Book draft

# Mathematical Foundations
# of
# Reinforcement Learning

Shiyu Zhao

March, 2024

# Contents