

Hadoop Operations And Cluster Management Cookbook

Getting the books Hadoop Operations And Cluster Management Cookbook now is not type of challenging means. You could not unaccompanied going subsequent to books buildup or library or borrowing from your connections to way in them. This is an certainly simple means to specifically get lead by on-line. This online broadcast Hadoop Operations And Cluster Management Cookbook can be one of the options to accompany you taking into account having supplementary time.

It will not waste your time. consent me, the e-book will certainly tell you new matter to read. Just invest little time to entry this on-line revelation Hadoop Operations And Cluster Management Cookbook as capably as review them wherever you are now.

Datenintensive Anwendungen designen Martin Kleppmann 2018-11-26 Daten stehen heute im Mittelpunkt vieler Herausforderungen im Systemdesign. Dabei sind komplexe Fragen wie Skalierbarkeit, Konsistenz, Zuverlässigkeit, Effizienz und Wartbarkeit zu klären. Darüber hinaus verfügen wir über eine überwältigende Vielfalt an Tools, einschließlich relationaler Datenbanken, NoSQL-Datenspeicher, Stream-und Batchprocessing und Message Broker. Aber was verbirgt sich hinter diesen Schlagworten? Und was ist die richtige Wahl für Ihre Anwendung? In diesem praktischen und umfassenden Leitfaden unterstützt Sie der Autor Martin Kleppmann bei der Navigation durch dieses schwierige Terrain, indem er die Vor- und Nachteile verschiedener Technologien zur Verarbeitung und Speicherung von Daten aufzeigt. Software verändert sich ständig, die Grundprinzipien bleiben aber gleich. Mit diesem Buch lernen Softwareentwickler und -architekten, wie sie die Konzepte in der Praxis umsetzen und wie sie Daten in modernen Anwendungen optimal nutzen können. Inspizieren Sie die Systeme, die Sie bereits verwenden, und erfahren Sie, wie Sie sie effektiver nutzen können Treffen Sie fundierte Entscheidungen, indem Sie die Stärken und Schwächen verschiedener Tools kennenlernen Steuern Sie die notwendigen Kompromisse in Bezug auf Konsistenz, Skalierbarkeit, Fehlertoleranz und Komplexität Machen Sie sich vertraut mit dem Stand der Forschung zu verteilten Systemen, auf denen moderne Datenbanken aufbauen Werfen Sie einen Blick hinter die Kulissen der wichtigsten Onlinedienste und lernen Sie von deren Architekturen

Learning Pyspark Tomasz Drabas 2017-04-28 Build real-time data intensive applications using the combined power of Python and Spark 2.0About This Book* Learn why and how you can efficiently use Python to implement various functionalities in Spark 2.0* Develop efficient, scalable real-time Spark solutions and deploy them* A comprehensive guide to take your understanding of implementing Spark with Python to the next levelWho This Book Is ForIf you are a Python developer who wants to learn about the Spark 2.0 ecosystem and how its functionalities can be implemented in Python, this book is for you. A firm understanding of Python is expected to get the best out of the book. Familiarity with Spark would be useful, but is not mandatory.What you will learn* Install, configure, and interact with Spark on a single machine* Build and interact with Spark DataFrames and Datasets using Spark SQL abstraction* Abstract various data sources with Blaze* Read, transform, and understand data and use it to train machine learning models* Familiarize yourself with the modeling pipeline capabilities of the machine learning module* Build machine learning models with MLlib* Package your application dependencies with spark-submit* Deploy locally built applications to clusterIn DetailApache Spark is an open source framework for efficient cluster computing with a strong interface for data parallelism and fault tolerance. This book will show you how you can leverage the power of Python and put it to use in the Spark ecosystem. You will start by getting a firm understanding of the Spark 2.0 architecture and how to set up a Python environment for Spark.Next you will get familiar with the PySpark packages and see how to get the data ready for processing. You will find out how to use the PySpark classes for RDD abstraction and Spark SQL abstraction, and understand streaming capabilities of Spark, machine learning using MLlib, polyglot persistence using Blaze, and graph processing using GraphX. Finally, you will see how you can configure Spark to deploy your applications on the cloud.By the end of this book, you will have established a firm understanding of the Spark Python API and how it can be used to build data-intensive applications.

Hadoop Real-World Solutions Cookbook Tanmay Deshpande 2016-03-31 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 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.

PySpark Recipes Raju Kumar Mishra 2017-12-09 Quickly find solutions to common programming problems encountered while processing big data. Content is presented in the popular problem-solution format. Look up the programming problem that you want to solve. Read the solution. Apply the solution directly in your own code. Problem solved! PySpark Recipes covers Hadoop and its shortcomings. The architecture of Spark, PySpark, and RDD are presented. You will learn to apply RDD to solve day-to-day big data problems. Python and NumPy are included and make it easy for new learners of PySpark to understand and adopt the model. What You Will Learn Understand the advanced features of PySpark2 and SparkSQL Optimize your code Program SparkSQL with Python Use Spark Streaming and Spark MLlib with Python Perform graph analysis with GraphFrames Who This Book Is For Data analysts, Python programmers, big data enthusiasts

PySpark Cookbook Denny Lee 2018-06-29 Combine the power of Apache Spark and Python to build effective big data applications Key Features Perform effective data processing, machine learning, and analytics using PySpark Overcome challenges in developing and deploying Spark solutions using Python Explore recipes for efficiently combining Python and Apache Spark to process data Book Description Apache Spark is an open source framework for efficient cluster computing with a strong interface for data parallelism and fault tolerance. The PySpark Cookbook presents effective and time-saving recipes for leveraging the power of Python and putting it to use in the Spark ecosystem. You'll start by learning the Apache Spark architecture and how to set up a Python environment for Spark. You'll then get familiar with the modules available in PySpark and start using them effortlessly. In addition to this, you'll discover how to abstract data with RDDs and DataFrames, and understand the streaming capabilities of PySpark. You'll then move on to using ML and MLlib in order to solve any problems related to the machine learning capabilities of PySpark and use GraphFrames to solve graph-processing problems. Finally, you will explore how to deploy your applications to the cloud using the spark-submit command. By the end of this book, you will be able to use the Python API for Apache Spark to solve any problems associated with building data-intensive applications. What you will learn Configure a local instance of PySpark in a virtual environment Install and configure Jupyter in local and multi-node environments Create DataFrames from JSON and a dictionary using pyspark.sql Explore regression and clustering models available in the ML module Use DataFrames to transform data used for modeling Connect to PubNub and perform aggregations on streams Who this book is for The PySpark Cookbook is for you if you are a Python developer looking for hands-on recipes for using the Apache Spark 2.x ecosystem in the best possible way. A thorough understanding of Python (and some familiarity with Spark) will help you get the best out of the book.

Applied OpenStack Design Patterns Uchit Vyas 2016-12-20 Learn practical and applied OpenStack cloud design solutions to gain maximum control over your infrastructure. You will achieve a complete controlled and customizable platform. Applied OpenStack Design Patterns teaches you how to map your application flow once you set up components and architectural design patterns. Also covered is storage management and computing to map user requests and allocations. Best practices of High Availability and Native Cluster Management are included. Solutions are presented to network components of OpenStack and to reduce latency and enable faster communication gateways between components of OpenStack and native applications. What You Will Learn: Design a modern cloud infrastructure Solve complex infrastructure application problems Understand OpenStack cloud infrastructure components Adopt a business impact analysis to support existing/new cloud infrastructure Use specific components to integrate an existing tool-chain set to gain agility and a quick, continuous delivery model Who This Book Is For: Seasoned solution architects, DevOps, and system engineers and analysts

Data Science für Dummies Lillian Pierson 2016-04-22 Daten, Daten, Daten? Sie haben schon Kenntnisse in Excel und Statistik, wissen aber noch nicht, wie all die Datensätze helfen sollen, bessere Entscheidungen zu treffen? Von Lillian Pierson bekommen Sie das dafür notwendige Handwerkszeug: Bauen Sie Ihre Kenntnisse in Statistik, Programmierung und Visualisierung aus. Nutzen Sie Python, R, SQL, Excel und KNIME. Zahlreiche Beispiele veranschaulichen die vorgestellten Methoden und Techniken. So können Sie die Erkenntnisse dieses Buches auf Ihre Daten übertragen und aus deren Analyse unmittelbare Schlüsse und Konsequenzen ziehen.

Big Data Management Peter Ghavami 2020-11-09 Data analytics is core to business and decision making. The rapid increase in data volume, velocity and variety offers both opportunities and challenges. While open source solutions to store big data, like Hadoop, offer platforms for exploring value and insight from big data, they were not originally developed with data security and governance in mind. Big Data Management discusses numerous policies, strategies and recipes for managing big data. It addresses data security, privacy, controls and life cycle management offering modern principles and open source architectures for successful governance of big data. The author has collected best practices from the world's leading organizations that have successfully implemented big data platforms. The topics discussed cover the entire data management life cycle, data quality, data stewardship, regulatory considerations, data council, architectural and operational models are presented for successful management of big data. The book is a must-read for data scientists, data engineers and corporate leaders who are implementing big data platforms in their organizations.

Apache ZooKeeper Essentials Saurav Haloi 2015-01-28 Whether you are a novice to ZooKeeper or already have some experience, you will be able to master the concepts of ZooKeeper and its usage with ease. This book assumes you to have some prior knowledge of distributed systems and high-level programming knowledge of C, Java, or Python, but no experience with Apache ZooKeeper is required.

Machine Learning Kochbuch Chris Albon 2019-03-22 Python-Programmierer finden in diesem Kochbuch nahezu 200 wertvolle und jeweils in sich abgeschlossene Anleitungen zu Aufgabenstellungen aus dem Bereich des Machine Learning, wie sie für die tägliche Arbeit typisch sind – von der Vorverarbeitung der Daten bis zum Deep Learning. Entwickler, die mit Python und seinen Bibliotheken einschließlich Pandas und Scikit-Learn vertraut sind, werden spezifische Probleme erfolgreich bewältigen – wie etwa Daten laden, Text und numerische Daten behandeln, Modelle auswählen, Dimensionalität reduzieren und vieles mehr. Jedes Rezept enthält Code, den Sie kopieren, zum Testen in eine kleine Beispieldatenmenge einfügen und dann anpassen können, um Ihre eigenen Anwendungen zu konstruieren. Darüber hinaus werden alle Lösungen diskutiert und wichtige Zusammenhänge hergestellt. Dieses Kochbuch unterstützt Sie dabei, den Schritt von der Theorie und den Konzepten hinein in die Praxis zu machen. Es liefert das praktische Rüstzeug, das Sie benötigen, um funktionierende Machine-Learning-Anwendungen zu entwickeln. In diesem Kochbuch finden Sie Rezepte für: Vektoren, Matrizen und Arrays den Umgang mit numerischen und kategorischen Daten, Texten, Bildern sowie Datum und Uhrzeit das Reduzieren der Dimensionalität durch Merkmalsextraktion oder Merkmalsauswahl Modellbewertung und -auswahl lineare und logistische Regression, Bäume und Wälder und k-nächste Nachbarn Support Vector Machine (SVM),

naive Bayes, Clustering und neuronale Netze das Speichern und Laden von trainierten Modellen

Einführung in SQL Alan Beaulieu 2009-08-31 SQL kann Spaß machen! Es ist ein erhebendes Gefühl, eine verworrene Datenmanipulation oder einen komplizierten Report mit einer einzigen Anweisung zu bewältigen und so einen Haufen Arbeit vom Tisch zu bekommen. Einführung in SQL bietet einen frischen Blick auf die Sprache, deren Grundlagen jeder Entwickler beherrschen muss. Die aktualisierte 2. Auflage deckt die Versionen MySQL 6.0, Oracle 11g und Microsoft SQL Server 2008 ab. Außerdem enthält sie neue Kapitel zu Views und Metadaten. SQL-Basics - in null Komma nichts durchstarten: Mit diesem leicht verständlichen Tutorial können Sie SQL systematisch und gründlich lernen, ohne sich zu langweilen. Es führt Sie rasch durch die Basics der Sprache und vermittelt darüber hinaus eine Reihe von häufig genutzten fortgeschrittenen Features. Mehr aus SQL-Befehlen herausholen: Alan Beaulieu will mehr vermitteln als die simple Anwendung von SQL-Befehlen: Er legt Wert auf ein tiefes Verständnis der SQL-Features und behandelt daher auch den Umgang mit Mengen, Abfragen innerhalb von Abfragen oder die überaus nützlichen eingebauten Funktionen von SQL. Die MySQL-Beispieldatenbank: Es gibt zwar viele Datenbankprodukte auf dem Markt, aber welches wäre zum Erlernen von SQL besser geeignet als MySQL, das weit verbreitete relationale Datenbanksystem? Der Autor hilft Ihnen, eine MySQL-Datenbank anzulegen, und nutzt diese für die Beispiele in diesem Buch. Übungen mit Lösungen: Zu jedem Thema finden Sie im Buch gut durchdachte Übungen mit Lösungen. So ist sichergestellt, dass Sie schnell Erfolgserlebnisse haben und das Gelernte auch praktisch umsetzen können.

PySpark SQL Recipes Raju Kumar Mishra 2019-03-18 Carry out data analysis with PySpark SQL, graphframes, and graph data processing using a problem-solution approach. This book provides solutions to problems related to dataframes, data manipulation summarization, and exploratory analysis. You will improve your skills in graph data analysis using graphframes and see how to optimize your PySpark SQL code. PySpark SQL Recipes starts with recipes on creating dataframes from different types of data source, data aggregation and summarization, and exploratory data analysis using PySpark SQL. You'll also discover how to solve problems in graph analysis using graphframes. On completing this book, you'll have ready-made code for all your PySpark SQL tasks, including creating dataframes using data from different file formats as well as from SQL or NoSQL databases. What You Will Learn Understand PySpark SQL and its advanced features Use SQL and HiveQL with PySpark SQL Work with structured streaming Optimize PySpark SQL Master graphframes and graph processing Who This Book Is For Data scientists, Python programmers, and SQL programmers.

Apache Sqoop Cookbook Kathleen Ting 2013-07-02 Integrating data from multiple sources is essential in the age of big data, but it can be a challenging and time-consuming task. This handy cookbook provides dozens of ready-to-use recipes for using Apache Sqoop, the command-line interface application that optimizes data transfers between relational databases and Hadoop. Sqoop is both powerful and bewildering, but with this cookbook's problem-solution-discussion format, you'll quickly learn how to deploy and then apply Sqoop in your environment. The authors provide MySQL, Oracle, and PostgreSQL database examples on GitHub that you can easily adapt for SQL Server, Netezza, Teradata, or other relational systems. Transfer data from a single database table into your Hadoop ecosystem Keep table data and Hadoop in sync by importing data incrementally Import data from more than one database table Customize transferred data by calling various database functions Export generated, processed, or backed-up data from Hadoop to your database Run Sqoop within Oozie, Hadoop's specialized workflow scheduler Load data into Hadoop's data warehouse (Hive) or database (HBase) Handle installation, connection, and syntax issues common to specific database vendors Learning PySpark Tomasz Drabas 2017-02-27 Build data-intensive applications locally and deploy at scale using the combined powers of Python and Spark 2.0 About This Book Learn why and how you can efficiently use Python to process data and build machine learning models in Apache Spark 2.0 Develop and deploy efficient, scalable real-time Spark solutions Take your understanding of using Spark with Python to the next level with this jump start guide Who This Book Is For If you are a Python developer who wants to learn about the Apache Spark 2.0 ecosystem, this book is for you. A firm understanding of Python is expected to get the best out of the book. Familiarity with Spark would be useful, but is not mandatory. What You Will Learn Learn about Apache Spark and the Spark 2.0 architecture Build and interact with Spark DataFrames using Spark SQL Learn how to solve graph and deep learning problems using GraphFrames and TensorFrames respectively Read, transform, and understand data and use it to train machine learning models Build machine learning models with MLlib and ML Learn how to submit your applications programmatically using spark-submit Deploy locally built applications to a cluster In Detail Apache Spark is an open source framework for efficient cluster computing with a strong interface for data parallelism and fault tolerance. This book will show you how to leverage the power of Python and put it to use in the Spark ecosystem. You will start by getting a firm understanding of the Spark 2.0 architecture and how to set up a Python environment for Spark. You will get familiar with the modules available in PySpark. You will learn how to abstract data with RDDs and DataFrames and understand the streaming capabilities of PySpark. Also, you will get a thorough overview of machine learning capabilities of PySpark using ML and MLlib, graph processing using GraphFrames, and polyglot persistence using Blaze. Finally, you will learn how to deploy your applications to the cloud using the spark-submit command. By the end of this book, you will have established a firm understanding of the Spark Python API and how it can be used to build data-intensive applications. Style and approach This book takes a very comprehensive, step-by-step approach so you understand how the Spark ecosystem can be used with Python to develop efficient, scalable solutions. Every chapter is standalone and written in a very easy-to-understand manner, with a focus on both the hows and the whys of each concept.

HBase High Performance Cookbook Ruchir Choudhry 2017-01-31 Exciting projects that will teach you how complex data can be exploited to gain maximum insights About This Book Architect a good HBase cluster for a very large distributed system Get to grips with the concepts of performance tuning with HBase A practical guide full of engaging recipes and attractive screenshots to enhance your system's performance Who This Book Is For This book is intended for developers and architects who want to know all about HBase at a hands-on level. This book is also for big data enthusiasts and database developers who have worked with other NoSQL databases and now want to explore HBase as another futuristic scalable database solution in the big data space. What You Will Learn Configure HBase from a high performance perspective Grab data from various RDBMS/Flat files into the HBASE systems Understand table design and perform CRUD operations Find out how the communication between the client and server happens in HBase Grasp when to use and avoid MapReduce and how to perform various tasks with it Get to know the concepts of scaling with HBase through practical examples Set up Hbase in the Cloud for a small scale environment Integrate HBase with other tools including Elasticsearch In Detail Apache HBase is a non-relational NoSQL database management system that runs on top of HDFS. It is an open source, distributed, versioned, column-oriented store and is written in Java to provide random real-time access to big Data. We'll start off by ensuring you have a solid understanding the basics of HBase, followed by giving you a thorough explanation of architecting a HBase cluster as per our project specifications. Next, we will explore the scalable structure

of tables and we will be able to communicate with the HBase client. After this, we'll show you the intricacies of MapReduce and the art of performance tuning with HBase. Following this, we'll explain the concepts pertaining to scaling with HBase. Finally, you will get an understanding of how to integrate HBase with other tools such as Elasticsearch. By the end of this book, you will have learned enough to exploit HBase for boost system performance. Style and approach This book is intended for software quality assurance/testing professionals, software project managers, or software developers with prior experience in using Selenium and Java to test web-based applications. This books also provides examples for C#, Python, and Ruby users.

Practical Data Analysis Hector Cuesta 2013-10-22 Each chapter of the book quickly introduces a key 'theme' of Data Analysis, before immersing you in the practical aspects of each theme. You'll learn quickly how to perform all aspects of Data Analysis. Practical Data Analysis is a book ideal for home and small business users who want to slice & dice the data they have on hand with minimum hassle.

Scala: Guide for Data Science Professionals Pascal Bugnion 2017-02-24 Scala will be a valuable tool to have on hand during your data science journey for everything from data cleaning to cutting-edge machine learning About This Book Build data science and data engineering solutions with ease An in-depth look at each stage of the data analysis process — from reading and collecting data to distributed analytics Explore a broad variety of data processing, machine learning, and genetic algorithms through diagrams, mathematical formulations, and source code Who This Book Is For This learning path is perfect for those who are comfortable with Scala programming and now want to enter the field of data science. Some knowledge of statistics is expected. What You Will Learn Transfer and filter tabular data to extract features for machine learning Read, clean, transform, and write data to both SQL and NoSQL databases Create Scala web applications that couple with JavaScript libraries such as D3 to create compelling interactive visualizations Load data from HDFS and HIVE with ease Run streaming and graph analytics in Spark for exploratory analysis Bundle and scale up Spark jobs by deploying them into a variety of cluster managers Build dynamic workflows for scientific computing Leverage open source libraries to extract patterns from time series Master probabilistic models for sequential data In Detail Scala is especially good for analyzing large sets of data as the scale of the task doesn't have any significant impact on performance. Scala's powerful functional libraries can interact with databases and build scalable frameworks — resulting in the creation of robust data pipelines. The first module introduces you to Scala libraries to ingest, store, manipulate, process, and visualize data. Using real world examples, you will learn how to design scalable architecture to process and model data — starting from simple concurrency constructs and progressing to actor systems and Apache Spark. After this, you will also learn how to build interactive visualizations with web frameworks. Once you have become familiar with all the tasks involved in data science, you will explore data analytics with Scala in the second module. You'll see how Scala can be used to make sense of data through easy to follow recipes. You will learn about Bokeh bindings for exploratory data analysis and quintessential machine learning with algorithms with Spark ML library. You'll get a sufficient understanding of Spark streaming, machine learning for streaming data, and Spark graphX. Armed with a firm understanding of data analysis, you will be ready to explore the most cutting-edge aspect of data science — machine learning. The final module teaches you the A to Z of machine learning with Scala. You'll explore Scala for dependency injections and implicits, which are used to write machine learning algorithms. You'll also explore machine learning topics such as clustering, dimensionality reduction, Naive Bayes, Regression models, SVMs, neural networks, and more. This learning path combines some of the best that Packt has to offer into one complete, curated package. It includes content from the following Packt products: Scala for Data Science, Pascal Bugnion Scala Data Analysis Cookbook, Arun Manivannan Scala for Machine Learning, Patrick R. Nicolas Style and approach A complete package with all the information necessary to start building useful data engineering and data science solutions straight away. It contains a diverse set of recipes that cover the full spectrum of interesting data analysis tasks and will help you revolutionize your data analysis skills using Scala.

Apache Spark Quick Start Guide Shrey Mehrotra 2019-01-31 A practical guide for solving complex data processing challenges by applying the best optimizations techniques in Apache Spark. Key Features Learn about the core concepts and the latest developments in Apache Spark Master writing efficient big data applications with Spark's built-in modules for SQL, Streaming, Machine Learning and Graph analysis Get introduced to a variety of optimizations based on the actual experience Book Description Apache Spark is a flexible framework that allows processing of batch and real-time data. Its unified engine has made it quite popular for big data use cases. This book will help you to get started with Apache Spark 2.0 and write big data applications for a variety of use cases. It will also introduce you to Apache Spark — one of the most popular Big Data processing frameworks. Although this book is intended to help you get started with Apache Spark, but it also focuses on explaining the core concepts. This practical guide provides a quick start to the Spark 2.0 architecture and its components. It teaches you how to set up Spark on your local machine. As we move ahead, you will be introduced to resilient distributed datasets (RDDs) and DataFrame APIs, and their corresponding transformations and actions. Then, we move on to the life cycle of a Spark application and learn about the techniques used to debug slow-running applications. You will also go through Spark's built-in modules for SQL, streaming, machine learning, and graph analysis. Finally, the book will lay out the best practices and optimization techniques that are key for writing efficient Spark applications. By the end of this book, you will have a sound fundamental understanding of the Apache Spark framework and you will be able to write and optimize Spark applications. What you will learn Learn core concepts such as RDDs, DataFrames, transformations, and more Set up a Spark development environment Choose the right APIs for your applications Understand Spark's architecture and the execution flow of a Spark application Explore built-in modules for SQL, streaming, ML, and graph analysis Optimize your Spark job for better performance Who this book is for If you are a big data enthusiast and love processing huge amount of data, this book is for you. If you are data engineer and looking for the best optimization techniques for your Spark applications, then you will find this book helpful. This book also helps data scientists who want to implement their machine learning algorithms in Spark. You need to have a basic understanding of any one of the programming languages such as Scala, Python or Java.

ETL with Azure Cookbook Christian Coté 2020-09-30 Explore the latest Azure ETL techniques both on-premises and in the cloud using Azure services such as SQL Server Integration Services (SSIS), Azure Data Factory, and Azure Databricks Key Features Understand the key components of an ETL solution using Azure Integration Services Discover the common and not-so-common challenges faced while creating modern and scalable ETL solutions Program and extend your packages to develop efficient data integration and data transformation solutions Book Description ETL is one of the most common and tedious procedures for moving and processing data from one database to another. With the help of this book, you will be able to speed up the process by designing effective ETL solutions using the Azure services available for handling and transforming any data to suit your requirements. With this cookbook, you'll become well versed in all the features of SQL Server Integration Services (SSIS) to

perform data migration and ETL tasks that integrate with Azure. You'll learn how to transform data in Azure and understand how legacy systems perform ETL on-premises using SSIS. Later chapters will get you up to speed with connecting and retrieving data from SQL Server 2019 Big Data Clusters, and even show you how to extend and customize the SSIS toolbox using custom-developed tasks and transforms. This ETL book also contains practical recipes for moving and transforming data with Azure services, such as Data Factory and Azure Databricks, and lets you explore various options for migrating SSIS packages to Azure. Toward the end, you'll find out how to profile data in the cloud and automate service creation with Business Intelligence Markup Language (BIML). By the end of this book, you'll have developed the skills you need to create and automate ETL solutions on-premises as well as in Azure. What you will learn

Explore ETL and how it is different from ELT
Move and transform various data sources with Azure ETL and ELT services
Use SSIS 2019 with Azure HDInsight clusters
Discover how to query SQL Server 2019 Big Data Clusters hosted in Azure
Migrate SSIS solutions to Azure and solve key challenges associated with it
Understand why data profiling is crucial and how to implement it in Azure Databricks
Get to grips with BIML and learn how it applies to SSIS and Azure Data Factory solutions
Who this book is for
This book is for data warehouse architects, ETL developers, or anyone who wants to build scalable ETL applications in Azure. Those looking to extend their existing on-premise ETL applications to use big data and a variety of Azure services or others interested in migrating existing on-premise solutions to the Azure cloud platform will also find the book useful. Familiarity with SQL Server services is necessary to get the most out of this book.

Programmieren mit Scala Dean Wampler 2010-10-31 Sie ist elegant, schlank, modern und flexibel: Die Rede ist von Scala, der neuen Programmiersprache für die Java Virtual Machine (JVM). Sie vereint die Vorzüge funktionaler und objektorientierter Programmierung, ist typischer als Java, lässt sich nahtlos in die Java-Welt integrieren – und eine in Scala entwickelte Anwendung benötigt oft nur einen Bruchteil der Codezeilen ihres Java-Pendants. Kein Wunder, dass immer mehr Firmen, deren große, geschäftskritische Anwendungen auf Java basieren, auf Scala umsteigen, um ihre Produktivität und die Skalierbarkeit ihrer Software zu erhöhen. Das wollen Sie auch? Dann lassen Sie sich von den Scala-Profis Dean Wampler und Alex Payne zeigen, wie es geht. Ihre Werkzeugkiste: Schon bevor Sie loslegen, sind Sie weiter, als Sie denken: Sie können Ihre Java-Programme weiter verwenden, Java-Bibliotheken nutzen, Java von Scala aus aufrufen und Scala von Java aus. Auch Ihre bevorzugten Entwicklungswerkzeuge wie NetBeans, IntelliJ IDEA oder Eclipse stehen Ihnen weiter zur Verfügung, dazu Kommandozeilen-Tools, Plugins für Editoren, Werkzeuge von Drittanbietern – und natürlich Ihre Programmiererfahrung. In Programmieren mit Scala erfahren Sie, wie Sie sich all das zunutze machen. Das Hybridmodell: Die Paradigmen "funktional" und "objektorientiert" sind keine Gegensätze, sondern ergänzen sich unter dem Scala-Dach zu einem sehr produktiven Ganzen. Nutzen Sie die Vorteile funktionaler Programmierung, wann immer sich das anbietet – und seien Sie so frei, auf die guten alten Seiteneffekte zu bauen, wenn Sie das für nötig halten. Futter für die Profis: Skalierbare Nebenläufigkeit mit Aktoren, Aufzucht und Pflege von XML mit Scala, Domainspezifische Sprachen, Tipps zum richtigen Anwendungsdesign – das sind nur ein paar der fortgeschrittenen Themen, in die Sie mit den beiden Autoren eintauchen. Danach sind Sie auch Profi im Programmieren mit Scala.

Hadoop MapReduce v2 Cookbook - Second Edition Thilina Gunarathne 2015-02-25 If you are a Big Data enthusiast and wish to use Hadoop v2 to solve your problems, then this book is for you. This book is for Java programmers with little to moderate knowledge of Hadoop MapReduce. This is also a one-stop reference for developers and system admins who want to quickly get up to speed with using Hadoop v2. It would be helpful to have a basic knowledge of software development using Java and a basic working knowledge of Linux.

Hadoop Operations and Cluster Management Cookbook Shumin Guo 2013 Solve specific problems using individual self-contained code recipes, or work through the book to develop your capabilities. This book is packed with easy-to-follow code and commands used for illustration, which makes your learning curve easy and quick. If you are a Hadoop cluster system administrator with Unix/Linux system management experience and you are looking to get a good grounding in how to set up and manage a Hadoop cluster, then this book is for you. It's assumed that you will have some experience in Unix/Linux command line already, as well as being familiar with network communication basics.

Einführung in Apache Solr Markus Klose 2014-02

SQL Server 2017 Integration Services Cookbook Christian Cote 2017-06-30 Harness the power of SQL Server 2017 Integration Services to build your data integration solutions with ease About This Book Acquaint yourself with all the newly introduced features in SQL Server 2017 Integration Services Program and extend your packages to enhance their functionality This detailed, step-by-step guide covers everything you need to develop efficient data integration and data transformation solutions for your organization Who This Book Is For This book is ideal for software engineers, DW/ETL architects, and ETL developers who need to create a new, or enhance an existing, ETL implementation with SQL Server 2017 Integration Services. This book would also be good for individuals who develop ETL solutions that use SSIS and are keen to learn the new features and capabilities in SSIS 2017. What You Will Learn Understand the key components of an ETL solution using SQL Server 2016-2017 Integration Services Design the architecture of a modern ETL solution Have a good knowledge of the new capabilities and features added to Integration Services Implement ETL solutions using Integration Services for both on-premises and Azure data Improve the performance and scalability of an ETL solution Enhance the ETL solution using a custom framework Be able to work on the ETL solution with many other developers and have common design paradigms or techniques Effectively use scripting to solve complex data issues In Detail SQL Server Integration Services is a tool that facilitates data extraction, consolidation, and loading options (ETL), SQL Server coding enhancements, data warehousing, and customizations. With the help of the recipes in this book, you'll gain complete hands-on experience of SSIS 2017 as well as the 2016 new features, design and development improvements including SCD, Tuning, and Customizations. At the start, you'll learn to install and set up SSIS as well other SQL Server resources to make optimal use of this Business Intelligence tools. We'll begin by taking you through the new features in SSIS 2016/2017 and implementing the necessary features to get a modern scalable ETL solution that fits the modern data warehouse. Through the course of chapters, you will learn how to design and build SSIS data warehouses packages using SQL Server Data Tools. Additionally, you'll learn to develop SSIS packages designed to maintain a data warehouse using the Data Flow and other control flow tasks. You'll also be demonstrated many recipes on cleansing data and how to get the end result after applying different transformations. Some real-world scenarios that you might face are also covered and how to handle various issues that you might face when designing your packages. At the end of this book, you'll get to know all the key concepts to perform data integration and transformation. You'll have explored on-premises Big Data integration processes to create a classic data warehouse, and will know how to extend the toolbox with custom tasks and transforms. Style and approach This cookbook follows a problem-solution approach and tackles all kinds of data integration scenarios by using the capabilities of SQL Server 2016 Integration Services. This book is well supplemented with

screenshots, tips, and tricks. Each recipe focuses on a particular task and is written in a very easy-to-follow manner.

Lernen mit Big Data Viktor Mayer-Schönberger 2014-08-15 Was heute noch undenkbar scheint, ist morgen schon Alltag – sprechende Übungsbücher, Schulaufgaben, die von den Schülern lernen. Schneller als gedacht, wird Big Data Einzug in Schulen und Klassenzimmer halten, so die These der beiden Experten und Erfolgsautoren Viktor Mayer-Schönberger und Kenneth Cukier. Und damit das Schulsystem und das Lernen von Grund auf verändern. Die beiden Autoren von Big Data erklären, welche Neuheiten uns erwarten. Und zeigen, dass es nicht nur positiv ist, den Fortschritt der Schüler und Studenten immer besser messen zu können. Vor lauter PISA und Rankings bleibt oft das Wesentliche auf der Strecke – eine gute Bildung. Die Gefahr ist, dass das Lernen von der Quantität der Daten dominiert wird, und nicht von der Qualität, von Kreativität oder von Ideen. Sie plädieren daher eindringlich dafür, unsere Bildungssysteme schnellstens zukunftsfähig zu machen.

Apache Mesos Cookbook David Blomquist 2017-04-28 Over 50 recipes on the core features of Apache Mesos and running big data frameworks in Mesos About This Book* Learn to install and configure Mesos to suit the needs of your organization* Follow step-by-step instructions to deploy application frameworks on top of Mesos, saving you many hours of research and trial and error* Use this practical guide packed with powerful recipes to implement Mesos and easily integrate it with other application frameworks Who This Book Is For This book is for system administrators, engineers, and big data programmers. Basic experience with big data technologies such as Hadoop or Spark would be useful but is not essential. A working knowledge of Apache Mesos is expected. What you will learn* Set up Mesos on different operating systems* Use the Marathon and Chronos frameworks to manage multiple applications* Work with Mesos and Docker* Integrate Mesos with Spark and other big data frameworks* Use networking features in Mesos for effective communication between containers* Configure Mesos for high availability using Zookeeper* Secure your Mesos clusters with SASL and Authorization ACLs* Solve everyday problems and discover the best practices In Detail Apache Mesos is open source cluster sharing and management software. Deploying and managing scalable applications in large-scale clustered environments can be difficult, but Apache Mesos makes it easier with efficient resource isolation and sharing across application frameworks. The goal of this book is to guide you through the practical implementation of the Mesos core along with a number of Mesos supported frameworks. You will begin by installing Mesos and then learn how to configure clusters and maintain them. You will also see how to deploy a cluster in a production environment with high availability using Zookeeper. Next, you will get to grips with using Mesos, Marathon, and Docker to build and deploy a PaaS. You will see how to schedule jobs with Chronos. We'll demonstrate how to integrate Mesos with big data frameworks such as Spark, Hadoop, and Storm. Practical solutions backed with clear examples will also show you how to deploy elastic big data jobs. You will find out how to deploy a scalable continuous integration and delivery system on Mesos with Jenkins. Finally, you will configure and deploy a highly scalable distributed search engine with ElasticSearch. Throughout the course of this book, you will get to know tips and tricks along with best practices to follow when working with Mesos.

Spring Recipes Daniel Rubio 2014-11-14 Spring Recipes: A Problem-Solution Approach, Third Edition builds upon the best-selling success of the previous editions and focuses on the latest Spring Framework features for building enterprise Java applications. This book provides code recipes for the following, found in the latest Spring: Spring fundamentals: Spring IoC container, Spring AOP/ AspectJ, and more. Spring enterprise: Spring Java EE integration, Spring Integration, Spring Batch, Spring Remoting, messaging, transactions, and working with big data and the cloud using Hadoop and MongoDB. Spring web: Spring MVC, other dynamic scripting, integration with the popular Grails Framework (and Groovy), REST/web services, and more This book guides you step-by-step through topics using complete and real-world code examples. When you start a new project, you can consider copying the code and configuration files from this book, and then modifying them for your needs. This can save you a great deal of work over creating a project from scratch!

Mastering Large Datasets with Python John Wolohan 2020-01-15 Summary Modern data science solutions need to be clean, easy to read, and scalable. In Mastering Large Datasets with Python, author J.T. Wolohan teaches you how to take a small project and scale it up using a functionally influenced approach to Python coding. You'll explore methods and built-in Python tools that lend themselves to clarity and scalability, like the high-performing parallelism method, as well as distributed technologies that allow for high data throughput. The abundant hands-on exercises in this practical tutorial will lock in these essential skills for any large-scale data science project. Purchase of the print book includes a free eBook in PDF, Kindle, and ePub formats from Manning Publications. About the technology Programming techniques that work well on laptop-sized data can slow to a crawl—or fail altogether—when applied to massive files or distributed datasets. By mastering the powerful map and reduce paradigm, along with the Python-based tools that support it, you can write data-centric applications that scale efficiently without requiring codebase rewrites as your requirements change. About the book Mastering Large Datasets with Python teaches you to write code that can handle datasets of any size. You'll start with laptop-sized datasets that teach you to parallelize data analysis by breaking large tasks into smaller ones that can run simultaneously. You'll then scale those same programs to industrial-sized datasets on a cluster of cloud servers. With the map and reduce paradigm firmly in place, you'll explore tools like Hadoop and PySpark to efficiently process massive distributed datasets, speed up decision-making with machine learning, and simplify your data storage with AWS S3.

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

Apache Sqoop Cookbook Kathleen Ting 2013-07-02 Integrating data from multiple sources is essential in the age of big data, but it can be a challenging and time-consuming task. This handy cookbook provides dozens of ready-to-use recipes for using Apache Sqoop, the command-line interface application that optimizes data transfers between relational databases and Hadoop. Sqoop is both powerful and bewildering, but with this cookbook's problem-solution-discussion format, you'll quickly learn how to deploy and then apply Sqoop in your environment. The authors provide MySQL, Oracle, and PostgreSQL database examples on GitHub that you can easily adapt for SQL Server, Netezza, Teradata, or other relational systems. Transfer data from a single database table

into your Hadoop ecosystem Keep table data and Hadoop in sync by importing data incrementally Import data from more than one database table Customize transferred data by calling various database functions Export generated, processed, or backed-up data from Hadoop to your database Run Sqoop within Oozie, Hadoop's specialized workflow scheduler Load data into Hadoop's data warehouse (Hive) or database (HBase) Handle installation, connection, and syntax issues common to specific database vendors Sieben Wochen, sieben Datenbanken Eric Redmond 2012

Arduino Kochbuch Michael Margolis 2012-08-31 Mit dem Arduino-Kochbuch, das auf der Version Arduino 1.0 basiert, erhalten Sie ein Füllhorn an Ideen und praktischen Beispielen, was alles mit dem Mikrocontroller gezaubert werden kann. Sie lernen alles über die Arduino-Softwareumgebung, digitale und analoge In- und Outputs, Peripheriegeräte, Motorensteuerung und fortgeschrittenes Arduino-Coding. Egal ob es ein Spielzeug, ein Detektor, ein Roboter oder ein interaktives Kleidungsstück werden soll: Elektronikbegeisterte finden über 200 Rezepte, Projekte und Techniken, um mit dem Arduino zu starten oder bestehende Arduino-Projekt mit neuen Features aufzupimpen.

Böse Julia Shaw 2018-09-24 Von Psychopathen wie Charles Manson oder Serienmördern wie Jack the Ripper geht eine unheimliche Faszination aus. Doch woher kommt sie? Und warum verdrängen wir so gern das alltäglichere Böse – von den eigenen Gewaltphantasien bis zum Machtmissbrauch im Büro? Die Kriminalpsychologin und Bestsellerautorin Julia Shaw taucht das Phänomen des Bösen in neues Licht. Shaw sucht und findet das Böse nicht nur in den Gehirnen von Massenmördern, sondern in jedem von uns. Und sie erläutert mithilfe psychologischer Fallstudien und neuester neurowissenschaftlicher Erkenntnisse, wie wir uns mit unserer dunklen Seite versöhnen. Ein augenöffnendes Buch, das die vertrauten Kategorien von Gut und Böse völlig über den Haufen wirft.

HDInsight Essentials - Second Edition Rajesh Nadipalli 2015-01-27 If you want to discover one of the latest tools designed to produce stunning Big Data insights, this book features everything you need to get to grips with your data. Whether you are a data architect, developer, or a business strategist, HDInsight adds value in everything from development, administration, and reporting.

Die Berechnung der Zukunft Nate Silver 2013-09-02 Zuverlässige Vorhersagen sind doch möglich! Nate Silver ist der heimliche Gewinner der amerikanischen Präsidentschaftswahlen 2012: ein begnadeter Statistiker, als »Prognose-Popstar« und »Wundernerd« weltberühmt geworden. Er hat die Wahlergebnisse aller 50 amerikanischen Bundesstaaten absolut exakt vorausgesagt – doch damit nicht genug: Jetzt zeigt Nate Silver, wie seine Prognosen in Zukunft Terroranschläge, Umweltkatastrophen und Finanzkrisen verhindern sollen. Gelingt ihm die Abschaffung des Zufalls? Warum werden Wettervorhersagen immer besser, während die Terrorattacken vom 11.09.2001 niemand kommen sah? Warum erkennen Ökonomen eine globale Finanzkrise nicht einmal dann, wenn diese bereits begonnen hat? Das Problem ist nicht der Mangel an Informationen, sondern dass wir die verfügbaren Daten nicht richtig deuten. Zuverlässige Prognosen aber würden uns helfen, Zufälle und Ungewissheiten abzuwehren und unser Schicksal selbst zu bestimmen. Nate Silver zeigt, dass und wie das geht. Erstmals wendet er seine Wahrscheinlichkeitsrechnung nicht nur auf Wahlprognosen an, sondern auf die großen Probleme unserer Zeit: die Finanzmärkte, Ratingagenturen, Epidemien, Erdbeben, den Klimawandel, den Terrorismus. In all diesen Fällen gibt es zahlreiche Prognosen von Experten, die er überprüft – und erklärt, warum sie meist falsch sind. Gleichzeitig schildert er, wie es gelingen kann, im Rauschen der Daten die wesentlichen Informationen herauszufiltern. Ein unterhaltsamer und spannender Augenöffner!

Mastering Python Scientific Computing Hemant Kumar Mehta 2015-09-23 A complete guide for Python programmers to master scientific computing using Python APIs and tools About This Book The basics of scientific computing to advanced concepts involving parallel and large scale computation are all covered. Most of the Python APIs and tools used in scientific computing are discussed in detail The concepts are discussed with suitable example programs Who This Book Is For If you are a Python programmer and want to get your hands on scientific computing, this book is for you. The book expects you to have had exposure to various concepts of Python programming. What You Will Learn Fundamentals and components of scientific computing Scientific computing data management Performing numerical computing using NumPy and SciPy Concepts and programming for symbolic computing using SymPy Using the plotting library matplotlib for data visualization Data analysis and visualization using Pandas, matplotlib, and IPython Performing parallel and high performance computing Real-life case studies and best practices of scientific computing In Detail In today's world, along with theoretical and experimental work, scientific computing has become an important part of scientific disciplines. Numerical calculations, simulations and computer modeling in this day and age form the vast majority of both experimental and theoretical papers. In the scientific method, replication and reproducibility are two important contributing factors. A complete and concrete scientific result should be reproducible and replicable. Python is suitable for scientific computing. A large community of users, plenty of help and documentation, a large collection of scientific libraries and environments, great performance, and good support makes Python a great choice for scientific computing. At present Python is among the top choices for developing scientific workflow and the book targets existing Python developers to master this domain using Python. The main things to learn in the book are the concept of scientific workflow, managing scientific workflow data and performing computation on this data using Python. The book discusses NumPy, SciPy, SymPy, matplotlib, Pandas and IPython with several example programs. Style and approach This book follows a hands-on approach to explain the complex concepts related to scientific computing. It details various APIs using appropriate examples.

High Performance Websites Steve Souders 2008

Apache Kafka 1.0 Cookbook Raúl Estrada 2017-12-22 Simplify real-time data processing by leveraging the power of Apache Kafka 1.0 Key Features Use Kafka 1.0 features such as Confluent platforms and Kafka streams to build efficient streaming data applications to handle and process your data Integrate Kafka with other Big Data tools such as Apache Hadoop, Apache Spark, and more Hands-on recipes to help you design, operate, maintain, and secure your Apache Kafka cluster with ease Book Description Apache Kafka provides a unified, high-throughput, low-latency platform to handle real-time data feeds. This book will show you how to use Kafka efficiently, and contains practical solutions to the common problems that developers and administrators usually face while working with it. This practical guide contains easy-to-follow recipes to help you set up, configure, and use Apache Kafka in the best possible manner. You will use Apache Kafka Consumers and Producers to build effective real-time streaming applications. The book covers the recently released Kafka version 1.0, the Confluent Platform and Kafka Streams. The programming aspect covered in the book will teach you how to perform important tasks such as message validation, enrichment and composition. Recipes focusing on optimizing the performance of your Kafka cluster, and integrate Kafka with a variety of third-party tools such as Apache Hadoop, Apache Spark, and Elasticsearch will help ease your day to day collaboration with Kafka greatly. Finally, we cover tasks related to monitoring and securing your Apache Kafka cluster using tools such as Ganglia and Graphite. If you're looking to become the go-to person in your organization when it comes to working with Apache Kafka, this book is the only resource you need to have.

What you will learn -Install and configure Apache Kafka 1.0 to get optimal performance -Create and configure Kafka Producers and Consumers -Operate your Kafka clusters efficiently by implementing the mirroring technique -Work with the new Confluent platform and Kafka streams, and achieve high availability with Kafka -Monitor Kafka using tools such as Graphite and Ganglia -Integrate Kafka with third-party tools such as Elasticsearch, Logstash, Apache Hadoop, Apache Spark, and more Who this book is for This book is for developers and Kafka administrators who are looking for quick, practical solutions to problems encountered while operating, managing or monitoring Apache Kafka. If you are a developer, some knowledge of Scala or Java will help, while for administrators, some working knowledge of Kafka will be useful.

Hadoop 2.x Administration Cookbook Gurmukh Singh 2017-05-26 Over 100 practical recipes to help you become an expert Hadoop administrator About This Book Become an expert Hadoop administrator and perform tasks to optimize your Hadoop Cluster Import and export data into Hive and use Oozie to manage workflow. Practical recipes will help you plan and secure your Hadoop cluster, and make it highly available Who This Book Is For If you are a system administrator with a basic understanding of Hadoop and you want to get into Hadoop administration, this book is for you. It's also ideal if you are a Hadoop administrator who wants a quick reference guide to all the Hadoop administration-related tasks and solutions to commonly occurring problems What You Will Learn Set up the Hadoop architecture to run a Hadoop cluster smoothly Maintain a Hadoop cluster on HDFS, YARN, and MapReduce Understand high availability with Zookeeper and Journal Node Configure Flume for data ingestion and Oozie to run various workflows Tune the Hadoop cluster for optimal performance Schedule jobs on a Hadoop cluster using the Fair and Capacity scheduler Secure your cluster and troubleshoot it for various common pain points In Detail Hadoop enables the distributed storage and processing of large datasets across clusters of computers. Learning how to administer Hadoop is crucial to exploit its unique features. With this book, you will be able to overcome common problems encountered in Hadoop administration. The book begins with laying the foundation by showing you the steps needed to set up a Hadoop cluster and its various nodes. You will get a better understanding of how to maintain Hadoop cluster, especially on the HDFS layer and using YARN and MapReduce. Further on, you will explore durability and high availability of a Hadoop cluster. You'll get a better understanding of the schedulers in Hadoop and how to configure and use them for your tasks. You will also get hands-on experience with the backup and recovery options and the performance tuning aspects of Hadoop. Finally, you will get a better understanding of troubleshooting, diagnostics, and best practices in Hadoop administration. By the end of this book, you will have a proper understanding of working with Hadoop clusters and will also be able to secure, encrypt it, and configure auditing for your Hadoop clusters. Style and approach This book contains short recipes that will help you run a Hadoop cluster efficiently. The recipes are solutions to real-life problems that administrators encounter while working with a Hadoop cluster

Raspberry Pi ??????(???) Richard Grimmett 2014-09-19 ??????????????????????
??Raspberry
Pi??
??GPS??
??GPS????
??GPS????
??#???? GOTOP Information Inc.

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