# Social Media Data Mining and Analytics

Gabor Szabo
Gungor Polatkan
Oscar Boykin
Antonios Chalkiopoulos

# Contents