

---

# Big Data Analytics With Spark Home Springer

---

Thank you extremely much for downloading **Big Data Analytics With Spark Home Springer**. Most likely you have knowledge that, people have seen numerous periods for their favorite books later this Big Data Analytics With Spark Home Springer, but stop taking place in harmful downloads.

Rather than enjoying a good ebook behind a mug of coffee in the afternoon, on the other hand they juggled subsequently some harmful virus inside their computer. **Big Data Analytics With Spark Home Springer** is easy to get to in our digital library an online entry to it is set as public fittingly you can download it instantly. Our digital library saves in fused countries, allowing you to acquire the most less latency times to download any of our books behind this one. Merely said, the Big Data Analytics With Spark Home Springer is universally compatible similar to any devices to read.

*Big Data Analytics With Spark Home  
Springer*

2021-02-07

---

## WIGGINS RAMOS

---

*Data Analytics with Hadoop* Packt Publishing Ltd

This IBM® Redbooks® publication provides topics to help the technical community take advantage of the resilience, scalability, and performance of the IBM Power Systems™ platform to implement or integrate an IBM Data Engine for Hadoop and Spark solution for analytics solutions to access, manage, and analyze data sets to improve business outcomes. This book documents topics to demonstrate and take advantage of the analytics strengths of the IBM POWER8® platform, the IBM analytics software portfolio, and selected third-party tools to help solve customer's data analytic workload requirements. This book describes how to plan, prepare, install, integrate, manage, and

show how to use the IBM Data Engine for Hadoop and Spark solution to run analytic workloads on IBM POWER8. In addition, this publication delivers documentation to complement available IBM analytics solutions to help your data analytic needs. This publication strengthens the position of IBM analytics and big data solutions with a well-defined and documented deployment model within an IBM POWER8 virtualized environment so that customers have a planned foundation for security, scaling, capacity, resilience, and optimization for analytics workloads. This book is targeted at technical professionals (analytics consultants, technical support staff, IT Architects, and IT Specialists) that are responsible for delivering analytics solutions and support on IBM Power Systems.

High Performance Spark Packt Publishing Ltd

Explore big data concepts, platforms, analytics, and their applications using the power of Hadoop 3 Key Features Learn

Hadoop 3 to build effective big data analytics solutions on-premise and on cloud Integrate Hadoop with other big data tools such as R, Python, Apache Spark, and Apache Flink Exploit big data using Hadoop 3 with real-world examples Book Description Apache Hadoop is the most popular platform for big data processing, and can be combined with a host of other big data tools to build powerful analytics solutions. Big Data Analytics with Hadoop 3 shows you how to do just that, by providing insights into the software as well as its benefits with the help of practical examples. Once you have taken a tour of Hadoop 3's latest features, you will get an overview of HDFS, MapReduce, and YARN, and how they enable faster, more efficient big data processing. You will then move on to learning how to integrate Hadoop with the open source tools, such as Python and R, to analyze and visualize data and perform statistical computing on big data. As you get acquainted with all this, you will explore how to use Hadoop 3 with Apache Spark and Apache Flink for real-time data analytics and stream processing. In addition to this, you will understand how to use Hadoop to build analytics solutions on the cloud and an end-to-end pipeline to perform big data analysis using practical use cases. By the end of this book, you will be well-versed with the analytical capabilities of the Hadoop ecosystem. You will be able to build powerful solutions to perform big data analytics and get insight effortlessly. What you will learn Explore the new features of Hadoop 3 along with HDFS, YARN, and MapReduce Get well-versed with the analytical capabilities of Hadoop ecosystem using practical examples Integrate Hadoop with R and Python for more efficient big data processing Learn to use Hadoop with Apache Spark and Apache

Flink for real-time data analytics Set up a Hadoop cluster on AWS cloud Perform big data analytics on AWS using Elastic Map Reduce Who this book is for Big Data Analytics with Hadoop 3 is for you if you are looking to build high-performance analytics solutions for your enterprise or business using Hadoop 3's powerful features, or you're new to big data analytics. A basic understanding of the Java programming language is required. Big Data Analytics with R Simon and Schuster Master alternative Big Data technologies that can do what Hadoop can't: real-time analytics and iterative machine learning. When most technical professionals think of Big Data analytics today, they think of Hadoop. But there are many cutting-edge applications that Hadoop isn't well suited for, especially real-time analytics and contexts requiring the use of iterative machine learning algorithms. Fortunately, several powerful new technologies have been developed specifically for use cases such as these. Big Data Analytics Beyond Hadoop is the first guide specifically designed to help you take the next steps beyond Hadoop. Dr. Vijay Srinivas Agneeswaran introduces the breakthrough Berkeley Data Analysis Stack (BDAS) in detail, including its motivation, design, architecture, Mesos cluster management, performance, and more. He presents realistic use cases and up-to-date example code for: Spark, the next generation in-memory computing technology from UC Berkeley Storm, the parallel real-time Big Data analytics technology from Twitter GraphLab, the next-generation graph processing paradigm from CMU and the University of Washington (with comparisons to alternatives such as Pregel and Piccolo) Halo also offers architectural and design guidance and code sketches for

scaling machine learning algorithms to Big Data, and then realizing them in real-time. He concludes by previewing emerging trends, including real-time video analytics, SDNs, and even Big Data governance, security, and privacy issues. He identifies intriguing startups and new research possibilities, including BDAS extensions and cutting-edge model-driven analytics. Big Data Analytics Beyond Hadoop is an indispensable resource for everyone who wants to reach the cutting edge of Big Data analytics, and stay there: practitioners, architects, programmers, data scientists, researchers, startup entrepreneurs, and advanced students.

[Hands-On Big Data Analytics with PySpark](#) Springer Nature

Get to grips with processing large volumes of data and presenting it as engaging, interactive insights using Spark and Python. Key Features Get a hands-on, fast-paced introduction to the Python data science stack Explore ways to create useful metrics and statistics from large datasets Create detailed analysis reports with real-world data Book Description Processing big data in real time is challenging due to scalability, information inconsistency, and fault tolerance. Big Data Analysis with Python teaches you how to use tools that can control this data avalanche for you. With this book, you'll learn practical techniques to aggregate data into useful dimensions for posterior analysis, extract statistical measurements, and transform datasets into features for other systems. The book begins with an introduction to data manipulation in Python using pandas. You'll then get familiar with statistical analysis and plotting techniques. With multiple hands-on activities in store, you'll be able to analyze data that is distributed on several computers by using Dask. As you progress,

you'll study how to aggregate data for plots when the entire data cannot be accommodated in memory. You'll also explore Hadoop (HDFS and YARN), which will help you tackle larger datasets. The book also covers Spark and explains how it interacts with other tools. By the end of this book, you'll be able to bootstrap your own Python environment, process large files, and manipulate data to generate statistics, metrics, and graphs. What you will learn Use Python to read and transform data into different formats Generate basic statistics and metrics using data on disk Work with computing tasks distributed over a cluster Convert data from various sources into storage or querying formats Prepare data for statistical analysis, visualization, and machine learning Present data in the form of effective visuals Who this book is for Big Data Analysis with Python is designed for Python developers, data analysts, and data scientists who want to get hands-on with methods to control data and transform it into impactful insights. Basic knowledge of statistical measurements and relational databases will help you to understand various concepts explained in this book.

**Big Data Processing with Apache Spark** Lulu.com

Solve Data Analytics Problems with Spark, PySpark, and Related Open Source Tools Spark is at the heart of today's Big Data revolution, helping data professionals supercharge efficiency and performance in a wide range of data processing and analytics tasks. In this guide, Big Data expert Jeffrey Aven covers all you need to know to leverage Spark, together with its extensions, subprojects, and wider ecosystem. Aven combines a language-agnostic introduction to foundational Spark concepts with extensive programming examples utilizing the popular and

intuitive PySpark development environment. This guide's focus on Python makes it widely accessible to large audiences of data professionals, analysts, and developers—even those with little Hadoop or Spark experience. Aven's broad coverage ranges from basic to advanced Spark programming, and Spark SQL to machine learning. You'll learn how to efficiently manage all forms of data with Spark: streaming, structured, semi-structured, and unstructured. Throughout, concise topic overviews quickly get you up to speed, and extensive hands-on exercises prepare you to solve real problems. Coverage includes:

- Understand Spark's evolving role in the Big Data and Hadoop ecosystems
- Create Spark clusters using various deployment modes
- Control and optimize the operation of Spark clusters and applications
- Master Spark Core RDD API programming techniques
- Extend, accelerate, and optimize Spark routines with advanced API platform constructs, including shared variables, RDD storage, and partitioning
- Efficiently integrate Spark with both SQL and nonrelational data stores
- Perform stream processing and messaging with Spark Streaming and Apache Kafka
- Implement predictive modeling with SparkR and Spark MLlib

*Scala Programming for Big Data Analytics* Apress

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.

*Essential PySpark for Scalable Data Analytics* Apress

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

*Learning PySpark* Packt Publishing Ltd

Design, process, and analyze large sets of complex data in real time About This Book Get acquainted with transformations and database-level interactions, and ensure the reliability of messages processed using Storm Implement strategies to solve the challenges of real-time data processing Load datasets, build queries, and make recommendations using Spark SQL Who This Book Is For If you are a Big Data architect, developer, or a programmer who wants to develop applications/frameworks to implement real-time analytics using open source technologies, then this book is for you. What You Will Learn Explore big data technologies and frameworks Work through practical challenges and use cases of real-time analytics versus batch analytics Develop real-world use cases for processing and analyzing data in real-time using the programming paradigm of Apache Storm Handle and process real-time transactional data Optimize and tune Apache Storm for varied workloads and production deployments Process and stream data with Amazon Kinesis and Elastic MapReduce Perform interactive and exploratory data analytics using Spark SQL Develop common enterprise architectures/applications for real-time and batch analytics In Detail Enterprise has been striving hard to deal with the challenges of data arriving in real time or near real time. Although there are technologies such as Storm and Spark (and

many more) that solve the challenges of real-time data, using the appropriate technology/framework for the right business use case is the key to success. This book provides you with the skills required to quickly design, implement and deploy your real-time analytics using real-world examples of big data use cases. From the beginning of the book, we will cover the basics of varied real-time data processing frameworks and technologies. We will discuss and explain the differences between batch and real-time processing in detail, and will also explore the techniques and programming concepts using Apache Storm. Moving on, we'll familiarize you with "Amazon Kinesis" for real-time data processing on cloud. We will further develop your understanding of real-time analytics through a comprehensive review of Apache Spark along with the high-level architecture and the building blocks of a Spark program. You will learn how to transform your data, get an output from transformations, and persist your results using Spark RDDs, using an interface called Spark SQL to work with Spark. At the end of this book, we will introduce Spark Streaming, the streaming library of Spark, and will walk you through the emerging Lambda Architecture (LA), which provides a hybrid platform for big data processing by combining real-time and precomputed batch data to provide a near real-time view of incoming data. Style and approach This step-by-step is an easy-to-follow, detailed tutorial, filled with practical examples of basic and advanced features. Each topic is explained sequentially and supported by real-world examples and executable code snippets.

**Big Data Analytics with Java** FT Press

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.

Knowledge Graphs and Big Data Processing John Wiley & Sons  
 Over insightful 90 recipes to get lightning-fast analytics with Apache Spark About This Book Use Apache Spark for data processing with these hands-on recipes Implement end-to-end, large-scale data analysis better than ever before Work with powerful libraries such as MLLib, SciPy, NumPy, and Pandas to gain insights from your data Who This Book Is For This book is for novice and intermediate level data science professionals and data analysts who want to solve data science problems with a distributed computing framework. Basic experience with data science implementation tasks is expected. Data science professionals looking to skill up and gain an edge in the field will find this book helpful. What You Will Learn Explore the topics of data mining, text mining, Natural Language Processing, information retrieval, and machine learning. Solve real-world analytical problems with large data sets. Address data science challenges with analytical tools on a distributed system like Spark (apt for iterative algorithms), which offers in-memory processing and more flexibility for data analysis at scale. Get hands-on experience with algorithms like Classification, regression, and recommendation on real datasets using Spark MLLib package. Learn about numerical and scientific computing using NumPy and SciPy on Spark. Use Predictive Model Markup Language (PMML) in



Spark for statistical data mining models. In Detail Spark has emerged as the most promising big data analytics engine for data science professionals. The true power and value of Apache Spark lies in its ability to execute data science tasks with speed and accuracy. Spark's selling point is that it combines ETL, batch analytics, real-time stream analysis, machine learning, graph processing, and visualizations. It lets you tackle the complexities that come with raw unstructured data sets with ease. This guide will get you comfortable and confident performing data science tasks with Spark. You will learn about implementations including distributed deep learning, numerical computing, and scalable machine learning. You will be shown effective solutions to problematic concepts in data science using Spark's data science libraries such as MLLib, Pandas, NumPy, SciPy, and more. These simple and efficient recipes will show you how to implement algorithms and optimize your work. Style and approach This book contains a comprehensive range of recipes designed to help you learn the fundamentals and tackle the difficulties of data science. This book outlines practical steps to produce powerful insights into Big Data through a recipe-based approach.

**Big Data Processing Using Spark in Cloud** Packt Publishing Ltd

Use PySpark to easily crush messy data at-scale and discover proven techniques to create testable, immutable, and easily parallelizable Spark jobs Key Features Work with large amounts of agile data using distributed datasets and in-memory caching Source data from all popular data hosting platforms, such as HDFS, Hive, JSON, and S3 Employ the easy-to-use PySpark API to deploy big data Analytics for production Book Description

Apache Spark is an open source parallel-processing framework that has been around for quite some time now. One of the many uses of Apache Spark is for data analytics applications across clustered computers. In this book, you will not only learn how to use Spark and the Python API to create high-performance analytics with big data, but also discover techniques for testing, immunizing, and parallelizing Spark jobs. You will learn how to source data from all popular data hosting platforms, including HDFS, Hive, JSON, and S3, and deal with large datasets with PySpark to gain practical big data experience. This book will help you work on prototypes on local machines and subsequently go on to handle messy data in production and at scale. This book covers installing and setting up PySpark, RDD operations, big data cleaning and wrangling, and aggregating and summarizing data into useful reports. You will also learn how to implement some practical and proven techniques to improve certain aspects of programming and administration in Apache Spark. By the end of the book, you will be able to build big data analytical solutions using the various PySpark offerings and also optimize them effectively. What you will learn Get practical big data experience while working on messy datasets Analyze patterns with Spark SQL to improve your business intelligence Use PySpark's interactive shell to speed up development time Create highly concurrent Spark programs by leveraging immutability Discover ways to avoid the most expensive operation in the Spark API: the shuffle operation Re-design your jobs to use reduceByKey instead of groupByKey Create robust processing pipelines by testing Apache Spark jobs Who this book is for This book is for developers, data scientists, business analysts, or anyone who needs to reliably

analyze large amounts of large-scale, real-world data. Whether you're tasked with creating your company's business intelligence function or creating great data platforms for your machine learning models, or are looking to use code to magnify the impact of your business, this book is for you.

*Big Data Analytics Beyond Hadoop* Packt Publishing Ltd

Big Data Analytics with Spark is a step-by-step guide for learning Spark, which is an open-source fast and general-purpose cluster computing framework for large-scale data analysis. You will learn how to use Spark for different types of big data analytics projects, including batch, interactive, graph, and stream data analysis as well as machine learning. In addition, this book will help you become a much sought-after Spark expert. Spark is one of the hottest Big Data technologies. The amount of data generated today by devices, applications and users is exploding. Therefore, there is a critical need for tools that can analyze large-scale data and unlock value from it. Spark is a powerful technology that meets that need. You can, for example, use Spark to perform low latency computations through the use of efficient caching and iterative algorithms; leverage the features of its shell for easy and interactive Data analysis; employ its fast batch processing and low latency features to process your real time data streams and so on. As a result, adoption of Spark is rapidly growing and is replacing Hadoop MapReduce as the technology of choice for big data analytics. This book provides an introduction to Spark and related big-data technologies. It covers Spark core and its add-on libraries, including Spark SQL, Spark Streaming, GraphX, and MLlib. Big Data Analytics with Spark is therefore written for busy professionals who prefer learning a new technology from a

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

Spark in Action O'Reilly Media

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.

*Frank Kane's Taming Big Data with Apache Spark and Python*  
Packt Publishing Ltd  
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. Foreword by Rob Thomas. 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

### **Big Data Analytics** IBM Redbooks

Summary Introducing Data Science teaches you how to accomplish the fundamental tasks that occupy data scientists. Using the Python language and common Python libraries, you'll experience firsthand the challenges of dealing with data at scale and gain a solid foundation in data science. Purchase of the print book includes a free eBook in PDF, Kindle, and ePub formats from Manning Publications. About the Technology Many companies need developers with data science skills to work on projects ranging from social media marketing to machine learning. Discovering what you need to learn to begin a career as a data scientist can seem bewildering. This book is designed to help you get started. About the Book Introducing Data Science

Introducing Data Science explains vital data science concepts and teaches you how to accomplish the fundamental tasks that occupy data scientists. You'll explore data visualization, graph databases, the use of NoSQL, and the data science process. You'll use the Python language and common Python libraries as you experience firsthand the challenges of dealing with data at scale. Discover how Python allows you to gain insights from data sets so big that they need to be stored on multiple machines, or from data moving so quickly that no single machine can handle it. This book gives you hands-on experience with the most popular Python data science libraries, Scikit-learn and StatsModels. After reading this book, you'll have the solid foundation you need to start a career in data science. What's Inside Handling large data Introduction to machine learning Using Python to work with data Writing data science algorithms About the Reader This book assumes you're comfortable reading code in Python or a similar language, such as C, Ruby, or JavaScript. No prior experience with data science is required. About the Authors Davy Cielen, Arno D. B. Meysman, and Mohamed Ali are the founders and managing partners of Optimately and Maiton, where they focus on developing data science projects and solutions in various sectors. Table of Contents Data science in a big data world The data science process Machine learning Handling large data on a single computer First steps in big data Join the NoSQL movement The rise of graph databases Text mining and text analytics Data visualization to the end user *Learning Spark* "O'Reilly Media, Inc." Get started with distributed computing using PySpark, a single unified framework to solve end-to-end data analytics at scale Key

Features Discover how to convert huge amounts of raw data into meaningful and actionable insights Use Spark's unified analytics engine for end-to-end analytics, from data preparation to predictive analytics Perform data ingestion, cleansing, and integration for ML, data analytics, and data visualization Book Description Apache Spark is a unified data analytics engine designed to process huge volumes of data quickly and efficiently. PySpark is Apache Spark's Python language API, which offers Python developers an easy-to-use scalable data analytics framework. Essential PySpark for Scalable Data Analytics starts by exploring the distributed computing paradigm and provides a high-level overview of Apache Spark. You'll begin your analytics journey with the data engineering process, learning how to perform data ingestion, cleansing, and integration at scale. This book helps you build real-time analytics pipelines that help you gain insights faster. You'll then discover methods for building cloud-based data lakes, and explore Delta Lake, which brings reliability to data lakes. The book also covers Data Lakehouse, an emerging paradigm, which combines the structure and performance of a data warehouse with the scalability of cloud-based data lakes. Later, you'll perform scalable data science and machine learning tasks using PySpark, such as data preparation, feature engineering, and model training and productionization. Finally, you'll learn ways to scale out standard Python ML libraries along with a new pandas API on top of PySpark called Koalas. By the end of this PySpark book, you'll be able to harness the power of PySpark to solve business problems. What you will learn Understand the role of distributed computing in the world of big data Gain an appreciation for Apache Spark as the de facto

go-to for big data processing Scale out your data analytics process using Apache Spark Build data pipelines using data lakes, and perform data visualization with PySpark and Spark SQL Leverage the cloud to build truly scalable and real-time data analytics applications Explore the applications of data science and scalable machine learning with PySpark Integrate your clean and curated data with BI and SQL analysis tools Who this book is for This book is for practicing data engineers, data scientists, data analysts, and data enthusiasts who are already using data analytics to explore distributed and scalable data analytics. Basic to intermediate knowledge of the disciplines of data engineering, data science, and SQL analytics is expected. General proficiency in using any programming language, especially Python, and working knowledge of performing data analytics using frameworks such as pandas and SQL will help you to get the most out of this book.

#### **Advanced Analytics with Spark** "O'Reilly Media, Inc."

In the second edition of this practical book, four Cloudera data scientists present a set of self-contained patterns for performing large-scale data analysis with Spark. The authors bring Spark, statistical methods, and real-world data sets together to teach you how to approach analytics problems by example. Updated for Spark 2.1, this edition acts as an introduction to these techniques and other best practices in Spark programming. You'll start with an introduction to Spark and its ecosystem, and then dive into patterns that apply common techniques—including classification, clustering, collaborative filtering, and anomaly detection—to fields such as genomics, security, and finance. If you have an entry-level understanding of machine learning and statistics, and

you program in Java, Python, or Scala, you'll find the book's patterns useful for working on your own data applications. With this book, you will: Familiarize yourself with the Spark programming model Become comfortable within the Spark ecosystem Learn general approaches in data science Examine complete implementations that analyze large public data sets Discover which machine learning tools make sense for particular problems Acquire code that can be adapted to many uses  
*Beginning Apache Spark Using Azure Databricks* Packt Publishing Ltd

A handy reference guide for data analysts and data scientists to help to obtain value from big data analytics using Spark on Hadoop clusters About This Book This book is based on the latest 2.0 version of Apache Spark and 2.7 version of Hadoop integrated with most commonly used tools. Learn all Spark stack components including latest topics such as DataFrames, DataSets, GraphFrames, Structured Streaming, DataFrame based ML Pipelines and SparkR. Integrations with frameworks such as HDFS, YARN and tools such as Jupyter, Zeppelin, NiFi, Mahout, HBase Spark Connector, GraphFrames, H2O and Hivemall. Who This Book Is For Though this book is primarily aimed at data analysts and data scientists, it will also help architects, programmers, and practitioners. Knowledge of either Spark or Hadoop would be beneficial. It is assumed that you have basic programming background in Scala, Python, SQL, or R programming with basic Linux experience. Working experience within big data environments is not mandatory. What You Will Learn Find out and implement the tools and techniques of big data analytics using Spark on Hadoop clusters with wide variety

of tools used with Spark and Hadoop Understand all the Hadoop and Spark ecosystem components Get to know all the Spark components: Spark Core, Spark SQL, DataFrames, DataSets, Conventional and Structured Streaming, MLLib, ML Pipelines and Graphx See batch and real-time data analytics using Spark Core, Spark SQL, and Conventional and Structured Streaming Get to grips with data science and machine learning using MLLib, ML Pipelines, H2O, Hivemall, Graphx, SparkR and Hivemall. In Detail Big Data Analytics book aims at providing the fundamentals of Apache Spark and Hadoop. All Spark components – Spark Core, Spark SQL, DataFrames, Data sets, Conventional Streaming, Structured Streaming, MLLib, Graphx and Hadoop core components – HDFS, MapReduce and Yarn are explored in greater depth with implementation examples on Spark + Hadoop clusters. It is moving away from MapReduce to Spark. So, advantages of Spark over MapReduce are explained at great depth to reap benefits of in-memory speeds. DataFrames API, Data Sources API and new Data set API are explained for building Big Data analytical applications. Real-time data analytics using Spark Streaming with Apache Kafka and HBase is covered to help building streaming applications. New Structured streaming concept is explained with an IOT (Internet of Things) use case. Machine learning techniques are covered using MLLib, ML Pipelines and SparkR and Graph Analytics are covered with GraphX and GraphFrames components of Spark. Readers will also get an opportunity to get started with web based notebooks such as Jupyter, Apache Zeppelin and data flow tool Apache NiFi to analyze and visualize data. Style and approach This step-by-step pragmatic guide will make life easy no matter what your level of

experience. You will deep dive into Apache Spark on Hadoop clusters through ample exciting real-life examples. Practical tutorial explains data science in simple terms to help programmers and data analysts get started with Data Science [Spark for Data Science](#) Packt Publishing Ltd

Learn the basics of analytics on big data using Java, machine learning and other big data tools About This Book Acquire real-world set of tools for building enterprise level data science applications Surpasses the barrier of other languages in data science and learn create useful object-oriented codes Extensive use of Java compliant big data tools like apache spark, Hadoop, etc. Who This Book Is For This book is for Java developers who are looking to perform data analysis in production environment. Those who wish to implement data analysis in their Big data applications will find this book helpful. What You Will Learn Start from simple analytic tasks on big data Get into more complex tasks with predictive analytics on big data using machine learning Learn real time analytic tasks Understand the concepts with examples and case studies Prepare and refine data for analysis Create charts in order to understand the data See various real-world datasets In Detail This book covers case studies such as sentiment analysis on a tweet dataset, recommendations on a movielens dataset, customer segmentation on an ecommerce dataset, and graph analysis on actual flights dataset. This book is an end-to-end guide to implement analytics on big data with Java. Java is the de facto language for major big data environments, including Hadoop. This book will teach you how to perform analytics on big data with production-friendly Java. This book basically divided into two sections. The first part is an

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

*Advanced Analytics with Spark* "O'Reilly Media, Inc."

Analyze your data and delve deep into the world of machine learning with the latest Spark version, 2.0 About This Book Perform data analysis and build predictive models on huge datasets that leverage Apache Spark Learn to integrate data science algorithms and techniques with the fast and scalable computing features of Spark to address big data challenges Work through practical examples on real-world problems with sample code snippets Who This Book Is For This book is for anyone who wants to leverage Apache Spark for data science and machine learning. If you are a technologist who wants to expand your

knowledge to perform data science operations in Spark, or a data scientist who wants to understand how algorithms are implemented in Spark, or a newbie with minimal development experience who wants to learn about Big Data Analytics, this book is for you! What You Will Learn Consolidate, clean, and transform your data acquired from various data sources Perform statistical analysis of data to find hidden insights Explore graphical techniques to see what your data looks like Use machine learning techniques to build predictive models Build scalable data products and solutions Start programming using the RDD, DataFrame and Dataset APIs Become an expert by improving your data analytical skills In Detail This is the era of Big Data. The words 'Big Data' implies big innovation and enables a competitive advantage for businesses. Apache Spark was designed to perform Big Data analytics at scale, and so Spark is

equipped with the necessary algorithms and supports multiple programming languages. Whether you are a technologist, a data scientist, or a beginner to Big Data analytics, this book will provide you with all the skills necessary to perform statistical data analysis, data visualization, predictive modeling, and build scalable data products or solutions using Python, Scala, and R. With ample case studies and real-world examples, Spark for Data Science will help you ensure the successful execution of your data science projects. Style and approach This book takes a step-by-step approach to statistical analysis and machine learning, and is explained in a conversational and easy-to-follow style. Each topic is explained sequentially with a focus on the fundamentals as well as the advanced concepts of algorithms and techniques. Real-world examples with sample code snippets are also included.