# Statistics and Data Visualization in Climate Science with R and Python

SAMUEL S. P. SHEN

San Diego State University

GERALD R. NORTH

Texas A&M University

# CAMBRIDGE
UNIVERSITY PRESS

# Contents

## Contents