# contents

**ix**

## 6 Diagnosing and tuning performance problems   194