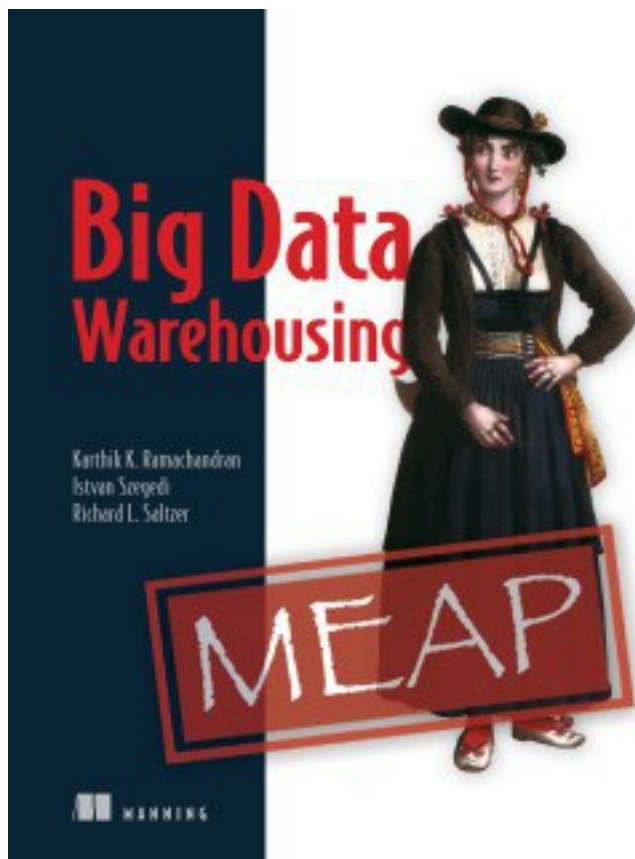# Big Data Warehousing

[Big Data Warehousing_下载链接1_](#)

著者:Karthik Ramachandran

出版者:

出版时间:2016-3-30

装帧:平装

isbn:9781633430280

Big Data Warehousing teaches you new techniques for common data warehousing tasks such as data ingest, SQL queries and report generation in a big data environment. You'll get a quick tour of using Hive and Impala to query and analyze large semi-structured datasets and learn how to build an Extract, Load, and Transform (ETL) workflow You'll explore data extraction with Sqoop and address the practical

question of schemas for modeling and transforming big data. As you progress through the book, you'll survey data governance with Falcon, how to build dataflows with Oozie, approaches to data processing, writing queries with SparkSQL, and data security using Apache Sentry and Knox.

作者介绍:

Karthik Ramachandran is a software engineer and Big Data expert who makes big data technologies and machine learning accessible to business users. He has extensive experience both with traditional enterprise data warehousing solutions as well as with the Hadoop ecosystem. Istvan Szegedi is a senior technical solutions architect working with enterprise data technologies and Hadoop. Richard Saltzer is a Software Engineer on Cloudera's internal data platform team where he builds scalable ingestion pipelines with Impala.

12. SPARK SQL
PART 4: OTHER CONSIDERATIONS
13. SECURITY
· · · · · · ([收起](#))

[Big Data Warehousing_下载链接1](#)

# 标签

hadoop

bigdata

# 评论

------------------------------
[Big Data Warehousing_下载链接1](#)

# 书评

------------------------------
[Big Data Warehousing_下载链接1](#)