# grokking
# Machine Learning

Luis G. Serrano

**Foreword by Sebastian Thrun**

# contents