# Multivariate Data Analysis
## 6th Edition

An introduction to Multivariate Analysis, Process Analytical Technology and Quality by Design

Kim H. Esbensen

and

Brad Swarbrick

with contributions from Frank Westad, Pat Whitcombe and Mark Anderson

**CAMO**
Bring data to life

# Contents