \left(\frac{dy(t)}{dt} + P(t)y(t) \right) = \frac{d}{dt} \left(y(t) \cdot e^{\int_a^t P(\tau) d\tau} \right)$ is a total differential, making it trivial to integrate (aka, anti-differentiate); and
- Then, as in Integration by Parts, we “hope” that we can find an antiderivative for the right side of the equation,

$$Q(t) \cdot e^{\int_a^t P(\tau) d\tau}.$$

Note: When computing the integrating factor, $e^{\int_a^t P(\tau) d\tau}$, you are free to choose a to simplify the integral, $\int_a^t P(\tau) d\tau$. Zero and minus infinity are often convenient values.

Example 9.7. Compute and plot a solution to the first-order linear inhomogeneous ODE.

$$m \frac{dv(t)}{dt} = -\kappa v(t) - mg, \quad v(t_0) = v_0. \tag{9.5}$$

Solution: Ans. $v(t) = v_0 e^{-\frac{\kappa}{m}(t-t_0)} - \frac{mg}{\kappa} \cdot \left(1 - e^{-\frac{\kappa}{m}(t-t_0)} \right)$.

Given the differential equation $m \frac{dv(t)}{dt} = -\kappa v(t) - mg$, we can rewrite it as,

$$\frac{dv(t)}{dt} + \frac{\kappa}{m} v(t) = -g.$$

The solution to this equation can be found using an **integrating factor**, as explained in Fact 9.6. In this case, $P(t) = \frac{\kappa}{m}$, which has antiderivative $\frac{\kappa}{m} \cdot t + a$. Taking $a = 0$ gives the integrating factor as

$$e^{\frac{\kappa}{m}t}.$$

As developed in Fact 9.6, the integrating factor has the “magical property” that when we multiply both sides of the ODE by $e^{\frac{\kappa}{m}t}$, giving us,

$$e^{\frac{\kappa}{m}t} \frac{dv(t)}{dt} + \frac{\kappa}{m} e^{\frac{\kappa}{m}t} v(t) = -ge^{\frac{\kappa}{m}t},$$

the left side of the equation is equal to $\frac{d}{dt}(e^{\frac{\kappa}{m}t}v(t))$. Therefore, integrating both sides with respect to t gives,

$$\int \frac{d}{dt}(e^{\frac{\kappa}{m}t}v(t)) dt = - \int ge^{\frac{\kappa}{m}t} dt.$$

Determining an antiderivative for the left side is now trivial, which was the job of the integrating factor. Determining an antiderivative for the right side is straightforward because it involves an exponential. Integrating and rearranging, we therefore obtain,

$$e^{\frac{\kappa}{m}t}v(t) = -\frac{mg}{\kappa}(e^{\frac{\kappa}{m}t} - 1) + C,$$

where C is the constant of integration. To find C , we use the initial condition $v(t_0) = v_0$, per

$$e^{\frac{\kappa}{m}t_0}v_0 = -\frac{mg}{\kappa}(e^{\frac{\kappa}{m}t_0} - 1) + C.$$

Solving for C and then substituting back into the equation for $v(t)$, yields the final solution,

$$v(t) = v_0 e^{-\frac{\kappa}{m}(t-t_0)} - \frac{mg}{\kappa}(1 - e^{-\frac{\kappa}{m}(t-t_0)}).$$

This gives us the general solution to the ODE as,

$$m \frac{dv(t)}{dt} = mg - \kappa v(t), \quad v(t_0) = v_0 \iff v(t) = v_0 e^{-\frac{\kappa}{m}(t-t_0)} - \frac{mg}{\kappa}(1 - e^{-\frac{\kappa}{m}(t-t_0)}). \quad (9.6)$$

The limit as $t \rightarrow \infty$ is called the “terminal velocity”. No, it’s not the speed that terminates a beginner parachutist. It’s the steady-state speed reached by the parachutist,

$$\lim_{t \rightarrow \infty} v(t) = \lim_{t \rightarrow \infty} v_0 e^{-\frac{\kappa}{m}(t-t_0)} - \frac{mg}{\kappa}(1 - e^{-\frac{\kappa}{m}(t-t_0)}) = -\frac{mg}{\kappa}.$$

From this, we see that if a 75 kg parachutist wants to land with a speed of -2 m/s, they need a parachute that provides a drag coefficient of $\kappa = 367.88$. A plot is given in Fig. 9.4-(a). ■

Example 9.8. Compute and plot a solution to the first-order nonlinear, inhomogeneous ordinary differential equation,

$$m \frac{dv(t)}{dt} = \kappa v^2(t) - mg, \quad v(t_0) = v_0. \quad (9.7)$$

Solution: Ans. A closed-form solution cannot be given in terms of elementary functions, as defined in [Secrets of the Arcane 7.39](#). To prove this fact requires [Differential Galois Theory](#), which few non-mathematicians acquire in their mathematical studies. Even though a closed-form solution cannot be found, we can still compute numerical solutions, as shown in Fig. 9.4-(b). Moreover, qualitative analysis of the equation can be done as in the video [Parachute Physics](#) by LabRat Scientific.

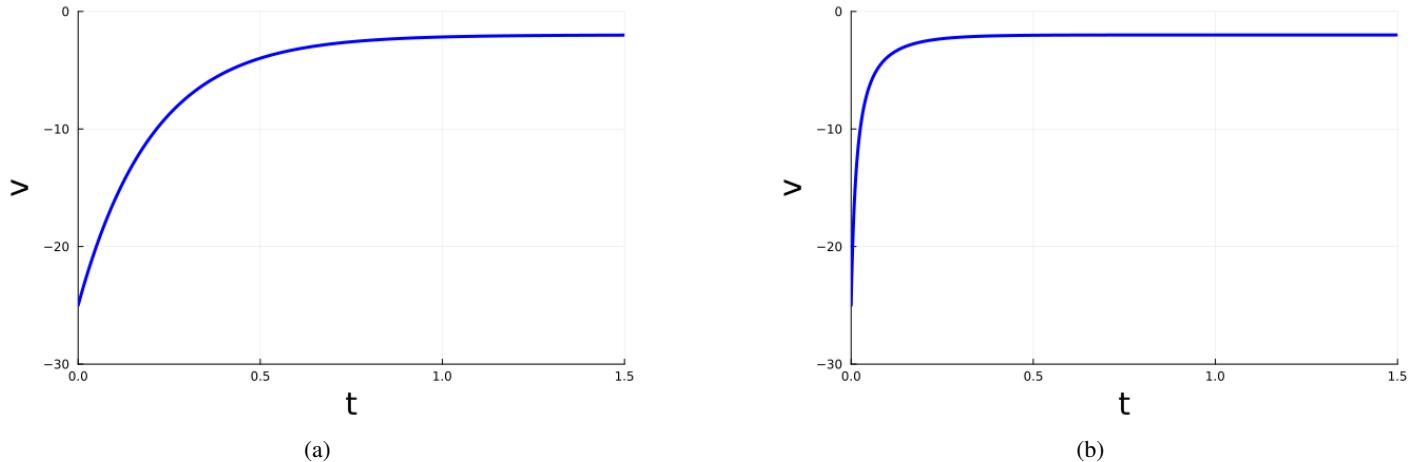


Figure 9.4: Velocity profiles of parachutists, with an emphasis on when each model “achieves a terminal velocity”. (a) Linear drag model. (b) Nonlinear drag model. The linear model predicts a more gentle deceleration to the terminal velocity. In practice, a chute does not open instantaneously, so the braking effect of a real parachute (with nonlinear drag) may not be as harsh as we are predicting!

So yes, we can still compute a terminal velocity, v_{ss} . At the terminal velocity, $\frac{dv(t)}{dt} \equiv 0$. Using this as a definition, we have

$$\frac{dv(t)}{dt} \equiv 0 \implies \kappa v_{ss}^2 = mg \implies v_{ss} = \sqrt{\frac{mg}{\kappa}}.$$

We can also turn it around to “size the parachute” because larger parachutes provide more air resistance. For a parachutist to land with a velocity of v_{ss} m/s implies

$$\kappa = \frac{mg}{v_{ss}^2}.$$

For a parachutist with a mass of 75 kg and desiring a terminal velocity of -2 m/s, $\kappa = 183.94$. A plot is given in Fig. 9.4-(b).

It's remarkable that, just by adding a constant term to the right side of an ODE with one derivative, one unknown, and one equation, we went from something solvable by an elementary method, namely separation of variables, to an ODE where it is provably impossible to write down an analytical solution. Now, go back and look at the ODEs we derived via Lagrange's method and ask yourself, what's the chance we can find an analytical solution to those equations? Ah, what is one in a googol (one in 10^{100})?

9.2.2 Numerical Solutions

Solving ODEs in Julia Using DifferentialEquations.jl

Among the many powerful tools in Julia for solving Ordinary Differential Equations (ODEs) is the `DifferentialEquations.jl` package. This package provides a comprehensive suite of solvers and utilities for dealing with various types of differential equations. To begin, you need to define the differential equation in a function. For example, consider an ODE of the form $\frac{dy}{dt} = f(y, t)$, where f is some function of y and t . In Julia, you would define this as function `f(y, t)`. Next, you set the initial condition and the time span for which you want to solve the ODE. This is done by creating an `ODEProblem` object with your function, initial condition, and time span. Finally, you solve the ODE using the `solve` function from the package. This function takes the problem object and optionally a solver algorithm and other parameters. The result is a solution object that can be analyzed or plotted to understand the behavior of the system described by the ODE. For more information and examples, visit the [DifferentialEquations.jl documentation](#).

Endeavor with a Linear Drag Model

```
# Import useful packages
```

```

2 using DifferentialEquations
3 using Plots
4 gr()
5
6 # Space Shuttle Endeavor with Linear Drag Model
7 struct Params
8     m::Float64      # Mass in kilograms
9     kappa::Float64   # Drag coefficient
10    g::Float64       # Acceleration due to gravity in m/s^2
11 end
12
13 # Initialize the Params structure with specific values
14 m = 78e3 # kg 78,000 kg or 39 metric tons
15 g = 9.81 # m/s^2
16 v0 = 120 # m/s
17 kappa = (0.9*v0*m)/3000; @show kappa
18 #
19 params = Params(m, kappa, g)
20
21 # define the ODE dv/dt = f(v, params, t)
22 function f(v, params, t)
23     dv = -(params.kappa / params.m) * v
24     return dv
25 end
26
27 # Set the initial condition as a vector!
28 v0 = 120 # m/s
29 v0= [v0]
30
31 # Set the time interval
32 T = (0.0, 100) # 100 second simulation period?
33
34 # Setup the ODE problem with out-of-place function
35 problem = ODEProblem{false}(f, v0, T, params)
36
37 # solve the ODE problem using the Runge Kutta Tsitouras 5/4 Integrator
38 sol = solve(problem, Tsit5())
39
40 # plot the solution
41 p1 = plot(sol, lw=3, guidefont = 20, xlabel="t", ylabel="v", legend=false, color=:blue)
42 ylims!(p1, (0, 125))
43
44 #  $e^{-\frac{\kappa}{m} \cdot t_f} = 0.1 \iff t_f = -\frac{m}{\kappa} \log(0.1)$ 
45 Tf = -(m/kappa)*log(0.1); @show Tf
46 p1 = scatter!(p1, [Tf], [0.1*v0], color=:red, markersize=6)

```

Output

```

kappa = 2808
Tf = 63.960697027612376

```

Endeavor with a Nonlinear Drag Model

```

1 # Import useful packages
2 using DifferentialEquations
3 using Plots

```

```

4 gr()
5
6 # Space Shuttle Endeavor with Nonlinear Drag Model
7 struct Params
8     m::Float64      # Mass in kilograms
9     kappa::Float64   # Drag coefficient
10    g::Float64       # Acceleration due to gravity in m/s^2
11 end
12
13 # Initialize the Params structure with specific values
14 m = 78e3 # 78,000 kg
15 g = 9.81
16 kappa = (m/3000)*log(10); @show kappa
17 params = Params(m, kappa, g)
18
19 # define the ODE dv/dt = f(v,params,t)
20 function f(v, params, t)
21     dv = -(params.kappa / params.m) * v.^2
22     return dv
23 end
24
25 # Set the initial condition as a vector!
26 v0 = 120 # m/s
27 v0= [v0]
28
29 # Set the time interval
30 T = (0.0, 100)
31
32 # Setup the ODE problem with out-of-place function
33 problem = ODEProblem{false}(f, v0, T, params)
34
35 # solve the ODE problem using the Runge Kutta Tsitouras 5/4 Integrator
36 sol = solve(problem, Tsit5())
37
38 # plot the solution
39 p1 = plot(sol, lw=3, guidefont = 20, xlabel="t", ylabel="v", legend=false, color=:blue)
40
41 #  $e^{-\frac{\kappa}{m} t_f} = 0.1 \iff t_f = -\frac{m}{\kappa} \log(0.1)$ 
42 Tf = 9*m/(kappa*v0[1]); @show Tf
43 p1 = scatter!(p1, [Tf], [0.1*v0], color=:red, markersize=6)
44 ylims!(p1, (0, 125))

```

Output

```

kappa = 59.867212417845195
Tf = 97.71625842823164

```

Parachutist with a Linear Drag Model

```

1 # Import useful packages
2 using DifferentialEquations
3 using Plots
4 gr()
5
6 # Parachutist with a linear drag model
7 struct Params

```

```

8   m::Float64      # Mass in kilograms
9   kappa::Float64  # Drag coefficient
10  g::Float64      # Acceleration due to gravity in m/s^2
11 end
12
13 #vTerminal = -\frac{m g}{\kappa}
14
15 # Initialize the Params structure with specific values
16 m = 75 # kg
17 g = 9.81 # m/s^2
18
19 vTerminal = -2 # m/s
20 kappa = -m*g/vTerminal; @show kappa
21 params = Params(m, kappa, g)
22
23 # define the ODE dv/dt = f(v,params,t)
24 function f(v, params, t)
25     dv = -params.g .- (params.kappa / params.m) * v
26     return dv
27 end
28
29 # Set the initial condition as a vector
30 v0 = -25 # m/s speed when chute is opened
31 v0= [v0]
32
33 # Set the time interval
34 T = (0.0, 1.5) #
35
36 # Setup the ODE problem with out-of-place function
37 problem = ODEProblem{false}(f, v0, T, params)
38
39 # solve the ODE problem using the Runge Kutta Tsitouras 5/4 Integrator
40 sol = solve(problem, Tsit5());
41
42 # plot the solution
43 p1 = plot(sol, lw=3, guidefont = 20, xlabel="t", ylabel="v", legend=false, color=:blue)
44 ylims!(p1, (-30, 0))

```

Output

kappa = 367.875

Parachutist with a Nonlinear Drag Model

```

1 # Import useful packages
2 using DifferentialEquations
3 using Plots
4 gr()
5
6 # Parachutist with a nonlinear drag model
7 struct Params
8     m::Float64      # Mass in kilograms
9     kappa::Float64  # Drag coefficient
10    g::Float64      # Acceleration due to gravity in m/s^2
11 end
12

```

```

13 #vTerminal = -\frac{m g}{\kappa}
14
15 # Initialize the Params structure with specific values
16 m = 75 # kg
17 g = 9.81 # m/s^2
18
19
20 vTerminal = -2 # m/s
21 kappa = m*g/(vTerminal^2); @show kappa
22 params = Params(m, kappa, g)
23
24 # define the ODE dv/dt = f(v,params,t)
25 function f(v, params, t)
26     dv = -params.g .- (params.kappa / params.m) .* abs.(v).*v
27     return dv
28 end
29
30 # Set the initial condition as a vector
31 v0 = -25 # m/s speed when chute is opened
32 v0= [v0]
33
34 # Set the time interval
35 T = (0.0, 1.5) #
36
37 # Setup the ODE problem with out-of-place function
38 problem = ODEProblem{false}(f, v0, T, params)
39
40 # solve the ODE problem using the Runge Kutta Tsitouras 5/4 Integrator
41 sol = solve(problem, Tsit5());
42
43 # plot the solution
44 p1 = plot(sol, lw=3, guidefont = 20, xlabel="t", ylabel="v", legend=false, color=:blue)
45 ylims!(p1, (-30, 0))

```

Output

kappa = 183.9375

9.2.3 Finite Escape Time

The ODEs we have investigated so far have solutions that exist for all $t \geq 0$. This is NOT always the case. Consider

$$\frac{dx(t)}{dt} = 1 + x^2(t), \quad (9.8)$$

with initial condition $x_0 = x(0) = 0$. Then, using the method of separation of variables, we arrive at

$$\begin{aligned}
 \frac{1}{1+x^2} dx &= d\tau \\
 \int_0^{x(t)} \frac{1}{1+x^2} dx &= \int_0^t d\tau \\
 \text{atan}(x(t)) &= t \\
 x(t) &= \tan(t).
 \end{aligned}$$

We know the function $\tan(t)$ has a vertical asymptote at $\frac{\pi}{2}$, that is,

$$\lim_{t \rightarrow \frac{\pi}{2}^-} \tan(t) = \infty.$$

When the solution of the ODE has a vertical asymptote, we say that the ODE itself has a **finite escape time**: in finite time, the solution of the ODE ceases to exist; said another way, the solution explodes to $\pm\infty$ at a finite time called t_{escape} . As you might imagine, finite escape times drive a numerical solver crazy.

```

1 # Import useful packages
2 using DifferentialEquations
3 using Plots
4 gr()
5
6 # Finite Escape Time
7
8 # define the ODE dv/dt = f(v,params,t)
9 function f(x,~, t) # If there are no parameters to pass, use a tilde
10    dx = 1 .+ x.^2 # Function is vectorized
11    return dx
12 end
13
14 # Set the initial condition as a vector
15 x0 = 0.0 # m/s
16 x0= [x0]
17
18 # Set the time interval
19 T = (0.0, 3)
20
21 # Setup the ODE problem with out-of-place function
22 problem = ODEProblem{false}(f, x0, T)
23
24 # solve the ODE problem using the Runge Kutta Tsitouras 5/4 Integrator
25 sol = solve(problem, Tsit5());
26
27 # plot the solution
28 plot(sol)

```

Output

Warning: dt(4.440892098500626e-16) <= dtmin(4.440892098500626e-16) at t=1.5707114767768096. Aborting. There is either an error in your model specification or the true solution is unstable.

By “unstable”, the solver means “numerically unstable”. For you, it means the solutions blew up, and you need to decide if that is a “feature” or a “bug”! In our case, it is not a bug in the code. It is a feature of the model we are trying to simulate. While it may not be a desirable feature, the finite escape time is inherent in the ODE itself and has nothing to do with the capabilities of the numerical solver. If we had set the time interval to any value $0 < T < \frac{\pi}{2}$, the solver would have found a very good approximation of the solution. In the output’s first line, the solver tells us it encountered a problem at $t = 1.5707114767768096$, which is very close to $\frac{\pi}{2}$.

9.2.4 One-Dimensional ODE with Multiple Solutions

So far, we have associated a single solution to each ODE we have solved. Some nonlinear ODEs have more than one solution as highlighted in this fun video [The Dome Paradox: A Loophole in Newton’s Laws](#) by Up and Atom. Here, we present a simple

example of a one-dimensional ordinary differential equation (ODE) that exhibits multiple solutions. **Spoiler Alert: it solves the dome paradox!**

Consider the differential equation

$$\frac{dx(t)}{dt} = (x(t))^{2/3}, \quad x(0) = 0. \quad (9.9)$$

The right side of this ODE, $x^{2/3} = \sqrt[3]{x^2}$, is notable because it is not differentiable at $x = 0$. Indeed,

$$\frac{d}{dx}(x^{2/3}) = \frac{2}{3}x^{-\frac{1}{3}} = \frac{2}{3}\frac{1}{\sqrt[3]{x}}$$

blows up at the origin.

In terms of solutions, the trivial solution,

$$x(t) \equiv 0 \quad (9.10)$$

satisfies the ODE and the initial condition. Surprisingly, there are other (non-trivial) solutions as well. For instance, the function $\varphi : [0, \infty) \rightarrow \mathbb{R}$ by

$$\varphi(t) = \frac{t^3}{27}, \quad (9.11)$$

also satisfies the initial condition and the differential equation. Indeed,

- $\frac{d}{dt}\varphi(t) = \frac{d}{dt}\frac{t^3}{27} = \frac{t^2}{9}$, and
- $(\varphi(t))^{2/3} = \left(\left(\frac{t}{3}\right)^3\right)^{2/3} = \left(\frac{t}{3}\right)^2 = \frac{t^2}{9}$.

Perhaps even more mind-blowing², for all constants $c > 0$,

$$\varphi_c(t) = \begin{cases} 0 & \text{if } t < c, \\ \frac{(t-c)^3}{27} & \text{if } t \geq c, \end{cases} \quad (9.12)$$

is also a solution. This example illustrates that some differential equations, particularly those where the right side is not everywhere continuously differentiable, can have multiple solutions. Understanding the criteria under which a unique solution exists is crucial in the study of ODEs. We address this in Chapter 9.4.3.

A related phenomenon is that the trajectories of some ODEs can be highly sensitive to their initial conditions. In such cases, “small” changes in the initial condition can lead to “large” variations in the resulting time trajectory. In both mathematics and popular science, this behavior is often referred to as “**Chaos**”. If you’d like an accessible introduction to this topic, we recommend the video [Chaos: The Science of the Butterfly Effect](#) by Veritasium.

9.2.5 (Optional Read:) Can an ODE Have No Solution?

You see this question a lot on the internet, so we thought to address it by giving various examples. If they are all a bit wacky, that’s a sign that most ODEs do have a solution.

- $\left(\frac{dx(t)}{dt}\right)^2 + 1 = 0$ does not have a real solution for any initial condition for the same reason that $y^2 + 1 = 0$ has no real solutions.
- $\frac{dx(t)}{dt} = \frac{1}{x(t)}$ does not have a solution for the initial condition $x(0) = 0$ because the right side of the ODE blows up at $x(t) = 0$.
- $\frac{dx(t)}{dt} = \begin{cases} 1 & x \in \mathbb{Q} \\ 0 & \text{otherwise,} \end{cases}$ does not have a solution for any initial condition because the function on the right side of the ODE is not Riemann integrable.

Why is that a problem? By the First Fundamental Theorem of Calculus, if the equation has a solution $x(t)$, then it must be an antiderivative of $x'(t)$, meaning that $x(t) = x(a) + \int_a^t x'(\tau) d\tau$. But for that to hold, we need to make sense of the integral. Because we only know Riemann Integration, we are stuck. Now, if you learn the Lebesgue Integral, you have a way out of this corner.

²Why is this mind-blowing? Suppose you are building a gadget and through much labor, you show that one of its variables satisfies (9.9). Then the variable behaves almost like something “random” because, at any moment in time, the “universe” rolls a die, chooses $c > 0$, and spookily produces a non-zero solution. Wild!

The following result comes from Functional Analysis and Topology, branches of mathematics that math majors take their senior year or in grad school. It is only included here to underline that any ODE that may be of interest to an engineer can be proven to have a solution, justifying the wackiness of our “counter examples.”

[Peano's Existence Theorem] Let G be an open set in $\mathbb{R} \times \mathbb{R}^n$ and let $f : G \rightarrow \mathbb{R}^n$ be a continuous function. If (t_0, x_0) is a point in G , then there exists an interval I containing t_0 and a function $x : I \rightarrow \mathbb{R}^n$ that is a solution to the initial value problem

$$\frac{dx}{dt} = f(t, x), \quad x(t_0) = x_0. \quad (9.13)$$

Proof. The proof of Peano's Existence Theorem is based on the application of the Arzelà–Ascoli theorem and Schauder's fixed point theorem. The idea is to construct a sequence of approximate solutions using successive approximations and show that this sequence has a convergent subsequence that converges to an actual solution of the differential equation. The continuity of f ensures that the sequence of approximate solutions is equicontinuous and uniformly bounded, which allows the application of the Arzelà–Ascoli theorem. \square

9.2.6 (Optional Read:) The Independent Variable in an ODE can be any Strictly Increasing Quantity

The study of heat conduction is a fundamental aspect of thermodynamics and continuum mechanics. In particular, one-dimensional steady-state heat conduction can be modeled using a **second-order** ordinary differential equation. This model provides insights into how temperature varies along a rod or a similar object when the system has reached a steady state. If time were also considered, we would have a PDE.

Consider a long, thin rod of length L , with thermal conductivity k , and wrapped in thermal insulation along its length, except at the ends. The rod is heated at one end and cooled at the other, creating a temperature gradient along its length. Let x be a position along the rod and $T(x)$ the temperature at the point x . The temperature distribution $T(x)$ along the rod is described by the steady-state heat equation,

$$-k \frac{d^2T}{dx^2} = 0, \quad (9.14)$$

where k is the thermal conductivity of the material. The boundary conditions might be, for example, $T(0) = T_0$ (temperature at the heated end) and $T(L) = T_L$ (temperature at the cooled end).

The general solution to this differential equation is an affine (aka, linear plus a constant) temperature profile,

$$T(x) = C_1 x + C_2, \quad (9.15)$$

where C_1 and C_2 are constants determined by the boundary conditions. For the given boundary conditions, the solution becomes,

$$T(x) = \frac{T_L - T_0}{L} x + T_0. \quad (9.16)$$

This equation describes a linear temperature gradient from the heated end to the cooled end of the rod. While the examples we study will mostly have time as the independent variable, you will encounter various independent variables in other engineering courses.

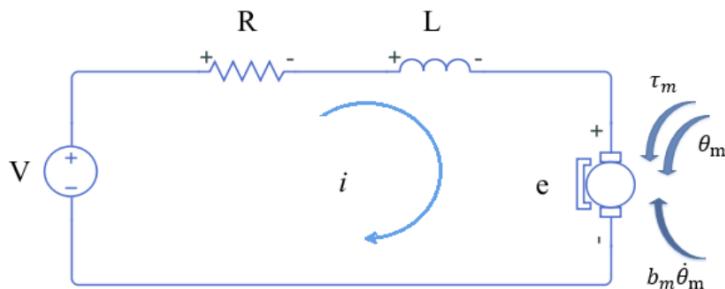
9.3 Higher-order ODEs and a Direct Current (DC) Motor Model

Most models of physical systems involve more than one derivative, or more than one variable, or both! Fig. 9.5-(b) shows the pendulum we studied in Chapter 6.4. Its differential equation is

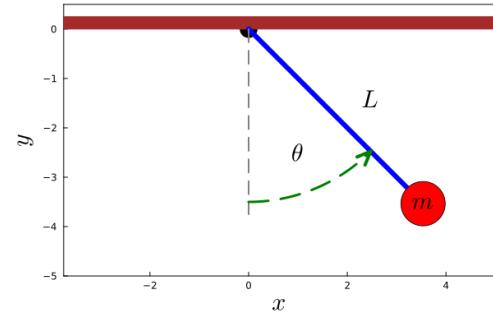
$$m_p \cdot L_p^2 \cdot \ddot{\theta} + m_p \cdot g \cdot L_p \cdot \sin(\theta) = \tau_p, \quad (9.17)$$

where τ_p is the torque applied to the pendulum at the revolute joint. This is an example of a **second-order ODE in a single variable**, θ . It is second order because the second derivative, $\ddot{\theta}$, appears in the model.

This section aims to show that real differential equations typically have higher order than one and contain multiple variables. In other words, higher-order ODEs are not just a mathematical curiosity invented to torture young learners. If you are already convinced of this, then go straight to Definition 9.10 and then move on. If you need more convincing, or you just want to see how motors can be modeled and integrated with robot arms, then read on.



(a)



(b)

Figure 9.5: (a) A [circuit model of DC motor](#) by Research Gate. (b) The classic pendulum formed by a massless rod with a point mass attached at its end from Fig. 6.10. The torque from the motor can be used to achieve a desired behavior, such as balancing the pendulum in the upright position. **Notice the clash in notation, where ME-types use L for length and EE-types use L for inductance. This is very common and you have to learn how to navigate it.** In the text, we add a subscript $_p$ to the parameters of the pendulum and a subscript $_m$ to the parameters of the motor.

Fig. 9.5-(a) shows a standard circuit representation of a DC motor; more comprehensive treatments can be found in [DC Motor Speed: System Modeling](#) by the MathWorks (creator of MATLAB) and [How does an Electric Motor work? DC Motor explained](#) by The Engineering Mindset. Using [Kirchoff's Current and Voltage Laws](#)³, its model can be written as

$$L_m \frac{di}{dt} + R_m i = V - b_m \dot{\theta}_m, \quad (9.18)$$

where θ_m is the angle of the motor's shaft, L_m is the inductance of the motor's windings and R_m is their resistance, V is the DC source voltage applied to the motor, and $b_m \dot{\theta}$ represents the “back emf” (Back Electromotive Force⁴). The output shaft of the motor is typically connected to the revolute joint of the pendulum through some sort of gearing that imposes

$$\theta_m = N\theta,$$

where N is the gear ratio, which could be in the range of 100. Gearing is needed because typical motors are efficient when spinning fast and producing low torque, while most robotics applications require high torque and low speed! **The current of the motor is now a second variable in the differential equations modeling the overall system.** The torque produced by a DC motor is proportional to the current, $\tau_m = K_m \cdot i$, where K_m is called the motor torque constant.

We have presented these details for two reasons:

- The differential equations you will encounter in engineering come from physical systems. Each discipline has its collection of important differential equation models. You typically do not learn how to develop the models in a mathematics course. Your author learned how to model through courses in Physics, Chemistry, EE, and ME (Robotics did not exist as an academic discipline back in the day), and through consulting work in industry. Pay attention in your courses to the models. Even if the models are simplified, they open doors to the more realistic models you will encounter in practice.
- The ODE model of the Cassie bipedal robot has 40 variables in it and fills up at least 100 pages in 10-point font. Through the application of appropriate software tools, the students who work with Cassie [manage the robot's ODE model](#)⁵ just fine!

Fact 9.9: A Third-Order ODE Model for a Motorized Pendulum

The overall system comprised of a DC motor connected to a pendulum, as shown in Fig. 9.5, obeys the following third-order

³You can learn these in Michigan’s EECS 215 Introduction to Electronic Circuits, EECS 314 Electrical Circuits, Systems, and Applications, or ROB 310 Robot Sensors and Signals.

⁴The back (or counter) electromotive force (emf) is the voltage generated by a running motor that acts to counter the supplied voltage. In the simple case of a permanent-magnet DC motor, it is proportional to the rotational speed of the motor.

⁵Y. Gong was a second-year MS student when designing the gait controller for Cassie. You do not see the model in the paper because 100 pages of printed formulas would serve no useful purpose.

ODE

$$m_p \cdot L_m \cdot L_p^2 \cdot \ddot{\theta} + (R_m \cdot m_p \cdot L_p^2) \cdot \ddot{\theta} + (m_p \cdot g \cdot L_m \cdot L_p \cdot \cos(\theta) + N \cdot K_m \cdot b_m) \cdot \dot{\theta} + R_m \cdot m_p \cdot g \cdot L_p \cdot \sin(\theta) = N \cdot K_m \cdot V. \quad (9.19)$$

where $V(t)$ is the voltage supplied to the motor. The derivation is given below as an optional read.

High-order ODEs of a Single Variable

Definition 9.10. An ODE of the form

$$y^{(n)} = f(y^{(n-1)}, \dots, \dot{y}, y, t) \quad (9.20)$$

is called an **n -th order (nonlinear) ODE**. While an ODE of the form,

$$y^{(n)} + a_{n-1}y^{(n-1)} + a_{n-2}y^{(n-2)} + \dots + a_1\dot{y} + a_0y = 0 \quad (9.21)$$

is called an **n -th order linear time-invariant ODE** when $a_k \in \mathbb{R}$ are constants, $0 \leq k \leq n - 1$.

- (a) Simple algebra is required to put (9.19) in the form of (9.20).
- (b) After which, replacing $\cos(\theta)$ by 1.0 and $\sin(\theta)$ by θ yields a linear time-invariant ODE model.
- (c) It is OK to have a term multiplying $y^{(n)}$ as long as it never vanishes.

Example 9.11. (Optional Read:) From the following equations for the pendulum driven by a DC motor, derive Fact 9.9.

$$\begin{aligned} m_p \cdot L_p^2 \cdot \ddot{\theta} + m_p \cdot g \cdot L_p \cdot \sin(\theta) &= \tau_p & (a) \\ L_m \frac{di}{dt} + R_m i &= V - b_m \dot{\theta}_m & (b) \\ \theta_m &= N\theta & (c) \\ \tau_m &= K_m i & (d) \\ \tau_p &= N\tau_m & (e) \end{aligned} \quad (9.22)$$

Solution: Substituting (9.22)-(d) and -(e) into (9.22)-(a) gives

$$m_p \cdot L_p^2 \cdot \ddot{\theta} + m_p \cdot g \cdot L_p \cdot \sin(\theta) = N \cdot K_m \cdot i.$$

Differentiating the above equation with respect to time gives

$$m_p \cdot L_p^2 \cdot \ddot{\theta} + m_p \cdot g \cdot L_p \cdot \cos(\theta) \cdot \dot{\theta} = N \cdot K_m \cdot \frac{di}{dt}, \quad (9.23)$$

where the dot notation for the derivative was not used on the motor current, i , for obvious reasons. To arrive at (9.19), we need to eliminate $\frac{di}{dt}$ from the above equation. To accomplish this, we solve (9.22)-(b) for $\frac{di}{dt}$, giving us

$$\frac{di}{dt} = -\frac{R_m}{L_m}i + \frac{1}{L_m}V - \frac{b_m}{L_m}\dot{\theta}_m.$$

Next, we replace the current i by

$$\begin{aligned} i &= \frac{1}{K_m}\tau_m \\ &= \frac{1}{N \cdot K_m}\tau_p \\ &= \frac{1}{N \cdot K_m}(m_p \cdot L_p^2 \cdot \ddot{\theta} + m_p \cdot g \cdot L_p \cdot \sin(\theta)) \\ &= \frac{m_p \cdot L_p^2}{N \cdot K_m} \cdot \ddot{\theta} + \frac{m_p \cdot g \cdot L_p}{N \cdot K_m} \cdot \sin(\theta), \end{aligned}$$

resulting in

$$\begin{aligned} \frac{di}{dt} &= -\frac{R_m}{L_m} \left(\frac{m_p \cdot L_p^2}{N \cdot K_m} \cdot \ddot{\theta} + \frac{m_p \cdot g \cdot L_p}{N \cdot K_m} \cdot \sin(\theta) \right) + \frac{1}{L_m} V - \frac{b_m}{L_m} \dot{\theta}_m \\ &= -\frac{R_m \cdot m_p \cdot L_p^2}{L_m \cdot N \cdot K_m} \cdot \ddot{\theta} - \frac{R_m \cdot m_p \cdot g \cdot L_p}{L_m \cdot N \cdot K_m} \cdot \sin(\theta) + \frac{1}{L_m} V - \frac{b_m}{L_m} \dot{\theta}_m. \end{aligned} \quad (9.24)$$

The final step is to substitute (9.24) into (9.23), giving

$$m_p \cdot L_p^2 \cdot \ddot{\theta} + m_p \cdot g \cdot L_p \cdot \cos(\theta) \cdot \dot{\theta} = N \cdot K_m \cdot \left(-\frac{R_m \cdot m_p \cdot L_p^2}{L_m \cdot N \cdot K_m} \cdot \ddot{\theta} - \frac{R_m \cdot m_p \cdot g \cdot L_p}{L_m \cdot N \cdot K_m} \cdot \sin(\theta) + \frac{1}{L_m} V - \frac{b_m}{L_m} \dot{\theta}_m \right). \quad (9.25)$$

Simplifying then gives (9.19).

Note: For a real project, one would use a Symbolic Package to combine and simplify the above equations. Hand calculations are error-prone and do not lend themselves to good documentation.

The following is inelegant and contains some hand calculations. However, the substitutions and simplifications are done symbolically, which is where the errors tend to occur.

```

1 using SymPy
2
3 # Define the symbolic variables
4 @syms th dth ddth mp Lp g taup Lm Rm i V bm N Km
5
6 # Define the expressions
7 #i = taup / (N * Km)
8 di = (V - bm * dth - Rm * i) / Lm
9
10 # Perform the substitution
11 di = di.subs(i, taup / (N * Km))
12 #
13 # (mp*Lp^2) * ddth = taup - mp*g*Lp*sin(th) and hence,
14 # taup = (mp*Lp^2) * ddth + mp*g*Lp*sin(th)
15 di = di.subs(taup, (mp*Lp^2)*ddth + mp*g*Lp*sin(th))
16 di = simplify(di)
17
18 # Display the substituted expression
19 println(di)

```

Output

```
(-Km*N*(-V + bm*dth) - Lp*Rm*mp*(Lp*ddth + g*sin(th))) / (Km*Lm*N)
```

```

1 #ddth = (N * Km * i - mp*g*Lp*sin(th)) / (mp*Lp^2)
2 ddth = (N * Km * di - mp*g*Lp*cos(th)*dth) / (mp*Lp^2)
3 dddth = simplify(dddth)
4
5 # Extracting numerator and denominator
6 num, den = dddth.as_numer_denom()
7
8 # print numerator and denominator of dddth, theta triple dot
9 println(num); println(" ")
10 println(den)
11 dendddth = simplify(dddth * den)
12
13 println(" ")

```

```
14 println("ODE is: $den * dddth = ", dendddth)
```

Output

```
-Km*N*(-V + bm*dth) - Lm*Lp*dth*g*mp*cos(th) - Lp*Rm*mp*(Lp*ddth + g*sin(th))  
Lm*Lp^2*mp
```

```
ODE is: Lm*Lp^2*mp * dddth = Km*N*(V - bm*dth) - Lm*Lp*dth*g*mp*cos(th)  
- Lp*Rm*mp*(Lp*ddth + g*sin(th))
```

All of the terms match up with the hand calculations. ■

9.4 Vector ODEs: Multidimensional Dynamics, Single Independent Variable

We first discuss an important source of “vector ODEs” with the goal of understanding what that even means! Then, we can discuss what a solution of the ODE means, when solutions are unique, and when solutions exist for an unbounded time interval (i.e., no finite escape time).

9.4.1 The Return of the Robot Equations

In Prop. 6.38, we introduced the Robot equations,

$$D(q) \cdot \ddot{q} + C(q, \dot{q}) \cdot \dot{q} + G(q) = B \cdot \tau, \quad (9.26)$$

where,

- $q \in \mathbb{R}^n$ is a vector of generalized position coordinates,
- $\dot{q} \in \mathbb{R}^n$ is the corresponding vector of velocities,
- $\ddot{q} \in \mathbb{R}^n$ is the corresponding vector of accelerations,
- and $\tau \in \mathbb{R}^m$ is a vector of motor torques.

Moreover,

- $D(q)$ is the $n \times n$ mass-inertia matrix and is invertible,
- $G(q) := \nabla V(q)$, is the $n \times 1$ gradient of the potential energy,
- the n vector $C(q, \dot{q}) \cdot \dot{q}$ contains all the remaining terms in (6.61),
- and B is the $n \times m$ torque distribution matrix.

Due to the appearance of $\ddot{q} = \frac{d^2}{dt^2}q$ in (9.26), the Robot Equations are a **vector of second-order ODEs in the independent variable t** . They look totally intimidating. Our first goal is to rewrite them as a vector of first-order ODEs,

$$\dot{x} = f(x, \tau),$$

which, beyond appearing “more approachable”, will greatly help us with both our analytical work and our numerical solvers. **It is important to note that so few of the Robot Equations have closed-form solutions that we routinely use numerical algorithms to compute solutions.**

Method 9.12: From Second-Order ODEs to First-Order ODEs

Consider the Robot Equations in (9.26) and let $x_1 := q$, $x_2 := \dot{q}$. In terms of these variables, we have that

$$\dot{x}_1 = \dot{q} = x_2$$

$$\begin{aligned}\dot{x}_2 &= \ddot{q} = -D^{-1}(q) \cdot [C(q, \dot{q}) \cdot \dot{q} + G(q)] + D^{-1}(q) \cdot B \cdot \tau \\ &= -D^{-1}(x_1) \cdot [C(x_1, x_2) \cdot x_2 + G(x_1)] + D^{-1}(x_1) \cdot B \cdot \tau.\end{aligned}\quad (9.27)$$

Once we stack x_1 and x_2 to form the $2n \times 1$ vector x , we have that (9.26) is equivalent to

$$\underbrace{\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix}}_{\dot{x}} = \underbrace{\begin{bmatrix} x_2 \\ -D^{-1}(x_1) \cdot [C(x_1, x_2) \cdot x_2 + G(x_1)] + D^{-1}(x_1) \cdot B \cdot \tau \end{bmatrix}}_{f(x, \tau)}. \quad (9.28)$$

An equation of the form

$$\dot{x} = f(x, \tau) \quad (9.29)$$

is called in various situations,

- a vector of first-order ODEs,
- a vector of controlled first-order ODEs (due to the presence of the motor torques, τ , which are used to control the robot),
- a **system of first-order ODEs**,
- or a **state-variable model**.

The components of the vector x are called state variables. They are important because you need to provide initial conditions for each state variable in order to uniquely specify a solution to the ODE.

Note: As you know from ROB 101 *Computational Linear Algebra*, real engineers abhor computing inverses of matrices. Instead, we recognize that \dot{x}_2 is the unique solution to the linear system of equations

$$D(x_1) \cdot \dot{x}_2 = [C(x_1, x_2) \cdot x_2 + G(x_1)] + B \cdot \tau,$$

which we would solve via LU Factorization or QR Factorization. In Julia, we have the backslash operator, which will select either LU or QR as it deems most appropriate for solving the equations, giving us

$$\dot{x}_2 = D(x_1) \setminus (-[C(x_1, x_2) \cdot x_2 + G(x_1)] + B \cdot \tau)$$

as the better way to implement in code the second row of (9.28).

Example 9.13. Here is the pendulum model from Example 6.37,

$$\underbrace{\frac{d}{dt} \frac{\partial \mathcal{L}(q, \dot{q})}{\partial \dot{q}}}^{m \cdot L^2 \cdot \ddot{\theta}} - \underbrace{(-m \cdot g \cdot L \cdot \sin(\theta))}_{\frac{\partial \mathcal{L}(q, \dot{q})}{\partial q}} = \underbrace{0}_{\Gamma}, \quad (9.30)$$

where $\Gamma = 0_{1 \times 1}$ because a motor has not been installed at the pivot and there is no friction. Express the model as a vector first-order ODE and provide a simulation with initial conditions $\theta_0 = \frac{3\pi}{8}$ and $\dot{\theta}_0 = 0$. Recall that the pendulum is pointing straight down at rest with $\theta_0 = 0$ and $\dot{\theta}_0 = 0$.

Solution: Following the template given in Method 9.12, we define $x_1 := \theta$ and $x_2 := \dot{\theta}$. After a bit of algebra, this gives

$$\begin{aligned}\dot{x}_1 &= x_2 \\ \dot{x}_2 &= -\frac{g}{L} \cdot \sin(x_1),\end{aligned}$$

because $\frac{m \cdot g \cdot L}{m \cdot L^2} = \frac{g}{L}$. It is perfectly fine to leave the equations in this form. It is equally fine to express them as

$$\underbrace{\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix}}_{\dot{x}} = \underbrace{\begin{bmatrix} x_2 \\ -\frac{g}{L} \cdot \sin(x_1) \end{bmatrix}}_{f(x)}.$$

The initial condition becomes

$$x(0) = \begin{bmatrix} \theta_0 \\ \dot{\theta}_0 \end{bmatrix} = \begin{bmatrix} \frac{3\pi}{8} \\ 0 \end{bmatrix}.$$

```

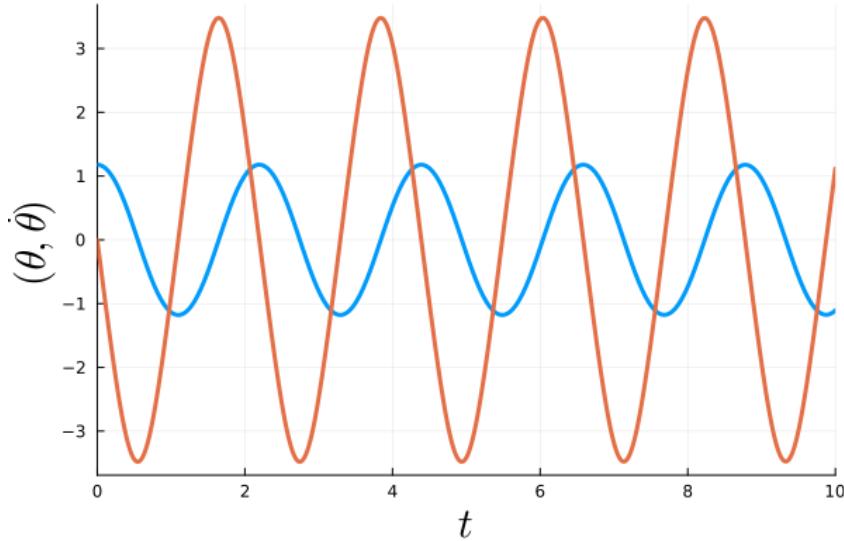
1 # Import useful packages
2 using DifferentialEquations
3 using Plots
4 using LaTeXStrings
5 gr()
6
7 # Pendulum Model
8
9 struct Params
10    m::Float64      # Mass in kilograms
11    L::Float64       # Length in meters
12    g::Float64       # Acceleration due to gravity in m/s^2
13 end
14
15
16 # Initialize the Params structure with specific values
17 m = 2
18 L = 1
19 g = 9.81
20 params = Params(m, L, g)
21
22
23 # define the ODE dx/dt = f(x,params,t)
24 function f(x, params, t)
25     dx = [x[2], -(params.g/params.L)*sin.(x[1])]
26     return dx
27 end
28
29 # Set the initial condition as a vector
30 x0 = [3*pi/8; 0.0]
31
32 # Set the time interval
33 T = (0.0, 10)
34
35 # Setup the ODE problem with out-of-place function
36 problem = ODEProblem{false}(f, x0, T, params)
37
```

```

38 # solve the ODE problem using the Runge Kutta Tsitouras 5/4 Integrator
39 sol = solve(problem, Tsit5());
40
41 # plot the solution
42 p1 = plot(sol, lw=3, guidefont = 20, xlabel=L"t", ylabel=L"(\theta, \dot{\theta})", legend=false)

```

Output The blue trace is the angle of the pendulum, while the red trace is its angular velocity. The pendulum oscillates as expected.



Fact 9.14: 2-Link Manipulator

The base angle q_1 is an absolute angle, with its zero position aligned with the x -axis. The angle of the second link, q_2 , is relative to the base (first) link. Hence, $q_2 = 0$ has the second link extending straight out from the base link. The angular velocities are denoted dq_1 and dq_2 in the code. Actuators have been placed at the base link and the second link, though we will not be using them in this Chapter.

```

1 function dyn_mod_2LinkManipulator(q, dq)
2 # DYN_MOD_2LINKMANIPULATOR
3 # 2023-11-20 17:26:11
4 #
5 # Author: Grizzle
6 #
7 # Model NOTATION: D(q)ddq + C(q,dq)*dq + G(q) = B*tau
8 # The Robot Equations: From Lagrange's Equations of Motion
9 #
10 g, L1, L2, m1, m2 = modelParameters()
11 #
12 # Variable names for the model
13 q1, q2 = q
14 dq1, dq2 = dq
15 #
16 D = zeros(2, 2)
17 D[1, 1] = L1^2*m1 + m2*(L1^2 + 2*L1*L2*cos(q2) + L2^2)
18 D[1, 2] = L2*m2*(L1*cos(q2) + L2)
19 D[2, 1] = L2*m2*(L1*cos(q2) + L2)

```

```

20 D[2, 2] = L2^2*m2
21 #
22 C = zeros(2, 2)
23 C[1, 1] = -L1*L2*dq2*m2*sin(q2)
24 C[1, 2] = -L1*L2*m2*(dq1 + dq2)*sin(q2)
25 C[2, 1] = L1*L2*dq1*m2*sin(q2)
26 #
27 G = zeros(2)
28 G[1] = g*m2*(L1*cos(q1) + L2*cos(q1 + q2)) + L1*g*m1*cos(q1)
29 G[2] = L2*g*m2*cos(q1 + q2)
30 #
31 B = zeros(2, 2)
32 B[1, 1] = 1
33 B[2, 2] = 1
34 #
35 JacG = zeros(2, 2)
36 JacG[1, 1] = g*m2*(-L1*sin(q1) - L2*sin(q1 + q2)) - L1*g*m1*sin(q1)
37 JacG[1, 2] = -L2*g*m2*sin(q1 + q2)
38 JacG[2, 1] = -L2*g*m2*sin(q1 + q2)
39 JacG[2, 2] = -L2*g*m2*sin(q1 + q2)
40 #
41 return (D=D, C=C, G=G, B=B, JacG=JacG)
42 end

```

To bring the model into the workspace, use

```

1 # bring the model into the workspace for later use
2 include("dyn_mod_2LinkManipulator.jl")

```

Example 9.15. A two-link manipulator model is given in Fact 9.14, following the same angle convention as in Fig. 1.6-(b). Express the model as a vector first-order ODE and provide a simulation with initial conditions

$$\begin{aligned} q_1(0) &= -\frac{3\pi}{8} \\ q_2(0) &= 0 \\ \dot{q}_1(0) &= 0 \\ \dot{q}_2(0) &= 0. \end{aligned}$$

Solution: Following the template given in Method 9.12, we define

$$x_1 := \begin{bmatrix} q_1 \\ q_2 \end{bmatrix} \text{ and } x_2 := \begin{bmatrix} \dot{q}_1 \\ \dot{q}_2 \end{bmatrix}.$$

We do not symbolically compute the inverse of the mass-inertia matrix because, in larger, more realistic models, it would make our simulation very slow. Hence, we write the model as

$$\dot{x}_1 = x_2$$

$$D(x_1)\dot{x}_2 = -C(x_1, x_2) \cdot x_2 - G(x_1),$$

and use the “backslash” operator to solve for \dot{x}_2 . In code, it looks quite clean.

```
# Import useful packages
```

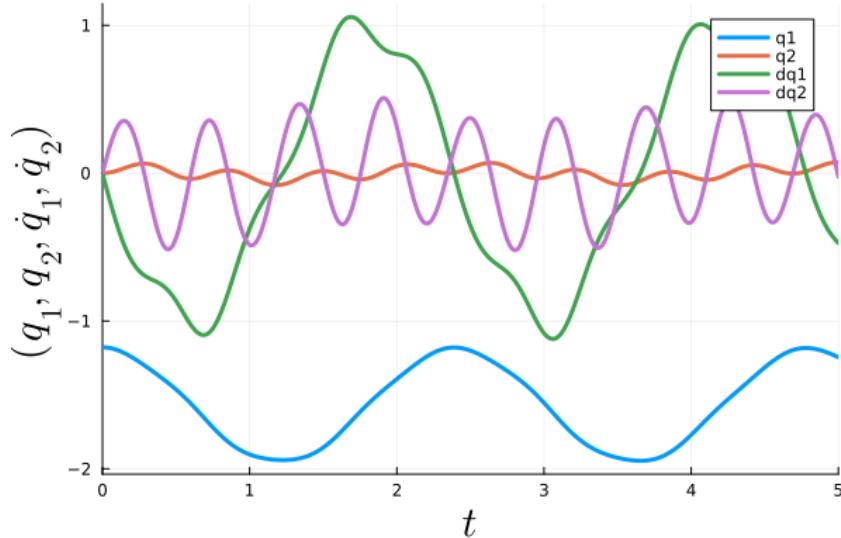
```

2 using DifferentialEquations
3 using Plots
4 using LaTeXStrings
5 using LinearAlgebra # for the backslash command
6 gr()
7
8 # 2-Link Manipulator Model
9 function modelParameters()
10    g = 9.81
11    L1 = 1.0
12    L2 = 0.5
13    m1 = 1.0
14    m2 = 3.0
15    return g, L1, L2, m1, m2
16 end
17
18 # bring the model into the workspace for later use
19 include("dyn_mod_2LinkManipulator.jl")
20
21 # define the ODE dx/dt = f(x,t)
22 #
23 # the parameters are passed via modelParameters() in the
24 # function dyn_mod_2LinkManipulator(q, dq)
25 #
26 function f(x, ~, t)
27    n = floor(Int, length(x) / 2) # In Julia n/2 is a Float64
28    q = x[1:n]
29    dq = x[n+1:end]
30    model = dyn_mod_2LinkManipulator(q, dq)
31    dx1 = dq
32    dx2 = (model.D) \ (-model.C*dq - model.G) # note the use of backslash
33    dx=[dx1;dx2]
34    return dx
35 end
36
37 # Set the initial condition as a vector
38 x0= [-3pi/8; 0.0; 0; 0]
39
40 # Set the time interval
41 T = (0.0, 5)
42
43 # Setup the ODE problem with out-of-place function
44 # params removed
45 problem = ODEProblem{false}(f, x0, T)
46
47 # solve the ODE problem using the Runge Kutta Tsitouras 5/4 Integrator
48 sol = solve(problem, Tsit5());
49
50 # plot the solution
51 p1 = plot(sol, lw=3, guidefont=20, xlabel=L"t", ylabel=L"(q_1, q_2, \dot{q}_1, \dot{q}_2)", label=["q1" "q2" "dq1" "dq2"], legend=:topright)

```

Output The initial condition has the manipulator arm hanging downward like a pendulum ($q_1(0) = -\frac{\pi}{2}$ would be straight down). Consequently, q_1 oscillates similar to a single-link downward pointing pendulum, with one important change, some of the “energy” is transferred to the second link, which begins oscillating as well, even though it started in perfect alignment with the base link. Robot models with even two planar links have very complicated dynamics. Can you imagine the dynamics of a 3D bipedal robot, such as

Cassie, with 14 links? Your answer should be an emphatic **NO!**



9.4.2 Solutions of First-order Vector ODEs

What is a Solution?

Definition 9.16. A function, $\varphi : [t_0, T) \rightarrow \mathbb{R}^n$, $T > t_0$, is a **solution** to the Ordinary Differential Equation (ODE)

$$\dot{x}(t) = f(x(t)), \quad x(t_0) = x_0, \quad (9.31)$$

if

- (a) $\varphi : [t_0, T) \rightarrow \mathbb{R}^n$ is continuous (meaning, each of its components is continuous);
- (b) $\varphi(t_0) = x_0$ (satisfies the initial condition); and
- (c) for all $t \in (t_0, T)$, the derivative $\dot{\varphi}(t) := \frac{d}{dt}\varphi(t)$ exists and $\dot{\varphi}(t) - f(\varphi(t)) = 0_{n \times 1}$ (satisfies the differential equation).

Note: It is very common to drop the time variable and write the ODE as $\dot{x} = f(x)$, $x(t_0) = x_0$. The “dot” notation for the derivative indicates the dependence on time.

Example 9.17. Show that for all constants $c > 0$,

$$\varphi_c(t) = \begin{cases} 0 & \text{if } t < c, \\ \frac{(t-c)^3}{27} & \text{if } t \geq c, \end{cases}$$

is a solution to the differential equation $\dot{x} = \sqrt[3]{x^2}$, $x(0) = 0$.

Solution We need to check the three conditions in Definition 9.16.

- (a) **Continuity:** The only point where continuity of $\varphi_c(t)$ may be in doubt is $t = c$. We check, therefore, the two one-sided limits

and compare them to the value of the function at c , viz,

$$\lim_{t \rightarrow c^-} \varphi_c(t) = \lim_{t \rightarrow c^-} 0 = 0 \quad (\text{the function is identically zero for } t < c)$$

$$\lim_{t \rightarrow c^+} \varphi_c(t) = \lim_{t \rightarrow c^+} \frac{(t-c)^3}{27} = 0, \quad (\text{using the function's definition for } t > c)$$

$$\varphi_c(t) \Big|_{t=c} = \frac{(t-c)^3}{27} \Big|_{t=c} = 0.$$

Hence, the function is continuous at c and therefore, at all points in its domain.

$$(b) \textbf{Satisfies the Initial Condition: } 0 = \frac{(t-c)^3}{27} \Big|_{t=c} = x_0.$$

(c) **Satisfies the Differential Equation:**

- For $t < c$, $\varphi_c(t)$ is identically zero, and hence its derivative is also zero. Also, $(0)^{2/3} = 0$. Hence, $\dot{\varphi}_c(t) - f(\varphi_c(t)) = 0$ for all $t < c$.
- For $t > c$, $\varphi_c(t) = \frac{(t-c)^3}{27}$, and hence by the Power Rule,

$$\frac{d}{dt} \varphi_c(t) = 3 \frac{(t-c)^2}{27} = \frac{(t-c)^2}{3^2} = \left(\frac{(t-c)}{3} \right)^2.$$

This takes care of the left side of the differential equation. Turning to the right side,

$$(\varphi_c(t))^{2/3} = \left(\frac{(t-c)^3}{27} \right)^{2/3} = \left(\left(\frac{(t-c)}{3} \right)^3 \right)^{2/3} = \left(\frac{(t-c)}{3} \right)^2.$$

Therefore, for all $t \neq c$, $\frac{d}{dt} \varphi_c(t) - (\varphi_c(t))^{2/3} = 0$. **What about $t = c$** , the point where the solution transitions from being identically zero to being a quadratic function of t ? We need to evaluate the derivative at that point using left and right limits. It all works out, but we skip the details because there is a more general notion of a solution given in Chapter 9.4.4 that obviates the calculation!

■.

9.4.3 Existence and Uniqueness of Solutions

We recall from ROB 101 *Computational Linear Algebra* that the **Euclidean norm** of a vector $x = [x_1 \ x_2 \ \dots \ x_n]$ is defined as

$$\|x\|_2 := \sqrt{(x_1)^2 + (x_2)^2 + \dots + (x_n)^2} = \sqrt{\sum_{i=1}^n (x_i)^2}. \quad (9.32)$$

There are also other norms, such

$$\|x\|_p := \sqrt[p]{\sum_{i=1}^n (x_i)^p}, \quad 1 \leq p < \infty \quad \text{and} \quad \|x\|_{\max} := \max_{1 \leq i \leq n} |x_i|.$$

Norm of a Matrix

Definition 9.18. Let A be an $n \times m$ matrix with entries a_{ij} . The **Euclidean norm** of A is defined as

$$\|A\|_2 := \sqrt{\sum_{i=1}^n \sum_{j=1}^m (a_{ij})^2}, \quad (9.33)$$

the square root of the sum of all squared entries of A . This is the same as stacking the columns of A into a large vector and applying the Euclidean norm to the resulting vector.

Note: The Euclidean or 2-norm is the default norm in Julia obtained with `norm(A)`. There are many other ways to define norms on matrices. It is often easiest to use the following matrix norm

$$\|A\|_\infty := \max_{i,j} |a_{ij}|, 1 \leq i, j \leq n, \quad (9.34)$$

which goes by the name **max norm** or the **infinity norm**.

Open Ball of Radius r

Open balls generalize to \mathbb{R}^n the notion of open intervals in \mathbb{R} .

Definition 9.19. Let x_0 be a point in \mathbb{R}^n and $r > 0$ a constant. The **open ball of radius r centered about x_0** is defined as

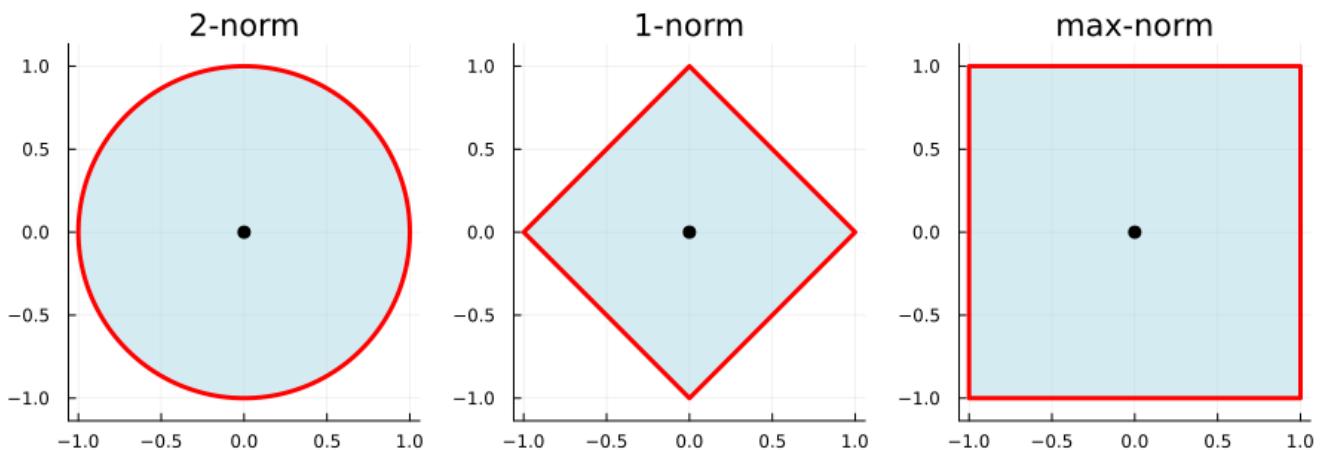
$$B_r(x_0) := \{x \in \mathbb{R}^n \mid \|x - x_0\| < r\}. \quad (9.35)$$

The ball is called “open” because it does not include the boundary points where $\|x - x_0\| = r$.

Note: Open balls are used to specify when some property holds **locally** (meaning near a given point x_0), versus **globally** (meaning it holds for all $x \in \mathbb{R}^n$).

Example 9.20. Sketch the open balls about the origin for the 2-norm, 1-norm, and max-norm.

Solution:



```

1 using Plots
2
3 # Function to plot the unit ball for a given norm
4 function plot_unit_ball(norm_type::String)
5     theta = range(0, stop=2π, length=1000) # parameter for the angle
6
7     if norm_type == "2-norm"
8         # Circle: x^2 + y^2 = 1
9         x = cos.(theta)
10        y = sin.(theta)
11    elseif norm_type == "1-norm"
12        # Diamond: |x| + |y| = 1, simplified representation
13        x = [0, 1, 0, -1, 0]
14        y = [1, 0, -1, 0, 1]
15    elseif norm_type == "max-norm"
16        # Square: max(|x|, |y|) = 1, simplified representation
17        x = [1, 1, -1, -1, 1]
18        y = [1, -1, -1, 1, 1]
19    else
20        error("Unknown norm type")
21    end
22
23    plot(x, y, label="", aspect_ratio=:equal, title="$(norm_type)", fill=(0, 0.5, :
24 lightblue), lw=3, color=:red)
25 end
26
27 # Plot unit balls for each norm
28 p1 = plot_unit_ball("2-norm")
29 p2 = plot_unit_ball("1-norm")
30 p3 = plot_unit_ball("max-norm")
31
32 # Display all plots in a single figure
33 pall= plot(p1, p2, p3, layout=(1, 3), size=(900, 300))

```

Output See the above plot.

Proposition 9.21: Existence and Uniqueness of Solutions

Consider a vector-valued ODE,

$$\dot{x} = f(x), \quad x(t_0) = x_0 \in \mathbb{R}^n. \quad (9.36)$$

- (a) **(Local Existence and Uniqueness of a Solution)** If for some $r > 0$, the Jacobian, $\frac{\partial f(x)}{\partial x}$, exists and is continuous for all $x \in B_r(x_0)$, then there exists $\delta > 0$ such that the ODE has a unique solution on the interval $[t_0, t_0 + \delta]$.
- (b) **(Global Existence and Uniqueness of a Solution)** If there exists a constant $0 \leq L < \infty$ such that, for all $x \in \mathbb{R}^n$, $\left\| \frac{\partial f(x)}{\partial x} \right\| \leq L$, then the ODE has a unique solution on the unbounded interval $[t_0, \infty)$.

Note: It is often easier to check the condition for global existence and uniqueness with the max (aka, infinity norm) because it can be computed element by element; indeed, $\left\| \frac{\partial f(x)}{\partial x} \right\|_{\infty} \leq L$ for all $x \in \mathbb{R}^n$ if, and only if, $\sup_{x \in \mathbb{R}^n} \left| \frac{\partial f_i(x)}{\partial x_j} \right| \leq L$.

Note 2: Most of the ODEs that appear in engineering practice satisfy the condition for local existence and uniqueness of a solution. However, the vector ODEs that model robots are typically quadratic in the angular velocities. Hence, their partial

derivatives will be linear in the angular velocities, which means the global bound $\|\frac{\partial f(x)}{\partial x}\| \leq L$ fails.

Note: This topic is treated in great detail in Michigan's EECS 562 Nonlinear Systems and Control; see it for proofs.

Exercise 9.22. Apply Prop. 9.21 to the ODEs given in Fact 9.2, Four Models of Chutes, to check for the local or global existence and uniqueness of their solutions. Comment on whether your conclusions match those of our previous calculations (of the solutions), or not.

For ease of reference, we repeat the ODEs here, rewritten slightly so that the left side only contains a derivative.

- (a) **Model 1: Parallel Motion, Linear Drag:** $\frac{dv(t)}{dt} = -\frac{\kappa}{m}v(t)$, $v(t_0) = v_0$.
- (b) **Model 2: Parallel Motion, Quadratic Drag:** $\frac{dv(t)}{dt} = -\frac{\kappa}{m}v^2(t)$, $v(t_0) = v_0$.
- (c) **Model 3: Vertical Motion, Linear Drag:** $\frac{dv(t)}{dt} = -\frac{\kappa}{m}v(t) - g$, $v(t_0) = v_0$.
- (d) **Model 4: Vertical Motion, Quadratic Drag:** $\frac{dv(t)}{dt} = \frac{\kappa}{m}v^2(t) - m$, $v(t_0) = v_0$.

Solutions:

We use the generic notation, $\dot{v} = f(v)$, for each of the models.

- (a) $\dot{v} = -\frac{\kappa}{m}v$, $v(t_0) = v_0$. **Ans.** Prop. 9.21 implies global existence and uniqueness of solutions.

We have $f(v) = -\frac{\kappa}{m}v$. Therefore, $\frac{df(v)}{dv} = -\frac{\kappa}{m}$ and $|\frac{df(v)}{dv}| = \left|-\frac{\kappa}{m}\right| = \frac{\kappa}{m}$ is finite for all $v \in \mathbb{R}$. Therefore, Prop. 9.21 implies a solution exists and is unique for all $t > t_0$ and all $v_0 \in \mathbb{R}$.

Comparison to our calculations in Example 9.3: we obtained a globally defined solution but had no means to check uniqueness.

- (b) $\dot{v} = -\frac{\kappa}{m}v^2$, $v(t_0) = v_0$. **Ans.** Prop. 9.21 implies local existence and uniqueness of solutions.

We have $f(v) = -\frac{\kappa}{m}v^2$. Therefore, $\frac{df(v)}{dv} = -2\frac{\kappa}{m}v$. The magnitude of the Jacobian, $|\frac{df(v)}{dv}|$, is finite on any open ball $B_r(v_0)$, $0 < r < \infty$, but grows unbounded as $|v| \rightarrow \infty$. Therefore, Prop. 9.21 implies the existence of $\delta > 0$ such that a solution exists and is unique for all $t \in [t_0, t_0 + \delta]$ and all $v_0 \in \mathbb{R}$, but does not imply global existence and uniqueness.

Comparison to our calculations in Example 9.5: we obtained a globally defined solution, even though the sufficient conditions of the Proposition only guaranteed a local solution. At the time we solved this problem, we had no means to check uniqueness.

Note: If we “flip the sign” on the ODE to $\frac{dv}{dt} = \frac{\kappa}{m}v^2$, then $|\frac{df(v)}{dv}| = 2\frac{\kappa}{m}|v|$ as before, again implying only local existence and uniqueness of solutions. This time, the method of separation of variables gives

$$v(t) = \frac{v_0}{1 - v_0 \frac{\kappa}{m}(t - t_0)},$$

which blows up at $t = t_0 + \frac{m}{\kappa v_0}$ because the denominator vanishes. By flipping the sign in front of the equation, we go from one with a globally defined solution to one that blows up in finite time. The test given in Prop. 9.21 is based on the norm of the Jacobian of the right side of the equation, and norms are invariant under sign changes.

- (c) $\dot{v} = -\frac{\kappa}{m}v(t) - g$, $v(t_0) = v_0$. **Ans.** Prop. 9.21 implies global existence and uniqueness of solutions.

We have $f(v) = -\frac{\kappa}{m}v - g$. Therefore, $\frac{df(v)}{dv} = -\frac{\kappa}{m}$ and $|\frac{df(v)}{dv}| = \left|-\frac{\kappa}{m}\right| = \frac{\kappa}{m} < \infty$. Therefore, Prop. 9.21 implies a solution exists and is unique for all $t > t_0$ and all $v_0 \in \mathbb{R}$.

Comparison to our calculations in Example 9.7: we obtained a globally defined solution but had no means to check uniqueness.

- (d) $\dot{v} = \frac{\kappa}{m} v^2(t) - g$, $v(t_0) = v_0$. **Ans.** Prop. 9.21 implies local existence and uniqueness of solutions.

We have $f(v) = -\frac{\kappa}{m}v^2 - g$. Therefore, $\frac{df(v)}{dv} = -2\frac{\kappa}{m}v$ and the magnitude of the Jacobian, $|\frac{df(v)}{dv}|$, is finite on any open ball $B_r(v_0)$, $0 < r < \infty$, but grows unbounded as $|v| \rightarrow \infty$. Therefore, Prop. 9.21 implies the existence of $\delta > 0$ such that a solution exists and is unique for all $t \in [t_0, t_0 + \delta]$ and all $v_0 \in \mathbb{R}$.

Comparison to our calculations in Example 9.8: we were unable to obtain an analytical solution. Our numerical calculations showed the existence of a solution over an adequately large time domain. Our previous work provided no means to check uniqueness. ■

Exercise 9.23. Revisit the ODE in (9.9), where we produced an infinite number of solutions (non-uniqueness on steroids),

$$\dot{x} = (x)^{2/3}, \quad x(0) = 0.$$

In what manner does it fail the conditions in Prop. 9.21?

Solution: The discussion in Chapter 9.2.4 gave us a heads up on the issue, but we'll continue as if that had not happened! Define $f(x) := x^{2/3}$. Prop. 9.21 tells us to evaluate

$$\frac{df(x)}{dx} = \frac{d}{dx}(x^{2/3}) = \frac{2}{3}x^{-\frac{1}{3}} = \frac{2}{3}\frac{1}{\sqrt[3]{x}},$$

which blows up at the initial condition, $x_0 = 0$. Hence, none of the conditions of Prop. 9.21 are met, so it has nothing to say about the solutions of the ODE. ■

Exercise 9.24. Differential equations are pivotal in modeling epidemics, showcasing how susceptible individuals ($x_1(t)$) and infectious individuals ($x_2(t)$) interact over time. The transition of susceptibles to infectious hinges on direct contact, captured by $\beta x_1 x_2$, where β denotes the infection rate. Furthermore, infectious individuals either recover, gaining immunity, or succumb to the illness, at a rate proportional to γx_2 , with γ reflecting the disease's time constant. Vaccination efforts, represented by a rate u , can reduce the number of susceptible individuals, illustrating the model's capacity to inform public health interventions.

The dynamics are governed by two differential equations,

$$\frac{dx_1}{dt} = -\beta x_1 x_2 - u, \tag{9.37}$$

$$\frac{dx_2}{dt} = \beta x_1 x_2 - \gamma x_2, \tag{9.38}$$

where, $x_1(t)$ and $x_2(t)$ represent susceptible and infectious individuals, respectively, integrating how the disease spreads and recedes within the community. Parameters β and γ quantify infection and recovery rates, pivotal in understanding the epidemic's trajectory. This model, by delineating interaction dynamics, aids in predicting and mitigating disease spread, underscoring the value of mathematical approaches in public health.

The epidemic model depends on two parameters $\beta > 0$ and $\gamma > 0$, and two states, x_1 and x_2 . Assuming u is constant,

- (a) determine if local existence and uniqueness of solutions can be guaranteed for all $x_1 \geq 0$, $x_2 \geq 0$; moreover;
- (b) determine if global existence and uniqueness of solutions can be guaranteed.

Solution:

Define $f(x) := \begin{bmatrix} -\beta x_1 x_2 - u \\ \beta x_1 x_2 - \gamma x_2 \end{bmatrix}$.

- (a) Local existence and uniqueness of solutions: **Ans.** Yes.

The Jacobian,

$$\frac{\partial f(x)}{\partial x} = \begin{bmatrix} -\beta x_2 & -\beta x_1 \\ \beta x_2 & \beta x_1 - \gamma \end{bmatrix},$$

is continuous everywhere, and hence, the condition for local existence and uniqueness is met for all $x_1 \geq 0, x_2 \geq 0$.

- (b) Global existence and uniqueness of solutions: **Ans.** No.

We use the max norm and note immediately that the $(1, 1)$ element of the Jacobian satisfies $\sup_{x_1 \geq 0, x_2 \geq 0} \left| \frac{f_1(x)}{\partial x_1} \right| = \sup_{x_2 \geq 0} |-\beta x_2| = \infty$, and hence

$$\left\| \frac{\partial f(x)}{\partial x} \right\|_\infty > L$$

for all finite L .

Note: Does this mean that solutions definitely blow up? No! Our calculations mean that we cannot prove global existence and uniqueness with the elementary means treated in the textbook. More powerful methods from Michigan's EECS 560 *Nonlinear Systems and Control* must be employed. Consider $V(x_1, x_2) := x_1 + x_2$, the total population of individuals. Then an easy calculation shows that $\frac{d}{dt}V(x_1, x_2) = \dot{x}_1 + \dot{x}_2 = -u - \gamma x_2 \leq 0$ when $u \geq 0$. Hence, the total population remains bounded, ruling out a finite escape time, and ensuring global existence and uniqueness of solutions. ■

ODE Examples from a Plethora of Engineering Domains: The (free) textbook, [State Models and Dynamical Systems](#) by Prof. Emeritus N. Harris McClamroch, aims to enhance undergraduate students' understanding of dynamic physical processes through applied mathematics and systems theory. It is written assuming that readers have a background in calculus, differential equations, Laplace transforms, physics, mechanics, and circuits. The book introduces mathematical systems theory, using case studies to clarify complex concepts.

Example 9.25. Show that the condition for global existence and uniqueness fails for the 2-link manipulator model from Fact 9.14 and Example 9.15.

Solution: Let $\dot{x} = f(x)$ be the state variable model as expressed in Example 9.15, after inverting $D(q)$. Using symbolic computations, we compute

$$J_{33} := \frac{\partial f_3(q_1, q_2, \dot{q}_1, \dot{q}_2)}{\partial \dot{q}_1} = \alpha(q_1, q_2)\dot{q}_1 + \beta(q_1, q_2)\dot{q}_2,$$

where

$$\alpha(q_1, q_2) = \frac{-4 \cdot m_2 \cdot (L_2 + L_1 \cdot \cos(q_2)) \cdot \sin(q_2)}{L_1 \cdot (-2 \cdot m_1 - m_2 + m_2 \cdot \cos(2 \cdot q_2))}$$

$$\beta(q_1, q_2) = \frac{-4 \cdot L_2 \cdot m_2 \cdot \sin(q_2)}{-2 \cdot L_1 \cdot m_1 - L_1 \cdot m_2 + L_1 \cdot m_2 \cdot \cos(2 \cdot q_2)}$$

Hence, for $q_2 \neq \pm \frac{\pi}{2}$ (modulo π), $\alpha(q_1, q_2) \neq 0$ and then

$$\sup_{\dot{q}_1} J_{33} = \infty,$$

and therefore $\left\| \frac{\partial f}{\partial x} \right\|$ cannot be uniformly bounded. More generally, if a robot model has two or more revolute joints, then it is pretty much guaranteed to fail the test we have given for global existence and uniqueness of solutions for the same reason underlined here: at least one component of the Jacobian will be linear (affine) in at least one component of the generalized velocity coordinates. ■

9.4.4 (Optional Read:) More General Notion of a Solution

Consider a rather general vector ODE that may even depend explicitly upon time,

$$\dot{x}(t) = f(x(t), t), \quad x(t_0) = x_0. \quad (9.39)$$

An ODE is an equation relating $x(t)$ and its derivative, $\dot{x}(t)$. We know that $x(t)$ is always an antiderivative of $\dot{x}(t)$. The First Fundamental Theorem of Calculus⁶ states that if $x(t)$ is an antiderivative of $x'(t)$ ($x'(t) = \dot{x}(t) = \frac{d}{dt}x(t)$ depending on the notation of the moment), the two functions must be related by

$$x(t) = x(t_0) + \int_{t_0}^t x'(\tau) d\tau \quad \text{and} \quad \frac{d}{dt} \left(x(t_0) + \int_{t_0}^t x'(\tau) d\tau \right) = x'(t). \quad (9.40)$$

We can then plug (9.39) into (9.40) and obtain that $x(t)$, $f(x(t), t)$, $x'(t)$ and $\dot{x}(t)$ are related by

$$x(t) = x(t_0) + \int_{t_0}^t \underbrace{f(x(\tau), \tau)}_{x'(\tau)} d\tau \quad \text{and} \quad \frac{d}{dt} \left(x(t_0) + \int_{t_0}^t \underbrace{f(x(\tau), \tau)}_{x'(\tau)} d\tau \right) = \underbrace{\dot{x}(t)}_{x'(t)}. \quad (9.41)$$

Therefore, the equation

$$x(t) = x_0 + \int_{t_0}^t f(x(\tau), \tau) d\tau, \quad (9.42)$$

relating $x(t)$ and an indefinite integral depending upon $x(t)$, is equivalent to the ODE (9.39), relating $x(t)$ and its derivative $\dot{x}(t)$. The interesting property of (9.42) is that the derivative of $x(t)$ does not show up explicitly, and hence functions that are only “piecewise differentiable”, instead of differentiable everywhere, can also be solutions of the ODE (9.39). **We care about this as engineers because, when we design feedback controllers for a system (e.g., a robot), we often send command signals to the closed-loop system that are only piecewise continuous; that is, the command signals have jumps in them.** Engineer to Robot: “Move forward at 1 m/s.” **The engineer sees an obstacle. Engineer screams to Robot: “Stop! Stop immediately...” as they slam the joystick in full reverse!** As engineers, we’d like to know that a solution to the differential equations that govern our robot does not cease to make sense as soon as we send a very abrupt (non-differentiable) change in the command signal.

9.5 Linear Systems of ODEs

When the right side of the ODE $\dot{x} = f(x)$ is linear in x , one would hope that it is simpler to compute analytical solutions and/or understand various properties of the solution. This is indeed the case and motivates the present section.

Homogeneous Linear Systems

Definition 9.26. For A an $n \times n$ real matrix,

$$\dot{x} = Ax, \quad x(t_0) = x_0 \quad (9.43)$$

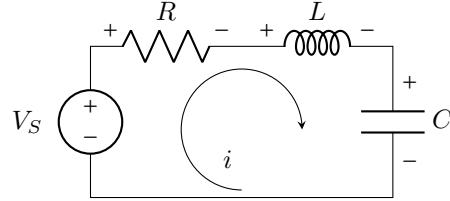
is called a **homogeneous linear system of differential equations**.

Notes: By Prop. 5.48-(a), $\frac{\partial}{\partial x} (Ax) = A$, which is independent of x . Hence, Prop. 9.21-(b) implies the global existence and uniqueness of solutions. In other words, for all $x_0 \in \mathbb{R}^n$, $t_0 \in \mathbb{R}$, there exists a solution on $[t_0, \infty)$ and the solution is unique.

9.5.1 Examples from Circuits and Mechanics

Example 9.27. Using Kirchoff’s Laws for linear RLC circuits, derive a differential equation model for the given circuit. The basics of Kirchoff’s laws will be given in the solution. For more details, you can see **Kirchhoff’s Voltage Law - KVL Circuits, Loop Rule & Ohm’s Law - Series Circuits, Physics** and **Kirchhoff’s Law, Junction & Loop Rule, Ohm’s Law - KCL & KVL Circuit Analysis - Physics**, both by the Organic Chemistry Tutor.

⁶While the Fundamental Theorems of Calculus apply to scalar-valued functions of a single variable, we can apply them to each component of a vector-valued function of a single variable, t .



Solution: V_s is a voltage source, which could be a battery or something more capable, such as a power supply. R denotes a resistor, L an inductor, and C a capacitor. The loop denotes the reference direction for the current flowing through the circuit, $i(t)$. **Kirchoff's Voltage Law** says that the sum of the voltages around any loop must equal zero. This gives us

$$V_s + V_R + V_L + V_C = 0. \quad (9.44)$$

So far, there are no derivatives in sight! The **constitutive relations** for the three circuit elements are

$$\begin{aligned} V_R &= R \cdot i \\ V_L &= L \frac{di}{dt} \\ i &= C \frac{dV_C}{dt}, \end{aligned} \quad (9.45)$$

which, in plain words, means that the voltage across a resistor is proportional to the current flowing through it, the voltage across the inductor is proportional to the time derivative of the current, and the current through the capacitor is proportional to the time derivative of the voltage across the capacitor, denoted here as $\frac{dV_C}{dt}$. Now that we have derivatives, we can build an ODE model. The rest is “algebra”.

If we assign state variables by

$$\begin{aligned} x_1 &:= V_C \\ x_2 &:= i, \end{aligned} \quad (9.46)$$

then, from line three of (9.45), we have

$$\dot{x}_1 = \frac{1}{C} x_2. \quad (9.47)$$

The equation for \dot{x}_2 comes from Kirchoff's Voltage Law as given by (9.44) in tandem with lines one and two of (9.45). Indeed, from (9.44), we obtain

$$V_L = -V_s - V_R - V_C,$$

which from (9.45) gives,

$$L \frac{di}{dt} = -V_s - r \cdot i - V_C.$$

From the definition of the state variables in (9.46), we obtain

$$\begin{aligned} L \dot{x}_2 &= -Cx_1 - Rx_2 - V_s \\ \Updownarrow \\ \dot{x}_2 &= -\frac{C}{L}x_1 - \frac{R}{L}x_2 - \frac{1}{L}V_s. \end{aligned} \quad (9.48)$$

Combining (9.47) and (9.48) gives

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} 0 & \frac{1}{C} \\ -\frac{C}{L} & -\frac{R}{L} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 0 \\ \frac{1}{L}V_s \end{bmatrix}. \quad (9.49)$$

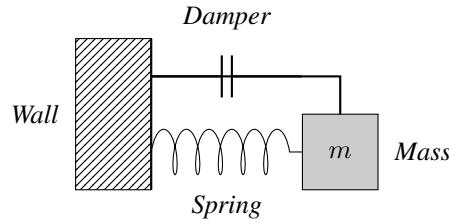
Equation (9.49) is an inhomogeneous system of linear equations due to the voltage source, V_s . If we assume $V_s \equiv 0$, we obtain

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} 0 & \frac{1}{C} \\ -\frac{C}{L} & -\frac{R}{L} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}. \quad (9.50)$$

The initial conditions for (9.50) come from the initial voltage across the capacitor and the initial current through the inductor. ■

Hence, RLC circuits give rise to linear systems of ODEs. In a similar manner, mechanical systems without “rotating elements” (aka, only have elements that move in a straight line) “mostly” give rise also to linear systems of ODEs when stiction and dead zones are ignored.

Example 9.28. Using either Newton’s Laws or Lagrange’s Equations of Motion, derive a linear differential equation model for the mass-spring-damper system below.



Solution: We assume the wall is immovable and let x be the distance of the mass from the wall, with positions to the right of the wall being positive. The velocity of the mass is denoted \dot{x} , and its Kinetic Energy is $\frac{1}{2}m\dot{x}^2$. Let the rest position of the spring be x_s (the position of the spring when it is neither compressed nor stretched). The Potential Energy of the spring is then $\frac{1}{2}k(x - x_s)^2$. The effect of the linear damper, $-b\dot{x}$, must be placed in Γ . With these assumptions, Lagrange’s equations become

$$L(x, \dot{x}) = \frac{1}{2}m\dot{x}^2 - \frac{1}{2}k(x - x_s)^2 \quad (9.51)$$

$$\frac{d}{dt} \frac{\partial L}{\partial \dot{x}} - \frac{\partial L}{\partial x} = \Gamma,$$

where k is the spring constant, b is the damping coefficient, and $\Gamma = -b\dot{x}$. Computing the indicated derivatives gives

$$\begin{aligned} \frac{d}{dt}(m\dot{x}) + k(x - x_s) &= -b\dot{x} \\ \Downarrow \\ m\ddot{x} + b\dot{x} + k(x - x_s) &= 0. \end{aligned} \quad (9.52)$$

Assigning state variables as

$$\begin{aligned} x_1 &:= x - x_s \\ x_2 &:= \dot{x}, \end{aligned} \quad (9.53)$$

gives the relationships

$$\begin{aligned} \dot{x}_1 &= x_2 \\ m\dot{x}_2 &= -kx_1 - bx_2. \end{aligned} \quad (9.54)$$

Therefore, the linear homogeneous model is

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -\frac{k}{m} & -\frac{b}{m} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}. \quad (9.55)$$

The initial conditions for (9.55) come from the initial position of the mass (relative to the rest position of the spring) and the initial velocity of the mass. ■

9.5.2 Expressing Higher-order Linear ODEs in First-order Form

In Definition 9.10, we introduced the notion of “higher order ODES of a single variable”, where the linear model takes the form

$$y^{(n)} + a_{n-1}y^{(n-1)} + a_{n-2}y^{(n-2)} + \cdots + a_1\dot{y} + a_0y = 0.$$

Such models can be written as a vector first-order system of linear ODEs by assigning state variables as follows

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{n-1} \\ x_n \end{bmatrix} = \begin{bmatrix} y \\ \dot{y} \\ \vdots \\ y^{(n-2)} \\ y^{(n-1)} \end{bmatrix}, \quad (9.56)$$

so that

$$\frac{d}{dt} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{n-1} \\ x_n \end{bmatrix} = \begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \vdots \\ \dot{x}_{n-1} \\ \dot{x}_n \end{bmatrix} = \begin{bmatrix} x_2 \\ x_3 \\ \vdots \\ x_n \\ -a_0x_1 - a_1x_2 - \cdots - a_{n-1}x_n \end{bmatrix}. \quad (9.57)$$

In matrix form, the system of ODEs becomes

$$\dot{\mathbf{x}} = \underbrace{\begin{bmatrix} 0 & 1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 1 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & 0 & 0 \\ 0 & 0 & 0 & 0 & \ddots & 1 & 0 \\ 0 & 0 & 0 & 0 & \ddots & 0 & 1 \\ -a_0 & -a_1 & -a_2 & -a_3 & \cdots & -a_{n-2} & -a_{n-1} \end{bmatrix}}_A \cdot \mathbf{x}. \quad (9.58)$$

Example 9.29. Rewrite the second-order linear mass-spring-damper model,

$$m\ddot{x} + b\dot{x} + k(x - x_s) = 0,$$

as a linear state-variable model.

Solutions: We first solve for \ddot{x} , yielding

$$\ddot{x} = -\frac{k}{m}(x - x_s) - \frac{b}{m}\dot{x}.$$

Assigning state variables as

$$\begin{aligned} x_1 &:= x - x_s \\ x_2 &:= \dot{x}, \end{aligned}$$

gives the relationships

$$\begin{aligned}\dot{x}_1 &= x_2 \\ \dot{x}_2 &= -\frac{k}{m}x_1 - \frac{b}{m}x_2.\end{aligned}$$

Therefore, the linear homogeneous model is

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -\frac{k}{m} & -\frac{b}{m} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix},$$

just as we derived in Example 9.28. ■

9.5.3 Linearization of Nonlinear ODE Models

In engineering practice, linear models predominantly arise from local approximations of nonlinear models around equilibrium points, that is, initial conditions where the system remains constant unless disturbed. The utility of linear approximations might seem counterintuitive, given our ability to compute solutions for complex nonlinear models numerically. However, often, the real challenge lies not in modeling systems but in modifying their behavior. Consider the example of a Segway: while creating and simulating its nonlinear dynamical model is straightforward, designing a feedback control system to maintain its balance is far more complex. Here, linearization becomes invaluable. By simplifying a nonlinear system to its linear counterpart around an equilibrium point, we gain a more manageable framework. This linear approximation facilitates the use of well-developed linear control techniques like PID control and state-space analysis, offering crucial insights into the system's stability and response near the equilibrium. Such simplifications are essential for effective control system design and analysis despite being local approximations that may not fully capture the system's behavior in more extreme conditions.

Equilibrium Points

Definition 9.30. A vector $x_e \in \mathbb{R}^n$ is an **equilibrium point** of the differential equation $\dot{x} = f(x)$ if $f(x_e) = 0_{n \times 1}$. Equivalently, x_e is an equilibrium point if the constant function $\varphi(t, t_0, x_e) := x_e$ is a solution of the ODE.

Notes:

- In simple terms, equilibrium points are initial conditions from which the solution of the ODE does not change; it remains constant and equal to the initial condition. For example, if a pendulum is initialized at rest in the downward hanging position, it stays there!
- Determining points x_e such that $f(x_e) = 0$ is a root finding problem; hence, Newton-Raphson is an effective means for finding equilibrium points in cases where the physics of a problem does not make them obvious.
- For a linear ODE, $\dot{x} = Ax$, x_e is an equilibrium point if, and only if, $x_e \in \text{null}(A)$. In particular, for a linear ODE, the origin is always an equilibrium point.

Example 9.31. Find equilibrium points for the following ordinary differential equations,

(a) $\frac{dv(t)}{dt} = \frac{\kappa}{m}v^2(t) - g$ (parachutist with nonlinear drag, from Example 9.8).

(b) $\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} x_2 \\ -\frac{g}{L} \cdot \sin(x_1) \end{bmatrix}$ (pendulum model from Examples 6.37 and 9.13).

(c) $\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} x_2 \\ -D^{-1}(x_1) \cdot [C(x_1, x_2) \cdot x_2 + G(x_1)] \end{bmatrix}$ (the robot equations from Method 9.12).

Solutions:

(a) $\dot{v} = \frac{\kappa}{m} v^2 - g$ **Ans.** $v_e = -\sqrt{\frac{mg}{\kappa}}$.

We set the right side to zero and solve for v . This gives $v^2 = \frac{mg}{\kappa}$. Solving for v gives $v_e = \pm\sqrt{\frac{mg}{\kappa}}$. However, only the negative velocity makes sense because the parachutist is not falling upward. This is where “engineering” judgment comes into play. If we knew nothing about the model, we’d report both solutions as potential equilibria.

Note: The equilibrium point is what we called the “terminal velocity” earlier when studying this problem. It corresponds to the speed at which the downward acceleration of gravity is in equilibrium with the upward drag force generated by the parachute.

(b)
$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} x_2 \\ -\frac{g}{L} \cdot \sin(x_1) \end{bmatrix}$$
 Ans. $x_e = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ and $x_e = \begin{bmatrix} \pi \\ 0 \end{bmatrix}$.

We set the right side to zero and solve for x . This gives $x_2 = 0$, the angular velocity of the pendulum, and $\sin(x_1) = 0$, which implies $x_1 = k\pi$, $k \in \mathbb{Z}$. $k = 0$ corresponds to the downward equilibrium and $k = 1$ gives the upward equilibrium.

(c)
$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} x_2 \\ -D^{-1}(x_1) \cdot [C(x_1, x_2) \cdot x_2 + G(x_1)] \end{bmatrix}$$
 Ans. $x_e = \left\{ \begin{bmatrix} q_e \\ 0 \end{bmatrix} \in \mathbb{R}^{2n} \mid G(q_e) = 0_{n \times 1} \right\}$

We set the right side to zero and solve for x . The top row gives $x_2 = 0$, the generalized velocity. Substituting this into the bottom row gives $-D^{-1}(x_1) \cdot G(x_1) = 0_{n \times 1} \iff G(x_1) = 0_{n \times 1}$. Because x_1 is the generalized position of the system, we’ll use q_e to denote the roots of the potential energy gradient vector, $G(q_e) = 0_{n \times 1}$. There could be no roots, one root, or multiple roots, as with any vector function.

■

Example 9.32. We consider the 2-link manipulator with springs added to both revolute joints. With this modification, and after setting the parameters to

$$\begin{bmatrix} g \\ L_1 \\ L_2 \\ m_1 \\ m_2 \\ k_1 \\ k_2 \end{bmatrix} = \begin{bmatrix} 9.81 \text{ m/s}^2 \\ 1.0 \text{ m} \\ 0.5 \text{ m} \\ 1.0 \text{ kg} \\ 3.0 \text{ kg} \\ 24 \text{ N/rad} \\ 24 \text{ N/rad} \end{bmatrix}$$

we have

$$G(q_1, q_2) = \begin{bmatrix} 14.715 \cdot \cos(q_1 + q_2) + 24q_1 + 39.24 \cdot \cos(q_1) \\ 14.715 \cdot \cos(q_1 + q_2) + 24q_2 \end{bmatrix}$$

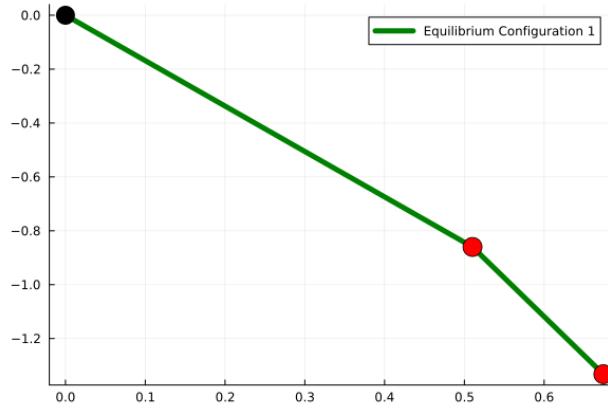
and

$$Jac_G(q_1, q_2) = \begin{bmatrix} 24.0 - 14.715 \cdot \sin(q_1 + q_2) - 39.24 \cdot \sin(q_1) & -14.715 \cdot \sin(q_1 + q_2) \\ -14.715 \cdot \sin(q_1 + q_2) & 24.0 - 14.715 \cdot \sin(q_1 + q_2) \end{bmatrix}.$$

Find the equilibrium point.

Solution: Ans. $q^* = \begin{bmatrix} -1.03536 \\ -0.20116 \end{bmatrix}$ radians.

The rest position of the springs corresponds to both links at $(\theta_1, \theta_2) = (0, 0)$, which corresponds to $p_1 = [0; 0]$ and $p_2 = [0; 0]$. With nothing to support the robot arm, gravity seeks to rotate its links downward, thereby engaging the springs. In other words, when hanging downward, the springs provide restorative (supporting) forces at the two joints. The stiffer the springs, the less motion is required by the links to be in equilibrium with gravity.



The following code implements a Newton-Raphson routine to find $G(q_1^*, q_2^*) = [0; 0]$.

```

1 using LinearAlgebra
2 include ("dyn_mod_2LinkManipulatorWithSprings.jl")
3 x0 = [pi/2, pi/4]
4 Model = dyn_mod_2LinkManipulatorWithSprings(x0, 0*x0)
5 G = Model.G
6 JacG = Model.JacG
7 k = 0
8 aTol = 1e-5
9 s = 0.1
10 xk = x0
11 while (k<1e5) && (norm(G)>aTol)
12     k = k+1
13     Delta_x = -JacG\G
14     xk = xk + s*Delta_x
15     Model = dyn_mod_2LinkManipulatorWithSprings(xk, 0*xk)
16     G = Model.G
17     JacG = Model.JacG
18 end
19 Model = dyn_mod_2LinkManipulatorWithSprings(xk, 0*xk)
20 G = Model.G
21 JacG = Model.JacG
22
23 display(G)
24 [xk x0]

```

Output

```
2-element Vector{Float64}:
-2.959527982682175e-6
 9.381671145725079e-6
```

```
2×2 Matrix{Float64}:
-1.03536  1.5708
-0.201158  0.785398
```

Linearization About an Equilibrium Point

Let $x_e \in \mathbb{R}^n$ be an **equilibrium point** of the differential equation $\dot{x} = f(x)$ and suppose that f is differentiable at x_e . Then from (5.34), the linearization of f about x_e is

$$f(x) = f(x_e) + \frac{\partial f(x)}{\partial x} \Big|_{x=x_e} (x - x_e) = \frac{\partial f(x_e)}{\partial x} \cdot (x - x_e),$$

where $f(x_e) = 0_{n \times 1}$ follows from x_e being an equilibrium point.

Definition 9.33. For a nonlinear vector ODE $\dot{x} = f(x)$, the **linearization about an equilibrium point** x_e is the linear vector ODE

$$\delta\dot{x} = \frac{\partial f(x_e)}{\partial x} \cdot \delta x, \quad (9.59)$$

where $\delta x := x - x_e$.

Notes:

- If $\delta x(t)$ is a solution to the linearized model (9.59), then $x_e + \delta x(t)$ is an **approximate solution** to the original nonlinear model.
- When $x_e = 0_{n \times 1}$, it common to write the linearized model as $\dot{x} = Ax$, for $A := \frac{\partial f(x_e)}{\partial x}$, thereby dropping the δ .
- Even when $x_e \neq 0_{n \times 1}$, it is still rather uncommon to use δx as the name of the state variables in the linear model. Some engineers abuse notation and reuse the variable x , others opt for a new variable name, such as z or χ .

Example 9.34. Linearize the following nonlinear models about the given equilibrium points.

(a) $\frac{dv(t)}{dt} = \frac{\kappa}{m} v^2(t) - g$ with $v_e = -\sqrt{\frac{mg}{\kappa}}$ (parachutist with nonlinear drag, from Example 9.8).

(b) $\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} x_2 \\ -\frac{g}{L} \cdot \sin(x_1) \end{bmatrix}$ with $x_e = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ and $x_e = \begin{bmatrix} \pi \\ 0 \end{bmatrix}$ (pendulum model from Examples 6.37 and 9.13).

Solutions:

(a) $\dot{v} = \frac{\kappa}{m} v^2 - g$ **Ans.** $\delta\dot{v} = -\left(2\sqrt{\frac{\kappa g}{m}}\right) \cdot \delta v$, where $\delta v = v - v_e$.

We identify $f(v) = \frac{\kappa}{m} v^2 - g$ and compute $\frac{df(v)}{dv} = \frac{d}{dv}\left(\frac{\kappa}{m} v^2 - g\right) = 2\frac{\kappa}{m} v$. Evaluating the “Jacobian” at $v_e = -\sqrt{\frac{mg}{\kappa}}$ gives $\frac{df(v_e)}{dv} = -2\sqrt{\frac{\kappa g}{m}}$. Hence, the linearized model is

$$\delta\dot{v} = -\left(2\sqrt{\frac{\kappa g}{m}}\right) \cdot \delta v, \quad \delta v = v - v_e.$$

Note: The scalar coefficient in front of δv is negative. Later, we’ll see that this imparts a “self-stabilizing action” to the model, meaning that if the velocity is initially greater than the equilibrium velocity, then it will decrease toward the equilibrium velocity, and conversely, if the velocity is initially lower than the equilibrium velocity, then it will increase toward the equilibrium velocity. If this were not the case, the parachutist would have to continually adjust the surface area of the parachute to regulate their speed, considerably augmenting the difficulty of the jump!

$$(b) \begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} x_2 \\ -\frac{g}{L} \cdot \sin(x_1) \end{bmatrix} \quad \text{Ans.} \quad \begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -\frac{g}{L} & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \text{ and } \begin{bmatrix} \delta \dot{x}_1 \\ \delta \dot{x}_2 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ \frac{g}{L} & 0 \end{bmatrix} \begin{bmatrix} \delta x_1 \\ \delta x_2 \end{bmatrix}.$$

We identify $f(x) = f(x_1, x_2) = \begin{bmatrix} x_2 \\ -\frac{g}{L} \cdot \sin(x_1) \end{bmatrix}$ and compute $\frac{\partial f(x_1, x_2)}{\partial x} = \begin{bmatrix} 0 & 1 \\ -\frac{g}{L} \cdot \cos(x_1) & 0 \end{bmatrix}$. Evaluating the Jacobian at the two equilibria gives

$$\frac{\partial f(0, 0)}{\partial x} = \begin{bmatrix} 0 & 1 \\ -\frac{g}{L} & 0 \end{bmatrix} \quad \text{and} \quad \frac{\partial f(\pi, 0)}{\partial x} = \begin{bmatrix} 0 & 1 \\ \frac{g}{L} & 0 \end{bmatrix}.$$

Hence, the two linear models are

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ -\frac{g}{L} & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \text{and} \quad \begin{bmatrix} \delta \dot{x}_1 \\ \delta \dot{x}_2 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ \frac{g}{L} & 0 \end{bmatrix} \begin{bmatrix} \delta x_1 \\ \delta x_2 \end{bmatrix},$$

where in the first model, because the equilibrium is at the origin, we kept x as the state variable, and in the second model, we are using $\begin{bmatrix} \delta x_1 \\ \delta x_2 \end{bmatrix} := \begin{bmatrix} x_1 - \pi \\ x_2 \end{bmatrix}$. There are no hard and fast rules governing these choices. You have a lot of freedom to do what you prefer.

Note: The two linear models differ by a minus sign in the (2,1) position of the “A-matrix”. Later, we’ll see that the “A-matrix” of the first model has purely imaginary eigenvalues $\lambda = \pm i\sqrt{\frac{g}{L}}$, indicating oscillatory behavior, while the second model has two real eigenvalues $\lambda = \pm \sqrt{\frac{g}{L}}$, with the positive eigenvalue indicating “instability”. This comports with what we know about an inverted (aka, upward pointing) pendulum: if we perturb it ever so slightly, it will fall over and not return to the upright position, while if we perturb a normal pendulum, it will oscillate. Segways, BallBots, bipedal robots, and humans each contain an (unstable) inverted pendulum that must be balanced! If this were not the case, babies would have a “self-righting” mechanism built in at birth. ■

Proposition 9.35: Linearization of the Robot Equations

We initially encountered the robot equations expressed as a set of nonlinear second-order differential equations,

$$D(q) \cdot \ddot{q} + C(q, \dot{q}) \cdot \dot{q} + G(q) = 0_{n \times 1},$$

where we are assuming here that there are no external forces acting on the robot. In Example 9.31, we determined that their equilibrium points correspond to $G(q_e) = 0_{n \times 1}$ and $\dot{q}_e = 0_{n \times 1}$. Their linearization about an equilibrium can be expressed in both second-order and first-order forms, namely,

$$D(q_e) \cdot \delta \ddot{q} + \frac{\partial G(q_e)}{\partial q} \cdot \delta q = 0_{n \times 1}, \quad (9.60)$$

where $\delta q = q - q_e$, and

$$\begin{bmatrix} \delta \dot{x}_1 \\ \delta \dot{x}_2 \end{bmatrix} = \begin{bmatrix} 0_n & I_n \\ A_{21} & 0_n \end{bmatrix} \begin{bmatrix} \delta x_1 \\ \delta x_2 \end{bmatrix}, \quad (9.61)$$

where $A_{21} = -D(q_e) \setminus \frac{\partial G(q_e)}{\partial q}$ because we know not to invert large matrices, I_n is the $n \times n$ identity matrix, $\delta x_1 = q - q_e$, and $\delta x_2 = \dot{q} - 0_{n \times 1}$.

Note: At first, it is surprising that the term $C(q, \dot{q}) \cdot \dot{q}$ does not show up in the linearization. The reason is that it consists of quadratic terms of the form $c_{ij}(q) \cdot \dot{q}_i \cdot \dot{q}_j$, which when linearized about $\dot{q}_e = 0_{n \times 1}$, $1 \leq i \leq n$, all vanish. By taking advantage of this property, the linearization of the robot equations is trivial to compute: you only have to compute one Jacobian, $\frac{\partial G(q_e)}{\partial q}$, and evaluate the mass inertia matrix at q_e .

Example 9.36. In Chapter 6.4.6, we derived Lagrange's equations for a 3-link manipulator. Using Julia, compute its linearized

model about the downward equilibrium, $q_e^d = \begin{bmatrix} -\frac{\pi}{2} \\ 0 \\ 0 \end{bmatrix}$ and the upward equilibrium $q_e^u = \begin{bmatrix} \frac{\pi}{2} \\ 0 \\ 0 \end{bmatrix}$.

Use the following function to define physical values to the parameters.

```

1 function modelParameters ()
2     g = 9.81 # m/s^2
3     L1 = 1 # m
4     L2 = 0.7
5     L3 = 0.5
6     m1 = 15 # kg
7     m2 = 10
8     m3 = 5
9     return g, L1, L2, L3, m1, m2, m3
10 end

```

Solutions:

Ans. Downward Equilibrium

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \\ \dot{x}_4 \end{bmatrix} = \begin{bmatrix} 0.0000 & 0.0000 & 0.0000 & 1.0000 & 0.0000 & 0.0000 \\ 0.0000 & 0.0000 & 0.0000 & 0.0000 & 1.0000 & 0.0000 \\ 0.0000 & 0.0000 & 0.0000 & 0.0000 & 0.0000 & 1.0000 \\ -9.8100 & 9.8100 & 0.0000 & 0.0000 & 0.0000 & 0.0000 \\ 9.8100 & -37.8386 & 7.0071 & 0.0000 & 0.0000 & 0.0000 \\ 0.0000 & 28.0286 & -36.4371 & 0.0000 & 0.0000 & 0.0000 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_4 \\ x_6 \end{bmatrix}.$$

Ans. Upward Equilibrium

$$\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \\ \dot{x}_4 \end{bmatrix} = \begin{bmatrix} 0.0000 & 0.0000 & 0.0000 & 1.0000 & 0.0000 & 0.0000 \\ 0.0000 & 0.0000 & 0.0000 & 0.0000 & 1.0000 & 0.0000 \\ 0.0000 & 0.0000 & 0.0000 & 0.0000 & 0.0000 & 1.0000 \\ 9.8100 & -9.8100 & 0.0000 & 0.0000 & 0.0000 & 0.0000 \\ -9.8100 & 37.8386 & -7.0071 & 0.0000 & 0.0000 & 0.0000 \\ 0.0000 & -28.0286 & 36.4371 & 0.0000 & 0.0000 & 0.0000 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_4 \\ x_6 \end{bmatrix}.$$

Note: The A_{21} block for the upward equilibrium is the negative of the A_{21} block for the downward equilibrium. This generalizes to three-links the observation we made for the single-link pendulum. In Example 9.37, we explore what this means for the eigenvalues of the respective linearized models.

The function `dyn_mod_3LinkManipulator(q, dq)` (see file `dyn_mod_3LinkManipulator.jl`), returns ($D=D$, $C=C$, $G=G$, $B=B$, $JacG=JacG$), where $JacG = \frac{\partial G}{\partial q}$. We will use D and $JacG$.

We only show the code for the downward equilibrium.

```

1 # Linear Model for the 3-link Manipulator
2 #
3 # bring the model into the workspace for later use
4 include ("dyn_mod_3LinkManipulator.jl")
5
6 #Build a function to define the parameters of the model
7
8 function modelParameters()
9     g = 9.81 # m/s^2
10    L1 = 1 # m
11    L2 = 0.7
12    L3 = 0.5
13    m1 = 15 # kg
14    m2 = 10
15    m3 = 5
16    return g, L1, L2, L3, m1, m2, m3
17 end
18
19 # Set equilibrium Point
20 if true
21     qe = [-pi/2, 0, 0]
22 else
23     qe = [pi/2, 0, 0]
24 end
25
26 F = dyn_mod_3LinkManipulator (qe, 0*qe)
27 D = F.D; display(D)
28 JacG = F.JacG; display(JacG)
29
30 # build the linearized model about the give equilibrium point
31 A21 = -D\JacG; A21 = cleanUp(A21); display(A21)
32 n, ~ = size(A21)
33 A = [zeros(n,n) I(n); A21 zeros(n,n) ]

```

Output

3×3 Matrix{Float64}:

68.1	25.1	5.5
25.1	12.1	3.0
5.5	3.0	1.25

3×3 Matrix{Float64}:

421.83	127.53	24.525
127.53	127.53	24.525
24.525	24.525	24.525

3×3 Matrix{Float64}:

-9.81	9.81	0.0
-------	------	-----

```

9.81 -37.8386 7.00714
0.0 28.0286 -36.4371

```

6×6 Matrix{Float64}:

```

0.0 0.0 0.0 1.0 0.0 0.0
0.0 0.0 0.0 0.0 1.0 0.0
0.0 0.0 0.0 0.0 0.0 1.0
-9.81 9.81 0.0 0.0 0.0 0.0
9.81 -37.8386 7.00714 0.0 0.0 0.0
0.0 28.0286 -36.4371 0.0 0.0 0.0

```

Example 9.37. Use Julia to compute the eigenvalues of both linearized models in Example 9.36.

$$\text{Solutions: } \begin{array}{ll} \text{Ans.} & \left[\begin{array}{c} \lambda_1 \\ \lambda_2 \\ \lambda_3 \\ \lambda_4 \\ \lambda_5 \\ \lambda_6 \end{array} \right]_d = \left[\begin{array}{c} -7.2389i \\ 7.2389i \\ -2.4528i \\ 2.4528i \\ -5.0664i \\ 5.0664i \end{array} \right] \text{ and } \text{Ans.} & \left[\begin{array}{c} \lambda_1 \\ \lambda_2 \\ \lambda_3 \\ \lambda_4 \\ \lambda_5 \\ \lambda_6 \end{array} \right]_u = \left[\begin{array}{c} -7.2389 \\ -5.0664 \\ -2.4528 \\ 2.4528 \\ 5.0664 \\ 7.2389 \end{array} \right]. \end{array}$$

Note: For the downward equilibrium, the eigenvalues are distinct and purely imaginary, which we will later show implies sustained periodic behavior. For the upward equilibrium, on the other hand, the eigenvalues are real and occur in positive-negative pairs, $\pm\lambda$. We'll see later that the positive real eigenvalues imply "instability", meaning solutions have norms that tend to infinity. Yikes!

(Optional Read:) Is there a logical explanation for why the eigenvalues change so drastically from one model to the other? Yes. There are two explanations. The first explanation is based on our physical experience with inverted pendula, as we discussed in Example 9.31-(b); one the system is disturbed, it topples over. The second and more mathematical explanation is that for the robot equations, we have

$$A = \begin{bmatrix} 0_n & I_n \\ A_{21} & 0_n \end{bmatrix}.$$

When A_{21} is an invertible $n \times n$ matrix, one can show that

$$\det(\lambda I_{2n} - A) = \det \begin{bmatrix} \lambda I_n & -I_n \\ -A_{21} & \lambda I_n \end{bmatrix} = \det(\lambda^2 I_n - A_{21}).$$

Therefore, the eigenvalues of A are the plus and minus square roots of the eigenvalues of A_{21} . That's pretty cool! The sign change on the entries of A_{21} for upward and downward equilibria changes purely imaginary eigenvalues into purely real eigenvalues. The reason behind the sign change of the entries of A_{21} is more technical. It arises because in the upward equilibrium, the potential energy is a maximum, while at the lower equilibrium, it is a minimum. Indeed, because $G(q) = \nabla PE(q)$, $\frac{\partial G(q)}{\partial q}$ is the Hessian of the potential energy. The Hessian is a positive definite matrix at a local minimum and a negative definite matrix at a local maximum. Linear Algebra is very powerful!

Code for the Downward Equilibrium

```

1 E=eigen(A)
2 evals = cleanUp(E.values); display(evals)

```

Output

```
6-element Vector{ComplexF64}:
0.0 - 7.23885032225449im
0.0 + 7.23885032225449im
0.0 - 2.452784229773624im
0.0 + 2.452784229773624im
0.0 - 5.066419822703628im
0.0 + 5.066419822703628im
```

CleanUp Function that handles complex numbers

```
1 function cleanUp (A, tol=1e-10)
2     B = copy (A)
3     for i in eachindex (B)
4         if isa (B [i], Complex)
5             # Clean up the real part
6             real_part = abs (real (B [i])) < tol ? 0.0 : real (B [i])
7             # Clean up the imaginary part
8             imag_part = abs (imag (B [i])) < tol ? 0.0 : imag (B [i])
9             # Reconstruct the complex number
10            B [i] = real_part + imag_part * im
11        else
12            # Original cleanup for non-complex numbers
13            if abs (B [i]) < tol
14                B [i] = 0.0
15            end
16        end
17    end
18    return B
19 end
```

Output

```
cleanUp (generic function with 2 methods)
```

9.5.4 The Matrix Exponential

We recall that for $a \in \mathbb{R}$, the Maclaurin/Taylor Expansions for e^{at} and $e^{a(t-t_0)}$ are

$$e^{at} = \sum_{k=0}^{\infty} a^k \cdot \frac{t^k}{k!} = 1 + a \cdot t + a^2 \cdot \frac{t^2}{2!} + \dots + a^k \cdot \frac{t^k}{k!} + \dots \text{ and}$$

$$e^{a(t-t_0)} = \sum_{k=0}^{\infty} a^k \cdot \frac{(t-t_0)^k}{k!} = 1 + a \cdot (t-t_0) + a^2 \cdot \frac{(t-t_0)^2}{2!} + \dots + a^k \cdot \frac{(t-t_0)^k}{k!} + \dots, \quad (9.62)$$

respectively⁷. In the following, we initially take $t_0 = 0$, so that the Maclaurin and Taylor expansions are the same.

Matrix Exponential and Exponential Function

In the following, I_n denotes the $n \times n$ identity matrix.

⁷The definition should be written as $e^{at} = 1 + \lim_{N \rightarrow \infty} \sum_{k=1}^N a^k \cdot \frac{t^k}{k!}$ so that the need for a limit is clear, and so that we do not have to use $0^0 = 1$. Normally, these “niceties” are assumed to be understood by the reader.

Definition 9.38. For A an $n \times n$ real matrix,

$$e^A := \sum_{k=0}^{\infty} A^k \cdot \frac{1}{k!} = I_n + A + A^2 \cdot \frac{1}{2!} + \cdots + A^k \cdot \frac{1}{k!} + \cdots \quad (9.63)$$

is called the **matrix exponential**. It is also denoted as $\exp(A)$.

Note: If $A = 0_n$, the $n \times n$ zero matrix, then $\exp(0_n) := I_n$. **Warning:** This shows that the matrix exponential is not the exponential of each term of the matrix. In other symbols, $[e^A]_{ij} \neq \exp(a_{ij})$. If it were, then $\exp(0_n)$ would equal an $n \times n$ matrix filled with ones, instead of being an identity matrix.

Definition 9.39. For A an $n \times n$ real matrix and $t \in \mathbb{R}$,

$$e^{At} := \sum_{k=0}^{\infty} A^k \cdot \frac{t^k}{k!} = I_n + A \cdot t + A^2 \cdot \frac{t^2}{2!} + \cdots + A^k \cdot \frac{t^k}{k!} + \cdots \quad (9.64)$$

is called the **matrix exponential function**. It is also denoted as $\exp(At)$.

Notes:

- The matrix exponential function is simply the matrix exponential evaluated at the matrix $A \cdot t = t \cdot A$, which depends upon time.
- One is often lazy and refers to e^{At} as the matrix exponential instead of the matrix exponential “function” because context usually makes the true meaning very clear.
- A more formal definition of the matrix exponential is $e^A := I_n + \lim_{N \rightarrow \infty} \sum_{k=1}^N A^k \cdot \frac{1}{k!}$, making the need to evaluate a limit clear. Using techniques from Michigan’s Math 451, the existence of the limit can be established for ALL $n \times n$ matrices with entries in \mathbb{R} or \mathbb{C} .

Example 9.40. Using the series definition, compute e^{At} for

$$(a) A_1 = \begin{bmatrix} a_{11} & 0 \\ 0 & a_{22} \end{bmatrix},$$

$$(b) A_2 = \begin{bmatrix} \lambda_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \lambda_n \end{bmatrix}, \text{ and}$$

$$(c) A_3 = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}.$$

Solutions: In each case, we need to find a pattern that makes evaluating the infinite sum feasible. The first two problems are remarkably easy.

$$(a) A_1 = \begin{bmatrix} a_{11} & 0 \\ 0 & a_{22} \end{bmatrix} \xrightarrow{\textcolor{blue}{\Rightarrow}} e^{A_1 t} = \begin{bmatrix} e^{a_{11}t} & 0 \\ 0 & e^{a_{22}t} \end{bmatrix}.$$

Before applying the definition, we note that raising a diagonal matrix to a power gives a diagonal matrix with entries raised to the same power, that is,

$$\underbrace{\begin{bmatrix} a_{11} & 0 \\ 0 & a_{22} \end{bmatrix}}_{k\text{-times}}^k = \underbrace{\begin{bmatrix} a_{11} & 0 \\ 0 & a_{22} \end{bmatrix} \times \cdots \times \begin{bmatrix} a_{11} & 0 \\ 0 & a_{22} \end{bmatrix}}_{k\text{-times}} = \begin{bmatrix} (a_{11})^k & 0 \\ 0 & (a_{22})^k \end{bmatrix},$$

as is easily checked for $k = 2$ and $k = 3$. The general case follows by induction (which would be good practice). Using this result, we have,

$$\begin{aligned} e^{A_1 t} &= I_2 + \begin{bmatrix} a_{11} & 0 \\ 0 & a_{22} \end{bmatrix} \cdot t + \begin{bmatrix} (a_{11})^2 & 0 \\ 0 & (a_{22})^2 \end{bmatrix} \cdot \frac{t^2}{2!} + \cdots + \begin{bmatrix} (a_{11})^k & 0 \\ 0 & (a_{22})^k \end{bmatrix} \cdot \frac{t^k}{k!} + \cdots \\ &= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + \begin{bmatrix} a_{11}t & 0 \\ 0 & a_{22}t \end{bmatrix} + \begin{bmatrix} (a_{11})^2 \frac{t^2}{2!} & 0 \\ 0 & (a_{22})^2 \frac{t^2}{2!} \end{bmatrix} + \cdots + \begin{bmatrix} (a_{11})^k \frac{t^k}{k!} & 0 \\ 0 & (a_{22})^k \frac{t^k}{k!} \end{bmatrix} + \cdots \\ &= \begin{bmatrix} e^{a_{11}t} & 0 \\ 0 & e^{a_{22}t} \end{bmatrix}. \end{aligned}$$

$$(b) A_2 = \begin{bmatrix} \lambda_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \lambda_n \end{bmatrix} \xrightarrow{\textcolor{blue}{\Rightarrow}} e^{A_2 t} = \begin{bmatrix} e^{\lambda_1 t} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & e^{\lambda_n t} \end{bmatrix}$$

The calculations follow exactly the same steps as part (a) and are left to the learner.

$$(c) A_3 = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix} \xrightarrow{\textcolor{blue}{\Rightarrow}} e^{A_3 t} = \begin{bmatrix} \cos(t) & \sin(t) \\ -\sin(t) & \cos(t) \end{bmatrix}.$$

Let's pause for just a moment and appreciate that a matrix exponential of a real-valued matrix can oscillate, while scalar exponentials of real variables are monotonic. Later, in Definition 9.52 and Prop. 9.53, we'll discover that this remarkable feature comes from Euler's formula, $e^{i\theta} = \cos(\theta) + i\sin(\theta)$, or, if we replace θ with ωt , $e^{i\omega t} = \cos(\omega t) + i\sin(\omega t)$.

To reduce clutter, we replace A_3 with A in the following.

The pattern occurring in the powers of A is more "subtle" this time. The first few powers are straightforward: by brute force multiplication, we obtain,

$$A^0 = I_2, \quad A = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}, \quad A^2 = \begin{bmatrix} -1 & 0 \\ 0 & -1 \end{bmatrix}, \quad A^3 = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}, \quad A^4 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = I_2,$$

where we notice that A^4 brings us back to where we started, the identity matrix. That is the pattern we are seeking!

We note immediately that

$$A^{4k} = \underbrace{A^4 \times \cdots \times A^4}_{k \text{ times}} = I_2 \times \cdots \times I_2 = I_2.$$

Next,

$$A^{4k+0} = I_2, \quad A^{4k+1} = A = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}, \quad A^{4k+2} = A^2 = \begin{bmatrix} -1 & 0 \\ 0 & -1 \end{bmatrix}, \quad A^{4k+3} = A^3 = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix},$$

which is kind of amazing the first time you see it.

To compute e^{At} , we need to substitute the above results into our power series for the matrix exponential,

$$e^{At} = I_2 + At + A^2 \frac{t^2}{2!} + A^3 \frac{t^3}{3!} + A^4 \frac{t^4}{4!} + A^5 \frac{t^5}{5!} \dots$$

By doing so, we obtain

$$[e^{At}]_{11} = 1 + 0 - \frac{t^2}{2} + 0 + \frac{t^4}{4!} + \cdots = \cos(t)$$

$$[e^{At}]_{12} = 0 + t + 0 - \frac{t^3}{3!} + 0 + \frac{t^5}{5!} + \cdots = \sin(t)$$

$$[e^{At}]_{21} = 0 - t + 0 + \frac{t^3}{3!} + 0 - \frac{t^5}{5!} + \cdots = -\sin(t)$$

$$[e^{At}]_{22} = 1 + 0 - \frac{t^2}{2} + 0 + \frac{t^4}{4!} + \cdots = \cos(t),$$

after recognizing the Maclaurin/Taylor expansions for cosine and sine. Therefore,

$$e^{At} = \begin{bmatrix} \cos(t) & \sin(t) \\ -\sin(t) & \cos(t) \end{bmatrix}.$$

■

9.5.5 Software Tools

As you should probably expect by now, there are software tools to help you with computing both e^A (the matrix exponential) and e^{At} (the matrix exponential function). We illustrate these next.

```

1 using LinearAlgebra
2 # Mass spring damper system
3 m = 5 # kg
4 k = 2 # N/m
5 b = 3 # Ns/m (Newton seconds per meter)
6 A = [0 1; -k/m -b/m]
7 expA = exp(A)

```

Output

```

2×2 Matrix{Float64}:
 0.839867  0.703132
 -0.281253  0.417988

```

Below is a comparison between the builtin function and a highly truncated implementation of the power series. You can see that the power series converges rapidly. This rapid convergence means that a finite number of terms well approximates the matrix exponential. We will differentiate with respect to time the matrix exponential **function** in the next subsection, using an argument based on the matrix exponential being well approximated by a finite number of terms. Its formal proof relies on analytically proving the absolute convergence of the Maclaurin/Taylor expansions for the matrix exponential, a step we will not take in this course. Bonjour, Math 451 Advanced Calculus!

```

1 # Compare to the builtin exp function
2 N=9
3 expA2 = I(2)
4 AkOverKfactorial = I(2)
5 for k = 1:N
6     AkOverKfactorial = A*AkOverKfactorial/k
7     expA2 = expA2 + AkOverKfactorial
8 end
9 [expA expA2]
```

Output

```
2×4 Matrix{Float64}:
 0.839867  0.703132  0.839867  0.703132
 -0.281253  0.417988  -0.281253  0.417988
```

Note: In Matlab, you would use `expA = expm(A)`, where the “m” appended to “exp” tells Matlab to compute the matrix exponential. Continuing with Matlab, the command, `expA = exp(A)` would compute the exponential of each term of the matrix, which is very much not the matrix exponential. Because Julia uses **broadcasting** to apply a scalar function to a vector- or matrix-valued object, it does not need the `expm()` command.

```

1 # Comparison of exp(A) and exp.(A)
2 [exp(A) exp.(A)] # Not the same
```

Output

```
2×4 Matrix{Float64}:
 0.839867  0.703132  1.0      2.71828
 -0.281253  0.417988  0.67032  0.548812
```

For large square matrices, the matrix exponential can be computed numerically. **The symbolic computation of the matrix exponential function, on the other hand, is limited to “smallish matrices with simple structure”. In practice, this tends not to be an issue, because of eigenvalues and eigenvectors, as we’ll highlight in the next section.**

We continue with “modest” results on symbolic calculations applied to Example 9.40.

```

1 using SymPy
2 #
3 @vars a11 a22 t
4 A1 = [a11 0; 0 a22]
5 A1t = A1*t # Define a 2x2 square matrix
6 expA1t =exp(A1t) # Compute the matrix exponential function
```

Output

$$\begin{bmatrix} e^{a11 \cdot t} & 0.0 \\ 0.0 & e^{a22 \cdot t} \end{bmatrix}$$

```

1 using SymPy
2 #
3 @vars lam1 lam2 lam3 lam4 t
4 A2 = Diagonal([lam1, lam2, lam3, lam4])
5 A2t = A2*t # Define a 2x2 square matrix
6 expA2t =exp(A2t) # Compute the matrix exponential function

```

Output

$$\begin{bmatrix} e^{lam1 \cdot t} & 0.0 & 0.0 & 0.0 \\ 0.0 & e^{lam2 \cdot t} & 0.0 & 0.0 \\ 0.0 & 0.0 & e^{lam3 \cdot t} & 0.0 \\ 0.0 & 0.0 & 0.0 & e^{lam4 \cdot t} \end{bmatrix}$$

```

1 using SymPy
2 #
3 @vars t
4 A3 = [0 1; -1 0]
5 A3t = A3*t # Define a 2x2 square matrix
6 expA3t =exp(A3t) # Compute the matrix exponential function

```

Output

$$\begin{bmatrix} \frac{e^{i \cdot t}}{2.0} + \frac{e^{(-i) \cdot t}}{2.0} & \frac{(-i) \cdot e^{i \cdot t}}{2.0} + \frac{i \cdot e^{(-i) \cdot t}}{2.0} \\ \frac{i \cdot e^{i \cdot t}}{2.0} - \frac{i \cdot e^{(-i) \cdot t}}{2.0} & \frac{e^{i \cdot t}}{2.0} + \frac{e^{(-i) \cdot t}}{2.0} \end{bmatrix}$$

You must know Euler's Formula from Chapter 1.11 to understand the answer! Moreover, for $\sin(t)$, you must recognize that

$$\frac{(-i) \cdot e^{i \cdot t}}{2.0} + \frac{i \cdot e^{(-i) \cdot t}}{2.0} = \frac{e^{i \cdot t} - e^{-i \cdot t}}{2i}$$

because $\frac{1}{i} = \frac{1}{i} \cdot \frac{i}{i} = -i$. Currently, applying `Sympy.simplify` to the answer does not help.

9.5.6 Solutions to Linear ODEs

Proposition 9.41: Derivative of the Matrix Exponential Function

For all $t \in \mathbb{R}$ and all $n \times n$ (constant) matrices,

$$\frac{d}{dt} e^{A \cdot t} = A \cdot e^{A \cdot t} = e^{A \cdot t} \cdot A. \quad (9.65)$$

Example 9.42. (Optional Read:) Give a non-rigorous but mostly persuasive argument that (9.65) holds.

Solution: The obstruction to being rigorous is that we do not have enough math to understand when the derivative of an infinite sum is the infinite sum of the derivatives. Hence, we approximate the matrix exponential by a finite sum,

$$e^{At} \approx e_N^{At} := I_n + \sum_{k=1}^N A^k \cdot \frac{t^k}{k!} = I_n + A \cdot t + A^2 \cdot \frac{t^2}{2!} + A^3 \cdot \frac{t^3}{3!} + \dots + A^{N-1} \cdot \frac{t^{N-1}}{(N-1)!} + A^N \cdot \frac{t^N}{N!},$$

where N sets the number of terms in the approximation. We can now use the fact that the derivative of a finite sum is the finite sum of the derivatives.

Computing derivatives term by term gives

$$\frac{d}{dt} I_n = 0_n \text{ and}$$

$$\begin{aligned} \frac{d}{dt} \left(A^k \cdot \frac{t^k}{k!} \right) &= A^k \cdot k \cdot \frac{t^{k-1}}{k!} \\ &= A^k \cdot \frac{t^{k-1}}{(k-1)!} \\ &= A \left(A^{k-1} \cdot \frac{t^{k-1}}{(k-1)!} \right) = \left(A^{k-1} \cdot \frac{t^{k-1}}{(k-1)!} \right) A, \end{aligned}$$

where we note that A can be factored to either side.

Using the above result,

$$\begin{aligned} \frac{d}{dt} \underbrace{\left(I_n + \sum_{k=1}^N A^k \cdot \frac{t^k}{k!} \right)}_{\approx e^{At}} &= 0_n + \sum_{k=1}^N A \left(A^{k-1} \cdot \frac{t^{k-1}}{(k-1)!} \right) \\ &= A \sum_{k=1}^N \left(A^{k-1} \cdot \frac{t^{k-1}}{(k-1)!} \right) \\ &= A \cdot \underbrace{\left(I_n + \sum_{k=1}^{N-1} \frac{t^k}{k!} \right)}_{\approx e^{At}}. \end{aligned}$$

Once again, we note that when A was factored out, it could have gone to either side. This completes our “mostly” persuasive argument.

Note: This is the same “proof” offered in EECS 560 Linear Systems Theory because it does not assume real analysis. ■

Proposition 9.43: Explicit Solution to a Linear System of ODEs

The unique solution to $\dot{x} = A \cdot x$, $x(t_0) = x_0$ is

$$\varphi(t, t_0, x_0) := e^{A(t-t_0)} \cdot x_0. \quad (9.66)$$

Note: It is perfectly fine to denote the solution by $\varphi(t) = e^{A(t-t_0)} \cdot x_0$, dropping the explicit dependence on the initial time, t_0 , and the initial state, x_0 , or to use $\varphi(t, x_0) = e^{At} \cdot x_0$, when $t_0 = 0$. You are in control of your own notation when communicating results to your colleagues.

Proof: All the conditions of Def. 9.16 are met; indeed:

- Prop. 9.41 states that the matrix exponential is differentiable; it follows that it must be continuous (meaning that each of the entries of the matrix exponential function is continuous in t).
- $\varphi(t_0, t_0, x_0) := e^{A(t_0-t_0)} \cdot x_0 = I_n \cdot x_0 = x_0$ (satisfies the initial condition).
- $\frac{\partial}{\partial t} \varphi(t, t_0, x_0) := \left(\frac{\partial}{\partial t} e^{A(t-t_0)} \right) \cdot x_0 = (A \cdot e^{A(t-t_0)}) \cdot x_0 = A \cdot e^{A(t-t_0)} \cdot x_0 = A \cdot \varphi(t, t_0, x_0)$ (satisfies the differential equation).

We already noted in Def. 9.16 why the solution is unique. Indeed, by Prop. 5.48-(a), $\frac{\partial}{\partial x} (Ax) = A$ is independent of x . Hence, Prop. 9.21-(b) implies that for all $x_0 \in \mathbb{R}^n$, $t_0 \in \mathbb{R}$, there exists a solution on $[t_0, \infty)$ and the solution is unique. ■

One rarely uses (9.66) to compute the solution to a real engineering problem. Instead, one uses solvers, such as "DifferentialEquations.jl". Equation (9.66) is used to deduce GENERAL PROPERTIES of SOLUTIONS to linear ODES, while a solver is used to COMPUTE SOLUTIONS for GIVEN INITIAL CONDITIONS.

Example 9.44. Compute the solution to the mass-spring-damper system in Example 9.28, using two methods:

- A numerical solver, such as *DifferentialEquations.jl*
- Prop. 9.43.

Take the parameters as mass $m = 5 \text{ Kg}$, spring constant $k = 2 \text{ N/m}$ (Newtons per meter), and damping coefficient $b = 3 \text{ Ns/m}$ (Newton seconds per meter).

Solution:

```

1 # Import useful packages
2 using DifferentialEquations
3 using Plots
4 using LaTeXStrings
5 using LinearAlgebra # for the backslash command
6 gr()
7
8 # Mass spring damper system
9 m = 5 # kg
10 k = 2 # N/m
11 b = 3 # Ns/m (Newton seconds per meter)
12 A = [0 1; -k/m -b/m]
13
14 # Define the ODE
15 function MassSpringDamper(du, x, A, t)
16     du .= A*x
17 end
18
19
20 # Set the initial condition as a vector
21 x0=[1; 0.0]
22
23 # Set the time interval
24 T = (0.0, 20)
25
26 # Setup the ODE problem
27 params = A
28 problem = ODEProblem(MassSpringDamper, x0, T, params)
29
30 # solve the ODE problem using the Runge Kutta Tsitouras 5/4 Integrator
31 sol = solve(problem, Tsit5())
32
33 # plot the continuous solution
34 p1 = plot(sol, lw=3, guidefont=20, xlabel=L"t", ylabel=L"x (x, v)",
35             label=[ "x" "v" ], legend=:topright)
36
37
38 # Discrete time vector
39 t=0:dt:T[2]
40
41 # use a for loop to compute the matrix exponential for A*t[k] and

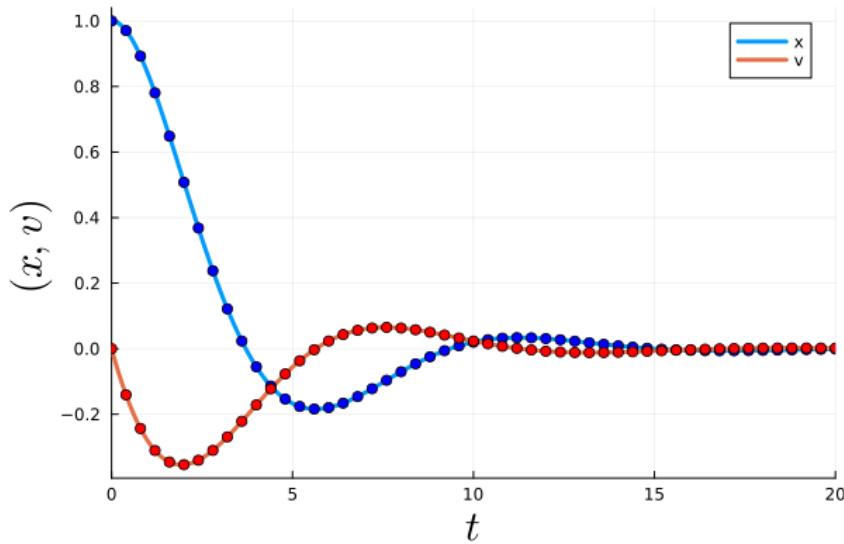
```

```

42 # then multiply by the initial condition to get the solution at t[k]
43 N = length(t)
44 x = zeros(2, N)
45 for k = 1:length(t)
46     x[:,k] = exp(A*t[k]) * x0
47 end
48
49 # Overlay the discrete-time solution with the continuous-time solution
50 scatter!(t, x[1,:], label = false, color=:blue)
51 scatter!(t, x[2,:], label = false, color=:red)

```

Output The solid lines are the solutions from the numerical solver, while the dots are the points computed by evaluating the matrix exponential function at $t_k = k \cdot 0.4$ seconds, for $0 \leq k \leq 50$. They are indistinguishable.



Proposition 9.45: (Optional Read:) Useful Properties of the Matrix Exponential

Here we highlight the similarities and differences between the matrix exponential function and the scalar exponential function. For A an $n \times n$ matrix,

- (a) $\det(e^A) = e^{\text{trace}(A)}$, where $\text{trace}(A)$ is the sum of its diagonal elements, and
- (b) $e^{A(s+t)} = e^{As} \cdot e^{At}$ for all real numbers s and t .

The first property says that, like the scalar exponential function, the matrix exponential is always invertible. The second property is similar to $e^{as} \cdot e^{at} = e^{a(s+t)}$, though it has a more profound interpretation in terms of solutions of linear ODEs. In particular, if we define $\varphi(t, t_0, x_0) := e^{A(t-t_0)} \cdot x_0$, then

$$\begin{aligned}
\varphi(s+t, t_0, x_0) &:= e^{A(s+t-t_0)} \cdot x_0 \\
&= e^{As} \cdot e^{A(t-t_0)} \cdot x_0 \\
&= e^{As} \cdot \varphi(t, t_0, x_0).
\end{aligned}$$

In other words, multiplying on the left by the matrix exponential $e^{A(t-t_0)}$ transitions the solution of $\dot{x} = Ax, x(t_0) = x_0$ from the value of the solution at the initial time, t_0 , to the value of the solution at time t , and multiplying on the left again by the matrix exponential e^{As} transitions the solution from time t to time to $s+t$.

Property (b) implies the following two additional properties,

- (c) $(e^{At})^{-1} = e^{-At}$ for all real numbers t , and
- (d) $(e^{At_d})^k = e^{A \cdot k \cdot t_d}$ for all real numbers t_d and integers k .

A key difference between the scalar exponential function and the matrix exponential function involves $e^{(A+B)t}$, for B an $n \times n$ matrix.

- (e) $e^{(A+B)t} = e^{At} \cdot e^{Bt} \iff A \cdot B - B \cdot A = 0_n$, which in words says the matrix exponential of a sum of matrices equals the product of the matrix exponentials if, and only if, A and B commute. Since most matrices do not commute, it is generally false that the matrix exponential of a sum is equal to the product of the matrix exponentials. Of course, if the matrices are 1×1 , meaning they are scalars, then they also commute.

Note: Even in graduate engineering courses, we typically do not cover the proofs of the above properties. Their main steps are related to proof by induction and binomial expansion. If you do want the proofs, we suggest [Matrix exponential](#) by Wikipedia. Another source is [Notes on the Matrix Exponential and Logarithm](#) by the Santa Cruz Institute for Particle Physics. Finally, some interesting comments are available on the Mathematical Stack Exchange, [Proof that \$\exp\(A+B\) = \exp\(A\)\exp\(B\)\$ iff \$AB = BA\$](#) .

9.5.7 Eigenvalues and Eigenvectors to the Rescue

Punchline for this Section: We've given several hints that eigenvalues are important. In this Section, we confirm it! From the eigenvalues of A , we can determine whether all solutions of $\dot{x} = Ax$ remain bounded and converge to zero as t goes to infinity, there are sustained oscillations, or some solutions tend to infinity. The reason behind this phenomenal fact? **If λ is an eigenvalue of A with eigenvector v , then $e^{\lambda t}$ is an eigenvalue of e^{At} with corresponding eigenvector v .** When λ is real, it is easy to understand whether $e^{\lambda t}$ converges to zero or explodes. Thanks to Euler's Formula, we can also handle the case of complex eigenvalues.

Review of Complex Numbers

If you wish a more extensive review of complex numbers, see ROB 101 *Computational Linear Algebra*, Appendix A.1. Here, we recall just a few facts. For $\alpha, \beta \in \mathbb{R}$, and $\mathbf{i}^2 = -1$,

$$z := \alpha + \mathbf{i}\beta$$

is a **complex number**. Its **complex conjugate** is

$$z^* := \alpha - \mathbf{i}\beta,$$

where the sign is flipped on the imaginary part. A **complex vector** is simply a vector with entries that are complex numbers. Hence, it can be written as

$$v := v_R + \mathbf{i}v_I,$$

where $v_R, v_I \in \mathbb{R}^n$; v_R is called the **real part** of the vector and v_I is called the **imaginary part** of the vector. It is important to understand that the imaginary part DOES NOT include \mathbf{i} , the imaginary unit number, and thus v_I is a real vector. The complex conjugate of the vector is

$$v^* := v_R - \mathbf{i}v_I,$$

where the sign is flipped on the imaginary part. Finally, the **magnitude** of the complex number $z^* := \alpha + \mathbf{i}\beta$ is defined to be

$$|z| := \sqrt{\alpha^2 + \beta^2} = \sqrt{z \cdot z^*}.$$

The **(Euclidean) norm of a complex vector** is

$$\|v\| := \sqrt{(v^*)^\top \cdot v} = \sqrt{\sum_{i=1}^n |v_i|^2},$$

where v_i are the components of v . Just as with norms of real vectors,

$$\|\alpha v\| = |\alpha| \cdot \|v\|,$$

for all $\alpha \in \mathbb{C}$ and $v \in \mathbb{C}^n$.

Eigenvalues and Eigenvectors of Matrices

We next recall from ROB 101 *Computational Linear Algebra*, Chapter 10.3 and Appendix A.2, the definition of eigenvalues and eigenvectors for $n \times n$ matrices.

Definition of Eigenvalues and Eigenvectors

Definition 9.46. Let A be an $n \times n$ matrix with real or complex coefficients. A scalar $\lambda \in \mathbb{C}$ is an **eigenvalue** of A , if there exists a non-zero vector $v \in \mathbb{C}^n$ such that $Av = \lambda v$. Any such vector v is called an **eigenvector** associated with λ .

Notes:

- Eigenvectors are not unique. If v is an eigenvector, then so is $\alpha \cdot v$ for any nonzero constant α .
- Most Linear Algebra packages normalize the eigenvector to have norm one.
- If A is a real matrix and $\lambda \in \mathbb{R}$, then we take $v \in \mathbb{R}^n$.
- If A is a real matrix and $\lambda \in \mathbb{C}$ satisfies $Av = \lambda v$, then λ^* , the complex conjugate of λ , is also an eigenvalue of A with corresponding eigenvector, v^* . In particular, $Av^* = \lambda^* v^*$. We say the eigenvalues and eigenvectors occur in **complex conjugate pairs**.

This next example helps to understand why eigenvalues tell us so much about solutions to linear systems of ODEs,

$$\dot{x} = Ax, \quad x(t_0 = 0) = x_0,$$

in particular, when the solutions converge to the origin, oscillate, or blow up.

Example 9.47. (Must Read:) The 2×2 matrix $A = \begin{bmatrix} a & \omega \\ -\omega & a \end{bmatrix}$ has eigenvalues $\lambda_1 = a + i\omega$ and $\lambda_2 = a - i\omega$. Use SymPy to compute e^{At} . On the basis of e^{At} , interpret the qualitative behavior of the solutions to the linear ODE $\dot{x} = Ax$ when the real parts of the eigenvalues are negative, zero, and positive.

Solution:

```

1 using SymPy
2
3 @vars a w t
4 A = [a w; -w a]*t # Define a 2x2 square matrix
5 expA = A.exp() # Another way to compute the matrix exponential

```

Output

$$\begin{bmatrix} \frac{e^{a \cdot t - i \cdot t \cdot \omega}}{2.0} + \frac{e^{a \cdot t + i \cdot t \cdot \omega}}{2.0} & \frac{i \cdot e^{a \cdot t - i \cdot t \cdot \omega}}{2.0} - \frac{i \cdot e^{a \cdot t + i \cdot t \cdot \omega}}{2.0} \\ \frac{(-i) \cdot e^{a \cdot t - i \cdot t \cdot \omega}}{2.0} + \frac{i \cdot e^{a \cdot t + i \cdot t \cdot \omega}}{2.0} & \frac{e^{a \cdot t - i \cdot t \cdot \omega}}{2.0} + \frac{e^{a \cdot t + i \cdot t \cdot \omega}}{2.0} \end{bmatrix} = \begin{bmatrix} e^{at} \cos(\omega t) & e^{at} \sin(\omega t) \\ -e^{at} \sin(\omega t) & e^{at} \cos(\omega t) \end{bmatrix}.$$

The simplification on the right was done by hand, based on

$$\begin{aligned} \frac{e^{a \cdot t - i \cdot t \cdot \omega}}{2.0} + \frac{e^{a \cdot t + i \cdot t \cdot \omega}}{2.0} &= e^{at} \left(\frac{e^{-i \cdot t \cdot \omega} + e^{i \cdot t \cdot \omega}}{2} \right) \\ &= e^{at} \left(\frac{\cos(t \cdot \omega) - i \sin(t \cdot \omega) + \cos(t \cdot \omega) + i \sin(t \cdot \omega)}{2} \right) \\ &= e^{at} \cos(\omega t) \end{aligned}$$

$$\begin{aligned} \frac{i \cdot e^{a \cdot t - i \cdot t \cdot \omega}}{2.0} - \frac{i \cdot e^{a \cdot t + i \cdot t \cdot \omega}}{2.0} &= e^{at} \left(\frac{i e^{-i \cdot t \cdot \omega} - i e^{i \cdot t \cdot \omega}}{2} \right) \\ &= e^{at} \left(\frac{i \cos(t \cdot \omega) + \sin(t \cdot \omega) - i \cos(t \cdot \omega) + \sin(t \cdot \omega)}{2} \right) \\ &= e^{at} \sin(\omega t). \end{aligned}$$

Thanks to eigenvalues and eigenvectors, **complex exponentials**, aka exponentials of the form, $e^{(a+i\omega)t}$, play a fundamental role in understanding the solutions of linear systems of ODEs.

For the present problem, we note that for all $x_0 \neq 0_{2 \times 1}$, the solution of $\dot{x} = Ax$, $x(t_0 = 0) = x_0$,

$$\varphi(t, t_0, x_0) := \begin{bmatrix} e^{at} \cos(\omega t) & e^{at} \sin(\omega t) \\ -e^{at} \sin(\omega t) & e^{at} \cos(\omega t) \end{bmatrix} \cdot \begin{bmatrix} x_{01} \\ x_{02} \end{bmatrix},$$

- converges to zero if $a < 0$ (decaying exponentially) when eigenvalues have a negative real part;
- oscillates forever if $a = 0$ (no exponential term) when eigenvalues are purely imaginary; and
- blows up if $a > 0$ (exploding exponentially) when eigenvalues have a positive real part.

■

Example 9.48. Use Julia to compute the eigenvalues and eigenvectors of the 3×3 matrix

$$A := \begin{bmatrix} 1.3916 & -0.4173 & -1.0494 \\ -0.4823 & -0.4498 & -0.2974 \\ 1.1709 & -0.2278 & 1.9801 \end{bmatrix}.$$

Moreover, verify the following properties:

- The eigenvalues and eigenvectors occur in complex conjugate pairs.
- $Av = \lambda v$.
- The eigenvectors are linearly independent and hence form a basis for \mathbb{C}^3 .
- If $v = v_R + i v_I$ is a complex eigenvector, then v and its complex conjugate v^* can be replaced by v_R and v_I , the real and imaginary parts of v , to build a basis for \mathbb{R}^3 .

Solutions:

- (a) The eigenvalues and eigenvectors occur in complex conjugate pairs.

```
1 using LinearAlgebra
2
3 A = [1.39164 -0.417337 -1.04938
4      -0.482258 -0.449806 -0.297397
5       1.1709   -0.2278    1.98006]
6
7 E = eigen(A)
```

Output

```
Eigen{ComplexF64, ComplexF64, Matrix{ComplexF64}, Vector{ComplexF64}}
values:
3-element Vector{ComplexF64}:
-0.5492714681705643 + 0.0im
 1.7355827340852823 - 1.0454627840167563im
 1.7355827340852823 + 1.0454627840167563im
vectors:
3×3 Matrix{ComplexF64}:
 0.206377+0.0im   0.17411+0.639401im   0.17411-0.639401im
 0.978444+0.0im   0.105157-0.0907931im   0.105157+0.0907931im
 -0.00741578+0.0im -0.735901-0.0im     -0.735901+0.0im
```

Hence, the eigenvalues are

$$\lambda_1 = -0.5493 + 0.0000i, \quad \lambda_2 = 1.7356 - 1.0455i, \quad \lambda_3 = 1.7356 + 1.0455i,$$

and the corresponding eigenvectors are

$$v_1 = \begin{bmatrix} 0.2064 + 0.0000i \\ 0.9784 + 0.0000i \\ -0.0074 + 0.0000i \end{bmatrix}, \quad v_2 = \begin{bmatrix} 0.1741 + 0.6394i \\ 0.1052 - 0.0908i \\ -0.7359 + 0.0000i \end{bmatrix}, \quad v_3 = \begin{bmatrix} 0.1741 - 0.6394i \\ 0.1052 + 0.0908i \\ -0.7359 + 0.0000i \end{bmatrix}.$$

λ_1 is real, with real eigenvector v_1 , while λ_2 and λ_3 are complex, with complex eigenvectors, v_2 and v_3 . Moreover, $\lambda_3 = \lambda_2^*$ and $v_3 = v_2^*$.

- (b) $Av = \lambda v$.

```
1 # verify A*v = lambda*v
2 for k = 1:3
3   v=E.vectors[:,k]
4   lam = E.values[k]
5   display([A*v lam*v])
6 end
```

The columns of $[Av \quad \lambda v]$ are identical.

Output

```

3×2 Matrix{ComplexF64}:
-0.113357+0.0im -0.113357+0.0im
-0.537432+0.0im -0.537432+0.0im
0.00407327+0.0im 0.00407327-0.0im

3×2 Matrix{ComplexF64}:
0.970653+0.927707im 0.970653+0.927707im
0.0875883-0.267517im 0.0875883-0.267517im
-1.27722+0.769357im -1.27722+0.769357im

3×2 Matrix{ComplexF64}:
0.970653-0.927707im 0.970653-0.927707im
0.0875883+0.267517im 0.0875883+0.267517im
-1.27722-0.769357im -1.27722-0.769357im

```

- (c) The eigenvectors are linearly independent and hence form a basis for \mathbb{C}^3 .

```

1 # Check for a basis
2 V=E.vectors
3
4 # Note V' computes the complex conjugate transpose of the complex matrix V
5 # V'*V is theoretically guaranteed to have a real determinant, though
6 # numerically, it may have a small imaginary part
7
8 # Check linear independence of the columns of V
9 det(V'*V)

```

Output

```
0.897060699129927 + 1.6064687371449844e-18im
```

We see that the determinant equals $0.8971 + \mathbf{i}1.6065 \times 10^{-18}$, that is, the imaginary part is the numerical equivalent of zero. Because the determinant is nonzero, the columns of V are linearly independent.

(Optional Read:) Note: $V' \cdot V := (V^*)^\top \cdot V$ is a **Hermitian Matrix**, that is a square matrix that is equal to the transpose of its (complex) conjugate matrix. The diagonal elements of a Hermitian matrix are all real numbers, and the element of the (i, j) position is equal to the (complex) conjugate of the element in the (j, i) position. In our case,

$$V' \cdot V := (V^*)^\top \cdot V = \begin{bmatrix} 1.0000 + 0.0000\mathbf{i} & 0.1443 + 0.0431\mathbf{i} & 0.1443 - 0.0431\mathbf{i} \\ 0.1443 - 0.0431\mathbf{i} & 1.0000 + 0.0000\mathbf{i} & 0.1658 - 0.2036\mathbf{i} \\ 0.1443 + 0.0431\mathbf{i} & 0.1658 + 0.2036\mathbf{i} & 1.0000 + 0.0000\mathbf{i} \end{bmatrix},$$

satisfies all of these properties.

- (d) If $v = v_R + \mathbf{i}v_I$ is a complex eigenvector, then v and its complex conjugate v^* can be replaced by v_R and v_I , the real and imaginary parts of v , to build a basis for \mathbb{R}^3 .

The vectors $V_{\text{new}} := [v_1, \text{real}(v_2), \text{imag}(v_2)]$ are linearly independent because the determinant of $V_{\text{new}}^\top \cdot V_{\text{new}}$ equals 0.224 and hence is nonzero.

```

1 # real(E.vectors[:,1]) removes the zero imaginary part and makes it
2 # a Float64 instead a ComplexF64
3

```

```

4 # real(E.vectors[:,2]) computes the real part of v2
5 # imag(E.vectors[:,2]) computes the imaginary part of v2
6
7 # Stack candidate basis vectors as columns of V_new
8 V_new = [real(E.vectors[:,1]) real(E.vectors[:,2]) imag(E.vectors[:,2])]
9
10 display(V_new)
11
12 # Check linear independence
13 det(V_new' * V_new)

```

Output

```

3×3 Matrix{Float64}:
 0.206377   0.17411   0.639401
 0.978444   0.105157  -0.0907931
 -0.00741578 -0.735901 -0.0

0.22426517478248179

```

■

Eigenvalues and Eigenvectors of the Matrix Exponential

Remarkably, eigenvectors of A are also eigenvectors of e^{At} , for all $t \in \mathbb{R}$, while the eigenvalues are transformed under the exponential map in a simple manner.

Proposition 9.49: Eigenvalues and Eigenvectors of the Matrix Exponential

Suppose that A is an $n \times n$ real matrix and that $Av = \lambda v$, with $v \neq 0_{n \times 1}$. Then, the following are true,

- (a) $A^k \cdot v = \lambda^k v$ for all integers $k \geq 1$,
- (b) $e^A \cdot v = e^\lambda v$, and
- (c) $e^{At} \cdot v = e^{\lambda t} v$ for all $t \in \mathbb{R}$.

Hence, if $v \neq 0_{n \times 1}$ is an eigenvector of A with corresponding eigenvalue λ , then v is an eigenvector of the matrix exponential function, e^{At} , with corresponding eigenvalue $e^{\lambda t}$. The surprising part is that only the eigenvalue is a function of time, and it (aka, $e^{\lambda t}$) is either a standard (real) exponential function of time or a complex exponential function of time. The eigenvectors of e^{At} are constant and identical to the eigenvectors of A .

Proof: (a) is proved by induction. We know the base case holds because we assumed that $Av = \lambda v$. Hence, assume $A^k \cdot v = \lambda^k v$ holds for some $k \geq 1$; we need to show it holds for $k + 1$.

$$\begin{aligned}
A^{k+1} \cdot v &= A \cdot (A^k \cdot v) \quad (\text{algebra}) \\
&= A \cdot (\lambda^k \cdot v) \quad (\text{by the induction hypothesis}) \\
&= \lambda^k \cdot (A \cdot v) \quad (\text{algebra, swapped the order of } A \text{ and } \lambda^k) \\
&= \lambda^k \cdot (\lambda v) \quad (\text{because } Av = \lambda v) \\
&= \lambda^{k+1} v \quad (\text{algebra}).
\end{aligned}$$

To prove (b), we use (a) with the power series for the matrix exponential, namely,

$$\begin{aligned}
e^A \cdot v &= \left(I_n + \sum_{k=1}^{\infty} \frac{t^k}{k!} A^k \right) \cdot v \quad (\text{definition of } e^A) \\
&= v + \sum_{k=1}^{\infty} \frac{t^k}{k!} A^k \cdot v \quad (\text{multiply through by } v) \\
&= v + \sum_{k=1}^{\infty} \frac{t^k}{k!} \lambda^k \cdot v \quad (\text{from part (a)}) \\
&= \left(1 + \sum_{k=1}^{\infty} \frac{t^k}{k!} \lambda^k \right) \cdot v \quad (\text{factor out } v) \\
&= e^{\lambda t} v \quad (\text{apply the definition of } e^{\lambda t}).
\end{aligned}$$

Part (c) is obtained from (b) by substituting At in place of A . ■.

Example 9.50. Use the eigenvectors of $A = \begin{bmatrix} -2.7690 & -4.1472 \\ 0.6061 & 1.2690 \end{bmatrix}$ as initial conditions for $\dot{x} = Ax$. Observe the “qualitative properties” of the corresponding solutions.

Solution: Let's find the eigenvalues and eigenvectors of A .

```

1 using LinearAlgebra
2
3 A = [-2.7699      -4.14723;  0.606141     1.26899]
4
5 E = eigen(A)

```

Output

```

Eigen{Float64, Float64, Matrix{Float64}, Vector{Float64}}
values:
2-element Vector{Float64}:
-2.0000057922545804
 0.50000579225458
vectors:
2×2 Matrix{Float64}:
-0.98324      0.785355
 0.182314    -0.619045

```

Hence, the eigenvalues

$$\lambda_1 = -2, \quad \lambda_2 = 0.5$$

are both real and distinct, and the corresponding eigenvectors are

$$v_1 = \begin{bmatrix} -0.983240 \\ 0.182314 \end{bmatrix}, \quad v_2 = \begin{bmatrix} 0.785355 \\ -0.619045 \end{bmatrix}$$

are real and linearly independent.

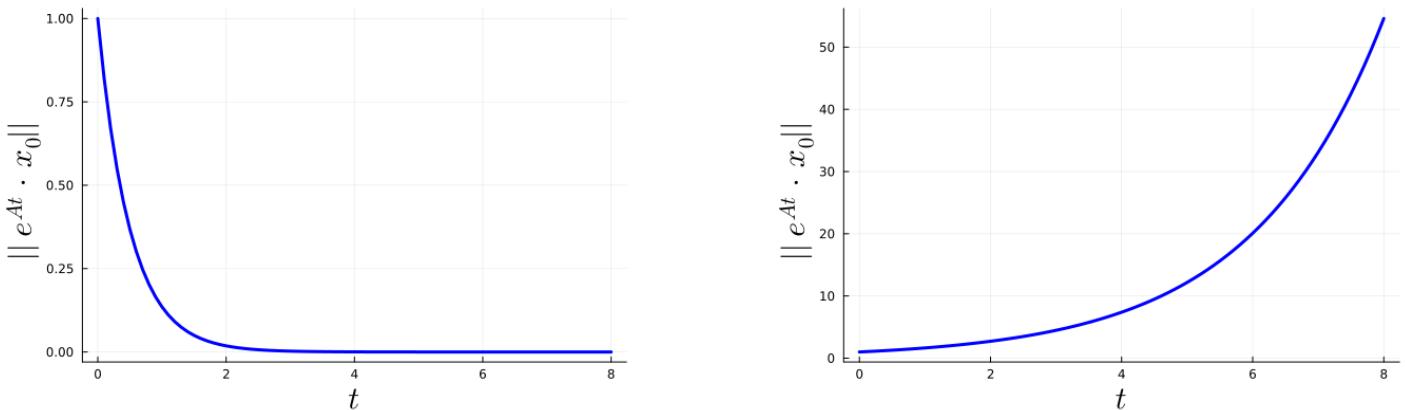
We note that $\lambda_1 < 0$ while $\lambda_2 > 0$. We next compute and plot solutions to $\dot{x} = Ax$, with $x_0 = v_1$ and $x_0 = v_2$. In fact, we are going to plot the norm of each solution, $\|e^{At}x_0\|$. We show only the code for the first eigenvector as the initial condition.

```

1 # Compute and plot ||x(t) = exp(At)*x0|| for x0 an eigenvector.
2 using Plots
3 using LinearAlgebra
4 using LaTeXStrings
5
6 function normSolution(t; A=A, x0=E.vectors[:, 1])
7     eAt = exp(A*t)
8     xt = eAt * x0
9     y = norm(xt)
10    return y
11 end
12
13 # time values where we compute the solution
14 t = 0:.1:8
15
16 y = normSolution.(t)
17
18 p1 = plot(t, y, guidefont = 20, legend = false, xlabel=L"t", ylabel=L"\|e^{At} \cdot x_0\|", lw=3,
19 color=:blue)

```

Output



The image on the left shows $\|x(t)\|$ decaying like $e^{\lambda_1 t} = e^{-2t}$ and the image on the right shows $\|x(t)\|$ exploding like $e^{\lambda_2 t} = e^{0.5t}$. This is because, when $x_0 = v$, an eigenvector, then $e^{At} \cdot v = e^{\lambda t}v$, and hence,

$$\|e^{At} \cdot v\| = \|e^{\lambda t} \cdot v\| = |e^{\lambda t}| \|v\|,$$

which implies for real eigenvalues,

$$\lim_{t \rightarrow \infty} \|e^{At} \cdot v\| = \begin{cases} 0 & \lambda < 0 \\ \|v\| & \lambda = 0 \\ \infty & \lambda > 0. \end{cases}$$

■

Our next challenge is to understand what happens in the case of complex eigenvalues and eigenvectors. Suppose that A is real with complex eigenvalue $\lambda = a + i\omega$ and corresponding eigenvector, $v = v_R + i v_I$. From Prop. 9.49,

$$e^{At}v = e^{\lambda t}v = e^{(a+i\omega)t}v.$$

Hence,

$$\|e^{At}v\| = \|e^{(a+i\omega)t}v\| = |e^{(a+i\omega)t}| \cdot \|v\| = |e^{at} \cdot e^{i\omega t}| \cdot \|v\| = |e^{at}| \cdot |e^{i\omega t}| \cdot \|v\| = e^{at} \cdot |e^{i\omega t}| \cdot \|v\|.$$

By Euler's Formula,

$$e^{i\omega t} = \cos(\omega t) + i \sin(\omega t),$$

and therefore,

$$|e^{i\omega t}| = |\cos(\omega t) + i\sin(\omega t)| = \sqrt{\cos^2(\omega t) + \sin^2(\omega t)} = 1.$$

After all of this, we obtain that

$$\|e^{At}v\| = e^{at}\|v\|,$$

where $a = \operatorname{real}(\lambda)$, the real part of the eigenvalue λ .

Example 9.51. Use the eigenvectors of $A = \begin{bmatrix} -1.8931 & 11.4911 & 2.8358 & -21.4883 \\ 1.9318 & -6.8883 & -0.6626 & 9.1683 \\ -1.4704 & 2.0461 & 0.1305 & -8.2359 \\ 1.5956 & -4.4190 & -0.9576 & 5.6509 \end{bmatrix}$ as initial conditions for $\dot{x} = Ax$. Observe the “qualitative properties” of the corresponding solutions.

Solution: Let's find the eigenvalues and eigenvectors of A .

```

1 # Compute eigenvalues and eigenvectors
2 A = [-1.89311 11.49111 2.83579 -21.4883
3   1.93182 -6.88831 -0.662625 9.16832
4   -1.47045 2.04612 0.130518 -8.23588
5   1.59557 -4.41903 -0.95763 5.65091]
6
7 E=eigen(A)
8
9 for k = 1:2
10    println("Pair of eigenvalues and vectors")
11    display(E.values[(2*k-1):2*k])
12    display(E.vectors[:,(2*k-1):2*k])
13    println(" ")
14 end

```

Output

```

Pair of eigenvalues and vectors
2-element Vector{ComplexF64}:
-2.000000213236534 - 1.9999961984135264im
-2.000000213236534 + 1.9999961984135264im
4×2 Matrix{ComplexF64}:
-0.314334+0.172478im -0.314334-0.172478im
-0.0323151-0.563835im -0.0323151+0.563835im
-0.654925-0.0im -0.654925+0.0im
-0.121327-0.329915im -0.121327+0.329915im

```

```

Pair of eigenvalues and vectors
2-element Vector{ComplexF64}:
0.5000042132365348 - 2.9999956474834977im
0.5000042132365348 + 2.9999956474834977im
4×2 Matrix{ComplexF64}:
0.77073-0.0im 0.77073+0.0im
-0.0107431+0.313272im -0.0107431-0.313272im
0.385818-0.320985im 0.385818+0.320985im
-0.0406637+0.232767im -0.0406637-0.232767im

```

Hence, the eigenvalues

$$\lambda_1 = -2 - i2, \quad \lambda_2 = -2 + i2, \quad \lambda_3 = 0.5 - i3, \quad \lambda_4 = 0.5 + i3$$

are complex and respect complex conjugate symmetry, namely, $\lambda_2 = \lambda_1^*$ and $\lambda_4 = \lambda_3^*$. The corresponding eigenvectors

$$v_1 = \begin{bmatrix} -0.314334 + 0.172478i \\ -0.0323151 - 0.563835i \\ -0.654925 - 0.000000i \\ -0.121327 - 0.329915i \end{bmatrix}, \quad v_2 = \begin{bmatrix} -0.314334 - 0.172478i \\ -0.0323151 + 0.563835i \\ -0.654925 + 0.000000i \\ -0.121327 + 0.329915i \end{bmatrix}$$

$$v_3 = \begin{bmatrix} 0.77073 - 0.000000i \\ -0.0107431 + 0.313272i \\ 0.385818 - 0.320985i \\ -0.0406637 + 0.232767i \end{bmatrix}, \quad v_4 = \begin{bmatrix} 0.77073 + 0.000000i \\ -0.0107431 - 0.313272i \\ 0.385818 + 0.320985i \\ -0.0406637 - 0.232767i \end{bmatrix}$$

are also complex and respect complex conjugate symmetry, $v_2 = v_1^*$ and $v_4 = v_3^*$.

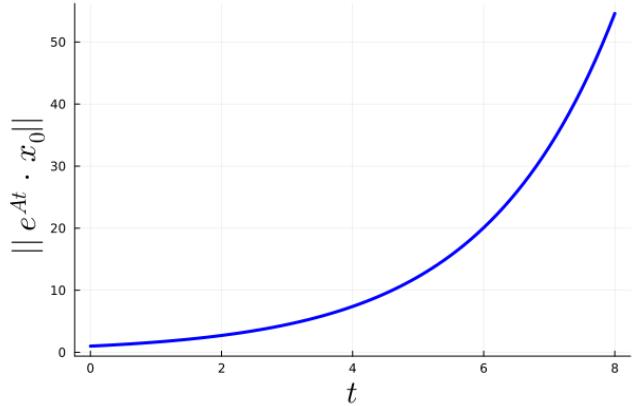
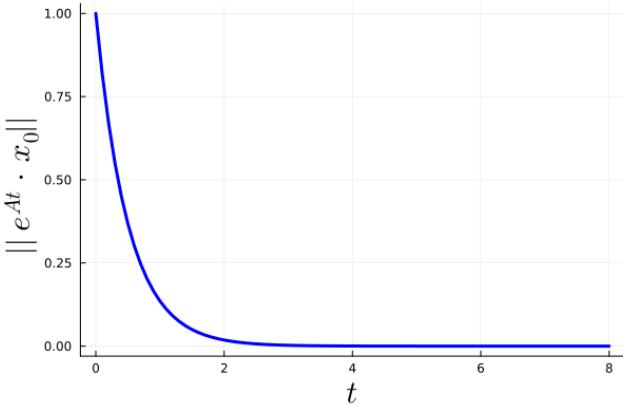
We note that $\text{real}(\lambda_1) = \text{real}(\lambda_2) = -2 < 0$ while $\text{real}(\lambda_3) = \text{real}(\lambda_4) = 0.5 > 0$. We next compute and plot solutions to $\dot{x} = Ax$, with $x_0 = v_1$ and $x_0 = v_3$. In fact, we are going to plot the norm of each solution, $\|e^{At}x_0\|$. We show only the code for the first eigenvector as the initial condition.

```

1 # Compute and plot ||x(t) = exp(At)*x0|| for x0 an eigenvector.
2 using Plots
3 using LinearAlgebra
4 using LaTeXStrings
5
6 function normSolution(t; A=A, x0=E.vectors[:, 1])
7     eAt = exp(A*t)
8     xt = eAt * x0
9     y = norm(xt)
10    return y
11 end
12
13 # time values where we compute the solution
14 t = 0:.1:8
15
16 y = normSolution.(t)
17
18 p1 = plot(t, y, guidefont = 20, legend = false, xlabel=L"t", ylabel=L"\|e^{At} \cdot x_0\|", lw=3,
19           color=:blue )

```

Output



The image on the left shows $\|x(t)\|$ decaying like $e^{\text{real}(\lambda_1)t} = e^{-2t}$ and the image on the right shows $\|x(t)\|$ exploding like $e^{\text{real}(\lambda_3)t} = e^{0.5t}$. This is because, when $x_0 = v$, an eigenvector, then $e^{At} \cdot v = e^{\lambda t} v$, and hence,

$$\|e^{At} \cdot v\| = \|e^{\lambda t} \cdot v\| = |e^{\lambda t}| \cdot \|v\| = e^{\text{real}(\lambda)t} \cdot \|v\|,$$

which implies for complex eigenvalues,

$$\lim_{t \rightarrow \infty} \|e^{At} \cdot v\| = \begin{cases} 0 & \text{real}(\lambda) < 0 \\ \|v\| & \text{real}(\lambda) = 0 \\ \infty & \text{real}(\lambda) > 0. \end{cases}$$

What does it mean when the initial condition is a complex vector? That seems very strange! You make a good point. We've made great progress in understanding the solutions of $\dot{x} = Ax$ for special initial conditions, namely, x_0 is an eigenvector of A . But yeah, on a real system, we cannot set x_0 to a complex value. So, did we just do a bunch of work for nothing?

Hardly. However, it does take a bit more Linear Algebra to make sense of our calculations.

Complex (scalar) Exponential Functions

Complex exponentials are a small, yet very powerful generalization of Euler's Formula. If you like to learn from videos, here is [The MAGIC of Complex Exponentials!](#) by FiguR3 iT ouT (sic).

Complex Exponentials

Definition 9.52. For $a, \omega \in \mathbb{R}$,

$$e^{(a+i\omega)t} := e^{at} \cdot (\cos(\omega t) + i \sin(\omega t)) \quad (9.67)$$

is called a **complex exponential**.

Equivalently, let $z := a + i\omega \in \mathbb{C}$ be a complex number, then e^{zt} is a **complex exponential**.

Note: We first encountered a complex exponential in the context of Example 9.50. Complex exponentials arise in the solution of linear ODEs when there are complex eigenvalues.

Proposition 9.53: Properties of Complex Exponentials

Complex exponentials are fundamental in various fields of mathematics, engineering, and physics. Below are some of the key properties of complex exponentials, for $a, \omega \in \mathbb{R}$ and $z := a + i\omega \in \mathbb{C}$:

- (a) **Periodicity:** The function $e^{i\omega t}$ is periodic with period $T = \frac{2\pi}{|\omega|}$, for $\omega \neq 0$.
- (b) **Multiplicative Property:** For any two complex numbers $z_1 = a_1 + i\omega_1$ and $z_2 = a_2 + i\omega_2$, the exponential function satisfies $e^{z_1 t} \cdot e^{z_2 t} = e^{(z_1+z_2)t}$.
- (c) **Inverse:** The inverse of e^{zt} is given by e^{-zt} .
- (d) **Differentiation and Integration:** The derivative and antiderivative of e^{zt} with respect to t are:
- $\frac{d}{dt} e^{zt} = ze^{zt}$ and
 - for $z \neq 0$, $\int e^{zt} dt = \frac{1}{z} e^{zt} + C$, where C is the constant of integration.
- (e) **Conjugation:** The complex conjugate of $e^{(a+i\omega)t}$ is $e^{(a-i\omega)t}$.
- (f) **Magnitude:** The magnitude of $e^{(a+i\omega)t}$ is $|e^{(a+i\omega)t}| = e^{at}$, which exhibits exponential growth or decay based on the sign of the real part, a .

Note: These properties make complex exponentials an indispensable tool in the analysis of oscillatory and wave phenomena, signal processing, control theory, and many other areas.

Understanding the Roles Played by the Real and Imaginary Parts of Complex Eigenvalues and Eigenvectors

Proposition 9.54: Real and Imaginary Parts of a Complex Eigenvector as Initial Conditions

Suppose that A is a real $n \times n$ matrix with a complex eigenvalue $\lambda = a + i\omega$ and corresponding eigenvector $v = v_R + iv_I$. Then

$$\begin{aligned} e^{At}v_R &= e^{at}\cos(\omega t)v_R - e^{at}\sin(\omega t)v_I \\ e^{At}v_I &= e^{at}\sin(\omega t)v_R + e^{at}\cos(\omega t)v_I. \end{aligned} \tag{9.68}$$

Note: Equation (9.68) helps us to see the oscillatory aspect of the solutions of linear ODEs when there are complex eigenvalues. The imaginary part of the eigenvalue, $\omega := \text{imag}(\lambda)$ sets the (angular) frequency of the oscillation, while the real part of the eigenvalue, $a := \text{real}(\lambda)$, determines the “amplification” or “attenuation” factor, depending on a is positive, negative, or zero.

When the initial condition is a real eigenvector, the solution of $\dot{x} = Ax$ evolves entirely along the eigenvector via $e^{\lambda t}v$. It is “living” in a one-dimensional subspace given by $\text{span}\{v\}$. For a complex eigenvector, when the ODE is initialized at v_R or v_I , the solution “lives” in a two-dimensional subspace given by $\text{span}\{v_R, v_I\}$. Within that subspace, the solution spirals to the origin if $a < 0$, spirals to infinity if $a > 0$, and forms a circle if $a = 0$. To be able to deduce such detailed information through Linear Algebra is pretty awesome, don’t you agree?

Example 9.55. Using our new insights, re-analyze the solutions to $\dot{x} = Ax$ for the matrix in Example 9.51. The key data are

$$A = \begin{bmatrix} -1.8931 & 11.4911 & 2.8358 & -21.4883 \\ 1.9318 & -6.8883 & -0.6626 & 9.1683 \\ -1.4704 & 2.0461 & 0.1305 & -8.2359 \\ 1.5956 & -4.4190 & -0.9576 & 5.6509 \end{bmatrix}$$

has eigenvalues,

$$\lambda_1 = -2 - i2, \quad \lambda_2 = \lambda_1^* = -2 + i2, \quad \lambda_3 = 0.5 - i3, \quad \lambda_4 = \lambda_3^* = 0.5 + i3$$

and corresponding eigenvectors

$$v_1 = \begin{bmatrix} -0.3143 \\ -0.0323 \\ -0.6549 \\ -0.1213 \end{bmatrix} + i \begin{bmatrix} 0.1725 \\ -0.5638 \\ -0.0000 \\ -0.3299 \end{bmatrix}, \quad v_2 = v_1^* = \begin{bmatrix} -0.3143 \\ -0.0323 \\ -0.6549 \\ -0.1213 \end{bmatrix} - i \begin{bmatrix} 0.1725 \\ -0.5638 \\ -0.0000 \\ -0.3299 \end{bmatrix}$$

$$v_3 = \begin{bmatrix} 0.7707 \\ -0.0107 \\ 0.3858 \\ -0.0407 \end{bmatrix} + i \begin{bmatrix} -0.0000 \\ 0.3133 \\ -0.3210 \\ 0.2328 \end{bmatrix}, \quad v_4 = v_3^* = \begin{bmatrix} 0.7707 \\ -0.0107 \\ 0.3858 \\ -0.0407 \end{bmatrix} - i \begin{bmatrix} -0.0000 \\ 0.3133 \\ -0.3210 \\ 0.2328 \end{bmatrix}.$$

Solutions: Compute the solution of $e^{At} \cdot v_{R,1}$ and plot its “projection” onto $\text{span}\{v_{R,1}, v_{I,1}\}$.

```

1 using LinearAlgebra
2 using Plots
3 using LaTeXStrings
4
5 A = [-1.89311 11.4911 2.83579 -21.4883
6 1.93182 -6.88831 -0.662625 9.16832
7 -1.47045 2.04612 0.130518 -8.23588
8 1.59557 -4.41903 -0.95763 5.65091]
9
10 E=eigen(A)
11 (lam1, lam2, lam3, lam4) = E.values
12 (v1, v2, v3, v4) = E.vectors
13
14 t = 0.0:0.02:100
15
16 # For x0 = real(v1)
17 x = exp.(real(lam1)*t).*cos.(imag(lam1)*t)
18 y = -exp.(real(lam1)*t).*sin.(imag(lam1)*t)
19
20 p1 = plot(x, y, guidefont = 15, lw=3, color=:blue, legend=false,
21 aspect_ratio=1, xlabel=L"e^{at} \cos(\omega t) v_R", ylabel=L"-e^{at} \sin(\omega t) v_I")
22
23 # mark the initial condition
24 scatter!([1], [0], markersize=8, color=:red)
25
26
27 display(p1)
```

Compute the solution of $e^{At} \cdot v_{I,3}$ and plot its “projection” onto $\text{span}\{v_{R,3}, v_{I,3}\}$.

```

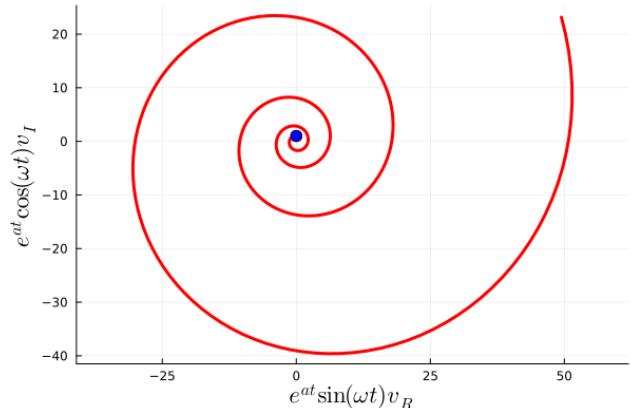
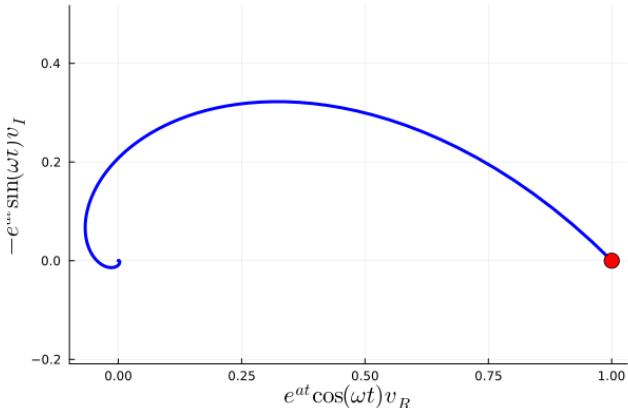
1 using LinearAlgebra
2 using Plots
3 using LaTeXStrings
4
5 A = [-1.89311 11.4911 2.83579 -21.4883
6 1.93182 -6.88831 -0.662625 9.16832
```

```

7 -1.47045   2.04612   0.130518   -8.23588
8    1.59557   -4.41903   -0.95763      5.65091]
9
10 E=eigen(A)
11 (lam1, lam2, lam3, lam4) = E.values
12 (v1, v2, v3, v4) = E.vectors
13
14 t = 0.0:0.02:8
15
16 # For x0 = imag(v3)
17 x = exp.(real(lam3)*t).*sin.(imag(lam3)*t)
18 y = exp.(real(lam3)*t).*cos.(imag(lam3)*t)
19
20 p1 = plot(x, y, guidefont = 15, lw=3, color=:red, legend=false,
21 aspect_ratio=1, xlabel=L"e^{at} \sin(\omega t) v_R", ylabel=L"e^{at} \cos(\omega t) v_I")
22
23 # mark the initial condition
24 scatter!([0], [1], markersize=6, color=:blue)
25
26 display(p1)

```

Output



On the left is the solution of $e^{At}v_{R,1}$, projected onto $\text{span}\{v_{R,1}, v_{I,1}\}$, with the red dot marking the initial condition. e^{-2t} decays so quickly that we cannot see with our eyes very much of the spiral into the origin. On the right is the solution of $e^{At}v_{I,3}$, projected onto $\text{span}\{v_{R,3}, v_{I,3}\}$, with the blue dot marking the initial condition. $e^{0.5t}$ expands slowly enough that we can see it spiral outward to infinity. ■.

Based on the results in this subsection, we understand very well how the solutions to $\dot{x} = Ax$ evolve for initial conditions tied to eigenvectors. The next subsection treats general initial conditions.

9.5.8 Exponential Stability of Linear Systems of ODEs with Implications for Nonlinear ODEs

Equilibrium Points and Stability Properties

Let A be a real $n \times n$ matrix and consider the linear system of ODEs $\dot{x} = Ax$. We specialize Def. 9.30 to the case of linear systems.

Definition 9.56. $x_e \in \mathbb{R}^n$ is an **equilibrium point** of $\dot{x} = Ax$ if $Ax_e = 0$. Equivalently, the constant vector $\varphi(t, t_0, x_e) := x_e$ is a solution of the ODE.

Note: For a linear ODE, the origin is always an equilibrium point.

Definition 9.57. $x_e \in \mathbb{R}^n$ is a **globally asymptotically stable equilibrium point** of $\dot{x} = Ax$ if for all $x_0 \in \mathbb{R}$,

$$\lim_{t \rightarrow \infty} e^{At} x_0 = 0.$$

Definition 9.58. $x_e \in \mathbb{R}^n$ is a **globally exponentially stable equilibrium point** of $\dot{x} = Ax$ if there exist $\kappa > 0$ and $\gamma < \infty$ such that, for all $x_0 \in \mathbb{R}$,

$$\|e^{At} x_0\| \leq \gamma \cdot e^{-\kappa t} \cdot \|x_0\|.$$

Note: For a linear system of ODEs, Global **exponential** stability of an equilibrium point and global **asymptotic** stability of an equilibrium point are equivalent concepts. For a nonlinear system of ODEs, they are different notions of stability. In this course, understanding that each notion implies that solutions converge to zero as time marches to infinity is enough. The rate at which the convergence takes place is less relevant.

We laid the seeds for guessing that the **eigenvalues** provide all the information needed to check for asymptotic stability. The **eigenvectors** were present because, when used as an initial condition, the solution to a linear ODE is very simple, $e^{At} \cdot v = e^{\lambda t} v$.

Proposition 9.59: Exponential Stability of Linear Systems of ODEs

Let A be a real $n \times n$ matrix with eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$, and consider the linear system of ODEs,

$$\dot{x} = Ax.$$

If $\text{real}(\lambda_i) < 0$ for $1 \leq i \leq n$ (aka, the real parts of all eigenvalues are negative), then the **origin** is globally exponentially stable.

Note: The proof is given in Chapter 9.9 for the case of distinct eigenvalues. The general case is part of EECS 560 Linear Systems and uses something called a **Jordan Normal Form** for square matrices. While it's an awesome thing to know, saving a few things for grad school is okay.

Example 9.60. Determine if the origin is a globally exponentially stable equilibrium point for the following linear systems of ODEs.

- (a) Consider the mass-spring-damper system in Example 9.28. When $m = 5 \text{ kg}$, $k = 2 \text{ N/m}$, and $b = 3 \text{ Ns/m}$ (Newton seconds per meter), the model becomes

$$\dot{x} = \begin{bmatrix} 0.0000 & 1.0000 \\ -0.4000 & -0.6000 \end{bmatrix} \cdot x.$$

- (b) Consider the RLC network in Example 9.27. For a consumer-grade RLC circuit in an older WiFi radio receiver, where discrete circuit elements were still used, typical component values might be:

- (R) 10 kilo-ohms ($k\Omega$)
- (L) 20 nanohenries (nH)
- (C) 10 nanofarads (nF),

yielding $\dot{x} = Ax$, with

$$A = \begin{bmatrix} 0.0000e+00 & 1.0000e+08 \\ -5.0000e-01 & -5.0000e+11 \end{bmatrix}.$$

Solutions:

- (a) Mass-spring-damper system: **Ans.** The origin is globally exponentially stable because the eigenvalues all have negative real parts.

```
1 using LinearAlgebra
2 # Mass spring damper system
3 m = 5 # kg
4 k = 2 # N/m
5 b = 3 # Ns/m (Newton seconds per meter)
6 A = [0 1; -k/m -b/m]
7 E = eigen(A)
```

Output

```
Eigen{ComplexF64, ComplexF64, Matrix{ComplexF64}, Vector{ComplexF64}}
values:
2-element Vector{ComplexF64}:
-0.2999999999999993 - 0.5567764362830022im
-0.2999999999999993 + 0.5567764362830022im

vectors:
2×2 Matrix{ComplexF64}:
 0.845154-0.0im      0.845154+0.0im
 -0.253546-0.470562im -0.253546+0.470562im
```

- (b) RLC network: **Ans.** The origin is globally exponentially stable because the eigenvalues all have negative real parts.

```
1 using LinearAlgebra
2 # RLC Circuit
3 R = 10e3 # ohms
4 L = 20e-9 # Henries
5 C = 10e-9 # farads
6 A = [0 1/C; -C/L -R/L]
7 E = eigen(A)
```

Output

```
Eigen{Float64, Float64, Matrix{Float64}, Vector{Float64}}
values:
2-element Vector{Float64}:
-4.99999999999999e11
-0.0001220703125

vectors:
2×2 Matrix{Float64}:
 -0.0002    1.0
  1.0       -1.0e-12
```

Example 9.61. Take a random initial condition for the mass-spring-damper system and plot the norm of its solution.

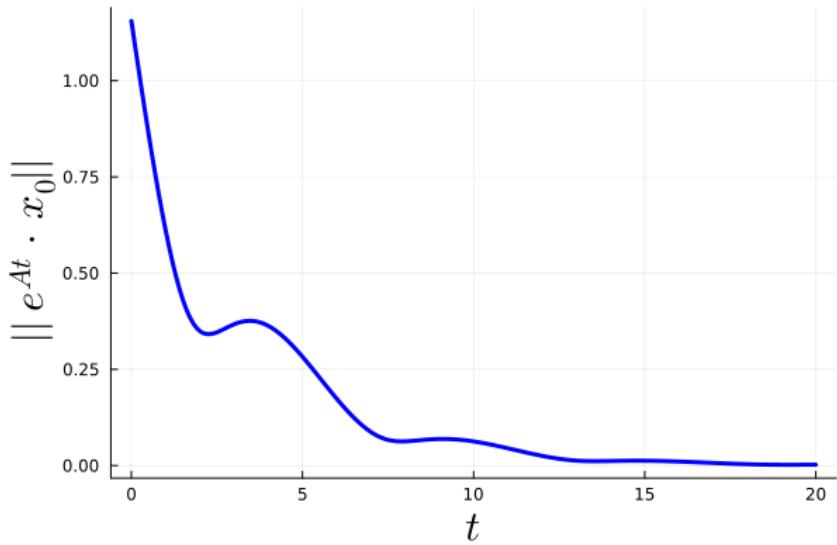
Solution:

```
1 # Take a random initial condition for the mass-spring-damper system and plot the norm of
  its solution.
2 using Plots
3 using LinearAlgebra
4 using Random
5 using LaTeXStrings
6
7 # Mass spring damper system
8 m = 5 # kg
9 k = 2 # N/m
10 b = 3 # Ns/m (Newton seconds per meter)
11 A = [0 1; -k/m -b/m]
12
13
14 Random.seed!(543210)
15 x0 = randn(2, 1); display(x0)
16 function normSolution(t; A=A, x0=x0)
17     eAt = exp(A*t)
18     xt = eAt * x0
19     y = norm(xt)
20     return y
21 end
22
23 # time values where we compute the solution
24 t = 0:0.01:20
25
26 y = normSolution.(t)
27
28 p1 = plot(t, y, guidefont = 20, legend = false, xlabel=L"t", ylabel=L"\|e^{At} \cdot x_0\|", lw=3,
29 color=:blue )
```

Output

```
2×1 Matrix{Float64}:
-0.9072986899852817
 0.7151631134658601
```

Because the eigenvalues are complex, we see oscillation in the norm as the solution decays to the origin, $x_e = [0; 0]$.



■

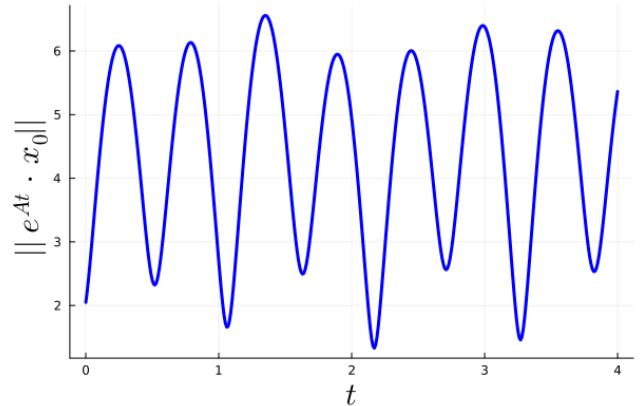
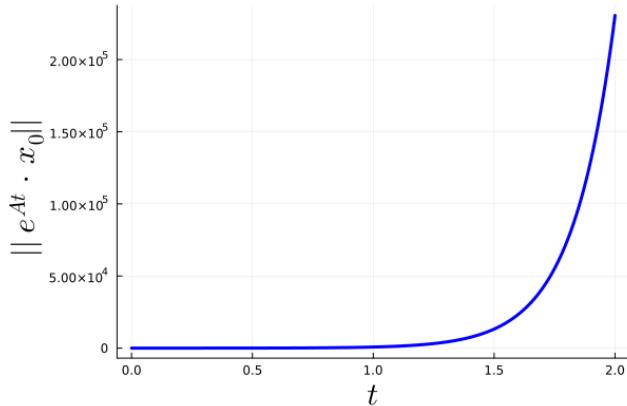
Example 9.62. Determine if the origin is a globally exponentially stable equilibrium point for the linearized model of the 2-link robot arm, for the (a) upward equilibrium and (b) downward equilibrium. We chose the 2-link arm so that we need to do everything from scratch:

- build the dynamic model,
- compute the linearization for each of the equilibria, and
- evaluate the eigenvalues.

In addition, for a random initial condition, compute solutions to each of the linearized models.

Solutions:

- (a) Upward equilibrium. **Ans.** The origin is not exponentially stable. There are four real eigenvalues, of which two are positive.
- (b) Downward equilibrium. **Ans.** The origin is not exponentially stable. There are four purely imaginary eigenvalues, meaning their real part is identically zero.



Using Lagrange's Method, we obtain,

```

1 function dyn_mod_2LinkManipulator(q, dq)
2 # DYN_MOD_2LINKMANIPULATOR
3 # 2023-12-05 14:45:06
4 #
5 # Author: Grizzle
6 #
7 # Model NOTATION: D(q)ddq + C(q,dq)*dq + G(q) = B*tau
8 # The Robot Equations: From Lagrange's Equations of Motion
9 #
10 g, L1, L2, m1, m2 = modelParameters()
11 #
12 # Variable names for the model
13 q1, q2= q
14 dq1, dq2 = dq
15 #
16 D = zeros(2, 2)
17 D[1, 1] = L1^2*m1 + m2*(L1^2 + 2*L1*L2*cos(q2) + L2^2)
18 D[1, 2] = L2*m2*(L1*cos(q2) + L2)
19 D[2, 1] = L2*m2*(L1*cos(q2) + L2)
20 D[2, 2] = L2^2*m2
21 #
22 C = zeros(2, 2)
23 C[1, 1] = -L1*L2*dq2*m2*sin(q2)
24 C[1, 2] = -L1*L2*m2*(dq1 + dq2)*sin(q2)
25 C[2, 1] = L1*L2*dq1*m2*sin(q2)
26 C[2, 2] = 0.0
27 #
28 G = zeros(2)
29 G[1] = g*m2*(L1*cos(q1) + L2*cos(q1 + q2)) + L1*g*m1*cos(q1)
30 G[2] = L2*g*m2*cos(q1 + q2)
31 #
32 B = zeros(2, 2)
33 B[1, 1] = 1
34 B[2, 2] = 1
35 #
36 JacG = zeros(2, 2)
37 JacG[1, 1] = g*m2*(-L1*sin(q1) - L2*sin(q1 + q2)) - L1*g*m1*sin(q1)
38 JacG[1, 2] = -L2*g*m2*sin(q1 + q2)
39 JacG[2, 1] = -L2*g*m2*sin(q1 + q2)
40 JacG[2, 2] = -L2*g*m2*sin(q1 + q2)
41 #
42 return (D=D, C=C, G=G, B=B, JacG=JacG)
43 end

```

Upward Equilibrium

```

1 # Linear Model for the 2-link Manipulator
2 #
3 # bring the model into the workspace for later use
4 include("dyn_mod_2LinkManipulator.jl")
5
6 #Build a function to define the parameters of the model
7
8 function modelParameters()
9     g = 9.81 # m/s^2

```

```

10   L1 = 1 # m
11   L2 = 0.7
12   m1 = 15 # kg
13   m2 = 10
14   return g, L1, L2, m1, m2
15 end
16
17 # Set equilibrium Point
18 if false
19     qe = [-pi/2, 0] #downward
20 else
21     qe = [pi/2, 0] # upward
22 end
23
24 F = dyn_mod_2LinkManipulator(qe, 0*qe)
25 D = F.D;
26 println("D matrix"); display(D)
27 JacG = F.JacG;
28 println("JacG matrix"); display(JacG)
29
30 # build the linearized model about the give equilibrium point
31 A21 = -D\JacG; A21 = cleanUp(A21);
32 println("A21 matrix"); display(A21)
33 n, ~ = size(A21)
34 A = [zeros(n, n) I(n); A21 zeros(n, n)];
35 println("A matrix"); display(A)
36
37 # Compute eigenvalues
38 E = eigen(A)

```

Output

```

D matrix
2×2 Matrix{Float64}:
 43.9  11.9
 11.9   4.9

JacG matrix
2×2 Matrix{Float64}:
 -313.92  -68.67
 -68.67  -68.67

A21 matrix
2×2 Matrix{Float64}:
 9.81  -6.54
 -9.81  29.8971

A matrix
4×4 Matrix{Float64}:
 0.0    0.0    1.0    0.0
 0.0    0.0    0.0    1.0
 9.81  -6.54    0.0    0.0
 -9.81  29.8971    0.0    0.0

Eigen{Float64, Float64, Matrix{Float64}, Vector{Float64}}
values:
4-element Vector{Float64}:

```

```

-5.718391382198321
-2.647100840002673
2.647100840002673
5.71839138219832

```

```

vectors:
4×4 Matrix{Float64}:
-0.0473235 -0.324822 0.324822 0.0473235
 0.165632 -0.139209 0.139209 -0.165632
 0.270614 0.859836 0.859836 0.270614
-0.947151 0.368501 0.368501 -0.947151

```

Downward Equilibrium

```

1 # Linear Model for the 2-link Manipulator
2 #
3 # bring the model into the workspace for later use
4 include("dyn_mod_2LinkManipulator.jl")
5
6 #Build a function to define the parameters of the model
7
8 function modelParameters()
9     g = 9.81 # m/s^2
10    L1 = 1 # m
11    L2 = 0.7
12    m1 = 15 # kg
13    m2 = 10
14    return g, L1, L2, m1, m2
15 end
16
17 # Set equilibrium Point
18 if true
19     qe = [-pi/2, 0] #downward
20 else
21     qe = [pi/2, 0] # upward
22 end
23
24 F = dyn_mod_2LinkManipulator(qe, 0*qe)
25 D = F.D;
26 println("D matrix"); display(D)
27 JacG = F.JacG;
28 println("JacG matrix"); display(JacG)
29
30 # build the linearized model about the given equilibrium point
31 A21 = -D\JacG; A21 = cleanUp(A21);
32 println("A21 matrix"); display(A21)
33 n, ~ = size(A21)
34 A = [zeros(n, n) I(n); A21 zeros(n, n)];
35 println("A matrix"); display(A)
36
37 # Compute eigenvalues
38 E = eigen(A)

```

Output

D matrix

```

2×2 Matrix{Float64}:
43.9  11.9
11.9   4.9

JacG matrix
2×2 Matrix{Float64}:
313.92  68.67
68.67   68.67

A21 matrix
2×2 Matrix{Float64}:
-9.81    6.54
 9.81  -29.8971

A matrix
4×4 Matrix{Float64}:
 0.0    0.0    1.0   0.0
 0.0    0.0    0.0   1.0
-9.81   6.54   0.0   0.0
 9.81  -29.8971  0.0   0.0

Eigen{ComplexF64, ComplexF64, Matrix{ComplexF64}, Vector{ComplexF64}}
values:
4-element Vector{ComplexF64}:
2.2204460492502973e-16 - 2.6471008400026728im
2.2204460492502973e-16 + 2.6471008400026728im
3.945433097643709e-16 - 5.718391382198318im
3.945433097643709e-16 + 5.718391382198318im

vectors:
4×4 Matrix{ComplexF64}:
 6.06496e-17-0.324822im      ...  1.65597e-17-0.0473235im
-2.13366e-17-0.139209im      2.56542e-17+0.165632im
-0.859836-0.0im              0.270614+2.991e-17im
-0.368501+2.66706e-16im      -0.947151+0.0im

```

Solutions are computed using the same code as in Example 9.61.

This next example illustrates the power of feedback control. The particular “type” of controller employed is linear state-variable feedback. The example uses material from Michigan’s EECS 560 *Linear Systems Theory* to design the feedback controller and illustrates two important points:

- an unstable equilibrium point can often be made into an exponentially stable equilibrium point with a properly designed feedback controller; and
- global exponential stability of a linearized model about an equilibrium point implies local exponential stability of the nonlinear model about the same equilibrium point.

Many feedback controllers in engineering practice are designed on the basis of these two facts; the second fact is proved in Michigan’s EECS 562 *Nonlinear Systems and Control*.

Example 9.63. Consider the following linearized model of the 2-link robot arm about the upward equilibrium, where motors (aka, actuators) have been added to the two joints,

$$\delta \dot{x} = \underbrace{\begin{bmatrix} 0.000 & 0.000 & 1.000 & 0.000 \\ 0.000 & 0.000 & 0.000 & 1.000 \\ 9.810 & -6.540 & 0.000 & 0.000 \\ -9.810 & 29.897 & 0.000 & 0.000 \end{bmatrix}}_A \delta x + \underbrace{\begin{bmatrix} 0.000 & 0.000 \\ 0.000 & 0.000 \\ 0.067 & -0.162 \\ -0.162 & 0.597 \end{bmatrix}}_B \delta \tau, \quad (9.69)$$

where $\delta x := \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}}_x - \underbrace{\begin{bmatrix} 1.571 \\ 0.000 \\ 0.000 \\ 0.000 \end{bmatrix}}_{x_e}$ and $\delta \tau := \underbrace{\begin{bmatrix} \tau_1 \\ \tau_2 \end{bmatrix}}_\tau - \underbrace{\begin{bmatrix} 0 \\ 0 \end{bmatrix}}_{\tau_e}$. τ_e is the value of the motor torque at the equilibrium point. In this case, it equals $0_{2 \times 1}$, because no torque is required at the equilibrium.

Because motors have been attached to both joints (aka, the joints are now actuated), we have the possibility to apply motor torques (τ_1, τ_2) to render the upward equilibrium exponentially stable. When the applied motor torques are made a linear function of the states of the model, this is called **linear state-variable feedback control**: $\tau = -K \cdot (x - x_e) = -K \delta x$, where, in our case, K is a 2×4 matrix of **controller gains**. Using methods from Michigan's EECS 560, a linear state-variable feedback controller has been designed for you,

$$\delta \tau = - \underbrace{\begin{bmatrix} 357.820 & 80.570 & 62.084 & 16.829 \\ 80.570 & 73.570 & 16.829 & 6.930 \end{bmatrix}}_K \delta x. \quad (9.70)$$

Your tasks are:

(a) Determine the eigenvalues of the closed-loop linearized system.

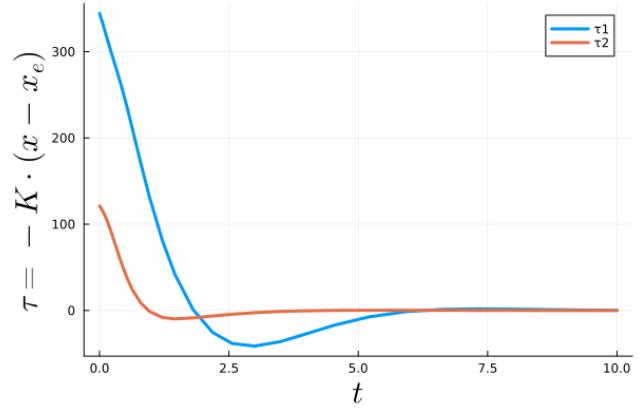
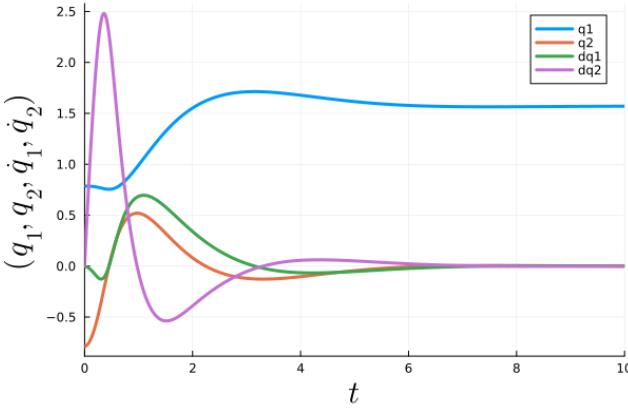
(b) Apply the linear state-variable feedback to the nonlinear model of the 2-link robot arm and produce a simulation for $x_0 =$

$$\begin{bmatrix} 0.785 \\ -0.785 \\ 0.000 \\ 0.000 \end{bmatrix}$$

, which has the base link at $\frac{\pi}{4}$ from the horizontal and the second link at $-\frac{\pi}{4}$ from the base link.

Will it topple over, or will it stand up and balance?

Solution: Ans. Here is a **GIF** generated from a simulation of the nonlinear 2-link robot arm, using the linear feedback controller. See the **Caveat** below.



Caveat: The model we have chosen is particularly simple because (1) each joint of the robot arm is actuated, and (2), we have assumed the motors can produce as much torque as the controller demands of them. Not all robots have as many actuators as joints and real motors have limitations. Chapter 10 will dig into these issues. For now, we ignore them and simply marvel at the fact that we were able to balance the arm in the inverted position. The BallBot used in ROB 311 *How to Build Robots and Make them Move* is an **underactuated robot**, meaning it has fewer actuators than it has position variables. We will confront the balance control of such super-challenging robots in Chapter 10. You may recall the Segway from Project 3 in ROB 101 *Computational Linear Algebra*; it too is an underactuated mobile robot.

Build the linearized model:

```

1 # Linear Model for the 2-link Manipulator
2 #
3 # bring the model into the workspace for later use
4 include("dyn_mod_2LinkManipulator.jl")
5
6 #Build a function to define the parameters of the model
7
8 function modelParameters()
9     g = 9.81 # m/s^2
10    L1 = 1 # m
11    L2 = 0.7
12    m1 = 15 # kg
13    m2 = 10
14    return g, L1, L2, m1, m2
15 end
16
17 # Set equilibrium Point
18 if false
19     qe = [-pi/2, 0] #downward
20 else
21     qe = [pi/2, 0] # upward
22 end
23
24 F = dyn_mod_2LinkManipulator(qe, 0*qe)
25 D = F.D;
26 JacG = F.JacG;
27 B = F.B;
28
29 # build the linearized model about the given equilibrium point
30 A21 = -D\JacG; A21 = cleanUp(A21);

```

```

31 n, ~ = size(A21)
32 A = [zeros(n,n) I(n); A21 zeros(n,n)];
33 B = [zeros(n,n); F.D \ B];
34 display(A)
35 display(B)

```

Output

```

4x4 Matrix{Float64}:
 0.0      0.0      1.0      0.0
 0.0      0.0      0.0      1.0
 9.81    -6.54     0.0      0.0
-9.81    29.8971   0.0      0.0

4x2 Matrix{Float64}:
 0.0      0.0
 0.0      0.0
 0.0666667 -0.161905
-0.161905   0.597279

```

Using the given feedback gain matrix, K, compute the closed-loop eigenvalues

Consider a linear model $\dot{x} = Ax + Bu$, where u is the traditional notation for the control input in the Feedback Control Community. Suppose that a linear feedback $u = -Kx$ is applied to the system. Then we have

$$\begin{aligned}\dot{x} &= Ax + Bu \\ &= Ax + B \cdot (-Kx) \quad (\text{substitute in the feedback}) \\ &= (A - B \cdot K)x \quad (\text{algebra}).\end{aligned}$$

Hence, the eigenvalues of $A - B \cdot K$ determine the asymptotic stability of the closed-loop system.

```

1 K=[
2 357.82 80.57 62.084 16.8291
3 80.57 73.57 16.8291 6.92965
4 ]
5 eigen(A-B*K)

```

Output

```

Eigen{ComplexF64, ComplexF64, Matrix{ComplexF64}, Vector{ComplexF64}}
values:
4-element Vector{ComplexF64}:
-0.7071159639815681 - 0.7070975982722726im
-0.7071159639815681 + 0.7070975982722726im
-0.7071061788755746 - 0.7071073834970084im
-0.7071061788755746 + 0.7071073834970084im

vectors:
4x4 Matrix{ComplexF64}:
-0.223141-2.18102e-11im -0.223141+2.18102e-11im ... -0.233087-5.00046e-11im
 0.670976-0.0im          0.670976+0.0im           -0.667586+0.0im
 0.157786+0.157782im    0.157786-0.157782im    0.164817-0.164818im
-0.474458-0.474445im    -0.474458+0.474445im     0.472054-0.472055im

```

The closed-loop eigenvalues have been placed at $-\frac{\sqrt{2}}{2} \pm i \frac{\sqrt{2}}{2}$. Hence, the upright equilibrium is now globally exponentially stable for the linearized model.

Simulate the Nonlinear Closed-loop System

```
1 # Import useful packages
2 using DifferentialEquations
3 using Plots
4 using LaTeXStrings
5 using LinearAlgebra # for the backslash command
6 gr()
7
8
9 # bring the model into the workspace for later use
10 include("dyn_mod_2LinkManipulator.jl")
11
12 # Parameters for the Feedback controller
13
14 struct Params
15     qe::Vector{Float64}          # equilibrium point
16     K ::Matrix{Float64}          # LQR Feedback gain matrix
17 end
18 # Initialize the Params structure with specific values
19 params = Params(qe, K)
20
21
22 # define the ODE dx/dt = f(x,t)
23 #
24 # the parameters are passed via modelParameters() in the
25 # function dyn_mod_2LinkManipulator(q, dq)
26 #
27 function f(x,params,t)
28     n = floor(Int, length(x) / 2) # In Julia n/2 is a Float64
29     q = x[1:n]
30     dq = x[n+1:end]
31     K = params.K
32     qe = params.qe
33     xe = [qe; 0*qe]
34     tau = -K*(x-xe)
35     model = dyn_mod_2LinkManipulator(q, dq)
36     dx1 = dq
37     dx2 = (model.D) \ (-model.C*dq-model.G+model.B*tau) # note the use of backslash
38     dx=[dx1;dx2]
39     return dx
40 end
41
42 # Set the initial condition as a vector
43 x0= [qe;0*qe] - [pi/4; pi/4; 0; 0]
44
45 # Set the time interval
46 T = (0.0, 10)
47
48 # Setup the ODE problem with out-of-place function
49 problem = ODEProblem{false}(f, x0, T, params)
50
51 # solve the ODE problem using the Runge Kutta Tsitouras 5/4 Integrator
52 sol = solve(problem, Tsit5());
53 #sol = solve(problem, Vern9());
```

```

55
56 # plot the solution
57 p1 = plot(sol, lw=3, guidefont=20, xlabel=L"t", ylabel=L"(q_1, q_2, \dot{q}_1, \dot{q}_2)",
58         label=["q1" "q2" "dq1" "dq2"], legend=:topright)
59 display(p1)
60
61
62 # Extract state variables from the solution
63 x_sol = Array(sol)
64
65 # Preallocate the array for control signals
66 tau_sol = zeros(2, size(x_sol, 2))
67
68 # Define terms in the feedback controller
69 K = params.K
70 qe = params.qe
71 xe = [qe; 0*qe]
72
73 # Compute the control signal at each time step
74 for i in 1:size(x_sol, 2)
75     x = x_sol[:, i]
76     tau_sol[:, i] = -K * (x - xe)
77 end
78
79 # Plot the control signal
80 p2 = plot(sol.t, tau_sol', lw=3, guidefont=20, xlabel=L"t", ylabel=L"\tau = -K \cdot \dot{x} - K \cdot (x - x_e)",
81         label=["tau_1" "tau_2"], legend=:topright)
82 display(p2)

```

Output The plots were moved to the beginning of the solution.

Make a GIF

```

1 using Plots
2
3 # Robot parameters
4 g, L1, L2, m1, m2 = modelParameters()
5
6 # Function to calculate forwardKinematics
7 function forwardKinematics(q1, q2)
8     p1 = L1 * [cos(q1), sin(q1)]
9     p2 = p1 + L2 * [cos(q1 + q2), sin(q1 + q2)]
10    return p1, p2
11 end
12
13 # Extract q1 and q2 from the solution
14 q1_vals = x_sol[1, :]
15 q2_vals = x_sol[2, :]
16
17 # Create a plot for each time step and save as a GIF
18 anim = @animate for i in 1:length(q1_vals)
19     p1, p2 = forwardKinematics(q1_vals[i], q2_vals[i])
20     plot([0, p1[1], p2[1]], [0, p1[2], p2[2]], legend=false, lw=3, color=:blue,
21           xlims=(-(L1+L2), L1+L2), ylims=(-(L1+L2)/10, 1.1*(L1+L2)), aspect_ratio=1)
22     scatter!([p1[1]], [p1[2]], color=:red, markersize = m1)
23     scatter!([p2[1]], [p2[2]], color=:red, markersize = m2)

```

```

24 scatter! ([0.0], [0.0], color=:black, markersize = m1+m2)
25 end
26
27 gif (anim, "inverted2LinkRobotArm.gif", fps = 10)

```

Output Here is the **GIF**; it is the same as the one provided at the beginning of the solution. **Just Kidding!** The GIF shows what happens if one messes up the SIGN on the feedback controller gain. Instead of implementing $u = -Kx$, the sign was flipped and the feedback implemented as $u = Kx$. The resulting eigenvalues are

$$\lambda_{1,\dots,4} = \begin{bmatrix} -7.472 \\ -3.232 \\ 4.646 \\ 8.886 \end{bmatrix}.$$

In the linearized model, the states head off to infinity (in norm). In the nonlinear model, there are rotational joints, where arithmetic modulo 2π comes into play ... at least in theory. In most real robots, the limbs cannot make a full rotation without having a self-collision with another part of the robot, or, hitting a hard stop that has been installed for safety reasons. Hitting the hard stop makes a loud noise and usually breaks something. If you are fortunate enough to be working on a modern commercial robot, the system likely has a piece of software supervising the commands to the robot joints and it will shut down the robot before anything dangerous has happened.

In a feedback loop, if the sign on the gains is wrong, the closed-loop system is usually toast. Many other things can cause a feedback controller to fail, but getting the sign wrong is the most egregious error. It seems like a silly mistake, one nearly impossible to make, but all it takes is one miscommunication with a team member for it to happen, as your author experienced early in his career. And it was on me as the lead control system designer. Talk about **egg on your face** ... it's like putting your hand on a hot stove ... you do not repeat the error.

Julia Packages for Feedback Control Design

- **ControlSystems.jl:** A comprehensive control systems toolbox in Julia for system modeling, analysis, controller design, and simulation. Supports state-space and frequency domain methods.
- **RobustAndOptimalControl.jl:** Focuses on robust and optimal control, extending the capabilities of ControlSystems.jl. Includes advanced techniques like H-infinity and robust control.
- **Systems.jl:** Offers tools for modeling, simulation, and analysis of dynamical systems, essential in control design.
- **LCMCore.jl:** Useful for control system applications, particularly in communication between system components.

9.6 When in a Bind, Euler's Method is Always There for You

Yes, him again, the Einstein of his era! Euler's Method is one of the most fundamental and historically significant algorithms for computing numerical solutions to ODEs. Named after the uber-prolific Swiss mathematician Leonhard Euler, this method provides a straightforward yet powerful approach for approximating solutions to ODEs, a cornerstone in mathematical modeling of natural and engineered processes. Euler's Method is particularly renowned for its simplicity and ease of implementation, making it an excellent starting point for understanding numerical techniques in differential equations. At its core, the method transforms the often daunting task of solving complex ODEs into a `for` loop. This is achieved by approximating the slope of the solution curve at a series of points, thereby constructing a piecewise linear approximation to the solution. While Euler's Method may not boast the precision of more advanced techniques, such as **Runge-Kutta** methods, its conceptual clarity offers invaluable insights into the behavior of

differential equations. It serves as a stepping stone to more sophisticated algorithms, providing a foundational understanding of how numerical integration can bridge the gap between analytical mathematics and practical computation.

Secrets of the Arcane 9.64: Euler's Method and Katherine Johnson, an American Hero!

A super interesting video: [Katherine Johnson and Euler's Method](#) by the UCLA Modeling Class.

In the bustling halls of NASA during the 1950s through 1980s, Katherine Johnson, an African American mathematician, emerged as a pivotal figure. Amidst the Space Race, Johnson's brilliance shone, especially in her application of Euler's method, a key mathematical technique. Euler's method, akin to a mathematical magic wand, simplifies complex equations into smaller, manageable steps. Johnson wielded this tool with precision, transforming the trajectories of spacecraft into a series of calculations that even non-mathematicians could fathom. Her most notable use of Euler's method was in calculating the trajectory for John Glenn's 1962 orbital mission. Johnson's calculations were so accurate that Glenn famously insisted on her verification before his flight, saying, "If she says they're good, then I'm ready to go." By breaking down the arcane into the understandable, Johnson not only propelled rockets into space but also carved a path for future mathematicians. Her legacy, intertwined with Euler's enduring method, reminds us that complex problems often have elegant, step-by-step solutions.

Katherine Johnson graduated with highest honors in Mathematics and French from West Virginia State College in 1937. Her exceptional talent in mathematics was evident from a young age and was the foundation of her illustrious career at NASA.

Key Dates:

- **1953:** Katherine Johnson joins NASA.
- **1961:** Calculates trajectory for Mercury mission.
- **1962:** Verifies calculations for John Glenn's orbit.
- **1969:** Leads team computing orbit for Apollo 11's lunar landing.

Katherine Johnson, a mathematician who made the moon reachable, and Euler's method, her trusty guide, together remind us of the power of simplicity in solving some of the most complex problems faced by humankind. Related Videos: [NASA Trailblazer: Katherine Johnson](#) by National Geographic and [I found Katherine Johnson's actual calculations at NASA](#) by Tibebs.

Method 9.65: Euler's Method is There When You do not Have Access to Packages

Consider a vector ordinary differential equation, $\dot{x} = f(x, t)$, with initial condition $x(t_0) = x_0$. If the ODE has a solution, $x(t, t_0, x_0)$ defined on $[t_0, T]$, then it can be numerically approximated as follows:

- Choose a step time $h > 0$.
- In a for loop, compute

$$x_{n+1} = x_n + h \cdot f(x_n, t_n), \quad (9.71)$$

from $n = 0$ to $n = \text{floop}\left(\frac{T}{h}\right)$.

The process starts at t_0 with x_0 and proceeds in steps of h to approximate the solution at subsequent points. The sequence of vectors x_n is an approximation of $x(t, t_0, x_0)$ at time t_n for $t_{n+1} = t_n + h$.

Note: The global error in an ODE Solver is the error in the numerical solution after a finite number of steps. For Euler's method, the global error is proportional to the step size, making it a first-order method. Specifically, the global error is of the order $O(h)$. This means that if you halve the step size, the total accumulated error over a fixed interval is also roughly halved.

The rate at which error accumulates in Euler's method can be problematic, especially for stiff equations or equations where the solution exhibits rapid changes. The method's simplicity and ease of implementation come at the cost of accuracy and stability, particularly for larger step sizes or over longer intervals. For more accurate results, especially over longer intervals

or with more complex ODEs, higher-order methods like the [Runge-Kutta](#) methods are often preferred, as they offer better accuracy for a given step size, albeit with increased computational complexity.

Example 9.66. Use Euler's Method to solve the following ODEs:

- (a) The single-link pendulum from Examples 6.37 and 9.13.
- (b) The 2-link robot arm stabilized in the upward equilibrium from Example 9.63.

Solutions:

- (a) The single-link pendulum from Examples 6.37 and 9.13.

```

1 # Import useful packages
2 using Plots
3 using LaTeXStrings
4 gr()
5
6 # Single-link Pendulum Model
7
8 struct ParamsPendulum
9     m::Float64      # Mass in kilograms
10    L::Float64       # Length in meters
11    g::Float64       # Acceleration due to gravity in m/s^2
12 end
13
14
15 # Initialize the Params structure with specific values
16 m = 2
17 L = 1
18 g = 9.81
19 params = ParamsPendulum(m, L, g)
20
21
22 # define the ODE dx/dt = f(x,params,t)
23 function f(x, params, t)
24     dx = [x[2], -(params.g/params.L)*sin.(x[1])]
25     return dx
26 end
27
28 # Set the initial condition as a vector
29 x0= [3*pi/8; 0.0]
30
31 # Set the time interval
32 T = (0.0, 10)
33
34 # Euler's method
35 dt = 1e-2
36 N = floor(Int64, T[2]/dt)
37 x=zeros(2, N+1)
38 x[:, 1]=x0
39 t=zeros(N+1)
40 t[1] = T[1]
41 for n=1:N
42     t[n+1] = t[n]+dt
43     x[:, n+1] = x[:, n] + dt*f(x[:, n], params, t[n])

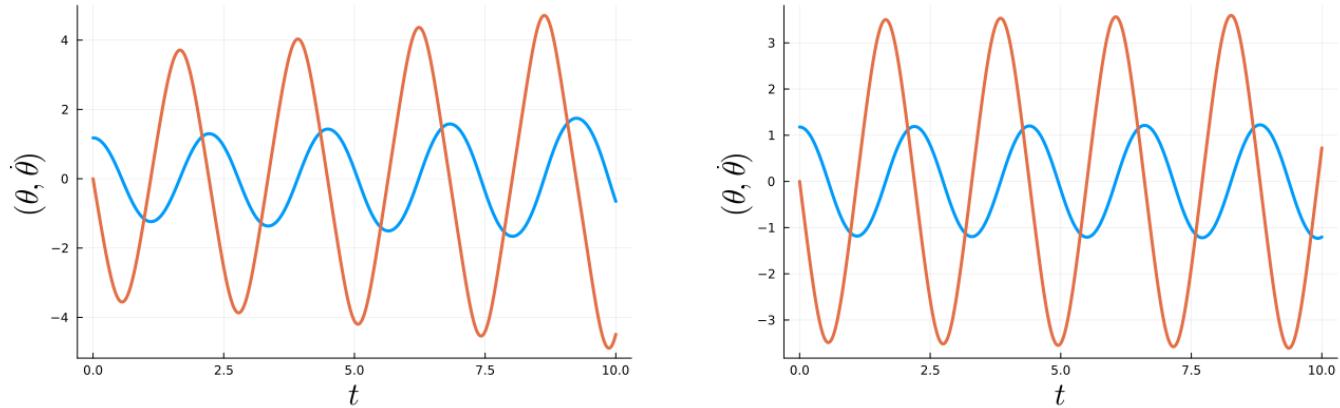
```

```

44 end
45
46 p1 = plot(t, x', lw=3, guidefont = 20, xlabel=L"t", ylabel=L"\dot{\theta}", \dot{theta}, legend=
  false)
47 display(p1)
48 png(p1, "PendulumSimulationChap09Eulerdt1eminus2")

```

Output In the left-side plot for $dt = 0.01$, the simulation does not faithfully represent a pendulum; we can see the velocity growing each swing of the pendulum. In the right-side plot, already for $dt = 0.001$, the simulation more faithfully represents a pendulum. If we look closely, we can still see the velocity growing each swing of the pendulum, but much slower. This is a drawback of Euler for “marginally stable” models, that is, models whose linearization has purely imaginary eigenvalues.



- (b) The 2-link robot arm stabilized in the upward equilibrium from Example 9.63.

```

1 # Import useful packages
2 using Plots
3 using LaTeXStrings
4 using LinearAlgebra # for the backslash command
5 gr()
6
7
8 # bring the model into the workspace for later use
9 include ("dyn_mod_2LinkManipulator.jl")
10
11 #Build a function to define the parameters of the model
12 function modelParameters()
13     g = 9.81 # m/s^2
14     L1 = 1 # m
15     L2 = 0.7
16     m1 = 15 # kg
17     m2 = 10
18     return g, L1, L2, m1, m2
19 end
20
21 # Parameters for the Feedback controller
22
23 struct Params
24     qe::Vector{Float64}          # equilibrium point
25     K ::Matrix{Float64}          # LQR Feedback gain matrix
26 end
27 # Initialize the Params structure with specific values

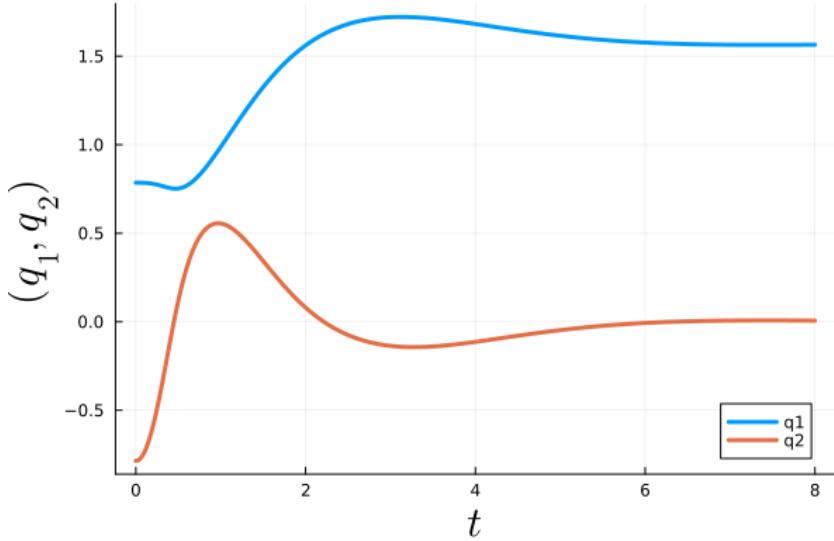
```

```

28
29 K=[ 357.82  80.57  62.084   16.8291
30      80.57  73.57  16.8291   6.92965]
31 qe = [pi/2, 0] # upward
32 params = Params(qe, K)
33
34
35 # define the ODE dx/dt = f(x,t)
36 #
37 # the parameters are passed via modelParameters() in the
38 # function dyn_mod_2LinkManipulator(q, dq)
39 #
40 function f(x,params,t)
41     n = floor(Int, length(x) / 2) # In Julia n/2 is a Float64
42     q = x[1:n]
43     dq = x[n+1:end]
44     K = params.K
45     qe = params.qe
46     xe = [qe; 0*qe]
47     tau = -K*(x-xe)
48     model = dyn_mod_2LinkManipulator(q, dq)
49     dq1 = dq
50     dx2 = (model.D) \ (-model.C*dq-model.G+model.B*tau) # note the use of backslash
51     dx=[dq1;dx2]
52     return dx
53 end
54
55 # Set the initial condition as a vector
56 x0= [qe;0*qe] - [pi/4; pi/4; 0; 0]
57
58 # Set the time interval
59 T = (0.0, 8)
60
61 # Euler's method
62 dt = 1e-2
63 N = floor(Int64,T[2]/dt)
64 x=zeros(4,N+1)
65 x[:,1]=x0
66 t=zeros(N+1)
67 t[1] = T[1]
68 for n=1:N
69     t[n+1] = t[n]+dt
70     x[:,n+1] = x[:,n] + dt*f(x[:,n],params,t[n])
71 end
72 #p1 = plot(t, x[1:2,:]', lw=3, guidefont = 20, xlabel=L"t", ylabel=L"(q_1, q_2)", legend = ("q_1", "q_2"))
73 p1 = plot(t, x[1, :], lw=3, label="q1", guidefont=20, xlabel=L"t", ylabel=L"(q_1, q_2)", )
74 plot!(p1, t, x[2, :], lw=3, label="q2")
75 display(p1)

```

Output The feedback controller in this system was designed so that the eigenvalues of the closed-loop linearized model have negative real parts. In such cases, Euler tends to do quite well. Here, we see the first link converging to the upright position, $q_1 = \frac{\pi}{2}$, and the second link aligning itself with the first link, that is, $q_2 = 0$.



Fact 9.67: Derivation of Euler's Method or How to Never Forget It

Consider an ODE $\dot{x} = f(x, t)$. Using a forward difference to approximate the derivative of $x(t)$ gives,

$$\begin{aligned} \frac{x(t+h) - x(t)}{h} &= f(x(t), t) \\ \Updownarrow \\ x(t+h) - x(t) &= h \cdot f(x(t), t) \\ \Updownarrow \\ x(t+h) &= x(t) + h \cdot f(x(t), t) \\ \Updownarrow \quad (\text{substituting in } t_{n+1} := t_n + h) \\ x(t_{n+1}) &= x(t_n) + h \cdot f(x(t_n), t_n). \end{aligned}$$

The derivation is amazingly simple and easy to re-derive at a moment's notice when you need it most.

9.7 Resonance in ODEs

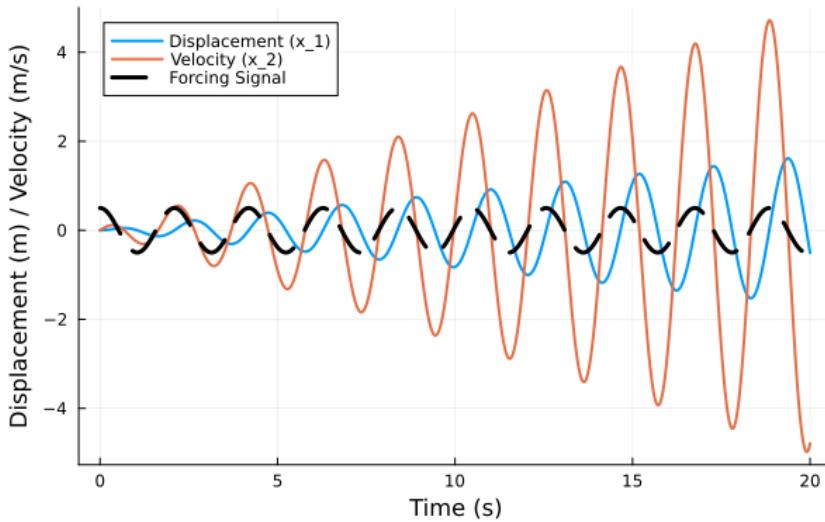
The collapse of the [Tacoma Narrows Bridge \(1940\)](#) (see also, the video [Tacoma Narrows Bridge Collapse ‘Gallopin’ Gert!](#)) in Washington State in 1940 serves as a stark reminder of the importance of considering dynamic effects in structural design. This event is often associated with the concept of resonance; however, the actual cause of the collapse is more accurately attributed to a phenomenon known as *aeroelastic flutter*, which differs from resonance.

Resonance occurs when a system oscillates at its natural frequency due to external periodic driving forces at that same frequency. Initially, it was believed that the Tacoma Narrows Bridge collapsed due to resonance caused by wind. Further analysis, however, has shown that the situation was more complex.

Aeroelastic Flutter is a self-exciting and self-sustaining oscillation caused by the interaction between aerodynamic forces, elastic properties, and inertial forces of a structure. For the Tacoma Narrows Bridge, wind provided an external force that initiated oscillations in the bridge. The bridge's design—particularly its narrow, flexible, and lightweight construction—made it exceptionally vulnerable to these aeroelastic oscillations.

Below we consider a mass-spring system, with no damping. It has a natural frequency of oscillation, $\omega_n = \sqrt{k/m}$, where k is the spring constant and m is the mass. If one gives the system a tap, it will simply oscillate at its natural frequency. In the code below, we force the system with $\cos(\omega_0 t)$, where ω_0 is chosen as the natural frequency of the system. This creates a **resonance response** that grows like $t \cos(\omega_0 t + \theta_0)$. The increasing amplitude with time will eventually cause any real-world system to reach its physical limits and break, similar to what aeroelastic flutter did to the [Tacoma Narrows Bridge](#).

Resonance in a Forced Oscillator



```

1 using DifferentialEquations, Plots
2
3 # Constants
4 m = 1.0    # Mass (kg)
5 k = 9.0    # Spring constant (N/m)
6 F0 = 0.5   # Amplitude of the external force (N)
7 w_n = sqrt(k / m)  # Natural frequency (rad/s)
8 w = w_n    # Forcing frequency = Natural frequency for resonance
9
10 # ODE Definition
11 function forced_oscillator!(du, u, p, t)
12     x, v = u
13     du[1] = v
14     du[2] = (F0*cos(w*t) - k*x) / m
15 end
16
17 # Initial conditions
18 u0 = [0.0, 0.0]  # Initial displacement and velocity
19 tspan = (0.0, 20.0)  # Time span
20
21 # Solve the ODE
22 prob = ODEProblem(forced_oscillator!, u0, tspan)
23 sol = solve(prob, Tsit5(), reltol=1e-8, abstol=1e-8)
24
25 # Create a finer time grid for interpolation
26 fine_t = range(tspan[1], tspan[2], length=10000) # Increased number of points for finer
27     interpolation
28
29 # Interpolate the solution at the finer time grid
30 fine_x = [sol(t)[1] for t in fine_t]
31 fine_v = [sol(t)[2] for t in fine_t]
32
33 # Forcing signal
34 forcing_signal = F0*cos.(w.*fine_t)
35
36 p1 = plot(fine_t, fine_x, label="Displacement (x_1)", xlabel="Time (s)", ylabel="Displacement (m) /
37     Velocity (m/s)",

```

```

37     title="Resonance in a Forced Oscillator", linewidth=2)
38 plot!(fine_t, fine_v, label="Velocity (x_2)", linewidth=2) # Smooth velocity curve
39 plot!(fine_t, forcing_signal, label="Forcing Signal", linestyle=:dash, color=:black, lw=3) # Forcing
    signal
40
41 # Display the plot
42 plot!(legend=:topleft)

```

Output See the above plot.

(Optional Read:) Videos for Enhancing Your Understanding and Intuition:

- [A better description of resonance](#) by Steve Mould. He illustrates wave patterns via the heights of flames! If you believe that fire makes anything better, then this video is for you. It's a bit slow initially, so you may want to jump ahead and then return to the beginning after you've seen the punchline.
- [Resonance and the Sounds of Music](#) by Walter Lewin. Very old-school set up by an absolute master with 1.64M subscribers.
- [Resonance Introduction using 9 Demonstrations](#) by Flipping Physics (so nerdy that the humor enhances understanding).

9.8 (Optional Read:) More on Julia Packages for Solving ODEs and PDEs

In Julia, the ecosystem for solving differential equations, including ordinary differential equations (ODEs) and partial differential equations (PDEs), is extensive and under continuous development. The primary effort that began under the umbrella of “JuliaDiffE” has since been integrated into a larger, more comprehensive initiative known as the [SciML \(Scientific Machine Learning\)](#) ecosystem. This initiative aims to provide a cohesive suite of tools for scientific computing, which includes solving differential equations, machine learning in the context of scientific problems, and much more.

- **For ODEs:**
 - **DifferentialEquations.jl:** A flagship package within the SciML ecosystem offering a unified interface for solving a broad spectrum of differential equations, including ODEs, PDEs, stochastic differential equations (SDEs), and more. It is highly flexible and features many algorithms for solving ODEs.
 - **OrdinaryDiffEq.jl:** A part of the DifferentialEquations.jl suite focused specifically on ODEs, providing a comprehensive array of solvers and analysis tools.
- **For PDEs:**
 - **DifferentialEquations.jl:** Also acts as the go-to tool for PDEs, offering methods to discretize and solve PDEs either directly or through conversion into ODEs or SDEs.
 - **FEniCS.jl:** A Julia wrapper for the FEniCS project, an acclaimed open-source computing platform for solving PDEs using finite element methods, enabling Julia users to utilize FEniCS’s potent PDE-solving capabilities.
 - **Gridap.jl:** Provides a comprehensive toolkit for the simulation of PDEs using the finite element method (FEM), offering a powerful and adaptable framework for complex multiphysics simulations.
 - **JuAFEM.jl:** Aimed at solving finite element method problems, focusing on a simple and intuitive interface for Julia users.
- **Other Notable Packages:**
 - **ModelingToolkit.jl:** Part of the SciML ecosystem designed to ease the modeling of scientific problems, including differential equations creation and manipulation, often used alongside DifferentialEquations.jl.
 - **ParameterizedFunctions.jl:** Facilitates defining systems of differential equations with additional features like symbolic derivatives and parameter estimation, enhancing DifferentialEquations.jl’s functionality.
 - **PDEOperators.jl:** Offers tools for defining differential operators for PDE discretization, simplifying work with complex differential operators.

- **ApproxFun.jl:** Though not exclusively for differential equations, this package enables function approximation and can be used for solving differential equations.

These packages are integral to a broader, highly integrated ecosystem that includes tools for optimization, machine learning, and beyond, allowing for advanced problem-solving approaches across various scientific fields. The Julia community actively develops and maintains these packages, ensuring regular updates and expansions to the ecosystem.

9.9 (Optional Read:) Proofs Associated with the Chapter

Proposition 9.68: Real Basis Vectors from Eigenvectors

Suppose that A is a real $n \times n$ matrix. Then we know its eigenvalues and eigenvectors are either real or they occur in complex conjugate pairs. Let $\lambda_1, \dots, \lambda_p, \lambda_1^*, \dots, \lambda_p^*$ be the complex eigenvalues and let $\lambda_{2p+1}, \dots, \lambda_n$ be the real eigenvalues, with a similar enumeration for the corresponding eigenvectors $\{v_1, \dots, v_p, v_1^*, \dots, v_p^*, v_{2p+1}, \dots, v_n\}$. If the eigenvalues are distinct, that is $\lambda_i \neq \lambda_j$ for $i \neq j$, then the set of n real vectors

$$\{v_{R1}, \dots, v_{Rp}, v_{I1}, \dots, v_{Ip}, v_{2p+1}, \dots, v_n\}$$

is linearly independent and forms a basis of \mathbb{R}^n .

Moreover, if $\lambda = \lambda_R + i\lambda_I$ and $v = v_R + iv_I$ are a complex eigenvalue and eigenvector, then

$$\begin{aligned} Av_R &= \lambda_R v_R - \lambda_I v_I \\ Av_I &= \lambda_I v_R + \lambda_R v_I. \end{aligned} \tag{9.72}$$

Proof: Chapter 2.7 of the ROB 501 textbook, **Mathematics for Robotics**, establishes that the eigenvalues being distinct implies that the set of eigenvectors

$$\{v_1, \dots, v_p, v_1^*, \dots, v_p^*, v_{2p+1}, \dots, v_n\}$$

forms a basis of $(\mathbb{C}^n, \mathbb{C})$. Hence,

$$\dim \text{span}\{v_1, \dots, v_p, v_1^*, \dots, v_p^*, v_{2p+1}, \dots, v_n\} = n,$$

the dimension of $(\mathbb{C}^n, \mathbb{C})$.

By the definition of the real and imaginary parts of vectors, for each $1 \leq j \leq p$, $v_j = v_{Rj} + iv_{Ij}$ and $v_j^* = v_{Rj} - iv_{Ij}$, and moreover, $v_{Rj} = \frac{v_j + v_j^*}{2}$ and $v_{Ij} = \frac{v_j - v_j^*}{2i}$. Hence,

$$\text{span}\{v_j, v_j^*\} = \text{span}\{v_{Rj}, v_{Ij}\},$$

and

$$\text{span}\{v_1, \dots, v_p, v_1^*, \dots, v_p^*, v_{2p+1}, \dots, v_n\} = \{v_{R1}, \dots, v_{Rp}, v_{I1}, \dots, v_{Ip}, v_{2p+1}, \dots, v_n\}.$$

It follows that

$$\dim \text{span}\{v_{R1}, \dots, v_{Rp}, v_{I1}, \dots, v_{Ip}, v_{2p+1}, \dots, v_n\} = n,$$

establishing the vectors are linearly independent in $(\mathbb{C}^n, \mathbb{C})$, and hence also in $(\mathbb{R}^n, \mathbb{R})$.

To establish (9.72), we use

$$\begin{aligned} Av &= A(v_R + iv_I) = (\lambda_R + i\lambda_I) \cdot (v_R + iv_I) = \lambda v \\ Av^* &= A(v_R - iv_I) = (\lambda_R - i\lambda_I) \cdot (v_R - iv_I) = \lambda^* v^*. \end{aligned}$$

This gives

$$\begin{aligned} Av_R &= A \frac{v + v^*}{2} = \lambda_R v_R - \lambda_I v_I \\ Av_I &= A \frac{v - v^*}{2i} = \lambda_I v_R + \lambda_R v_I. \end{aligned}$$

An alternative proof is

$$\begin{aligned} Av_R &= \text{real}(Av) = \text{real}((\lambda_R + i\lambda_I) \cdot (v_R + iv_I)) = \lambda_R v_R - \lambda_I v_I \\ Av_I &= \text{imag}(Av) = \text{imag}((\lambda_R + i\lambda_I) \cdot (v_R + iv_I)) = \lambda_I v_R + \lambda_R v_I. \end{aligned}$$

■

Proposition 9.69: $e^{At}v_R$ and $e^{At}v_I$

$$\begin{aligned} e^{At}v_R &= e^{\lambda_R t} \cos(\lambda_I t) v_R - e^{\lambda_R t} \sin(\lambda_I t) v_I \\ e^{At}v_I &= e^{\lambda_R t} \sin(\lambda_I t) v_R + e^{\lambda_R t} \cos(\lambda_I t) v_I. \end{aligned} \quad (9.73)$$

Proof: As in Prop. 9.68, we write $\lambda = \lambda_R + i\lambda_I$ and $v = v_R + iv_I$. Then,

$$\begin{aligned} e^{At}v_R &= \text{real}(e^{At}v) \\ &= \text{real}(e^{\lambda t}v) \\ &= \text{real}(e^{\lambda_R t}(\cos(\lambda_I t) + i\sin(\lambda_I t))(v_R + iv_I)) \\ &= e^{\lambda_R t} \cos(\lambda_I t) v_R - e^{\lambda_R t} \sin(\lambda_I t) v_I \end{aligned}$$

Similarly,

$$\begin{aligned} e^{At}v_I &= \text{imag}(e^{At}v) \\ &= \text{imag}(e^{\lambda t}v) \\ &= \text{imag}(e^{\lambda_R t}(\cos(\lambda_I t) + i\sin(\lambda_I t))(v_R + iv_I)) \\ &= e^{\lambda_R t} \sin(\lambda_I t) v_R + e^{\lambda_R t} \cos(\lambda_I t) v_I \end{aligned}$$

■

From Prop. 9.35 Linearization of the Robot Equations, we have that the linearized model about any equilibrium point has the form

$$\begin{bmatrix} \delta\dot{x}_1 \\ \delta\dot{x}_2 \end{bmatrix} = \begin{bmatrix} 0_n & I_n \\ A_{21} & 0_n \end{bmatrix} \begin{bmatrix} \delta x_1 \\ \delta x_2 \end{bmatrix}. \quad (9.74)$$

Proposition 9.70: Eigen-Facts about the Robot Equations

Suppose the matrices A and B are invertible and that matrix C is defined as

$$C = \begin{bmatrix} 0 & B \\ A & 0 \end{bmatrix}.$$

Then,

$$\det(C - \lambda I) = \det \begin{bmatrix} -\lambda I & B \\ A & -\lambda I \end{bmatrix} = \det(\lambda^2 I - AB).$$

Therefore, the eigenvalues of C are the square roots of the eigenvalues of AB . So, they depend on the eigenvalues of AB rather than the individual eigenvalues of the matrices A and B .

Note: If AB has a negative real eigenvalue, then C will have a pair of purely imaginary eigenvalues. If AB has a positive

real eigenvalue, then C will have both a positive and negative eigenvalue. If AB has a pair of complex conjugate eigenvalues, then C will have two complex conjugate pairs of eigenvalues.

Proof: Based on [Eigenvalues of block matrices with zero diagonal blocks](#). Given two vectors $u, v \in \mathbb{R}^n$, let $w = \begin{bmatrix} u \\ v \end{bmatrix} \in \mathbb{R}^{2n}$ be the vector whose first n components are those of u and whose last n are those of v . The matrix-vector product $C \cdot w$ gives

$$\begin{bmatrix} 0 & B \\ A & 0 \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} Av \\ Bu \end{bmatrix}.$$

Therefore, solving $Cw = \lambda w$ reduces to solving the following two matrix equations,

$$\begin{cases} Av = \lambda u \\ Bu = \lambda v. \end{cases}$$

Since A and B are invertible, it follows that $v = \lambda A^{-1}u$ and $u = \lambda B^{-1}v$:

$$\begin{cases} Av = \lambda \cdot \lambda B^{-1}v \\ Bu = \lambda \cdot \lambda A^{-1}u \end{cases} \rightarrow \begin{cases} BA v = \lambda^2 v \\ AB u = \lambda^2 u \end{cases}$$

We conclude that $\lambda \in \mathbb{C}$ is an eigenvalue of C if, and only if, λ^2 is an eigenvalue of AB (or BA , since they have the same eigenvalues). ■

Proof of Prop. 9.53 (Properties of Complex Exponentials)

Complex exponentials are fundamental in various fields of mathematics, engineering, and physics. Below are some of the key properties of complex exponentials, for $a, \omega \in \mathbb{R}$ and $z := a + i\omega \in \mathbb{C}$:

- (a) **Periodicity:** The function $e^{i\omega t}$ is periodic with period $T = \frac{2\pi}{|\omega|}$, for $\omega \neq 0$.
- (b) **Multiplicative Property:** For any two complex numbers $z_1 = a_1 + i\omega_1$ and $z_2 = a_2 + i\omega_2$, the exponential function satisfies $e^{z_1 t} \cdot e^{z_2 t} = e^{(z_1+z_2)t}$.
- (c) **Inverse:** The inverse of e^{zt} is given by e^{-zt} .
- (d) **Differentiation and Integration:** The derivative and antiderivative of e^{zt} with respect to t are:
 - $\frac{d}{dt} e^{zt} = ze^{zt}$ and
 - for $z \neq 0$, $\int e^{zt} dt = \frac{1}{z} e^{zt} + C$, where C is the constant of integration.
- (e) **Conjugation:** The complex conjugate of $e^{(a+i\omega)t}$ is $e^{(a-i\omega)t}$.
- (f) **Magnitude:** The magnitude of $e^{(a+i\omega)t}$ is $|e^{(a+i\omega)t}| = e^{at}$, which exhibits exponential growth or decay based on the sign of the real part, a .

Proof:

- (a) **Periodicity:** By Euler's formula, $e^{i\omega t} = \cos(\omega t) + i\sin(\omega t)$. Because $\cos(\omega(t + \frac{2\pi}{|\omega|})) = \cos(\omega t + \text{sign}(\omega)2\pi) = \cos(\omega t)$ and $\sin(\omega(t + \frac{2\pi}{|\omega|})) = \sin(\omega t + \text{sign}(\omega)2\pi) = \sin(\omega t)$, it follows that $e^{i\omega(t + \frac{2\pi}{|\omega|})} = e^{i\omega t}$, proving periodicity with period $\frac{2\pi}{|\omega|}$.

(b) **Multiplicative Property:** Let $z_1 = a_1 + i\omega_1$ and $z_2 = a_2 + i\omega_2$. Then

$$\begin{aligned} e^{z_1 t} &= e^{a_1 t} \cos(\omega_1 t) + i e^{a_1 t} \sin(\omega_1 t) \\ e^{z_2 t} &= e^{a_2 t} \cos(\omega_2 t) + i e^{a_2 t} \sin(\omega_2 t). \end{aligned}$$

Therefore,

$$\begin{aligned} e^{z_1 t} \cdot e^{z_2 t} &= (e^{a_1 t} \cos(\omega_1 t) + i e^{a_1 t} \sin(\omega_1 t)) \cdot (e^{a_2 t} \cos(\omega_2 t) + i e^{a_2 t} \sin(\omega_2 t)) \\ &= e^{a_1 t} \cdot e^{a_2 t} (\cos(\omega_1 t) + i \sin(\omega_1 t)) \cdot (\cos(\omega_2 t) + i \sin(\omega_2 t)) \\ &= e^{a_1 t} \cdot e^{a_2 t} (\cos(\omega_1 t) \cos(\omega_2 t) - \sin(\omega_1 t) \sin(\omega_2 t)) + i e^{a_1 t} \cdot e^{a_2 t} (\cos(\omega_1 t) \sin(\omega_2 t) + \sin(\omega_1 t) \cos(\omega_2 t)) \\ &= e^{a_1 t} \cdot e^{a_2 t} \cos((\omega_1 + \omega_2)t) + i e^{a_1 t} \cdot e^{a_2 t} \sin((\omega_1 + \omega_2)t) \\ &= e^{(a_1 + a_2)t + (\omega_1 + \omega_2)t} \\ &= e^{(z_1 + z_2)t}. \end{aligned}$$

(c) **Inverse:** From (b), $e^{zt} \cdot e^{-zt} = e^{(z-z)t} = e^0 = 1$, and hence, e^{-zt} is the multiplicative inverse of e^{zt} .

(d) **Differentiation and Integration:** The most important observation is that by using Euler's formula, we are only differentiating or integrating real functions.

Suppose $z = a + i\omega$, and write $e^{zt} = e^{at} (\cos(\omega t) + i \sin(\omega t)) = e^{at} \cos(\omega t) + i e^{at} \sin(\omega t)$.

- By the product rule and chain rule,

$$\begin{aligned} \frac{d}{dt} e^{zt} &= \frac{d}{dt} (e^{at} \cos(\omega t)) + i \frac{d}{dt} (e^{at} \sin(\omega t)) \\ &= (a e^{at} \cos(\omega t) - \omega e^{at} \sin(\omega t)) + i(a e^{at} \sin(\omega t) + \omega e^{at} \cos(\omega t)) \\ &= a e^{at} (\cos(\omega t) + i \sin(\omega t)) + i \omega e^{at} (\cos(\omega t) + i \sin(\omega t)) \\ &= (a + i\omega) \cdot e^{at} (\cos(\omega t) + i \sin(\omega t)) \\ &= z e^{zt}. \end{aligned}$$

- We assume $z \neq 0$. $\int e^{zt} dt = \int e^{at} \cos(\omega t) dt + i \int e^{at} \sin(\omega t) dt$. We note that if one of a or ω is equal to zero, then the two integrals are trivial to evaluate; we leave this case to the learner and assume, therefore, that both a and ω are non-zero. Each of the two integrals requires two applications of integration by parts! We show all the steps for the first integral.

The integration by parts formula is given by,

$$\int u dv = uv - \int v du.$$

We set $u = \cos(\omega t)$ and $dv = e^{at} dt$ and compute

- 1) $du = -\omega \sin(\omega t) dt$, and
- 2) $v = \int e^{at} dt = \frac{1}{a} e^{at}$, given $a \neq 0$.

Using the integration by parts formula, we have,

$$\begin{aligned} \int e^{at} \cos(\omega t) dt &= uv - \int v du \\ &= \left(\cos(\omega t) \cdot \frac{1}{a} e^{at} \right) - \int \left(\frac{1}{a} e^{at} \right) (-\omega \sin(\omega t)) dt \\ &= \frac{1}{a} e^{at} \cos(\omega t) + \frac{\omega}{a} \int e^{at} \sin(\omega t) dt. \end{aligned}$$

The second integral $\int e^{at} \sin(\omega t) dt$ can be solved similarly by integration by parts. Let's set $u = \sin(\omega t)$ and $dv = e^{at} dt$, then,

3) $du = \omega \cos(\omega t) dt$, and

4) $v = \frac{1}{a} e^{at}$.

Applying integration by parts again, we have,

$$\int e^{at} \sin(\omega t) dt = \frac{1}{a} e^{at} \sin(\omega t) - \frac{\omega}{a} \int e^{at} \cos(\omega t) dt,$$

where we recognize $\int e^{at} \cos(\omega t) dt$ as the original integral.

Solving for the Original Integral: Let $I = \int e^{at} \cos(\omega t) dt$. From the steps above, we have,

$$\begin{aligned} I &= \frac{1}{a} e^{at} \cos(\omega t) + \frac{\omega}{a} \left(\frac{1}{a} e^{at} \sin(\omega t) - \frac{\omega}{a} I \right) \\ &= \frac{1}{a} e^{at} \cos(\omega t) + \frac{\omega}{a^2} e^{at} \sin(\omega t) - \frac{\omega^2}{a^2} I. \end{aligned}$$

Solving for I , we find,

$$\begin{aligned} I + \frac{\omega^2}{a^2} I &= \frac{1}{a} e^{at} \cos(\omega t) + \frac{\omega}{a^2} e^{at} \sin(\omega t) \\ &\Downarrow \\ I \left(1 + \frac{\omega^2}{a^2} \right) &= \frac{e^{at}}{a} \cos(\omega t) + \frac{\omega}{a^2} e^{at} \sin(\omega t) \\ &\Downarrow \\ I &= \frac{e^{at}}{a + \frac{\omega^2}{a}} \left(\cos(\omega t) + \frac{\omega}{a} \sin(\omega t) \right) \\ &\Downarrow \\ I &= \frac{ae^{at}}{a^2 + \omega^2} \cos(\omega t) + \frac{\omega e^{at}}{a^2 + \omega^2} \sin(\omega t). \end{aligned}$$

Therefore, the antiderivative of $e^{at} \cos(\omega t)$ is,

$$\int e^{at} \cos(\omega t) dt = \frac{ae^{at}}{a^2 + \omega^2} \cos(\omega t) + \frac{\omega e^{at}}{a^2 + \omega^2} \sin(\omega t) + C,$$

where C is the constant of integration.

In a similar manner, the antiderivative of $e^{at} \sin(\omega t)$ is given by,

$$\int e^{at} \sin(\omega t) dt = \frac{e^{at}}{a^2 + \omega^2} (a \sin(\omega t) - \omega \cos(\omega t)) + C,$$

where C is the constant of integration.

Now, we must combine the two antiderivatives to obtain the antiderivative for the original function.

For the function $f(t) = e^{at} \cos(\omega t) + i e^{at} \sin(\omega t)$, the antiderivative is the sum of the two antiderivatives computed above. Thus,

$$\begin{aligned} \int f(t) dt &= \int e^{at} \cos(\omega t) dt + i \int e^{at} \sin(\omega t) dt \\ &= \frac{e^{at}}{a^2 + \omega^2} (a \cos(\omega t) + \omega \sin(\omega t)) + i \frac{e^{at}}{a^2 + \omega^2} (a \sin(\omega t) - \omega \cos(\omega t)) + C \\ &= \frac{e^{at}}{a^2 + \omega^2} a (\cos(\omega t) + i \sin(\omega t)) - i \frac{e^{at}}{a^2 + \omega^2} \omega (\cos(\omega t) + i \sin(\omega t)) + C \\ &= \frac{e^{at}}{a^2 + \omega^2} (a - i \omega) (\cos(\omega t) + i \sin(\omega t)) + C \text{ (after combining like terms)}. \end{aligned}$$

Recalling Euler's formula, $e^{i\omega t} = \cos(\omega t) + i \sin(\omega t)$. Substituting this into our expression gives,

$$\int f(t) dt = \frac{e^{at}}{a^2 + \omega^2} (a - i\omega) e^{i\omega t} + C.$$

Noticing that $a^2 + \omega^2 = (a + i\omega)(a - i\omega)$, $a + i\omega = z$, and $e^{zt} = e^{at} \cdot e^{i\omega t}$, we can rewrite the antiderivative as,

$$\int f(t) dt = \frac{e^{at} \cdot e^{i\omega t} \cdot (a - i\omega)}{(a + i\omega)(a - i\omega)} + C = \frac{1}{z} e^{zt} + C,$$

when z is not zero.

(e) **Conjugation:** Because $e^{(a+i\omega)t} = e^{at} (\cos(\omega t) + i \sin(\omega t))$, the complex conjugate of $e^{(a+i\omega)t}$ is $e^{at} (\cos(\omega t) - i \sin(\omega t)) = e^{at} (\cos(-\omega t) + i \sin(-\omega t)) = e^{(a-i\omega)t}$.

(f) **Magnitude:** The magnitude of $e^{(a+i\omega)t} = e^{at} (\cos(\omega t) + i \sin(\omega t))$ is

$$|e^{at} \cos(\omega t) + i e^{at} \sin(\omega t)| = \sqrt{(e^{at} \cos(\omega t))^2 + (e^{at} \sin(\omega t))^2} = \sqrt{(e^{at})^2 \cdot (\cos^2(\omega t) + \sin^2(\omega t))} = e^{at}. \blacksquare$$

Proof of Prop. 9.59 (Exponential Stability of Linear Systems of ODEs)

Let A be a real $n \times n$ matrix with eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$, and consider the linear system of ODEs,

$$\dot{x} = Ax.$$

If $\text{real}(\lambda_i) < 0$ for $1 \leq i \leq n$ (aka, the real parts of all eigenvalues are negative), then the **origin** is globally exponentially stable.

Proof: We give the proof for the case of distinct eigenvalues. By Prop. 9.68, set of n real vectors

$$\{v_{R1}, \dots, v_{Rp}, v_{I1}, \dots, v_{Ip}, v_{2p+1}, \dots, v_n\}$$

is linearly independent and forms a basis of \mathbb{R}^n . Hence, an arbitrary $x_0 \in \mathbb{R}^n$ can be written as a linear combination

$$x_0 = \sum_{j=1}^p (\alpha_j v_{Rj} + \beta_j v_{Ij}) + \sum_{j=2p+1}^n \gamma_j v_j.$$

Next, applying Prop. 9.69, we have

$$\begin{aligned} e^{At} x_0 &= \sum_{j=1}^p (\alpha_j e^{At} v_{Rj} + \beta_j e^{At} v_{Ij}) + \sum_{j=2p+1}^n \gamma_j e^{At} v_j \\ &= \sum_{j=1}^p (\alpha_j e^{At} v_{Rj} + \beta_j e^{At} v_{Ij}) + \sum_{j=2p+1}^n \gamma_j e^{At} v_j \\ &= \sum_{j=1}^p (\alpha_j (e^{\lambda_{Rj}t} \cos(\lambda_{Ij}t) v_{Rj} - e^{\lambda_{Rj}t} \sin(\lambda_{Ij}t) v_{Ij}) + \beta_j (e^{\lambda_{Rj}t} \sin(\lambda_{Ij}t) v_{Rj} + e^{\lambda_{Rj}t} \cos(\lambda_{Ij}t) v_{Ij})) + \sum_{j=2p+1}^n \gamma_j e^{\lambda_j t} v_j \\ &= \sum_{j=1}^p e^{\lambda_{Rj}t} (\alpha_j (\cos(\lambda_{Ij}t) v_{Rj} - \sin(\lambda_{Ij}t) v_{Ij}) + \beta_j (\sin(\lambda_{Ij}t) v_{Rj} + \cos(\lambda_{Ij}t) v_{Ij})) + \sum_{j=2p+1}^n \gamma_j e^{\lambda_j t} v_j. \end{aligned}$$

From this, we see that if all eigenvalues have negative real parts, then $x(t) := e^{At} x_0 \rightarrow 0_{n \times 1}$ as $t \rightarrow \infty$, completing the proof. ■

Note: Suppose one of the eigenvalues has real part that is non-negative. Then by taking x_0 to be the (real part of) the corresponding eigenvector, we then have a solution that does not converge to zero. It will either oscillate, stay constant, or explode. ■

Chapter 10

Laplace Transforms through the Lens of Feedback Control

Learning Objectives

- To introduce the concept and mathematical theory of Laplace Transforms.
- To become familiar with common Laplace transforms and their properties.
- To understand the process and methods of finding inverse Laplace transforms.
- To apply Laplace transforms for solving linear ordinary differential equations (ODEs).
- To define and derive transfer functions for linear time-invariant systems.
- To understand the importance and application of transfer functions in the design of Single-Input Single-Output (SISO) feedback loops.
- To use appropriate software tools for each of the above topics.

Outcomes

- Gaining a foundational understanding of Laplace transforms and their application in engineering problems.
- Mastery of standard Laplace transforms of common functions and their key properties.
- Ability to understand inverse Laplace transforms using techniques such as partial fraction decomposition and tools in Julia.
- Proficiency in applying Laplace transforms to simplify and solve linear ODEs with initial conditions.
- Comprehensive understanding of transfer functions, their derivation, and significance in system analysis.
- Applying knowledge of transfer functions to design and analyze SISO feedback loops in engineering systems, particularly in robotics.

The Laplace Transform Magically Turns Linear Differential Equations into Algebraic Equations. One of the most remarkable aspects of Laplace Transforms lies in their ability to transmute the complexities of linear ODEs into the more tractable realm of algebra. This metamorphosis is not merely a theoretical curiosity but an immensely practical tool that unlocks a new perspective on ODEs. By translating the differential equations into algebraic ones, Laplace Transforms simplify the process of understanding solutions of ODEs, showing how to modify the equations to achieve asymptotic stability. This technique allows for a more straightforward analysis and design of systems governed by ODEs, thereby bridging the gap between complex mathematical theory and real-world application.

10.1 Introduction

Thanks to Chapter 9, we now understand what an ODE is, how to compute its solution for given initial conditions, and whether its solutions asymptotically converge to an equilibrium point or not. That is already a big accomplishment. We now make a radical change in focus from analyzing solutions of ODEs to actively modifying their solutions so that they exhibit properties that we deem desirable. To achieve this, we introduce controlled inputs into our ODE models of physical systems. These inputs serve as a mechanism to influence and steer the evolution of the system's state. By strategically manipulating these inputs, we can guide the system's trajectory toward our objectives, such as achieving exponential stability, driving an important signal to a desired value (e.g., setting a robot's speed), or attenuating the effects of perturbations. This approach marks a significant shift from passive analysis to active control, opening the door to designing systems that respond effectively to commanded inputs and achieve specific, desired outcomes.

The fascinating possibilities of feedback control are introduced within the context of linear Single-Input Single-Output (SISO) systems, utilizing Laplace Transforms and transfer functions. We employ the Segway model as a practical example to demonstrate the foundational principles of feedback control, while also touching upon concepts like PD controllers and pre-compensators. Through our exploration of Laplace Transforms and transfer functions, we aim to provide an initial framework for understanding how these tools can be used to design basic yet effective feedback control strategies. The chapter builds upon the concepts of stability and feedback introduced in Chapter 9 and Example 9.63, illustrating the impact of these tools in engineering. Our principal goal is to provide the concepts and methods required for stabilizing the BallBot mobile robot used in ROB 311 Build Robots and Make Them Move.

It's important to note, however, that the material in this Chapter, while illustrative and intriguing, remains an introductory overview and is not a substitute for more advanced courses such as Michigan's EECS 460 Control Systems Analysis and Design, ME 461 Automatic Control, ROB 415 Robot Controls, or any of the other amazing Control Courses listed [here](#).



Figure 10.1: An Input-Output System. The signal $u(t)$ represents the action of an **actuator**, such as a DC motor, a rocket engine, or a voltage supply. The output is a signal that is important to the system designer, such as the speed of a robot or the angle of a joint. If $u(t)$ and $y(t)$ are scalar-valued, then one has a Single-Input, Single-Output (SISO) System. If $u(t)$ or $y(t)$ are vectors, then one has a Multi-Input Multi-Output (MIMO) System.

10.2 Single-Input Single-Output (SISO) Linear Systems

A system is an ODE with both inputs and outputs. A drone (aka, quadrotor) is a system. A system's inputs are signals that cause changes in the motion of the entire system or parts of the system. A quadrotor-drone has four DC motors with propellers attached to them. These are the actuators. The output of the drone could be its 3D position and its 3D orientation in \mathbb{R}^3 . A simpler output would be its speed in a single direction, say along the x -axis. It would be up to you (or your team, or your boss) to define the important quantities to be regulated in the system. Here, we define linear systems with a single input and output. It's THE PLACE to start. The planar Segway model that will guide our discussion has one input (a motor at the wheels) and two outputs (the angle of the Segway's body and its speed) that we will learn to treat successively. Controlling two outputs with one input is more challenging than controlling one output. C'est la vie, bébé! That's Robotics in a nutshell.

10.2.1 Input-Output and State-Variable Models

Every domain has its “secret codes and rituals”. In **Feedback Control**, the generic notation for an input is $u(t)$, and for an output, $y(t)$. We will use this notation so that you are not caught off guard in a follow-on course.

SISO Linear Systems

Definition 10.1.

$$y^{(n)} + a_{n-1}y^{(n-1)} + a_{n-2}y^{(n-2)} + \dots + a_1\dot{y} + a_0y = b_0u + b_1\dot{u} + \dots + b_mu^{(m)} \quad (10.1)$$

is a **linear input-output system**. It is called an input-output system because the only variables present are the input, the output, and their derivatives. Because we will always assume that both $u(t)$ and $y(t)$ are scalars, it is a **SISO linear system**.

Note: We are taking the coefficients a_k to be constants. In this case, the model is called a **Linear Time-Invariant (LTI)** input-output system. If the coefficients depend on time, then one has a **Linear Time-Varying (LTV)** input-output system. Even grad students struggle with LTV models, so don’t worry, we’re not going there.

Definition 10.2. Consider a linear system of ODEs, $\dot{x} = Ax$, with $x \in \mathbb{R}^n$. It becomes a **control system** if we add an input and an output to the model. In particular, for b , an $n \times 1$ column vector and c , a $1 \times n$ row vector,

$$\begin{aligned} \dot{x} &= Ax + bu \\ y &= cx \end{aligned} \quad (10.2)$$

is a **SISO linear state variable model**. While studying MIMO linear state variable models is super interesting, we will not do that.

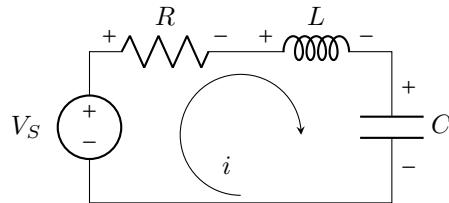
Note: In a course totally dedicated to feedback control, you learn that the models (10.1) and (10.2) are two sides of the same coin. So why do we give both forms? Because the first one makes it very easy to motivate **transfer functions**, and the second one is how we will obtain linear control models for the Segway and the BallBot.

Note: An input-output system is **causal** if its output at time t does not depend upon future values of the input signal. It is a non-trivial fact that the system in (10.1) is causal if, and only if, $m \leq n$, that is, the highest derivative appearing in the input is less than or equal to the highest derivative appearing in the output. The state-variable model (10.2) is always causal.

We illustrate how these models arise in circuits and mechanics.

10.2.2 Input-Output Models from Circuits

Example 10.3. We rework the RLC circuit from Example 9.27 to derive a SISO model instead of a state-variable model, using $y := V_C$ and $u := V_S$.



Solution: Ans. $\frac{d^2y}{dt^2} + \frac{R}{L}\frac{dy}{dt} + \frac{1}{LC}y = \frac{1}{LC}u$.

As before, V_S is a voltage source, R denotes a resistor, L an inductor, and C a capacitor. The loop denotes the reference direction for the current flowing through the circuit, $i(t)$. **Kirchoff's Voltage Law** says that the sum of the voltages around any loop must equal zero, yielding

$$-V_S + V_R + V_L + V_C = 0. \quad (10.3)$$

As before, the derivatives arise from two of the three **constitutive relations** for the circuit elements, namely

$$\begin{aligned} V_R &= R \cdot i \\ V_L &= L \frac{di}{dt} \\ i &= C \frac{dV_C}{dt}. \end{aligned} \tag{10.4}$$

From lines one and three of (10.4), we obtain

$$V_R = RC \frac{dV_C}{dt}, \tag{10.5}$$

and from lines two and three of (10.4), we obtain

$$\begin{aligned} V_L &= L \frac{d}{dt} \left(C \frac{dV_C}{dt} \right) \\ &= LC \frac{d^2V_C}{dt^2}. \end{aligned} \tag{10.6}$$

Substituting these into (10.3) yields

$$\begin{aligned} -V_S + RC \frac{dV_C}{dt} + LC \frac{d^2V_C}{dt^2} + V_C &= 0 \\ \Downarrow (\text{rearranging terms}) \\ LC \frac{d^2V_C}{dt^2} + RC \frac{dV_C}{dt} + V_C &= V_S. \end{aligned} \tag{10.7}$$

Upon defining $y := V_C$ and $u := V_S$, we have

$$\begin{aligned} LC \frac{d^2y}{dt^2} + RC \frac{dy}{dt} + y &= u \\ \Downarrow (\text{dividing through by } LC) \\ \frac{d^2y}{dt^2} + \frac{R}{L} \frac{dy}{dt} + \frac{1}{LC} y &= \frac{1}{LC} u. \end{aligned} \tag{10.8}$$

Note: Resistors heat up as the current increases, and this can change the “local resistance” of the element. Modeling this effect could lead to either a nonlinear or LTV model. ■

Example 10.4. Re-do the above example, but this time, take i , the loop current, as the output. The input remains the applied voltage, V_S .

Solution: Ans. $\frac{d^2y}{dt^2} + \frac{R}{L} \frac{dy}{dt} + \frac{1}{LC} y = \frac{1}{L} \frac{du}{dt}$, where we note that the input is differentiated.

By differentiating (10.3), we obtain

$$-\dot{V}_S + \dot{V}_R + \dot{V}_L + \dot{V}_C = 0. \tag{10.9}$$

This opening move makes the remaining computations easier. How would you know to do it? Experience, like that obtained in Michigan’s EECS 215 Introduction to Electronic Circuits. Remember, this Calculus course is not attempting to teach you electronics, mechanics, or physics; it is, however, illustrating how Calculus shows up in these subjects. To acquire domain knowledge, you need to carefully construct your course plan.

After the brilliant opening move in (10.9), we next apply the constitutive relationships in (10.4). By differentiating its first and second lines, we obtain

$$\begin{aligned} \frac{dV_R}{dt} &= R \frac{di}{dt} \\ \frac{dV_L}{dt} &= L \frac{d^2i}{dt^2}, \end{aligned} \tag{10.10}$$

and from the third line, we obtain directly,

$$\frac{dV_C}{dt} = \frac{1}{C} i. \tag{10.11}$$

Substituting these into (10.9) yields

$$-\frac{dV_S}{dt} + R \frac{di}{dt} + L \frac{d^2i}{dt^2} + \frac{1}{C} i = 0. \quad (10.12)$$

Dividing through by $\frac{1}{L}$ and rearranging terms gives

$$\frac{d^2i}{dt^2} + \frac{R}{L} \frac{di}{dt} + \frac{1}{LC} i = \frac{1}{L} \frac{dV_S}{dt}. \quad (10.13)$$

Setting $y := i$ and $u := V_s$ then gives the answer. ■



Figure 10.2: Segways take many forms. The one we are modeling is closer to image (a), which shows a classic Segway transporter. Do you think the designer of the balance algorithm for the devices knew ahead of time the exact mass, moments of inertia, etc, for the riders? Or did they have to make the balance algorithm robust to a wide range of human and mechanical riders? To accomplish the latter takes more than our brief introduction to feedback control. You can read about the [algorithm used by Cassie](#) to achieve the very meta stunt of a robot riding another robot.

10.2.3 A State-Variable Model of a (Planar) Segway Transporter

A Segway Transporter is a Mobile Robot Ridden by a Human or another Robot. In the following, we detail some of its key components.

- **Mechanical Components:**

- *Wheels and Tires:* Two large wheels with durable tires designed for traction.
- *Frame and Body:* Sturdy frame supporting the rider's platform, handlebar, and other components.
- *Gear System:* Gears connect the electric motors to the wheels, enabling precise speed and movement control.

- **Electrical Components:**

- *Batteries:* Rechargeable Lithium-ion or Nickel-Metal Hydride batteries provide power to the motors and control systems.
- *Electric Motors:* Two independent brushless DC electric motors, one for each wheel, providing propulsion and control.
- *Control Circuitry:* Electronic control units for managing power distribution, motor control, and balancing algorithms.

- **Sensing Components:**

- *Gyroscopic Sensors:* Multiple gyroscopes provide data on the tilt and orientation of the transporter for balance control.

- **Accelerometers:** Sensors measuring acceleration forces to maintain balance and control motion.
- **Encoders:** Encoders in the wheels measure rotation, which, when differentiated, aid in speed control and stability.

A planar set of **robot equations** (aka, Lagrange's Equations of Motion) for the Segway Transporter is derived in Chapter 10.11. The model lumps the rider and the Segway platform into a single object with a combined mass, center of mass, and moment of inertia about its center of mass. The lean angle of the combined platform and rider is denoted by θ and is measured with respect to the vertical, with zero being upright. Each wheel of the Segway has a given mass, radius, and moment of inertia. Furthermore, each wheel is connected to a DC motor through gearing. The mass of the motors is lumped into the mass of the platform because they are physically bolted to it. The model includes the moment of inertia of each motor's rotor and the gear ratio. The control variable u in the model is the torque generated by the motor and applied to its rotor, typically called "motor torque". The torque applied between the platform and each wheel is multiplied by the gear ratio. The model does not include the electrical dynamics of the DC motors. Hence, one could build an even more realistic model that used the voltage applied to the motor as the control input.

How does a Segway accelerate to a desired speed or brake to a stop? If you say it uses the motors, you are only half right. Why only half? Because if the lean angle is zero, that is, the platform+rider's center of mass is directly over the wheels, then the linear acceleration is zero. To accelerate forward, the rider must lean forward. Then, to prevent a face plant, the motor plus control system must accelerate the wheels forward so as to drive the wheels under the center of mass of the platform+rider. Once the center of mass returns to being over the wheels, the linear acceleration becomes zero, and the Segway rolls along at a constant speed. To come to a stop, the rider leans opposite to the direction of travel. In a real 3D Segway, steering is accomplished by leaning the platform's handle to the left or right. In our planar (2D) Segway, there is no steering. In some sense, we are on train tracks.

Fact 10.5: Linearized (State-Variable) Model of the Segway Transporter

The linearization of the Segway Transporter about the upright equilibrium yields the state-variable control model

$$\frac{d}{dt} \begin{bmatrix} \theta \\ \varphi \\ \dot{\theta} \\ \dot{\varphi} \end{bmatrix} = \begin{bmatrix} 0.000 & 0.000 & 1.000 & 0.000 \\ 0.000 & 0.000 & 0.000 & 1.000 \\ 50.937 & 0.000 & 0.000 & 0.000 \\ -162.687 & 0.000 & 0.000 & 0.000 \end{bmatrix} \begin{bmatrix} \theta \\ \varphi \\ \dot{\theta} \\ \dot{\varphi} \end{bmatrix} + \begin{bmatrix} 0.000 \\ 0.000 \\ -21.468 \\ 80.227 \end{bmatrix} u, \quad (10.14)$$

where u is the motor torque. The coefficients in the matrix are constants; thus, the system is LTI. The model can also be written as

$$\underbrace{\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \\ \dot{x}_3 \\ \dot{x}_4 \end{bmatrix}}_{\dot{x}} = \underbrace{\begin{bmatrix} 0.000 & 0.000 & 1.000 & 0.000 \\ 0.000 & 0.000 & 0.000 & 1.000 \\ 50.937 & 0.000 & 0.000 & 0.000 \\ -162.687 & 0.000 & 0.000 & 0.000 \end{bmatrix}}_A \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}}_x + \underbrace{\begin{bmatrix} 0.000 \\ 0.000 \\ -21.468 \\ 80.227 \end{bmatrix}}_b u, \quad (10.15)$$

where $x_1 := \theta$, $x_2 := \varphi$, $x_3 := \dot{\theta}$, and $x_4 := \dot{\varphi}$.

We next define interesting outputs for the model.

- **Lean Angle of the Segway Platform plus Rider:** $y_\theta = c_\theta \cdot x = \underbrace{\begin{bmatrix} 1 & 0 & 0 & 0 \end{bmatrix}}_{c_\theta} x$. We need this quantity to be near zero for the Segway to be upright.
- **Three different means to define speed:**

- **Horizontal Speed Level 1:** $y_{v1} = c_{v1} \cdot x = \underbrace{\begin{bmatrix} 0.000 & 0.000 & 0.000 & 0.250 \end{bmatrix}}_{c_{v1}} x$ multiplies the radius of the wheel (0.25 m) by the relative angular velocity of the wheel with respect to the angle of the platform (in rad/s) to obtain a speed in m/s. The relative angular velocity is different from the wheel's angular velocity in the world frame when the platform also has angular velocity.
- **Horizontal Speed Level 2:** $y_{v2} = c_{v2} \cdot x = \underbrace{\begin{bmatrix} 0.000 & 0.000 & 0.250 & 0.250 \end{bmatrix}}_{c_{v2}} x$, uses wheel radius and the absolute angular velocity of the wheel. The latter is the wheel's angular velocity in the world frame, even when the platform also has angular velocity.
- **Horizontal Speed Level 3:** $y_{v3} = c_{v3} \cdot x = \underbrace{\begin{bmatrix} 0.000 & 0.000 & 0.863 & 0.250 \end{bmatrix}}_{c_{v3}} x$, is the horizontal velocity of the center of mass of the overall Segway. This seems like the correct quantity to control—**except**—the center of mass of the real system will vary with the mass and geometry of the rider. So many decisions. Engineering is hard!

For each of these outputs, $\dot{x} = Ax + bu$, $y = cx$ creates a SISO state-variable model for the Segway.

Note: When the Segway is upright and not wobbling wildly, $\dot{\theta}$, the angular velocity of the platform plus rider, will be very small compared to the angular velocity of the wheels. Hence, one may posit that the three horizontal speeds above are the same for all intents and purposes.

10.3 The Laplace Transform and an Algebraic Look at Linear ODEs

The Laplace Transform is a powerful mathematical tool that plays a crucial role in the analysis and solution of linear time-invariant (LTI) ordinary differential equations (ODEs). Its ability to transform equations involving derivatives with respect to time into algebraic equations in the Laplace variable, s (AKA, the complex frequency domain), simplifies the operations of differentiation and integration in the time domain to mere multiplication and division in the frequency domain. Furthermore, the solution can be brought back from the frequency domain to the original time domain by taking the inverse Laplace Transform. This technique is especially useful in engineering disciplines such as control systems, circuits, signal processing, and communications, where LTI systems are common.

The Laplace Transform

A function of time $f : \mathbb{R} \rightarrow \mathbb{R}$ is said to be **right-sided** if $f(t) = 0$ for all $t < 0$.

Definition 10.6. For a right-sided piecewise continuous function, $f(t)$, its **Laplace transform**^a is the improper integral,

$$\begin{aligned} F(s) := \mathcal{L}\{f(t)\} &:= \int_{0^-}^{\infty} e^{-st} f(t) dt \\ &:= \lim_{a \rightarrow 0^-, b \rightarrow \infty} \int_a^b e^{-st} f(t) dt, \end{aligned} \tag{10.16}$$

where $s = \sigma + i\omega \in \mathbb{C}$ is called the **Laplace transform variable**.

Note: We use the notation “ $f(t) \leftrightarrow F(s)$ ” to denote that $F(s)$ is the Laplace transform of $f(t)$ and $f(t)$ is the **inverse Laplace transform** of $F(s)$. The two functions $f(t)$ and $F(s)$ are also called **Laplace transform pairs**.

We'll see soon that the lower limit, 0^- (zero minus), is there to handle initial conditions for ODEs. Because the upper limit on the Laplace transform is infinite, the integral does not always exist. We take care of that next.

Definition 10.7. The **region of convergence (ROC)** of $F(s)$ is

$$\begin{aligned} ROC &:= \{s \in \mathbb{C} \mid e^{-st} f(t) \text{ is absolutely integrable}\} \\ &= \{s \in \mathbb{C} \mid \int_{0^-}^{\infty} |e^{-st} f(t)| dt < \infty\}. \end{aligned} \quad (10.17)$$

Note: Because $s = \sigma + i\omega$, and because $|e^{-(\sigma+i\omega)t}| = |e^{-\sigma t} \cdot e^{-i\omega t}| = e^{-\sigma t}$, the ROC can also be written as

$$ROC := \{s \in \mathbb{C} \mid \int_{0^-}^{\infty} e^{-\sigma t} |f(t)| dt < \infty\}. \quad (10.18)$$

^aTechnically, this is the **unilateral Laplace transform**. There exists also a bilateral Laplace transform; because it is not used much in engineering, we do not introduce it.

Remark 10.8. (Optional Read:) The Laplace transform is a doubly improper integral due to the two limits. In theory, this is a bit painful. In practice, for all the functions that are important to engineers, one uses Table 10.11 (or a larger version of it) to compute the Laplace transform and its inverse. As stated above, the use of 0^- (zero minus) as the lower limit is mainly due to the way that the Laplace transform interacts with the initial conditions of a linear ODE. In fact, the initial conditions are the limits from the left of $y(t)$ and its derivatives at $t = 0$. A second reason is the **unit impulse “function”**, the first function listed in Table 10.11. A course such as Michigan’s EECS 216 Signals and Systems delves more deeply into all of this. Our goal is to provide more of a “Quick-Reference Guide” so that you can rapidly see the power and insight unleashed by the Laplace transform before dedicating an entire course to its mastery.

Proposition 10.9: Two Key Properties of the Laplace Transform

- (a) For all functions $f(t)$ of the form $t^n e^{at} u_{\text{stp}}(t)$, $t^n e^{at} \sin(\omega t) u_{\text{stp}}(t)$, and $t^n e^{at} \cos(\omega t) u_{\text{stp}}(t)$, the region of convergence of $F(s) = \mathcal{L}\{f(t)\}$ is

$$ROC = \{s \in \mathbb{C} \mid \text{real}(s - a) > 0\} = \{\sigma + i\omega \in \mathbb{C} \mid \sigma > a\}.$$

- (b) Suppose that $f(t)$ is continuously differentiable and both $f(t)$ and $f'(t)$ have Laplace transforms. Then,

$$\mathcal{L}\{f'(t)\} = sF(s) - f(0^-),$$

for $F(s) = \mathcal{L}\{f(t)\}$ and s in the intersection of their respective regions of convergence. In other words, differentiation in the Laplace domain corresponds to multiplication by the Laplace variable s .

The Laplace transforms of high-order derivatives are given in Table 10.10.

Note: For a SISO linear system, 0^- comes from the initial conditions on the output signal, $y(t)$. For the input signal, $u(t)$, its value at 0^- is ALWAYS zero. We have consistently multiplied input signals by the unit step function, $u_{\text{stp}}(t)$, so that $\lim_{t \rightarrow 0^-} u^{(k)}(t) = 0$, for $k = 0, 1, 2, \dots$

(Optional Read:) Proof: (a) is immediate from Example 8.13.

To establish (b), we use our old friend, integration by parts: $\int u \, dv = u \cdot v - \int v \, du$.

$$\begin{aligned}
\mathcal{L}\{f'(t)\} &= \int_{0^-}^{\infty} f'(t)e^{-st} dt \\
&= \int_{0^-}^{\infty} \underbrace{e^{-st}}_u \cdot \underbrace{f'(t) dt}_v \\
&= \underbrace{e^{-st}}_u \cdot \underbrace{f(t)}_v \Big|_{t=0^-}^{t=\infty} - \int_{0^-}^{\infty} \underbrace{f(t)}_v \cdot \underbrace{(-s) \cdot e^{-st} dt}_{du} \\
&= \underbrace{e^{-st} \cdot f(t)}_{\text{must be careful here}} \Big|_{t=0^-}^{t=\infty} + s \underbrace{\int_{0^-}^{\infty} f(t) \cdot e^{-st} dt}_{s \cdot F(s)}.
\end{aligned}$$

The second term of the last row is $s \cdot F(s)$, the Laplace variable s times the Laplace transform of $f(t)$. As indicated, the first term takes a bit of care, but we covered the key technical part in Prop. 8.12-(b). For $s \in ROC$ of both Laplace transforms,

$$\lim_{t \rightarrow \infty} e^{-st} \cdot f(t) = 0.$$

Therefore, $e^{-st} \cdot f(t) \Big|_{t=0^-}^{t=\infty} = 0 - f(0^-)$, yielding

$$\mathcal{L}\{f'(t)\} = sF(s) - f(0^-).$$

Let's do at least one of the higher derivatives. If we use the notation, $f''(t) = (f'(t))'$, then we have that

$$\mathcal{L}\{f''(t)\} = \mathcal{L}\{(f'(t))'\} = s\mathcal{L}\{f'(t)\} - f'(0^-) = s(sF(s) - f(0^-)) - f'(0^-) = s^2F(s) - sf(0^-) - f'(0^-),$$

assuming all indicated Laplace transforms exist. From here, a proof by induction gives the remainder of Table 10.10. ■

Table 10.10: Laplace Transforms of Derivatives with Non-zero Initial Conditions, assuming all indicated Laplace transforms exist.

Time Domain	Laplace Domain
$f'(t)$	$sF(s) - f(0^-)$
$f''(t)$	$s^2F(s) - sf(0^-) - f'(0^-)$
$f'''(t)$	$s^3F(s) - s^2f(0^-) - sf'(0^-) - f''(0^-)$
\vdots	\vdots
$f^{(n)}(t)$	$s^nF(s) - s^{n-1}f(0^-) - s^{n-2}f'(0^-) - \dots - f^{(n-1)}(0^-)$
$f^{(n)}(t)$	$s^nF(s)$ when all initial conditions are zero.

Very Important: An input signal $u(t)$ never has (non-zero) initial conditions. Only output signals $y(t)$ can have non-zero initial conditions. In other symbols, for all input signals, $u(0^-) := \lim_{t \rightarrow 0^-} u(t) = 0$ and for all $k = 1, 2, \dots$, $u^{(k)}(0^-) :=$

$$\lim_{t \rightarrow 0^-} u^{(k)}(t) = 0.$$

10.3.1 Common Laplace Transform Pairs

The traditional path through Laplace transforms is to spend an hour or so computing many Laplace transform pairs, even though that is not how engineers use the Laplace transform in practice. We recommend you look at Table 10.11, skim Chapter 10.3.2, and then move on to Chapter 10.3.4.

Table 10.11: Common Laplace Transform Pairs. The unit step function $u_{\text{stp}}(t)$ is included to emphasize that the Laplace transforms are only valid for signals that start at or after time $t = 0$; in other words, the signals must be zero for $t < 0$. The unit impulse, $\delta(t)$, is presented as an optional read in Chapter 10.10. Its ubiquity in other courses makes it hard to ignore completely.

Function $f(t), t \geq 0$	Laplace Transform $F(s)$	ROC (can be ignored)
$\delta(t)$	1	all s
$u_{\text{stp}}(t)$	$\frac{1}{s}$	$\text{real}(s) > 0$
$t^n u_{\text{stp}}(t)$	$\frac{n!}{s^{n+1}}$	$\text{real}(s) > 0$
$e^{at} u_{\text{stp}}(t)$	$\frac{1}{s-a}$	$\text{real}(s-a) > 0$
$t^n e^{at} u_{\text{stp}}(t)$	$\frac{n!}{(s-a)^{n+1}}$	$\text{real}(s-a) > 0$
$\sin(\omega t) u_{\text{stp}}(t)$	$\frac{\omega}{s^2 + \omega^2}$	$\text{real}(s) > 0$
$\cos(\omega t) u_{\text{stp}}(t)$	$\frac{s}{s^2 + \omega^2}$	$\text{real}(s) > 0$
$e^{at} \sin(\omega t) u_{\text{stp}}(t)$	$\frac{\omega}{(s-a)^2 + \omega^2}$	$\text{real}(s-a) > 0$
$e^{at} \cos(\omega t) u_{\text{stp}}(t)$	$\frac{s-a}{(s-a)^2 + \omega^2}$	$\text{real}(s-a) > 0$
$t \sin(\omega t) u_{\text{stp}}(t)$	$\frac{2\omega s}{(s^2 + \omega^2)^2}$	$\text{real}(s) > 0$
$t \cos(\omega t) u_{\text{stp}}(t)$	$\frac{s^2 - \omega^2}{(s^2 + \omega^2)^2}$	$\text{real}(s) > 0$
$f(t-b) \cdot u_{\text{stp}}(t-b), b \geq 0$	$e^{-bs} F(s)$	ROC of $F(s)$

10.3.2 (Optional Read:) How to Compute Laplace Transform Pairs by Hand

This section is an optional read because, in practice, one NEVER computes Laplace transforms by hand, and in fact, even though they can be computed symbolically, that is rarely done either! Instead, one simply uses a table.

Example 10.12. Determine for which values of the complex variable $s = \sigma + i\omega$ the following limits exist and are zero. Use the notation $\text{real}(s)$ to denote σ , the real part of the complex number, s , and $\text{imag}(s)$ to denote ω , the imaginary part of the complex number, s .

- (a) $\lim_{t \rightarrow \infty} e^{st}$
- (b) $\lim_{t \rightarrow \infty} e^{-st}$
- (c) $\lim_{t \rightarrow \infty} e^{-(s+a)t}$
- (d) $\lim_{t \rightarrow \infty} t e^{-st}$

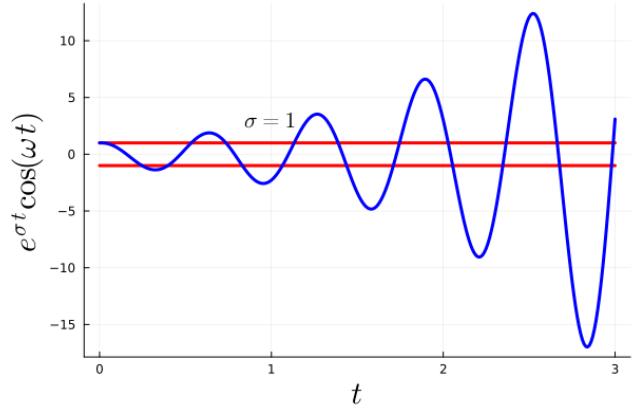
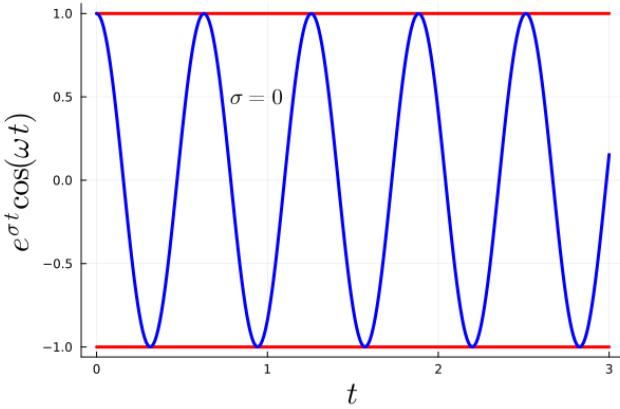
Solutions: For use in the following solutions, we note that

$$e^{st} = e^{(\sigma+i\omega)t} = e^{\sigma t} \cdot e^{i\omega t} = e^{\sigma t} \cdot \cos(\omega t) + i e^{\sigma t} \cdot \sin(\omega t).$$

Moreover, $\lim_{t \rightarrow \infty} e^{st}$ exists if, and only if, the two limits $\lim_{t \rightarrow \infty} e^{\sigma t} \cdot \cos(\omega t)$ and $\lim_{t \rightarrow \infty} e^{\sigma t} \cdot \sin(\omega t)$ both exist.

$$(a) \text{ Ans. } \lim_{t \rightarrow \infty} e^{st} = \lim_{t \rightarrow \infty} e^{\sigma t} \cdot e^{i\omega t} = \begin{cases} 0 & \text{real}(s) < 0 \\ 1 & s = 0 \\ \text{does not exist} & \text{real}(s) \geq 0, \text{imag}(s) \neq 0 \\ \infty & \text{real}(s) > 0, \text{imag}(s) = 0. \end{cases}$$

The limit does not exist when $\text{real}(s) \geq 0, \text{imag}(s) \neq 0$ because $\lim_{t \rightarrow \infty} e^{\sigma t} \cdot \cos(\omega t)$ “oscillates” as $t \rightarrow \infty$; more precisely, for all $T > 0$, there exist $t_1, t_2 \geq T$ such that $e^{\sigma t_1} \cdot \cos(\omega t_1) = +1$ and $e^{\sigma t_2} \cdot \cos(\omega t_2) = -1$ as illustrated in the figures below.



$$(b) \text{ Ans. } \lim_{t \rightarrow \infty} e^{-st} = \lim_{t \rightarrow \infty} e^{-\sigma t} \cdot e^{-i\omega t} = \begin{cases} 0 & \text{real}(s) > 0 \\ 1 & s = 0 \\ \text{does not exist} & \text{real}(s) \leq 0, \text{imag}(s) \neq 0 \\ \infty & \text{real}(s) < 0, \text{imag}(s) = 0. \end{cases}$$

$$(c) \text{ Ans. } \lim_{t \rightarrow \infty} e^{-(\sigma+a)t} \cdot e^{-i\omega t} = \begin{cases} 0 & \text{real}(s) + a > 0 \\ 1 & s + a = 0 \\ \text{does not exist} & \text{real}(s) + a \leq 0, \text{imag}(s) \neq 0 \\ \infty & \text{real}(s) + a < 0, \text{imag}(s) = 0 \end{cases} = \begin{cases} 0 & \text{real}(s) > -a \\ 1 & s = -a \\ \text{does not exist} & \text{real}(s) \leq -a, \text{imag}(s) \neq 0 \\ \infty & \text{real}(s) < -a, \text{imag}(s) = 0 \end{cases}$$

$$(d) \text{ Ans. } \lim_{t \rightarrow \infty} t e^{-st} = \lim_{t \rightarrow \infty} t e^{-\sigma t} \cdot e^{-i\omega t} = \begin{cases} 0 & \text{real}(s) > 0 \\ \infty & s = 0 \\ \text{does not exist} & \text{real}(s) \leq 0, \text{imag}(s) \neq 0 \\ \infty & \text{real}(s) \leq 0, \text{imag}(s) = 0. \end{cases}$$

Here, L'Hôpital's Rule is useful. Suppose that $\sigma < 0$. Then, $t e^{\sigma t}$ presents us with infinity times zero in the limit, which we can rewrite as infinity over infinity, $\frac{t}{e^{-\sigma t}} = \frac{t}{e^{|\sigma|t}}$, because $\sigma < 0$. Hence, applying L'Hôpital's Rule gives

$$\lim_{t \rightarrow \infty} t e^{\sigma t} = \lim_{t \rightarrow \infty} \frac{t}{e^{|\sigma|t}} = \lim_{t \rightarrow \infty} \frac{1}{|\sigma| e^{|\sigma|t}} = 0.$$

■

Example 10.13. Compute the Laplace transforms of the following signals, taking into account the limits computed in Example 10.12.

$$(a) u_{\text{stp}}(t)$$

$$(b) e^{-at}u_{\text{stp}}(t)$$

$$(c) tu_{\text{stp}}(t).$$

Solutions:

$$(a) \text{ Ans. } u_{\text{stp}}(t) \longleftrightarrow \frac{1}{s}$$

Using antiderivatives,

$$\begin{aligned}\mathcal{L}\{u_{\text{stp}}(t)\} &= \int_{0^-}^{\infty} e^{-st}u_{\text{stp}}(t) dt \\ &= \int_0^{\infty} e^{-st} dt \\ &= \frac{-1}{s} e^{-st} \Big|_{t=0}^{t=\infty} \\ &= \frac{-1}{s} (0 - 1) \quad \text{real}(s) > 0 \\ &= \frac{1}{s} \quad \text{real}(s) > 0.\end{aligned}$$

The above calculation is saying that $\mathcal{L}\{u_{\text{stp}}(t)\}$ exists and is equal to $\frac{1}{s}$ for all $\{s \in \mathbb{C} \mid \text{real}(s) > 0\}$. The set of all values in the complex plane for which the Laplace transform of a function exists is called its **region of convergence**. This set is very important for the bilateral transform, while it can be ignored completely for the unilateral Laplace transform. Why is that? It's because of Prop. 10.9-(a), which essentially says that the ROC can be computed by inspection for functions that commonly occur in engineering and STEM.

$$(b) e^{-at}u_{\text{stp}}(t) \quad \text{Ans.} \quad \frac{1}{s+a}$$

Using antiderivatives,

$$\begin{aligned}\mathcal{L}\{u_{\text{stp}}(t)\} &= \int_{0^-}^{\infty} e^{-st}e^{-at}u_{\text{stp}}(t) dt \\ &= \int_0^{\infty} e^{-(s+a)t} dt \\ &= \frac{-1}{s+a} e^{-(s+a)t} \Big|_{t=0}^{t=\infty} \\ &= \frac{-1}{s+a} (0 - 1) \quad \text{real}(s) > -a \\ &= \frac{1}{s+a} \quad \text{real}(s) > -a.\end{aligned}$$

$$(c) tu_{\text{stp}}(t) \quad \text{Ans.} \quad \frac{1}{s^2}$$

Using integration by parts, with $u = t$ and $dv = e^{-st} dt \implies v = -\frac{1}{s}e^{-st}$,

$$\begin{aligned}
\mathcal{L}\{u_{\text{stp}}(t)\} &= \int_{0^-}^{\infty} e^{-st} t u_{\text{stp}}(t) dt \\
&= \underbrace{\int_0^{\infty} t e^{-st} dt}_{\int u dv} \\
&= \underbrace{-\frac{t}{s} e^{-st} \Big|_{t=0}^{t=\infty}}_{u \cdot v} - \underbrace{\int_0^{\infty} \frac{-1}{s} e^{-st} dt}_{\int v du} \\
&= -\frac{t}{s} e^{-st} \Big|_{t=0}^{t=\infty} + \frac{1}{s} \cdot \int_0^{\infty} e^{-st} dt \\
&= (0 - 0) + \frac{-1}{s^2} e^{-st} \Big|_{t=0}^{t=\infty} \\
&= \frac{-1}{s^2} (0 - 1) \quad \text{real}(s) > 0 \\
&= \frac{1}{s^2} \quad \text{real}(s) > 0.
\end{aligned}$$

■

10.3.3 Computing Inverse Laplace Transforms by Hand

While one uses tables to compute Laplace transforms, the task of computing inverse Laplace transforms requires more work. For example, the Laplace transform

$$F(s) = \frac{s+1}{(s+2)(s+3)} \tag{10.19}$$

is nowhere to be found in Table 10.11. However, it is easy to verify that

$$\frac{s+1}{(s+2)(s+3)} = \frac{2}{s+3} - \frac{1}{s+2}, \tag{10.20}$$

and hence, by Table 10.11, the inverse Laplace transform is

$$\begin{aligned}
f(t) &= \mathcal{L}^{-1}\left\{\frac{2}{s+3}\right\} - \mathcal{L}^{-1}\left\{\frac{1}{s+2}\right\} \\
&= 2e^{-3t}u_{\text{stp}}(t) - e^{-2t}u_{\text{stp}}(t).
\end{aligned} \tag{10.21}$$

The process of going from (10.19) to (10.20) is called a partial fraction expansion, or PFE for short. We first encountered PFEs in Chapter 7.3.5.

Proposition 10.14: Partial Fraction Expansion

A proper rational function $F(s)$ with distinct roots r_1, r_2, \dots, r_n can be expanded as

$$F(s) = \sum_{j=1}^n \frac{k_j}{s - r_j},$$

where the coefficients k_j are called the residues of $F(s)$ at the poles r_j . The residues k_j can be computed using:

$$k_j = \lim_{s \rightarrow r_j} F(s)(s - r_j).$$

Remark: If r_j is complex, the associated residue k_j will generally be a complex number. Its complex conjugate k_j^* corresponds to the residue at the conjugate pole r_j^* .

Example 10.15. Compute the inverse Laplace transforms of the following rational functions, using the method of partial fraction expansion, where necessary.

$$(a) \frac{1}{s^2+2s+2}$$

$$(b) \frac{s+3}{(s+1)(s-2)}$$

$$(c) \frac{s^2+4}{(s^2+1)(s+3)}$$

Solutions:

$$(a) \text{ Ans. } \frac{1}{s^2+2s+2} \longleftrightarrow e^{-t} \sin(t) u_{\text{stp}}(t)$$

Rewrite the denominator as a completed square:

$$s^2 + 2s + 2 = (s + 1)^2 + 1.$$

The Laplace transform becomes:

$$\frac{1}{s^2 + 2s + 2} = \frac{1}{(s + 1)^2 + 1}.$$

Using the table of Laplace transforms:

$$\mathcal{L}^{-1} \left\{ \frac{1}{(s + 1)^2 + 1} \right\} = e^{-t} \sin(t) u_{\text{stp}}(t).$$

$$(b) \text{ Ans. } \frac{s+3}{(s+1)(s-2)} \longleftrightarrow \left(-\frac{2}{3}e^{-t} + \frac{5}{3}e^{2t} \right) u_{\text{stp}}(t)$$

Using partial fraction expansions,

$$\frac{s+3}{(s+1)(s-2)} = \frac{k_1}{s+1} + \frac{k_2}{s-2}.$$

Solving for k_1 and k_2 ,

$$k_1 = \lim_{s \rightarrow -1} \frac{(s+3)(s+1)}{(s+1)(s-2)} = \lim_{s \rightarrow -1} \frac{s+3}{s-2} = -\frac{2}{3}, \quad k_2 = \lim_{s \rightarrow 2} \frac{(s+3)(s-2)}{(s+1)(s-2)} = \lim_{s \rightarrow 2} \frac{s+3}{s+1} = \frac{5}{3}.$$

Substituting back,

$$\frac{s+3}{(s+1)(s-2)} = \frac{-\frac{2}{3}}{s+1} + \frac{\frac{5}{3}}{s-2}.$$

Using the inverse Laplace transform:

$$\mathcal{L}^{-1} \left\{ \frac{-\frac{2}{3}}{s+1} \right\} = -\frac{2}{3}e^{-t} u_{\text{stp}}(t), \quad \mathcal{L}^{-1} \left\{ \frac{\frac{5}{3}}{s-2} \right\} = \frac{5}{3}e^{2t} u_{\text{stp}}(t).$$

Therefore:

$$f(t) = \left(-\frac{2}{3}e^{-t} + \frac{5}{3}e^{2t} \right) u_{\text{stp}}(t).$$

$$(c) \text{ Ans. } \frac{s^2+4}{(s^2+1)(s+3)} \longleftrightarrow \left(-\frac{3}{10} \cos(t) - \frac{9}{10} \sin(t) + \frac{13}{10} e^{-3t} \right) u_{\text{stp}}(t)$$

This problem is the most complicated, so we show all of the intermediate steps. **By the time we are done, you'll be ready for software to compute inverse Laplace transforms.** We express

$$\frac{s^2+4}{(s^2+1)(s+3)} = \frac{k_1}{s-i} + \frac{k_2}{s+i} + \frac{k_3}{s+3}$$

as a partial fraction expansion.

Step 1: Solve for k_1 , k_2 , and k_3

$$\begin{aligned} k_1 &= \lim_{s \rightarrow i} \frac{(s^2+4)(s-i)}{(s^2+1)(s+3)} \\ &= \lim_{s \rightarrow i} \frac{(s^2+4)(s-i)}{(s+i)(s-i)(s+3)} \\ &= \lim_{s \rightarrow i} \frac{(s^2+4)}{(s+i)(s+3)} \\ &= \frac{i^2+4}{(i+i)(i+3)} \\ &= \frac{-1+4}{2i(i+3)} \\ &= \frac{3}{2i(i+3)} \\ &= \frac{3}{-2+6i} \\ &= \frac{3(-2-6i)}{(-2+6i)(-2-6i)} \\ &= \frac{-6-18i}{40} \\ &= -\frac{3}{20} - \frac{9i}{20}. \end{aligned}$$

$$\begin{aligned}
k_2 &= \lim_{s \rightarrow -i} \frac{(s^2 + 4)(s + i)}{(s^2 + 1)(s + 3)} \\
&= \lim_{s \rightarrow -i} \frac{(s^2 + 4)(s + i)}{(s + i)(s - i)(s + 3)} \\
&= \lim_{s \rightarrow -i} \frac{(s^2 + 4)}{(s - i)(s + 3)} \\
&= \frac{(-i)^2 + 4}{(-i - i)(-i + 3)} \\
&= \frac{-1 + 4}{-2i(-i + 3)} \\
&= \frac{3}{-2i(-i + 3)} \\
&= \frac{3}{-2 - 6i} \\
&= \frac{3(-2 + 6i)}{(-2 - 6i)(-2 + 6i)} \\
&= \frac{-6 + 18i}{40} \\
&= -\frac{3}{20} + \frac{9i}{20}.
\end{aligned}$$

$$\begin{aligned}
k_3 &= \lim_{s \rightarrow -3} \frac{(s^2 + 4)(s + 3)}{(s^2 + 1)(s + 3)} \\
&= \lim_{s \rightarrow -3} \frac{s^2 + 4}{s^2 + 1} \\
&= \frac{((-3)^2 + 4)}{((-3)^2 + 1)} \\
&= \frac{(9 + 4)}{(9 + 1)} \\
&= \frac{13}{10}.
\end{aligned}$$

Step 2: Substitute Back The partial fraction expansion becomes,

$$\frac{s^2 + 4}{(s^2 + 1)(s + 3)} = \frac{-\frac{3}{20} - \frac{9i}{20}}{s - i} + \frac{-\frac{3}{20} + \frac{9i}{20}}{s + i} + \frac{\frac{13}{10}}{s + 3}.$$

Step 3: Compute the Inverse Laplace Transform Using the standard inverse Laplace transforms,

$$\mathcal{L}^{-1} \left\{ \left(-\frac{3}{20} - \frac{9i}{20} \right) \frac{1}{s-i} \right\} = \left(-\frac{3}{20} - \frac{9i}{20} \right) e^{it} u_{\text{stp}}(t),$$

$$\mathcal{L}^{-1} \left\{ \left(-\frac{3}{20} + \frac{9i}{20} \right) \frac{1}{s+i} \right\} = \left(-\frac{3}{20} + \frac{9i}{20} \right) e^{-it} u_{\text{stp}}(t),$$

$$\mathcal{L}^{-1} \left\{ \frac{13}{10} \frac{1}{s+3} \right\} = \frac{13}{10} e^{-3t} u_{\text{stp}}(t).$$

In the following, we drop the unit step function to save space. We add it back at the end of the computations.

Step 4: Combine Complex Exponentials Combine the terms with complex exponentials using Euler's formula:

$$e^{it} = \cos(t) + i \sin(t), \quad e^{-it} = \cos(t) - i \sin(t).$$

Consider,

$$\left(-\frac{3}{20} - \frac{9i}{20} \right) e^{it} + \left(-\frac{3}{20} + \frac{9i}{20} \right) e^{-it}.$$

Using Euler's formula gives,

$$\left(-\frac{3}{20} - \frac{9i}{20} \right) (\cos(t) + i \sin(t)) + \left(-\frac{3}{20} + \frac{9i}{20} \right) (\cos(t) - i \sin(t)).$$

Expanding terms,

$$\left(-\frac{3}{20} \cos(t) - \frac{9i}{20} \cos(t) - \frac{3i}{20} \sin(t) - \frac{9}{20} \sin(t) \right) + \left(-\frac{3}{20} \cos(t) + \frac{9i}{20} \cos(t) + \frac{3i}{20} \sin(t) - \frac{9}{20} \sin(t) \right).$$

Simplifying,

$$-\frac{6}{20} \cos(t) - \frac{18}{20} \sin(t),$$

and thus

$$-\frac{3}{10} \cos(t) - \frac{9}{10} \sin(t).$$

Step 5: Final Result The time-domain function is

$$f(t) = -\frac{3}{10} \cos(t) - \frac{9}{10} \sin(t) + \frac{13}{10} e^{-3t}.$$

Thus, the final result is:

$$f(t) = \left(-\frac{3}{10} \cos(t) - \frac{9}{10} \sin(t) + \frac{13}{10} e^{-3t} \right) u_{\text{stp}}(t).$$

■

10.3.4 Algebraic Approach to Solving Linear ODEs

Our principal goal is to transform a linear time-invariant ODE into an algebraic equation and then underline the utility of this process. While in engineering practice, one rarely uses the Laplace transform to compute solutions to ODEs, doing so is the most direct way to see derivatives disappear in front of your eyes. In that light, we give two examples that reduce the act of solving an ODE to dealing with rational functions of the Laplace variable, s .

Example 10.16. Use the Laplace transform to compute the solution to the linear time-invariant ODE

$$\ddot{y}(t) + 5\dot{y}(t) + 6y(t) = u(t) + 7\dot{u}(t), \quad (10.22)$$

for $u(t) = e^{-t} u_{\text{stp}}(t)$ and subject to the initial conditions $y(0^-) = 2$ and $\dot{y}(0^-) = 3$.

Solution: **Ans.** $Y(s) = \frac{2(s^2+11s+7)}{s^3+6s^2+11s+6}$ and $y(t) = -3e^{-t}u_{\text{stp}}(t) + 22e^{-2t}u_{\text{stp}}(t) - 17e^{-3t}u_{\text{stp}}(t)$.

Using Table 10.10, we take the Laplace transform of both sides of (10.22), yielding

$$(s^2Y(s) - sy(0^-) - \dot{y}(0^-)) + 5(sY(s) - y(0^-)) + 6Y(s) = U(s) + 7sU(s). \quad (10.23)$$

At this point, we could use Table 10.11 to obtain $U(s) = \mathcal{L}\{e^{-t}u_{\text{stp}}(t)\} = \frac{1}{s+1}$, substitute in the given values for $y(0^-)$ and $\dot{y}(0^-)$, and do the algebra with one of the various symbolic manipulation packages (see below) to obtain a closed-form expression for the Laplace transform of the output, $Y(s)$. However, because this is our first example, we'll also show how the calculations would go by hand because, in other Engineering courses, you'll be expected to solve the problem that way.

Continuing with a manual solution,

$$\begin{aligned} (s^2Y(s) - sy(0^-) - \dot{y}(0^-)) + 5(sY(s) - y(0^-)) + 6Y(s) &= U(s) + 7sU(s) \\ &\Downarrow \text{(move initial conditions to the right side)} \\ (s^2Y(s) + 5sY(s) + 6Y(s)) &= (U(s) + 7sU(s)) + (sy(0^-) + \dot{y}(0^-) + 5y(0^-)) \\ &\Downarrow \text{(group terms)} \\ (s^2 + 5s + 6)Y(s) &= (1 + 7s)U(s) + (s + 5)y(0^-) + \dot{y}(0^-) \\ &\Downarrow \text{(divide by } (s^2 + 5s + 6)) \\ Y(s) &= \frac{1 + 7s}{s^2 + 5s + 6}U(s) + \frac{(s + 5)y(0^-) + \dot{y}(0^-)}{s^2 + 5s + 6}, \end{aligned}$$

giving us the Laplace transform¹ of the output $Y(s)$ in terms of the input signal and the initial conditions

$$Y(s) = \underbrace{\frac{1 + 7s}{s^2 + 5s + 6}U(s)}_{\text{forced response}} + \underbrace{\frac{(s + 5)y(0^-) + \dot{y}(0^-)}{s^2 + 5s + 6}}_{\text{initial condition response}} \quad (10.24)$$

As you may imagine, $\frac{1+7s}{s^2+5s+6}U(s)$ is called the **forced response** because it is the part of the output's response that depends on the input signal, $u(t)$. Similarly, $\frac{(s+5)y(0^-)+\dot{y}(0^-)}{s^2+5s+6}$ is called the **initial condition response** because it is the part of the output's response that depends on the initial conditions, $y(0^-)$ and $\dot{y}(0^-)$. The fact that the **total response** is the sum of the forced response and the initial condition response is called the **superposition principle**.

Soon, we'll call $G(s) := \frac{7s+1}{s^2+5s+6}$ the **transfer function** of the input-output system. Note that you can write it down without doing any computations whatsoever. The denominator comes directly from the left side of the ODE, while the numerator comes from the right side. The name "transfer function" comes from $G(s)$ "transferring the effect of the input, $U(s)$, to the output, $Y(s)$ ", via

$$Y(s) = G(s)U(s).$$

The next step in the process is to use Table 10.11 to obtain $U(s) = \mathcal{L}\{e^{-t}u_{\text{stp}}(t)\} = \frac{1}{s+1}$ and to substitute in actual values for $y(0^-)$ and $\dot{y}(0^-)$. Because it's not that bad, we'll write it out,

$$\begin{aligned} Y(s) &= \frac{1 + 7s}{s^2 + 5s + 6}U(s) + \frac{(s + 5)y(0^-) + \dot{y}(0^-)}{s^2 + 5s + 6} \\ &= \frac{1 + 7s}{s^2 + 5s + 6} \cdot \frac{1}{s + 1} + \frac{(s + 5)2 + 3}{s^2 + 5s + 6} \\ &= \frac{2(s^2 + 11s + 7)}{s^3 + 6s^2 + 11s + 6}, \end{aligned}$$

¹Recall, all of the tedious algebra can be done with symbolic manipulation.

after a bit of algebra.

Here comes a key part of the solution process: To compute $y(t)$, we need to evaluate the inverse Laplace transform of $Y(s)$. A quick perusal of Table 10.11 shows that we are out of luck. Such a complicated Laplace transform is nowhere to be seen. **The “trick” is to perform a partial fraction expansion (PFE) of $Y(s)$** , just as we did in Chapter 7.3.5 when determining antiderivatives for rational functions of x .

The key step in computing a PFE is finding the denominator’s roots. Here, we have

$$s^3 + 6 \cdot s^2 + 11 \cdot s + 6 = (s + 3)(s + 2)(s + 1).$$

Hence, the PFE takes the form

$$\frac{2 \cdot (s^2 + 11 \cdot s + 7)}{s^3 + 6 \cdot s^2 + 11 \cdot s + 6} = \frac{k_1}{s + 3} + \frac{k_2}{s + 2} + \frac{k_3}{s + 1}.$$

Instead of computing the coefficients by hand, we appeal to SymPy, yielding

$$Y(s) = \frac{-17}{s + 3} + \frac{22}{s + 2} - \frac{3}{s + 1}.$$

Now, we can compute inverse Laplace transforms of each of these terms via Table 10.11, because

$$\frac{-17}{s + 3} \longleftrightarrow -17e^{-3t}u_{\text{stp}}(t)$$

$$\frac{22}{s + 2} \longleftrightarrow 22e^{-2t}u_{\text{stp}}(t)$$

$$\frac{-3}{s + 1} \longleftrightarrow -3e^{-t}u_{\text{stp}}(t).$$

Therefore,

$$y(t) = -3e^{-t}u_{\text{stp}}(t) + 22e^{-2t}u_{\text{stp}}(t) - 17e^{-3t}u_{\text{stp}}(t).$$

Here is how the entire problem can be solved in code.

```
1 using SymPy
2
3 # Define the symbolic variable
4 @SymPy.syms s
5
6 # Define the input's Laplace transform and the initial conditions
7 U = 1/(s+1)
8 yzero = 2.0
9 dyzero = 3.0
10
11 # D(s) Y(s) = RHS
12
13 # Define the right-hand side (RHS) expression
14 RHS = (1+7s)*U + (s+5)*yzero + dyzero
15
16 # Define D(s)
17 D = s^2 + 5*s + 6
18
19 # Define Y as the RHS divided by D
20 Y = RHS / D
21
22 # Simplify Y
23 Y = SymPy.expand(Y) # Forces terms to be multiplied out
24 Y = SymPy.simplify(Y)
25 display(Y)
```

```

26
27 # Compute the partial fraction expansion
28 PFE = apart(Y, s)

```

Output $Y(s) = \frac{2 \cdot (s^2 + 11 \cdot s + 7)}{s^3 + 6 \cdot s^2 + 11 \cdot s + 6}$

$$PFE = -\frac{17}{s+3} + \frac{22}{s+2} - \frac{3}{s+1}$$

It is also possible to compute the inverse Laplace transform in code without first computing a partial fraction expansion. How sweet is that?

1 using PyCall

```

2
3 # Import Python's SymPy
4 sympy = pyimport("sympy")
5
6 # Define the symbols
7 s, t = sympy.symbols("s t")
8
9 # Define the expression Y(s)
10 Y = (2.0*s^2 + 22.0*s + 14.0)/(1.0*s^3 + 6.0*s^2 + 11.0*s + 6.0)
11
12 # Compute the inverse Laplace transform
13 y = sympy.inverse_laplace_transform(Y, s, t)
14
15 println(y)

```

Output

$$-17.0 \cdot \exp(-3.0 \cdot t) \cdot \text{Heaviside}(t) + 22.0 \cdot \exp(-2.0 \cdot t) \cdot \text{Heaviside}(t) - 3.0 \cdot \exp(-t) \cdot \text{Heaviside}(t)$$

Note: $\text{Heaviside}(t)$ is what we are denoting as $u_{\text{stp}}(t)$. Sometimes `sympy` will express the answer in a funky way, with $e^{-3t} u_{\text{stp}}$ factored out, such as

$$(-3 \cdot \exp(2 \cdot t) + 22 \cdot \exp(t) - 17) \cdot \exp(-3 \cdot t) \cdot \text{Heaviside}(t)$$

When $e^{-3t} u_{\text{stp}}(t)$ is multiplied through the expression in parentheses, the answer we gave above is recovered. ■

Amazing Observation: Suppose we did not really care about the exact form of the output $y(t)$, but only wanted to know that

$$\lim_{t \rightarrow \infty} y(t) = 0.$$

Then, we do not need the exact values of the coefficients k_i . Once we know that the roots of the denominator polynomial have negative real parts, Table 10.11 implies that we have decaying exponentials, and the output converges to zero. **Laplace must have been a genius!**

Remark 10.17. In your studies, you will often encounter complex mathematical problems that require symbolic computation, such as solving equations, performing integration, or calculating inverse Laplace transforms. While Julia offers powerful computational capabilities, its native symbolic mathematics packages, `SymPy.jl` and `Symbolics.jl`, sometimes lack certain advanced features or specific functions found in Python's `Sympy` library. To bridge this gap, you can use `PyCall`, a Julia package that allows one to call Python functions directly from Julia. This approach enables you to leverage the full power and extensive functionality of Python's `Sympy`, which is a mature library for symbolic mathematics. The integration of Python's capabilities within Julia's environment exemplifies the flexibility and interoperability of modern programming languages, allowing one to choose the best tool for each specific task and seamlessly combine them to enhance your computational experience.

Example 10.18. Use Laplace transforms to compute the solution to the SISO linear state-variable control system

$$\underbrace{\begin{bmatrix} \dot{x}_1 \\ \dot{x}_2 \end{bmatrix}}_{\dot{x}} = \underbrace{\begin{bmatrix} -5 & 1 \\ -6 & 0 \end{bmatrix}}_A \underbrace{\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}}_x + \underbrace{\begin{bmatrix} 7 \\ 1 \end{bmatrix}}_b u$$

$$y = \underbrace{\begin{bmatrix} 1 & 0 \end{bmatrix}}_c x,$$
(10.25)

for $u(t) = e^{-t}u_{\text{stp}}(t)$ and subject to the initial conditions

$$x(0^-) := \begin{bmatrix} x_1(0^-) \\ x_2(0^-) \end{bmatrix} = \begin{bmatrix} 2 \\ 13 \end{bmatrix}.$$

Solution: **Ans.** $Y(s) = \frac{2(s^2+11s+7)}{s^3+6s^2+11s+6}$ and $y(t) = -3e^{-t}u_{\text{stp}}(t) + 22e^{-2t}u_{\text{stp}}(t) - 17e^{-3t}u_{\text{stp}}(t)$.

We do the solution by hand to explain the “tricky parts”, such as, **what could it even mean to apply the Laplace transform to vectors and matrices?** Answer: because everything is linear, applying the Laplace transform row by row works great. **After this one time, never again will we grind it out by hand! So pay attention.** You’ll see that all the computations can be done in matrix-vector form, as any good citizen of ROB 101 would want to do them. Once we have everything in matrix-vector form, we can code it up and do it in a snap!

To apply the Laplace transform row by row, we rewrite (10.25) as

$$\begin{aligned} \dot{x}_1 &= -5x_1 + x_2 + 7u \\ \dot{x}_2 &= -6x_1 + u \\ y &= x_1 \end{aligned} \tag{10.26}$$

This sends us back to ROB 101 Computational Linear Algebra, Lecture 1, before we knew how to express systems of linear equations in vector-matrix form. Taking Laplace transforms line by line gives,

$$\begin{aligned} sX_1(s) - x_1(0^-) &= -5X_1(s) + X_2(s) + 7U(s) \\ sX_2(s) - x_2(0^-) &= -6X_1(s) + U(s) \\ Y(s) &= X_1(s). \end{aligned} \tag{10.27}$$

Moving the initial conditions to the right side yields,

$$\begin{aligned} sX_1(s) &= -5X_1(s) + X_2(s) + 7U(s) + x_1(0^-) \\ sX_2(s) &= -6X_1(s) + U(s) + x_2(0^-) \\ Y(s) &= X_1(s). \end{aligned} \tag{10.28}$$

We now rewrite the above in matrix-vector form. It’s not hard if you work through it step by step. Define

$$X(s) := \begin{bmatrix} X_1(s) \\ X_2(s) \end{bmatrix}. \tag{10.29}$$

Then, we can write the first two rows of (10.28) as

$$\begin{aligned}
 sX(s) &= AX(s) + bU(s) + x(0^-) \\
 &\Downarrow \\
 sX(s) - AX(s) &= bU(s) + x(0^-) \\
 &\Downarrow \\
 sI_2 X(s) - AX(s) &= bU(s) + x(0^-) \\
 &\Downarrow \\
 (sI_2 - A) X(s) &= bU(s) + x(0^-),
 \end{aligned} \tag{10.30}$$

for A and b as defined in (10.25). Moreover, the last row becomes

$$Y(s) = cX(s), \tag{10.31}$$

with c as defined in (10.25). Solving for $X(s)$ in (10.30) gives

$$X(s) = (sI - A)^{-1} (bU(s) + x(0^-)). \tag{10.32}$$

Substituting (10.32) into (10.31) gives

$$Y(s) = c (sI - A)^{-1} (bU(s) + x(0^-)). \tag{10.33}$$

The initial conditions look like a modification of the input. We write the final answer as

$$Y(s) = \underbrace{c (sI - A)^{-1} bU(s)}_{\text{forced response}} + \underbrace{c (sI - A)^{-1} x(0^-)}_{\text{initial condition response}}. \tag{10.34}$$

The term multiplying the input, $G(s) := c (sI - A)^{-1} b$, will soon be called the **transfer function**. While $G(s)$ may not look like a rational function of s , it can be proven to be one using the **adjoint method** of computing the matrix inverse. We'll compute the inverse symbolically to show the desired property for small examples.

Now that we have the general formula, we'll do the computation of the above in code!

```

1 using SymPy, LinearAlgebra
2
3 # Define the symbolic variable
4 @SymPy.syms s
5
6 # Linear state variable model
7 A=[-5 1;-6 0.0]
8 b=[7.0; 1]
9 c=[1.0 0]
10
11 # Initial Conditions
12 x0=[2.0; 13]
13
14 # Laplace transform of the input
15 U = 1/(s+1)
16
17 # Compute Y(s)
18 Y = c*inv(s*I(2) - A)*b*U + c*inv(s*I(2) - A)*x0
19
20 # Y is a 1 x 1 vector, so in Julia, we extract it
21 Y = Y[1]
22
23 # Simplify Y
24 Y = SymPy.expand(Y) # Forces terms to be multiplied out

```

```

25 Y = SymPy.simplify(Y)
26 display(Y)
27 println(Y)
28
29
30 # Compute the partial fraction expansion
31 PFE = apart(Y, s)

```

Output

$$Y(s) = \frac{72 \cdot s^2 + 792 \cdot s + 504}{36 \cdot s^3 + 216 \cdot s^2 + 396 \cdot s + 216}$$

$$(72.0 \cdot s^2 + 792.0 \cdot s + 504.0) / (36.0 \cdot s^3 + 216.0 \cdot s^2 + 396.0 \cdot s + 216.0)$$

$$PFE = \frac{-3.000}{s+1.000} + \frac{11.000}{0.500 \cdot s + 1.000} - \frac{5.667}{0.333 \cdot s + 1.000}$$

If you look carefully, you can see that SymPy has not removed a factor of 36 from the numerator and denominator of $Y(s)$. Moreover, in the PFE, if we multiplied the second term by $\frac{2}{2}$ and the third term by $\frac{3}{3}$ we'd obtain simpler answers, but we leave them as such because we just don't care.

Next, we obtain the inverse Laplace transform.

```

1 using PyCall
2
3 # Import Python's SymPy
4 sympy = pyimport("sympy")
5
6 # Define the symbols
7 s, t = sympy.symbols("s t")
8
9 # Define the expression Y(s)
10 Y = (72.0*s^2 + 792.0*s + 504.0) / (36.0*s^3 + 216.0*s^2 + 396.0*s + 216.0)
11
12 # Compute the inverse Laplace transform
13 y = sympy.inverse_laplace_transform(Y, s, t)
14
15 println(y)

```

Output

$$-17.0 \cdot \exp(-3.0 \cdot t) \cdot \text{Heaviside}(t) + 22.0 \cdot \exp(-2.0 \cdot t) \cdot \text{Heaviside}(t) - 3.0 \cdot \exp(-t) \cdot \text{Heaviside}(t)$$

the same answer we obtained in the previous example. **Ta-da! Input-output and state-variable methods are really two sides of the same coin once you take Michigan's EECS 560 Linear Systems Theory.**

10.3.5 Transfer Functions

A transfer function for an LTI (Linear Time-Invariant) input-output system is a Laplace (aka, frequency) domain representation that describes the relationship between the input and output of a system **when the initial conditions are identically zero**. It is a ratio of two polynomials in terms of the Laplace variable s . The transfer function allows for the study of system stability without directly solving differential equations in the time domain. By examining the transfer function's "poles and zeros" (defined shortly), one can infer important properties about the system's response to a wide variety of inputs.

Transfer Function of a SISO Linear System

Definition 10.19. *The transfer function of the linear input output model*

$$y^{(n)} + a_{n-1}y^{(n-1)} + a_{n-2}y^{(n-2)} + \dots + a_1\dot{y} + a_0y = b_0u + b_1\dot{u} + \dots + b_mu^{(m)}$$

is the **rational function**

$$G(s) := \frac{Y(s)}{U(s)} = \frac{b_ms^m + \dots + b_1s + b_0}{s^n + a_{n-1}s^{n-1} + \dots + a_1s + a_0}. \quad (10.35)$$

Note that $Y(s) = G(s) \cdot U(s)$ is precisely the Laplace transform of the linear system (10.1) when all **initial conditions are zero**.

Note: When you need the transfer function of an input-output system, you simply write it down by inspection. No calculations are required.

Definition 10.20. *The transfer function of the SISO linear state-variable model $\dot{x} = Ax + bu, y = cx$ is*

$$G(s) := \frac{Y(s)}{U(s)} = c(sI - A)^{-1}b. \quad (10.36)$$

Note that $Y(s) = G(s) \cdot U(s)$ is precisely the Laplace transform of the linear state-variable control system (10.2) when all **initial conditions are zero**.

Note: While the transfer function of a SISO linear state-variable model usually cannot be determined by inspection, it does not matter because we use code to compute it.

10.3.6 The Segway Transporter à la Laplace

We derive a transfer function representation for the Segway from the state-variable model given in Fact 10.5.

Example 10.21. *For the Segway transporter, use one of the symbolic computation packages in Julia to compute the Laplace transform of the state-variable model $\dot{x} = Ax + bu$, assuming zero initial conditions,*

$$H(s) := \frac{X(s)}{U(s)} = (sI_4 - A)^{-1}b.$$

In addition, interpret the components of $H(s)$.

Solution: Ans. $H(s) = \begin{bmatrix} \frac{-21.468}{s^2.000-50.937} \\ \frac{80.227 \cdot s^2-594.021}{s^2 \cdot (s^2-50.937)} \\ \frac{-21.468 \cdot s}{s^2-50.937} \\ \frac{80.227 \cdot s^2-594.021}{s \cdot (s^2-50.937)} \end{bmatrix} = \begin{bmatrix} \frac{\theta(s)}{U(s)} \\ \frac{\varphi(s)}{U(s)} \\ \frac{\dot{\theta}(s)}{U(s)} \\ \frac{\dot{\varphi}(s)}{U(s)} \end{bmatrix}.$

The components of $H(s)$ are as labeled because the components of the state vector x are $\theta, \varphi, \dot{\theta}$, and $\dot{\varphi}$. Moreover, because $\mathcal{L}\{\dot{\theta}\} = s \cdot \mathcal{L}\{\theta\}$ and $\mathcal{L}\{\dot{\varphi}\} = s \cdot \mathcal{L}\{\varphi\}$, we expect to see this reflected in the transfer functions that comprise $\frac{X(s)}{U(s)}$. Indeed, $\frac{\dot{\theta}(s)}{U(s)} = s \cdot \frac{\theta(s)}{U(s)}$ and $\frac{\dot{\varphi}(s)}{U(s)} = s \cdot \frac{\varphi(s)}{U(s)}$. Equivalently,

$$\frac{\theta(s)}{U(s)} = \frac{1}{s} \cdot \frac{\dot{\theta}(s)}{U(s)} \text{ and } \frac{\varphi(s)}{U(s)} = \frac{1}{s} \cdot \frac{\dot{\varphi}(s)}{U(s)}.$$

In other words, the Laplace variable s represents differentiation, and $\frac{1}{s}$ represents integration.

The following code computes the linearized state-variable model and the “transfer vector”, $H(s) := \frac{X(s)}{U(s)} = (sI_4 - A)^{-1}b$. Here, we are using SymPy, but it could just as easily be done with Symbolics or PyCall. So many decisions....

```
1 using SymPy, LinearAlgebra
2
3 # bring the model into the workspace for later use
4 include("dyn_mod_planarSegway.jl")
5
6 # Set equilibrium Point
7 qe = [0.0; 0]
8
9 F = dyn_mod_planarSegway(qe, 0*qe)
10 D = F.D;
11 JacG = F.JacG;
12 JacComx = F.JacComx
13 B = F.B;
14
15 # build the linearized model about the given equilibrium point
16 A21 = -D\JacG;
17 n, m = size(B)
18 A = [zeros(n, n) I(n); A21 zeros(n, n)];
19 b = [zeros(n, m); F.D \ B];
20 display(A)
21 display(b)
22
23
24 # Define the symbolic variable
25 @SymPy.syms s
26 H = inv(s*I(4) - A)*b
27
28 # Simplify H
29 H = SymPy.expand.(H) # Forces terms to be multiplied out
30 H = SymPy.simplify.(H)
31 #display(H)
32 println(H)
```

Output

```
4×4 Matrix{Float64}:
 0.0      0.0    1.0    0.0
 0.0      0.0    0.0    1.0
 50.9373  0.0    0.0    0.0
-162.687  0.0    0.0    0.0

4×1 Matrix{Float64}:
 0.0
 0.0
-21.467741053946913
 80.22714380349518

Sym[-1093.50791254434 / (50.9372602267016*s^2 - 2594.60447940271);
(4086.55090116364*s^2 - 30257.7781854544) / (s^2*(50.9372602267016*s^2 - 2594.60447940271));
-21.4677410539469*s / (s^2 - 50.9372602267016);
(4086.55090116364*s^2 - 30257.7781854544) / (s*(50.9372602267016*s^2 - 2594.60447940271));]
```

The primary control task for a Segway is to keep it upright. This means we need to keep the lean angle of the platform plus rider near zero. In fact, for a Segway you would ride yourself, maintaining the Segway upright is the only task the feedback controller faces because it is assumed that the rider will regulate speed through their own body lean angle. Later in this Chapter, we will also design a secondary controller that adjusts the overall lean angle so as to achieve a desired “speed”.

Example 10.22. (Segway’s Transfer Functions) Compute the SISO transfer functions $G_\theta(s) := \frac{Y_\theta(s)}{U(s)}$ and $G_v(s) := \frac{Y_v(s)}{U(s)}$, where the latter uses the level-1 notion of speed defined in Fact 10.5.

Solution: **Ans.** The code below carries out the computations required to obtain the results in (10.37); they are then presented pictorially in Figure 10.3.

$$G_\theta(s) = \frac{-21.468}{s^2 - 50.937} \quad (10.37)$$

$$G_v(s) = \frac{-20.057(s^2 - 7.404)}{s(s^2 - 50.937)}.$$

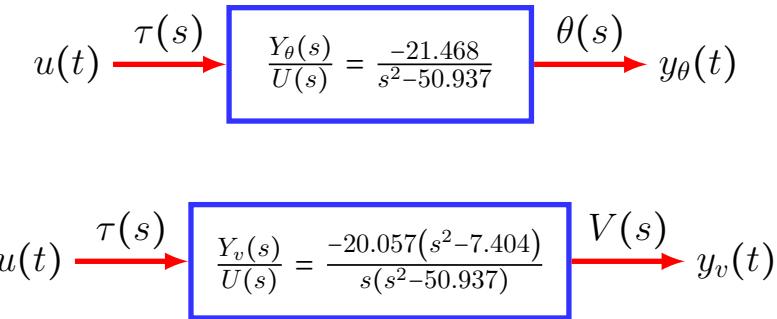


Figure 10.3: The Segway’s transfer functions for the lean angle and cart speed are as given, where $u(t)$ is the motor torque. At the end of the Chapter, we will seek to control **simultaneously** these two key quantities. Can you see any special structure that may help us with that?

```

1 using LinearAlgebra, Symbolics
2
3 # Set equilibrium Point
4 qe = [0.0; 0]
5
6 F = dyn_mod_planarSegway(qe, 0*qe)
7 D = F.D;
8 JacG = F.JacG;
9 B = F.B;
10
11 # build the linearized model about the given equilibrium point
12 A21 = -D\JacG; # A21 = cleanUp(A21);
13 n, m = size(B)
14 A = [zeros(n,n) I(n); A21 zeros(n,n)];
15 b = [zeros(n,m); F.D \ B];
16 display(A)
17 display(b)
18
19 # Define the various outputs
20 c_th = [1.0 0 0 0]
21
22 # model params

```

```

23 g, m_wh, m_pend, L, r_wh, J_pend, J_wh, J_rotor, N = modelParameters()
24
25 c_v1 = [0 0 0 r_wh]
26
27
28 # Define the symbolic variables and transfer functions
29 @Symbolics.variables s
30 G_th = c_th*inv(s*I(4) - A)*b
31 G_v1 = c_v1*inv(s*I(4) - A)*b
32
33 # Simplify
34 G_th = Symbolics.simplify(G_th[1])
35 println(G_th)
36
37 println(" ")
38
39 G_v1 = Symbolics.simplify(G_v1[1])
40 println(G_v1)

```

Output

```

4×4 Matrix{Float64}:
 0.0      0.0    1.0    0.0
 0.0      0.0    0.0    1.0
 50.9373  -0.0   0.0    0.0
-162.687  -0.0   0.0    0.0

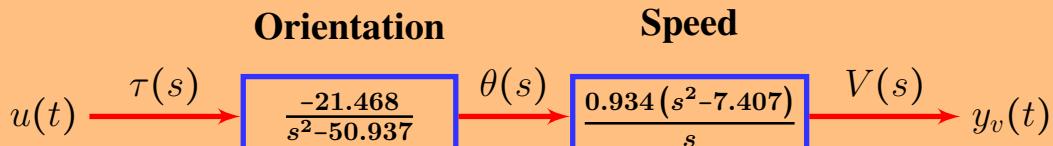
4×1 Matrix{Float64}:
 0.0
 0.0
-21.467741053946913
 80.22714380349518

-21.467741053946913 / (-50.93726022670157 + s^2)
(-148.50513185627278 + 20.056785950873795(s^2)) / (-50.93726022670157s + s^3)

```

Fact 10.23: Factoring the Segway's Transfer Function

Because transfer functions are rational functions of the Laplace variable, s , the rules of ordinary algebra apply to them. While factoring an ODE seems far-fetched at best, expressing polynomials and rational functions in terms of factors is ordinary. Applying this idea to the Segway, its transfer function from motor torque to cart speed can be factored into two terms, as shown below. Indeed, $\frac{-20.057}{-21.468} = 0.934$ to three significant digits, while the other terms are “self-evident”.



The above factorization shows mathematically how the orientation of the Segway's body directly affects its speed, something that is much harder to see from the ODE itself. It suggests a two-pronged attack on how to control the Segway: first, stabilize the orientation and, second, use the orientation as a pseudo-actuator for the speed.

10.4 Poles, Zeros, and BIBO Stability

In Chapter 9.5.8, we studied linear ODEs of the form $\dot{x} = Ax$, $x(t_0) = x_0$, and learned that all solutions asymptotically converged to the origin if, and only if, the real parts of the eigenvalues of A have negative real parts. **The main message of this section is that the roots of the denominator of a transfer function work just like eigenvalues: if their real parts are negative, then bounded input signals give rise to bounded output signals.**

10.4.1 Poles and Zeros

A transfer function is **rational** if it can be expressed as the ratio of two polynomials.

Poles and Zeros

Definition 10.24. For a rational transfer function $G(s) = \frac{N(s)}{D(s)}$, the polynomials $N(s)$ and $D(s)$ are said to be **coprime** if they have no common roots. In this case, the roots of $N(s)$ are the **zeros** of $G(s)$ and the roots of $D(s)$ are the **poles**.

Note: If $N(s)$ and $D(s)$ have common factors, hence yielding common roots, the common factors must be canceled before determining the zeros and poles of the transfer function.

Example 10.25. Determine the zeros and poles of the transfer functions

$$(a) G_1(s) = \frac{2}{s^2 - 1}$$

$$(b) G_2(s) = \frac{s+2}{s^2 - 1}$$

$$(c) G_3(s) = \frac{s+1}{s^2 - 1}$$

$$(d) G_4(s) = \frac{2s-4}{s^2 + 2s + 2}.$$

Solution:

(a) $N_1(s) = 2$ has no roots and $D_1(s) = s^2 - 1 = (s - 1)(s + 1)$ has roots ± 1 . $N_1(s)$ and $D_1(s)$ are therefore coprime; it follows that $G_1(s)$ has no zeros and two poles at ± 1 .

(b) $N_2(s) = s + 2$ has one root -2 which is distinct from the roots of $D_2(s) = s^2 - 1$, which are ± 1 . Thus $N_2(s)$ and $D_2(s)$ are coprime and it follows that $G_2(s)$ has one zero at -2 and two poles at ± 1 .

(c) $N_3(s) = s + 1$ has one root -1 which is also a root of $D_3(s) = s^2 - 1$. Thus $N_3(s)$ and $D_3(s)$ are *not coprime*. To determine the zeros and poles of $G_3(s)$, the common factor of $(s + 1)$ is removed, yielding

$$G_3(s) = \frac{s + 1}{(s + 1)(s - 1)} = \frac{1}{s - 1}.$$

It follows that $G_3(s)$ has no zeros and one pole at 1.

- (d) $N_4(s) = 2(s - 2)$ has a root at 2, which is distinct from the roots of $D_4(s) = s^2 + 2s + 2 = (s + 1 + i)(s + 1 - i)$, which are $-1 \pm i$. $N_4(s)$ and $D_4(s)$ are therefore coprime; it follows that $G_4(s)$ has a zero at 2 and two poles at $-1 \pm i$.

10.4.2 BIBO Stability

We treated bounded functions way back in Chapter 4.6. An **input or output signal** $s : [0, \infty) \rightarrow \mathbb{R}$ is **bounded** if there exists a constant M such that $|s(t)| \leq M$ for all $t \geq 0$. Hence, signals such as constants, $\sin(\omega t) u_{\text{stp}}(t)$, $\cos(\omega t) u_{\text{stp}}(t)$, $t^k e^{at} u_{\text{stp}}(t)$ for $a < 0$, and products of such functions are bounded. Functions such as $e^{at} u_{\text{stp}}(t)$ and $t^k u_{\text{stp}}(t)$ for $a > 0$ and $k > 0$ are unbounded.

BIBO Stability

Definition 10.26. A system is bounded-input bounded-output *stable*, or **BIBO stable** for short, if the output is bounded whenever the input is bounded.

If the system is **not BIBO stable**, it is said to be **unstable**.

Note: To be extra clear, the word “whenever” here means the output must remain bounded for “all possible” bounded inputs, $u : [0, \infty) \rightarrow \mathbb{R}$. That sounds like it could be hard to check!

Recall that a system is **causal** if present values of the output do not depend on future values of the input. If the system has a rational transfer function $G(s)$, then the system is causal if, and only if, the degree of the numerator is less than or equal to the degree of the denominator. **Transfer functions that satisfy this property have a special name; they are said to be proper.**

Proposition 10.27: Poles and BIBO Stability

A system with a proper rational transfer function is BIBO stable if, and only if, its poles have negative real parts. In other words, the stability of a rational transfer function

$$G(s) = \frac{N(s)}{D(s)},$$

is completely determined by its poles as long as the degree of the numerator polynomial is less than or equal to the degree of the denominator polynomial.

Note: The zeros do not play a role in stability.

To illustrate the connections between time-domain models, transfer functions, poles, and BIBO stability, consider a system represented by a first-order differential equation with input $u(t)$ and output $y(t)$,

$$\dot{y}(t) - y(t) = u(t).$$

In the s-domain, the model becomes $sY(s) - Y(s) = U(s)$, yielding

$$Y(s) = G(s)U(s),$$

where the transfer function is

$$G(s) = \frac{1}{s - 1}.$$

The pole is $s = +1$, and thus the system is unstable. Is it easy to give a bounded input that results in an unbounded output? Absolutely.

For definiteness, suppose the input is a unit step, $u_{\text{stp}}(t)$. Then, from Table 10.11, $U(s) = \frac{1}{s}$ giving

$$Y(s) = \frac{1}{s - 1} \frac{1}{s}.$$

Applying the method of partial fraction expansion,

$$Y(s) = -\frac{1}{s} + \frac{1}{s - 1}.$$

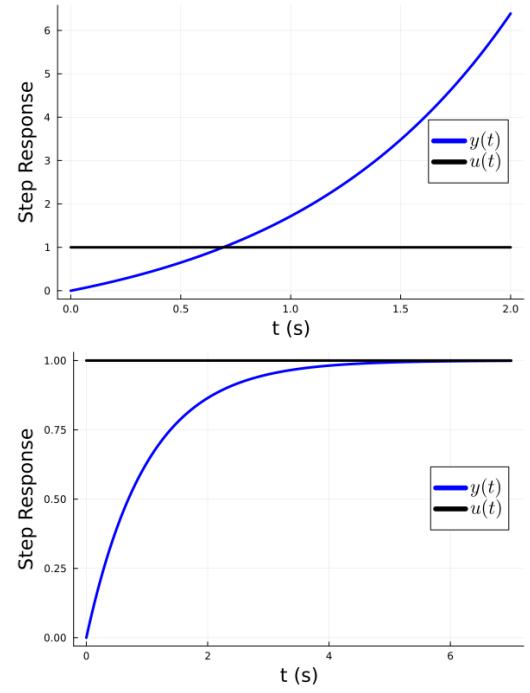
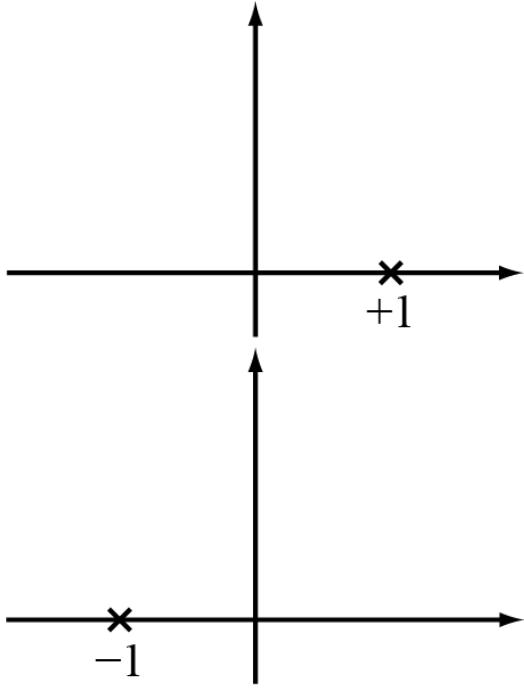


Figure 10.4: The positions of the poles are marked with a \times . For a positive real pole, bounded inputs can give rise to unbounded outputs, whereas for a negative real pole, bounded inputs always give rise to bounded outputs.

Taking the inverse Laplace transform yields

$$y(t) = (e^t - 1) u_{\text{stp}}(t). \quad (10.38)$$

From this equation, the output $y(t)$ is unbounded, diverging to infinity as t goes to infinity, as shown in the upper half of Fig. 10.4.

On the other hand, for the model

$$\dot{y}(t) + y(t) = u(t),$$

the transfer function is

$$G(s) = \frac{1}{s+1}.$$

The pole is $s = -1$, and hence the system is BIBO stable. Computing the output in response to a unit-step input gives

$$y(t) = (1 - e^{-t}) u_{\text{stp}}(t). \quad (10.39)$$

From this equation, the output converges to 1.0 as t goes to infinity, and, in particular, it is bounded for all $t \geq 0$, as illustrated in the lower half of Fig. 10.4.

Example 10.28. Determine the stability of the transfer functions

$$(a) G_1(s) = \frac{2}{s^2 - 1}$$

$$(b) G_2(s) = \frac{s+2}{s^2 - 1}$$

$$(c) G_3(s) = \frac{s+1}{s^2 - 1}$$

$$(d) G_4(s) = \frac{2s-4}{s^2 + 2s + 2}$$

$$(e) G_5(s) = \frac{s-1}{s},$$

where the first four transfer functions are taken from the previous example, and the last one is new.

Solutions: $G_1(s)$, $G_2(s)$ and $G_3(s)$ each has a pole at $+1$, and thus they are unstable. The poles of $G_4(s)$ are $-1 \pm i$, which have negative real parts, and hence the transfer function is BIBO stable. $G_5(s)$ has a pole at the origin, and hence it is unstable. ■

At this point, the role of zeros seems somewhat mysterious. The zeros will come into the picture when closed-loop systems are analyzed.

Example 10.29. (Optional Read:) If a rational transfer function is not proper, do the poles still determine its stability? The short answer is no.

Solution: If you are curious as to why not, consider $G(s) = s$, which corresponds to a differentiator, that is, the system

$$y(t) = \frac{du(t)}{dt},$$

where $u(t)$ is the input. The input $\sin(t^2)$ is bounded, since it takes values between -1 and 1 , but the corresponding output $y(t) = 2t \cos(t^2)$ is unbounded.

The hardcore among you may think that the transfer function has no poles, and therefore, it is untrue that all of its poles have negative real parts. Instead of arguing semantics, consider the transfer function,

$$\begin{aligned} G(s) &= G_a(s) + G_b(s) \\ &= s + \frac{1}{s+2} \\ &= \frac{s(s+2)}{s+2} + \frac{1}{s+2} \\ &= \frac{s^2 + 2s + 1}{s+2}, \end{aligned}$$

which is improper. A quick calculation shows that it has two zeros at -1 and a pole at -2 . But the system is not BIBO stable because the bounded input $u(t) = \sin(t^2)$ produces the output

$$y(t) = y_a(t) + y_b(t),$$

where $y_a(t) = 2t \cos(t^2)$, which is unbounded, and the second term, being the output of a system with BIBO stable transfer function $G_b(s)$, must be bounded. The sum of an unbounded signal and a bounded signal is unbounded. Thus, a bounded input has produced an unbounded output, showing the system is unstable. ■

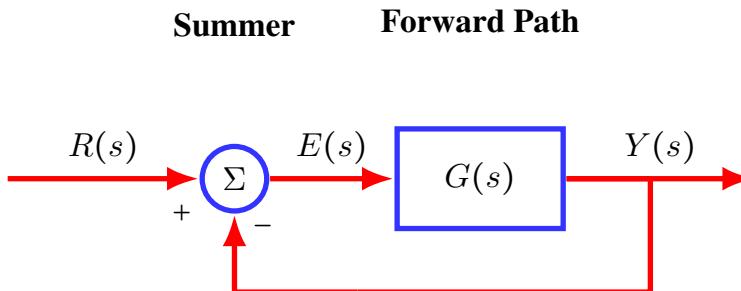


Figure 10.5: Unity negative feedback system.

10.5 Unity Feedback Systems

A common control system configuration, and the easiest one to study when first learning the subject of feedback control, is the **unity feedback system** shown in Fig. 10.5. It is also called a **negative feedback system** due to the minus sign on the feedback loop at

the summer. In this figure, $Y(s)$ is the system output, $R(s)$ is the reference command, and $E(s)$ is the error between the output and the reference. The reference command $R(s)$ normally specifies the desired output of the system. When the reference command is a constant, the reference value is also called a **setpoint**. The transfer function in the **forward path** is $G(s)$, whereas the transfer function in the **feedback path** (the bottom loop) is one, the source of the name, **unity feedback system**.

10.5.1 Closed-loop Transfer Functions

The block diagram in Fig. 10.5 is equivalent to the following set of linear equations

$$E(s) = R(s) - Y(s) \quad (10.40)$$

$$Y(s) = G(s)E(s). \quad (10.41)$$

Solving these two equations for $Y(s)$ as a function of $R(s)$ results in

$$\begin{aligned} Y(s) &= G(s)(R(s) - Y(s)) \\ &= G(s)R(s) - G(s)Y(s), \end{aligned}$$

and thus

$$Y(s) + G(s)Y(s) = G(s)R(s),$$

which is equivalent to

$$(1 + G(s))Y(s) = G(s)R(s).$$

Hence

$$Y(s) = \frac{G(s)}{1 + G(s)}R(s), \quad (10.42)$$

from which it follows that the transfer function from $R(s)$ to $Y(s)$ is

$$\frac{Y(s)}{R(s)} = \frac{G(s)}{1 + G(s)}. \quad (10.43)$$

The same steps can be repeated to determine the transfer function from $R(s)$ to $E(s)$,

$$E(s) = R(s) - G(s)Y(s)$$

and

$$(1 + G(s))E(s) = R(s).$$

Thus,

$$E(s) = \frac{1}{1 + G(s)}R(s), \quad (10.44)$$

which yields the transfer function as

$$\frac{E(s)}{R(s)} = \frac{1}{1 + G(s)}. \quad (10.45)$$

Example 10.30. For the unity feedback system in Fig. 10.5, compute the poles of $\frac{Y(s)}{R(s)}$ and $\frac{E(s)}{R(s)}$ for each of the following forward-path transfer functions,

$$(a) G_1(s) = \frac{2}{s^2 - 1}$$

$$(b) G_2(s) = \frac{s+2}{s^2 - 1}$$

$$(c) G_3(s) = \frac{s+1}{s^2 - 1}.$$

Solutions: Before starting the solutions, suppose you had the following ratio of rational numbers

$$\frac{\frac{3}{4}}{1 + \frac{3}{4}},$$

and you wanted to reduce it to a standard ratio of two integers. You would be unlikely to find the problem intimidating because you know how to clear the denominator as follows,

$$\begin{aligned} \frac{\frac{3}{4}}{1 + \frac{3}{4}} &= \frac{3}{4 + 3} \quad (\text{multiply top and bottom by 4}) \\ &= \frac{3}{7}. \end{aligned}$$

Because there are no common factors, this is the final answer. If you somehow arrived at $\frac{6}{14}$, you would remove the common factor of two (aka, 6 and 14 are not **coprime**). We are going to do the same for rational functions. **What is mindbogglingly cool is that these rational functions are transfer functions and, THEREFORE, correspond to ODEs!** Never in a million years would you have imagined manipulating derivatives like this. Laplace was a genius.

- (a) $G_1(s) = \frac{2}{s^2 - 1}$ **Ans.** Poles of $\frac{Y(s)}{R(s)}$ and $\frac{E(s)}{R(s)}$ are the same and equal $\pm i$. Hence, the closed-loop system is not BIBO stable.

$$\begin{aligned} \frac{Y(s)}{R(s)} &= \frac{G_1(s)}{1 + G_1(s)} \\ &= \frac{\frac{2}{s^2 - 1}}{1 + \frac{2}{s^2 - 1}} \quad (\text{multiply top and bottom by } (s^2 - 1)) \\ &= \frac{2}{(s^2 - 1) + 2} \\ &= \frac{2}{s^2 + 1}. \end{aligned}$$

Hence, the poles are $\pm i$. We did not have to worry about a common factor because the numerator is a constant and has no roots. In a similar fashion, we compute

$$\begin{aligned} \frac{E(s)}{R(s)} &= \frac{1}{1 + G_1(s)} \\ &= \frac{1}{1 + \frac{2}{s^2 - 1}} \quad (\text{multiply top and bottom by } (s^2 - 1)) \\ &= \frac{2(s^2 - 1)}{(s^2 - 1) + 2} \\ &= \frac{2(s^2 - 1)}{s^2 + 1}. \end{aligned}$$

The numerator has roots ± 1 , while the denominator has roots $\pm i$. There are no common factors and hence the poles are $\pm i$.

- (b) $G_2(s) = \frac{s+2}{s^2 - 1}$ **Ans.** Poles of $\frac{Y(s)}{R(s)}$ and $\frac{E(s)}{R(s)}$ are the same and equal $-\frac{1}{2} \pm i\frac{\sqrt{3}}{2}$. Hence, the closed-loop system is BIBO stable.

$$\begin{aligned} \frac{Y(s)}{R(s)} &= \frac{G_2(s)}{1 + G_2(s)} \\ &= \frac{\frac{s+2}{s^2 - 1}}{1 + \frac{s+2}{s^2 - 1}} \quad (\text{multiply top and bottom by } (s^2 - 1)) \\ &= \frac{s+2}{(s^2 - 1) + (s+2)} \\ &= \frac{s+2}{s^2 + s + 1}. \end{aligned}$$

The root of the numerator is -2 and the roots of the denominator are $-\frac{1}{2} \pm i\frac{\sqrt{3}}{2}$. Because there are no common factors, the poles are $-\frac{1}{2} \pm i\frac{\sqrt{3}}{2}$. In a similar fashion, we compute

$$\begin{aligned}\frac{E(s)}{R(s)} &= \frac{1}{1 + G_2(s)} \\ &= \frac{1}{1 + \frac{s+2}{s^2-1}} \quad (\text{multiply top and bottom by } (s^2 - 1)) \\ &= \frac{(s^2 - 1)}{(s^2 - 1) + (s + 2)} \\ &= \frac{(s^2 - 1)}{s^2 + s + 1}.\end{aligned}$$

The numerator has roots ± 1 , while the denominator has roots $-\frac{1}{2} \pm i\frac{\sqrt{3}}{2}$. There are no common factors and hence the poles are $-\frac{1}{2} \pm i\frac{\sqrt{3}}{2}$.

- (c) $G_3(s) = \frac{s+1}{s^2-1}$ **Ans.** The pole of $\frac{Y(s)}{R(s)}$ and $\frac{E(s)}{R(s)}$ is the same and equals $s = 0$. Hence, the closed-loop system is not BIBO stable.

$$\begin{aligned}\frac{Y(s)}{R(s)} &= \frac{G_3(s)}{1 + G_3(s)} \\ &= \frac{\frac{s+1}{s^2-1}}{1 + \frac{s+1}{s^2-1}} \quad (\text{multiply top and bottom by } (s^2 - 1)) \\ &= \frac{s+1}{(s^2 - 1) + (s + 1)} \\ &= \frac{s+1}{s^2 + s} \\ &= \frac{s+1}{s(s+1)} \\ &= \frac{1}{s}\end{aligned}$$

The numerator has no roots (because we saw the cancellation and took care of it) and the root of the denominator is 0. Because there are no common factors, the pole is $s = 0$. In a similar fashion, we compute

$$\begin{aligned}\frac{E(s)}{R(s)} &= \frac{1}{1 + G_3(s)} \\ &= \frac{1}{1 + \frac{s+1}{s^2-1}} \quad (\text{multiply top and bottom by } (s^2 - 1)) \\ &= \frac{(s^2 - 1)}{(s^2 - 1) + (s + 1)} \\ &= \frac{(s^2 - 1)}{s^2 + s} \\ &= \frac{(s+1)(s-1)}{s(s+1)} \\ &= \frac{s-1}{s}.\end{aligned}$$

The numerator has a root at $+1$, while the denominator has a root at 0. There are no common factors and hence the pole is $s = 0$.

Note: The common factor in the numerator and denominator can be traced back to the common factor in the open-loop transfer function,

$$G_3(s) = \frac{s+1}{s^2-1} = \frac{s+1}{(s-1)(s+1)}.$$

If the common factor had been removed BEFORE computing the closed-loop transfer functions, the resulting polynomials would not have had a common factor; that is, they, too, would have been coprime. As part of feedback design, we sometimes like to cancel a “STABLE zero” by a “STABLE pole”. We’ll come back to this point later. ■

We deliberately went through every step in each problem to underline the repetitive nature of the computations.

Proposition 10.31: Closed-Loop Transfer Functions

Suppose that $G(s) = \frac{N(s)}{D(s)}$ is the forward-path transfer function for the unity feedback system in Fig. 10.5. Then, the closed-loop transfer functions are

$$\frac{Y(s)}{R(s)} = \frac{N(s)}{D(s) + N(s)} \quad (10.46)$$

$$\frac{E(s)}{R(s)} = \frac{D(s)}{D(s) + N(s)}.$$

Moreover, if $N(s)$ and $D(s)$ are coprime, then so are the numerators and denominators of $\frac{Y(s)}{R(s)}$ and $\frac{E(s)}{R(s)}$.

Proof: Simple calculations show that the two transfer functions associated with a unity feedback system are as stated. Indeed,

$$\frac{Y(s)}{R(s)} = \frac{\frac{N(s)}{D(s)}}{1 + \frac{N(s)}{D(s)}} = \frac{N(s)}{D(s) + N(s)},$$

and

$$\frac{E(s)}{R(s)} = \frac{1}{1 + \frac{N(s)}{D(s)}} = \frac{D(s)}{D(s) + N(s)}.$$

In each case, the numerator and denominator have been multiplied by $D(s)$ to simplify the expressions to a ratio of polynomials.

Next, we show that if $N(s)$ and $D(s)$ are coprime (no common roots), then the polynomials $N(s)$ and $D(s) + N(s)$ are coprime as well. To show this, suppose that s_1 is a root of $N(s)$, that is, $N(s_1) = 0$. Can it also be a root of $D(s) + N(s)$? The answer is no, because if $N(s_1) = 0$, then $D(s_1) + N(s_1) = D(s_1)$. Hence, if $D(s_1) + N(s_1) = 0$, then $D(s_1) = 0$, which shows that s_1 would then have to be a root of $D(s)$, which is not possible due to $N(s)$ and $D(s)$ being coprime. The identical argument applies to the second transfer function as well. ■

10.5.2 Closed-loop Poles and Zeros

Proposition 10.31 puts us in a position to define the closed-loop poles and zeros of a unity feedback system.

Closed-Loop Poles and Zeros

Definition 10.32. For the unity feedback system in Fig. 10.5, with forward-path transfer function $G(s) = \frac{N(s)}{D(s)}$ coprime, the **closed-loop poles** are the roots of $D(s) + N(s)$ and the **closed-loop zeros** are the roots of $N(s)$. It follows that the zeros of the forward-path transfer function are also the closed-loop zeros of the corresponding unity feedback loop

Notes: (a) If $N(s)$ and $D(s)$ have common factors, hence yielding common roots, the common factors must be canceled before determining the zeros and poles of the closed-loop system. (b) Because we are normally interested in shaping the transient response of the transfer function from reference command to output, the zeros of interest to us are those of $\frac{Y(s)}{R(s)}$.

Even though it is a slight abuse of terminology, the zeros of $G(s)$ are often called the open-loop zeros, and the poles of $G(s)$ are often called the open-loop poles. This terminology makes sense because if one *opens the loop* by disconnecting the feedback path from the summer in Fig. 10.5, the transfer function from $R(s)$ to $Y(s)$ is $G(s)$.

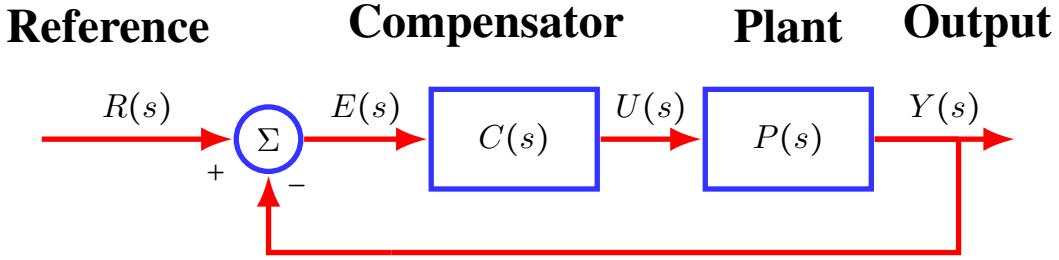


Figure 10.6: Cascade control system configuration, where $C(s)$ is the transfer function of the compensator and $P(s)$ is the transfer function of the plant. Their product, $G(s) = C(s)P(s)$ is the overall forward-path transfer function.

Example 10.33. Give the open-loop (OL) and closed-loop (CL) poles and zeros for the forward-path (FP) transfer functions taken from Example 10.30.

$$(a) G_1(s) = \frac{2}{s^2 - 1}$$

$$(b) G_2(s) = \frac{s+2}{s^2 - 1}$$

$$(c) G_3(s) = \frac{s+1}{s^2 - 1}.$$

Solutions: Most of the work has been done in Examples 10.25, 10.28, and 10.30. The purpose here is to collect it all together and practice the new vocabulary.

FP $G(s)$	OL Zeros	OL Poles	CL Zeros	CL Poles
$\frac{2}{s^2 - 1}$	None	± 1	None	$\pm i$
$\frac{s+2}{s^2 - 1}$	-2	± 1	-2	$-\frac{1}{2} \pm i \frac{\sqrt{3}}{2}$
$\frac{s+1}{s^2 - 1}$	None	+1	None	0

When there are no open-loop zeros, there are no closed-loop zeros. Because the numerator and denominator of $G_3(s)$ are not coprime, we had to first remove their common factor of $(s + 1)$ before completing the table. ■

10.6 Cascade Control Architecture

Figure 10.6 is a refinement of the unity feedback system, where the forward path is written as a **cascade** connection of a **compensator** $C(s)$ and **plant** $P(s)$. Recall that “plant” is the generic term for the system that is to be controlled. The input to the plant is denoted $U(s)$. The plant input is the means by which the plant output is changed. For the Segway, the plant input was the torque at the wheels provided by the motors, while for the RLC circuit, it was the applied voltage.

The compensator is the algorithm that determines how to vary the input to the plant to regulate the plant output. In the cascade control configuration of Fig. 10.6, the compensator acts on $E(s) := R(s) - Y(s)$, the **error** between the reference command $R(s)$ and the measured output $Y(s)$, to produce the control input by

$$U(s) = C(s)E(s). \quad (10.47)$$

When the error is zero, the plant output is exactly following or **tracking** the reference command. One way to think of the compensator is that it attempts to adjust the plant input so that the plant output tracks the reference command. Equivalently, the compensator seeks to drive the error between the reference and output signals to zero.

The closed-loop transfer functions, (10.43) and (10.45) in Chapter 10.5.1 easily updated to this new situation by identifying $G(s) := C(s)P(s)$. Doing so yields,

$$\frac{Y(s)}{R(s)} = \frac{C(s)P(s)}{1 + C(s)P(s)} \quad (10.48)$$

and

$$\frac{E(s)}{R(s)} = \frac{1}{1 + C(s)P(s)}. \quad (10.49)$$

Because $U(s) = C(s)E(s)$, it follows that

$$\frac{U(s)}{R(s)} = \frac{C(s)}{1 + C(s)P(s)}. \quad (10.50)$$

All three of these transfer functions become important in most control system designs.

10.6.1 Two Common Compensators: Proportional (P) and Proportional-Derivative (PD)

We introduce here the two most commonly used compensators. A more complete list is available in Appendix ??.

Proportional compensator

$$C(s) = K_P, \quad (10.51)$$

where K_P is called a **proportional gain**. With a proportional compensator, the control signal is given by

$$u(t) = K_P e(t) = K_P(r(t) - y(t)).$$

Adjusting K_P allows the control signal to react more or less “vigorously” to the error between the commanded value of the output (aka, reference signal), $r(t)$, and the measured value of the output, $y(t)$. It adds neither zeros nor poles to the forward path of the feedback loop. For a first-order plant, $P(s)$, gain adjustment alone can achieve stability, as we illustrate next.

Example 10.34. (Proportional feedback applied to an unstable first-order plant) Suppose the plant has transfer function $P(s) = \frac{1}{s-1}$ and the compensator is a proportional feedback $C(s) = K_P$, where K_P is a gain to be chosen. Determine the closed-loop transfer function $\frac{Y(s)}{R(s)}$ and evaluate its BIBO stability as a function of K_P .

Solution:

$$\begin{aligned} \frac{Y(s)}{R(s)} &= \frac{C(s)P(s)}{1 + C(s)P(s)} \\ &= \frac{K_P \frac{1}{s-1}}{1 + K_P \frac{1}{s-1}} \quad (\text{multiply top and bottom by } (s-1)) \\ &= \frac{K_P}{s-1+K_P}. \end{aligned}$$

The pole of the transfer function is the root of the denominator polynomial, $s - 1 + K_P = 0$, yielding a single pole $1 - K_P$. The closed-loop system is, therefore, BIBO stable if, and only if, $K_P > 1$. ■

Feedback control design would be way too easy if we could stabilize all plants with proportional feedback! The following result will make it easier to show a limitation of proportional feedback control.

Fact 10.35: Simplified Stability Test

For real numbers a_1 and a_0 , the roots of the quadratic equation

$$s^2 + a_1 s + a_0 = 0 \quad (10.52)$$

have negative real parts if, and only if, $a_1 > 0$ and $a_0 > 0$.

Note: To apply the above result to $\bar{a}_2 s^2 + \bar{a}_1 s + \bar{a}_0 = 0$, with $\bar{a}_2 \neq 0$, you simply divide through by \bar{a}_2 , yielding $a_1 = \frac{\bar{a}_1}{\bar{a}_2}$ and $a_0 = \frac{\bar{a}_0}{\bar{a}_2}$.

For a cubic equation, the result is, for real numbers a_2 , a_1 , and a_0 , the roots of the equation

$$s^3 + a_2 s^2 + a_1 s + a_0 = 0$$

have negative real parts if, and only if, $a_2 > 0$, $a_0 > 0$, and $a_1 - \frac{a_0}{a_2} > 0$. Similar results can be given for polynomials of arbitrary degree. The computations are typically organized in a table called the **Routh array**, which is typically covered in Michigan's EECS 460 Control Systems Analysis and Design.

Example 10.36. (Proportional feedback applied to an unstable second-order plant) Suppose the plant is instead $P(s) = \frac{1}{s(s-2)}$, with the previous proportional feedback, $C(s) = K_P$. Determine the transfer function $\frac{Y(s)}{R(s)}$ and evaluate its stability as a function of K_P .

Solution: Evaluating the closed-loop transfer function gives

$$\frac{Y(s)}{R(s)} = \frac{K_P}{s^2 - 2s + K_P}.$$

The poles of the transfer function are the roots of the polynomial $s^2 - 2s + K_P = 0$. Applying the result of Fact 10.35, it is seen that the coefficient $a_1 = -2 < 0$ is independent of the value of the gain K_P , and hence at least one root will have a non-negative real part no matter the value of K_P . The closed-loop system is, therefore, unstable for all values of the gain K_P , and **something more sophisticated than proportional feedback is required to stabilize the closed-loop system**. ■

PD or Proportional-Derivative compensator

$$C(s) = K_P + K_D s = K_D \left(s + \frac{K_P}{K_D} \right) \quad (10.53)$$

has two gains, the proportional gain, K_P , which acts as before, and the **derivative gain**, K_D . With a proportional-derivative (PD) compensator, the control signal is given by

$$u(t) = K_P e(t) + K_D \dot{e}(t) = K_P(r(t) - y(t)) + K_D(\dot{r}(t) - \dot{y}(t)).$$

The derivative term allows the controller to react with “anticipation” to changes in the error signal. Whoa! Does the controller become a psychic? Not exactly. If we approximate the derivative as

$$\dot{e}(t) \approx \frac{e(t + \delta t) - e(t)}{\delta t},$$

the control signal approximately depends on the error term at the future time, $t + \delta t$, giving it a head start on making a correction to the closed-loop system. Another way to look at it is that the linear approximation of the error signal at a point t_0 is

$$e(t) \approx e(t_0) + \underbrace{\frac{de(t_0)}{dt}}_{\dot{e}(t_0)}(t - t_0),$$

showing again that the derivative of e allows the prediction of future values of the error signal over small intervals of time. This discussion is meant to provide intuition, so if it seems hand-wavy or confusing to you, just ignore it. The Laplace transform will soon clear up lingering doubts.

A PD compensator adds a zero in the forward path at $-\frac{K_P}{K_D}$, but has no poles. The derivative action of a PD compensator will prove useful for stabilizing a system and for reducing oscillations.

Example 10.37. (PD compensation applied to an unstable second-order plant) Suppose the plant is $P(s) = \frac{1}{s(s-2)}$ and $C(s) = K_P + K_D s$. Determine the transfer function $\frac{Y(s)}{R(s)}$ and evaluate its stability as a function of the controller gains.

Solution: Evaluating the closed-loop transfer function gives

$$\frac{Y(s)}{R(s)} = \frac{K_P + K_D s}{(s^2 - 2s) + (K_P + K_D s)} = \frac{K_P + K_D s}{s^2 + (K_D - 2)s + K_P}.$$

The poles of the transfer function are the roots of the polynomial $s^2 + (K_D - 2)s + K_P = 0$. After identifying the coefficients $a_1 = (K_D - 2)$ and $a_0 = K_P$, Fact 10.35 yields that the closed-loop system's poles have negative real parts if, and only if, $K_P > 0$ and $K_D > 2$. ■

Is stability all we care about in a feedback loop? Hardly. We'll next look at the steady-state error, that is, the long-term "accuracy" of the control system.

10.6.2 Steady-State Error

The **steady-state error**, or e_{SS} , is the asymptotic² value of the error signal defined in Figures 10.5 and 10.6. The steady-state error depends on the reference command and the forward-path transfer function. The expression for the error signal was derived in (10.44) and (10.49). When e_{SS} is small, the control system is accurately tracking the reference command as the system reaches a steady state. What happens between the closed-loop system being "excited" by an input signal and its reaching a steady state is called its transient response; this will be the subject of a separate section.

While the term "steady-state" can be defined for a range of reference signals, we will focus on the most common one, the unit-step function, that is, $R(s) = \frac{1}{s}$. Appendix ?? also looks at steady-state error for ramp inputs.

Proposition 10.38: Steady-State Error for Step Inputs Applied to BIBO Stable Closed-Loop Systems

Assume the forward path transfer function is $G(s) = \frac{N(s)}{D(s)}$ and the closed-loop system is BIBO stable. Then the steady-state error for a unit-step command is well-defined and is given by

$$e_{SS} = \lim_{s \rightarrow 0} \frac{D(s)}{D(s) + N(s)}. \quad (10.54)$$

Moreover, if $G(s)$ does not have a pole at the origin, then

$$e_{SS} = \frac{D(0)}{D(0) + N(0)} \neq 0,$$

while if $G(s)$ has one or more poles at the origin, then $D(0) = 0$ and the steady-state error is

$$e_{SS} = 0.$$

Note: To be clear, the result also holds for $G(s) = C(s)P(s)$ arising from a Cascade Controller Architecture.

Note: What happens to e_{SS} if the reference is $r(t) = r_0 \cdot u_{\text{stp}}(t)$, a scaled unit-step? Because $E(s) = \frac{1}{1+G(s)} R(s)$, it follows that if the reference is scaled by r_0 , then so is the error. Hence, e_{SS} gives the *relative error*, or, said another way, the *percent error* in the output is $100 e_{SS}\%$.

e_{SS}	% Error
0.10	10 %
0.25	25 %
1.00	100 %

Example 10.39. (Steady-state error) Determine the steady-state error for a unit-step input when a unity feedback system is specified by

²That is, its value as $t \rightarrow \infty$.

- (a) $C_1(s) = 10$ and $P_1(s) = \frac{4}{s+4}$.
 (b) $C_2(s) = 10$ and $P_2(s) = \frac{4}{s(s+4)}$.
 (c) $C_3(s) = (10 + 2s)$ and $P_3(s) = \frac{1}{s(s-4)}$.

Solution: The step responses are shown in Fig. 10.7. The calculations for the time-domain expressions are given in Example 10.40 below. In engineering practice, if time-domain signals are required, they are usually computed numerically by a computer-aided control design program rather than by tedious hand calculations, but in other courses, you may need to do them by hand.

- (a) $C_1(s) = 10$ and $P_1(s) = \frac{4}{s+4}$ **Ans.** $e_{SS} \approx 9.1\%$.

The forward-path transfer function is

$$C_1(s)P_1(s) = \frac{40}{s+4} =: \frac{N_1(s)}{D_1(s)}.$$

The pole of the closed-loop system satisfies

$$D_1(s) + N_1(s) = (s+4) + 40 = 0 \implies s = -44,$$

and thus, the closed-loop system is BIBO stable. For a step input, the steady-state error is

$$e_{SS} = \frac{D_1(0)}{D_1(0) + N_1(0)} = \frac{4}{40+4} = \frac{1}{11} \approx 0.091.$$

- (b) $C_2(s) = 10$ and $P_2(s) = \frac{4}{s(s+4)}$. **Ans.** $e_{SS} = 0.0$.

The forward-path transfer function is

$$C_2(s)P_2(s) = \frac{40}{s(s+4)} =: \frac{N_2(s)}{D_2(s)}.$$

The pole of the closed-loop system satisfies

$$D_2(s) + N_2(s) = s(s+4) + 40 = s^2 + 4s + 40 = 0,$$

and thus by Fact 10.35, the closed-loop system is BIBO stable. For a step input, the steady-state error is

$$e_{SS} = \frac{D_2(0)}{D_2(0) + N_2(0)} = \frac{0}{0+40} = 0.0.$$

We emphasize that because the denominator of the forward-path transfer function has a pole at the origin, $D_2(0) = 0$. Recall that because s is a differentiator, $\frac{1}{s}$ is an integrator. You can read about **integral control** in Appendix ??.

- (c) $C_3(s) = (10 + 2s)$ and $P_3(s) = \frac{1}{s(s-4)}$. **Ans.** e_{SS} is undefined because the closed-loop is unstable. In practice, you would see the error wildly oscillating with increasing amplitude because the closed-loop poles are complex, with a positive real part (as you can check with the quadratic formula).

The forward-path transfer function is

$$C_3(s)P_3(s) = \frac{(10 + 2s)}{s(s-4)} =: \frac{N_3(s)}{D_3(s)}.$$

The pole of the closed-loop system satisfies

$$D_3(s) + N_3(s) = s(s-4) + (10 + 2s) = s^2 - 2s + 10 = 0,$$

and thus by Fact 10.35, the closed-loop system is unstable. Therefore, the steady-state error is undefined.

We show how the time domain signals in Fig. 10.7 are computed for part (a). The others are similar.

```

1 using ControlSystems, Plots, LaTeXStrings
2 # ControlSystems is slow to load
3
4 # Specify forward path transfer functions
5 C = 10.0
6 P = tf([4], [1.0, 4.0])
7
8 # build closed-loop tf
9 sys = feedback(C*P, 1)
10
11 # Compute and Plot the Step Response
12 t = 0:0.01:1.0
13
14 # Time trajectory
15 y, t = step(sys, t); y = y'
16
17 # Plot
18 p1 = plot(t, y, guidefont = 15, lw=3, color=:blue, label=L"$y(t)$", legendfontsize=15,
19 xlabel="t (s)", ylabel="Step Response", legend=:right)
20 Ref = ustep.(t)
21 Error = Ref - y
22 p1 = plot!(p1, t, Ref, lw=3, color=:black, legend=:right, label=L"$u(t)$")
23 p1 = plot!(p1, t, Error, lw=3, color=:red, legend=:right, label=L"$e(t)$")
24
25 png(p1, "C1P1StepResponse")
26 display(p1)

```

Output See Fig. 10.7.

Example 10.40. (Optional Read:) Use code to determine time-domain expressions for $e(t)$ and then take limits to confirm the results in Example 10.39.

Solutions:

In each case, we use code to compute the inverse Laplace transform of $E(s) := \frac{1}{1+C(s)P(s)} \cdot \frac{1}{s}$, where $\frac{1}{s} = \mathcal{L}\{u_{\text{stp}}(t)\}$.

(a) $C_1(s) = 10$ and $P_1(s) = \frac{4}{s+4}$ yield $e_1(t) = \frac{1}{11}u_{\text{stp}}(t) + \frac{10}{11}e^{-4t}u_{\text{stp}}(t)$ and hence $\lim_{t \rightarrow \infty} e_1(t) = \frac{1}{11}$.

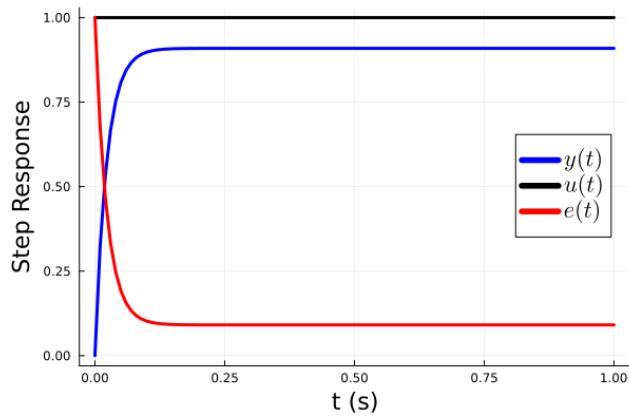
(b) $C_2(s) = 10$ and $P_2(s) = \frac{4}{s(s+4)}$ yield $e_2(t) = \frac{1}{3}e^{-2t} \sin(6t)u_{\text{stp}}(t) + e^{-2t} \cos(6t)u_{\text{stp}}(t)$ and hence $\lim_{t \rightarrow \infty} e_2(t) = 0$.

(c) $C_3(s) = (10 + 2s)$ and $P_3(s) = \frac{1}{s(s-4)}$ yield $e_3(t) = e^t \cos(3t)u_{\text{stp}}(t) - e^t \sin(3t)u_{\text{stp}}(t)$ and hence $\lim_{t \rightarrow \infty} e_3(t)$ is undefined. Indeed, $e_3(t)$ grows unbounded while oscillating between positive and negative values, never making up its mind between $+\infty$ and $-\infty$!

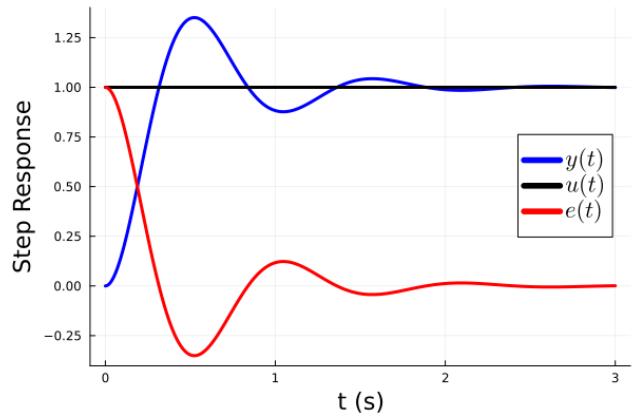
```

1 using PyCall
2
3 # Import Python's SymPy
4 sympy = pyimport("sympy")
5
6 # Define the symbols
7 s, t = sympy.symbols("s t")
8

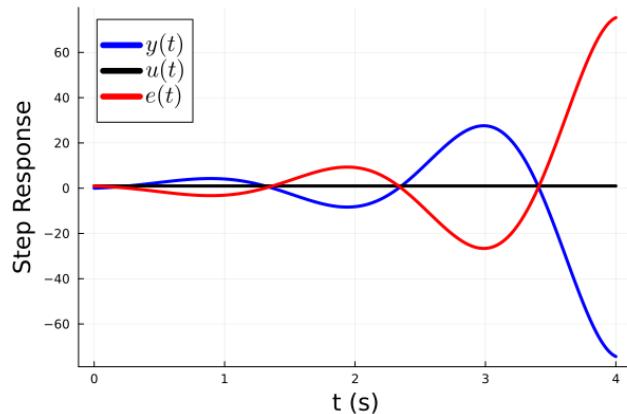
```



(a)



(b)



(c)

Figure 10.7: The step responses for Examples 10.39 and 10.40, with the error in red. In (a) and (b), a steady-state error is achieved, while in (c), the error does not have a well-defined limit due to the closed-loop system being unstable.

```

9 Ustep = 1/s
10
11 # Define the expression E1(s)
12 C1 = 10 + 0*s
13 P1 = 4/(s+4)
14 E1 = (1/(1 + C1*P1))*Ustep
15
16 # Compute the inverse Laplace transform
17 e1 = sympy.inverse_laplace_transform(E1, s, t)
18
19 println("Part 1")
20 println(e1)
21 println(" ")
22
23 # Define the expression E2(s)
24 C2 = 10 + 0*s
25 P2 = 4/(s*(s+4))
26 E2 = (1/(1 + C2*P2))*Ustep
27
28 # Compute the inverse Laplace transform
29 e2 = sympy.inverse_laplace_transform(E2, s, t)
30
31 println("Part 2")

```

```

32 println(e2)
33 println(" ")
34
35 # Define the expression E3(s)
36 C3 = 10 + 2*s
37 P3 = 1/(s*(s-4))
38 E3 = (1/(1 + C3*P3))*Ustep
39
40 # Compute the inverse Laplace transform
41 e3 = sympy.inverse_laplace_transform(E3, s, t)
42
43 println("Part 3")
44 println(e3)

```

Output Recall that Heaviside(t) is the unit-step function, $u_{\text{stp}}(t)$. Also note that $\lim_{t \rightarrow \infty} u_{\text{stp}}(t) = 1.0$

Part 1
 $\text{Heaviside}(t)/11 + 10 * \exp(-44*t) * \text{Heaviside}(t)/11$

Part 2
 $(\exp(-2*t) * \sin(6*t)/3 + \exp(-2*t) * \cos(6*t)) * \text{Heaviside}(t)$

Part 3
 $(-\exp(t) * \sin(3*t) + \exp(t) * \cos(3*t)) * \text{Heaviside}(t)$

10.7 Transient Response of First- and Second-order Systems

Thus far, we have learned how to compute closed-loop transfer functions, evaluate their stability, and analyze steady-state tracking performance for step inputs. This section looks at the transient response of a system to an abrupt change in input, that is, the **step response**. The primary objective is to relate the qualitative features of the step response to the poles and zeros of the transfer function.

10.7.1 Common Performance Specifications

In this section, whether the system being considered is open-loop or closed-loop is irrelevant. Let the Laplace-transform representation of the system be

$$Y(s) = G(s) \frac{1}{s},$$

where $\frac{1}{s}$ is the Laplace transform of a unit-step input, $Y(s)$ is the output, and $G(s)$ is a BIBO stable transfer function. Let $y(t)$ be the resulting time response of the output, such as the one shown in Fig. 10.8. Section 10.6.2, because the transfer function is BIBO stable, $y(t)$ asymptotically converges to a steady-state value, denoted here by y_∞ .

The key qualitative features of Fig. 10.8 are listed below:

1. **Peak value:** y_{\max} is the peak value of the output.

In the case that the output monotonically approaches its steady-state value y_∞ , as in Fig. 10.9 for example, then y_{\max} is defined to be y_∞ .

2. **Percent overshoot:** the difference $y_{\max} - y_\infty$ is the amount the output exceeds its steady-state value, which could be zero,

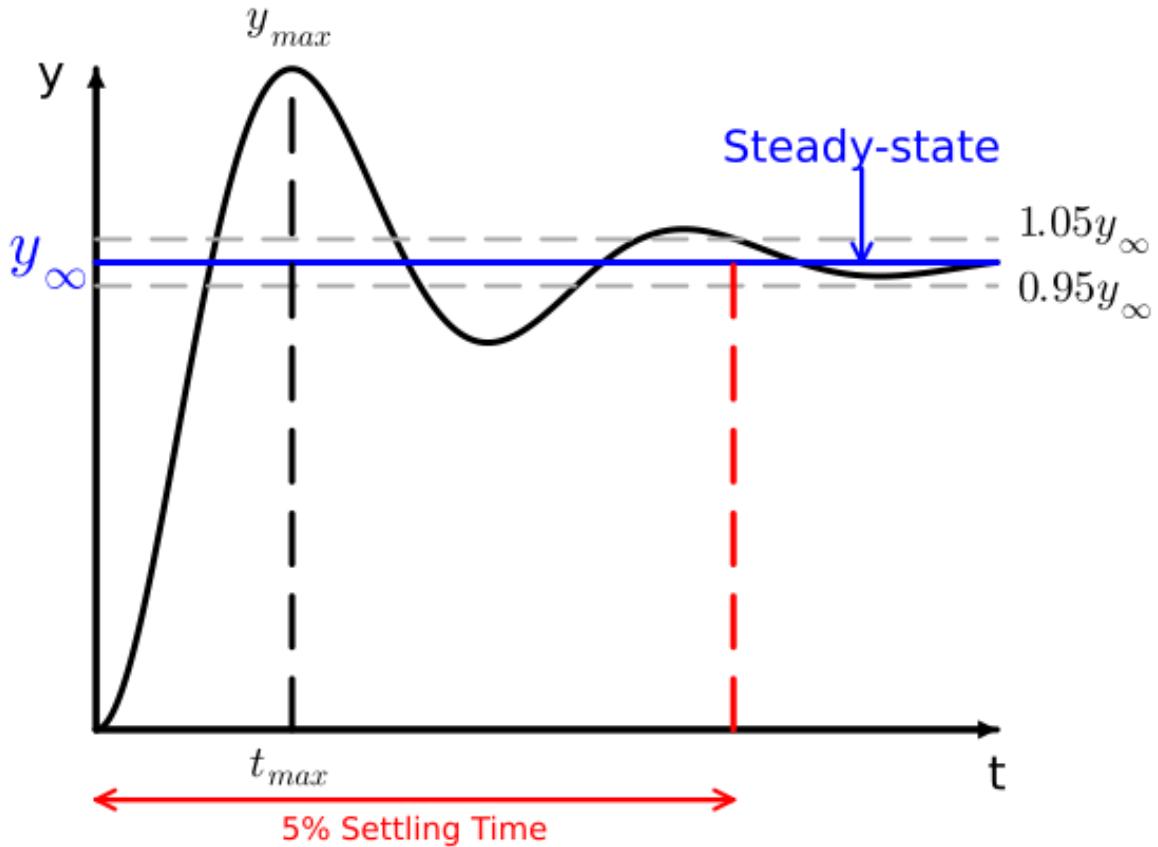


Figure 10.8: Typical step response of an underdamped system. The 5% settling time is the time it takes for the system to enter and then remain within the interval $[0.95y_\infty, 1.05y_\infty]$.

as indicated above. The percent overshoot is given by

$$\text{Percent overshoot} = \frac{y_{max} - y_\infty}{y_\infty} \times 100\%$$

3. **5% settling time** = $\min\{T \geq 0 : |y(t) - y_\infty| \leq 0.05y_\infty, \forall t \geq T\}$. In other words, the 5% settling time T_s is the smallest $T > 0$ such that

$$0.95y_\infty \leq y(t) \leq 1.05y_\infty \quad (10.55)$$

for all $t \geq T$. Though not used here, the 2% settling time is fairly common as well.

4. **Peak Time:** Time to reach y_{max} ; it could be infinite.

5. **(Not shown) Rise time:** time to go from 10% to 90% of final value. Most commonly used for first-order systems or for over-damped second-order systems.

We will see in Section 10.8 that feedback design requires good rules of thumb for relating the overall system's performance to the zeros and poles of the transfer function. These rules of thumb are developed next.

[JWG: Use Figure 10.8 as a template for transient response plots.]

10.7.2 First-order System without a Zero

Consider a first-order system with unity DC gain, meaning that $G(0) = 1.0$, namely

$$G(s) = \frac{1}{\tau s + 1}. \quad (10.56)$$

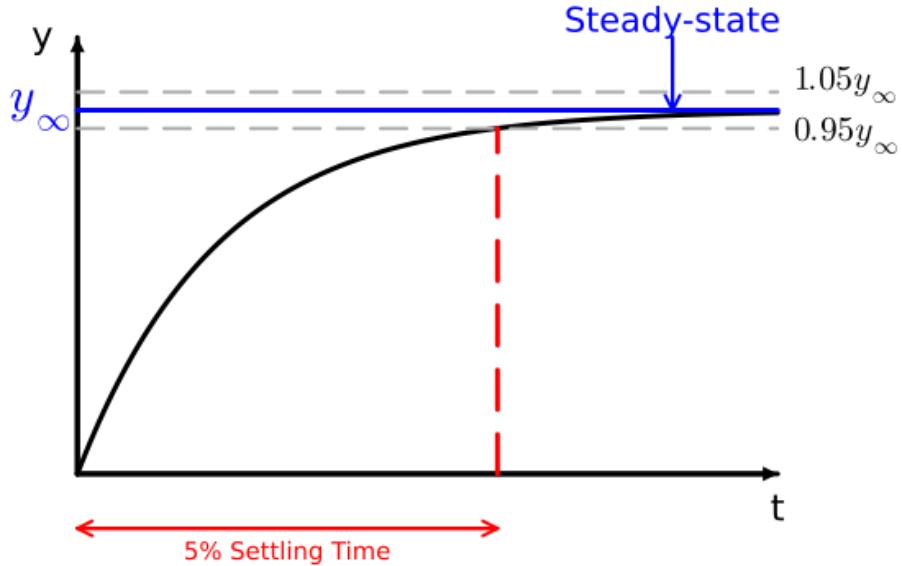


Figure 10.9: (**First-order Systems:**) The step response of a first-order system $\frac{k_0}{\tau s + 1}$ has $10\% - 90\%$ rise time approximately equal to 2τ , no overshoot, and its 5% settling time is approximately 3τ . The peak value is the same as y_∞ , meaning it is not achieved for any finite value of t .

When $\tau > 0$, it is called the **time constant** of the system. The pole is then given by $-1/\tau$. The step response can be computed to be

$$y(t) = \left(1 - e^{-\frac{t}{\tau}}\right) u_{\text{stp}}(t),$$

which is shown in Fig. 10.9. The time constant is the time it takes the step response to go from zero to 63.21% of its final value. The step response has no overshoot, and thus both the peak time and the 0% to 100% rise time are infinite.

Because it is easier to remember round numbers, it is common to estimate the time constant by the time it takes the step response to go from zero to 60% of its final value. In the spirit of the zero to 60% approximation for the time constant, the $10\% - 90\%$ rise time is approximately equal to 2τ and the 5% settling time is approximately equal to 3τ . These are summarized below for later use:

$0\% \text{ to } 60\% \text{ of final value} = \tau,$ $10\% \text{ to } 90\% \text{ rise time} = 2\tau,$ $5\% \text{ settling time} = 3\tau.$	(10.57)
---	---------

If the transfer function is instead $K \frac{1}{\tau s + 1}$, so that the DC gain is K , the above approximations remain valid because they depend only on percentage changes in the output, and not on absolute values. This is an important advantage in expressing quantities in relative terms, or percentages.

10.7.3 Second-order System without a Zero

Consider a second-order system with unity **DC gain**, that is, $G(0) = 1.0$,

$$G(s) = \frac{a_0}{s^2 + a_1 s + a_0}, \quad (10.58)$$

where $a_0 > 0$ and $a_1 > 0$ are assumed for BIBO stability. The step response of the system becomes straightforward to understand when the model is reparameterized in terms of the **damping ratio** ζ and **undamped natural frequency** ω_n , defined as

$$\omega_n^2 = a_0 \quad (10.59)$$

$$2\zeta\omega_n = a_1. \quad (10.60)$$

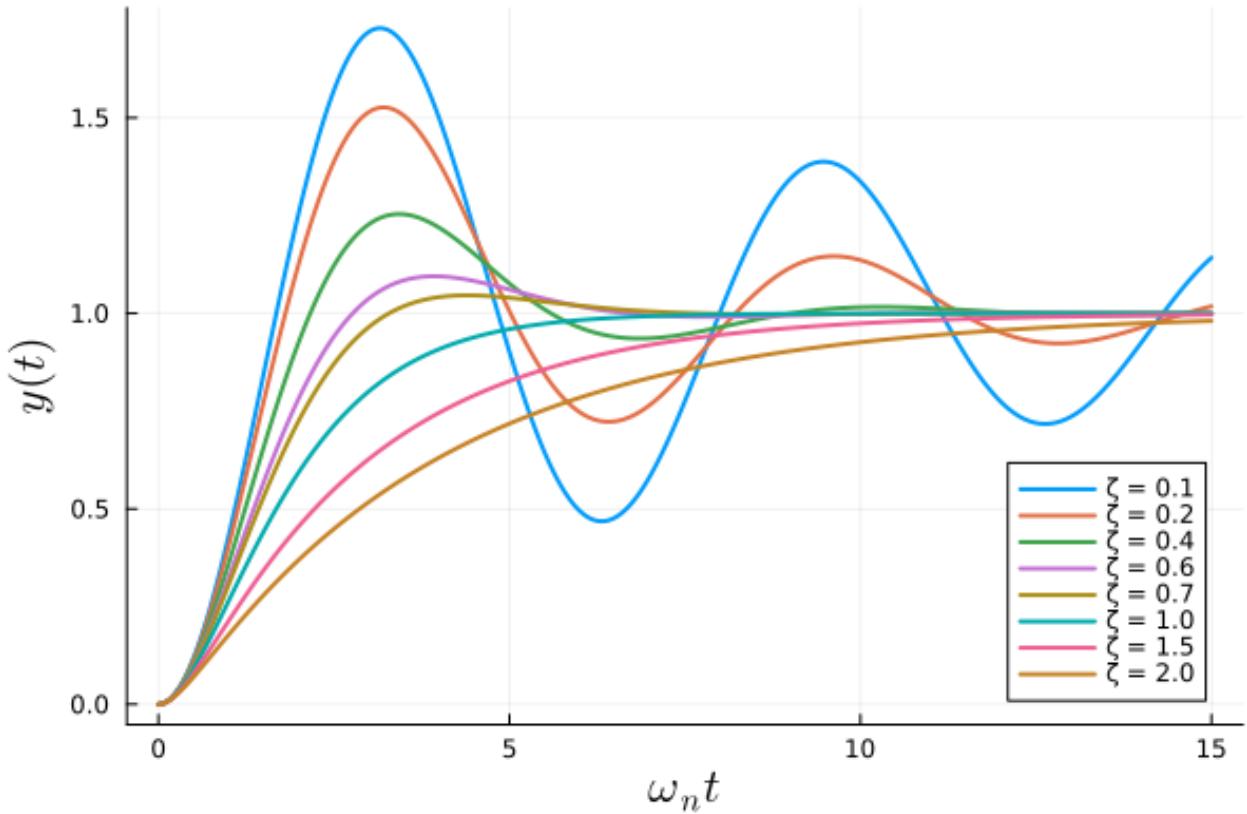


Figure 10.10: Step response of a standard second-order system, $G(s) = \frac{\omega_n^2}{s^2 + 2\zeta\omega_n s + \omega_n^2}$. There is overshoot for $0 < \zeta < 1$, while there is no overshoot for $\zeta \geq 1$. There is a trade-off between fast response and overshoot. **Note that the x -axis is scaled and equal to $\omega_n t$.**

This gives,

$$G(s) = \frac{\omega_n^2}{s^2 + 2\zeta\omega_n s + \omega_n^2}, \quad (10.61)$$

where

$$\omega_n = \sqrt{a_0} \quad (10.62)$$

$$\zeta = \frac{a_1}{2\omega_n} = \frac{a_1}{2\sqrt{a_0}}. \quad (10.63)$$

The step response for a range of damping ratios is shown in Fig. 10.10. When the damping ratio ζ is between zero and one, there is overshoot and the output oscillates before settling to its final value. When the damping ratio is one or greater, there is no overshoot. The speed with which the step response converges to its steady-state value is determined by the undamped natural frequency. When ω_n is small, the response is sluggish, meaning the settling time and rise time are long, and when ω_n is large, the response is rapid. For example, if $\omega_n = 1000$, the time scale in Fig. 10.10 becomes milliseconds, whereas, if $\omega_n = 1/60$, the time scale is minutes.

The system in (10.61) is said to be in **standard form**. Note that if the model has instead the form

$$\tilde{G}(s) = \frac{k}{s^2 + a_1 s + a_0}, \quad (10.64)$$

then it can be re-expressed as

$$\tilde{G}(s) = \frac{k}{a_0} \frac{a_0}{s^2 + a_1 s + a_0} = K \frac{\omega_n^2}{s^2 + 2\zeta\omega_n s + \omega_n^2}, \quad (10.65)$$

where $K = \frac{k}{a_0}$. The gain K introduces a change of scale in the system's output, but it does not change the percent overshoot, settling time, rise time, etc., because these are all relative quantities.

Poles

Applying the quadratic formula, the poles of the transfer function are computed to be

$$s_{1,2} = -\zeta\omega_n \pm i\omega_n\sqrt{\zeta^2 - 1}.$$

Depending on the sign of $\zeta^2 - 1$, three cases are possible:

- **underdamped** $0 < \zeta < 1$: a pair of complex conjugate poles;
- **critically damped** $\zeta = 1$: two repeated real poles; and
- **overdamped** $\zeta > 1$: two distinct real poles.

Step Response

The step responses for the three different types of standard second-order systems are described in detail. The derivations are available in Appendix ??.

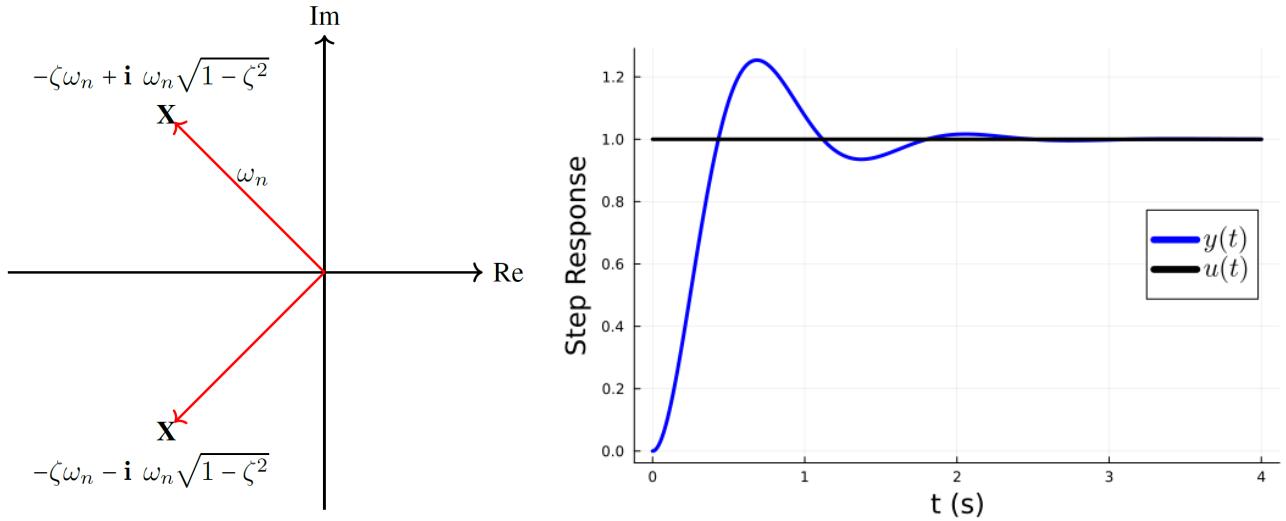


Figure 10.11: The complex poles of an underdamped system ($0 < \zeta < 1$) give rise to overshoot and oscillations in the step response.

Underdamped ($0 < \zeta < 1$): The poles are complex and equal to

$$s_{1,2} = -\zeta\omega_n \pm i\omega_n\sqrt{1 - \zeta^2}. \quad (10.66)$$

Using the method of partial fraction expansion, the step response is computed to be,

$$y(t) = 1 - \frac{1}{\sqrt{1 - \zeta^2}} e^{-\zeta\omega_n t} \sin(\omega_n\sqrt{1 - \zeta^2}t + \theta) u_{\text{stp}}(t), \quad (10.67)$$

where, from Fig. 10.11,

$$\cos(\theta) = \frac{\zeta\omega_n}{\omega_n} = \zeta.$$

The decaying, sinusoidal nature of the output is clearly seen.

The peak value is determined by computing the time derivative of the output, setting it to zero, and evaluating the relative maxima

and minima. Doing so gives

$$y_{\max} = 1 + e^{\frac{-\zeta\pi}{\sqrt{1-\zeta^2}}}, \quad (10.68)$$

$$\text{percent overshoot} = 100e^{\frac{-\zeta\pi}{\sqrt{1-\zeta^2}}} \% \quad (10.69)$$

$$\text{peak time } t_{\max} = \frac{\pi}{\omega_n \sqrt{1 - \zeta^2}}. \quad (10.70)$$

The 0 – 100% rise time is given by

$$t_r = \frac{\pi - \theta}{\omega_n \sqrt{1 - \zeta^2}}. \quad (10.71)$$

Though the exact calculation of the settling time is difficult, for $0.3 \leq \zeta \leq 0.9$, the following approximation is usually sufficient,

$$5\% \text{ settling time} \approx \frac{3}{\zeta \omega_n}. \quad (10.72)$$

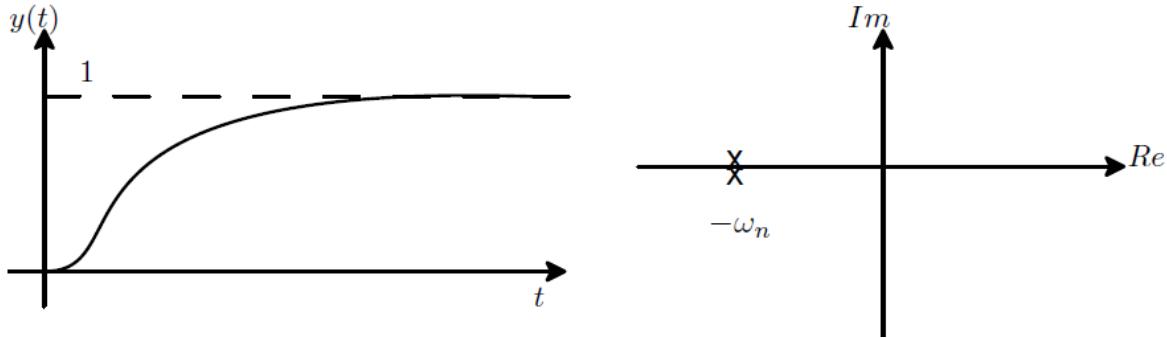


Figure 10.12: The step response of a critically damped ($\zeta = 1$), second-order system with no zeros is monotonic. The repeated real poles at $-\omega_n$ are denoted with a double \times

Critically damped ($\zeta = 1$): The transfer function has two, repeated real poles and can be written as

$$G(s) = \frac{\omega_n^2}{(s + \omega_n)^2}. \quad (10.73)$$

Following the method of partial fraction expansion, the step response is computed to be

$$y(t) = (1 - e^{-\omega_n t} - \omega_n t e^{-\omega_n t}) u_{\text{stp}}(t). \quad (10.74)$$

By computing the time derivative of the output, it can be shown that for all $t \geq 0$, $\dot{y}(t) > 0$, and therefore, there are neither oscillations nor overshoot; moreover, the peak time is infinite. An approximation of the settling time is given by

$$5\% \text{ settling time} \approx \frac{4.8}{\omega_n}. \quad (10.75)$$

Overdamped ($\zeta > 1$): The transfer function has two distinct real poles and can be written as

$$G(s) = \frac{p_1}{s + p_1} \frac{p_2}{s + p_2} = \frac{1}{\frac{1}{p_1}s + 1} \frac{1}{\frac{1}{p_2}s + 1},$$

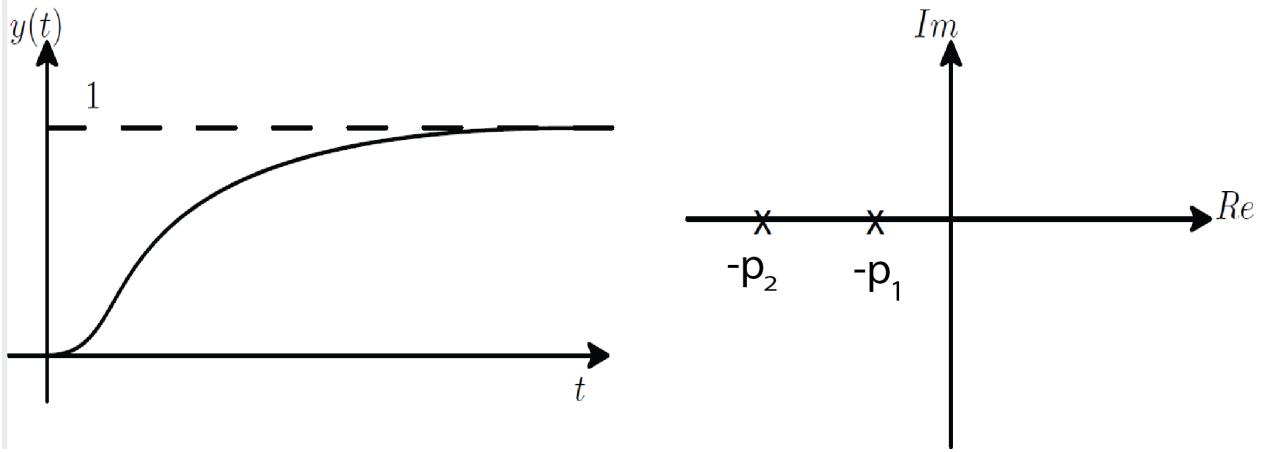


Figure 10.13: Step response of an overdamped, second-order system with no zeros.

where

$$\begin{aligned} p_1 &= \omega_n \left[\zeta - \sqrt{\zeta^2 - 1} \right] \\ p_2 &= \omega_n \left[\zeta + \sqrt{\zeta^2 - 1} \right]. \end{aligned} \quad (10.76)$$

Following the method of partial fraction expansion, the step response is computed to be

$$y(t) = \left(1 - \frac{p_2}{p_2 - p_1} e^{-p_1 t} - \frac{p_1}{p_1 - p_2} e^{-p_2 t} \right) u_{\text{stp}}(t). \quad (10.77)$$

By computing the time derivative of the output, it can also be shown that for all $t \geq 0$, $\dot{y}(t) > 0$, which implies that the output monotonically approaches its final value, and thus there is no overshoot and the peak time is infinite, as shown in Fig. 10.13.

Because there is no overshoot, the primary qualitative feature of the output is its settling time.

- For $1 < \zeta < 1.5$, the poles are close to one another and the terms in (10.77) decay at similar rates. An exact value for the settling time is difficult to compute, but a good approximation can be shown to be

$$5\% \text{ settling time} \approx \frac{6.4}{\omega_n} (\zeta - 1) + \frac{4.8}{\omega_n}. \quad (10.78)$$

- For $\zeta \geq 2$, the poles are far apart because $p_2 \geq 10p_1$. In this case, the last term in (10.77) decays much more rapidly than the second term. Moreover,

$$\frac{p_2}{p_2 - p_1} \approx 1,$$

and thus,

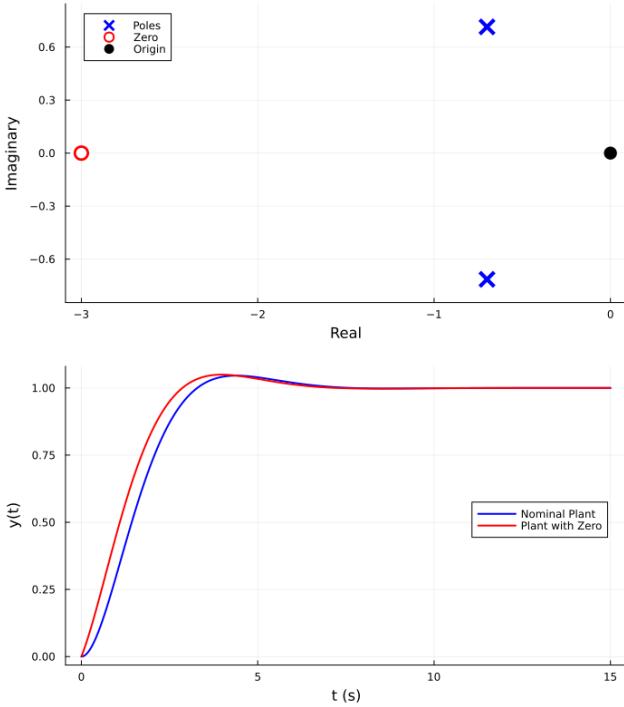
$$y(t) \approx (1 - e^{-p_1 t}) u_{\text{stp}}(t), \quad (10.79)$$

which is the step response of a first-order system with time constant $\tau = 1/p_1$, that is,

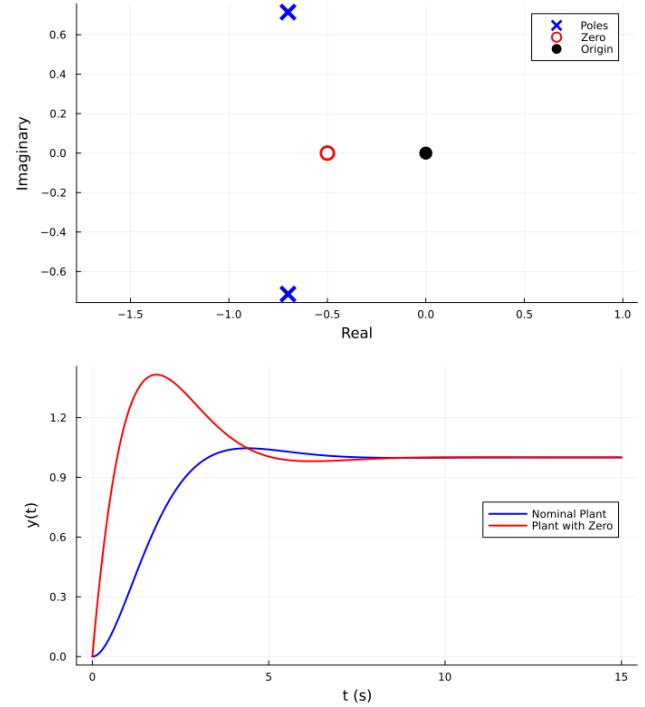
$$G(s) \approx \frac{1}{\frac{1}{p_1} s + 1}. \quad (10.80)$$

Consequently, the 5% settling time can be estimated from the corresponding formula for a first-order system, namely

$$5\% \text{ settling time} \approx \frac{3}{p_1}.$$



(a)



(b)

Figure 10.14: (**Stable Zeros Increase Overshoot:**) In (a), the zero is far to the left of the poles, and its effect if small, while in (b), the zero is near the poles and greatly increases the overshoot.

10.7.4 Effects of Zeros

Zeros in a transfer function can profoundly alter the transient response. The magnitude of their effect depends on where they are located with respect to the poles. This will be illustrated by considering a single zero in a second-order system,

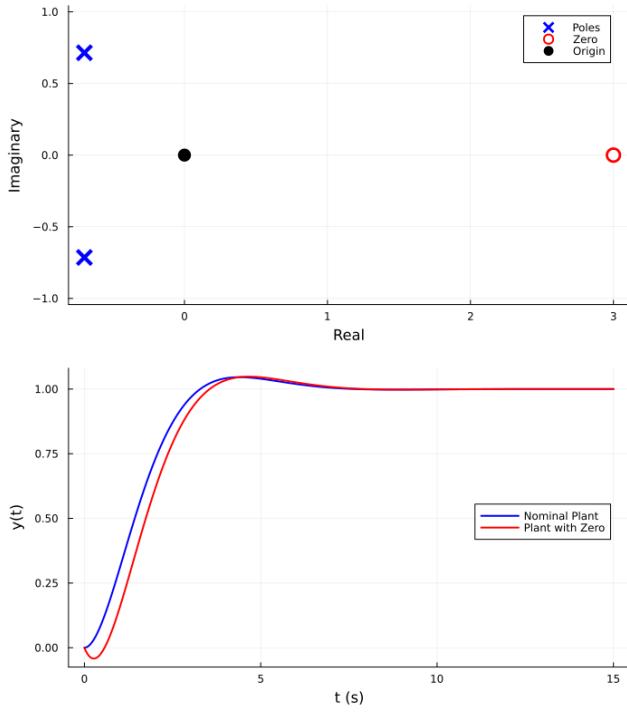
$$\bar{G}(s) = \frac{\left(\frac{1}{z}s + 1\right)\omega_n^2}{s^2 + 2\zeta\omega_n s + \omega_n^2}, \quad (10.81)$$

which has a zero at $-z$. For the purposes of our analysis, the zero has been introduced in such a way that it does not change the DC gain of the transfer function. While this excludes $z = 0$, it is fine from a practical point of view. The value of z can be either positive or negative, placing the zero in the left-half or right-half plane, respectively.

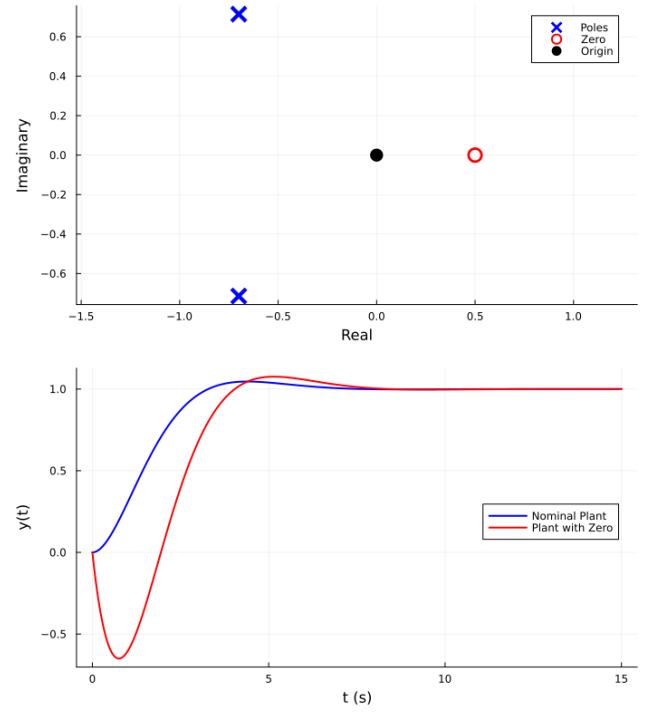
As shown in Figures 10.14 and 10.15, the effects of a zero can be quite startling. It depends on the location of the zero.

Underdamped $0 < \zeta \leq 1$

- $z < 0$, then the zero is in the right-half plane, which is officially called a **non-minimum phase zero**, though we will call it informally an “**unstable zero**”. A very special phenomenon happens: the step response exhibits **undershooting**, meaning that the output initially moves in a direction opposite to the commanded value. Zeros in the left-half plane do not cause undershooting.
- $|z| > 10\zeta\omega_n$, the zero can be neglected in the sense that the step response is essentially indistinguishable from that of a system without a zero.
- $|z| \approx 5\zeta\omega_n$ the zero can be neglected because the step response exhibits only a slight increase in overshoot and settling time from a system without a zero.
- $|z| \approx 2\zeta\omega_n$, the zero has a significant effect on the step response step response. For $z > 0$, the 0-100% rise time is reduced by



(a)



(b)

Figure 10.15: (Unstable Zeros Cause Undershooting:) In (a), the zero is far to the right of the origin, and its effect is small, while in (b), the zero is closer to the origin and causes a lot of undershooting.

50%, and the overshoot is approximately doubled. For $z < 0$, the undershoot is roughly double the nominal overshoot.

- $|z| \approx \zeta\omega_n$, the zero has a profound effect on the step response step response. For $z > 0$, the 0-100% rise time is reduced by 75%, and the overshoot is quadrupled. For $z < 0$, the undershoot may be 100% or more.
- $|z| < 0.5\zeta\omega_n$, the step response is hardly recognizable as that of a second-order system!

Overdamped and Critically Damped, $\zeta \geq 1$

Zeros can also strongly affect the behavior of overdamped systems. When the zero is in the right-half plane, undershooting occurs. When the zero is in the left-half plane and is to the right of the two real poles, overshooting occurs. Even though there is overshooting, the response does not oscillate. The output reaches a peak value and then decays monotonically to its steady-state value. When the zero is between the two poles or to the left of both poles, no overshooting occurs.

10.7.5 (Optional Read:) Why Zeros Alter the Transient Response

Let $Y'(s)$ be the output of the transfer function in (10.81), when the input is a unit step. $Y'(s)$ can be written as

$$Y'(s) = Y(s) + \frac{1}{z} s Y(s),$$

where $Y(s)$ is the step response of a standard second-order system without a zero. In the time domain, recalling that $\frac{d}{dt} \leftrightarrow s$, we have

$$y'(t) = y(t) + \frac{1}{z} \frac{dy(t)}{dt}.$$

When z is large in magnitude, its reciprocal is small, and the added derivative term has little effect, meaning that $y'(t) \approx y(t)$. When z is small in magnitude, its reciprocal is large, and the derivative term can dominate the sum, in which case, $y'(t)$ is very different

Compensator Plant

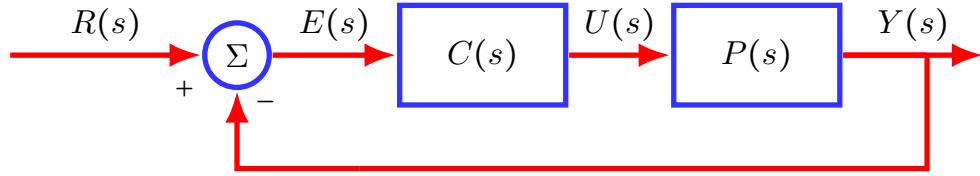


Figure 10.16: Reminder on a cascade compensator and plant in a unity feedback loop. This section seeks to design $C(s)$. It's easy to forget that we are manipulating ODEs, but thanks to Laplace, $s \leftrightarrow \frac{d}{dt}$.

from $y(t)$. For zeros in the left-half plane, $z > 0$. It follows that when $y(t)$ is increasing, its derivative is positive, and therefore $y'(t) = y(t) + \frac{1}{z} \frac{dy(t)}{dt} > y(t)$, leading to faster response and increased overshoot. When $z < 0$, it can be shown that

$$\left. \frac{dy'(t)}{dt} \right|_{t=0} < 0,$$

and thus $y'(t)$ initially decreases, giving rise to undershooting.

10.8 Design of Cascade Compensators and Pre-compensators

Control system design is in some sense the opposite of control system analysis. In analysis, one seeks to predict the response of a system to inputs. In design, one seeks to select the components in a system so that it will perform in a desired manner. This section tackles the problem of designing closed-loop systems to achieve stability, rapid response, low oscillation, and low steady-state error. It will draw upon all of the results we have developed throughout the chapter.

10.8.1 First-order Systems

First-order systems provide an opportunity to explore feedback concepts without too many complications. Consider, therefore, a first-order plant, possibly unstable and without a zero,

$$P(s) = \frac{k_0}{s + p_0}. \quad (10.82)$$

The initial design objective is to create a stable, closed-loop system with a 5% settling time of T_s seconds for a step input.

The design process begins by translating the specification for the transient response to the desired location of the closed-loop pole, or, equivalently for a first-order system, the desired time constant. From our analysis of first-order systems, we obtain

$$\tau_{des} = \frac{T_s}{3}.$$

(10.83)

Proportional Feedback for Stabilization of First-order Plants

The next step in the design process is to select the compensator and solve for the unknown gain or gains. The simplest compensator is proportional feedback in a unity feedback configuration, as in Fig. 10.16; hence we start with

$$C(s) = K_P. \quad (10.84)$$

The closed-loop system is then

$$G_{cl}(s) = \frac{C(s)P(s)}{1 + C(s)P(s)} = \frac{K_P k_0}{s + p_0 + K_p k_0},$$

which has time constant

$$\tau_{cl} = \frac{1}{p_0 + K_p k_0}.$$

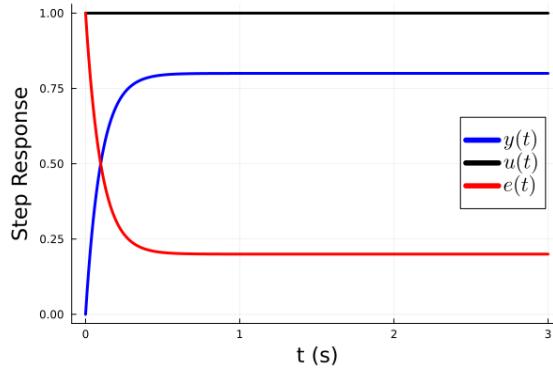


Figure 10.17: Step response and steady-state error for $P(s) = \frac{2}{s+2}$ and $C(s) = 4$. Black is the unit step input signal, blue is the closed-loop system's step response, and red is the tracking error of 20%.

Setting τ_{cl} equal to τ_{des} and solving for the compensator gain yields

$$K_P = \frac{\frac{1}{\tau_{des}} - p_0}{k_0}. \quad (10.85)$$

The final step is to compute the step response and verify that the design meets the specifications.

Example 10.41. (Stable first-order system) Suppose that $P(s) = \frac{2}{s+2}$, which has a time constant of 500 ms, and hence a 5% settling time of 1.5 s. It is desired that the closed-loop system be five times as fast, with a 5% settling time of 300 ms. Design an appropriate controller.

Solution: A proportional compensator is selected. From (10.83), the desired time constant is 100 ms. Following the steps outlined above, or immediately applying (10.85), yields

$$K_P = 4.$$

The step response is shown in Fig. 10.17 along with the error between the reference command and the output. It is observed that the steady-state error is 0.2, in other words, the error is 20%.

A system that is open-loop unstable is illustrated next.

Example 10.42. (Unstable first-order system) Suppose that $P(s) = \frac{2}{s-1}$ and the objective is closed-loop stability with a 5% settling time of 300 ms. Design an appropriate controller.

Solution: A proportional compensator is selected. From (10.83), the desired time constant is 100 ms and the required compensator gain is computed to be

$$K_P = 5.5$$

The step response is shown in Fig. 10.18 along with the error between the reference command and the output. It is observed that the steady-state error is -0.1.

In most applications, the sign of the error is irrelevant, only its magnitude counts. Hence, even in the above case when $e_{SS} = -0.1$, one still says the error is 10%, meaning that its magnitude is 10%.

Fixing the Steady-state Error

In Examples 10.41 and 10.42, the steady-state error is 20% and 10%, respectively. Can it be reduced? There are two ways to accomplish this. Consider once again the general first-order plant given in (10.82) and the proportional compensator in (10.84). Applying Chapter 10.6.2 to the system yields the transfer function of the error,

$$\frac{E(s)}{R(s)} = \frac{1}{s + p_0 + K_P k_0}.$$

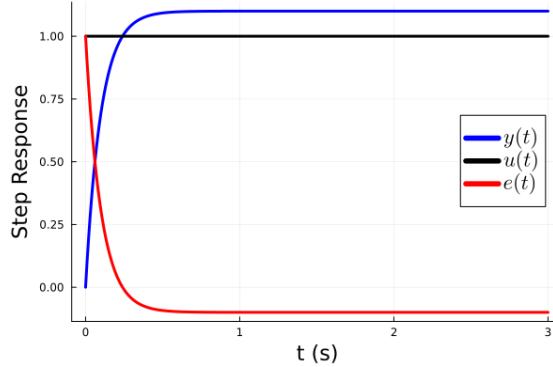


Figure 10.18: Step response and steady-state error for $P(s) = \frac{2}{s-1}$ and $C(s) = 5.5$. Black is the unit step input signal, blue is the closed-loop system's step response, and red is the tracking error of 10%.

Assuming that K_P has been chosen to be consistent with BIBO stability (namely, $p_0 + K_P k_0 > 0$), then the steady-state error for a unit-step input is

$$e_{SS} = \frac{p_0}{p_0 + K_P k_0}. \quad (10.86)$$

It follows that

$$\lim_{|K_P| \rightarrow \infty} e_{ss} = 0.$$

The feedback gain can never be set to infinity, but the above analysis demonstrates that the magnitude of the steady-state error can be made as small as desired by taking $|K_P|$ to be sufficiently large, assuming, once again, that closed-loop stability can be assured for large values of the feedback gain AND assuming you have a powerful enough actuator to handle the resulting control signal. Powerful motors on mobile devices require heavy batteries. Engineering is all about tradeoffs.

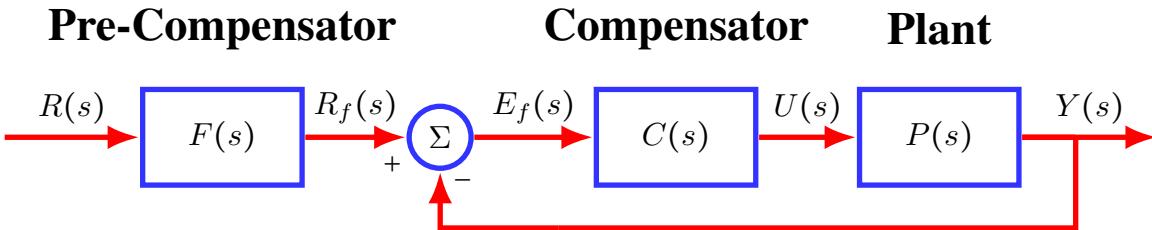


Figure 10.19: A unity feedback loop with a pre-compensator (aka, pre-filter) $F(s)$ added to the overall system. The pre-filter can be used to ameliorate the effects of left-half plane zeros or to “fix” steady-state error. The overall structure of the diagram is called a two-degree of freedom controller: one degree of freedom is the forward-path compensator $C(s)$ that can now focus more on stabilization, while the second degree of freedom is the pre-compensator, $F(s)$, that can take care of some of the transient and steady-state tasks that were not handled by $C(s)$.

There is a second means of reducing the steady-state error that does not involve high controller gains. Figure 10.19 shows a pre-compensator inserted in the system. The true error remains as

$$E(s) := R(s) - Y(s),$$

which is different than the “filtered error” $E_f(s)$. After a bit of algebra, and assuming closed-loop stability, one computes for a unit-step input,

$$e_{SS} = \frac{1 + (1 - F(0)) C(0) P(0)}{1 + C(0) P(0)}. \quad (10.87)$$

Solving the above for $F(0)$, the DC-gain for the pre-compensator, to achieve $e_{SS} = 0$ yields

$$F(0) = \frac{1 + C(0) P(0)}{C(0) P(0)}. \quad (10.88)$$

For Example 10.41, we can let the pre-compensator be the constant gain,

$$F(s) := K_f = \frac{5}{4},$$

and for Example 10.42 we can let the pre-compensator be the constant gain,

$$F(s) := K_f = \frac{10}{9}.$$

All we are doing is exploiting linearity: by inserting a scale factor, a “dummy” reference signal, R_f , is created that achieves the desired steady-state error between the real reference signal and the measured output.

10.8.2 Second-order Systems

We now consider a second-order system without a zero,

$$P(s) = \frac{k_0}{s^2 + a_1 s + a_0}. \quad (10.89)$$

As in the case of first-order systems, we will first address the transient response, which includes stability. When using a PD controller, it will insert a (stable) zero that will “mess up” our design considerations. In a second design phase, will address the zero with a pre-compensator. We will also take up steady-state error.

The design process begins by translating the specifications for the transient response into the desired location of the closed-loop poles. We will continue to assume that one of the specifications is the 5% settling time. For simplicity and definiteness, we will take the second specification to be the percent overshoot. From Section 10.7.3, the poles are uniquely determined by the damping ratio ζ and the undamped natural frequency ω_n . If we further assume that $0 < \zeta < 1$, then the transient specifications relate to the desired closed-loop poles via

$$\begin{aligned} \zeta &\leftrightarrow \text{percent overshoot} = 100e^{\frac{-\zeta\pi}{\sqrt{1-\zeta^2}}} \% \\ \omega_n &\leftrightarrow \text{5\% settling time} \approx \frac{3}{\zeta\omega_n}. \end{aligned} \quad (10.90)$$

In the above expression, the limit as $\zeta \rightarrow 1$ yields no overshoot. Common numbers to keep in mind for the damping ratio are given in the table below. Only novice control system designers use the exact formula for the damping ratio from (10.90); experienced designers realize that some tweaking will be required in the end no matter what because the model is never exact in practice (yes, it can be in HW problems).

ζ	\approx Overshoot	True Overshoot
1.0	0 %	0.00 %
0.9	0 %	0.15 %
0.8	2 %	1.5 %
0.7	5 %	4.6 %
0.6	10 %	9.5 %
0.5	15 %	16.3 %

PD Control

The next step in the design process is to select the compensator and solve for the unknown gain or gains. A PD compensator will be assumed in a unity feedback configuration, as in Fig. 10.16. Hence we start with

$$C(s) = K_P + K_D s = K_D \left(s + \frac{K_P}{K_D} \right).$$

The closed-loop system is then

$$\begin{aligned} G_{cl}(s) &= \frac{C(s)P(s)}{1 + C(s)P(s)} \\ &= k_0 \frac{K_P + K_D s}{s^2 + (a_1 + K_D k_0)s + (a_0 + K_P k_0)} \\ &= k_0 \frac{K_D(s + K_P/K_D)}{s^2 + (a_1 + K_D k_0)s + (a_0 + K_P k_0)}, \end{aligned}$$

which shows that the characteristic equation is

$$s^2 + \underbrace{(a_1 + K_D k_0)}_{2\zeta\omega_n} s + \underbrace{(a_0 + K_P k_0)}_{\omega_n^2} = 0.$$

Identifying the various terms and performing elementary algebra show the relationship between the compensator gains and the closed-loop poles to be

$$\begin{aligned} K_P &= \frac{\omega_n^2 - a_0}{k_0} \\ K_D &= \frac{2\zeta\omega_n - a_1}{k_0}. \end{aligned} \tag{10.91}$$

Due to the PD compensator, the closed-loop transfer function has a zero at $-K_P/K_D$. Recall that Section 10.7.4 showed that a zero could significantly alter the step response of a system. This will motivate a modification to our unity feedback configuration.

Before moving on, this video on a “poorly tuned” PD controller shows the importance of what we are studying: [Control System Tuning #rocket #engineering gyat](#) by BPS.shorts.

Example 10.43. (Stabilizing an oscillator) Given that $P(s) = \frac{1}{0.5s^2 + 2}$, which has poles $\pm i2$, it is desired that the closed-loop system have less than 5% overshoot and a 5% settling time of 700 ms. Show that the step response of the open-loop system has sustained oscillations and then design a controller to meet the objectives for the closed-loop system.

Solution: In the s-domain (aka, Laplace domain), the response of the open-loop system to a step input is

$$\begin{aligned} Y(s) &= P(s) \frac{1}{s} \\ &= \frac{2}{(s + i2)(s - i2)} \frac{1}{s} \\ &= \frac{A_1}{s} + \frac{A_2}{s + i2} + \frac{A_3}{s - i2}, \end{aligned}$$

where $A_1 = \frac{1}{2}$, $A_2 = -\frac{1}{4}$, $A_3 = -\frac{1}{4}$. Taking the inverse Laplace transform term-by-term and applying Euler’s formula leads to

$$\begin{aligned} y(t) &= \frac{1}{2}u(t) - \frac{1}{4}e^{-i2t}u(t) - \frac{1}{4}e^{i2t}u(t) \\ &= \frac{1}{2}(1 - \cos(2t))u(t). \end{aligned}$$

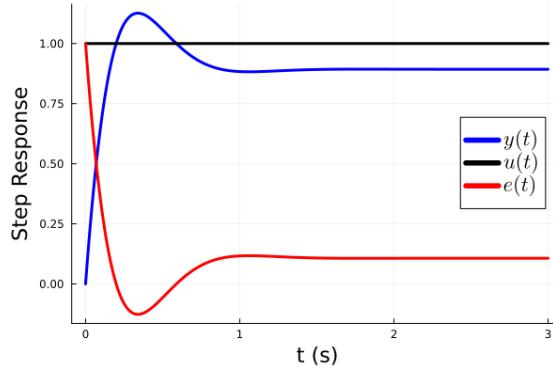


Figure 10.20: Step response and steady-state error for $P(s) = \frac{1}{0.5s^2+2}$ and $C(s) = 16.7 + 4.3s$. The zero in the forward path induced by the PD compensator (controller) causes more overshoot than is desired. This can be fixed with a pre-compensator.

Hence, the step response of the open-loop system is a sustained oscillation at a frequency of 2 rad/s.

For the controller, a PD compensator is assumed. From the overshoot specification, the damping ratio is selected to be $\zeta = 0.7$. From the settling time specification and (10.90), the undamped natural frequency is computed to be $\omega_n = 6.1$. Following the control design steps outlined above, or immediately applying (10.91), yields

$$K_P = 16.7 \text{ and } K_D = 4.3.$$

The step response is shown in Fig. 10.20 along with the error between the reference command and the output. The steady-state error is 11%. The 5% settling time is 1.2 s, almost twice the design value. The percent overshoot is evaluated to be

$$\frac{y_{max} - y_\infty}{y_\infty} \times 100\% = \frac{1.13 - 0.89}{0.89} \times 100\% = 27\%,$$

more than five times the design value.

To understand why the transient performance differs from the design specifications, we have to recall that the PD compensator introduced a zero in the forward path of the feedback loop. Where did the zero of the PD compensator end up in relation to the closed-loop poles? The zero is at

$$-\frac{K_P}{K_D} = -3.88,$$

which places it to the right of the poles, which have real part

$$-\zeta\omega_n = -4.27.$$

Based on Fig. 10.14, a large overshoot is expected. We can address this through the two-degree-of-freedom design in Figure 10.19, where this time, $F(s)$ will be more than just a gain. The transfer function from reference to output is

$$\frac{Y(s)}{R(s)} = F(s) \frac{C(s)P(s)}{1 + C(s)P(s)},$$

which is simply the product of the transfer functions of the pre-compensator and the standard unity feedback loop.

The pre-compensator can be used to remove unwanted left-half plane zeros in the transfer function by the deliberate creation of **stable pole-zero cancellations**. If there are no unstable pole-zero cancellations, the control loop with the pre-compensator is BIBO stable if, and only if, both the pre-compensator and the control loop without the pre-compensator are BIBO stable. Moreover, if $F(0) = 1$, that is, the pre-compensator has unity DC gain, then the steady-state error for a unit-step input is the same, with or without the pre-compensator, namely

$$e_{SS} = \frac{1}{1 + C(0)P(0)}.$$

Example 10.44. (Stabilizing an oscillator continued) Consider again the system in Example 10.43, where it was desired that the closed-loop system have less than 5% overshoot and a 5% settling time of 700 ms. Use a two-degree-of-freedom controller to remove the troublesome zero introduced by the PD compensator.

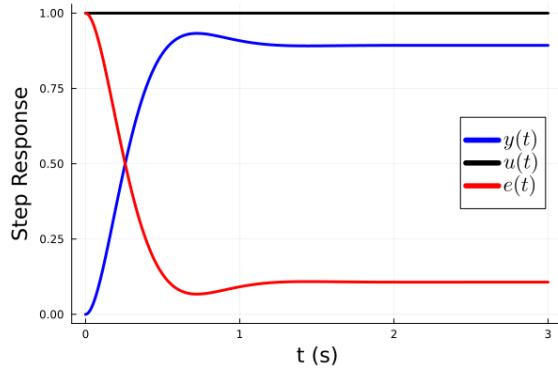


Figure 10.21: Step response when using a pre-compensator. The plant is $P(s) = \frac{1}{0.5s^2+2}$, $C(s) = 16.7 + 4.3s = 4.3(s + 3.88)$ and $F(s) = \frac{3.88}{s+3.88}$.

Solution: The PD compensator $C(s)$ is retained. The pre-compensator is selected as

$$F(s) = \frac{z}{s+z} = \frac{1}{\frac{s}{z} + 1}, \quad (10.92)$$

where $z = \frac{K_P}{K_D} = 3.88$. This cancels the zero with a stable pole. The step response of the system is shown in Fig. 10.21. The overshoot and settling time specifications are now met.

Remark: We leave it to the reader to adjust the DC gain of the pre-compensator, $F(s)$, using (10.88) to zero the steady-state error.

10.8.3 Dealing with More Complicated Systems

The methods we have been following have allowed us to assign all of the poles of first- and second-order systems in order to meet a set of transient and steady-state design specifications. Assigning all the poles breaks down rapidly beyond the cases we have considered. To be sure, a few more tricks can be brought to bear, thereby slightly extending the range of systems that can be treated using the approach developed in this chapter, and these extensions will be outlined in the design examples of the next section and in the homework exercises. Nevertheless, faced with a plant with four poles and a zero, even these extensions do not work very often. So, what is an aspiring control designer to do?

At the junior or senior undergraduate level, most universities offer a course in feedback control. There, you may be exposed to the Root Locus Method of Evans and the Frequency Domain Methods of Nyquist and Bode. These methods will allow you to treat single-input, single-output, linear plants of considerable complexity. You will also be able to consider more comprehensive design specifications. At the MS level, many universities offer courses in multivariable linear systems, nonlinear systems, stochastic systems, adaptive systems, discrete-event systems, and much more. Control Systems is a very active research area, rich with theoretical questions and exciting, novel applications.

10.9 Feedback Design for a Linearized Model of a Planar Segway Transporter

The key transfer functions for the Segway are given in Example 10.22 and Fig. 10.3. We will stabilize the body lean angle before moving on to regulating speed.

10.9.1 Controlling Body Lean Angle

The relevant transfer function is

$$\frac{\theta(s)}{\tau(s)} := G_\theta(s) = \frac{-21.468}{s^2 - 50.937}. \quad (10.93)$$

The transfer function comes from the linearized model of the Segway developed in Chapter 10.3.6. It is a simple second-order transfer function with poles at $\pm\sqrt{50.937} = \pm 7.137$. Hence, the system is unstable and requires feedback control to stay upright.

Iteration # 1

We'll use the time-domain specs from Example 10.44, namely less than 5% overshoot and a 5% settling time of 700 ms. We'll also use a two-degree-of-freedom controller to remove the troublesome zero that we will introduce via a PD compensator.

We need to turn the specs into gains for the PD compensator, using (10.91), where for this plant, $k_0 = -21.468$, $a_0 = -50.937$ and $a_1 = 0$. Turning the crank, we obtain,

$$\begin{aligned}\zeta &= 0.7 \\ \omega_n &= \frac{3}{\zeta \cdot T_s} = 6.122 \\ K_P &= \frac{\omega_n^2 - a_0}{k_0} = -4.12 \\ K_D &= \frac{2\zeta\omega_n - a_1}{k_0} = -0.40,\end{aligned}$$

and hence the PD compensator is

$$C_{PD}(s) = K_D \left(s - \frac{K_P}{K_D} \right) = -0.28(s + 11.55).$$

With these parameters, we can compute the steady-state error,

$$e_{SS} = \frac{D(0)}{D(0) + N(0)} = \frac{a_0}{a_0 + k_0 K_P} = -1.36 \text{ or } 136\%,$$

which is unacceptably large. We choose to take a different approach to the performance specs.

```

1 # Define plant transfer function params
2 # P1 = k0/(s^2 + als + a0)
3 k0 = -21.468
4 a1 = 0
5 a0 = -50.937
6
7 # Set time domain specs for the closed-loop system
8 zeta = 0.7 # 5 % overshoot
9 Ts = 0.7
10
11 # Ts = 3.0/(zeta*wn)
12 wn = 3.0/(zeta*Ts)
13
14 # PD gains
15 Kp= (wn^2 - a0)/k0
16 Kd= (2*zeta*wn - a1)/k0
17
18 # Zero location for PD controller
19 # Kd*(s+z)
20 z = Kp/Kd
21
22 # Steady state error
23 # N(s)/D(S) = C(s)*P(s)
24 #           = (Kp + Kds)* k0/(s^2 + als + a0)
25 #
26 # eSS = D(0) / (D(0) + N(0))
27 eSS = a0 / (a0 + Kp*k0)
28
29 # Print relevant data
30 println("zeta = $zeta")
31 println("Ts = $Ts")
```

```

32 println ("wn = $wn")
33 println ("Kp = $Kp")
34 println ("Kd = $Kd")
35 println ("z = $z")
36 println ("eSS = $eSS")

```

Output

```

zeta = 0.7
Ts = 0.7
wn = 6.122448979591837
Kp = -4.118752632183022
Kd = -0.3992653517527749
z = 10.315827842565596
eSS = -1.3588859666666666

```

Iteration #2

We now know that steady-state error is a big challenge for this system. We can flip the script and design for a given e_{SS} and overshoot (related to damping ratio, ζ), leaving the settling time (related to ω_n) as a free parameter.

For the design, we retain $\zeta = 0.7$ and specify $e_{SS} = 20\%$. Both $e_{SS} = 0.20$ and $e_{SS} = -0.20$ satisfy the spec. Which one to choose? The one that gives real numbers in the calculations below! It turns out that we need $e_{SS} = -0.20$ (try the other one, and you'll agree). Following this plan, from the expression for steady-state error, we obtain the proportional gain,

$$K_P = \frac{a_0 - a_0 e_{SS}}{k_0 e_{SS}} = -14.236.$$

The undamped natural frequency is then

$$\omega_n = \sqrt{K_P k_0 + a_0} = 15.959,$$

and the derivative gain,

$$K_D = \frac{2\zeta \omega_n - a_1}{k_0} = -1.041,$$

giving the PD compensator

$$C(s) = -14.236 + -1.041 s = -1.041(s + 13.679). \quad (10.94)$$

The 5% settling time is

$$T_s \approx \frac{3}{\zeta \omega_n} = 0.269.$$

In the code, we also design the pre-compensator, $F(s)$, using (10.92) to cancel the stable zero introduced by the PD compensator and we adjust its DC gain using (10.88) to set the steady-state error to zero. The resulting transfer function from the commanded value for the lean angle, $\theta^{des}(s)$, to $\theta(s)$ is then

$$\frac{\theta(s)}{\theta^{des}(s)} = F(s) \frac{C(s) P(s)}{1 + C(s) P(s)} = \frac{254.685}{s^2 + 22.342 s + 254.685}, \quad (10.95)$$

after removing the stable pole-zero cancellation. We note that the transfer function has a DC gain of one.

```

1 # Define plant transfer function params
2 # P = k0/(s^2 + als + a0)
3 k0= -21.468
4 a1 = 0
5 a0 = -50.937
6
7 # Specifications

```

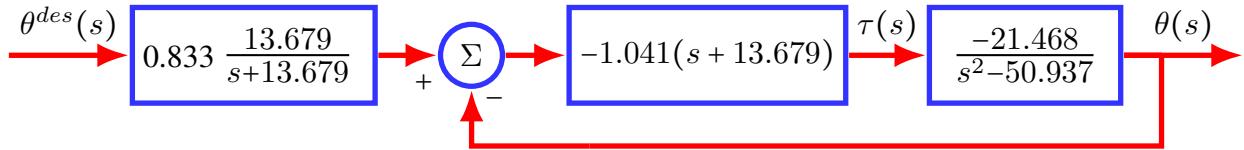


Figure 10.22: Closed-loop system to stabilize the lean angle or orientation of the Segway's body. The overall transfer function is given in (10.95), where the stable pole-zero cancellation has been removed. The result is a standard second-order system with a DC gain of one.

```

8 zeta = 0.7 # 5% overshoot
9 eSS = -0.20
10
11 # Proportional Gain
12 Kp = (a0 - a0*eSS) / (k0*eSS)
13
14 # Undamped natural frequency
15 # becomes complex with eSS = +0.20
16 wn = sqrt(Kp*k0 + a0)
17
18 # Derivative gain
19 Kd= (2*zeta*wn - a1) / k0
20
21 # Zero location for PD controller
22 # Kd*(s+z)
23 z = Kp/Kd
24
25 # Settling time
26 Ts = 3 / (zeta*wn)
27
28 # Steady state error to double check
29 # N(s)/D(S) = C(s)*P(s)
30 # = (Kp + Kds) * k0 / (s^2 + als + a0)
31 #
32 # eSS = D(0) / (D(0) + N(0))
33 eSS = a0 / (a0 + Kp*k0)
34
35 # Pre-compensator DC gain
36 kF = (1 + Kp*k0/a0) / (Kp*k0/a0)
37
38 # Print relevant data
39 println("zeta = $zeta")
40 println("Ts = $Ts")
41 println("wn = $wn")
42 println("Kp = $Kp")
43 println("Kd = $Kd")
44 println("z = $z")
45 println("eSS = $eSS")
46 println("kF = $kF")

```

Output

```

zeta = 0.7
Ts = 0.2685477580249118
wn = 15.958853342267417
Kp = -14.236165455561764
Kd = -1.0407301415676533

```

```

z   = 13.679017150514934
eSS = -0.20000000000000004
kF = 0.8333333333333333

```

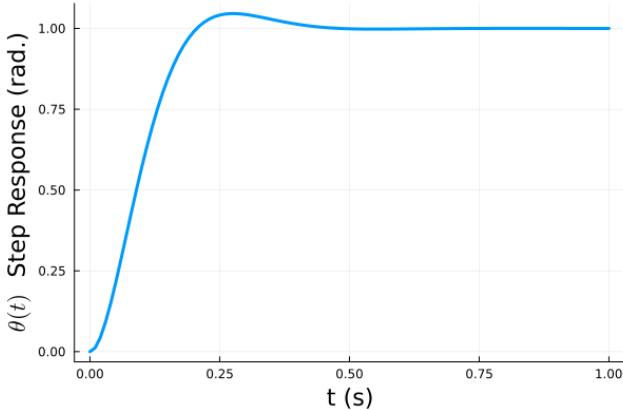
Next, we check the step response to see that the specs are actually met.

```

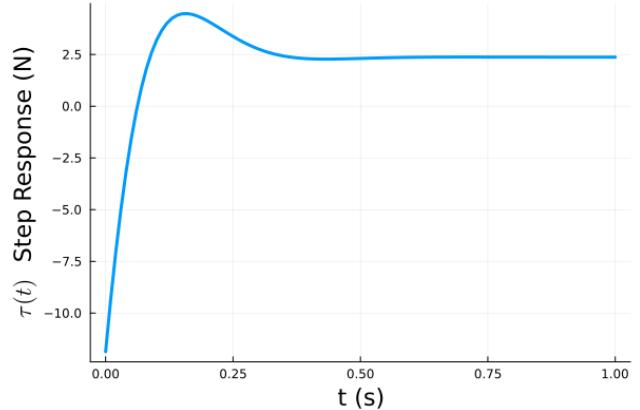
1 using ControlSystems, Plots, LaTeXStrings
2
3 # Define the Plant, Controller, and Precompensator Transfer Functions
4 P = tf([k0], [1, a1, a0]) # Plant transfer function
5 C = tf([Kd, Kp], [1]) # Controller transfer function
6 F = tf([kF], [1/z, 1]) # Precompensator transfer function
7
8
9 # Combine the Transfer Functions
10 G = C * P
11
12 # Create a Unity Feedback Loop
13 CL = feedback(G, 1)
14
15 # Multiply by the precompensator
16 sys = F * CL
17
18 # Compute and Plot the Step Response
19 t = 0:0.01:1 # Time vector
20 y, t = step(sys, t)
21 p1 = plot(t, y', lw=3, label=false, guidefont = 15,
22           xlabel="t (s)", ylabel=L"\theta(t) Step Response (rad.)")
23
24
25
26 tfU = F*C/(1+C*P)
27 # Compute and Plot the Step Response
28 t = 0:0.01:1 # Time vector
29 u, t = step(tfU, t)
30
31 p2 = plot(t, u', lw=3, label=false, guidefont = 15,
32           xlabel="t (s)", ylabel=L"\tau(t) Step Response (N)")
33
34 display(p1)
35 display(p2)

```

Output See Fig. 10.23.



(a)

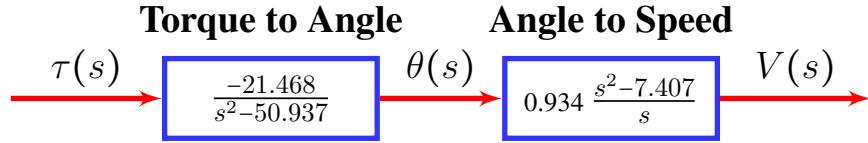


(b)

Figure 10.23: Step responses: (a) Lean angle of the Segway and rider. The overshoot is 5% and the steady-state error is zero. (b) The torque required to achieve the step response. The peak of -10 N is a lot of torque, but you have to realize that a “unit step” in θ is one radian (approximately 60°). Because the system is linear, the torque required for 10° would be -1.67 N, which is realistic. The sign just specifies the direction the motor is spinning.

10.9.2 Controlling Speed and Lean Angle

Our objective now is to control speed in addition to orientation. Of course, we cannot just add another actuator to the overall system. While we could perhaps switch back and forth from controlling one quantity for a few tens of milliseconds and then the other, that would make for a really interesting ride! We need a better idea. It’s time to exploit the factorization in Fact 10.23, namely



Recall that this comes from the two transfer functions for the Segway given in (10.37),

$$G_\theta(s) = \frac{\theta(s)}{\tau(s)} = -21.468 \frac{1}{s^2 - 50.937}$$

$$G_v(s) = \frac{V(s)}{\tau(s)} = -20.057 \frac{s^2 - 7.404}{s(s^2 - 50.937)},$$

which share a common term. Taking their ratio to eliminate the common factor yields

$$\frac{V(s)}{\theta(s)} = 0.934 \frac{s^2 - 7.407}{s};$$

that is,

$$V(s) = \left(0.934 \frac{s^2 - 7.407}{s} \right) \theta(s). \quad (10.96)$$

Hence, appending the transfer function $\left(0.934 \frac{s^2 - 7.407}{s} \right)$ to the output of the feedback system in Fig. 10.22 results in the new output being speed, and the input (which we treat as a pseudo actuator) being θ^{des} .

While Fig. 10.24 may look complicated, from our work on regulating the body orientation, namely (10.95), we have that

$$\frac{\theta(s)}{\theta^{des}(s)} = \frac{254.685}{s^2 + 22.342 s + 254.685},$$

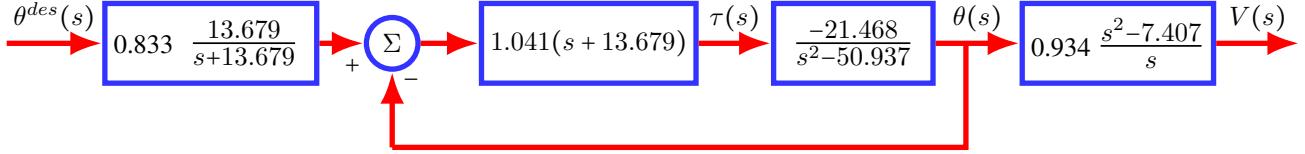


Figure 10.24: The speed Portion of the Segway’s transfer function has been appended to the closed-loop system that stabilizes the orientation of the Segway’s body.

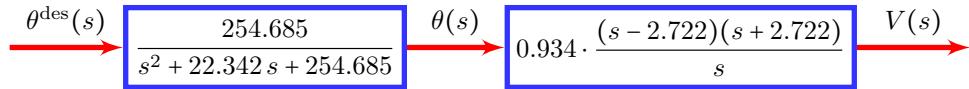
which allows us to replace the first three blocks of the feedback diagram with a standard second-order transfer function, yielding

$$V(s) = \left(0.934 \frac{s^2 - 7.407}{s} \right) \left(\frac{254.685}{s^2 + 22.342s + 254.685} \right) \theta^{des}(s),$$

or equivalently,

$$V(s) = \left(0.934 \frac{(s - 2.722)(s + 2.722)}{s} \right) \left(\frac{254.685}{s^2 + 22.342s + 254.685} \right) \theta^{des}(s). \quad (10.97)$$

Hence, we can now focus on the transfer function,



We need to stabilize the pole at the origin and mitigate the effects of the zeros. We cannot do anything about the unstable zero, but we can cancel the stable zero with the transfer function

$$\frac{2.722}{s + 2.722}.$$

Because the poles of $\frac{254.685}{s^2 + 22.342s + 254.685}$ are $-11.171 \pm 11.397i$, which are already nicely into the left half plane, it turns out that, similar to a pole of a first-order transfer function, the pole at the origin can be stabilized with proportional feedback. Hence, for our second compensator, we propose the simple transfer function,

$$C_2(s) = K_2 \frac{2.722}{s + 2.722}, \quad (10.98)$$

where the gain K_2 will be tuned to achieve a “nice step response”, that is, one that is sufficiently fast, but does not overshoot. We do not need to worry about the steady-state error because it will be zero due to the pole at the origin.

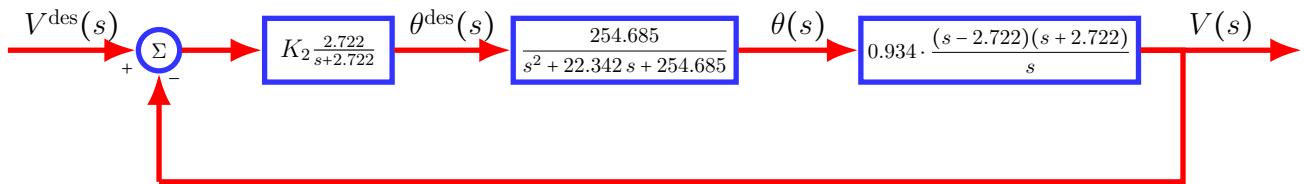


Figure 10.25: K_2 must be adjusted to obtain a nice step response. This can be done in simulation. Later, as an optional read, we show that the appropriate range for BIBO stability is $-0.227 < K_2 < 0$. The negative value for K_2 is due to the term, $(s - 2.722)$.

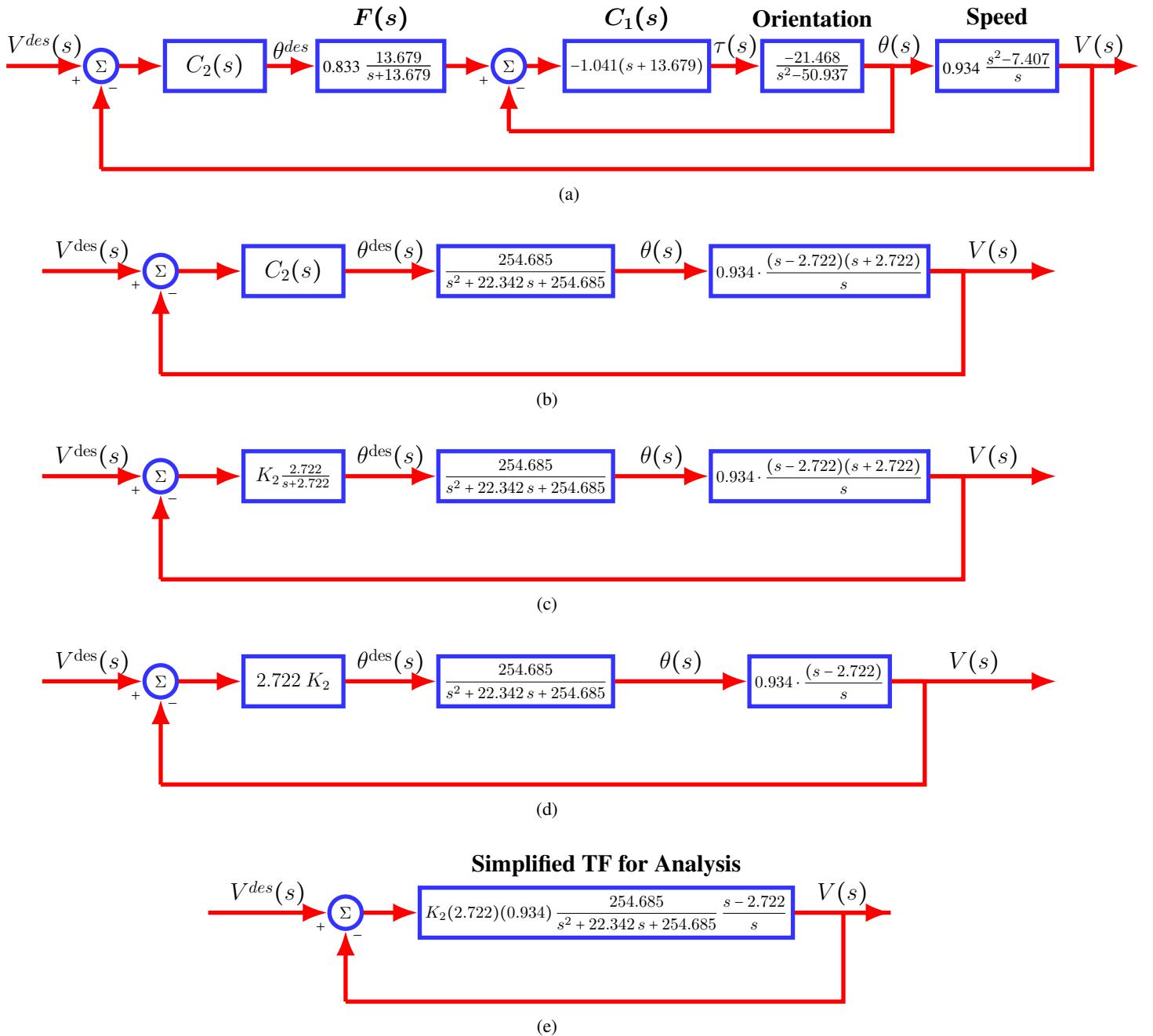


Figure 10.26: (**Steps for Understanding the Design of the Speed Controller**) This part of the design is harder because the transfer function involved is more complicated. However, because of how we designed $F(s)$ and $C_1(s)$ it simplifies greatly. The transfer function from θ^{des} to θ is a standard second-order transfer function; moreover, the numerator of the speed transfer function factors into two real zeros. This makes the design of $C_2(s)$ look much more approachable in (b) and (c). We hence propose that $C_2(s)$ as a proportional gain K and a “pre-compensator-like” term that cancels the “stable” zero. The overall transfer function simplifies in (e) to a third-order system with a zero. Its stability can be analyzed in closed-form.

(Optional Read:) Analysis Required to Stabilize the Speed

To analyze how to regulate the Segway's forward speed using feedback control, we begin by constructing the full open-loop transfer function of the speed control system. The proposed control architecture consists of a compensator $C_2(s)$ in series with the system's dynamics, all placed in a unity feedback loop.

The three main components in the forward path are:

- The speed transfer function,

$$\frac{V(s)}{\theta(s)} = 0.934 \cdot \frac{(s - 2.722)(s + 2.722)}{s}.$$

- The closed-loop lean angle dynamics,

$$\frac{\theta(s)}{\theta_{\text{des}}(s)} = \frac{254.685}{s^2 + 22.342 s + 254.685}.$$

- The proposed compensator for speed control that deliberately cancels a stable zero,

$$C_2(s) = K_2 \cdot \frac{2.722}{s + 2.277}.$$

The full open-loop transfer function is the product of these three components,

$$G(s) = C_2(s) \cdot \left(\frac{\theta(s)}{\theta_{\text{des}}(s)} \right) \cdot \left(\frac{V(s)}{\theta(s)} \right) = C_2(s) \cdot \left(\frac{V(s)}{\theta(s)} \right) \cdot \left(\frac{\theta(s)}{\theta_{\text{des}}(s)} \right).$$

Substituting each part in, we obtain,

$$\begin{aligned} G(s) &= \left(K_2 \cdot \frac{2.722}{s + 2.277} \right) \cdot \left(0.934 \cdot \frac{(s - 2.722)(s + 2.722)}{s} \right) \cdot \left(\frac{254.685}{s^2 + 22.342 s + 254.685} \right) \\ &= (2.722 K_2) \cdot \left(0.934 \cdot \frac{(s - 2.722)}{s} \right) \cdot \left(\frac{254.685}{s^2 + 22.342 s + 254.685} \right) \\ &= 2.542 K_2 \cdot \left(\frac{s - 2.722}{s} \right) \cdot \left(\frac{254.685}{s^2 + 22.342 s + 254.685} \right), \end{aligned}$$

where the zero at $s = 2.722$ in the speed transfer function is canceled by the pole at $s = -2.722$ in the compensator, $C_2(s)$. This is an intentional cancellation of the stable zero, leaving only the unstable zero at $s = +2.722$.

The terms in the standard second-order system came from setting overshoot and 5%-settling time. For a reason that will become clear in a second, we rewrite $G(s)$ as,

$$G(s) = K_3 \cdot \left(\frac{s - z_2}{s} \right) \cdot \left(\frac{\omega_n^2}{s^2 + 2\zeta\omega_n s + \omega_n^2} \right),$$

where $K_3 = 2.542 K_2$, $z_2 = 2.722$, $\zeta = 0.7$, and $\omega_n = 15.959$. The closed-loop transfer function, with unity feedback, is then given by,

$$\frac{V(s)}{V^{\text{des}}(s)} = \frac{G(s)}{1 + G(s)} = \frac{K_3 \omega_n^2 (s - z_2)}{s^3 + \underbrace{2\zeta\omega_n}_{a_2} s^2 + \underbrace{\omega_n^2(1 + K_3)}_{a_1} s + \underbrace{(-K_3 \omega_n^2 z_2)}_{a_0}}.$$

This expression forms the basis for analyzing the stability and performance of the system from our simplified stability rules for third-order polynomials, namely,

$$a_0 > 0 \iff K_3 < 0 \text{ because } z_2 > 0$$

$$a_2 > 0 \iff \text{automatic from the definitions of } \zeta \text{ and } \omega_n$$

$$a_1 - \frac{a_0}{a_2} > 0 \iff 1 + K_3 \left(1 + \frac{z_2}{2\zeta\omega_n} \right) > 0 \iff 1 + K_3 \left(1 + \frac{z_2 T_s}{6} \right) > 0 \iff K_3 > -\frac{6}{6 + z_2 T_s},$$

where $0.269 = T_s \approx \frac{3}{\zeta\omega_n}$ is the value we used for the 5% settling time. This gives

$$-0.891 < K_3 < 0 \iff -0.351 < K_2 < 0$$

for the range of gains that will stabilize the speed loop, once we substitute in for T_s and $K_3 = 2.542 K_2$.

Step Responses

The final feedback control diagram is given in Fig. 10.27 and a family of step responses for various values of K is given in Fig. 10.28.

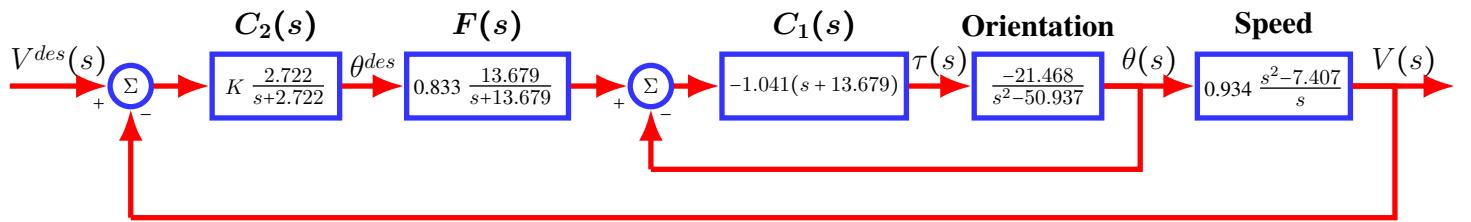


Figure 10.27: (**Inner-Outer Feedback Loops for the Segway**) There are two unity feedback loops. The inner loop is from Fig. 10.22, our controller for stabilizing the orientation of the Segway. It includes everything from $\theta^{des}(s)$ (the input to the pre-compensator, $F(s)$) all the way to $\theta(s)$, the output of the orientation block. The other blocks and the left-most summing junction form the outer unity feedback loop for stabilizing the Segway's speed.

```
1 using ControlSystems, Plots, LaTeXStrings
```

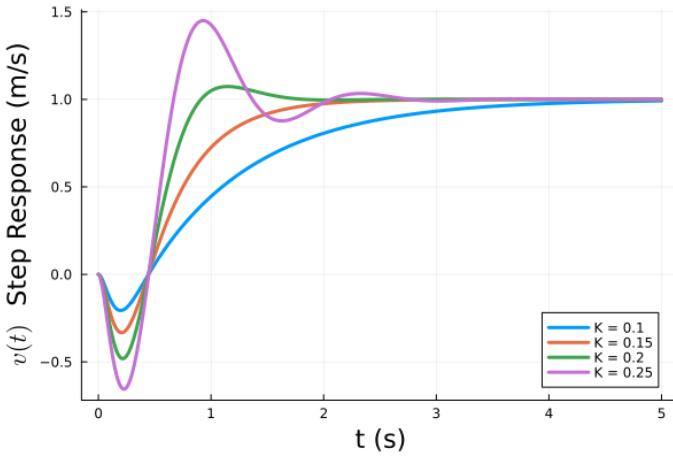
```
2
3
4 # Plant transfer function for the inner feedback loop
5 # from theta desired to theta
6 Pa = tf([254.685], [1, 22.342, 254.685])
7
8 # remaining part of the Segway's transfer fuction from torque to speed
9 Pb = tf([0.934, 0, -0.934*7.407], [1.0, 0])
10
11 # design transfer function for the outer feedbck loop
12 P2 = Pa*Pb # Final Plant TF
13
14 # Search over gains to find a good one
15 # remainder of the controller is inside the for loop
16 K = -0.1*[1, 1.5, 2, 2.5]
17
18 # Create empty plot that we can add to in the loop
19 p1 = plot(guidefont = 15)
20
21 # For loop to compute multiple step responses
22 for i = 1:length(K)
23
24 # Compensator depending on K
25 # It cancels a stable zero
26 C2 = tf([K[i]*2.722], [1, 2.722])
27 G2 = C2*P2
28
29 # Outer Unity feedback loop
30 sys2 = feedback(G2, 1)
31
32 # Compute and Plot the Step Response
33 t = 0:0.01:5.0
34
35 # Time vector
36 v, t = step(sys2, t)
37
38 # Create a LaTeX string for the legend
```

```

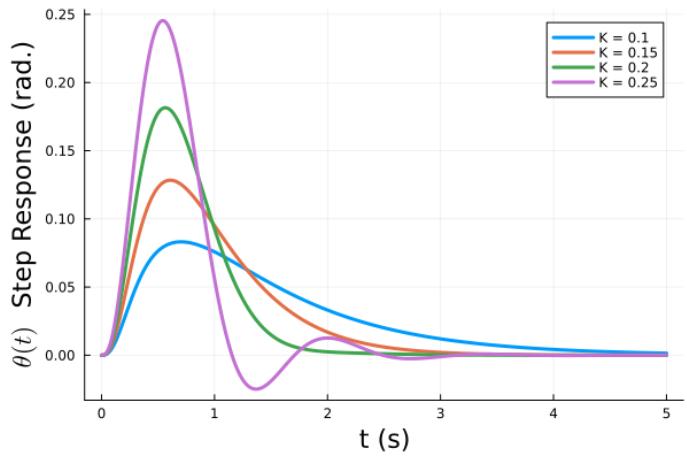
39 K_label = string("K = ", round(K[i], digits=2))
40
41 # Add to plot
42 p1 = plot!(p1, t, v', lw=3, label=K_label)
43
44 end
45
46
47 # Make the plot
48 p1 = plot!(p1, xlabel="t (s)", ylabel=L"v(t) \u2248 Step Response (m/s)")

```

Output See Fig. 10.28.



(a)



(b)

Figure 10.28: **Step responses:** (a) Shows the step response in speed for a range of gains K in $C_2(s)$, while (b) shows the corresponding lean angle trajectories. It is worth looking at both plots. Here, a unit step in speed means 1 m/s, a comfortable walking speed, hence somewhat slow in Segway terms. For $K = 0.2$, there is some overshoot in the speed response, but perhaps not too much. However, in the lean angle, we see approximately 0.17 radians, or 19.5° of undershoot, which is a lot. The safer bet is likely $K = 0.15$ or even 0.1.

10.9.3 The Real Deal: Implementing the Controller on the Nonlinear Model

It's time for the rubber to meet the road, literally. Our goal is to implement the controller on the Segway's nonlinear model. We break the overall goal into simpler steps:

- Separate the controller terms in Fig. 10.27 from the model terms.
- Transform the controller from the Laplace domain into ODEs in the time domain.
- Verify the controller on the linearized ODE model of the Segway before evaluating it on the nonlinear model. Because the controller worked well on the transfer function, it "has to work" on the ODE version of the linearized model, **unless our simulation code has a bug**.
- Replace the linear model with the nonlinear model. For "small" initial conditions, the response should be nearly identical to the linearized model. For "large" initial conditions, all bets are off.

Isolating the Controller

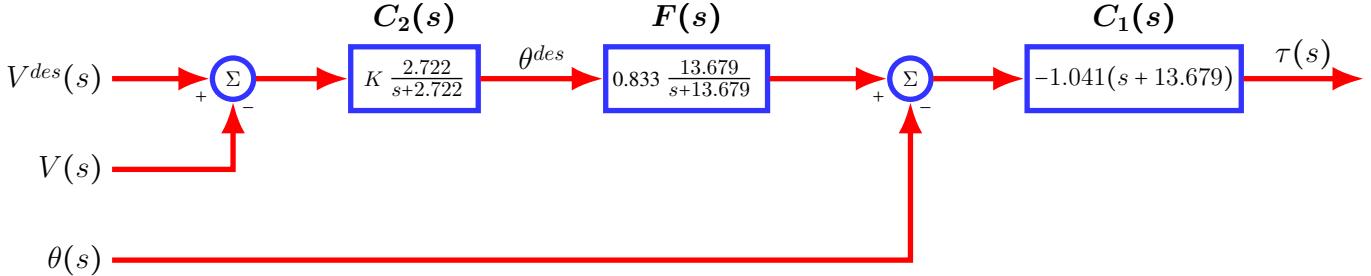


Figure 10.29: (**Just the Controller for the Segway**) Once the transfer function blocks for the Segway's model are removed from Fig. 10.27, we are left with just the controller blocks. The controller has three input signals: the desired speed, v^{des} , the measured speed, $V(s)$, and the measured orientation or lean angle, θ . It has one output signal, the motor torque, τ .

In Fig. 10.27, two of the transfer functions generate the Segway's orientation, $\theta(s)$, and cart speed, $V(s)$. These input signals to the controller will come from the nonlinear state-variable model of the Segway. On a real Segway or BallBot, they would come from sensors on the robots. When we remove model transfer functions from the block diagram, we are left with Fig. 10.29.

Where the feedback blocks required $\theta(s)$ and $V(s)$ from model blocks, we have moved them to the left of the diagram to show that they are now inputs to the controller along with the desired speed, $V^{des}(s)$. The block diagram in Fig. 10.29 is a representation of the following linear equation

$$\begin{aligned} \tau(s) &= C_2(s) F(s) C_1(s) (V^{des}(s) - V(s)) - C_1(s) \theta(s) \\ &= \begin{bmatrix} C_2(s) F(s) C_1(s) & -C_1(s) \end{bmatrix} \cdot \begin{bmatrix} V^{des}(s) - V(s) \\ \theta(s) \end{bmatrix}, \end{aligned} \quad (10.99)$$

where we have (temporarily) highlighted the input and output signals to make them easier to locate in the expressions.

Turning the Transfer Function Representation of the Controller into a State-variable Model

To merge the control calculation in (10.99) with the ODEs of the Segway model (linear or nonlinear), we need to transform the transfer function computations back into ODEs. This is possible because each s in (10.99) corresponds to a $\frac{d}{dt}$, thanks to the Laplace transform. While the conversion is not so hard to do by hand, it is better to do it in software. Here, we use the package `ControlSystems.jl`.

There is one sticky point: the PD controller $K_D \left(s + \frac{K_P}{K_D} \right)$ is not proper. The practical solution is to use a “dirty derivative” (yes, that is its name; you cannot make that up). The controller is rewritten as

$$K_D \left(\frac{s}{\tau_d s + 1} + \frac{K_P}{K_D} \right),$$

where τ_d is the time constant of the dirty derivative; in the code below, it is chosen to be 10 ms. With this modification, the PD controller is now proper and can be changed into an ODE.

```

1 using ControlSystems
2
3 # Build function to simplify process of going from Transfer Function
4 # Representation to State-variable model
5
6 function transferFunction2LinearODE(Sys_tf)
7     # Convert to state-variable representation
8     Sys_ss = ss(Sys_tf)
9     # dx/dt = A x + B u
10    # y = C x + D u

```

```

11    # Reduce to minimal and balanced state-space representation
12    # Explained in Michigan EECS 565 Linear Feedback Control Systems
13    #
14    tol = 1e-3 # for stable pole-zero cancellations
15    Sys_minimal = minreal(Sys_ss, tol)
16    # Improve numerical characteristics
17    Sys_ODE, ~, ~ = balreal(Sys_minimal) # the tildes are important
18
19    # Sys_ODE matrices
20    #     A = Sys_ODE.A
21    #     B = Sys_ODE.B
22    #     C = Sys_ODE.C
23    #     D = Sys_ODE.D
24
25
26    return Sys_ODE
27 end

```

Output

transferFunction2LinearODE (generic function with 1 method)

```

1 using ControlSystems
2
3
4 # Create a causal (proper) PD controller that admits a state variable model
5 tau_d = 1e-2 # Dirty derivative time constant
6 s = tf("s")
7 s_dirty = s / (tau_d * s + 1) # dirty derivative
8
9 # Controllers and Filter
10    # gains
11 K1 = -1.041 # Gain for C1(s)
12 K2 = 0.15 # Gain for C2(s)
13 K3 = 0.833 # gain for F(s)
14    # TFs
15 C1 = K1 * (s_dirty + 13.679)
16 C2 = K2 * tf([2.722], [1, 2.722])
17 F = K3 * tf([13.679], [1, 13.679])
18
19 # Open-loop multi-input transfer function from [V^{des}, V, theta] to tau
20 # Based on the diagram that isolates the feedback controller
21 controller_TF = [C2 * F * C1 -C1]
22
23 # Convert to state-variable representation
24 controller_ODE = transferFunction2LinearODE(controller_TF)
25 Ac = controller_ODE.A
26 Bc = controller_ODE.B
27 Cc = controller_ODE.C
28 Dc = controller_ODE.D
29
30 println("State-Variable Representation of the Controller:")
31 display(Ac)
32 display(Bc)
33 display(Cc)
34 display(Dc)

```

Output

State-Variable Representation of the Controller:

3×3 Matrix{Float64}:

```
-99.9969      -4.37006     1.09924
 -0.0695394    -2.91994     1.44874
 -0.00121461   1.44793    -13.4842
```

3×2 Matrix{Float64}:

```
 0.0517967  102.029
 -2.22881    0.073081
  0.560787   0.000209314
```

1×3 Matrix{Float64}:

```
-102.029  -2.23001  0.560787
```

1×2 Matrix{Float64}:

```
 0.0  118.34
```

Function to implement the linearized Segway model

```
1 # Linear Segway Model
2
3 function SegwayLinearModel(dx, x, tau, t)
4     A = [
5         0.0      0.0  1.0  0.0
6         0.0      0.0  0.0  1.0
7         50.9373 -0.0  0.0  0.0
8        -162.687 -0.0  0.0  0.0]
9     b=[ 
10        0.0
11        0.0
12        -21.467741053946913
13        80.22714380349518]
14     dx = A*x + b*tau
15     return dx
16 end
```

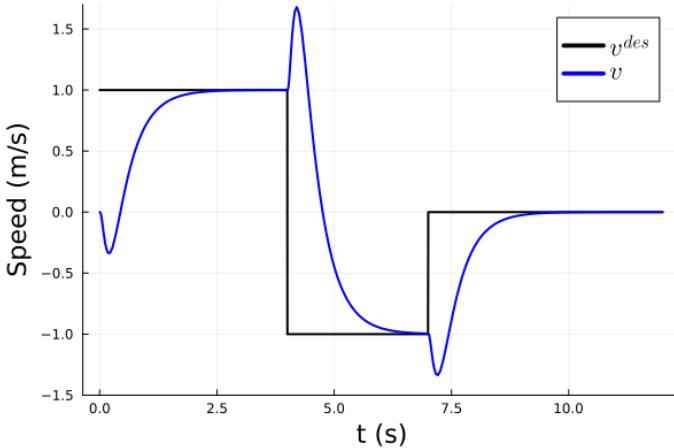
Output

```
SegwayLinearModel (generic function with 1 method)
```

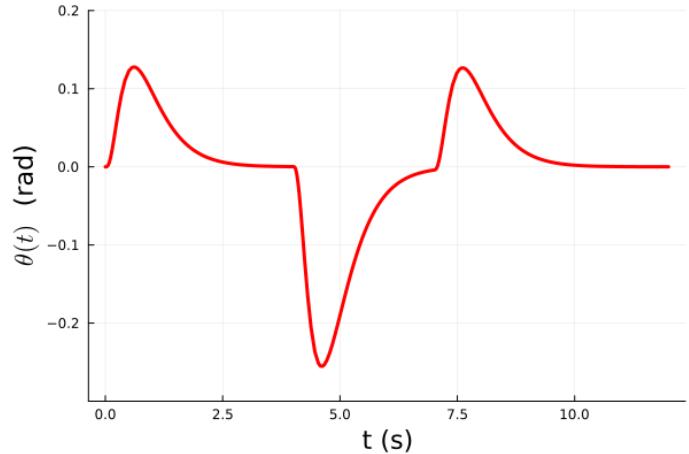
Put the Controller in Closed-loop with the Linear Model

Custom function to make it easy to form a closed-loop system

```
1 # Custom function to make it easy to form a closed-loop system
2
3 function simulate_control_system(plant_model, controller_ODE, K_IC, referenceCommand, xm0, xc0, tspan, p)
4     # Combined state-space representation
5     function combined_system(dx, x, p, t)
6         # Extract model and controller states
7         xm, xc = x[1:length(xm0)], x[(length(xm0)+1):end]    # Split the state variables
8
9         # Check if referenceCommand is a function or a constant
10        if isa(referenceCommand, Function)
11            # referenceCommand is a function
```



(a)



(b)

Figure 10.30: **Tracking a Speed Profile with the Linearized Segway Model:** (a) Shows a time-varying speed profile as the reference command and the closed-loop linear system's response, while (b) shows the corresponding lean angle trajectory.

```

12     Ref = referenceCommand(t)
13
14     else
15         # referenceCommand is a constant
16         # Define behavior when referenceCommand is a constant
17         Ref = referenceCommand
18
19     end
20
21     # Input to the controller
22     # Calculate uc from K_IC, uRef, and xm
23     uc = K_IC * [xm; Ref]
24     # uc[1] = max(min(uc[1], 0.5), -0.5)
25     # Output of the controller
26     tau = controller_ODE.C * xc + controller_ODE.D * uc
27     # tau is a 1 x 1 matrix. Must extract its value
28     # Plant dynamics
29     dxm = plant_model(dx[1:length(xm0)], xm, tau[1], t)
30     dx[1 : length(xm0)] = dxm
31     # Linear controller dynamics
32     dx[(length(xm0)+1):end] = controller_ODE.A * xc + controller_ODE.B * uc
33     # return for clarity
34     return dx
35
36 end
37
38 # Form full set of initial conditions
39 x0 = [xm0; xc0] # Combined initial state
40 # Form and solve the differential equations
41 prob = ODEProblem(combined_system, x0, tspan, p)
42 sol = solve(prob, Tsit5())
43 return sol
44
45 end

```

Output

simulate_control_system (generic function with 1 method)

Function to implement the linearized Segway model

```

1 # Linear Segway Model
2

```

```

3 function SegwayLinearModel(dx, x, tau, t)
4     A = [
5         0.0      0.0   1.0   0.0
6         0.0      0.0   0.0   1.0
7     50.9373 -0.0   0.0   0.0
8    -162.687 -0.0   0.0   0.0]
9     b=[
10        0.0
11        0.0
12    -21.467741053946913
13     80.22714380349518]
14 dx = A*x + b*tau
15 return dx
16 end

```

Output

SegwayLinearModel (generic function with 1 method)

Define a speed profile and simulate the closed-loop system

```

1 function vDes(t)
2 # Define a speed profile
3 if t > 7
4     vdes = 0.0
5 elseif t > 4
6     vdes = -1.0
7 else
8     vdes = 1.0
9 end
10 return vdes
11 end
12 using DifferentialEquations, Plots
13
14 function modelParameters(outSelect=0)
15     g = 9.81 # m/s^2
16     m_wh = 6 # kg mass of 2 wheels
17     m_pend = 45+75 # kg mass of the person, platform, and the handlebar
18     L = (45*.3 + 75*0.85)/(45+75) # m CoM for the person plus platform measured from
19     # the axels of the wheels
20     r_wh = 0.25 # m, Wheel radius
21     J_pend = m_pend*L^2/6 # kg m^2 moment of inertia about the center of mass
22     J_wh = 0.1 # kg m^2 moment of intertia of each wheel
23     J_rotor = 1e-4 # kg m^2 per rotor, one for each wheel
24     N = 50 # gear ratio between each motor rotor and wheel
25     if outSelect > 0
26         # named tuple
27         return (g=g, m_wh=m_wh, m_pend=m_pend, L=L, r_wh=r_wh, J_pend=J_pend,
28                 J_wh=J_wh, J_rotor=J_rotor, N=N)
29     else # standard tuple
30         return g, m_wh, m_pend, L, r_wh, J_pend, J_wh, J_rotor, N
31     end
32 end
33
34 params = modelParameters(1)
35 xm0 = zeros(4,1)
36 xc0 = zeros(3,1)

```

```

37
38 # Interconnection matrix between model states, reference, and the inputs to the
39 # controller. Vdes assumed appended to the bottom of the vmoel's state vector
40 K_IC = zeros(2,length(xm0)+1)
41 # Row 1 forms Vdes - v
42 K_IC[1,5] = 1.0;
43 K_IC[1,4] = -params.r_wh # v = dot_phi times wheel radius
44 # Row 2 picks off theta (orientation angle in radians)
45 K_IC[2,1] = 1.0
46
47 tspan = (0.0,12.0)
48 referenceCommand = 1.0
49 p=0.0
50 sol = simulate_control_system(SegwayLinearModel, contoller_ODE, K_IC, vDes, xm0, xc0, tspan, p)
51
52 # Plot everything as a first check before making nice plots
53 pSol = plot(sol)
54 display(pSol)
55
56 # Reshape sol to be easier to work with
57 function niceFormat(sol)
58     t = sol.t
59     states = sol.u
60     # Reshape each matrix to a row vector and stack them vertically
61     x = vcat([reshape(mat, 1, :) for mat in states]...)
62     return (t=t, x=x)
63 end
64
65 mySol=niceFormat(sol);
66
67 using Plots
68
69 t = mySol.t
70 x = mySol.x
71 theta = x[:,1]
72 phi = x[:,2]
73 dtheta = x[:,3]
74 dphi = x[:,4]
75
76 params = modelParameters(1)
77 # Level-1 Speed
78 v = params.r_wh*dphi
79
80 Ref = vDes.(t)
81
82 p1 = plot(t, Ref, lw=2, guidefont = 15, label=L"$v^{des}$", legendfontsize=15,
83 xlabel="t (s)", ylabel="Speed (m/s)", color=:black)
84 plot!(t, v, label=L"$v$ ", lw=2, color=:blue)
85
86 p2 = plot(t, theta, lw=3, label=false, guidefont = 15, legendfontsize=15,
87 xlabel="t (s)", ylabel=L"\theta(t) (rad)", color=:red)
88
89
90 display(p1)
91 display(p2)

```

Output The results are shown in Fig. 10.30.

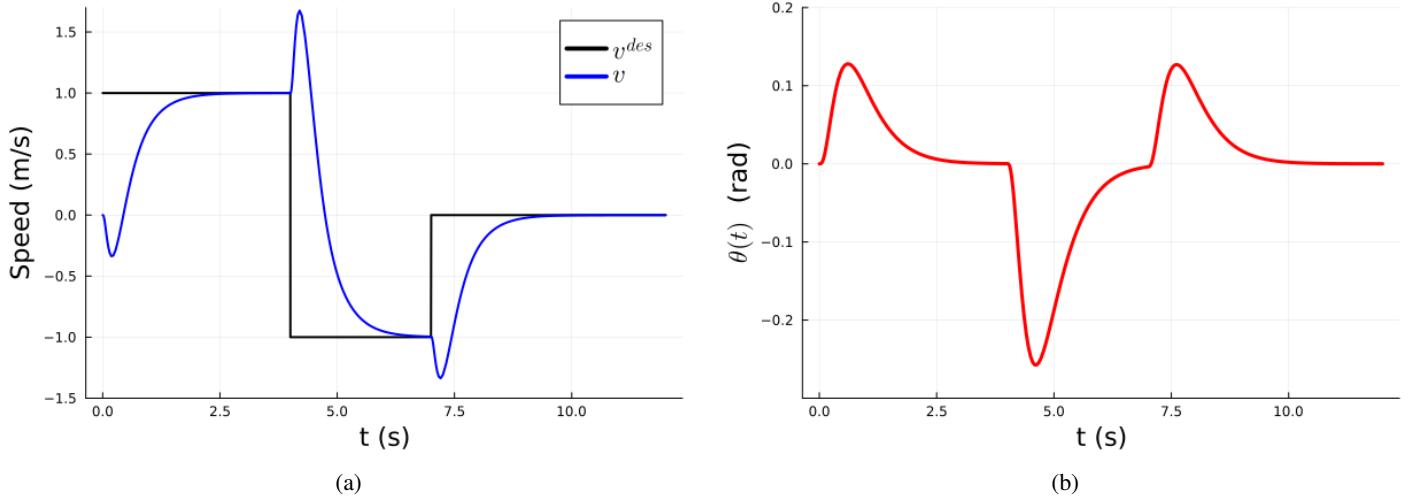


Figure 10.31: **Tracking a Speed Profile with the Nonlinear Segway Model:** (a) Shows a time-varying speed profile as the reference command and the closed-loop linear system's response, while (b) shows the corresponding lean angle trajectory. The results are indistinguishable from Fig. 10.30. Why is that? Because the controller maintains a small lean angle, the linear model is valid.

Put the Controller in Closed-loop with the Nonlinear Model

We simply replace the linearized plant model with the nonlinear Lagrangian model of the Segway and re-run the previous closed-loop simulation code.

```

1 # bring the model into the workspace for later use
2 include ("dyn_mod_planarSegway.jl")
3
4
5 function SegwayNonlinearModel(dx, x, tau, t)
6     n = floor(Int, length(x) / 2) # In Julia n/2 is a Float64
7     q = x[1:n]
8     dq = x[n+1:end]
9     model = dyn_mod_planarSegway(q, dq)
10    dx1 = dq
11    dx2 = (model.D) \ (-model.C*dq - model.G + model.B*tau) # note the use of backslash
12    dx=[dx1;dx2]
13    return dx
14 end

```

Output The results are shown in Fig. 10.31

Our final trick of the trade: We add one line to the function that simulates the closed-loop system

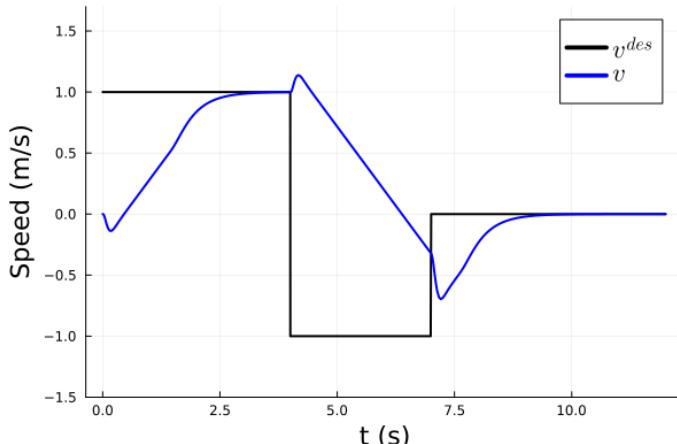
```

1 # Original line
2         uc = K_IC * [xm; Ref]
3 # Additional line
4         uc[1] = max(min(uc[1], 0.5), -0.5)

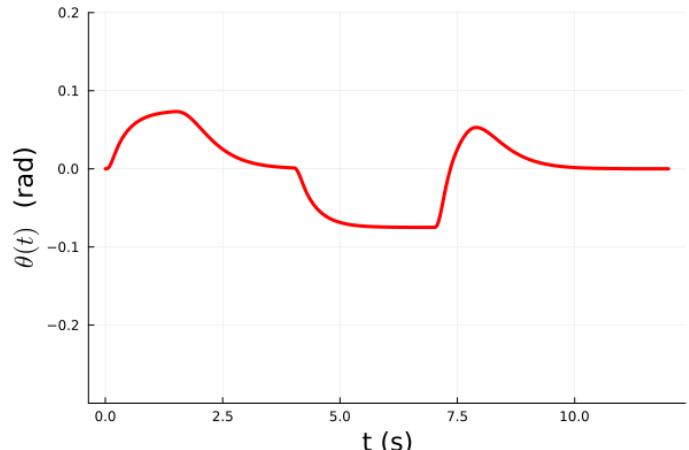
```

which implements the following saturation on the error in the speed. Hence, if the speed change is huge, the controller does not over-respond. The downside is that the system is less responsive for both acceleration and deceleration. The upside is that the rider may do fewer face plants. This idea can be very helpful for the BallBot.

$$uc[1] = \max(\min(uc[1], 0.5), -0.5) = \begin{cases} 0.5 & uc[1] > 0.5 \\ uc[1] & -0.5 \leq uc[1] \leq 0.5 \\ -0.5 & \text{otherwise.} \end{cases}$$



(a)



(b)

Figure 10.32: **Tracking a Speed Profile with the Nonlinear Segway Model after a Saturation Element was added to the Feedback Loop:** (a) Shows a time-varying speed profile as the reference command and the closed-loop nonlinear system's response, while (b) shows the corresponding lean angle trajectory. The saturation element limits the controller's response to errors in tracking the reference speed to ± 0.5 m/s. This makes the system less responsive on the one hand, while drastically reducing lean angle on the other. Which is better? Depends on the planned use of the Segway.

10.10 (Optional Read:) The Impulse “Function”

To do or not to do, that is the question.

10.11 (Optional Read:) Robot Equations for a Planar Segway

```

1 function dyn_mod_planarSegway(q, dq)
2 # DYN_MOD_PLANARSEGWAY
3 # 2023-12-20 17:49:53
4 #
5 # Author: Grizzle
6 #
7 # Model NOTATION: D(q)ddq + C(q,dq)*dq + G(q) = B*tau
8 # The Robot Equations: From Lagrange's Equations of Motion
9 #
10 g, m_wh, m_pend, L, r_wh, J_pend, J_wh, J_rotor, N = modelParameters()
11 #
12 # Variable names for the model
13 theta, phi = q
14 dtheta, dphi = dq
15 #
16 D = zeros(2, 2)
17 D[1, 1] = (J_pend*m_pend + J_pend*m_wh + (m_pend^2)*(r_wh^2) + (m_wh^2)*
18 (r_wh^2) + 2.0J_wh*m_pend + 2.0J_wh*m_wh + (L^2)*(m_pend^2)*
19 (cos(theta)^2) + (L^2)*(m_pend^2)*(sin(theta)^2) + 2.0m_pend*

```

```

20     m_wh*(r_wh^2) + 2.0L*r_wh*(m_pend^2)*cos(theta) + 2.0L*m_pend*
21     m_wh*r_wh*cos(theta)) / (m_pend + m_wh)
22 D[1, 2] = 2.0J_wh + m_pend*(r_wh^2) + m_wh*(r_wh^2) + L*m_pend*r_wh*
23     cos(theta)
24 D[2, 1] = 2.0J_wh + m_pend*(r_wh^2) + m_wh*(r_wh^2) + L*m_pend*r_wh*
25     cos(theta)
26 D[2, 2] = 2.0J_wh + m_pend*(r_wh^2) + m_wh*(r_wh^2) + 2.0J_rotor*(N^2)
27 #
28 C = zeros(2, 2)
29 C[1, 1] = (dtheta*(0.5(L^2)*(m_pend^2)*sin(2theta) - L*r_wh*(m_pend^2)*
30     sin(theta) - (L^2)*(m_pend^2)*cos(theta)*sin(theta) - L*m_pend*
31     m_wh*r_wh*sin(theta))) / (m_pend + m_wh)
32 C[2, 1] = -L*dtheta*m_pend*r_wh*sin(theta)
33 #
34 G = zeros(2)
35 G[1] = -L*g*m_pend*sin(theta)
36 G[2] = 0
37 #
38 B = zeros(2, 1)
39 B[2, 1] = 50
40 #
41 JacG = zeros(2, 2)
42 JacG[1, 1] = -L*g*m_pend*cos(theta)
43 #
44 JacComx = zeros(2)
45 JacComx[1] = (m_pend*(r_wh + L*cos(theta)) + m_wh*r_wh) / (m_pend + m_wh)
46 JacComx[2] = (m_pend*r_wh + m_wh*r_wh) / (m_pend + m_wh)
47 #
48 return (D=D, C=C, G=G, B=B, JacG=JacG, JacComx=JacComx)
49 end

```

10.12 (Optional Read:) Proofs Associated with the Chapter

We start with a remark on Prop. 10.9: it is possible for a right-sided function $f(t)$ to have a Laplace transform and a derivative $f'(t)$ that is continuous on $[0, \infty)$, but yet $f'(t)$ does not have a Laplace transform. We give a counterexample below. Importantly, $f(t) = \sin(e^{t^2})$, is not a function that occurs in engineering practice.

(a) **f has a Laplace transform:** For all $t \geq 0$ we have

$$|f(t)| = |\sin(e^{t^2})| \leq 1.$$

Therefore, for any $s > 0$ the Laplace transform

$$F(s) = \int_0^\infty f(t)e^{-st} dt$$

satisfies

$$|F(s)| \leq \int_0^\infty e^{-st} dt = \frac{1}{s} < \infty.$$

Hence, f is Laplace transformable.

(b) **f is continuously differentiable:** A straightforward application of the chain rule gives that

$$f'(t) = \cos(e^{t^2}) \cdot \frac{d}{dt}(e^{t^2}) = \cos(e^{t^2}) \cdot 2te^{t^2},$$

so that

$$f'(t) = 2te^{t^2} \cos(e^{t^2})$$

exists and is continuous on $[0, \infty)$.

- (c) **f' does not have a Laplace transform:** One has to show that the integral $\int_{0^-}^{\infty} f'(t)e^{-st} dt$ is divergent for all $s \in \mathbb{C}$. You can ask your favorite LLM to do this for you. The key issue is that the term e^{t^2} grows too rapidly and cannot be dominated by $e^{-\sigma t}$ for any $\sigma > 0$ finite. The oscillations due to $\cos(e^{t^2})$ pose technical challenges in forming the bounds, but these can be overcome with methods developed in a course on Real Analysis.

■

Proof of Prop. 10.14 (Partial Fraction Expansion)

A proper rational function $F(s)$ with distinct roots r_1, r_2, \dots, r_n can be expanded as

$$F(s) = \sum_{j=1}^n \frac{k_j}{s - r_j},$$

where the coefficients k_j are called the residues of $F(s)$ at the poles r_j . The residues k_j can be computed using

$$k_j = \lim_{s \rightarrow r_j} F(s)(s - r_j).$$

Remark: If r_j is complex, the associated residue k_j will generally be a complex number. Its complex conjugate k_j^* corresponds to the residue at the conjugate pole r_j^* .

Proof: Start with the partial fraction form

$$F(s) = \sum_{j=1}^n \frac{k_j}{s - r_j}.$$

Multiply both sides by $(s - r_i)$, where r_i is a distinct root

$$F(s)(s - r_i) = \sum_{j=1}^n k_j \frac{s - r_i}{s - r_j}.$$

At $s = r_i$, note that all terms in the summation vanish except the term where $j = i$, because for $j \neq i$, the numerator $(s - r_i)$ is zero, and hence

$$F(s)(s - r_i) \Big|_{s=r_i} = k_i.$$

It follows that

$$k_i = \lim_{s \rightarrow r_i} F(s)(s - r_i).$$

■

Proof of the Remark: Suppose $r_{j+1} = r_j^*$, where $r_j = a_j + i\omega_j$ is a complex pole of $F(s)$, which means that $\omega_j \neq 0$. We start by writing the denominator associated with the poles r_j and r_j^* as the real polynomial,

$$(s - r_j)(s - r_j^*) = (s - a_j)^2 + \omega_j^2.$$

Define $G(s)$ by:

$$F(s) = G(s) \cdot \frac{1}{(s - a_j)^2 + \omega_j^2} = G(s) \cdot \frac{1}{(s - r_j)(s - r_j^*)},$$

where $G(s)$ is a ratio of real polynomials because $F(s)$ is.

Step 1: Partial Fraction Expansion

Using a partial fraction expansion for the term $\frac{1}{(s - r_j)(s - r_j^*)}$, we can write:

$$F(s) = G(s) \cdot \left(\frac{k_j}{s - r_j} + \frac{k_{j+1}}{s - r_j^*} \right),$$

where k_j and k_{j+1} are the residues of $F(s)$ at r_j and r_j^* , respectively.

Step 2: Symmetry Property of $G(s)$

Because $F(s)$ has real coefficients, $G(s)$ also has real coefficients and satisfies the symmetry property:

$$G(s^*) = (G(s))^*.$$

Step 3: Relationship Between k_j and k_{j+1}

From the residue formula, we know,

$$k_j = \lim_{s \rightarrow r_j} F(s)(s - r_j), \quad k_{j+1} = \lim_{s \rightarrow r_j^*} F(s)(s - r_j^*).$$

Substitute the expression for $F(s)$ into the formula for k_j ,

$$k_j = \lim_{s \rightarrow r_j} G(s) \cdot \frac{1}{(s - r_j^*)} = G(r_j) \cdot \frac{1}{(r_j - r_j^*)} = G(r_j) \cdot \frac{1}{(2i\omega_j)}.$$

Similarly, substitute the expression for $F(s)$ into the formula for k_{j+1} ,

$$k_{j+1} = \lim_{s \rightarrow r_j^*} G(s) \cdot \frac{1}{(s - r_j)} = G(r_j^*) \cdot \frac{1}{(r_j^* - r_j)} = G(r_j^*) \cdot \frac{1}{(-2i\omega_j)}.$$

Step 3: Relationship Between k_j and k_{j+1}

Using the symmetry property of $G(s)$, we know,

$$G(r_j^*) = G(r_j)^*.$$

Substitute this into the expression for k_{j+1} ,

$$k_{j+1} = G(r_j^*) \cdot \frac{1}{(-2i\omega_j)} = G(r_j)^* \cdot \left(\frac{1}{(2i\omega_j)} \right)^*.$$

Simplifying gives

$$k_{j+1} = \left(G(r_j) \cdot \frac{1}{(2i\omega_j)} \right)^* = k_j^*.$$

Conclusion: The residue k_{j+1} at the pole r_j^* is the complex conjugate of the residue k_j at the pole r_j , as promised. ■

The purpose of partial fraction expansion is to express $F(s)$ as a sum of simple terms. This allows the inverse Laplace transform to be computed term-by-term, significantly simplifying the process of converting frequency-domain functions to their time-domain equivalents.

Proof of Prop. 10.27 (Poles and BIBO Stability)

A system with a proper rational transfer function is BIBO stable if, and only if, its poles have negative real parts. In other words, the stability of a rational transfer function

$$G(s) = \frac{N(s)}{D(s)},$$

is completely determined by its poles as long as the degree of the numerator polynomial is less than or equal to the degree of the denominator polynomial.

Proof: The proof assumes the poles are distinct, and the degree of the numerator polynomial is at least one less than the degree of the denominator polynomial. Moreover, we focus on showing that if there is a pole with a zero or positive real part, there is always a

bounded input that results in an unbounded output. The complete general proof is covered in EECS 216 Signals and Systems.

Without loss of generality, we assume that the denominator polynomial is monic and write $D(s) = (s - r_1)(s - r_2)\cdots(s - r_n)$ as a product of its roots. Its partial fraction expansion becomes, therefore,

$$G(s) = \sum_{j=1}^n \frac{k_j}{s - r_j},$$

where r_j are the poles of the system, and k_j are the “residues”, the official name of the coefficients of the partial fraction expansion! We make two observations:

- All of the terms k_j must be non-zero. This is because, otherwise, when the partial fraction expansion is recombined by placing the terms over a common denominator, the term corresponding to $(s - r_j)$ will not show up in the denominator, reducing its degree and eliminating r_j as a pole.
- If a pole is complex, its complex conjugate must also be a pole.

The response to a unit step input, $u_{\text{stp}}(t)$, is given by,

$$y(t) = \mathcal{L}^{-1} \left\{ G(s) \cdot \frac{1}{s} \right\}.$$

We analyze two cases: real poles and complex poles.

Case 1: Real poles. If r_j is real, the corresponding term in $G(s)$ is $\frac{k_j}{s - r_j}$. For a unit step input,

$$Y_j(s) = \frac{k_j}{s - r_j} \frac{1}{s}.$$

- If $r_j \neq 0$, the Laplace transform becomes

$$\begin{aligned} \frac{k_j}{s - r_j} \frac{1}{s} &= \alpha_0 \frac{1}{s} + \alpha_j \frac{1}{s - r_j} \quad (\text{by PFE}) \\ &= -\frac{k_j}{r_j} \frac{1}{s} + \frac{k_j}{r_j} \frac{1}{s - r_j} \quad (\text{solving for the coefficients}), \end{aligned}$$

and hence

$$y_j(t) = \frac{k_j}{r_j} (e^{r_j t} - 1) u_{\text{stp}}(t).$$

If $r_j > 0$, the exponential term $e^{r_j t}$ grows unbounded as $t \rightarrow \infty$, causing $y(t)$ to blow up. If $r_j < 0$, the exponential term $e^{r_j t}$ decays to zero as $t \rightarrow \infty$.

- If $r_j = 0$, the term $\frac{1}{s - r_j}$ becomes $\frac{1}{s}$, resulting in

$$y_j(t) = k_j \mathcal{L}^{-1} \left\{ \frac{1}{s^2} \right\} = k_j t u_{\text{stp}}(t),$$

a ramp function, which also grows unbounded as $t \rightarrow \infty$.

Hence, for stability, all real poles must satisfy $r_j < 0$.

Case 2: Complex poles. Let $r_j = a_j + i\omega_j$ with $\omega_j \neq 0$ (otherwise, it is not complex!), and assume r_j and its complex conjugate $r_{j+1} = r_j^* = a_j - i\omega_j$ appear sequentially. The corresponding terms in $G(s)$ are $\frac{k_j}{s - r_j} + \frac{k_j^*}{s - r_j^*}$, and for a unit-step input,

$$Y_j(s) = \left(\frac{k_j}{s - r_j} + \frac{k_j^*}{s - r_j^*} \right) \frac{1}{s}.$$

Therefore,

$$\begin{aligned} \left(\frac{k_j}{s - r_j} + \frac{k_j^*}{s - r_j^*} \right) \frac{1}{s} &= \alpha_0 \frac{1}{s} + \alpha_j \frac{1}{s - r_j} + \alpha_j^* \frac{1}{s - r_j^*} \quad (\text{by PFE}) \\ &= -2\operatorname{real}\left(\frac{k_j}{r_j}\right) \frac{1}{s} + \frac{k_j}{r_j} \frac{1}{s - r_j} + \frac{k_j^*}{r_j^*} \frac{1}{s - r_j^*} \quad (\text{solving for the coefficients}). \end{aligned}$$

After some algebra, the inverse Laplace transform becomes

$$y_j(t) = 2 \left[\operatorname{real}\left(\frac{k_j}{r_j}\right) e^{a_j t} \cos(\omega_j t) - \operatorname{imag}\left(\frac{k_j}{r_j}\right) e^{a_j t} \sin(\omega_j t) - \operatorname{real}\left(\frac{k_j}{r_j}\right) \right] u_{\text{stp}}(t).$$

- If $a_j > 0$, the exponential term $e^{a_j t}$ grows unbounded as $t \rightarrow \infty$, causing $y(t)$ to blow up.
- If $a_j < 0$, the exponential term $e^{a_j t}$ decays to zero as $t \rightarrow \infty$.
- If $a_j = 0$, we replace the unit step input with another bounded input, $\sin(\omega_j t)u_{\text{stp}}(t)$, where the frequency of the sinusoid, ω_j , is carefully chosen to match the natural frequency of the system, causing resonance.

If $a_j = 0$, the poles are purely imaginary and the terms in the partial fraction expansion become

$$\frac{k_j}{s - i\omega_j} + \frac{k_j^*}{s + i\omega_j}.$$

Combining these terms over the common denominator $(s^2 + \omega_j^2)$, the numerator becomes,

$$k_j(s + i\omega_j) + k_j^*(s - i\omega_j).$$

Expanding the numerator,

$$k_j(s + i\omega_j) + k_j^*(s - i\omega_j) = (k_j + k_j^*)s + i\omega_j(k_j - k_j^*).$$

$$k_j + k_j^* = 2\operatorname{real}(k_j), \quad k_j - k_j^* = 2i\operatorname{imag}(k_j).$$

Thus, the numerator simplifies to:

$$2\operatorname{real}(k_j)s + 2i\operatorname{imag}(k_j)\omega_j.$$

The transfer function term becomes:

$$\frac{k_j}{s - i\omega_j} + \frac{k_j^*}{s + i\omega_j} = \frac{2\operatorname{real}(k_j)s + 2i\operatorname{imag}(k_j)\omega_j}{s^2 + \omega_j^2}.$$

Now, consider the input $\sin(\omega_j t)u_{\text{stp}}(t)$, whose Laplace transform is:

$$\mathcal{L}\{\sin(\omega_j t)u_{\text{stp}}(t)\} = \frac{\omega_j}{s^2 + \omega_j^2}.$$

The output is:

$$y_j(t) = \mathcal{L}^{-1} \left\{ \frac{(2\operatorname{real}(k_j)s + 2i\operatorname{imag}(k_j)\omega_j)\omega_j}{(s^2 + \omega_j^2)^2} \right\}.$$

Split the numerator into two parts,

$$2\operatorname{real}(k_j)\omega_j \frac{s}{(s^2 + \omega_j^2)^2} \quad \text{and} \quad 2i\operatorname{imag}(k_j)\omega_j \frac{\omega_j}{(s^2 + \omega_j^2)^2}.$$

Using standard inverse Laplace transforms,

$$-\mathcal{L}^{-1} \left\{ \frac{s}{(s^2 + \omega_j^2)^2} \right\} = t \cos(\omega_j t),$$

$$-\mathcal{L}^{-1}\left\{\frac{\omega_j}{(s^2+\omega_j^2)^2}\right\} = t \sin(\omega_j t).$$

The output becomes,

$$\begin{aligned} y_j(t) &= 2 \operatorname{real}(k_j) \omega_j t \cos(\omega_j t) + 2 \operatorname{imag}(k_j) \omega_j t \sin(\omega_j t) \\ &= t (2 \operatorname{real}(k_j) \omega_j \cos(\omega_j t) + 2 \operatorname{imag}(k_j) \omega_j \sin(\omega_j t)) u_{\text{stp}}(t). \end{aligned}$$

This clearly shows that the output grows linearly with time, indicating resonance and instability.

Thus, for BIBO Stability, all complex poles must have $\operatorname{real}(r_j) = a_j < 0$.

Conclusion: For BIBO stability, the real parts of all poles must be strictly negative ($\operatorname{real}(r_j) < 0$). Otherwise, it is always possible to excite the system with a bounded input that yields an unbounded output. ■

Proof of Prop. 10.38 (Steady-State Error for Step Inputs Applied to BIBO Stable Closed-Loop Systems)

Assume the forward path transfer function is $G(s) = \frac{N(s)}{D(s)}$ and the closed-loop system is BIBO stable. Then the steady-state error for a unit-step command is well-defined and is given by

$$e_{SS} = \lim_{s \rightarrow 0} \frac{D(s)}{D(s) + N(s)}. \quad (10.100)$$

Moreover, if $G(s)$ does not have a pole at the origin, then

$$e_{SS} = \frac{D(0)}{D(0) + N(0)} \neq 0,$$

while if $G(s)$ has one or more poles at the origin, then $D(0) = 0$ and the steady-state error is

$$e_{SS} = 0.$$

Proof: We give the proof for the case that the roots of $D(s) + N(s)$ are distinct. From the given information, $\frac{E(s)}{R(s)} = \frac{D(s)}{D(s) + N(s)}$. If we assume that $r(t) = u_{\text{stp}}(t)$, then $R(s) = \frac{1}{s}$ and $E(s) = \frac{D(s)}{D(s) + N(s)} \frac{1}{s}$. Let r_1, r_2, \dots, r_n be the roots of $D(s) + N(s)$. Then the partial fraction expansion of $E(s)$ becomes

$$E(s) = \frac{k_0}{s} + \sum_{i=1}^n \frac{k_i}{s - r_i},$$

and from Prop. 10.14,

$$k_0 = \lim_{s \rightarrow 0} \frac{s D(s)}{D(s) + N(s)} \frac{1}{s} = \frac{D(0)}{D(0) + N(0)}.$$

Taking the inverse Laplace transform term-by-term gives

$$e(t) = k_0 u_{\text{stp}}(t) + \sum_{i=1}^n e^{r_i t} u_{\text{stp}}(t).$$

By BIBO stability of the closed-loop system, the real parts of the poles r_i are all negative. Hence,

$$\lim_{t \rightarrow \infty} e(t) = k_0 = \frac{D(0)}{D(0) + N(0)}$$

as claimed.

The poles of the open-loop transfer function, $\frac{Y(s)}{U(s)} = G(s) = \frac{N(s)}{D(s)}$, are the zeros of the closed-loop transfer function $\frac{E(s)}{R(s)} = \frac{D(s)}{D(s) + N(s)}$. Hence, $0 = k_0 = \frac{D(0)}{D(0) + N(0)}$ if, and only if, $D(0) = 0$. This completes the proof. ■