# Demand Prediction in Retail

A Practical Guide to Leverage Data and Predictive Analytics

Maxime C. Cohen, Paul-Emile Gras, Arthur Pentecoste, Renyu Zhang