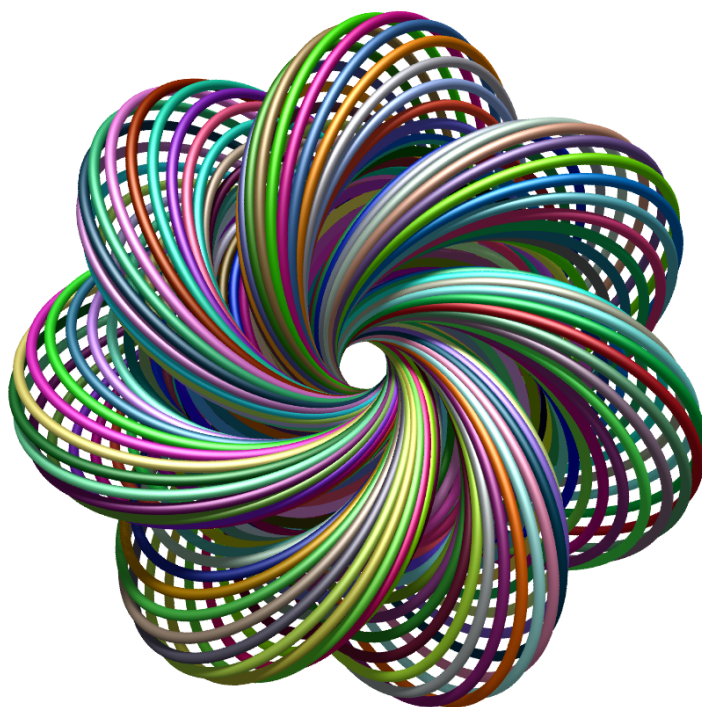# Essentials of Mathematical Methods:

## *Foundations, Principles, and Algorithms*

**Yuguang Yang**

version 3.0

*God used beautiful mathematics in creating the world. –Paul Dirac*

*Dedicated to*

*those who appreciate the power of mathematical methods*
*and enjoy learning it.*

# Preface

## Objective

Today, mathematical methods, models, and computational algorithms are playing increasingly significant roles in addressing major challenges arising from scientific research and technological developments. Although many novel methods and algorithms, such as deep learning and artificial intelligence, are emerging and reshaping various areas at an unprecedented pace, their core ideas and working mechanisms are inherently related to and deeply rooted in some essential mathematical foundations and principles. By performing an in-depth survey on the underlying foundations, principles, and algorithms, this book aims to navigate the vast landscape of mathematical methods widely used in diverse scientific and engineering domains.

This book starts with a survey of mathematical foundations, including essential concepts and theorems in real analysis, linear algebra, and related fundamentals. Then it examines a broad spectrum of applied mathematical methods, ranging from traditional ones such as optimizations and dynamical system modeling, to state-of-the-art such as machine learning, deep learning, and reinforcement learning. The emphasis is placed on methods for stochastic and dynamical system modeling, optimal decision-making, and statistical learning. For each topic, this book organizes fundamental definitions, theorems, methods, and algorithms in a logical and illuminating way.

## Features and Highlights

- Comprehensive, essential, and self-contained.
- Concepts, theorems, and discussions are developed to suit real-world applications.
- Key references and resources are provided on each topic.
- Comparisons and discussions on similar definitions and theorems.
- An evolving book with regular updates on Github.

## Acknowledgment

This book evolved from my study notes during my PhD studies at the Johns Hopkins University (JHU). I want to thank the following professors at JHU for their courses and valuable discussion: Daniel Robinson, Teresa Lebair, Andrea Prosperetti, Gregory Chirikjian, Michael Kahadan, James C. Spall, Marin Kobilarov, Suchi Saria, Michael Dimitz, Sean Sun, Ari Turner, Gregory Eyink, Amitabh Basu, John Miller and David Audley. I also want to thank Rachael Zhang for her editorial assistance.

Yuguang Yang, Fall 2019
yangyutu123@gmail.com

# Notations

- $\mathbb{R}$: real numbers.
- $\mathbb{R}_+$: nonnegative real numbers.
- $\mathbb{R}_{++}$: positive real numbers.
- $\bar{\mathbb{R}}$: extended real numbers.
- $\mathbb{C}$: complex numbers.
- $\mathbb{F}$: real or complex numbers.
- $\mathbb{Q}$: rational numbers.
- $\mathbb{Z}$: integer numbers.
- $\mathbb{P}$: positive numbers.
- $\mathcal{P}_n$: polynomial of degree of n.
- $\mathbb{N}$: natural numbers.
- $\mathcal{R}(A)$: the range of matrix $A$.
- $\mathcal{N}(A)$: the null space of matrix $A$.
- $V$: vector space.
- $det(A)$: the determinant of matrix $A$.
- $rank(A)$: the rank of matrix $A$.
- $\|\cdot\|_2$: Euclidean 2 norm of a vector of a matrix.
- $\|\cdot\|_F$: Frobenius norm of a matrix.
- $\rho(A)$: the spectral radius of matrix $A$.
- $Tr(A)$: the trace of matrix $A$.
- $L^2[a,b]$: Lebesgue integrable function on $[a,b]$.
- $L^1[a,b]$: Lebesgue integrable function on $[a,b]$.
- $N(0,1)$: standard Gaussian distribution.
- $N(\mu,\sigma^2)$: Gaussian distribution with mean $\mu$ and variance $\sigma^2$.
- $MN(\mu,\Sigma)$: multivariate Gaussian distribution with mean vector $\mu$ and covariance matrix $\Sigma$.
- $\mathbf{1}(x), I(x)$ indicator function.
- $E[X], \mathbf{E}[X], \mathbb{E}[X]$ expectation of random variable $X$.
- $Var[X]$ variance of random variable $X$.

# CONTENTS

# v statistical learning methods

# vi optimal control and reinforcement learning methods

# vii appendix

# LIST OF ALGORITHMS

# LIST OF FIGURES

# LIST OF TABLES

# CONTENTS

# ii mathematical optimization methods

## iii   classical statistical methods

# iv　dynamics modeling methods

## vi optimal control and reinforcement learning methods

## vii appendix

Part I

MATHEMATICAL FOUNDATIONS