# Contents

## III  Special topics     263

## 10 Ensemble learning     265

## 11 From online learning to bandits     291