# Hadoop In Action Chuck Lam

Summary Machine Learning in Action is
unique book that blends the
foundational theories of machine
learning with the practical realities
of building tools for everyday data
analysis. You'll use the flexible
Python programming language to build
programs that implement algorithms for
data classification, forecasting,
recommendations, and higher-level
features like summarization and
simplification. About the Book A
machine is said to learn when its
performance improves with experience.
Learning requires algorithms and
programs that capture data and ferret
out the interestingor useful patterns.
Once the specialized domain of analysts
and mathematicians, machine learning is
becoming a skill needed by many.
Machine Learning in Action is a clearly
written tutorial for developers. It
avoids academic language and takes you
straight to the techniques you'll use
in your day-to-day work. Many (Python)

examples present the core algorithms of statistical data processing, data analysis, and data visualization in code you can reuse. You'll understand the concepts and how they fit in with tactical tasks like classification, forecasting, recommendations, and higher-level features like summarization and simplification. Readers need no prior experience with machine learning or statistical processing. Familiarity with Python is helpful. Purchase of the print book comes with an offer of a free PDF, ePub, and Kindle eBook from Manning. Also available is all code from the book. What's Inside A no-nonsense introduction Examples showing common ML tasks Everyday data analysis Implementing classic algorithms like Apriori and Adaboos Table of Contents PART 1 CLASSIFICATION Machine learning basics Classifying with k-Nearest Neighbors Splitting datasets one feature at a time: decision trees Classifying with probability theory: naïve Bayes Logistic regression Support vector machines Improving classification with the AdaBoost meta

algorithm PART 2 FORECASTING NUMERIC VALUES WITH REGRESSION Predicting numeric values: regression Tree-based regression PART 3 UNSUPERVISED LEARNING Grouping unlabeled items using k-means clustering Association analysis with the Apriori algorithm Efficiently finding frequent itemsets with FP-growth PART 4 ADDITIONAL TOOLS Using principal component analysis to simplify data Simplifying data with the singular value decomposition Big data and MapReduce

Until now, design patterns for the MapReduce framework have been scattered among various research papers, blogs, and books. This handy guide brings together a unique collection of valuable MapReduce patterns that will save you time and effort regardless of the domain, language, or development framework you're using. Each pattern is explained in context, with pitfalls and caveats clearly identified to help you avoid common design mistakes when modeling your big data architecture. This book also provides a complete overview of MapReduce that explains its origins and implementations, and why

design patterns are so important. All code examples are written for Hadoop. Summarization patterns: get a top-level view by summarizing and grouping data Filtering patterns: view data subsets such as records generated from one user Data organization patterns: reorganize data to work with other systems, or to make MapReduce analysis easier Join patterns: analyze different datasets together to discover interesting relationships Metapatterns: piece together several patterns to solve multi-stage problems, or to perform several analytics in the same job Input and output patterns: customize the way you use Hadoop to load or store data "A clear exposition of MapReduce programs for common data processing patterns—this book is indespensible for anyone using Hadoop." --Tom White, author of Hadoop: The Definitive Guide Barack Rosenshine's Principles of Instruction are widely recognised for their clarity and simplicity and their potential to support teachers seeking to engage with cognitive science and the wider world of education research. In this concise new guide, Rosenshine

fan Tom Sherrington amplifies and augments the principles and further demonstrates how they can be put into practice in everyday classrooms.The second half of the book contain Rosenshine's original paper Principles of Instruction, as published in 2O1O by the International Academy of Education (IAE) - a paper with a superb worldwide reputation for relating research findings to classroom practice. Filled with practical, step-by-step instructions and clear explanations for the most important and useful tasks.This book provides quick recipes for using Hive to read data in various formats, efficiently querying this data, and extending Hive with any custom functions you may need to insert your own logic into the data pipeline.This book is written for data analysts and developers who want to use their current knowledge of SQL to be more productive with Hadoop. It assumes that readers are comfortable writing SQL queries and are familiar with Hadoop at the level of the classic WordCount example. Streaming Data

Essential techniques to help you
process, and get unique insights from,
big data, 2nd Edition
Mahout in Action
Apache Hive Cookbook
Building Effective Algorithms and
Analytics for Hadoop and Other Systems

*Summary Streaming Data introduces the concepts and
requirements of streaming and real-time data systems. The book is
an idea-rich tutorial that teaches you to think about how to
efficiently interact with fast-flowing data. Purchase of the print
book includes a free eBook in PDF, Kindle, and ePub formats
from Manning Publications. About the Technology As humans,
we're constantly filtering and deciphering the information
streaming toward us. In the same way, streaming data applications
can accomplish amazing tasks like reading live location data to
recommend nearby services, tracking faults with machinery in real
time, and sending digital receipts before your customers leave the
shop. Recent advances in streaming data technology and
techniques make it possible for any developer to build these
applications if they have the right mindset. This book will let you
join them. About the Book Streaming Data is an idea-rich tutorial
that teaches you to think about efficiently interacting with fast-
flowing data. Through relevant examples and illustrated use cases,
you'll explore designs for applications that read, analyze, share,
and store streaming data. Along the way, you'll discover the roles
of key technologies like Spark, Storm, Kafka, Flink, RabbitMQ,
and more. This book offers the perfect balance between big-picture
thinking and implementation details. What's Inside The right way
to collect real-time data Architecting a streaming pipeline
Analyzing the data Which technologies to use and when About the
Reader Written for developers familiar with relational database*

*concepts. No experience with streaming or real-time applications required. About the Author Andrew Psaltis is a software engineer focused on massively scalable real-time analytics. Table of Contents PART 1 - A NEW HOLISTIC APPROACH Introducing streaming data Getting data from clients: data ingestion Transporting the data from collection tier: decoupling the data pipeline Analyzing streaming data Algorithms for data analysis Storing the analyzed or collected data Making the data available Consumer device capabilities and limitations accessing the data PART 2 - TAKING IT REAL WORLD Analyzing Meetup RSVPs in real time*

*The massive datasets required for most modern businesses are too large to safely store and efficiently process on a single server. Hadoop is an open source data processing framework that provides a distributed file system that can manage data stored across clusters of servers and implements the MapReduce data processing model so that users can effectively query and utilize big data. The new Hadoop 2.0 is a stable, enterprise-ready platform supported by a rich ecosystem of tools and related technologies such as Pig, Hive, YARN, Spark, Tez, and many more. Hadoop in Action, Second Edition, provides a comprehensive introduction to Hadoop and shows how to write programs in the MapReduce style. It starts with a few easy examples and then moves quickly to show how Hadoop can be used in more complex data analysis tasks. It covers how YARN, new in Hadoop 2, simplifies and supercharges resource management to make streaming and real-time applications more feasible. Included are best practices and design patterns of MapReduce programming. The book expands on the first edition by enhancing coverage of important Hadoop 2 concepts and systems, and by providing new chapters on data management and data science that reinforce a practical understanding of Hadoop. Purchase of the print book includes a free eBook in PDF, Kindle, and ePub formats from Manning Publications.*

*You receive an e-mail. It contains an offer for a complete personal computer system. It seems like the retailer read your mind since you were exploring computers on their web site just a few hours prior…. As you drive to the store to buy the computer bundle, you get an offer for a discounted coffee from the coffee shop you are getting ready to drive past. It says that since you're in the area, you can get 10% off if you stop by in the next 20 minutes…. As you drink your coffee, you receive an apology from the manufacturer of a product that you complained about yesterday on your Facebook page, as well as on the company's web site…. Finally, once you get back home, you receive notice of a special armor upgrade available for purchase in your favorite online video game. It is just what is needed to get past some spots you've been struggling with…. Sound crazy? Are these things that can only happen in the distant future? No. All of these scenarios are possible today! Big data. Advanced analytics. Big data analytics. It seems you can't escape such terms today. Everywhere you turn people are discussing, writing about, and promoting big data and advanced analytics. Well, you can now add this book to the discussion. What is real and what is hype? Such attention can lead one to the suspicion that perhaps the analysis of big data is something that is more hype than substance. While there has been a lot of hype over the past few years, the reality is that we are in a transformative era in terms of analytic capabilities and the leveraging of massive amounts of data. If you take the time to cut through the sometimes-over-zealous hype present in the media, you'll find something very real and very powerful underneath it. With big data, the hype is driven by genuine excitement and anticipation of the business and consumer benefits that analyzing it will yield over time. Big data is the next wave of new data sources that will drive the next wave of analytic innovation in business, government, and academia. These innovations have the potential to radically change how organizations view their business. The analysis that big data enables will lead to decisions*

*that are more informed and, in some cases, different from what they are today. It will yield insights that many can only dream about today. As you'll see, there are many consistencies with the requirements to tame big data and what has always been needed to tame new data sources. However, the additional scale of big data necessitates utilizing the newest tools, technologies, methods, and processes. The old way of approaching analysis just won't work. It is time to evolve the world of advanced analytics to the next level. That's what this book is about. Taming the Big Data Tidal Wave isn't just the title of this book, but rather an activity that will determine which businesses win and which lose in the next decade. By preparing and taking the initiative, organizations can ride the big data tidal wave to success rather than being pummeled underneath the crushing surf. What do you need to know and how do you prepare in order to start taming big data and generating exciting new analytics from it? Sit back, get comfortable, and prepare to find out!*

*Hadoop in ActionSimon and Schuster*

*A Learner's Guide to Big Numbers, Statistics, and Good Decisions*

*Programming Hive*

*Spring Security in Action*

*The Definitive Guide*

*Taming The Big Data Tidal Wave*

If you are ready to dive into the MapReduce framework for processing large datasets, this practical book takes you step by step through the algorithms and tools you need to build distributed MapReduce applications with Apache Hadoop or Apache Spark. Each chapter provides a recipe for solving a massive computational problem, such as building a recommendation system. You'll learn how to implement the appropriate MapReduce solution with code that you can use in your projects. Dr. Mahmoud Parsian covers basic design patterns,

optimization techniques, and data mining and machine learning solutions for problems in bioinformatics, genomics, statistics, and social network analysis. This book also includes an overview of MapReduce, Hadoop, and Spark. Topics include: Market basket analysis for a large set of transactions Data mining algorithms (K-means, KNN, and Naive Bayes) Using huge genomic data to sequence DNA and RNA Naive Bayes theorem and Markov chains for data and market prediction Recommendation algorithms and pairwise document similarity Linear regression, Cox regression, and Pearson correlation Allelic frequency and mining DNA Social network analysis (recommendation systems, counting triangles, sentiment analysis)
Summary Mahout in Action is a hands-on introduction to machine learning with Apache Mahout. Following real-world examples, the book presents practical use cases and then illustrates how Mahout can be applied to solve them. Includes a free audio- and video-enhanced ebook. About the Technology A computer system that learns and adapts as it collects data can be really powerful. Mahout, Apache's open source machine learning project, captures the core algorithms of recommendation systems, classification, and clustering in ready-to-use, scalable libraries. With Mahout, you can immediately apply to your own projects the machine learning techniques that drive Amazon, Netflix, and others. About this Book This book covers machine learning using Apache Mahout. Based on experience with real-world applications, it introduces practical use cases and illustrates how Mahout can be applied to solve them. It

places particular focus on issues of scalability and how to apply these techniques against large data sets using the Apache Hadoop framework. This book is written for developers familiar with Java -- no prior experience with Mahout is assumed. Owners of a Manning pBook purchased anywhere in the world can download a free eBook from manning.com at any time. They can do so multiple times and in any or all formats available (PDF, ePub or Kindle). To do so, customers must register their printed copy on Manning's site by creating a user account and then following instructions printed on the pBook registration insert at the front of the book. What's Inside Use group data to make individual recommendations Find logical clusters within your data Filter and refine with on-the-fly classification Free audio and video extras Table of Contents Meet Apache Mahout PART 1 RECOMMENDATIONS Introducing recommenders Representing recommender data Making recommendations Taking recommenders to production Distributing recommendation computations PART 2 CLUSTERING Introduction to clustering Representing data Clustering algorithms in Mahout Evaluating and improving clustering quality Taking clustering to production Real-world applications of clustering PART 3 CLASSIFICATION Introduction to classification Training a classifier Evaluating and tuning a classifier Deploying a classifier Case study: Shop It To Me Individual self-contained code recipes. Solve specific problems using individual recipes, or work through the book to develop your capabilities. If you are a big data enthusiast and striving to use Hadoop to solve your

problems, this book is for you. Aimed at Java programmers with some knowledge of Hadoop MapReduce, this is also a comprehensive reference for developers and system admins who want to get up to speed using Hadoop.

Our world is being revolutionized by data-driven methods: access to large amounts of data has generated new insights and opened exciting new opportunities in commerce, science, and computing applications. Processing the enormous quantities of data necessary for these advances requires large clusters, making distributed computing paradigms more crucial than ever. MapReduce is a programming model for expressing distributed computations on massive datasets and an execution framework for large-scale data processing on clusters of commodity servers. The programming model provides an easy-to-understand abstraction for designing scalable algorithms, while the execution framework transparently handles many system-level details, ranging from scheduling to synchronization to fault tolerance. This book focuses on MapReduce algorithm design, with an emphasis on text processing algorithms common in natural language processing, information retrieval, and machine learning. We introduce the notion of MapReduce design patterns, which represent general reusable solutions to commonly occurring problems across a variety of problem domains. This book not only intends to help the reader "think in MapReduce", but also discusses limitations of the programming model as well. This volume is a printed version of a work that appears in

the Synthesis Digital Library of Engineering and
Computer Science. Synthesis Lectures provide concise,
original presentations of important research and
development topics, published quickly, in digital and
print formats. For more information visit
www.morganclaypool.com
Hello, Android
Hadoop in Practice
Hadoop: The Definitive Guide
Understanding Big Data: Analytics for Enterprise Class
Hadoop and Streaming Data
Machine Learning in Action

Big Data represents a new era in data
exploration and utilization, and IBM is
uniquely positioned to help clients navigate
this transformation. This book reveals how
IBM is leveraging open source Big Data
technology, infused with IBM technologies, to
deliver a robust, secure, highly available,
enterprise-class Big Data platform. The three
defining characteristics of Big Data--volume,
variety, and velocity--are discussed. You'll
get a primer on Hadoop and how IBM is
hardening it for the enterprise, and learn
when to leverage IBM InfoSphere BigInsights
(Big Data at rest) and IBM InfoSphere Streams
(Big Data in motion) technologies. Industry
use cases are also included in this practical
guide. Learn how IBM hardens Hadoop for
enterprise-class scalability and reliability
Gain insight into IBM's unique in-motion and
at-rest Big Data analytics platform Learn
tips and tricks for Big Data use cases and

solutions Get a quick Hadoop primer
Summary Modern data science solutions need to
be clean, easy to read, and scalable. In
Mastering Large Datasets with Python, author
J.T. Wolohan teaches you how to take a small
project and scale it up using a functionally
influenced approach to Python coding. You'll
explore methods and built-in Python tools
that lend themselves to clarity and
scalability, like the high-performing
parallelism method, as well as distributed
technologies that allow for high data
throughput. The abundant hands-on exercises
in this practical tutorial will lock in these
essential skills for any large-scale data
science project. Purchase of the print book
includes a free eBook in PDF, Kindle, and
ePub formats from Manning Publications. About
the technology Programming techniques that
work well on laptop-sized data can slow to a
crawl—or fail altogether—when applied to
massive files or distributed datasets. By
mastering the powerful map and reduce
paradigm, along with the Python-based tools
that support it, you can write data-centric
applications that scale efficiently without
requiring codebase rewrites as your
requirements change. About the book Mastering
Large Datasets with Python teaches you to
write code that can handle datasets of any
size. You'll start with laptop-sized datasets
that teach you to parallelize data analysis
by breaking large tasks into smaller ones
that can run simultaneously. You'll then

scale those same programs to industrial-sized datasets on a cluster of cloud servers. With the map and reduce paradigm firmly in place, you'll explore tools like Hadoop and PySpark to efficiently process massive distributed datasets, speed up decision-making with machine learning, and simplify your data storage with AWS S3. What's inside An introduction to the map and reduce paradigm Parallelization with the multiprocessing module and pathos framework Hadoop and Spark for distributed computing Running AWS jobs to process large datasets About the reader For Python programmers who need to work faster with more data. About the author J. T. Wolohan is a lead data scientist at Booz Allen Hamilton, and a PhD researcher at Indiana University, Bloomington. Table of Contents: PART 1 1 ¦ Introduction 2 ¦ Accelerating large dataset work: Map and parallel computing 3 ¦ Function pipelines for mapping complex transformations 4 ¦ Processing large datasets with lazy workflows 5 ¦ Accumulation operations with reduce 6 ¦ Speeding up map and reduce with advanced parallelization PART 2 7 ¦ Processing truly big datasets with Hadoop and Spark 8 ¦ Best practices for large data with Apache Streaming and mrjob 9 ¦ PageRank with map and reduce in PySpark 10 ¦ Faster decision-making with machine learning and PySpark PART 3 11 ¦ Large datasets in the cloud with Amazon Web Services and S3 12 ¦ MapReduce in the cloud with Amazon's Elastic MapReduce

Now in its second edition, this book focuses on practical algorithms for mining data from even the largest datasets.

Ready to unlock the power of your data? With this comprehensive guide, you'll learn how to build and maintain reliable, scalable, distributed systems with Apache Hadoop. This book is ideal for programmers looking to analyze datasets of any size, and for administrators who want to set up and run Hadoop clusters. You'll find illuminating case studies that demonstrate how Hadoop is used to solve specific problems. This third edition covers recent changes to Hadoop, including material on the new MapReduce API, as well as MapReduce 2 and its more flexible execution model (YARN). Store large datasets with the Hadoop Distributed File System (HDFS) Run distributed computations with MapReduce Use Hadoop's data and I/O building blocks for compression, data integrity, serialization (including Avro), and persistence Discover common pitfalls and advanced features for writing real-world MapReduce programs Design, build, and administer a dedicated Hadoop cluster—or run Hadoop in the cloud Load data from relational databases into HDFS, using Sqoop Perform large-scale data processing with the Pig query language Analyze datasets with Hive, Hadoop's data warehousing system Take advantage of HBase for structured and semi-structured data, and ZooKeeper for building distributed systems

*Summary HBase in Action has all the knowledge you need to design, build, and run applications using HBase. First, it introduces you to the fundamentals of distributed systems and large scale data handling. Then, you'll explore real-world applications and code samples with just enough theory to understand the practical techniques. You'll see how to build applications with HBase and take advantage of the MapReduce processing framework. And along the way you'll learn patterns and best practices. About the Technology HBase is a NoSQL storage system designed for fast, random access to large volumes of data. It runs on commodity hardware and scales smoothly from modest datasets to billions of rows and millions of columns. About this Book HBase in Action is an experience-driven guide that shows you how to design, build, and run applications using HBase. First, it introduces you to the fundamentals of handling big data. Then, you'll explore HBase with the help of real applications and code samples and with just enough theory to back up the*

*practical techniques. You'll take advantage of the MapReduce processing framework and benefit from seeing HBase best practices in action. Purchase of the print book comes with an offer of a free PDF, ePub, and Kindle eBook from Manning. Also available is all code from the book. What's Inside When and how to use HBase Practical examples Design patterns for scalable data systems Deployment, integration, and design Written for developers and architects familiar with data storage and processing. No prior knowledge of HBase, Hadoop, or MapReduce is required. Table of Contents PART 1 HBASE FUNDAMENTALS Introducing HBase Getting started Distributed HBase, HDFS, and MapReduce PART 2 ADVANCED CONCEPTS HBase table design Extending HBase with coprocessors Alternative HBase clients PART 3 EXAMPLE APPLICATIONS HBase by example: OpenTSDB Scaling GIS on HBase PART 4 OPERATIONALIZING HBASE Deploying HBase Operations*

*Describes the features and functions of Apache Hive, the data infrastructure for Hadoop.*

*Data is arriving faster than you can process it and the overall volumes keep growing at a rate that keeps you awake at night. Hadoop can help you tame the data*

*beast. Effective use of Hadoop however
requires a mixture of programming, design,
and system administration skills. "Hadoop
Beginner's Guide" removes the mystery from
Hadoop, presenting Hadoop and related
technologies with a focus on building
working systems and getting the job done,
using cloud services to do so when it
makes sense. From basic concepts and
initial setup through developing
applications and keeping the system
running as the data grows, the book gives
the understanding needed to effectively
use Hadoop to solve real world problems.
Starting with the basics of installing and
configuring Hadoop, the book explains how
to develop applications, maintain the
system, and how to use additional products
to integrate with other systems. While
learning different ways to develop
applications to run on Hadoop the book
also covers tools such as Hive, Sqoop, and
Flume that show how Hadoop can be
integrated with relational databases and
log collection. In addition to examples on
Hadoop clusters on Ubuntu uses of cloud
services such as Amazon, EC2 and Elastic
MapReduce are covered.
You've heard the hype about Hadoop: it
runs petabyte-scale data mining tasks
insanely fast, it runs gigantic tasks on*

*clouds for absurdly cheap, it's been heavily committed to by tech giants like IBM, Yahoo!, and the Apache Project, and it's completely open-source (thus free). But what exactly is it, and more importantly, how do you even get a Hadoop cluster up and running? From Apress, the name you've come to trust for hands-on technical knowledge, Pro Hadoop brings you up to speed on Hadoop. You learn the ins and outs of MapReduce; how to structure a cluster, design, and implement the Hadoop file system; and how to build your first cloud-computing tasks using Hadoop. Learn how to let Hadoop take care of distributing and parallelizing your software—you just focus on the code, Hadoop takes care of the rest. Best of all, you'll learn from a tech professional who's been in the Hadoop scene since day one. Written from the perspective of a principal engineer with down-in-the-trenches knowledge of what to do wrong with Hadoop, you learn how to avoid the common, expensive first errors that everyone makes with creating their own Hadoop system or inheriting someone else's. Skip the novice stage and the expensive, hard-to-fix mistakes...go straight to seasoned pro on the hottest cloud-computing framework with Pro Hadoop.*

*Your productivity will blow your managers away.*
*Big Data Computing*
*Understanding the real-time pipeline*
*The Complete Works of William Shakespeare*
*Hadoop in Action*
*Mastering Large Datasets with Python*
*Data Warehousing in the Age of the Big Data will help you and your organization make the most of unstructured data with your existing data warehouse. As Big Data continues to revolutionize how we use data, it doesn't have to create more confusion. Expert author Krish Krishnan helps you make sense of how Big Data fits into the world of data warehousing in clear and concise detail. The book is presented in three distinct parts. Part 1 discusses Big Data, its technologies and use cases from early adopters. Part 2 addresses data warehousing, its shortcomings, and new architecture options, workloads, and integration techniques for Big Data and the data warehouse. Part 3 deals with data governance, data visualization, information life-cycle management, data scientists, and implementing a Big Data–ready data warehouse. Extensive appendixes include case studies from vendor implementations and a special segment on how we can build a healthcare information factory. Ultimately, this book will help you navigate through the complex layers of Big Data and data warehousing while providing you information on how to effectively think about using all these technologies and the architectures to design the next-generation data warehouse. Learn how to leverage Big Data by effectively integrating it into your data warehouse.*

*Includes real-world examples and use cases that clearly demonstrate Hadoop, NoSQL, HBASE, Hive, and other Big Data technologies Understand how to optimize and tune your current data warehouse infrastructure and integrate newer infrastructure matching data processing workloads and requirements*

*Discover how Apache Hadoop can unleash the power of your data. This comprehensive resource shows you how to build and maintain reliable, scalable, distributed systems with the Hadoop framework -- an open source implementation of MapReduce, the algorithm on which Google built its empire. Programmers will find details for analyzing datasets of any size, and administrators will learn how to set up and run Hadoop clusters. This revised edition covers recent changes to Hadoop, including new features such as Hive, Sqoop, and Avro. It also provides illuminating case studies that illustrate how Hadoop is used to solve specific problems. Looking to get the most out of your data? This is your book. Use the Hadoop Distributed File System (HDFS) for storing large datasets, then run distributed computations over those datasets with MapReduce Become familiar with Hadoop's data and I/O building blocks for compression, data integrity, serialization, and persistence Discover common pitfalls and advanced features for writing real-world MapReduce programs Design, build, and administer a dedicated Hadoop cluster, or run Hadoop in the cloud Use Pig, a high-level query language for large-scale data processing Analyze datasets with Hive, Hadoop's data warehousing system Take advantage of HBase, Hadoop's database for structured and semi-structured*

*data Learn ZooKeeper, a toolkit of coordination primitives for building distributed systems "Now you have the opportunity to learn about Hadoop from a master -- not only of the technology, but also of common sense and plain talk." --Doug Cutting, Cloudera*

*The go-to guidebook for deploying Big Data solutions withHadoop Today's enterprise architects need to understand how the Hadoopframeworks and APIs fit together, and how they can be integrated todeliver real-world solutions. This book is a practical, detailedguide to building and implementing those solutions, with code-levelinstruction in the popular Wrox tradition. It covers storing datawith HDFS and Hbase, processing data with MapReduce, and automatingdata processing with Oozie. Hadoop security, running Hadoop withAmazon Web Services, best practices, and automating Hadoopprocesses in real time are also covered in depth. With in-depth code examples in Java and XML and the latest onrecent additions to the Hadoop ecosystem, this complete resourcealso covers the use of APIs, exposing their inner workings andallowing architects and developers to better leverage and customizethem. The ultimate guide for developers, designers, and architectswho need to build and deploy Hadoop applications Covers storing and processing data with various technologies,automating data processing, Hadoop security, and deliveringreal-time solutions Includes detailed, real-world examples and code-levelguidelines Explains when, why, and how to use these tools effectively Written by a team of Hadoop experts in theprogrammer-to-programmer Wrox style Professional Hadoop Solutions is the*

*reference enterprisearchitects and developers need to maximize the power of Hadoop.*

*Special Features: · Introduction to MapReduce· Examples illustrating ideas in practice· Hadoop's Streaming API· Other related tools, like Pig and Hive About The Book: Hadoop in Action introduces the subject and teaches you how to write programs in the MapReduce style. It starts with a few easy examples and then moves quickly to show Hadoop use in more complex data analysis tasks. Included are best practices and design patterns of MapReduce programming.This book requires basic Java skills. Knowing basic statistical concepts can help with the more advanced examples.*

*Distributed Data at Web Scale*

*Pro Hadoop*

*Professional Hadoop Solutions*

*Cassandra: The Definitive Guide*

*HADOOP IN ACTION*

**This book takes you on a fantastic journey to discover the attributes of big data using Apache Hive. Key Features Grasp the skills needed to write efficient Hive queries to analyze the Big Data Discover how Hive can coexist and work with other tools within the Hadoop ecosystem Uses practical, example-oriented scenarios to cover all the newly released features of Apache Hive 2.3.3 Book Description In this book, we prepare you for your journey into big data by frstly introducing you to backgrounds in the big data domain, alongwith the process of setting up and getting familiar with your Hive working environment. Next, the book guides you through discovering and transforming the values of big**

**data with the help of examples. It also hones your skills in using the Hive language in an effcient manner. Toward the end, the book focuses on advanced topics, such as performance, security, and extensions in Hive, which will guide you on exciting adventures on this worthwhile big data journey. By the end of the book, you will be familiar with Hive and able to work effeciently to find solutions to big data problems What you will learn Create and set up the Hive environment Discover how to use Hive's definition language to describe data Discover interesting data by joining and filtering datasets in Hive Transform data by using Hive sorting, ordering, and functions Aggregate and sample data in different ways Boost Hive query performance and enhance data security in Hive Customize Hive to your needs by using user-defined functions and integrate it with other tools Who this book is for If you are a data analyst, developer, or simply someone who wants to quickly get started with Hive to explore and analyze Big Data in Hadoop, this is the book for you. Since Hive is an SQL-like language, some previous experience with SQL will be useful to get the most out of this book.
Summary Hadoop in Practice, Second Edition provides over 100 tested, instantly useful techniques that will help you conquer big data, using Hadoop. This revised new edition covers changes and new features in the Hadoop core architecture, including MapReduce 2. Brand new chapters cover YARN and integrating**

**Kafka, Impala, and Spark SQL with Hadoop. You'll also get new and updated techniques for Flume, Sqoop, and Mahout, all of which have seen major new versions recently. In short, this is the most practical, up-to-date coverage of Hadoop available anywhere. Purchase of the print book includes a free eBook in PDF, Kindle, and ePub formats from Manning Publications. About the Book It's always a good time to upgrade your Hadoop skills! Hadoop in Practice, Second Edition provides a collection of 104 tested, instantly useful techniques for analyzing real-time streams, moving data securely, machine learning, managing large-scale clusters, and taming big data using Hadoop. This completely revised edition covers changes and new features in Hadoop core, including MapReduce 2 and YARN. You'll pick up hands-on best practices for integrating Spark, Kafka, and Impala with Hadoop, and get new and updated techniques for the latest versions of Flume, Sqoop, and Mahout. In short, this is the most practical, up-to-date coverage of Hadoop available. Readers need to know a programming language like Java and have basic familiarity with Hadoop. What's Inside Thoroughly updated for Hadoop 2 How to write YARN applications Integrate real-time technologies like Storm, Impala, and Spark Predictive analytics using Mahout and RR Readers need to know a programming language like Java and have basic familiarity with Hadoop. About the Author Alex Holmes works on tough big-data problems. He is a software engineer, author, speaker, and**

**blogger specializing in large-scale Hadoop projects. Table of Contents PART 1 BACKGROUND AND FUNDAMENTALS Hadoop in a heartbeat Introduction to YARN PART 2 DATA LOGISTICS Data serialization—working with text and beyond Organizing and optimizing data in HDFS Moving data into and out of Hadoop PART 3 BIG DATA PATTERNS Applying MapReduce patterns to big data Utilizing data structures and algorithms at scale Tuning, debugging, and testing PART 4 BEYOND MAPREDUCE SQL on Hadoop Writing a YARN application**

**Spring Security in Action shows you how to prevent cross-site scripting and request forgery attacks before they do damage. You'll start with the basics, simulating password upgrades and adding multiple types of authorization. As your skills grow, you'll adapt Spring Security to new architectures and create advanced OAuth2 configurations. By the time you're done, you'll have a customized Spring Security configuration that protects against threats both common and extraordinary. Summary While creating secure applications is critically important, it can also be tedious and time-consuming to stitch together the required collection of tools. For Java developers, the powerful Spring Security framework makes it easy for you to bake security into your software from the very beginning. Filled with code samples and practical examples, Spring Security in Action teaches you how to secure your apps from the most common threats, ranging from injection**

**attacks to lackluster monitoring. In it, you'll learn how to manage system users, configure secure endpoints, and use OAuth2 and OpenID Connect for authentication and authorization. Purchase of the print book includes a free eBook in PDF, Kindle, and ePub formats from Manning Publications. About the technology Security is non-negotiable. You rely on Spring applications to transmit data, verify credentials, and prevent attacks. Adopting "secure by design" principles will protect your network from data theft and unauthorized intrusions. About the book Spring Security in Action shows you how to prevent cross-site scripting and request forgery attacks before they do damage. You'll start with the basics, simulating password upgrades and adding multiple types of authorization. As your skills grow, you'll adapt Spring Security to new architectures and create advanced OAuth2 configurations. By the time you're done, you'll have a customized Spring Security configuration that protects against threats both common and extraordinary. What's inside Encoding passwords and authenticating users Securing endpoints Automating security testing Setting up a standalone authorization server About the reader For experienced Java and Spring developers. About the author Laurentiu Spilca is a dedicated development lead and trainer at Endava, with over ten years of Java experience. Table of Contents PART 1 - FIRST STEPS 1 Security Today 2 Hello Spring Security PART 2 - IMPLEMENTATION 3 Managing users 4 Dealing with passwords 5**

**Implementing authentication 6 Hands-on: A small secured web application 7 Configuring authorization: Restricting access 8 Configuring authorization: Applying restrictions 9 Implementing filters 10 Applying CSRF protection and CORS 11 Hands-on: A separation of responsibilities 12 How does OAuth 2 work? 13 OAuth 2: Implementing the authorization server 14 OAuth 2: Implementing the resource server 15 OAuth 2: Using JWT and cryptographic signatures 16 Global method security: Pre- and postauthorizations 17 Global method security: Pre- and postfiltering 18 Hands-on: An OAuth 2 application 19 Spring Security for reactive apps 20 Spring Security testing**

**If you've been asked to maintain large and complex Hadoop clusters, this book is a must. Demand for operations-specific material has skyrocketed now that Hadoop is becoming the de facto standard for truly large-scale data processing in the data center. Eric Sammer, Principal Solution Architect at Cloudera, shows you the particulars of running Hadoop in production, from planning, installing, and configuring the system to providing ongoing maintenance. Rather than run through all possible scenarios, this pragmatic operations guide calls out what works, as demonstrated in critical deployments. Get a high-level overview of HDFS and MapReduce: why they exist and how they work Plan a Hadoop deployment, from hardware and OS selection to network requirements Learn setup and configuration**

**details with a list of critical properties Manage resources by sharing a cluster across multiple groups Get a runbook of the most common cluster maintenance tasks Monitor Hadoop clusters—and learn troubleshooting with the help of real-world war stories Use basic tools and techniques to handle backup and catastrophic failure**
**Instant Apache Hive Essentials How-to**
**Hadoop MapReduce Cookbook**
**Finding Opportunities in Huge Data Streams with Advanced Analytics**
**R For Dummies**
**Library Resources & Technical Services**

A guide for data managers and analyzers shares guidelines for identifying patterns, predicting future outcomes, and presenting findings to others; drawing on current research in cognitive science and learning theory while covering such additional topics as assessing data quality, handling ambiguous information, and organizing data within market groups. Original.

Google Android dominates the mobile market, and by targeting Android, your apps can run on most of the phones and tablets in the world. This new fourth edition of the #1 book for learning Android covers all modern Android versions from Android 4.1 through Android 5.0. Freshly added material covers new Android features such as Fragments and

Google Play Services. Android is a platform you can't afford not to learn, and this book gets you started. Android is a software toolkit for mobile phones and tablets, created by Google. It's inside more than a billion devices, making Android the number one platform for application developers. Your own app could be running on all those devices! Getting started developing with Android is easy. You don't even need access to an Android phone, just a computer where you can install the Android SDK and the emulator that comes with it. Within minutes, Hello, Android gets you creating your first working application: Android's version of "Hello, World." From there, you'll build up a more substantial example: an Ultimate Tic-Tac-Toe game. By gradually adding features to the game, you'll learn about many aspects of Android programming, such as creating animated user interfaces, playing music and sound effects, building location-based services (including GPS and cell-tower triangulation), and accessing web services. You'll also learn how to publish your applications to the Google Play Store. This fourth edition of the bestselling Android classic has been revised for Android 4.1-4.3 (Jelly Bean), 4.4 (KitKat), and Android 5.0 (Lollipop). Topics have been streamlined

and simplified based on reader feedback, and every page and example has been reviewed and updated for compatibility with the latest versions of Android. If you'd rather be coding than reading about coding, this book is for you.

Hadoop in Action teaches readers how to use Hadoop and write MapReduce programs. The intended readers are programmers, architects, and project managers who have to process large amounts of data offline. Hadoop in Action will lead the reader from obtaining a copy of Hadoop to setting it up in a cluster and writing data analytic programs. The book begins by making the basic idea of Hadoop and MapReduce easier to grasp by applying the default Hadoop installation to a few easy-to-follow tasks, such as analyzing changes in word frequency across a body of documents. The book continues through the basic concepts of MapReduce applications developed using Hadoop, including a close look at framework components, use of Hadoop for a variety of data analysis tasks, and numerous examples of Hadoop in action. Hadoop in Action will explain how to use Hadoop and present design patterns and practices of programming MapReduce. MapReduce is a complex idea both

conceptually and in its implementation, and Hadoop users are challenged to learn all the knobs and levers for running Hadoop. This book takes you beyond the mechanics of running Hadoop, teaching you to write meaningful programs in a MapReduce framework. This book assumes the reader will have a basic familiarity with Java, as most code examples will be written in Java. Familiarity with basic statistical concepts (e.g. histogram, correlation) will help the reader appreciate the more advanced data processing examples. Purchase of the print book comes with an offer of a free PDF, ePub, and Kindle eBook from Manning. Also available is all code from the book. This book is an example-based tutorial that deals with Optimizing Hadoop for MapReduce job performance. If you are a Hadoop administrator, developer, MapReduce user, or beginner, this book is the best choice available if you wish to optimize your clusters and applications. Having prior knowledge of creating MapReduce applications is not necessary, but will help you better understand the concepts and snippets of MapReduce class template code. Recipes for Scaling Up with Hadoop and Spark Optimizing Hadoop for MapReduce

MapReduce Design Patterns
Rosenshine's Principles in Action
Introducing Google's Mobile Development
Platform

**If your organization is looking for a storage solution to accommodate a virtually endless amount of data, this book will show you how Apache HBase can fulfill your needs. As the open source implementation of Google's BigTable architecture, HBase scales to billions of rows and millions of columns, while ensuring that write and read performance remain constant.HBase: The Definitive Guideprovides the details you require, whether you simply want to evaluate this high-performance, non-relational database, or put it into practice right away. HBase's adoption rate is beginning to climb, and several IT executives are asking pointed questions about this high-capacity database. This is the only book available to give you meaningful answers. Learn how to distribute large datasets across an inexpensive cluster of commodity servers Develop HBase clients in many programming languages, including Java, Python, and Ruby Get details on HBase's primary storage system, HDFS—Hadoop's distributed and replicated filesystem Learn how HBase's native interface to Hadoop's MapReduce framework enables easy development and execution of batch jobs that can scan entire tables Discover the integration between HBase and other facets of the Apache Hadoop project**
**Due to market forces and technological evolution,**

**Big Data computing is developing at an increasing rate. A wide variety of novel approaches and tools have emerged to tackle the challenges of Big Data, creating both more opportunities and more challenges for students and professionals in the field of data computation and analysis. Presenting a mix of industry cases and theory, Big Data Computing discusses the technical and practical issues related to Big Data in intelligent information management. Emphasizing the adoption and diffusion of Big Data tools and technologies in industry, the book introduces a broad range of Big Data concepts, tools, and techniques. It covers a wide range of research, and provides comparisons between state-of-the-art approaches. Comprised of five sections, the book focuses on: What Big Data is and why it is important Semantic technologies Tools and methods Business and economic perspectives Big Data applications across industries Master the programming language of choice among statisticians and data analysts worldwide Coming to grips with R can be tough, even for seasoned statisticians and data analysts. Enter R For Dummies, the quick, easy way to master all the R you'll ever need. Requiring no prior programming experience and packed with practical examples, easy, step-by-step exercises, and sample code, this extremely accessible guide is the ideal introduction to R for complete beginners. It also covers many concepts that intermediate-level programmers will find extremely useful. Master your R ABCs ? get up**

**to speed in no time with the basics, from installing and configuring R to writing simple scripts and performing simultaneous calculations on many variables Put data in its place ? get to know your way around lists, data frames, and other R data structures while learning to interact with other programs, such as Microsoft Excel Make data dance to your tune ? learn how to reshape and manipulate data, merge data sets, split and combine data, perform calculations on vectors and arrays, and much more Visualize it ? learn to use R's powerful data visualization features to create beautiful and informative graphical presentations of your data Get statistical ? find out how to do simple statistical analysis, summarize your variables, and conduct classic statistical tests, such as t-tests Expand and customize R ? get the lowdown on how to find, install, and make the most of add-on packages created by the global R community for a wide variety of purposes Open the book and find: Help downloading, installing, and configuring R Tips for getting data in and out of R Ways to use data frames and lists to organize data How to manipulate and process data Advice on fitting regression models and ANOVA Helpful hints for working with graphics How to code in R What R mailing lists and forums can do for you**
**Imagine what you could do if scalability wasn't a problem. With this hands-on guide, you'll learn how the Cassandra database management system handles hundreds of terabytes of data while**

**remaining highly available across multiple data centers. This expanded second edition—updated for Cassandra 3.0—provides the technical details and practical examples you need to put this database to work in a production environment. Authors Jeff Carpenter and Eben Hewitt demonstrate the advantages of Cassandra's non-relational design, with special attention to data modeling. If you're a developer, DBA, or application architect looking to solve a database scaling issue or future-proof your application, this guide helps you harness Cassandra's speed and flexibility. Understand Cassandra's distributed and decentralized structure Use the Cassandra Query Language (CQL) and cqlsh—the CQL shell Create a working data model and compare it with an equivalent relational model Develop sample applications using client drivers for languages including Java, Python, and Node.js Explore cluster topology and learn how nodes exchange data Maintain a high level of performance in your cluster Deploy Cassandra on site, in the Cloud, or with Docker Integrate Cassandra with Spark, Hadoop, Elasticsearch, Solr, and Lucene Data Algorithms HBase in Action**

**Hadoop Beginner's Guide**
**Hadoop Operations**

Easy, hands-on recipes to help you understand Hive and its integration with frameworks that are used widely in today's big data world About This Book Grasp a

complete reference of different Hive topics. Get to know the latest recipes in development in Hive including CRUD operations Understand Hive internals and integration of Hive with different frameworks used in today's world. Who This Book Is For The book is intended for those who want to start in Hive or who have basic understanding of Hive framework. Prior knowledge of basic SQL command is also required What You Will Learn Learn different features and offering on the latest Hive Understand the working and structure of the Hive internals Get an insight on the latest development in Hive framework Grasp the concepts of Hive Data Model Master the key concepts like Partition, Buckets and Statistics Know how to integrate Hive with other frameworks such as Spark, Accumulo, etc In Detail Hive was developed by Facebook and later open sourced in Apache community. Hive provides SQL like interface to run queries on Big Data frameworks. Hive provides SQL like syntax also called as HiveQL that includes all SQL capabilities like analytical functions which are the need of the hour in today's Big Data world. This book provides you easy installation steps with different types of metastores supported by Hive. This book has simple and easy to learn recipes for configuring Hive clients and services. You would also learn different Hive optimizations including Partitions and Bucketing. The book also covers the source code explanation of latest Hive version. Hive Query Language is being used by other frameworks including spark. Towards the end you will cover integration of Hive with these frameworks. Style and approach Starting with the basics and covering

the core concepts with the practical usage, this book is a complete guide to learn and explore Hive offerings. Let Hadoop For Dummies help harness the power of your data and rein in the information overload Big data has become big business, and companies and organizations of all sizes are struggling to find ways to retrieve valuable information from their massive data sets with becoming overwhelmed. Enter Hadoop and this easy-to-understand For Dummies guide. Hadoop For Dummies helps readers understand the value of big data, make a business case for using Hadoop, navigate the Hadoop ecosystem, and build and manage Hadoop applications and clusters. Explains the origins of Hadoop, its economic benefits, and its functionality and practical applications Helps you find your way around the Hadoop ecosystem, program MapReduce, utilize design patterns, and get your Hadoop cluster up and running quickly and easily Details how to use Hadoop applications for data mining, web analytics and personalization, large-scale text processing, data science, and problem-solving Shows you how to improve the value of your Hadoop cluster, maximize your investment in Hadoop, and avoid common pitfalls when building your Hadoop cluster From programmers challenged with building and maintaining affordable, scaleable data systems to administrators who must deal with huge volumes of information effectively and efficiently, this how-to has something to help you with Hadoop.
Data-intensive Text Processing with MapReduce
Parallelize and Distribute Your Python Code

HBase
A Guide for Developers and Administrators