

Beginning Apache Pig: Big Data Processing Made Easy

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

Read Free Beginning Apache Pig: Big Data Processing Made Easy

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

Read Free Beginning Apache Pig: Big Data Processing Made Easy

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

A handy reference guide for data analysts and data scientists to help to obtain value from big data analytics using Spark on Hadoop clusters About This Book This book is based on the latest 2.0 version of Apache Spark and 2.7 version of Hadoop integrated with most commonly used tools. Learn all Spark stack components including latest topics such as DataFrames, DataSets, GraphFrames, Structured Streaming, DataFrame based ML Pipelines and SparkR. Integrations with frameworks such as HDFS, YARN and tools such as Jupyter, Zeppelin, NiFi, Mahout, HBase Spark Connector, GraphFrames, H2O and Hivemall. Who This Book Is For Though this book is primarily aimed at data analysts and data scientists, it will also help architects, programmers, and practitioners. Knowledge of either Spark or Hadoop would be beneficial. It is assumed that you have basic

Read Free Beginning Apache Pig: Big Data Processing Made Easy

programming background in Scala, Python, SQL, or R programming with basic Linux experience. Working experience within big data environments is not mandatory. What You Will Learn Find out and implement the tools and techniques of big data analytics using Spark on Hadoop clusters with wide variety of tools used with Spark and Hadoop Understand all the Hadoop and Spark ecosystem components Get to know all the Spark components: Spark Core, Spark SQL, DataFrames, DataSets, Conventional and Structured Streaming, MLlib, ML Pipelines and Graphx See batch and real-time data analytics using Spark Core, Spark SQL, and Conventional and Structured Streaming Get to grips with data science and machine learning using MLlib, ML Pipelines, H2O, Hivemall, Graphx, SparkR and Hivemall. In Detail Big Data Analytics book aims at providing the fundamentals of Apache Spark and Hadoop. All Spark components - Spark Core, Spark SQL, DataFrames, Data sets, Conventional Streaming, Structured Streaming, MLlib, Graphx and Hadoop core components - HDFS, MapReduce and

Read Free Beginning Apache Pig: Big Data Processing Made Easy

Yarn are explored in greater depth with implementation examples on Spark + Hadoop clusters. It is moving away from MapReduce to Spark. So, advantages of Spark over MapReduce are explained at great depth to reap benefits of in-memory speeds. DataFrames API, Data Sources API and new Data set API are explained for building Big Data analytical applications. Real-time data analytics using Spark Streaming with Apache Kafka and HBase is covered to help building streaming applications. New Structured streaming concept is explained with an IOT (Internet of Things) use case. Machine learning techniques are covered using MLlib, ML Pipelines and SparkR and Graph Analytics are covered with GraphX and GraphFrames components of Spark. Readers will also get an opportunity to get started with web based notebooks such as Jupyter, Apache Zeppelin and data flow tool Apache NiFi to analyze and visualize data. Style and approach This step-by-step pragmatic guide will make life easy no matter what your level of experience. You will deep dive into Apache Spark on Hadoop clusters

Read Free Beginning Apache Pig: Big Data Processing Made Easy

through ample exciting real-life examples. Practical tutorial explains data science in simple terms to help programmers and data analysts get started with Data Science

This thoroughly revised second edition of "Big Data" introduces application of big data to various domains from farming to healthcare to managing traffic and many more. The book takes a big leap with introduction of three new primer on Data Modeling and Management, Artificial Intelligence and careers in Data Science. Important topics like Big Data Programming languages are simplified and areas like MongoDB have been expanded. The key concepts and technological developments are explained with illustrations. This simple and easy to understand book is aimed for the final year students of Computer Science, professionals and big data enthusiasts. With a series of pictures at the beginning of every chapter from nature and human interaction with it, the book tells a parallel story about life cycle and the many aspects of big data applications in primary education, water resource

Read Free Beginning Apache Pig: Big Data Processing Made Easy

management, precision farming, finance, etc. Few Highlights: • A new chapter on Data Science careers and job roles • A primer on Artificial Intelligence, and its advantages and threats • A primer on Data Modeling and Management • New section on General Data Protection Rights (GDPR) regime in Europe

In *Beginning Big Data with Power BI and Excel 2013*, you will learn to solve business problems by tapping the power of Microsoft's Excel and Power BI to import data from NoSQL and SQL databases and other sources, create relational data models, and analyze business problems through sophisticated dashboards and data-driven maps. While *Beginning Big Data with Power BI and Excel 2013* covers prominent tools such as Hadoop and the NoSQL databases, it recognizes that most small and medium-sized businesses don't have the Big Data processing needs of a Netflix, Target, or Facebook. Instead, it shows how to import data and use the self-service analytics available in Excel with Power BI. As you'll see through the book's numerous case examples, these tools—which you already know how

Read Free Beginning Apache Pig: Big Data Processing Made Easy

to use—can perform many of the same functions as the higher-end Apache tools many people believe are required to carry out in Big Data projects. Through instruction, insight, advice, and case studies, *Beginning Big Data with Power BI and Excel 2013* will show you how to: Import and mash up data from web pages, SQL and NoSQL databases, the Azure Marketplace and other sources. Tap into the analytical power of PivotTables and PivotCharts and develop relational data models to track trends and make predictions based on a wide range of data. Understand basic statistics and use Excel with PowerBI to do sophisticated statistical analysis—including identifying trends and correlations. Use SQL within Excel to do sophisticated queries across multiple tables, including NoSQL databases. Create complex formulas to solve real-world business problems using Data Analysis Expressions (DAX). "This course will teach to smoothly handle big data sets using Hadoop 3. The course starts by covering basic commands used by big data developers on a daily basis. Then, you'll focus on

Read Free Beginning Apache Pig: Big Data Processing Made Easy

HDFS architecture and command lines that a developer uses frequently. Next, you'll use Flume to import data from other ecosystems into the Hadoop ecosystem, which plays a crucial role in the data available for storage and analysis using MapReduce. Also, you'll learn to import and export data from RDBMS to HDFS and vice-versa using SQOOP. Then, you'll learn about Apache Pig, which is used to deal with data using Flume and SQOOP. Here you'll also learn to load, transform, and store data in Pig relation. Finally, you'll dive into Hive functionality and learn to load, update, delete content in Hive. By the end of the course, you'll have gained enough knowledge to work with big data using Hadoop. So, grab the course and handle big data sets with ease."--Resource description page.

A comprehensive guide to mastering the most advanced Hadoop 3 concepts Key Features Get to grips with the newly introduced features and capabilities of Hadoop 3 Crunch and process data using MapReduce, YARN, and a host of tools within the Hadoop ecosystem Sharpen your Hadoop skills with real-world case

Read Free Beginning Apache Pig: Big Data Processing Made Easy

studies and code Book 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

Read Free Beginning Apache Pig: Big Data Processing Made Easy

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 learn

Gain an in-depth understanding of distributed computing using Hadoop 3

Develop enterprise-grade applications using Apache Spark, Flink, and more

Build scalable and high-performance Hadoop data pipelines with security, monitoring, and data governance

Explore batch data processing patterns and how to model data in Hadoop

Master 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

Read Free Beginning Apache Pig: Big Data Processing Made Easy

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.

This guide is an ideal learning tool and reference for Apache Pig, the programming language that helps programmers describe and run large data projects on Hadoop. With Pig, they can analyze data without having to create a full-fledged application--making it easy for them to experiment with new data sets.

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 IBMs 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

Read Free Beginning Apache Pig: Big Data Processing Made Easy

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.

[Apache Hadoop 3 Quick Start Guide](#)

[BIG DATA AND HADOOP](#)

[Big Data: Second Edition](#)

[Beginning Apache Cassandra Development](#)

[Big Data for Chimps](#)

[Learning Apache Pig](#)

[Big Data Processing Made Easy](#)

[A Working Guide to the Complete Hadoop Toolset](#)

[Moving beyond MapReduce and Batch](#)

[Processing with Apache Hadoop 2](#)

[Big Data For Dummies](#)

[Building Real-World Big Data Systems on Azure HDInsight Using the Hadoop Ecosystem](#)

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

Read Free Beginning Apache Pig: Big Data Processing Made Easy

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. Big data is a term that describes the large volume of data – both structured and unstructured – that inundates a business on a day-to-day basis. But it ' s not the amount of data that ' s important. It ' s

Read Free Beginning Apache Pig: Big Data Processing Made Easy

what organizations do with the data that matters. Big data can be analyzed for insights that lead to better decisions and strategic business moves. The use of Big Data is becoming common these days by the companies to outperform their peers. In most industries, existing competitors and new entrants alike will use the strategies resulting from the analyzed data to compete, innovate and capture value. Big Data helps the organizations to create new growth opportunities and entirely new categories of companies that can combine and analyze industry data. These companies have ample information about the products and services, buyers and suppliers, consumer preferences that can be captured and analyzed. While the term “ big data ” is relatively new, the act of gathering and storing large amounts of information for eventual analysis is ages old. The concept gained momentum in the early 2000s when industry analyst Doug Laney articulated the now-mainstream definition of big data as the three Vs: Volume.

Organizations collect data from a variety of sources, including business transactions, social media and information from sensor or machine-to-machine data. In the past, storing it would 've been a problem – but new technologies (such as Hadoop) have eased the burden. The name 'Big Data' itself is related to a size which is enormous. Size of data plays very crucial role in determining value out of data. Also, whether a particular data can actually be considered as a Big Data or not, is dependent upon volume of data. Hence, 'Volume' is one characteristic which needs to be considered while dealing with 'Big Data'. Velocity. Data streams in at an unprecedented speed and must be dealt with in a timely manner. RFID tags, sensors and smart metering are driving the need to deal with torrents of data in near-real time. The term 'velocity' refers to the speed of generation of data. How fast the data is generated and processed to meet the demands, determines real potential in the data. Big Data Velocity deals with the speed at which data flows in from sources like business processes, application logs, networks and social media sites, sensors, Mobile devices, etc. The flow of data is massive and continuous. Variety. Data comes in all types of formats – from structured datasets numeric data in traditional databases to

Read Free Beginning Apache Pig: Big Data Processing Made Easy

unstructured text documents, email, video, audio, stock ticker data and financial transactions. Variety refers to heterogeneous sources and the nature of data, both structured and unstructured. During earlier days, spreadsheets and databases were the only sources of data considered by most of the applications. Now days, data in the form of emails, photos, videos, monitoring devices, PDFs, audio, etc. is also being considered in the analysis applications. This variety of unstructured data poses certain issues for storage, mining and analysing data.

Big Data Simplified blends technology with strategy and delves into applications of big data in specialized areas, such as recommendation engines, data science and Internet of Things (IoT) and enables a practitioner to make the right technology choice. The steps to strategize a big data implementation are also discussed in detail. This book presents a holistic approach to the topic, covering a wide landscape of big data technologies like Hadoop 2.0 and package implementations, such as Cloudera. In-depth discussion of associated technologies, such as MapReduce, Hive, Pig, Oozie, ApacheZookeeper, Flume, Kafka, Spark, Python and NoSQL databases like Cassandra, MongoDB, GraphDB, etc., is also included. Features: 1. Important concepts are backed by code snippets enabling step-by-step practical implementation 2. Includes case study with complete code and detailing the concepts are discussed 3. Numerous objective and subjective-type questions added for readers to evaluate their learning Table of Contents: Chapter 1) A Closer Look at Data Chapter 2) Introducing Big Data Chapter 3) Introducing Hadoop Chapter 4) Introducing MapReduce Chapter 5) Introducing NoSQL Chapter 6) Introducing Spark and Kafka Chapter 7) Other BigData Tools and Technologies Chapter 8) Working with Big Data in R Chapter 9) Working with Big Data in Python Chapter 10) Big Data Applied Chapter 11) Big Data Strategy Chapter 12) Case Study: Retail Near Real-time Analytics

“ This book is a critically needed resource for the newly released Apache Hadoop 2.0, highlighting YARN as the significant

Read Free Beginning Apache Pig: Big Data Processing Made Easy

breakthrough that broadens Hadoop beyond the MapReduce paradigm. ” —From the Foreword by Raymie Stata, CEO of Altiscale

The Insider ’ s Guide to Building Distributed, Big Data Applications with Apache Hadoop™ YARN 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. YARN project founder Arun Murthy and project lead Vinod Kumar Vavilapalli demonstrate how YARN increases scalability and cluster utilization, enables new programming models and services, and opens new options beyond Java and batch processing. They walk you through the entire YARN project lifecycle, from installation through deployment. You ’ ll find many examples drawn from the authors ’ cutting-edge experience—first as Hadoop ’ s earliest developers and implementers at Yahoo! and now as Hortonworks developers moving the platform forward and helping customers succeed with it. Coverage includes YARN ’ s goals, design, architecture, and components—how it expands the Apache Hadoop ecosystem

Exploring YARN on a single node
Administering YARN clusters and Capacity Scheduler
Running existing MapReduce applications
Developing a large-scale clustered YARN application
Discovering new open source frameworks that run under YARN

Over 90 hands-on recipes to help you learn and master the intricacies of Apache Hadoop 2.X, YARN, Hive, Pig, Oozie, Flume, Sqoop, Apache Spark, and Mahout

About This Book Implement outstanding Machine Learning use cases on your own analytics models and processes. Solutions to common problems when working with the Hadoop ecosystem. Step-by-step implementation of end-to-end big data use cases. Who This Book Is For Readers who have a basic knowledge of big data systems and want to advance their knowledge with hands-on recipes. What You Will Learn Installing and

Read Free Beginning Apache Pig: Big Data Processing Made Easy

maintaining Hadoop 2.X cluster and its ecosystem. Write advanced Map Reduce programs and understand design patterns. Advanced Data Analysis using the Hive, Pig, and Map Reduce programs. Import and export data from various sources using Sqoop and Flume. Data storage in various file formats such as Text, Sequential, Parquet, ORC, and RC Files. Machine learning principles with libraries such as Mahout Batch and Stream data processing using Apache Spark In Detail Big data is the current requirement. Most organizations produce huge amount of data every day. With the arrival of Hadoop-like tools, it has become easier for everyone to solve big data problems with great efficiency and at minimal cost. Grasping Machine Learning techniques will help you greatly in building predictive models and using this data to make the right decisions for your organization. Hadoop Real World Solutions Cookbook gives readers insights into learning and mastering big data via recipes. The book not only clarifies most big data tools in the market but also provides best practices for using them. The book provides recipes that are based on the latest versions of Apache Hadoop 2.X, YARN, Hive, Pig, Sqoop, Flume, Apache Spark, Mahout and many more such ecosystem tools. This real-world-solution cookbook is packed with handy recipes you can apply to your own everyday issues. Each chapter provides in-depth recipes that can be referenced easily. This book provides detailed practices on the latest technologies such as YARN and Apache Spark. Readers will be able to consider themselves as big data experts on completion of this book. This guide is an invaluable tutorial if you are planning to implement a big data warehouse for your business. Style and approach An easy-to-follow guide that walks you through world of big data. Each tool in the Hadoop ecosystem is explained in detail and the recipes are placed in such a manner that readers can implement them sequentially. Plenty of reference links are provided for advanced reading.

A comprehensive practical guide that walks you through the multiple stages of data management in enterprise and gives you numerous design patterns with appropriate code examples to solve frequent problems in each of these stages. The chapters are organized to mimic

Read Free Beginning Apache Pig: Big Data Processing Made Easy

the sequential data flow evidenced in Analytics platforms, but they can also be read independently to solve a particular group of problems in the Big Data life cycle. If you are an experienced developer who is already familiar with Pig and is looking for a use case standpoint where they can relate to the problems of data ingestion, profiling, cleansing, transforming, and egressing data encountered in the enterprises. Knowledge of Hadoop and Pig is necessary for readers to grasp the intricacies of Pig design patterns better.

Find the right big data solution for your business or organization Big data management is one of the major challenges facing business, industry, and not-for-profit organizations. Data sets such as customer transactions for a mega-retailer, weather patterns monitored by meteorologists, or social network activity can quickly outpace the capacity of traditional data management tools. If you need to develop or manage big data solutions, you'll appreciate how these four experts define, explain, and guide you through this new and often confusing concept. You'll learn what it is, why it matters, and how to choose and implement solutions that work. Effectively managing big data is an issue of growing importance to businesses, not-for-profit organizations, government, and IT professionals. Authors are experts in information management, big data, and a variety of solutions. Explains big data in detail and discusses how to select and implement a solution, security concerns to consider, data storage and presentation issues, analytics, and much more. Provides essential information in a no-nonsense, easy-to-understand style that is empowering. Big Data For Dummies cuts through the confusion and helps you take charge of big data solutions for your organization.

This comprehensive book focuses on better big-data security for healthcare organizations. Following an extensive introduction to the Internet of Things (IoT) in healthcare including challenging topics and scenarios, it offers an in-depth analysis of medical body area networks with the 5th generation of IoT communication technology along with its nanotechnology. It also describes a novel strategic framework and computationally intelligent model to measure possible security

Read Free Beginning Apache Pig: Big Data Processing Made Easy

vulnerabilities in the context of e-health. Moreover, the book addresses healthcare systems that handle large volumes of data driven by patients' records and health/personal information, including big-data-based knowledge management systems to support clinical decisions. Several of the issues faced in storing/processing big data are presented along with the available tools, technologies and algorithms to deal with those problems as well as a case study in healthcare analytics. Addressing trust, privacy, and security issues as well as the IoT and big-data challenges, the book highlights the advances in the field to guide engineers developing different IoT devices and evaluating the performance of different IoT techniques. Additionally, it explores the impact of such technologies on public, private, community, and hybrid scenarios in healthcare. This book offers professionals, scientists and engineers the latest technologies, techniques, and strategies for IoT and big data.

[Hadoop in Practice](#)

[Resilience in the Digital Age](#)

[An Introduction for Data Scientists](#)

[Learn about big data processing and analytics](#)

[Proceedings of ICMDE 2020, Volume 2](#)

[Mastering Hadoop 3](#)

[Big Data Simplified.1e](#)

[Beginning Apache Spark 2](#)

[With Resilient Distributed Datasets, Spark SQL, Structured Streaming and Spark Machine Learning library](#)

[Hands-on Beginner's Guide on Big Data and Hadoop 3](#)

Beginning Apache Cassandra Development introduces you to one of the most robust and best-performing NoSQL database platforms on the planet. Apache Cassandra is a document database following the JSON document model. It is specifically designed to manage large amounts of data across many commodity servers

Read Free Beginning Apache Pig: Big Data Processing Made Easy

without there being any single point of failure. This design approach makes Apache Cassandra a robust and easy-to-implement platform when high availability is needed. Apache Cassandra can be used by developers in Java, PHP, Python, and JavaScript—the primary and most commonly used languages. In *Beginning Apache Cassandra Development*, author and Cassandra expert Vivek Mishra takes you through using Apache Cassandra from each of these primary languages. Mishra also covers the Cassandra Query Language (CQL), the Apache Cassandra analog to SQL. You'll learn to develop applications sourcing data from Cassandra, query that data, and deliver it at speed to your application's users. Cassandra is one of the leading NoSQL databases, meaning you get unparalleled throughput and performance without the sort of processing overhead that comes with traditional proprietary databases. *Beginning Apache Cassandra Development* will therefore help you create applications that generate search results quickly, stand up to high levels of demand, scale as your user base grows, ensure operational simplicity, and—not least—provide delightful user experiences. Finding patterns in massive event streams can be difficult, but learning how to find them doesn't have to be. This unique hands-on guide shows you how to solve this and many other problems in large-scale data processing with simple, fun, and elegant tools that leverage Apache Hadoop. You'll gain a

Read Free Beginning Apache Pig: Big Data Processing Made Easy

practical, actionable view of big data by working with real data and real problems. Perfect for beginners, this book's approach will also appeal to experienced practitioners who want to brush up on their skills. Part I explains how Hadoop and MapReduce work, while Part II covers many analytic patterns you can use to process any data. As you work through several exercises, you'll also learn how to use Apache Pig to process data. Learn the necessary mechanics of working with Hadoop, including how data and computation move around the cluster Dive into map/reduce mechanics and build your first map/reduce job in Python Understand how to run chains of map/reduce jobs in the form of Pig scripts Use a real-world dataset—baseball performance statistics—throughout the book Work with examples of several analytic patterns, and learn when and where you might use them 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

Read Free Beginning Apache Pig: Big Data Processing Made Easy

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

Read Free Beginning Apache Pig: Big Data Processing Made Easy

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

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

Read Free Beginning Apache Pig: Big Data Processing Made Easy

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.

Develop applications for the big data landscape with Spark and Hadoop. This book also explains the role of Spark in developing scalable machine learning and analytics applications with Cloud technologies.

Beginning Apache Spark 2 gives you an introduction to Apache Spark and shows you how to work with it. Along the way, you'll discover resilient distributed datasets (RDDs); use Spark SQL for structured data; and learn stream processing and build real-time applications with Spark Structured Streaming. Furthermore, you'll learn the fundamentals of Spark ML for machine learning and much more. After you read this book, you will have the fundamentals to become proficient in using Apache Spark and know when and how to apply it to your big data applications. What You Will Learn Understand Spark unified data processing platform How to run Spark in Spark Shell or Databricks Use and manipulate RDDs Deal with structured data using Spark SQL through its operations and advanced functions Build real-time applications using Spark Structured Streaming Develop intelligent applications with the Spark Machine Learning library Who This Book Is For Programmers and developers active in big data, Hadoop, and Java but who are new to

Read Free Beginning Apache Pig: Big Data Processing Made Easy

the Apache Spark platform.

Big Data Analytics with R and Hadoop is a tutorial style book that focuses on all the powerful big data tasks that can be achieved by integrating R and Hadoop. This book is ideal for R developers who are looking for a way to perform big data analytics with Hadoop. This book is also aimed at those who know Hadoop and want to build some intelligent applications over Big data with R packages. It would be helpful if readers have basic knowledge of R.

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

Read Free Beginning Apache Pig: Big Data Processing Made Easy

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.

"In this Learning Apache Pig training course, expert author Tom Hanlon will teach you how to explore, manipulate, and analyze data stored on a Hadoop cluster. This course is designed for the absolute beginner, meaning no experience with Pig is required. You will start by learning how to use Pig, then jump into learning about Pig and HCatalog. From there, Tom will teach you about advanced Pig, including Pig scripts, parameters in Pig scripts, and Pig and Oozie. Finally, this video tutorial will teach you about Pig user defined functions and streaming."--Resource description page.

[Computational Methods and Data Engineering \(Pig, Zookeeper and HBase\)](#)

[Big data processing made easy](#)

[Beginning Big Data with Power BI and Excel](#)

Read Free Beginning Apache Pig: Big Data Processing Made Easy

2013

[Hadoop in Action](#)

[Learn the Essentials of Big Data Computing in the Apache Hadoop 2 Ecosystem](#)

[Big Data Analytics](#)

[Hadoop in 24 Hours, Sams Teach Yourself](#)

[A Guide to Massive-Scale Data Processing in Practice](#)

[Processing Big Data with Azure HDInsight](#)

[All You Need to Know About Big Data](#)

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

Read Free Beginning Apache Pig: Big Data Processing Made Easy

*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*

*The book contains the latest trend in IT industry 'BigData
and Hadoop'. It explains how big is 'Big Data' and why
everybody is trying to implement this into their IT project. It
includes research work on various topics, theoretical and
practical approach, each component of the architecture is
described along with current industry trends. Big Data and
Hadoop have taken together are a new skill as per the
industry standards. Readers will get a compact book along
with the industry experience and would be a reference to
help readers. KEY FEATURES Overview Of Big Data,
Basics of Hadoop, Hadoop Distributed File System, HBase,
MapReduce, HIVE: The Dataware House Of Hadoop, PIG:
The Higher Level Programming Environment, SQOOP:
Importing Data From Heterogeneous Sources, Flume, Ozzie,
Zookeeper & Big Data Stream Mining, Chapter-wise
Questions & Previous Years Questions*

*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*

Read Free Beginning Apache Pig: Big Data Processing Made Easy

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 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

Read Free Beginning Apache Pig: Big Data Processing Made Easy

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.

Big data analytics emerged as a revolution in the field of information technology. It is the ability of the organization

Read Free Beginning Apache Pig: Big Data Processing Made Easy

to stay agile which gives it a competitive edge over its competitors. Data harvesting and data analytics enable the organization identify new opportunities which in turn results in efficient operations, leads to smarter business moves and higher business turnovers. All these issues are addressed by big data analytics and its initiatives. Chapter 4 focuses on architecture of Pig, Apache Pig execution modes, Pig data types and operators. Apache Pig Latin data model is based on nested relations. The chapter provides description of different components of Pig Latin data model. The lab session includes installing Pig over Hadoop and exploring different Pig Latin operators. Chapter 5 deals with common services provides by zookeeper, architecture and components of zookeeper and zookeeper operation modes. The salient feature of the chapter is exploration of leader election algorithm and security of ZNodes through access control list. The chapter concludes with the hands-on lab sessions on installation of zookeeper and exposure to zookeeper command-line interface. Chapter 6 discusses different types of NoSQL databases, transformation rules from one data model to another and performs in-depth analysis of HBase data model. The features which are difficult to comprehend such as data compaction, data locality, HBase read and write operations are simplified with easy to understand figures and explanation. As a part of hands-on lab sessions, installation of HBase over Hadoop and exercises based on HBase general commands, DDL commands and DML commands are dealt with.

Learn to use Apache Pig to develop lightweight big data

Read Free Beginning Apache Pig: Big Data Processing Made Easy

applications easily and quickly. This book shows you many optimization techniques and covers every context where Pig is used in big data analytics. Beginning Apache Pig shows you how Pig is easy to learn and requires relatively little time to develop big data applications. The book is divided into four parts: the complete features of Apache Pig; integration with other tools; how to solve complex business problems; and optimization of tools. You'll discover topics such as MapReduce and why it cannot meet every business need; the features of Pig Latin such as data types for each load, store, joins, groups, and ordering; how Pig workflows can be created; submitting Pig jobs using Hue; and working with Oozie. You'll also see how to extend the framework by writing UDFs and custom load, store, and filter functions. Finally you'll cover different optimization techniques such as gathering statistics about a Pig script, joining strategies, parallelism, and the role of data formats in good performance.

What You Will Learn

- Use all the features of Apache Pig
- Integrate Apache Pig with other tools
- Extend Apache Pig
- Optimize Pig Latin code
- Solve different use cases for Pig Latin

Who This Book Is For

All levels of IT professionals: architects, big data enthusiasts, engineers, developers, and big data administrators

Big Data represents a new era in data exploration and utilization, and IBM is uniquely positioned to help clients navigate this transformation. This book reveals how IBM is leveraging open source Big Data technology, infused with IBM technologies, to deliver a robust, secure, highly available, enterprise-class Big Data platform. The three

Read Free Beginning Apache Pig: Big Data Processing Made Easy

defining characteristics of Big Data--volume, variety, and velocity--are discussed. You'll get a primer on Hadoop and how IBM is hardening it for the enterprise, and learn when to leverage IBM InfoSphere BigInsights (Big Data at rest) and IBM InfoSphere Streams (Big Data in motion) technologies. Industry use cases are also included in this practical guide. Learn how IBM hardens Hadoop for enterprise-class scalability and reliability Gain insight into IBM's unique in-motion and at-rest Big Data analytics platform Learn tips and tricks for Big Data use cases and solutions Get a quick Hadoop primer

For many organizations, Hadoop is the first step for dealing with massive amounts of data. The next step? Processing and analyzing datasets with the Apache Pig scripting platform. With Pig, you can batch-process data without having to create a full-fledged application, making it easy to experiment with new datasets. Updated with use cases and programming examples, this second edition is the ideal learning tool for new and experienced users alike. You'll find comprehensive coverage on key features such as the Pig Latin scripting language and the Grunt shell. When you need to analyze terabytes of data, this book shows you how to do it efficiently with Pig. Delve into Pig's data model, including scalar and complex data types Write Pig Latin scripts to sort, group, join, project, and filter your data Use Grunt to work with the Hadoop Distributed File System (HDFS) Build complex data processing pipelines with Pig's macros and modularity features Embed Pig Latin in Python for iterative processing and other advanced tasks Use Pig with Apache

Read Free Beginning Apache Pig: Big Data Processing Made Easy

Tez to build high-performance batch and interactive data processing applications Create your own load and store functions to handle data formats and storage mechanisms

[*Data Analytics with Hadoop*](#)

[*BIG DATA*](#)

[*Big Data Processing and Analysis Using PowerBI in Excel 2013*](#)

[*Big data processing at scale to unlock unique business insights*](#)

[*Understanding Big Data: Analytics for Enterprise Class Hadoop and Streaming Data*](#)

[*Programming Hive*](#)

[*Hadoop Application Architectures*](#)

[*Dataflow Scripting with Hadoop*](#)

[*Pig Design Patterns*](#)

[*Big Data and Hadoop*](#)

[*Hadoop For Dummies*](#)

Get a jump start on using Azure HDInsight and Hadoop Ecosystem components. As most Hadoop and Big Data projects are written in either Java, Scala, or Python, this book minimizes the effort to learn another language and is written from the perspective of a .NET developer. Hadoop components are covered, including Hive, Pig, HBase, Storm, and Spark on Azure HDInsight, and code samples are written in .NET only. Processing Big Data with Azure HDInsight covers the fundamentals of big data, how businesses are using it to their advantage, and how Azure HDInsight fits into the big data world. This book introduces Hadoop and big data concepts and then dives

Read Free Beginning Apache Pig: Big Data Processing Made Easy

into creating different solutions with HDInsight and the Hadoop Ecosystem. It covers concepts with real-world scenarios and code examples, making sure you get hands-on experience. The best way to utilize this book is to practice while reading. After reading this book you will be familiar with Azure HDInsight and how it can be utilized to build big data solutions, including batch processing, stream analytics, interactive processing, and storing and retrieving data in an efficient manner. What You'll Learn Understand the fundamentals of HDInsight and Hadoop Work with HDInsight cluster Query with Apache Hive and Apache Pig Store and retrieve data with Apache HBase Stream data processing using Apache Storm Work with Apache Spark Who This Book Is For Software developers, technical architects, data scientists/analysts, and Hadoop administrators who want to develop on Microsoft's managed Hadoop offering, HDInsight

Beginning Apache Pig Big data processing made easy
Apress

Apache Hadoop is the technology at the heart of the Big Data revolution, and Hadoop skills are in enormous demand. Now, in just 24 lessons of one hour or less, you can learn all the skills and techniques you'll need to deploy each key component of a Hadoop platform in your local environment or in the cloud, building a fully functional Hadoop cluster and using it with real programs and datasets. Each short, easy lesson builds on all that's come before, helping you master all of Hadoop's essentials, and extend it to meet your unique challenges. Apache Hadoop

Read Free Beginning Apache Pig: Big Data Processing Made Easy

in 24 Hours, Sams Teach Yourself covers all this, and much more: Understanding Hadoop and the Hadoop Distributed File System (HDFS) Importing data into Hadoop, and process it there Mastering basic MapReduce Java programming, and using advanced MapReduce API concepts Making the most of Apache Pig and Apache Hive Implementing and administering YARN Taking advantage of the full Hadoop ecosystem Managing Hadoop clusters with Apache Ambari Working with the Hadoop User Environment (HUE) Scaling, securing, and troubleshooting Hadoop environments Integrating Hadoop into the enterprise Deploying Hadoop in the cloud Getting started with Apache Spark Step-by-step instructions walk you through common questions, issues, and tasks; Q-and-As, Quizzes, and Exercises build and test your knowledge; "Did You Know?" tips offer insider advice and shortcuts; and "Watch Out!" alerts help you avoid pitfalls. By the time you're finished, you'll be comfortable using Apache Hadoop to solve a wide spectrum of Big Data problems. Ready to use statistical and machine-learning techniques across large data sets? This practical guide shows you why the Hadoop ecosystem is perfect for the job. Instead of deployment, operations, or software development usually associated with distributed computing, you'll focus on particular analyses you can build, the data warehousing techniques that Hadoop provides, and higher order data workflows this framework can produce. Data scientists and analysts will learn how to perform a wide range of techniques, from writing MapReduce and Spark

Read Free Beginning Apache Pig: Big Data Processing Made Easy

applications with Python to using advanced modeling and data management with Spark MLlib, Hive, and HBase. You'll also learn about the analytical processes and data systems available to build and empower data products that can handle—and actually require—huge amounts of data. Understand core concepts behind Hadoop and cluster computing Use design patterns and parallel analytical algorithms to create distributed data analysis jobs Learn about data management, mining, and warehousing in a distributed context using Apache Hive and HBase Use Sqoop and Apache Flume to ingest data from relational databases Program complex Hadoop and Spark applications with Apache Pig and Spark DataFrames Perform machine learning techniques such as classification, clustering, and collaborative filtering with Spark's MLlib

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

[Big Data Analytics with R and Hadoop](#)

[Apache Hadoop YARN](#)

[Programming Pig](#)

[Hadoop Real-World Solutions Cookbook](#)

[Hadoop: The Definitive Guide](#)

[Beginning Apache Pig](#)

[Big Data Made Easy](#)

[Hadoop 2 Quick-Start Guide](#)

[Internet of Things and Big Data Technologies for Next Generation Healthcare](#)

[Explore, Manipulate, and Analyze Big Data in the Hadoop](#)

Read Free Beginning Apache Pig: Big Data Processing Made Easy

[Ecosystem](#)

[Big Data Tools – Which, When and How? \(Volume– II\)](#)