Book draft

# Mathematical Foundation
# of
# Reinforcement Learning

Shiyu Zhao

August 2022

# Contents