*Agile Tools for Real-World Data*

# Python for Data Analysis

*Wes McKinney*

# Python for Data Analysis

*Wes McKinney*

**Python for Data Analysis**
by Wes McKinney

| | |
|---|---|
| **Editors:** Julie Steele and Meghan Blanchette | **Indexer:** BIM Publishing Services |
| **Production Editor:** Melanie Yarbrough | **Cover Designer:** Karen Montgomery |
| **Copyeditor:** Teresa Exley | **Interior Designer:** David Futato |
| **Proofreader:** BIM Publishing Services | **Illustrator:** Rebecca Demarest |

# Table of Contents