# Contents

7   Machine Learning                                                                        143
*Rayid Ghani and Malte Schierholz*

## 8 Text Analysis

*Evgeny Klochikhin and Jordan Boyd-Graber*