# PRACTICAL DATA SCIENCE

## with

# Hadoop®

## and

# Spark

Designing
and Building
**Effective Analytics**
at Scale

OFER MENDELEVITCH
CASEY STELLA
DOUGLAS EADLINE

# Practical Data Science with Hadoop® and Spark

# Practical Data Science with Hadoop and Spark: Designing and Building Effective Analytics at Scale

# <u>Table of Contents</u>

Pearson

# Table of Contents

# Table of Contents

# Table of Contents

# Table of Contents

# Table of Contents

Pearson

# Table of Contents

# Table of Contents

# **Table of Contents**

Pearson