# THE SUPERVISED LEARNING WORKSHOP

## A NEW, INTERACTIVE APPROACH TO UNDERSTANDING SUPERVISED LEARNING ALGORITHMS

# The Supervised Learning Workshop

## Second Edition

A New, Interactive Approach to Understanding Supervised Learning Algorithms

Blaine Bateman, Ashish Ranjan Jha,
Benjamin Johnston, and Ishita Mathur

**Packt>**

**The Supervised Learning Workshop**

**Second Edition**

**Authors:** Blaine Bateman, Ashish Ranjan Jha, Benjamin Johnston, and Ishita Mathur

# Table of Contents

# Chapter 2: Exploratory Data Analysis and Visualization    39

# Chapter 6: Ensemble Modeling | 297

# Chapter 7: Model Evaluation | 333