

Hadoop 2 Quick Start Guide Learn The Essentials Of Big Data Computing In The Apache Hadoop 2 Ecosystem Addison Wesley Data Analytics Series

When somebody should go to the book stores, search commencement by shop, shelf by shelf, it is truly problematic. This is why we offer the books compilations in this website. It will no question ease you to see guide **Hadoop 2 Quick Start Guide Learn The Essentials Of Big Data Computing In The Apache Hadoop 2 Ecosystem Addison Wesley Data Analytics Series** as you such as.

By searching the title, publisher, or authors of guide you in point of fact want, you can discover them rapidly. In the house, workplace, or perhaps in your method can be every best area within net connections. If you aspire to download and install the Hadoop 2 Quick Start Guide Learn The Essentials Of Big Data Computing In The Apache Hadoop 2 Ecosystem Addison Wesley Data Analytics Series , it is enormously simple then, since currently we extend the link to buy and create bargains to download and install Hadoop

2 Quick Start Guide Learn The Essentials Of Big Data Computing In The Apache Hadoop 2 Ecosystem
Addison Wesley Data Analytics Series hence simple!

Big Data Processing With Hadoop - Revathi, T.
2018-11-16

Due to the increasing availability of affordable internet services, the number of users, and the need for a wider range of multimedia-based applications, internet usage is on the rise. With so many users and such a large amount of data, the requirements of analyzing large data sets leads to the need for further advancements to information processing. Big Data Processing With Hadoop is an essential reference source that discusses possible solutions for millions of users working with a variety of data applications, who expect fast turnaround responses, but encounter issues with processing data at the rate it comes in. Featuring research on topics such as

market basket analytics, scheduler load simulator, and writing YARN applications, this book is ideally designed for IoT professionals, students, and engineers seeking coverage on many of the real-world challenges regarding big data.

Machine Learning with Scala Quick Start Guide -
Md. Rezaul Karim 2019-04-30

Supervised and unsupervised machine learning made easy in Scala with this quick-start guide. Key Features Construct and deploy machine learning systems that learn from your data and give accurate predictions Unleash the power of Spark ML along with popular machine learning algorithms to solve complex tasks in Scala. Solve hands-on problems by combining popular neural network architectures

such as LSTM and CNN using Scala with DeepLearning4j library Book Description Scala is a highly scalable integration of object-oriented nature and functional programming concepts that make it easy to build scalable and complex big data applications. This book is a handy guide for machine learning developers and data scientists who want to develop and train effective machine learning models in Scala. The book starts with an introduction to machine learning, while covering deep learning and machine learning basics. It then explains how to use Scala-based ML libraries to solve classification and regression problems using linear regression, generalized linear regression, logistic regression, support vector machine, and Naïve Bayes algorithms. It also covers tree-based ensemble techniques for solving both classification and regression problems. Moving ahead, it covers unsupervised learning techniques, such as

dimensionality reduction, clustering, and recommender systems. Finally, it provides a brief overview of deep learning using a real-life example in Scala. What you will learn Get acquainted with JVM-based machine learning libraries for Scala such as Spark ML and Deeplearning4j Learn RDDs, DataFrame, and Spark SQL for analyzing structured and unstructured data Understand supervised and unsupervised learning techniques with best practices and pitfalls Learn classification and regression analysis with linear regression, logistic regression, Naïve Bayes, support vector machine, and tree-based ensemble techniques Learn effective ways of clustering analysis with dimensionality reduction techniques Learn recommender systems with collaborative filtering approach Delve into deep learning and neural network architectures Who this book is for This book is for machine learning developers looking to train machine

learning models in Scala without spending too much time and effort. Some fundamental knowledge of Scala programming and some basics of statistics and linear algebra is all you need to get started with this book.

Apache Hadoop 3 Quick Start Guide - Hrishikesh Vijay Karambelkar 2018-10-31

A fast paced guide that will help you learn about Apache Hadoop 3 and its ecosystem Key Features Set up, configure and get started with Hadoop to get useful insights from large data sets Work with the different components of Hadoop such as MapReduce, HDFS and YARN Learn about the new features introduced in Hadoop 3 Book Description Apache Hadoop is a widely used distributed data platform. It enables large datasets to be efficiently processed instead of using one large computer to store and process the data. This book will get you started with the Hadoop ecosystem,

and introduce you to the main technical topics, including MapReduce, YARN, and HDFS. The book begins with an overview of big data and Apache Hadoop. Then, you will set up a pseudo Hadoop development environment and a multi-node enterprise Hadoop cluster. You will see how the parallel programming paradigm, such as MapReduce, can solve many complex data processing problems. The book also covers the important aspects of the big data software development lifecycle, including quality assurance and control, performance, administration, and monitoring. You will then learn about the Hadoop ecosystem, and tools such as Kafka, Sqoop, Flume, Pig, Hive, and HBase. Finally, you will look at advanced topics, including real time streaming using Apache Storm, and data analytics using Apache Spark. By the end of the book, you will be well versed with different configurations of the

Hadoop 3 cluster. What you will learn Store and analyze data at scale using HDFS, MapReduce and YARN Install and configure Hadoop 3 in different modes Use Yarn effectively to run different applications on Hadoop based platform Understand and monitor how Hadoop cluster is managed Consume streaming data using Storm, and then analyze it using Spark Explore Apache Hadoop ecosystem components, such as Flume, Sqoop, HBase, Hive, and Kafka Who this book is for Aspiring Big Data professionals who want to learn the essentials of Hadoop 3 will find this book to be useful. Existing Hadoop users who want to get up to speed with the new features introduced in Hadoop 3 will also benefit from this book. Having knowledge of Java programming will be an added advantage.

Practical Data Science with Hadoop and Spark - Ofer Mendeleevitch 2016-12-08

The Complete Guide to Data Science with Hadoop—For Technical Professionals, Businesspeople, and Students Demand is soaring for professionals who can solve real data science problems with Hadoop and Spark. Practical Data Science with Hadoop® and Spark is your complete guide to doing just that. Drawing on immense experience with Hadoop and big data, three leading experts bring together everything you need: high-level concepts, deep-dive techniques, real-world use cases, practical applications, and hands-on tutorials. The authors introduce the essentials of data science and the modern Hadoop ecosystem, explaining how Hadoop and Spark have evolved into an effective platform for solving data science problems at scale. In addition to comprehensive application coverage, the authors also provide useful guidance on the important steps of data ingestion, data munging, and visualization. Once the

groundwork is in place, the authors focus on specific applications, including machine learning, predictive modeling for sentiment analysis, clustering for document analysis, anomaly detection, and natural language processing (NLP). This guide provides a strong technical foundation for those who want to do practical data science, and also presents business-driven guidance on how to apply Hadoop and Spark to optimize ROI of data science initiatives. Learn What data science is, how it has evolved, and how to plan a data science career How data volume, variety, and velocity shape data science use cases Hadoop and its ecosystem, including HDFS, MapReduce, YARN, and Spark Data importation with Hive and Spark Data quality, preprocessing, preparation, and modeling Visualization: surfacing insights from huge data sets Machine learning: classification, regression, clustering, and anomaly detection Algorithms and Hadoop tools for

predictive modeling Cluster analysis and similarity functions Large-scale anomaly detection NLP: applying data science to human language HADOOP - TOM. WHITE 2015

Hadoop in Action - Chuck Lam 2010-11-30 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.

[Hadoop and Spark Fundamentals](#) - Doug Eadline 2018

"Hadoop and Spark Fundamentals LiveLessons provides 9+ hours of video introduction to the Apache Hadoop Big Data ecosystem. The tutorial includes background information and explains the core components of Hadoop, including Hadoop Distributed File Systems (HDFS), MapReduce, the YARN resource manager, and YARN Frameworks. In addition, it demonstrates how to use Hadoop at several levels, including the native Java interface, C++ pipes, and the universal streaming program interface. Examples include how to use benchmarks and high-level tools, including the Apache Pig scripting language, Apache Hive "SQL-like" interface, Apache Flume for streaming input, Apache Sqoop for import and

export of relational data, and Apache Oozie for Hadoop workflow management. In addition, there is comprehensive coverage of Spark, PySpark, and the Zeppelin web-GUI. The steps for easily installing a working Hadoop/Spark system on a desktop/laptop and on a local stand-alone cluster using the powerful Ambari GUI are also included. All software used in these LiveLessons is open source and freely available for your use and experimentation. A bonus lesson includes a quick primer on the Linux command line as used with Hadoop and Spark."-- Resource description page.

[Expert Hadoop 2 Administration](#) - Sam R. Alapati
2016-11-29

This is the eBook of the printed book and may not include any media, website access codes, or print supplements that may come packaged with the bound book. The Comprehensive, Up-to-Date Apache Hadoop Administration Handbook and

Reference “Sam Alapati has worked with production Hadoop clusters for six years. His unique depth of experience has enabled him to write the go-to resource for all administrators looking to spec, size, expand, and secure production Hadoop clusters of any size.” —Paul Dix, Series Editor In Expert Hadoop® Administration, leading Hadoop administrator Sam R. Alapati brings together authoritative knowledge for creating, configuring, securing, managing, and optimizing production Hadoop clusters in any environment. Drawing on his experience with large-scale Hadoop administration, Alapati integrates action-oriented advice with carefully researched explanations of both problems and solutions. He covers an unmatched range of topics and offers an unparalleled collection of realistic examples. Alapati demystifies complex Hadoop environments, helping you understand exactly what happens behind the

scenes when you administer your cluster. You'll gain unprecedented insight as you walk through building clusters from scratch and configuring high availability, performance, security, encryption, and other key attributes. The high-value administration skills you learn here will be indispensable no matter what Hadoop distribution you use or what Hadoop applications you run. Understand Hadoop's architecture from an administrator's standpoint

Create simple and fully distributed clusters
Run MapReduce and Spark applications in a Hadoop cluster
Manage and protect Hadoop data and high availability
Work with HDFS commands, file permissions, and storage management
Move data, and use YARN to allocate resources and schedule jobs
Manage job workflows with Oozie and Hue
Secure, monitor, log, and optimize Hadoop
Benchmark and troubleshoot Hadoop

Apache Hadoop YARN - Arun C. Murthy 2014

"Apache Hadoop is helping drive the Big Data revolution. Now, its data processing has been completely overhauled: Apache Hadoop YARN provides resource management at data center scale and easier ways to create distributed applications that process petabytes of data. And now in Apache Hadoop™ YARN, two Hadoop technical leaders show you how to develop new applications and adapt existing code to fully leverage these revolutionary advances." -- From the Amazon

Hadoop: The Definitive Guide - Tom White
2012-05-10

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

Hadoop 2.x Administration Cookbook - Gurmukh Singh 2017-05-26

Over 100 practical recipes to help you become an expert Hadoop administrator About This Book Become an expert Hadoop administrator and perform tasks to optimize your Hadoop Cluster Import and export data into Hive and use Oozie to manage workflow. Practical recipes will help you plan and secure your Hadoop cluster, and make it highly available Who This Book Is For If you are a system administrator with a basic understanding of Hadoop and you want to get into Hadoop administration, this book is for you. It's also ideal if you are a Hadoop administrator who wants a quick reference guide to all the Hadoop administration-related tasks and solutions to commonly occurring problems What You Will Learn Set up the Hadoop architecture to run a Hadoop cluster smoothly

Maintain a Hadoop cluster on HDFS, YARN, and MapReduce Understand high availability with Zookeeper and Journal Node Configure Flume for data ingestion and Oozie to run various workflows Tune the Hadoop cluster for optimal performance Schedule jobs on a Hadoop cluster using the Fair and Capacity scheduler Secure your cluster and troubleshoot it for various common pain points In Detail Hadoop enables the distributed storage and processing of large datasets across clusters of computers. Learning how to administer Hadoop is crucial to exploit its unique features. With this book, you will be able to overcome common problems encountered in Hadoop administration. The book begins with laying the foundation by showing you the steps needed to set up a Hadoop cluster and its various nodes. You will get a better understanding of how to maintain Hadoop cluster, especially on the HDFS layer and using YARN and MapReduce.

Further on, you will explore durability and high availability of a Hadoop cluster. You'll get a better understanding of the schedulers in Hadoop and how to configure and use them for your tasks. You will also get hands-on experience with the backup and recovery options and the performance tuning aspects of Hadoop. Finally, you will get a better understanding of troubleshooting, diagnostics, and best practices in Hadoop administration. By the end of this book, you will have a proper understanding of working with Hadoop clusters and will also be able to secure, encrypt it, and configure auditing for your Hadoop clusters. Style and approach This book contains short recipes that will help you run a Hadoop cluster efficiently. The recipes are solutions to real-life problems that administrators encounter while working with a Hadoop cluster
Hadoop Operations - Eric Sammer 2012-09-26
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

Hadoop Application Architectures - Mark Grover
2015-06-30

Get expert guidance on architecting end-to-end data management solutions with Apache Hadoop. While many sources explain how to use various components in the Hadoop ecosystem, this practical book takes you through architectural considerations necessary to tie those components together into a complete tailored application, based on your particular use case. To reinforce those lessons, the book's second section provides detailed examples of architectures used in some of the most commonly found Hadoop applications. Whether you're designing a new Hadoop application, or planning to integrate Hadoop into your existing data infrastructure, **Hadoop Application Architectures**

will skillfully guide you through the process. This book covers: Factors to consider when using Hadoop to store and model data Best practices for moving data in and out of the system Data processing frameworks, including MapReduce, Spark, and Hive Common Hadoop processing patterns, such as removing duplicate records and using windowing analytics Giraph, GraphX, and other tools for large graph processing on Hadoop Using workflow orchestration and scheduling tools such as Apache Oozie Near-real-time stream processing with Apache Storm, Apache Spark Streaming, and Apache Flume Architecture examples for clickstream analysis, fraud detection, and data warehousing

Learning Hadoop 2 - Garry Turkington 2015-02-13

If you are a system or application developer interested in learning how to solve practical problems using the Hadoop framework, then this

book is ideal for you. You are expected to be familiar with the Unix/Linux command-line interface and have some experience with the Java programming language. Familiarity with Hadoop would be a plus.

Hadoop 2 Quick-start Guide - Doug Eadline 2016

Advanced Intelligent Systems for Sustainable Development (AI2SD'2018) - Mostafa Ezziyyani
2019-03-06

This book includes the outcomes of the International Conference on Advanced Intelligent Systems for Sustainable Development (AI2SD-2018), held in Tangier, Morocco on July 12–14, 2018. Presenting the latest research in the field of computing sciences and information technology, it discusses new challenges and provides valuable insights into the field, the goal being to stimulate debate, and to promote closer interaction and interdisciplinary

collaboration between researchers and practitioners. Though chiefly intended for researchers and practitioners in advanced information technology management and networking, the book will also be of interest to those engaged in emerging fields such as data science and analytics, big data, internet of things, smart networked systems, artificial intelligence, expert systems and cloud computing.

Big Data Analytics with Hadoop 3 - Sridhar Alla

2018-05-31

Explore big data concepts, platforms, analytics, and their applications using the power of Hadoop 3
Key Features Learn Hadoop 3 to build effective big data analytics solutions on-premise and on cloud
Integrate Hadoop with other big data tools such as R, Python, Apache Spark, and Apache Flink
Exploit big data using Hadoop 3 with real-world examples
Book Description Apache Hadoop is the most popular platform for big data processing, and

can be combined with a host of other big data tools to build powerful analytics solutions. **Big Data Analytics with Hadoop 3** shows you how to do just that, by providing insights into the software as well as its benefits with the help of practical examples. Once you have taken a tour of Hadoop 3's latest features, you will get an overview of HDFS, MapReduce, and YARN, and how they enable faster, more efficient big data processing. You will then move on to learning how to integrate Hadoop with the open source tools, such as Python and R, to analyze and visualize data and perform statistical computing on big data. As you get acquainted with all this, you will explore how to use Hadoop 3 with Apache Spark and Apache Flink for real-time data analytics and stream processing. In addition to this, you will understand how to use Hadoop to build analytics solutions on the cloud and an end-to-end pipeline to perform big data analysis using practical

use cases. By the end of this book, you will be well-versed with the analytical capabilities of the Hadoop ecosystem. You will be able to build powerful solutions to perform big data analytics and get insight effortlessly. What you will learn Explore the new features of Hadoop 3 along with HDFS, YARN, and MapReduce Get well-versed with the analytical capabilities of Hadoop ecosystem using practical examples Integrate Hadoop with R and Python for more efficient big data processing Learn to use Hadoop with Apache Spark and Apache Flink for real-time data analytics Set up a Hadoop cluster on AWS cloud Perform big data analytics on AWS using Elastic Map Reduce Who this book is for Big Data Analytics with Hadoop 3 is for you if you are looking to build high-performance analytics solutions for your enterprise or business using Hadoop 3's powerful features, or you're new to big data analytics. A basic understanding of the Java

programming language is required.

Intelligent Systems and Applications - Kohei Arai
2020-08-25

The book *Intelligent Systems and Applications - Proceedings of the 2020 Intelligent Systems Conference* is a remarkable collection of chapters covering a wider range of topics in areas of intelligent systems and artificial intelligence and their applications to the real world. The Conference attracted a total of 545 submissions from many academic pioneering researchers, scientists, industrial engineers, students from all around the world. These submissions underwent a double-blind peer review process. Of those 545 submissions, 177 submissions have been selected to be included in these proceedings. As intelligent systems continue to replace and sometimes outperform human intelligence in decision-making processes, they have enabled a larger number of problems to be

tackled more effectively. This branching out of computational intelligence in several directions and use of intelligent systems in everyday applications have created the need for such an international conference which serves as a venue to report on up-to-the-minute innovations and developments. This book collects both theory and application based chapters on all aspects of artificial intelligence, from classical to intelligent scope. We hope that readers find the volume interesting and valuable; it provides the state of the art intelligent methods and techniques for solving real world problems along with a vision of the future research.

Big Data and Hadoop - VK Jain 2017-01-01

This book introduces you to the Big Data processing techniques addressing but not limited to various BI (business intelligence) requirements, such as reporting, batch analytics, online analytical processing (OLAP), data mining and Warehousing,

and predictive analytics. The book has been written on IBM's Platform of Hadoop framework. IBM Infosphere BigInsight has the highest amount of tutorial matter available free of cost on Internet which makes it easy to acquire proficiency in this technique. This therefore becomes highly vulnerable coaching materials in easy to learn steps. The book optimally provides the courseware as per MCA and M. Tech Level Syllabi of most of the Universities. All components of big Data Platform like Jaql, Hive Pig, Sqoop, Flume, Hadoop Streaming, Oozie: HBase, HDFS, FlumeNG, Whirr, Cloudera, Fuse, Zookeeper and Mahout: Machine learning for Hadoop has been discussed in sufficient Detail with hands on Exercises on each.

Big Data Fundamentals - Thomas Erl 2015-12-29

“This text should be required reading for everyone in contemporary business.” --Peter Woodhull, CEO, Modus21 “The one book that clearly describes and

links Big Data concepts to business utility.” --Dr. Christopher Starr, PhD “Simply, this is the best Big Data book on the market!” --Sam Rostam, Cascadian IT Group “...one of the most contemporary approaches I’ve seen to Big Data fundamentals...” --Joshua M. Davis, PhD The Definitive Plain-English Guide to Big Data for Business and Technology Professionals Big Data Fundamentals provides a pragmatic, no-nonsense introduction to Big Data. Best-selling IT author Thomas Erl and his team clearly explain key Big Data concepts, theory and terminology, as well as fundamental technologies and techniques. All coverage is supported with case study examples and numerous simple diagrams. The authors begin by explaining how Big Data can propel an organization forward by solving a spectrum of previously intractable business problems. Next, they demystify key analysis techniques and technologies and show how a Big

Data solution environment can be built and integrated to offer competitive advantages. Discovering Big Data’s fundamental concepts and what makes it different from previous forms of data analysis and data science Understanding the business motivations and drivers behind Big Data adoption, from operational improvements through innovation Planning strategic, business-driven Big Data initiatives Addressing considerations such as data management, governance, and security Recognizing the 5 “V” characteristics of datasets in Big Data environments: volume, velocity, variety, veracity, and value Clarifying Big Data’s relationships with OLTP, OLAP, ETL, data warehouses, and data marts Working with Big Data in structured, unstructured, semi-structured, and metadata formats Increasing value by integrating Big Data resources with corporate performance monitoring Understanding how Big Data leverages

distributed and parallel processing Using NoSQL and other technologies to meet Big Data's distinct data processing requirements Leveraging statistical approaches of quantitative and qualitative analysis Applying computational analysis methods, including machine learning

Cloud Computing for Science and Engineering - Ian Foster 2017-09-29

A guide to cloud computing for students, scientists, and engineers, with advice and many hands-on examples. The emergence of powerful, always-on cloud utilities has transformed how consumers interact with information technology, enabling video streaming, intelligent personal assistants, and the sharing of content. Businesses, too, have benefited from the cloud, outsourcing much of their information technology to cloud services. Science, however, has not fully exploited the advantages of the cloud. Could scientific discovery be accelerated

if mundane chores were automated and outsourced to the cloud? Leading computer scientists Ian Foster and Dennis Gannon argue that it can, and in this book offer a guide to cloud computing for students, scientists, and engineers, with advice and many hands-on examples. The book surveys the technology that underpins the cloud, new approaches to technical problems enabled by the cloud, and the concepts required to integrate cloud services into scientific work. It covers managing data in the cloud, and how to program these services; computing in the cloud, from deploying single virtual machines or containers to supporting basic interactive science experiments to gathering clusters of machines to do data analytics; using the cloud as a platform for automating analysis procedures, machine learning, and analyzing streaming data; building your own cloud with open source software; and cloud security. The book is

accompanied by a website, Cloud4SciEng.org, that provides a variety of supplementary material, including exercises, lecture slides, and other resources helpful to readers and instructors.

Hadoop Security - Ben Spivey 2015-06-29

As more corporations turn to Hadoop to store and process their most valuable data, the risk of a potential breach of those systems increases exponentially. This practical book not only shows Hadoop administrators and security architects how to protect Hadoop data from unauthorized access, it also shows how to limit the ability of an attacker to corrupt or modify data in the event of a security breach. Authors Ben Spivey and Joey Echeverria provide in-depth information about the security features available in Hadoop, and organize them according to common computer security concepts. You'll also get real-world examples that demonstrate how you can apply these concepts to

your use cases. Understand the challenges of securing distributed systems, particularly Hadoop Use best practices for preparing Hadoop cluster hardware as securely as possible Get an overview of the Kerberos network authentication protocol Delve into authorization and accounting principles as they apply to Hadoop Learn how to use mechanisms to protect data in a Hadoop cluster, both in transit and at rest Integrate Hadoop data ingest into enterprise-wide security architecture Ensure that security architecture reaches all the way to end-user access

Cloudera Administration Handbook - Rohit Menon 2014-07-18

An easy-to-follow Apache Hadoop administrator's guide filled with practical screenshots and explanations for each step and configuration. This book is great for administrators interested in setting up and managing a large Hadoop cluster. If you are an administrator, or want to be an administrator, and

you are ready to build and maintain a production-level cluster running CDH5, then this book is for you.

Hadoop 2 Quick-Start Guide - Douglas Eadline

2015-10-28

Get Started Fast with Apache Hadoop® 2, YARN, and Today's Hadoop Ecosystem With Hadoop 2.x and YARN, Hadoop moves beyond MapReduce to become practical for virtually any type of data processing. Hadoop 2.x and the Data Lake concept represent a radical shift away from conventional approaches to data usage and storage. Hadoop 2.x installations offer unmatched scalability and breakthrough extensibility that supports new and existing Big Data analytics processing methods and models. Hadoop® 2 Quick-Start Guide is the first easy, accessible guide to Apache Hadoop 2.x, YARN, and the modern Hadoop ecosystem. Building on his unsurpassed experience teaching Hadoop and Big

Data, author Douglas Eadline covers all the basics you need to know to install and use Hadoop 2 on personal computers or servers, and to navigate the powerful technologies that complement it. Eadline concisely introduces and explains every key Hadoop 2 concept, tool, and service, illustrating each with a simple “beginning-to-end” example and identifying trustworthy, up-to-date resources for learning more. This guide is ideal if you want to learn about Hadoop 2 without getting mired in technical details. Douglas Eadline will bring you up to speed quickly, whether you're a user, admin, devops specialist, programmer, architect, analyst, or data scientist. Coverage Includes Understanding what Hadoop 2 and YARN do, and how they improve on Hadoop 1 with MapReduce Understanding Hadoop-based Data Lakes versus RDBMS Data Warehouses Installing Hadoop 2 and core services on Linux machines, virtualized

sandboxes, or clusters Exploring the Hadoop Distributed File System (HDFS) Understanding the essentials of MapReduce and YARN application programming Simplifying programming and data movement with Apache Pig, Hive, Sqoop, Flume, Oozie, and HBase Observing application progress, controlling jobs, and managing workflows Managing Hadoop efficiently with Apache Ambari—including recipes for HDFS to NFSv3 gateway, HDFS snapshots, and YARN configuration Learning basic Hadoop 2 troubleshooting, and installing Apache Hue and Apache Spark *Apache Spark Quick Start Guide* - Shrey Mehrotra 2019-01-31 A practical guide for solving complex data processing challenges by applying the best optimizations techniques in Apache Spark. Key Features Learn about the core concepts and the latest developments in Apache Spark Master writing

efficient big data applications with Spark's built-in modules for SQL, Streaming, Machine Learning and Graph analysis Get introduced to a variety of optimizations based on the actual experience Book Description Apache Spark is a flexible framework that allows processing of batch and real-time data. Its unified engine has made it quite popular for big data use cases. This book will help you to get started with Apache Spark 2.0 and write big data applications for a variety of use cases. It will also introduce you to Apache Spark – one of the most popular Big Data processing frameworks. Although this book is intended to help you get started with Apache Spark, but it also focuses on explaining the core concepts. This practical guide provides a quick start to the Spark 2.0 architecture and its components. It teaches you how to set up Spark on your local machine. As we move ahead, you will be introduced to resilient distributed datasets (RDDs)

and DataFrame APIs, and their corresponding transformations and actions. Then, we move on to the life cycle of a Spark application and learn about the techniques used to debug slow-running applications. You will also go through Spark's built-in modules for SQL, streaming, machine learning, and graph analysis. Finally, the book will lay out the best practices and optimization techniques that are key for writing efficient Spark applications. By the end of this book, you will have a sound fundamental understanding of the Apache Spark framework and you will be able to write and optimize Spark applications. What you will learn

Learn core concepts such as RDDs, DataFrames, transformations, and more

Set up a Spark development environment

Choose the right APIs for your applications

Understand Spark's architecture and the execution flow of a Spark application

Explore built-in modules for SQL,

streaming, ML, and graph analysis

Optimize your Spark job for better performance

Who this book is for

If you are a big data enthusiast and love processing huge amount of data, this book is for you. If you are data engineer and looking for the best optimization techniques for your Spark applications, then you will find this book helpful. This book also helps data scientists who want to implement their machine learning algorithms in Spark. You need to have a basic understanding of any one of the programming languages such as Scala, Python or Java.

Mastering Spark with R - Javier Luraschi

2019-10-07

If you're like most R users, you have deep knowledge and love for statistics. But as your organization continues to collect huge amounts of data, adding tools such as Apache Spark makes a lot of sense. With this practical book, data scientists and

professionals working with large-scale data applications will learn how to use Spark from R to tackle big data and big compute problems. Authors Javier Luraschi, Kevin Kuo, and Edgar Ruiz show you how to use R with Spark to solve different data analysis problems. This book covers relevant data science topics, cluster computing, and issues that should interest even the most advanced users. Analyze, explore, transform, and visualize data in Apache Spark with R Create statistical models to extract information and predict outcomes; automate the process in production-ready workflows Perform analysis and modeling across many machines using distributed computing techniques Use large-scale data from multiple sources and different formats with ease from within Spark Learn about alternative modeling frameworks for graph processing, geospatial analysis, and genomics at scale Dive into advanced topics including custom

transformations, real-time data processing, and creating custom Spark extensions
Hadoop For Dummies - Dirk deRoos 2014-04-14
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 - Jimmy Lin 2010

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

MongoDB 4 Quick Start Guide - Doug Bierer
2018-09-28

A fast paced guide that will help you to create, read, update and delete data using MongoDB Key Features
Create secure databases with MongoDB
Manipulate and maintain your database
Model and use data in a No SQL environment with MongoDB
Book Description

MongoDB has grown to become the de facto NoSQL database with millions of users, from small start-ups to Fortune 500 companies. It can solve problems that are considered difficult, if not impossible, for aging RDBMS technologies. Written for version 4 of MongoDB, this book is the easiest way to get started with MongoDB. You will start by getting a MongoDB installation up and running in a safe and secure manner. You will learn how to perform mission-critical create, read, update, and delete operations, and set up database security. You will also learn about advanced features of MongoDB such as the aggregation pipeline, replication, and sharding. You will learn how to build a simple web application that uses MongoDB to respond to AJAX queries, and see how to make use of the MongoDB programming language driver for PHP. The examples incorporate new features available in MongoDB version 4 where appropriate. What you

will learn
Get a standard MongoDB database up and running quickly
Perform simple CRUD operations on the database using the MongoDB command shell
Set up a simple aggregation pipeline to return subsets of data grouped, sorted, and filtered
Safeguard your data via replication and handle massive amounts of data via sharding
Publish data from a web form to the database using a program language driver
Explore the basic CRUD operations performed using the PHP MongoDB driver
Who this book is for Web developers, IT professionals and Database Administrators (DBAs) who want to learn how to create and manage MongoDB databases.

Hadoop For Dummies - Dirk deRoos 2014-03-21

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 without 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, scalable data systems to administrators who must deal with huge volumes of information effectively and efficiently, this show-to has something to help you with Hadoop.

Kafka: The Definitive Guide - Neha Narkhede
2017-08-31

Every enterprise application creates data, whether it's log messages, metrics, user activity, outgoing messages, or something else. And how to move all of this data becomes nearly as important as the data itself. If you're an application architect, developer, or production engineer new to Apache Kafka, this practical guide shows you how to use this open source streaming platform to handle real-time data feeds. Engineers from Confluent and LinkedIn who are responsible for developing Kafka explain how to deploy production Kafka clusters, write reliable event-driven microservices, and build scalable

stream-processing applications with this platform. Through detailed examples, you'll learn Kafka's design principles, reliability guarantees, key APIs, and architecture details, including the replication protocol, the controller, and the storage layer.

Understand publish-subscribe messaging and how it fits in the big data ecosystem. Explore Kafka producers and consumers for writing and reading messages Understand Kafka patterns and use-case requirements to ensure reliable data delivery Get best practices for building data pipelines and applications with Kafka Manage Kafka in production, and learn to perform monitoring, tuning, and maintenance tasks Learn the most critical metrics among Kafka's operational measurements Explore how Kafka's stream delivery capabilities make it a perfect source for stream processing systems

Machine Learning with Apache Spark Quick Start

Guide - Jillur Quddus 2018-12-26

Combine advanced analytics including Machine Learning, Deep Learning Neural Networks and Natural Language Processing with modern scalable technologies including Apache Spark to derive actionable insights from Big Data in real-time Key Features Make a hands-on start in the fields of Big Data, Distributed Technologies and Machine Learning Learn how to design, develop and interpret the results of common Machine Learning algorithms Uncover hidden patterns in your data in order to derive real actionable insights and business value Book Description Every person and every organization in the world manages data, whether they realize it or not. Data is used to describe the world around us and can be used for almost any purpose, from analyzing consumer habits to fighting disease and serious organized crime. Ultimately, we manage data in order to derive value from it, and

many organizations around the world have traditionally invested in technology to help process their data faster and more efficiently. But we now live in an interconnected world driven by mass data creation and consumption where data is no longer rows and columns restricted to a spreadsheet, but an organic and evolving asset in its own right. With this realization comes major challenges for organizations: how do we manage the sheer size of data being created every second (think not only spreadsheets and databases, but also social media posts, images, videos, music, blogs and so on)? And once we can manage all of this data, how do we derive real value from it? The focus of Machine Learning with Apache Spark is to help us answer these questions in a hands-on manner. We introduce the latest scalable technologies to help us manage and process big data. We then introduce advanced analytical algorithms applied to real-world

use cases in order to uncover patterns, derive actionable insights, and learn from this big data. What you will learn Understand how Spark fits in the context of the big data ecosystem Understand how to deploy and configure a local development environment using Apache Spark Understand how to design supervised and unsupervised learning models Build models to perform NLP, deep learning, and cognitive services using Spark ML libraries Design real-time machine learning pipelines in Apache Spark Become familiar with advanced techniques for processing a large volume of data by applying machine learning algorithms Who this book is for This book is aimed at Business Analysts, Data Analysts and Data Scientists who wish to make a hands-on start in order to take advantage of modern Big Data technologies combined with Advanced Analytics.

Practical Hadoop Ecosystem - Deepak Vohra

2016-09-30

Learn how to use the Apache Hadoop projects, including MapReduce, HDFS, Apache Hive, Apache HBase, Apache Kafka, Apache Mahout, and Apache Solr. From setting up the environment to running sample applications each chapter in this book is a practical tutorial on using an Apache Hadoop ecosystem project. While several books on Apache Hadoop are available, most are based on the main projects, MapReduce and HDFS, and none discusses the other Apache Hadoop ecosystem projects and how they all work together as a cohesive big data development platform. What You Will Learn: Set up the environment in Linux for Hadoop projects using Cloudera Hadoop Distribution CDH 5 Run a MapReduce job Store data with Apache Hive, and Apache HBase Index data in HDFS with Apache Solr Develop a Kafka messaging system Stream Logs to HDFS with Apache Flume

Transfer data from MySQL database to Hive, HDFS, and HBase with Sqoop Create a Hive table over Apache Solr Develop a Mahout User Recommender System Who This Book Is For: Apache Hadoop developers. Pre-requisite knowledge of Linux and some knowledge of Hadoop is required.

Software Architecture for Big Data and the Cloud - Ivan Mistrik 2017-06-12

Software Architecture for Big Data and the Cloud is designed to be a single resource that brings together research on how software architectures can solve the challenges imposed by building big data software systems. The challenges of big data on the software architecture can relate to scale, security, integrity, performance, concurrency, parallelism, and dependability, amongst others. Big data handling requires rethinking architectural solutions to meet functional and non-functional requirements related to volume, variety and velocity. The book's

editors have varied and complementary backgrounds in requirements and architecture, specifically in software architectures for cloud and big data, as well as expertise in software engineering for cloud and big data. This book brings together work across different disciplines in software engineering, including work expanded from conference tracks and workshops led by the editors. Discusses systematic and disciplined approaches to building software architectures for cloud and big data with state-of-the-art methods and techniques Presents case studies involving enterprise, business, and government service deployment of big data applications Shares guidance on theory, frameworks, methodologies, and architecture for cloud and big data

R Web Scraping Quick Start Guide - Olgun Aydin 2018-10-31

Web Scraping techniques are getting more popular,

since data is as valuable as oil in 21st century. Through this book get some key knowledge about using XPath, regEX; web scraping libraries for R like rvest and RSelenium technologies. Key FeaturesTechniques, tools and frameworks for web scraping with RScrape data effortlessly from a variety of websites Learn how to selectively choose the data to scrape, and build your datasetBook Description Web scraping is a technique to extract data from websites. It simulates the behavior of a website user to turn the website itself into a web service to retrieve or introduce new data. This book gives you all you need to get started with scraping web pages using R programming. You will learn about the rules of RegEx and Xpath, key components for scraping website data. We will show you web scraping techniques, methodologies, and frameworks. With this book's guidance, you will become comfortable with the tools to write and

test RegEx and XPath rules. We will focus on examples of dynamic websites for scraping data and how to implement the techniques learned. You will learn how to collect URLs and then create XPath rules for your first web scraping script using rvest library. From the data you collect, you will be able to calculate the statistics and create R plots to visualize them. Finally, you will discover how to use Selenium drivers with R for more sophisticated scraping. You will create AWS instances and use R to connect a PostgreSQL database hosted on AWS. By the end of the book, you will be sufficiently confident to create end-to-end web scraping systems using R. What you will learnWrite and create regEX rulesWrite XPath rules to query your dataLearn how web scraping methods workUse rvest to crawl web pagesStore data retrieved from the webLearn the key uses of Rselenium to scrape dataWho this book is for This book is for R

programmers who want to get started quickly with web scraping, as well as data analysts who want to learn scraping using R. Basic knowledge of R is all you need to get started with this book.

Hadoop: The Definitive Guide - Tom White

2010-09-24

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

Big Data Made Easy - Michael Frampton 2014-12-31

Many corporations are finding that the size of their data sets are outgrowing the capability of their systems to store and process them. The data is becoming too big to manage and use with traditional tools. The solution: implementing a big data system. As Big Data Made Easy: A Working Guide to the Complete Hadoop Toolset shows, Apache Hadoop offers a scalable, fault-tolerant system for storing and processing data in parallel. It has a very rich toolset that allows for storage (Hadoop), configuration (YARN and ZooKeeper), collection (Nutch and Solr), processing (Storm, Pig, and Map Reduce), scheduling (Oozie), moving (Sqoop and Avro), monitoring (Chukwa, Ambari, and Hue), testing (Big Top), and analysis (Hive). The problem is that the Internet offers IT pros wading into big data

many versions of the truth and some outright falsehoods born of ignorance. What is needed is a book just like this one: a wide-ranging but easily understood set of instructions to explain where to get Hadoop tools, what they can do, how to install them, how to configure them, how to integrate them, and how to use them successfully. And you need an expert who has worked in this area for a decade—someone just like author and big data expert Mike Frampton. Big Data Made Easy approaches the problem of managing massive data sets from a systems perspective, and it explains the roles for each project (like architect and tester, for example) and shows how the Hadoop toolset can be used at each system stage. It explains, in an easily understood manner and through numerous examples, how to use each tool. The book also explains the sliding scale of tools available depending upon data size and when and how to use

them. Big Data Made Easy shows developers and architects, as well as testers and project managers, how to: Store big data Configure big data Process big data Schedule processes Move data among SQL and NoSQL systems Monitor data Perform big data analytics Report on big data processes and projects Test big data systems Big Data Made Easy also explains the best part, which is that this toolset is free. Anyone can download it and—with the help of this book—start to use it within a day. With the skills this book will teach you under your belt, you will add value to your company or client immediately, not to mention your career.

Hadoop in Practice - Alex Holmes 2014-10-12
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

Spark: The Definitive Guide - Bill Chambers
2018-02-08

Learn how to use, deploy, and maintain Apache Spark with this comprehensive guide, written by the creators of the open-source cluster-computing framework. With an emphasis on improvements and new features in Spark 2.0, authors Bill Chambers and Matei Zaharia break down Spark topics into distinct sections, each with unique goals. You'll explore the basic operations and common functions of Spark's structured APIs, as well as Structured Streaming, a new high-level API for

building end-to-end streaming applications.

Developers and system administrators will learn the fundamentals of monitoring, tuning, and debugging Spark, and explore machine learning techniques and scenarios for employing MLlib, Spark's scalable machine-learning library. Get a gentle overview of big data and Spark Learn about DataFrames, SQL, and Datasets—Spark's core APIs—through worked examples Dive into Spark's low-level APIs, RDDs, and execution of SQL and DataFrames Understand

how Spark runs on a cluster Debug, monitor, and tune Spark clusters and applications Learn the power of Structured Streaming, Spark's stream-processing engine Learn how you can apply MLlib to a variety of problems, including classification or recommendation

- Edward Capriolo 2012-09-26

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

Programming Hive