# SUPERVISED MACHINE LEARNING

## OPTIMIZATION FRAMEWORK AND APPLICATIONS WITH SAS AND R

Tanya Kolosova and Samuel Berestizhevsky

# Supervised Machine Learning
## Optimization Framework and Applications with SAS and R

Tanya Kolosova PhD
Associates in Analytics Inc., Boca Raton, Florida

Samuel Berestizhevsky MSc
Associates in Analytics Inc., Boca Raton, Florida

# Contents

## Part I

## Part III