

Apache Spark Hands On Session Uniroma2

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 the Apache Spark platform.

HadoopExam Learning Resources (www.HadoopExam.com). Provides many learning resources for Hadoop , BigData , Data Science and Analytics certifications as well as technical Books. We have following training's and books. 1. Hadoop Professional Training with Hands On sessions. 2. Apache Spark Professional Training with Hands On sessions. 3. Apache Pig Professional Training and Books. 4. Apache Hive Professional Training 5. Apache HBase training and Book

Hands-On Deep Learning with Apache SparkBuild and deploy distributed deep learning applications on Apache SparkPackt Publishing Ltd

Apache Spark is one of the fastest growing technology in BigData computing world. It support multiple programming languages like Java, Scala, Python and R. Hence, many existing and new framework started to integrate Spark platform as well in their platform e.g. Hadoop, Cassandra, EMR etc. While creating Spark certification material HadoopExam technical team found that there is no proper material and book is available for the Spark SQL (version 2.x) which covers the concepts as well as use of various features and found difficulty in creating the material. Therefore, they decided to create full length book for Spark SQL and outcome of that is this book. In this book technical team try to cover both fundamental concepts of Spark SQL engine and many exercises approx. 35+ so that most of the programming features can be covered. There are approximately 35 exercises and total 15 chapters which covers the programming aspects of SparkSQL. All the exercises given in this book are written using Scala. However, concepts remain same even if you are using different programming language.

A solution-based guide to put your deep learning models into production with the power of Apache Spark Key Features Discover practical recipes for distributed deep learning with Apache Spark Learn to use libraries such as Keras and TensorFlow Solve problems in order to train your deep learning models on Apache Spark Book Description With deep learning gaining rapid

mainstream adoption in modern-day industries, organizations are looking for ways to unite popular big data tools with highly efficient deep learning libraries. As a result, this will help deep learning models train with higher efficiency and speed. With the help of the Apache Spark Deep Learning Cookbook, you'll work through specific recipes to generate outcomes for deep learning algorithms, without getting bogged down in theory. From setting up Apache Spark for deep learning to implementing types of neural net, this book tackles both common and not so common problems to perform deep learning on a distributed environment. In addition to this, you'll get access to deep learning code within Spark that can be reused to answer similar problems or tweaked to answer slightly different problems. You will also learn how to stream and cluster your data with Spark. Once you have got to grips with the basics, you'll explore how to implement and deploy deep learning models, such as Convolutional Neural Networks (CNN) and Recurrent Neural Networks (RNN) in Spark, using popular libraries such as TensorFlow and Keras. By the end of the book, you'll have the expertise to train and deploy efficient deep learning models on Apache Spark. What you will learn

- Set up a fully functional Spark environment
- Understand practical machine learning and deep learning concepts
- Apply built-in machine learning libraries within Spark
- Explore libraries that are compatible with TensorFlow and Keras
- Explore NLP models such as Word2vec and TF-IDF on Spark
- Organize dataframes for deep learning evaluation
- Apply testing and training modeling to ensure accuracy
- Access readily available code that may be reusable

Who this book is for If you're looking for a practical and highly useful resource for implementing efficiently distributed deep learning models with Apache Spark, then the Apache Spark Deep Learning Cookbook is for you. Knowledge of the core machine learning concepts and a basic understanding of the Apache Spark framework is required to get the best out of this book. Additionally, some programming knowledge in Python is a plus.

Data is bigger, arrives faster, and comes in a variety of formats—and it all needs to be processed at scale for analytics or machine learning. But how can you process such varied workloads efficiently? Enter Apache Spark. Updated to include Spark 3.0, this second edition shows data engineers and data scientists why structure and unification in Spark matters. Specifically, this book explains how to perform simple and complex data analytics and employ machine learning algorithms. Through step-by-step walk-throughs, code snippets, and notebooks, you'll be able to:

- Learn Python, SQL, Scala, or Java high-level Structured APIs
- Understand Spark operations and SQL Engine
- Inspect, tune, and debug Spark operations with Spark configurations and Spark UI
- Connect to data sources: JSON, Parquet, CSV, Avro, ORC, Hive, S3, or Kafka
- Perform analytics on batch and streaming data using Structured Streaming
- Build reliable data pipelines with open source Delta Lake and Spark
- Develop machine learning pipelines with MLlib and productionize models using MLflow

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.

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.

Apache Spark is a flexible in-memory framework that allows processing of both batch and real-time data. Its unified engine has made it quite popular for big data use cases. This book will help you to quickly get started with Apache Spark 2.0 and write efficient big data applications for a variety of use cases.

Apache® Spark is one of the fastest growing technology in BigData computing world. It supports multiple programming languages like Java, Scala, Python and R. Hence, many existing and new framework started to integrate Spark platform as well in their platform e.g. Hadoop, Cassandra, EMR etc. While creating Spark certification material HadoopExam technical team found that there is no proper material and book is available for the Spark (version 2.x) which covers the concepts as well as use of various features and found difficulty in creating the material. Therefore, they decided to create full length book for Spark (Databricks® CRT020 Spark Scala/Python or PySpark Certification) and outcome of that is this book. In this book technical team try to cover both fundamental concepts of Spark 2.x topics which are part of the certification syllabus as well as add as many exercises as possible and in current version we have around 46 hands on exercises added which you can execute on the Databricks community edition, because each of this exercises tested on that platform as well, as this book is focused on the Scala version of the certification, hence all the exercises and their solution provided in the Scala. We have divided the entire book in the 13 chapters, as you move ahead chapter by chapter you would be comfortable with the Databricks Spark Scala certification (CRT020). All the exercises given in this book are written using Scala. However, concepts remain same even if you are using different programming language. This book explains the key feature to develop a complex and stable network that helps to gather the data to optimize the asset performance and maximize the production in the Industries leveraging on the cloud infrastructure and services. By the end, you can design the Industrial IoT network and the architecture for processing its data in the cloud.

Frank Kane's hands-on Spark training course, based on his bestselling Taming Big Data with Apache Spark and Python video, now available in a book. Understand and analyze large data sets using Spark on a single system or on a cluster. About This Book Understand how Spark can be distributed across computing clusters Develop and run Spark jobs efficiently using Python A hands-on tutorial by Frank Kane with over 15 real-world examples teaching you Big Data processing with Spark Who This Book Is For If you are a data scientist or data analyst who wants to learn Big Data

processing using Apache Spark and Python, this book is for you. If you have some programming experience in Python, and want to learn how to process large amounts of data using Apache Spark, Frank Kane's *Taming Big Data with Apache Spark and Python* will also help you. What You Will Learn Find out how you can identify Big Data problems as Spark problems Install and run Apache Spark on your computer or on a cluster Analyze large data sets across many CPUs using Spark's Resilient Distributed Datasets Implement machine learning on Spark using the MLlib library Process continuous streams of data in real time using the Spark streaming module Perform complex network analysis using Spark's GraphX library Use Amazon's Elastic MapReduce service to run your Spark jobs on a cluster In Detail Frank Kane's *Taming Big Data with Apache Spark and Python* is your companion to learning Apache Spark in a hands-on manner. Frank will start you off by teaching you how to set up Spark on a single system or on a cluster, and you'll soon move on to analyzing large data sets using Spark RDD, and developing and running effective Spark jobs quickly using Python. Apache Spark has emerged as the next big thing in the Big Data domain – quickly rising from an ascending technology to an established superstar in just a matter of years. Spark allows you to quickly extract actionable insights from large amounts of data, on a real-time basis, making it an essential tool in many modern businesses. Frank has packed this book with over 15 interactive, fun-filled examples relevant to the real world, and he will empower you to understand the Spark ecosystem and implement production-grade real-time Spark projects with ease. Style and approach Frank Kane's *Taming Big Data with Apache Spark and Python* is a hands-on tutorial with over 15 real-world examples carefully explained by Frank in a step-by-step manner. The examples vary in complexity, and you can move through them at your own pace.

If you're like most R users, you have deep knowledge and love for statistics. But as your organization continues to collect huge amounts of data, adding tools such as Apache Spark makes a lot of sense. With this practical book, data scientists and professionals working with large-scale data applications will learn how to use Spark from R to tackle big data and big compute problems. Authors Javier Luraschi, Kevin Kuo, and Edgar Ruiz show you how to use R with Spark to solve different data analysis problems. This book covers relevant data science topics, cluster computing, and issues that should interest even the most advanced users. Analyze, explore, transform, and visualize data in Apache Spark with R Create statistical models to extract information and predict outcomes; automate the process in production-ready workflows Perform analysis and modeling across many machines using distributed computing techniques Use large-scale data from multiple sources and different formats with ease from within Spark Learn about alternative modeling frameworks for graph processing, geospatial analysis, and genomics at scale Dive into advanced topics including custom transformations, real-time data processing, and creating custom Spark extensions

Master scala's advanced techniques to solve real-world problems in data analysis and gain valuable insights from your data Key Features A beginner's guide for performing data analysis loaded with numerous rich, practical examples Access to popular Scala libraries such as Breeze, Saddle for efficient data manipulation and exploratory analysis Develop applications in Scala for real-time analysis and machine learning in Apache Spark Book Description Efficient business decisions with an accurate sense of business data helps in delivering better performance across products and services. This book helps you to leverage the popular Scala libraries and tools for performing core data analysis tasks with ease. The book begins with a quick overview of the building blocks of a standard data analysis process. You will learn to perform basic tasks like Extraction, Staging, Validation, Cleaning, and Shaping of datasets. You will later deep dive into the data exploration and visualization areas of the data analysis life cycle. You will make use of popular Scala libraries like Saddle, Breeze, Vegas, and PredictionIO for processing your datasets. You will learn statistical methods for deriving meaningful insights from data. You will also learn to create applications for Apache Spark 2.x on complex data analysis, in real-time. You will discover traditional machine learning techniques for doing data analysis. Furthermore, you will also be introduced to neural networks and deep learning from a data analysis standpoint. By the end of this book, you will be capable of handling large sets of structured and unstructured data, perform exploratory analysis, and building efficient Scala applications for discovering and delivering insights What you will learn Techniques to determine the validity and confidence level of data Apply quartiles and n-tiles to datasets to see how data is distributed into many buckets Create data pipelines that combine multiple data lifecycle steps Use built-in features to gain a deeper understanding of the data Apply Lasso regression analysis method to your data Compare Apache Spark API with traditional Apache Spark data analysis Who this book is for If you are a data scientist or a data analyst who wants to learn how to perform data analysis using Scala, this book is for you. All you need is knowledge of the basic fundamentals of Scala programming.

CCA175 , CCP DE575

Gain the key language concepts and programming techniques of Scala in the context of big data analytics and Apache Spark. The book begins by introducing you to Scala and establishes a firm contextual understanding of why you should learn this language, how it stands in comparison to Java, and how Scala is related to Apache Spark for big data analytics. Next, you'll set up the Scala environment ready for examining your first Scala programs. This is followed by sections on Scala fundamentals including mutable/immutable variables, the type hierarchy system, control flow expressions and code blocks. The author discusses functions at length and highlights a number of associated concepts such as functional programming and anonymous functions. The book then delves deeper into Scala's powerful collections system because many of Apache Spark's APIs bear a strong resemblance to Scala collections. Along the

way you'll see the development life cycle of a Scala program. This involves compiling and building programs using the industry-standard Scala Build Tool (SBT). You'll cover guidelines related to dependency management using SBT as this is critical for building large Apache Spark applications. Scala Programming for Big Data Analytics concludes by demonstrating how you can make use of the concepts to write programs that run on the Apache Spark framework. These programs will provide distributed and parallel computing, which is critical for big data analytics. What You Will Learn See the fundamentals of Scala as a general-purpose programming language Understand functional programming and object-oriented programming constructs in Scala Use Scala collections and functions Develop, package and run Apache Spark applications for big data analytics Who This Book Is For Data scientists, data analysts and data engineers who intend to use Apache Spark for large-scale analytics. /div

Build efficient data flow and machine learning programs with this flexible, multi-functional open-source cluster-computing framework Key Features Master the art of real-time big data processing and machine learning Explore a wide range of use-cases to analyze large data Discover ways to optimize your work by using many features of Spark 2.x and Scala Book Description Apache Spark is an in-memory, cluster-based data processing system that provides a wide range of functionalities such as big data processing, analytics, machine learning, and more. With this Learning Path, you can take your knowledge of Apache Spark to the next level by learning how to expand Spark's functionality and building your own data flow and machine learning programs on this platform. You will work with the different modules in Apache Spark, such as interactive querying with Spark SQL, using DataFrames and datasets, implementing streaming analytics with Spark Streaming, and applying machine learning and deep learning techniques on Spark using MLlib and various external tools. By the end of this elaborately designed Learning Path, you will have all the knowledge you need to master Apache Spark, and build your own big data processing and analytics pipeline quickly and without any hassle. This Learning Path includes content from the following Packt products: Mastering Apache Spark 2.x by Romeo Kienzler Scala and Spark for Big Data Analytics by Md. Rezaul Karim, Sridhar Alla Apache Spark 2.x Machine Learning Cookbook by Siamak Amirghodsi, Meenakshi Rajendran, Broderick Hall, Shuen Mei Cookbook What you will learn Get to grips with all the features of Apache Spark 2.x Perform highly optimized real-time big data processing Use ML and DL techniques with Spark MLlib and third-party tools Analyze structured and unstructured data using SparkSQL and GraphX Understand tuning, debugging, and monitoring of big data applications Build scalable and fault-tolerant streaming applications Develop scalable recommendation engines Who this book is for If you are an intermediate-level Spark developer looking to master the advanced capabilities and use-cases of Apache Spark 2.x, this Learning Path is ideal for you. Big data professionals who want to learn how to integrate and use the features of Apache Spark and build a strong big data

pipeline will also find this Learning Path useful. To grasp the concepts explained in this Learning Path, you must know the fundamentals of Apache Spark and Scala.

In this book, you'll learn to implement some practical and proven techniques to improve aspects of programming and administration in Apache Spark. Techniques are demonstrated using practical examples and best practices. You will also learn how to use Spark and its Python API to create performant analytics with large-scale data.

Simplify machine learning model implementations with Spark About This Book Solve the day-to-day problems of data science with Spark This unique cookbook consists of exciting and intuitive numerical recipes Optimize your work by acquiring, cleaning, analyzing, predicting, and visualizing your data Who This Book Is For This book is for Scala developers with a fairly good exposure to and understanding of machine learning techniques, but lack practical implementations with Spark. A solid knowledge of machine learning algorithms is assumed, as well as hands-on experience of implementing ML algorithms with Scala. However, you do not need to be acquainted with the Spark ML libraries and ecosystem. What You Will Learn Get to know how Scala and Spark go hand-in-hand for developers when developing ML systems with Spark Build a recommendation engine that scales with Spark Find out how to build unsupervised clustering systems to classify data in Spark Build machine learning systems with the Decision Tree and Ensemble models in Spark Deal with the curse of high-dimensionality in big data using Spark Implement Text analytics for Search Engines in Spark Streaming Machine Learning System implementation using Spark In Detail Machine learning aims to extract knowledge from data, relying on fundamental concepts in computer science, statistics, probability, and optimization. Learning about algorithms enables a wide range of applications, from everyday tasks such as product recommendations and spam filtering to cutting edge applications such as self-driving cars and personalized medicine. You will gain hands-on experience of applying these principles using Apache Spark, a resilient cluster computing system well suited for large-scale machine learning tasks. This book begins with a quick overview of setting up the necessary IDEs to facilitate the execution of code examples that will be covered in various chapters. It also highlights some key issues developers face while working with machine learning algorithms on the Spark platform. We progress by uncovering the various Spark APIs and the implementation of ML algorithms with developing classification systems, recommendation engines, text analytics, clustering, and learning systems. Toward the final chapters, we'll focus on building high-end applications and explain various unsupervised methodologies and challenges to tackle when implementing with big data ML systems. Style and approach This book is packed with intuitive recipes supported with line-by-line explanations to help you understand how to optimize your work flow and resolve problems when working with complex data modeling tasks and predictive algorithms. This is a valuable resource for data scientists and those working on large scale data projects.

Before you can build analytics tools to gain quick insights, you first need to know how to process data in real time. With this practical guide, developers familiar with Apache Spark will learn how to put this in-memory framework to use for streaming data. You'll discover how Spark enables you to write streaming jobs in almost the same way you write batch jobs. Authors Gerard Maas and François Garillot help you explore the theoretical underpinnings of Apache Spark. This comprehensive guide features two sections that compare and contrast the streaming APIs Spark now supports: the original Spark Streaming library and the newer Structured Streaming API. Learn fundamental stream processing concepts and examine different streaming architectures Explore Structured Streaming through practical examples; learn different aspects of stream processing in detail Create and operate streaming jobs and applications with Spark Streaming; integrate Spark Streaming

with other Spark APIs Learn advanced Spark Streaming techniques, including approximation algorithms and machine learning algorithms Compare Apache Spark to other stream processing projects, including Apache Storm, Apache Flink, and Apache Kafka Streams Learn how to use, deploy, and maintain Apache Spark with this comprehensive guide, written by the creators of the open-source cluster-computing framework. With an emphasis on improvements and new features in Spark 2.0, authors Bill Chambers and Matei Zaharia break down Spark topics into distinct sections, each with unique goals. You'll explore the basic operations and common functions of Spark's structured APIs, as well as Structured Streaming, a new high-level API for building end-to-end streaming applications. Developers and system administrators will learn the fundamentals of monitoring, tuning, and debugging Spark, and explore machine learning techniques and scenarios for employing MLlib, Spark's scalable machine-learning library. Get a gentle overview of big data and Spark Learn about DataFrames, SQL, and Datasets—Spark's core APIs—through worked examples Dive into Spark's low-level APIs, RDDs, and execution of SQL and DataFrames Understand how Spark runs on a cluster Debug, monitor, and tune Spark clusters and applications Learn the power of Structured Streaming, Spark's stream-processing engine Learn how you can apply MLlib to a variety of problems, including classification or recommendation

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 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 Deep Learning is a subset of Machine Learning where data sets with several layers of complexity can be processed. This book teaches you the different techniques using which deep learning solutions can be implemented at scale, on Apache Spark. This will help you gain experience of implementing your deep learning models in many real-world use cases.

Explore and work with various Microsoft Azure services for real-time Data Analytics KEY FEATURES Understanding what Azure can do with your data Understanding the analytics services offered by Azure Understand how data can be transformed to generate more data Understand what is done after a Machine Learning model is built Go through some Data Analytics real-world use cases DESCRIPTION Data is the key input for Analytics. Building and implementing data platforms such as Data Lakes, modern Data Marts, and Analytics at scale require the right cloud platform that Azure provides through its services. The book starts by sharing how analytics has evolved and continues to evolve. Following the introduction, you will deep dive into ingestion technologies. You will learn about Data processing services in Azure. You will next learn about what is meant by a Data Lake and understand how Azure Data Lake Storage is used for analytical workloads. You will then learn about critical services that will provide actual Machine Learning capabilities in Azure. The book also talks about Azure Data Catalog for cataloging, Azure AD for Access Management, Web Apps and PowerApps for cloud web applications, Cognitive services for Speech, Vision,

Search and Language, Azure VM for computing and Data Science VMs, Functions as serverless computing, Kubernetes and Containers as deployment options. Towards the end, the book discusses two use cases on Analytics. WHAT WILL YOU LEARN Explore and work with various Azure services Orchestrate and ingest data using Azure Data Factory Learn how to use Azure Stream Analytics Get to know more about Synapse Analytics and its features Learn how to use Azure Analysis Services and its functionalities WHO THIS BOOK IS FOR This book is for anyone who has basic to intermediate knowledge of cloud and analytics concepts and wants to use Microsoft Azure for Data Analytics. This book will also benefit Data Scientists who want to use Azure for Machine Learning. TABLE OF CONTENTS 1. Data and its power 2. Evolution of Analytics and its Types 3. Internet of Things 4. AI and ML 5. Why cloud 6. What are a data lake and a modern datamart 7. Introduction to Azure services 8. Types of data 9. Azure Data Factory 10. Stream Analytics 11. Azure Data Lake Store and Azure Storage 12. Cosmos DB 13. Synapse Analytics 14. Azure Databricks 15. Azure Analysis Services 16. Power BI 17. Azure Machine Learning 18. Sample Architectures and synergies - Real-Time and Batch 19. Azure Data Catalog 20. Azure Active Directory 21. Azure Webapps 22. Power apps 23. Time Series Insights 24. Azure Cognitive Services 25. Azure Logicapps 26. Azure VM 27. Azure Functions 28. Azure Containers 29. Azure Kubernetes Service 30. Use Case 1 31. Use Case 2

If you want to build an enterprise-quality application that uses natural language text but aren't sure where to begin or what tools to use, this practical guide will help get you started. Alex Thomas, principal data scientist at Wisecube, shows software engineers and data scientists how to build scalable natural language processing (NLP) applications using deep learning and the Apache Spark NLP library. Through concrete examples, practical and theoretical explanations, and hands-on exercises for using NLP on the Spark processing framework, this book teaches you everything from basic linguistics and writing systems to sentiment analysis and search engines. You'll also explore special concerns for developing text-based applications, such as performance. In four sections, you'll learn NLP basics and building blocks before diving into application and system building: Basics: Understand the fundamentals of natural language processing, NLP on Apache Spark, and deep learning Building blocks: Learn techniques for building NLP applications—including tokenization, sentence segmentation, and named-entity recognition—and discover how and why they work Applications: Explore the design, development, and experimentation process for building your own NLP applications Building NLP systems: Consider options for productionizing and deploying NLP models, including which human languages to support

Apache Spark is amazing when everything clicks. But if you haven't seen the performance improvements you expected, or still don't feel confident enough to use Spark in production, this practical book is for you. Authors Holden Karau and Rachel Warren demonstrate performance optimizations to help your Spark queries run faster and handle larger data sizes, while using fewer resources. Ideal for software engineers, data engineers, developers, and system administrators working with large-scale data applications, this book describes techniques that can reduce data infrastructure costs and developer hours. Not only will you gain a more comprehensive understanding of Spark, you'll also learn how to make it sing. With this book, you'll explore: How Spark SQL's new interfaces improve performance over SQL's RDD data structure The choice between data joins in Core Spark and Spark SQL Techniques for getting the most out of standard RDD transformations How to work around performance issues in Spark's key/value pair paradigm Writing high-performance Spark code without Scala or the JVM How to test for functionality and performance when applying suggested improvements Using Spark MLlib and Spark ML machine learning libraries Spark's Streaming components and external community packages

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.

Powerful smart applications using deep learning algorithms to dominate numerical computing, deep learning, and functional programming. Key Features Explore machine learning techniques with prominent open source Scala libraries such as Spark ML, H2O, MXNet, Zeppelin, and DeepLearning4j Solve real-world machine learning problems by delving complex numerical computing with Scala functional programming in a scalable and faster way Cover all key aspects such as collection, storing, processing, analyzing, and evaluation required to build and deploy machine models on computing clusters using Scala Play framework. Book Description Machine learning has had a huge impact on academia and industry by turning data into actionable information. Scala has seen a steady rise in adoption over the past few years, especially in the fields of data science and analytics. This book is for data scientists, data engineers, and deep learning enthusiasts who have a background in complex numerical computing and want to know more hands-on machine learning application development. If you're well versed in machine learning concepts and want to expand your knowledge by delving into the practical implementation of these concepts using the power of Scala, then this book is what you need! Through 11 end-to-end projects, you will be acquainted with popular machine learning libraries such as Spark ML, H2O, DeepLearning4j, and MXNet. At the end, you will be able to use numerical computing and functional programming to carry out complex numerical tasks to develop, build, and deploy research or commercial projects in a production-ready environment. What you will learn Apply advanced regression techniques to boost the performance of predictive models Use different classification algorithms for business analytics Generate trading strategies for Bitcoin and stock trading using ensemble techniques Train Deep Neural Networks (DNN) using H2O and Spark ML Utilize NLP to build scalable machine learning models Learn how to apply reinforcement learning algorithms such as Q-learning for developing ML application Learn how to use autoencoders to develop a fraud detection application Implement LSTM and CNN models using DeepLearning4j and MXNet Who this book is for If you want to leverage the power of both Scala and Spark to make sense of Big Data, then this book is for you. If you are well versed with machine learning concepts and wants to expand your knowledge by delving into the practical implementation using the power of Scala, then this book is what you need!

Strong understanding of Scala Programming language is recommended. Basic familiarity with machine Learning techniques will be more helpful.

Develop robust, Scala-powered projects with the help of machine learning libraries such as SparkML to harvest meaningful insight
Key Features Gain hands-on experience in building data science projects with Scala Exploit powerful functionalities of machine learning libraries
Use machine learning algorithms and decision tree models for enterprise apps
Book Description Scala, together with the Spark Framework, forms a rich and powerful data processing ecosystem. Modern Scala Projects is a journey into the depths of this ecosystem. The machine learning (ML) projects presented in this book enable you to create practical, robust data analytics solutions, with an emphasis on automating data workflows with the Spark ML pipeline API. This book showcases or carefully cherry-picks from Scala's functional libraries and other constructs to help readers roll out their own scalable data processing frameworks. The projects in this book enable data practitioners across all industries gain insights into data that will help organizations have strategic and competitive advantage. Modern Scala Projects focuses on the application of supervisory learning ML techniques that classify data and make predictions. You'll begin with working on a project to predict a class of flower by implementing a simple machine learning model. Next, you'll create a cancer diagnosis classification pipeline, followed by projects delving into stock price prediction, spam filtering, fraud detection, and a recommendation engine. By the end of this book, you will be able to build efficient data science projects that fulfil your software requirements. What you will learn Create pipelines to extract data or analytics and visualizations Automate your process pipeline with jobs that are reproducible Extract intelligent data efficiently from large, disparate datasets Automate the extraction, transformation, and loading of data Develop tools that collate, model, and analyze data Maintain the integrity of data as data flows become more complex Develop tools that predict outcomes based on "pattern discovery" Build really fast and accurate machine-learning models in Scala Who this book is for Modern Scala Projects is for Scala developers who would like to gain some hands-on experience with some interesting real-world projects. Prior programming experience with Scala is necessary.

Implement machine learning, cognitive services, and artificial intelligence solutions by leveraging Azure cloud technologies
Key Features Learn advanced concepts in Azure ML and the Cortana Intelligence Suite architecture Explore ML Server using SQL Server and HDInsight capabilities Implement various tools in Azure to build and deploy machine learning models
Book Description Implementing Machine learning (ML) and Artificial Intelligence (AI) in the cloud had not been possible earlier due to the lack of processing power and storage. However, Azure has created ML and AI services that are easy to implement in the cloud. Hands-On Machine Learning with Azure teaches you how to perform advanced ML projects in the cloud in a cost-effective way. The book begins by covering the benefits of ML and AI in the cloud. You will then explore Microsoft's Team Data Science Process to establish a repeatable process for successful AI development and implementation. You will also gain an understanding of AI technologies available in Azure and the Cognitive Services APIs to integrate them into bot applications. This book lets you explore prebuilt templates with Azure Machine Learning Studio and build a model using canned algorithms that can be deployed as web services. The book then takes you through a preconfigured series of virtual machines in Azure targeted at AI development scenarios. You will get to grips with the ML Server and its capabilities in SQL and HDInsight. In the concluding chapters, you'll integrate patterns with other non-AI services in Azure. By the end of this book, you will be fully equipped to implement smart cognitive actions in your models. What you will learn Discover the benefits of leveraging the cloud for ML and AI Use Cognitive Services APIs to build intelligent bots Build a model using canned algorithms from Microsoft and deploy it as a web service Deploy virtual machines in AI development scenarios Apply R, Python, SQL Server, and Spark in Azure Build and deploy deep learning solutions with CNTK, MMLSpark,

and TensorFlow Implement model retraining in IoT, Streaming, and Blockchain solutions Explore best practices for integrating ML and AI functions with ADLA and logic apps Who this book is for If you are a data scientist or developer familiar with Azure ML and cognitive services and want to create smart models and make sense of data in the cloud, this book is for you. You'll also find this book useful if you want to bring powerful machine learning services into your cloud applications. Some experience with data manipulation and processing, using languages like SQL, Python, and R, will aid in understanding the concepts covered in this book

Data in all domains is getting bigger. How can you work with it efficiently? Recently updated for Spark 1.3, this book introduces Apache Spark, the open source cluster computing system that makes data analytics fast to write and fast to run. With Spark, you can tackle big datasets quickly through simple APIs in Python, Java, and Scala. This edition includes new information on Spark SQL, Spark Streaming, setup, and Maven coordinates. Written by the developers of Spark, this book will have data scientists and engineers up and running in no time. You'll learn how to express parallel jobs with just a few lines of code, and cover applications from simple batch jobs to stream processing and machine learning. Quickly dive into Spark capabilities such as distributed datasets, in-memory caching, and the interactive shell Leverage Spark's powerful built-in libraries, including Spark SQL, Spark Streaming, and MLlib Use one programming paradigm instead of mixing and matching tools like Hive, Hadoop, Mahout, and Storm Learn how to deploy interactive, batch, and streaming applications Connect to data sources including HDFS, Hive, JSON, and S3 Master advanced topics like data partitioning and shared variables Design, implement, and deliver successful streaming applications, machine learning pipelines and graph applications using Spark SQL API About This Book Learn about the design and implementation of streaming applications, machine learning pipelines, deep learning, and large-scale graph processing applications using Spark SQL APIs and Scala. Learn data exploration, data munging, and how to process structured and semi-structured data using real-world datasets and gain hands-on exposure to the issues and challenges of working with noisy and "dirty" real-world data. Understand design considerations for scalability and performance in web-scale Spark application architectures. Who This Book Is For If you are a developer, engineer, or an architect and want to learn how to use Apache Spark in a web-scale project, then this is the book for you. It is assumed that you have prior knowledge of SQL querying. A basic programming knowledge with Scala, Java, R, or Python is all you need to get started with this book. What You Will Learn Familiarize yourself with Spark SQL programming, including working with DataFrame/Dataset API and SQL Perform a series of hands-on exercises with different types of data sources, including CSV, JSON, Avro, MySQL, and MongoDB Perform data quality checks, data visualization, and basic statistical analysis tasks Perform data munging tasks on publically available datasets Learn how to use Spark SQL and Apache Kafka to build streaming applications Learn key performance-tuning tips and tricks in Spark SQL applications Learn key architectural components and patterns in large-scale Spark SQL applications In Detail In the past year, Apache Spark has been increasingly adopted for the development of distributed applications. Spark SQL APIs provide an optimized interface that helps developers build such applications quickly and easily. However, designing web-scale production applications using Spark SQL APIs can be a complex task. Hence, understanding the design and implementation best practices before you start your project will help you avoid these problems. This book gives an insight into the engineering practices used to design and build real-world, Spark-based applications. The book's hands-on examples will give you the required confidence to work on any future projects you encounter in Spark SQL. It starts by familiarizing you with data exploration and data munging tasks using Spark SQL and Scala. Extensive code examples will help you understand the methods used to implement typical use-cases for various types of applications. You will get a walkthrough of the key concepts and terms that are common to streaming, machine learning, and graph applications. You will also learn key

performance-tuning details including Cost Based Optimization (Spark 2.2) in Spark SQL applications. Finally, you will move on to learning how such systems are architected and deployed for a successful delivery of your project. **Style and approach** This book is a hands-on guide to designing, building, and deploying Spark SQL-centric production applications at scale.

Analyze vast amounts of data in record time using Apache Spark with Databricks in the Cloud. Learn the fundamentals, and more, of running analytics on large clusters in Azure and AWS, using Apache Spark with Databricks on top. Discover how to squeeze the most value out of your data at a mere fraction of what classical analytics solutions cost, while at the same time getting the results you need, incrementally faster. This book explains how the confluence of these pivotal technologies gives you enormous power, and cheaply, when it comes to huge datasets. You will begin by learning how cloud infrastructure makes it possible to scale your code to large amounts of processing units, without having to pay for the machinery in advance. From there you will learn how Apache Spark, an open source framework, can enable all those CPUs for data analytics use. Finally, you will see how services such as Databricks provide the power of Apache Spark, without you having to know anything about configuring hardware or software. By removing the need for expensive experts and hardware, your resources can instead be allocated to actually finding business value in the data. **This book guides you through some advanced topics such as analytics in the cloud, data lakes, data ingestion, architecture, machine learning, and tools, including Apache Spark, Apache Hadoop, Apache Hive, Python, and SQL. Valuable exercises help reinforce what you have learned.** **What You Will Learn** Discover the value of big data analytics that leverage the power of the cloud Get started with Databricks using SQL and Python in either Microsoft Azure or AWS Understand the underlying technology, and how the cloud and Apache Spark fit into the bigger picture See how these tools are used in the real world Run basic analytics, including machine learning, on billions of rows at a fraction of a cost or free **Who This Book Is For** Data engineers, data scientists, and cloud architects who want or need to run advanced analytics in the cloud. It is assumed that the reader has data experience, but perhaps minimal exposure to Apache Spark and Azure Databricks. The book is also recommended for people who want to get started in the analytics field, as it provides a strong foundation.

Work with Apache Spark using Scala to deploy and set up single-node, multi-node, and high-availability clusters. This book discusses various components of Spark such as Spark Core, DataFrames, Datasets and SQL, Spark Streaming, Spark MLlib, and R on Spark with the help of practical code snippets for each topic. **Practical Apache Spark** also covers the integration of Apache Spark with Kafka with examples. You'll follow a learn-to-do-by-yourself approach to learning – learn the concepts, practice the code snippets in Scala, and complete the assignments given to get an overall exposure. On completion, you'll have knowledge of the functional programming aspects of Scala, and hands-on expertise in various Spark components. You'll also become familiar with machine learning algorithms with real-time usage. **What You Will Learn** Discover the functional programming features of Scala Understand the complete architecture of Spark and its components Integrate Apache Spark with Hive and Kafka Use Spark SQL, DataFrames, and Datasets to process data using traditional SQL queries Work with different machine learning concepts and libraries using Spark's MLlib packages **Who This Book Is For** Developers and professionals who deal with batch and stream data processing.

This book contains the questions answers and some FAQ about the Databricks Spark Certification for version 2.x, which is the latest release from Apache Spark. In this book we will be having in total 75 practice questions. Almost all required question would have in detail explanation to the questions and answers, wherever required. Don't consider this book as a guide, it is more of question and answer practice book. This book also give some references as well like how to prepare further to ensure that you clear the certification exam. This book will particularly

focus on the Python version of the certification preparation material. Please note these are practice questions and not dumps, hence just memorizing the question and answers will not help in the real exam. You need to understand the concepts in detail as well as you should be able to solve the programming questions at the end in real world's work you should be able to write code using PySpark whether you are Data Engineer, Data Analytics Engineer, Data Scientists or Programmer. Hence, take the opportunity to learn each question and also go through the explanation of the questions.

This book constitutes the thoroughly refereed post-conference proceedings of the International Conference for Smart Health, ICSH 2016, held in Haikou, Hainan, China, in December 2016. The 23 full papers presented were carefully reviewed and selected from 52 submissions. They are organized around the following topics: big data and smart health; health data analysis and management; healthcare intelligent systems and clinical practice; medical monitoring and information extraction; clinical and medical data mining.

Apache® Spark is one of the fastest growing technology in BigData computing world. It supports multiple programming languages like Java, Scala, Python and R. Hence, many existing and new framework started to integrate Spark platform as well in their platform e.g. Hadoop, Cassandra, EMR etc. While creating Spark certification material HadoopExam technical team found that there is no proper material and book is available for the Spark (version 2.x) which covers the concepts as well as use of various features and found difficulty in creating the material. Therefore, they decided to create full length book for Spark (HDPSCD Spark Scala Certification) and outcome of that is this book. In this book technical team try to cover both fundamental concepts of Spark 2.x topics which are part of the certification syllabus as well as add as many exercises as possible and in current version we have around 10 hands on exercises added which you can execute on the Hortonworks sandbox, as this book is focused on the Scala version of the certification, hence all the exercises and their solution provided in the Scala. We have divided the entire book in the 7 chapters, as you move ahead chapter by chapter you would be comfortable with the HDPSCD Spark Scala certification. All the exercises given in this book are written using Scala. However, concepts remain same even if you are using different programming language.

This Book is published by www.HadoopExam.com (HadoopExam Learning Resources). Where you can find material and training's for preparing for Big-data, Cloud Computing, Analytics, Data Science and popular Programming Language. This Book will contain 130+ frequent interview questions for Spark 2.0 framework, which also covers the YARN framework, Spark streaming, Core Spark and SparkSQL, PySpark, these questions will not only help you in clearing interview process, but also you can understand various underline concepts, which Spark engine uses internally. Also, it is recommended that you go through the Spark Hands On Training provided by HadoopExam. In training we have created concepts as well as practicals by creating simple and complex problems with the use of Spark framework API. While publishing this book there are 32 modules available, which are in-line with Spark technology to be used on Hadoop Framework. As you know, Spark is one the most popular computing framework used and very well integrate with the Hadoop framework. You can see previously professionals were using MapReduce framework as a computing engine, but since Spark developed it is almost replaced by Spark engine, because Spark can give you rich API as well as it do most of the time data processing by having data in memory. Having data in-memory can save lot of disk I/O and drastically

improve the performance of submitted application. If you see now a days IOT and Machine learning are catching up and most of the professional started using higher level API created using Spark framework like MLlib, Graphx etc. Spark technology is now a days an exclusive skill, which most of developer want to learn. So to fulfill this need HadoopExam.com has many learning resources for learning Spark and doing certifications. Currently we have following products available to make you master in Apache Framework, visit HadoopExam.com for more detail. 1. Apache Spark Professional Training with Hands On Lab Sessions 2. O'Reilly Databricks Apache Spark Developer Certification Simulator 3. Hortonworks Spark Developer Certification 4. Cloudera CCA175 Hadoop and Spark Developer Certification 5. MapR Spark Certification preparation material This book has collection of questions, which are usually asked by the interviewer while filtering the candidates who had really worked on Spark framework which is well integrated with the Hadoop Framework.

This book covers the fundamentals of machine learning with Python in a concise and dynamic manner. It covers data mining and large-scale machine learning using Apache Spark. About This Book Take your first steps in the world of data science by understanding the tools and techniques of data analysis Train efficient Machine Learning models in Python using the supervised and unsupervised learning methods Learn how to use Apache Spark for processing Big Data efficiently Who This Book Is For If you are a budding data scientist or a data analyst who wants to analyze and gain actionable insights from data using Python, this book is for you. Programmers with some experience in Python who want to enter the lucrative world of Data Science will also find this book to be very useful, but you don't need to be an expert Python coder or mathematician to get the most from this book. What You Will Learn Learn how to clean your data and ready it for analysis Implement the popular clustering and regression methods in Python Train efficient machine learning models using decision trees and random forests Visualize the results of your analysis using Python's Matplotlib library Use Apache Spark's MLlib package to perform machine learning on large datasets In Detail Join Frank Kane, who worked on Amazon and IMDb's machine learning algorithms, as he guides you on your first steps into the world of data science. Hands-On Data Science and Python Machine Learning gives you the tools that you need to understand and explore the core topics in the field, and the confidence and practice to build and analyze your own machine learning models. With the help of interesting and easy-to-follow practical examples, Frank Kane explains potentially complex topics such as Bayesian methods and K-means clustering in a way that anybody can understand them. Based on Frank's successful data science course, Hands-On Data Science and Python Machine Learning empowers you to conduct data analysis and perform efficient machine learning using Python. Let Frank help you unearth the value in your data using the various data mining and data analysis techniques available in Python, and to develop efficient predictive models to predict future results. You will also learn how to perform large-scale machine learning on Big Data using Apache Spark. The book covers preparing your data for analysis, training machine learning models, and visualizing the final data analysis. Style and approach This comprehensive book is a perfect blend of theory and hands-on code examples in Python which can be used for your reference at any time.

Get command of your organizational Big Data using the power of data science and analytics Key Features A perfect companion to

boost your Big Data storing, processing, analyzing skills to help you take informed business decisions Work with the best tools such as Apache Hadoop, R, Python, and Spark for NoSQL platforms to perform massive online analyses Get expert tips on statistical inference, machine learning, mathematical modeling, and data visualization for Big Data Book Description Big Data analytics relates to the strategies used by organizations to collect, organize and analyze large amounts of data to uncover valuable business insights that otherwise cannot be analyzed through traditional systems. Crafting an enterprise-scale cost-efficient Big Data and machine learning solution to uncover insights and value from your organization's data is a challenge. Today, with hundreds of new Big Data systems, machine learning packages and BI Tools, selecting the right combination of technologies is an even greater challenge. This book will help you do that. With the help of this guide, you will be able to bridge the gap between the theoretical world of technology with the practical ground reality of building corporate Big Data and data science platforms. You will get hands-on exposure to Hadoop and Spark, build machine learning dashboards using R and R Shiny, create web-based apps using NoSQL databases such as MongoDB and even learn how to write R code for neural networks. By the end of the book, you will have a very clear and concrete understanding of what Big Data analytics means, how it drives revenues for organizations, and how you can develop your own Big Data analytics solution using different tools and methods articulated in this book. What you will learn - Get a 360-degree view into the world of Big Data, data science and machine learning - Broad range of technical and business Big Data analytics topics that caters to the interests of the technical experts as well as corporate IT executives - Get hands-on experience with industry-standard Big Data and machine learning tools such as Hadoop, Spark, MongoDB, KDB+ and R - Create production-grade machine learning BI Dashboards using R and R Shiny with step-by-step instructions - Learn how to combine open-source Big Data, machine learning and BI Tools to create low-cost business analytics applications - Understand corporate strategies for successful Big Data and data science projects - Go beyond general-purpose analytics to develop cutting-edge Big Data applications using emerging technologies Who this book is for The book is intended for existing and aspiring Big Data professionals who wish to become the go-to person in their organization when it comes to Big Data architecture, analytics, and governance. While no prior knowledge of Big Data or related technologies is assumed, it will be helpful to have some programming experience.

Summary The Spark distributed data processing platform provides an easy-to-implement tool for ingesting, streaming, and processing data from any source. In Spark in Action, Second Edition, you'll learn to take advantage of Spark's core features and incredible processing speed, with applications including real-time computation, delayed evaluation, and machine learning. Spark skills are a hot commodity in enterprises worldwide, and with Spark's powerful and flexible Java APIs, you can reap all the benefits without first learning Scala or Hadoop. Purchase of the print book includes a free eBook in PDF, Kindle, and ePub formats from Manning Publications. About the technology Analyzing enterprise data starts by reading, filtering, and merging files and streams from many sources. The Spark data processing engine handles this varied volume like a champ, delivering speeds 100 times faster than Hadoop systems. Thanks to SQL support, an intuitive interface, and a straightforward multilanguage API, you

can use Spark without learning a complex new ecosystem. About the book Spark in Action, Second Edition, teaches you to create end-to-end analytics applications. In this entirely new book, you'll learn from interesting Java-based examples, including a complete data pipeline for processing NASA satellite data. And you'll discover Java, Python, and Scala code samples hosted on GitHub that you can explore and adapt, plus appendixes that give you a cheat sheet for installing tools and understanding Spark-specific terms. What's inside Writing Spark applications in Java Spark application architecture Ingestion through files, databases, streaming, and Elasticsearch Querying distributed datasets with Spark SQL About the reader This book does not assume previous experience with Spark, Scala, or Hadoop. About the author Jean-Georges Perrin is an experienced data and software architect. He is France's first IBM Champion and has been honored for 12 consecutive years. Table of Contents PART 1 - THE THEORY CRIPPLED BY AWESOME EXAMPLES 1 So, what is Spark, anyway? 2 Architecture and flow 3 The majestic role of the dataframe 4 Fundamentally lazy 5 Building a simple app for deployment 6 Deploying your simple app PART 2 - INGESTION 7 Ingestion from files 8 Ingestion from databases 9 Advanced ingestion: finding data sources and building your own 10 Ingestion through structured streaming PART 3 - TRANSFORMING YOUR DATA 11 Working with SQL 12 Transforming your data 13 Transforming entire documents 14 Extending transformations with user-defined functions 15 Aggregating your data PART 4 - GOING FURTHER 16 Cache and checkpoint: Enhancing Spark's performances 17 Exporting data and building full data pipelines 18 Exploring deployment

[Copyright: 8b3c8562c1a72876ae70d620b178ce00](#)