

Getting Started With Impala Interactive SQL For Apache Hadoop

If you ally infatuation such a referred **Getting Started With Impala Interactive SQL For Apache Hadoop** books that will find the money for you worth, get the categorically best seller from us currently from several preferred authors. If you desire to witty books, lots of novels, tale, jokes, and more fictions collections are along with launched, from best seller to one of the most current released.

You may not be perplexed to enjoy every books collections Getting Started With Impala Interactive SQL For Apache Hadoop that we will unquestionably offer. It is not regarding the costs. Its practically what you dependence currently. This Getting Started With Impala Interactive SQL For Apache Hadoop , as one of the most enthusiastic sellers here will agreed be among the best options to review.

The Semantic Web - ISWC 2014 - Peter Mika
2014-10-09
The two-volume set LNCS 8796 and 8797

constitutes the refereed proceedings of the 13th International Semantic Web Conference, ISWC 2014, held in Riva del Garda, in October 2014.

The International Semantic Web Conference is the premier forum for Semantic Web research, where cutting edge scientific results and technological innovations are presented, where problems and solutions are discussed, and where the future of this vision is being developed. It brings together specialists in fields such as artificial intelligence, databases, social networks, distributed computing, Web engineering, information systems, human-computer interaction, natural language processing, and the social sciences. Part 1 (LNCS 8796) contains a total of 38 papers which were presented in the research track. They were carefully reviewed and selected from 180 submissions. Part 2 (LNCS 8797) contains 15 papers from the 'semantic Web in use' track which were accepted from 46 submissions. In addition, it presents 16 contributions of the RBDS track and 6 papers of the doctoral consortium.

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

Introducing Data Science - Davy Cielen

2016-05-02

Summary Introducing Data Science teaches you how to accomplish the fundamental tasks that occupy data scientists. Using the Python language and common Python libraries, you'll experience firsthand the challenges of dealing with data at scale and gain a solid foundation in data science. Purchase of the print book includes a free eBook in PDF, Kindle, and ePub formats from Manning Publications. About the Technology Many companies need developers

with data science skills to work on projects ranging from social media marketing to machine learning. Discovering what you need to learn to begin a career as a data scientist can seem bewildering. This book is designed to help you get started. About the Book Introducing Data Science Introducing Data Science explains vital data science concepts and teaches you how to accomplish the fundamental tasks that occupy data scientists. You'll explore data visualization, graph databases, the use of NoSQL, and the data science process. You'll use the Python language and common Python libraries as you experience firsthand the challenges of dealing with data at scale. Discover how Python allows you to gain insights from data sets so big that they need to be stored on multiple machines, or from data moving so quickly that no single machine can handle it. This book gives you hands-on experience with the most popular Python data science libraries, Scikit-learn and StatsModels. After reading this book, you'll have the solid

foundation you need to start a career in data science. What's Inside Handling large data Introduction to machine learning Using Python to work with data Writing data science algorithms About the Reader This book assumes you're comfortable reading code in Python or a similar language, such as C, Ruby, or JavaScript. No prior experience with data science is required. About the Authors Davy Cielen, Arno D. B. Meysman, and Mohamed Ali are the founders and managing partners of Optimately and Maiton, where they focus on developing data science projects and solutions in various sectors. Table of Contents Data science in a big data world The data science process Machine learning Handling large data on a single computer First steps in big data Join the NoSQL movement The rise of graph databases Text mining and text analytics Data visualization to the end user

Learning Spark - Jules S. Damji 2020-07-16
Data is bigger, arrives faster, and comes in a

variety of formats—and it all needs to be processed at scale for analytics or machine learning. But how can you process such varied workloads efficiently? Enter Apache Spark. Updated to include Spark 3.0, this second edition shows data engineers and data scientists why structure and unification in Spark matters. Specifically, this book explains how to perform simple and complex data analytics and employ machine learning algorithms. Through step-by-step walk-throughs, code snippets, and notebooks, you'll be able to: Learn Python, SQL, Scala, or Java high-level Structured APIs Understand Spark operations and SQL Engine Inspect, tune, and debug Spark operations with Spark configurations and Spark UI Connect to data sources: JSON, Parquet, CSV, Avro, ORC, Hive, S3, or Kafka Perform analytics on batch and streaming data using Structured Streaming Build reliable data pipelines with open source Delta Lake and Spark Develop machine learning pipelines with MLlib and productionize models

using MLflow

QlikView Your Business - Oleg Troyansky

2015-07-22

Unlock the meaning of your data with QlikView
The Qlik platform was designed to provide a fast and easy data analytics tool, and QlikView Your Business is your detailed, full-color, step-by-step guide to understanding Qlikview's powerful features and techniques so you can quickly start unlocking your data's potential. This expert author team brings real-world insight together with practical business analytics, so you can approach, explore, and solve business intelligence problems using the robust Qlik toolset and clearly communicate your results to stakeholders using powerful visualization features in QlikView and Qlik Sense. This book starts at the basic level and dives deep into the most advanced QlikView techniques, delivering tangible value and knowledge to new users and experienced developers alike. As an added benefit, every topic presented is enhanced with

tips, tricks, and insightful recommendations that the authors accumulated through years of developing QlikView analytics. This is the book for you: If you are a developer whose job is to load transactional data into Qlik BI environment, and who needs to understand both the basics and the most advanced techniques of Qlik data modelling and scripting If you are a data analyst whose job is to develop actionable and insightful QlikView visualizations to share within your organization If you are a project manager or business person, who wants to get a better understanding of the Qlik Business Intelligence platform and its capabilities
What You Will Learn: The book covers three common business scenarios - Sales, Profitability, and Inventory Analysis. Each scenario contains four chapters, covering the four main disciplines of business analytics: Business Case, Data Modeling, Scripting, and Visualizations. The material is organized by increasing levels of complexity. Following our comprehensive tutorial, you will

learn simple and advanced QlikView and Qlik Sense concepts, including the following: Data Modeling: Transforming Transactional data into Dimensional models Building a Star Schema Linking multiple fact tables using Link Tables Combing multiple tables into a single fact able using Concatenated Fact models Managing slowly changing dimensions Advanced date handling, using the As of Date table Calculating running balances Basic and Advanced Scripting: How to use the Data Load Script language for implementing data modeling techniques How to build and use the QVD data layer Building a multi-tier data architectures Using variables, loops, subroutines, and other script control statements Advanced scripting techniques for a variety of ETL solutions Building Insightful Visualizations in QlikView: Introduction into QlikView sheet objects — List Boxes, Text Objects, Charts, and more Designing insightful Dashboards in QlikView Using advanced calculation techniques, such as Set Analysis and

Advanced Aggregation Using variables for What-If Analysis, as well as using variables for storing calculations, colors, and selection filters Advanced visualization techniques - normalized and non-normalized Mekko charts, Waterfall charts, Whale Tail charts, and more Building Insightful Visualizations in Qlik Sense: Introducing Qlik Sense - how it is different from QlikView and what is similar? Creating Sense sheet objects Building and using the Library of Master Items Exploring Qlik Sense unique features — Storytelling, Geo Mapping, and using Extensions Whether you are just starting out with QlikView or are ready to dive deeper, QlikView Your Business is your comprehensive guide to sharpening your QlikView skills and unleashing the power of QlikView in your organization.

Mastering Hadoop 3 - Chanchal Singh

2019-02-28

A comprehensive guide to mastering the most advanced Hadoop 3 concepts Key FeaturesGet to

grips with the newly introduced features and capabilities of Hadoop 3Crunch and process data using MapReduce, YARN, and a host of tools within the Hadoop ecosystemSharpen your Hadoop skills with real-world case studies and codeBook Description Apache Hadoop is one of the most popular big data solutions for distributed storage and for processing large chunks of data. With Hadoop 3, Apache promises to provide a high-performance, more fault-tolerant, and highly efficient big data processing platform, with a focus on improved scalability and increased efficiency. With this guide, you'll understand advanced concepts of the Hadoop ecosystem tool. You'll learn how Hadoop works internally, study advanced concepts of different ecosystem tools, discover solutions to real-world use cases, and understand how to secure your cluster. It will then walk you through HDFS, YARN, MapReduce, and Hadoop 3 concepts. You'll be able to address common challenges like using Kafka efficiently, designing low latency,

reliable message delivery Kafka systems, and handling high data volumes. As you advance, you'll discover how to address major challenges when building an enterprise-grade messaging system, and how to use different stream processing systems along with Kafka to fulfil your enterprise goals. By the end of this book, you'll have a complete understanding of how components in the Hadoop ecosystem are effectively integrated to implement a fast and reliable data pipeline, and you'll be equipped to tackle a range of real-world problems in data pipelines. What you will learnGain an in-depth understanding of distributed computing using Hadoop 3Develop enterprise-grade applications using Apache Spark, Flink, and moreBuild scalable and high-performance Hadoop data pipelines with security, monitoring, and data governanceExplore batch data processing patterns and how to model data in HadoopMaster best practices for enterprises using, or planning to use, Hadoop 3 as a data

platform Understand security aspects of Hadoop, including authorization and authentication Who this book is for If you want to become a big data professional by mastering the advanced concepts of Hadoop, this book is for you. You'll also find this book useful if you're a Hadoop professional looking to strengthen your knowledge of the Hadoop ecosystem.

Fundamental knowledge of the Java programming language and basics of Hadoop is necessary to get started with this book.

Hadoop: The Definitive Guide - Tom White
2015-03-25

Get ready to unlock the power of your data. With the fourth edition of 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. Using Hadoop 2 exclusively, author Tom White presents new

chapters on YARN and several Hadoop-related projects such as Parquet, Flume, Crunch, and Spark. You'll learn about recent changes to Hadoop, and explore new case studies on Hadoop's role in healthcare systems and genomics data processing. Learn fundamental components such as MapReduce, HDFS, and YARN Explore MapReduce in depth, including steps for developing applications with it Set up and maintain a Hadoop cluster running HDFS and MapReduce on YARN Learn two data formats: Avro for data serialization and Parquet for nested data Use data ingestion tools such as Flume (for streaming data) and Sqoop (for bulk data transfer) Understand how high-level data processing tools like Pig, Hive, Crunch, and Spark work with Hadoop Learn the HBase distributed database and the ZooKeeper distributed configuration service

Stream Processing with Apache Flink -

Fabian Hueske 2019-04-11

Get started with Apache Flink, the open source

framework that powers some of the world's largest stream processing applications. With this practical book, you'll explore the fundamental concepts of parallel stream processing and discover how this technology differs from traditional batch data processing. Longtime Apache Flink committers Fabian Hueske and Vasia Kalavri show you how to implement scalable streaming applications with Flink's DataStream API and continuously run and maintain these applications in operational environments. Stream processing is ideal for many use cases, including low-latency ETL, streaming analytics, and real-time dashboards as well as fraud detection, anomaly detection, and alerting. You can process continuous data of any kind, including user interactions, financial transactions, and IoT data, as soon as you generate them. Learn concepts and challenges of distributed stateful stream processing Explore Flink's system architecture, including its event-time processing mode and fault-tolerance model

Understand the fundamentals and building blocks of the DataStream API, including its time-based and stateful operators Read data from and write data to external systems with exactly-once consistency Deploy and configure Flink clusters Operate continuously running streaming applications

Apache Superset Quick Start Guide -
Shashank Shekhar 2018-12-19

Integrate open source data analytics and build business intelligence on SQL databases with Apache Superset. The quick, intuitive nature for data visualization in a web application makes it easy for creating interactive dashboards. Key Features Work with Apache Superset's rich set of data visualizations Create interactive dashboards and data storytelling Easily explore data Book Description Apache Superset is a modern, open source, enterprise-ready business intelligence (BI) web application. With the help of this book, you will see how Superset integrates with popular databases like Postgres,

Google BigQuery, Snowflake, and MySQL. You will learn to create real time data visualizations and dashboards on modern web browsers for your organization using Superset. First, we look at the fundamentals of Superset, and then get it up and running. You'll go through the requisite installation, configuration, and deployment. Then, we will discuss different columnar data types, analytics, and the visualizations available. You'll also see the security tools available to the administrator to keep your data safe. You will learn how to visualize relationships as graphs instead of coordinates on plain orthogonal axes. This will help you when you upload your own entity relationship dataset and analyze the dataset in new, different ways. You will also see how to analyze geographical regions by working with location data. Finally, we cover a set of tutorials on dashboard designs frequently used by analysts, business intelligence professionals, and developers. What you will learn Get to grips with the fundamentals of data exploration using

Superset Set up a working instance of Superset on cloud services like Google Compute Engine Integrate Superset with SQL databases Build dashboards with Superset Calculate statistics in Superset for numerical, categorical, or text data Understand visualization techniques, filtering, and grouping by aggregation Manage user roles and permissions in Superset Work with SQL Lab Who this book is for This book is for data analysts, BI professionals, and developers who want to learn Apache Superset. If you want to create interactive dashboards from SQL databases, this book is what you need. Working knowledge of Python will be an advantage but not necessary to understand this book.

Moving Hadoop to the Cloud - Bill Havanki
2017-07-14

Until recently, Hadoop deployments existed on hardware owned and run by organizations. Now, of course, you can acquire the computing resources and network connectivity to run Hadoop clusters in the cloud. But there's a lot

more to deploying Hadoop to the public cloud than simply renting machines. This hands-on guide shows developers and systems administrators familiar with Hadoop how to install, use, and manage cloud-born clusters efficiently. You'll learn how to architect clusters that work with cloud-provider features—not just to avoid pitfalls, but also to take full advantage of these services. You'll also compare the Amazon, Google, and Microsoft clouds, and learn how to set up clusters in each of them. Learn how Hadoop clusters run in the cloud, the problems they can help you solve, and their potential drawbacks Examine the common concepts of cloud providers, including compute capabilities, networking and security, and storage Build a functional Hadoop cluster on cloud infrastructure, and learn what the major providers require Explore use cases for high availability, relational data with Hive, and complex analytics with Spark Get patterns and practices for running cloud clusters, from

designing for price and security to dealing with maintenance

Apache Hive Cookbook - Hanish Bansal
2016-04-29

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.

[MySQL Cookbook](#) - Paul DuBois 2003

DuBois organizes his cookbook's recipes into sections on the problem, the solution stated simply, and the solution implemented in code and discussed. The implementation and discussion sections are the most valuable, as they contain the command sequences, code listings, and design explanations that can be transferred to outside projects.

Apache Hive Essentials - Dayong Du

2015-02-26

If you are a data analyst, developer, or simply someone who wants to use Hive to explore and analyze data in Hadoop, this is the book for you. Whether you are new to big data or an expert, with this book, you will be able to master both the basic and the advanced features of Hive. Since Hive is an SQL-like language, some previous experience with the SQL language and

databases is useful to have a better understanding of this book.

[Building a Data Warehouse](#) - Vincent Rainardi
2008-03-11

Here is the ideal field guide for data warehousing implementation. This book first teaches you how to build a data warehouse, including defining the architecture, understanding the methodology, gathering the requirements, designing the data models, and creating the databases. Coverage then explains how to populate the data warehouse and explores how to present data to users using reports and multidimensional databases and how to use the data in the data warehouse for business intelligence, customer relationship management, and other purposes. It also details testing and how to administer data warehouse operation.

[Cloudera Impala](#) - John Russell 2013-11-25
Learn about Cloudera Impala--an open source project that's opening up the Apache Hadoop

software stack to a wide audience of database analysts, users, and developers. The Impala massively parallel processing (MPP) engine makes SQL queries of Hadoop data simple enough to be accessible to analysts familiar with SQL and to users of business intelligence tools--and it's fast enough to be used for interactive exploration and experimentation.

[Handbook of Research on Big Data Storage and Visualization Techniques](#) - Segall, Richard S.
2018-01-05

The digital age has presented an exponential growth in the amount of data available to individuals looking to draw conclusions based on given or collected information across industries. Challenges associated with the analysis, security, sharing, storage, and visualization of large and complex data sets continue to plague data scientists and analysts alike as traditional data processing applications struggle to adequately manage big data. The Handbook of Research on Big Data Storage and Visualization

Techniques is a critical scholarly resource that explores big data analytics and technologies and their role in developing a broad understanding of issues pertaining to the use of big data in multidisciplinary fields. Featuring coverage on a broad range of topics, such as architecture patterns, programming systems, and computational energy, this publication is geared towards professionals, researchers, and students seeking current research and application topics on the subject.

Microsoft Power BI Quick Start Guide -

Devin Knight 2018-07-30

Bring your data to life with Power BI Key Features Get to grips with the fundamentals of Microsoft Power BI and its Business Intelligence capabilities Build accurate analytical models, reports and dashboards Get faster and more intuitive insights from your data using Microsoft Power BI Book Description Microsoft Power BI is a cloud-based service that helps you easily visualize and share insights using your

organization's data. This book will get you started with business intelligence using the Power BI toolset, covering essential concepts such as installation, designing effective data models, as well as building basic dashboards and visualizations to make your data come to life You will learn how to get your data the way you want - connecting to data sources sources and how to clean your data with the Power BI Query Editor. You will next learn how to properly design your data model to make your data easier to work with.. You will next learn how to properly design your data model to navigate table relationships and build DAX formulas to make your data easier to work with. Visualizing your data is another key element of this book, and you will learn how to follow proper data visualization styles and enhanced digital storytelling techniques. By the end of this book, you will understand how to administer your organization's Power BI environment so deployment can be made seamless, data refreshes can run properly, and

security can be fully implemented What you will learn Connect to data sources using both import and DirectQuery options Use the Query Editor to apply data transformations and data cleansing processes, including learning how to write M and R scripts Design optimized data models by designing relationships and DAX calculations Leverage built-in and custom visuals to design effective reports Use the Power BI Desktop and Power BI Service to implement Row Level Security on your model Administer a Power BI cloud tenant for your organization Deploy your Power BI Desktop files into the Power BI Report Server Who this book is for This book is for aspiring Business Intelligence professionals who want to get up and running with Microsoft Power BI. If you have a basic understanding of BI concepts and want to learn how to apply them using Microsoft Power BI, this book is for you.

Disruptive Analytics - Thomas W. Dinsmore
2016-08-27

Learn all you need to know about seven key

innovations disrupting business analytics today. These innovations—the open source business model, cloud analytics, the Hadoop ecosystem, Spark and in-memory analytics, streaming analytics, Deep Learning, and self-service analytics—are radically changing how businesses use data for competitive advantage. Taken together, they are disrupting the business analytics value chain, creating new opportunities. Enterprises who seize the opportunity will thrive and prosper, while others struggle and decline: disrupt or be disrupted. Disruptive Business Analytics provides strategies to profit from disruption. It shows you how to organize for insight, build and provision an open source stack, how to practice lean data warehousing, and how to assimilate disruptive innovations into an organization. Through a short history of business analytics and a detailed survey of products and services, analytics authority Thomas W. Dinsmore provides a practical explanation of the most compelling

innovations available today. What You'll Learn Discover how the open source business model works and how to make it work for you See how cloud computing completely changes the economics of analytics Harness the power of Hadoop and its ecosystem Find out why Apache Spark is everywhere Discover the potential of streaming and real-time analytics Learn what Deep Learning can do and why it matters See how self-service analytics can change the way organizations do business Who This Book Is For Corporate actors at all levels of responsibility for analytics: analysts, CIOs, CTOs, strategic decision makers, managers, systems architects, technical marketers, product developers, IT personnel, and consultants.

Next-Generation Big Data - Butch Quinto
2018-06-12

Utilize this practical and easy-to-follow guide to modernize traditional enterprise data warehouse and business intelligence environments with next-generation big data technologies. Next-

Generation Big Data takes a holistic approach, covering the most important aspects of modern enterprise big data. The book covers not only the main technology stack but also the next-generation tools and applications used for big data warehousing, data warehouse optimization, real-time and batch data ingestion and processing, real-time data visualization, big data governance, data wrangling, big data cloud deployments, and distributed in-memory big data computing. Finally, the book has an extensive and detailed coverage of big data case studies from Navistar, Cerner, British Telecom, Shopzilla, Thomson Reuters, and Mastercard. What You'll Learn Install Apache Kudu, Impala, and Spark to modernize enterprise data warehouse and business intelligence environments, complete with real-world, easy-to-follow examples, and practical advice Integrate HBase, Solr, Oracle, SQL Server, MySQL, Flume, Kafka, HDFS, and Amazon S3 with Apache Kudu, Impala, and Spark Use

StreamSets, Talend, Pentaho, and CDAP for real-time and batch data ingestion and processing
Utilize Trifacta, Alteryx, and Datameer for data wrangling and interactive data processing
Turbocharge Spark with Alluxio, a distributed in-memory storage platform
Deploy big data in the cloud using Cloudera Director
Perform real-time data visualization and time series analysis using Zoomdata, Apache Kudu, Impala, and Spark
Understand enterprise big data topics such as big data governance, metadata management, data lineage, impact analysis, and policy enforcement, and how to use Cloudera Navigator to perform common data governance tasks
Implement big data use cases such as big data warehousing, data warehouse optimization, Internet of Things, real-time data ingestion and analytics, complex event processing, and scalable predictive modeling
Study real-world big data case studies from innovative companies, including Navistar, Cerner, British Telecom, Shopzilla, Thomson Reuters, and Mastercard

Who This Book Is For
BI and big data warehouse professionals interested in gaining practical and real-world insight into next-generation big data processing and analytics using Apache Kudu, Impala, and Spark; and those who want to learn more about other advanced enterprise topics
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

Getting Started with Impala - John Russell
2014-09-25

Learn how to write, tune, and port SQL queries and other statements for a Big Data environment, using Impala—the massively parallel processing SQL query engine for Apache Hadoop. The best practices in this practical

guide help you design database schemas that not only interoperate with other Hadoop components, and are convenient for administrators to manage and monitor, but also accommodate future expansion in data size and evolution of software capabilities. Written by John Russell, documentation lead for the Cloudera Impala project, this book gets you working with the most recent Impala releases quickly. Ideal for database developers and business analysts, the latest revision covers analytics functions, complex types, incremental statistics, subqueries, and submission to the Apache incubator. *Getting Started with Impala* includes advice from Cloudera's development team, as well as insights from its consulting engagements with customers. Learn how Impala integrates with a wide range of Hadoop components Attain high performance and scalability for huge data sets on production clusters Explore common developer tasks, such as porting code to Impala and optimizing performance Use tutorials for

working with billion-row tables, date- and time-based values, and other techniques Learn how to transition from rigid schemas to a flexible model that evolves as needs change Take a deep dive into joins and the roles of statistics

SQL on Big Data - Sumit Pal 2016-11-17

Learn various commercial and open source products that perform SQL on Big Data platforms. You will understand the architectures of the various SQL engines being used and how the tools work internally in terms of execution, data movement, latency, scalability, performance, and system requirements. This book consolidates in one place solutions to the challenges associated with the requirements of speed, scalability, and the variety of operations needed for data integration and SQL operations. After discussing the history of the how and why of SQL on Big Data, the book provides in-depth insight into the products, architectures, and innovations happening in this rapidly evolving space. SQL on Big Data discusses in detail the

innovations happening, the capabilities on the horizon, and how they solve the issues of performance and scalability and the ability to handle different data types. The book covers how SQL on Big Data engines are permeating the OLTP, OLAP, and Operational analytics space and the rapidly evolving HTAP systems. You will learn the details of: Batch Architectures—Understand the internals and how the existing Hive engine is built and how it is evolving continually to support new features and provide lower latency on queries Interactive Architectures—Understanding how SQL engines are architected to support low latency on large data sets Streaming Architectures—Understanding how SQL engines are architected to support queries on data in motion using in-memory and lock-free data structures Operational Architectures—Understanding how SQL engines are architected for transactional and operational systems to support transactions on Big Data

platforms Innovative Architectures—Explore the rapidly evolving newer SQL engines on Big Data with innovative ideas and concepts Who This Book Is For: Business analysts, BI engineers, developers, data scientists and architects, and quality assurance professionals/div

Hadoop: The Definitive Guide - Tom White
2015-03-25

Get ready to unlock the power of your data. With the fourth edition of 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. Using Hadoop 2 exclusively, author Tom White presents new chapters on YARN and several Hadoop-related projects such as Parquet, Flume, Crunch, and Spark. You'll learn about recent changes to Hadoop, and explore new case studies on Hadoop's role in healthcare systems and genomics data

processing. Learn fundamental components such as MapReduce, HDFS, and YARN Explore MapReduce in depth, including steps for developing applications with it Set up and maintain a Hadoop cluster running HDFS and MapReduce on YARN Learn two data formats: Avro for data serialization and Parquet for nested data Use data ingestion tools such as Flume (for streaming data) and Sqoop (for bulk data transfer) Understand how high-level data processing tools like Pig, Hive, Crunch, and Spark work with Hadoop Learn the HBase distributed database and the ZooKeeper distributed configuration service

-
Readings in Database Systems - Joseph M. Hellerstein 2005

The latest edition of a popular text and reference on database research, with substantial new material and revision; covers classical literature and recent hot topics. Lessons from database

research have been applied in academic fields ranging from bioinformatics to next-generation Internet architecture and in industrial uses including Web-based e-commerce and search engines. The core ideas in the field have become increasingly influential. This text provides both students and professionals with a grounding in database research and a technical context for understanding recent innovations in the field. The readings included treat the most important issues in the database area--the basic material for any DBMS professional. This fourth edition has been substantially updated and revised, with 21 of the 48 papers new to the edition, four of them published for the first time. Many of the sections have been newly organized, and each section includes a new or substantially revised introduction that discusses the context, motivation, and controversies in a particular area, placing it in the broader perspective of database research. Two introductory articles, never before published, provide an organized,

current introduction to basic knowledge of the field; one discusses the history of data models and query languages and the other offers an architectural overview of a database system. The remaining articles range from the classical literature on database research to treatments of current hot topics, including a paper on search engine architecture and a paper on application servers, both written expressly for this edition. The result is a collection of papers that are seminal and also accessible to a reader who has a basic familiarity with database systems.

Study on Data Placement Strategies in Distributed RDF Stores - D.D. Janke

2020-03-18

The distributed setting of RDF stores in the cloud poses many challenges, including how to optimize data placement on the compute nodes to improve query performance. In this book, a novel benchmarking methodology is developed for data placement strategies; one that overcomes these limitations by using a data-

placement-strategy-independent distributed RDF store to analyze the effect of the data placement strategies on query performance. Frequently used data placement strategies have been evaluated, and this evaluation challenges the commonly held belief that data placement strategies which emphasize local computation lead to faster query executions. Indeed, results indicate that queries with a high workload can be executed faster on hash-based data placement strategies than on, for example, minimal edge-cut covers. The analysis of additional measurements indicates that vertical parallelization (i.e., a well-distributed workload) may be more important than horizontal containment (i.e., minimal data transport) for efficient query processing. Two such data placement strategies are proposed: the first, found in the literature, is entitled overpartitioned minimal edge-cut cover, and the second is the newly developed molecule hash cover. Evaluation revealed a balanced query

workload and a high horizontal containment, which lead to a high vertical parallelization. As a result, these strategies demonstrated better query performance than other frequently used data placement strategies. The book also tests the hypothesis that collocating small connected triple sets on the same compute node while balancing the amount of triples stored on the different compute nodes leads to a high vertical parallelization.

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.

Getting Started with Kudu - Jean-Marc

Spaggiari 2018-07-09

Fast data ingestion, serving, and analytics in the Hadoop ecosystem have forced developers and architects to choose solutions using the least common denominator—either fast analytics at the cost of slow data ingestion or fast data ingestion at the cost of slow analytics. There is an answer to this problem. With the Apache Kudu column-oriented data store, you can easily perform fast analytics on fast data. This practical guide shows you how. Begun as an internal project at Cloudera, Kudu is an open source solution compatible with many data processing frameworks in the Hadoop environment. In this book, current and former solutions professionals from Cloudera provide use cases, examples, best practices, and sample code to help you get up to speed with Kudu. Explore Kudu's high-level design, including how it spreads data across servers Fully administer a Kudu cluster, enable security, and add or remove nodes Learn Kudu's client-side APIs, including how to integrate

Apache Impala, Spark, and other frameworks for data manipulation Examine Kudu's schema design, including basic concepts and primitives necessary to make your project successful Explore case studies for using Kudu for real-time IoT analytics, predictive modeling, and in combination with another storage engine

Big Data Analytics with Spark - Mohammed Guller 2015-12-29

Big Data Analytics with Spark is a step-by-step guide for learning Spark, which is an open-source fast and general-purpose cluster computing framework for large-scale data analysis. You will learn how to use Spark for different types of big data analytics projects, including batch, interactive, graph, and stream data analysis as well as machine learning. In addition, this book will help you become a much sought-after Spark expert. Spark is one of the hottest Big Data technologies. The amount of data generated today by devices, applications and users is exploding. Therefore, there is a

critical need for tools that can analyze large-scale data and unlock value from it. Spark is a powerful technology that meets that need. You can, for example, use Spark to perform low latency computations through the use of efficient caching and iterative algorithms; leverage the features of its shell for easy and interactive Data analysis; employ its fast batch processing and low latency features to process your real time data streams and so on. As a result, adoption of Spark is rapidly growing and is replacing Hadoop MapReduce as the technology of choice for big data analytics. This book provides an introduction to Spark and related big-data technologies. It covers Spark core and its add-on libraries, including Spark SQL, Spark Streaming, GraphX, and MLlib. Big Data Analytics with Spark is therefore written for busy professionals who prefer learning a new technology from a consolidated source instead of spending countless hours on the Internet trying to pick bits and pieces from different sources. The book

also provides a chapter on Scala, the hottest functional programming language, and the program that underlies Spark. You'll learn the basics of functional programming in Scala, so that you can write Spark applications in it. What's more, Big Data Analytics with Spark provides an introduction to other big data technologies that are commonly used along with Spark, like Hive, Avro, Kafka and so on. So the book is self-sufficient; all the technologies that you need to know to use Spark are covered. The only thing that you are expected to know is programming in any language. There is a critical shortage of people with big data expertise, so companies are willing to pay top dollar for people with skills in areas like Spark and Scala. So reading this book and absorbing its principles will provide a boost—possibly a big boost—to your career.

[Learning SQL on SQL Server 2005](#) - Sikha Saha Bagui 2006-04-26

Anyone who interacts with today's modern

databases needs to know SQL (Structured Query Language), the standard language for generating, manipulating, and retrieving database information. In recent years, the dramatic rise in the popularity of relational databases and multi-user databases has fueled a healthy demand for application developers and others who can write SQL code efficiently and correctly. If you're new to databases, or need a SQL refresher, Learning SQL on SQL Server 2005 is an ideal step-by-step introduction to this database query tool, with everything you need for programming SQL using Microsoft's SQL Server 2005—one of the most powerful and popular database engines used today. Plenty of books explain database theory. This guide lets you apply the theory as you learn SQL. You don't need prior database knowledge, or even prior computer knowledge. Based on a popular university-level course designed by authors Sikha Saha Bagui and Richard Walsh Earp, Learning SQL on SQL Server 2005 starts with

very simple SQL concepts, and slowly builds into more complex query development. Every topic, concept, and idea comes with examples of code and output, along with exercises to help you gain proficiency in SQL and SQL Server 2005. With this book, you'll learn: Beginning SQL commands, such as how and where to type an SQL query, and how to create, populate, alter and delete tables How to customize SQL Server 2005's settings and about SQL Server 2005's functions About joins, a common database mechanism for combining tables Query development, the use of views and other derived structures, and simple set operations Subqueries, aggregate functions and correlated subqueries, as well as indexes and constraints that can be added to tables in SQL Server 2005 Whether you're an undergraduate computer science or MIS student, a self-learner who has access to the new Microsoft database, or work for your company's IT department, Learning SQL on SQL Server 2005 will get you up to speed on

SQL in no time.

Mindstorms - Seymour A. Papert 2020-10-06

In this revolutionary book, a renowned computer scientist explains the importance of teaching children the basics of computing and how it can prepare them to succeed in the ever-evolving tech world. Computers have completely changed the way we teach children. We have Mindstorms to thank for that. In this book, pioneering computer scientist Seymour Papert uses the invention of LOGO, the first child-friendly programming language, to make the case for the value of teaching children with computers. Papert argues that children are more than capable of mastering computers, and that teaching computational processes like debugging in the classroom can change the way we learn everything else. He also shows that schools saturated with technology can actually improve socialization and interaction among students and between students and teachers. Technology changes every day, but the basic ways that

computers can help us learn remain. For thousands of teachers and parents who have sought creative ways to help children learn with computers, Mindstorms is their bible.

Programming Hive - Edward Capriolo

2012-09-26

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

Big Data Analytics with Java - Rajat Mehta

2017-07-31

Learn the basics of analytics on big data using Java, machine learning and other big data tools
About This Book Acquire real-world set of tools for building enterprise level data science applications Surpasses the barrier of other languages in data science and learn create useful object-oriented codes Extensive use of Java compliant big data tools like apache spark, Hadoop, etc. Who This Book Is For This book is for Java developers who are looking to perform data analysis in production environment. Those who wish to implement data analysis in their Big

data applications will find this book helpful.
What You Will Learn Start from simple analytic tasks on big data Get into more complex tasks with predictive analytics on big data using machine learning Learn real time analytic tasks Understand the concepts with examples and case studies Prepare and refine data for analysis Create charts in order to understand the data See various real-world datasets In Detail This book covers case studies such as sentiment analysis on a tweet dataset, recommendations on a movielens dataset, customer segmentation on an ecommerce dataset, and graph analysis on actual flights dataset. This book is an end-to-end guide to implement analytics on big data with Java. Java is the de facto language for major big data environments, including Hadoop. This book will teach you how to perform analytics on big data with production-friendly Java. This book basically divided into two sections. The first part is an introduction that will help the readers get acquainted with big data environments, whereas

the second part will contain a hardcore discussion on all the concepts in analytics on big data. It will take you from data analysis and data visualization to the core concepts and advantages of machine learning, real-life usage of regression and classification using Naive Bayes, a deep discussion on the concepts of clustering, and a review of simple neural networks on big data using deepLearning4j or plain Java Spark code. This book is a must-have book for Java developers who want to start learning big data analytics and want to use it in the real world. Style and approach The approach of book is to deliver practical learning modules in manageable content. Each chapter is a self-contained unit of a concept in big data analytics. Book will step by step builds the competency in the area of big data analytics. Examples using real world case studies to give ideas of real applications and how to use the techniques mentioned. The examples and case studies will be shown using both theory and code.

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

Learning Apache Drill - Charles Givre

2018-11-02

Get up to speed with Apache Drill, an extensible distributed SQL query engine that reads massive datasets in many popular file formats such as Parquet, JSON, and CSV. Drill reads data in HDFS or in cloud-native storage such as S3 and works with Hive metastores along with distributed databases such as HBase, MongoDB, and relational databases. Drill works everywhere: on your laptop or in your largest cluster. In this practical book, Drill committers Charles Givre and Paul Rogers show analysts and data scientists how to query and analyze raw data using this powerful tool. Data scientists today spend about 80% of their time just

gathering and cleaning data. With this book, you'll learn how Drill helps you analyze data more effectively to drive down time to insight. Use Drill to clean, prepare, and summarize delimited data for further analysis Query file types including logfiles, Parquet, JSON, and other complex formats Query Hadoop, relational databases, MongoDB, and Kafka with standard SQL Connect to Drill programmatically using a variety of languages Use Drill even with challenging or ambiguous file formats Perform sophisticated analysis by extending Drill's functionality with user-defined functions Facilitate data analysis for network security, image metadata, and machine learning

Reasoning Web. Learning, Uncertainty, Streaming, and Scalability - Claudia d'Amato

2018-09-14

This volume contains lecture notes of the 14th Reasoning Web Summer School (RW 2018), held in Esch-sur-Alzette, Luxembourg, in September 2018. The research areas of Semantic Web,

Linked Data, and Knowledge Graphs have recently received a lot of attention in academia and industry. Since its inception in 2001, the Semantic Web has aimed at enriching the existing Web with meta-data and processing methods, so as to provide Web-based systems with intelligent capabilities such as context awareness and decision support. The Semantic Web vision has been driving many community efforts which have invested a lot of resources in developing vocabularies and ontologies for annotating their resources semantically. Besides ontologies, rules have long been a central part of the Semantic Web framework and are available as one of its fundamental representation tools, with logic serving as a unifying foundation. Linked Data is a related research area which studies how one can make RDF data available on the Web and interconnect it with other data with the aim of increasing its value for everybody. Knowledge Graphs have been shown useful not only for Web search (as demonstrated by

Google, Bing, etc.) but also in many application domains.

Data Lake for Enterprises - Tomcy John
2017-05-31

A practical guide to implementing your enterprise data lake using Lambda Architecture as the base About This Book Build a full-fledged data lake for your organization with popular big data technologies using the Lambda architecture as the base Delve into the big data technologies required to meet modern day business strategies A highly practical guide to implementing enterprise data lakes with lots of examples and real-world use-cases Who This Book Is For Java developers and architects who would like to implement a data lake for their enterprise will find this book useful. If you want to get hands-on experience with the Lambda Architecture and big data technologies by implementing a practical solution using these technologies, this book will also help you. What You Will Learn Build an enterprise-level data lake using the

relevant big data technologies Understand the core of the Lambda architecture and how to apply it in an enterprise Learn the technical details around Sqoop and its functionalities Integrate Kafka with Hadoop components to acquire enterprise data Use flume with streaming technologies for stream-based processing Understand stream-based processing with reference to Apache Spark Streaming Incorporate Hadoop components and know the advantages they provide for enterprise data lakes Build fast, streaming, and high-performance applications using Elasticsearch Make your data ingestion process consistent across various data formats with configurability Process your data to derive intelligence using machine learning algorithms In Detail The term "Data Lake" has recently emerged as a prominent term in the big data industry. Data scientists can make use of it in deriving meaningful insights that can be used by businesses to redefine or transform the way they

operate. Lambda architecture is also emerging as one of the very eminent patterns in the big data landscape, as it not only helps to derive useful information from historical data but also correlates real-time data to enable business to take critical decisions. This book tries to bring these two important aspects — data lake and lambda architecture—together. This book is divided into three main sections. The first introduces you to the concept of data lakes, the importance of data lakes in enterprises, and getting you up-to-speed with the Lambda architecture. The second section delves into the principal components of building a data lake using the Lambda architecture. It introduces you to popular big data technologies such as Apache Hadoop, Spark, Sqoop, Flume, and Elasticsearch. The third section is a highly practical demonstration of putting it all together, and shows you how an enterprise data lake can be implemented, along with several real-world use-cases. It also shows you how other

peripheral components can be added to the lake to make it more efficient. By the end of this book, you will be able to choose the right big data technologies using the lambda architectural patterns to build your enterprise data lake. Style and approach The book takes a pragmatic approach, showing ways to leverage big data technologies and lambda architecture to build an enterprise-level data lake.

Getting Started with Impala - John Russell
2014-09-25

Learn how to write, tune, and port SQL queries and other statements for a Big Data environment, using Impala—the massively parallel processing SQL query engine for Apache Hadoop. The best practices in this practical guide help you design database schemas that not only interoperate with other Hadoop components, and are convenient for administrators to manage and monitor, but also accommodate future expansion in data size and evolution of software capabilities. Written by John Russell,

documentation lead for the Cloudera Impala project, this book gets you working with the most recent Impala releases quickly. Ideal for database developers and business analysts, the latest revision covers analytics functions, complex types, incremental statistics, subqueries, and submission to the Apache incubator. Getting Started with Impala includes advice from Cloudera’s development team, as well as insights from its consulting engagements with customers. Learn how Impala integrates with a wide range of Hadoop components Attain high performance and scalability for huge data sets on production clusters Explore common developer tasks, such as porting code to Impala and optimizing performance Use tutorials for working with billion-row tables, date- and time-based values, and other techniques Learn how to transition from rigid schemas to a flexible model that evolves as needs change Take a deep dive into joins and the roles of statistics

Hadoop in Practice - Alex Holmes 2014-09-29

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

Architecting Modern Data Platforms - Jan Kunigk 2018-12-05

There's a lot of information about big data technologies, but splicing these technologies into an end-to-end enterprise data platform is a daunting task not widely covered. With this practical book, you'll learn how to build big data infrastructure both on-premises and in the cloud and successfully architect a modern data

platform. Ideal for enterprise architects, IT managers, application architects, and data engineers, this book shows you how to overcome the many challenges that emerge during Hadoop projects. You'll explore the vast landscape of tools available in the Hadoop and big data realm in a thorough technical primer before diving into: Infrastructure: Look at all component layers in a modern data platform, from the server to the data center, to establish a solid foundation for data in your enterprise Platform: Understand aspects of deployment, operation, security, high availability, and disaster recovery, along with everything you need to know to integrate your platform with the rest of your enterprise IT Taking Hadoop to the cloud: Learn the important architectural aspects of running a big data platform in the cloud while maintaining enterprise security and high availability