

Introducción a Apache Spark

para empezar a programar
el *big data*

Mario Macías
Mauro Gómez
Rubén Tous
Jordi Torres



Director de la colección:

Diseño de la colección: Editorial UOC

Primera edición en lengua castellana: Noviembre 2015

© Mario Macías, Mauro Gómez, Rubèn Tous, Jordi Torres, del texto

© Diseño de la cubierta: Natàlia Serrano

© Editorial UOC (Oberta UOC Publishing, SL), de esta edición, AÑO
Rambla del Poblenou 156, 08018 Barcelona
<http://www.editorialuoc.com>

Realización editorial: Oberta UOC Publishing, SL

Maquetación: Maria García

Impresión: Service Point FMI, S.A

ISBN: 978-84-9116-037-3

Depósito legal: B 25326-2015

Ninguna parte de esta publicación, incluyendo el diseño general y de la cubierta, no puede ser copiada, reproducida, almacenada o transmitida de ninguna forma ni por ningún medio, ya sea eléctrico, químico, mecánico, óptico, de grabación, de fotocopia o por otros métodos, sin la autorización previa por escrito de los titulares del *copyright*.

Mario Macías

Doctor en Arquitectura de Computadores por la Universidad Politécnica de Cataluña, donde compagina su trabajo de profesor con el de investigador en el Barcelona Supercomputing Center-Centro Nacional de Supercomputación (BSC-CNS). Durante la última década ha trabajado en proyectos de investigación relacionados con *cloud computing* y eficiencia energética. Como actividad independiente, ha autopublicado dos libros de divulgación científica.

Mauro Gómez Parada

Graduado en Ingeniería Informática por la Universidad de Vigo. Actualmente está finalizando el máster de Ingeniería informática en la UPC (Universidad Politécnica de Barcelona) y trabaja como investigador sobre Spark con el grupo de Autonomic Systems del BSC-CNS. Como actividad independiente, ha colaborado en la creación de los Premios Galegos da Musica, nacidos en el año 2013 en Galicia.

Rubèn Tous

Doctor en Informática por la Universidad Pompeu Fabra. En la actualidad es profesor contratado doctor en el Departamento de Arquitectura de Computadores de la Universidad Politécnica de Cataluña e investigador colaborador en el Barcelona Supercomputing Center. Es experto en indexación, búsqueda y clasificación de información multimedia. Ha publicado más de cincuenta artículos de investigación en revistas y conferencias internacionales. Es coeditor de múltiples estándares de los grupos MPEG y JPEG de ISO, y ha sido codirector del Metadata Subgroup de JPEG.

Jordi Torres

Catedrático de la UPC y lidera un grupo de investigación en el BSC. Actualmente su investigación se centra en la convergencia de la computación de altas prestaciones con el *big data* y su aplicación a los retos que plantea la analítica del *big data* o la computación cognitiva. Dada su extensa carrera profesional en diferentes roles, también realiza actividades de consultoría y estrategia relacionadas con las tecnologías de próxima generación y su impacto, y actúa como experto para varias organizaciones y empresas o mentorizando a emprendedores. Una de sus pasiones es la divulgación científica, que lo ha llevado a escribir un par de libros, dar conferencias y colaborar con medios de comunicación como *La Vanguardia*. Mantiene un blog sobre tecnología en www.JordiTorres.eu.

El verdadero progreso es el que pone la tecnología al alcance de todos.

Henry Ford

Índice

Prólogo	15
Agradecimientos	17
Prefacio	21
Convenciones de formato	23
Capítulo I. ¿Qué es Apache Spark?	27
1. El <i>big data</i> ya está aquí	27
2. Extendiendo el software <i>stack</i> del <i>big data</i>	30
3. Mejorando la eficiencia del <i>big data</i>	31
4. Una pila de software <i>big data</i> unificada y evolucionada	33
5. Spark y Python	35
6. Spark y sus alternativas.....	37
Capítulo II. Descargar y empezar con Apache Spark	39
1. Descargar Apache Spark.....	39
2. Introducción al <i>shell</i> de Python.....	40
3. Conceptos esenciales de Spark	44
4. Aplicaciones autocontenidas	46
4.1. Ejecutando la aplicación mediante <i>spark-submit</i>	47
5. Configurando Spark.....	50

5.1. Spark properties: configuración a nivel de aplicación	50
5.2. Otras variables de configuración.....	52
Capítulo III. Conceptos básicos de Spark	53
1. Conjuntos de datos resilientes y distribuidos	53
1.1. Creación de RDD.....	54
1.2. Acciones sobre RDD.....	56
1.3. Transformaciones de RDD.....	61
1.4. Persistencia de RDD	69
2. Variables compartidas.....	72
2.1. Variables difundidas	72
2.2. Acumuladores.....	73
3. Caso práctico: palabras más frecuentes	75
Capítulo IV. Acceso a datos	79
1. Formatos de archivos	80
1.1. Ficheros de texto	80
1.2. Ficheros JSON.....	84
1.3. Ficheros CSV.....	86
1.4. Ficheros <i>SequenceFiles</i>	90
2. Bases de datos.....	94
2.1. JDBC	94
2.2. Hive.....	96
Capítulo V. SQL en Spark.....	99
1. <i>Data frames</i>	99
1.1. Creación de <i>data frames</i>	100
1.2. Operaciones básicas	103
2. Consultas SQL.....	112
2.1. Funciones definidas por el usuario.....	113

Capítulo VI. Procesando flujos de datos con Spark	115
1. Un ejemplo sencillo	116
2. Receptores	121
3. Transformaciones.....	122
3.1. Transformaciones básicas sobre DStreams.....	122
3.2. Transformaciones de tipo join	123
3.3. Operación <i>UpdateStateByKey</i>	124
3.4. Operación <i>transform</i>	126
3.5. Transformaciones en ventana (<i>windowed</i>)	127
4. Operaciones de salida con DStreams	130
5. DataFrames y operaciones SQL con Spark Streaming.....	131
 Capítulo VII. Aprendizaje automático con Spark.....	 135
1. El módulo MLlib.....	135
2. Uso de MLlib.....	137
2.1. Pasos a seguir	138
2.2. Ejemplo de agrupamiento.....	139
2.3. Crear los RDD	140
2.4. Convertir texto a valores numéricos.....	141
2.5. Entrenar el algoritmo.....	141
2.6. Evaluar el modelo.....	142
3. Funcionalidades de MLlib	144
3.1. Tipos de datos.....	144
3.2. Estadística	145
3.3. Clasificación y regresión.....	146
3.4. Extracción de características y transformación.....	148
3.5. Otros algoritmos y funcionalidades.....	149
4. Caso práctico: agrupamiento de imágenes etiquetadas..	149
5. Para saber más	154

Apéndices	155
Apéndice A. Análisis de grafos con Spark	157
Apéndice B. Breve introducción a Python	183

Prólogo

Nos encontramos en un momento muy emocionante a la hora de trabajar en computación paralela y *big data*. El gran volumen de datos que hoy en día se genera en todos los campos de la industria y la ciencia está revolucionando la forma como interactuamos con las aplicaciones, creamos productos y estudiamos el mundo a nuestro alrededor. Al mismo tiempo, las herramientas necesarias para trabajar con estos datos se han vuelto más fáciles de usar que nunca, puesto que los desarrolladores las han hecho accesibles a más y más usuarios, requiriéndoles menos y menos esfuerzo para adoptarlas. Espero que Apache Spark termine siendo una de estas herramientas para ti, que te aporte un nuevo medio para trabajar con datos de manera fácil, potente, e incluso a veces divertida de usar.

Por ello estoy encantado de ver este primer libro sobre Spark escrito en lengua española –hasta ahora todos los libros eran en lengua inglesa–, escrito por un fantástico equipo de autores. Mario, Mauro, Rubén y Jordi son destacados miembros de la comunidad con gran experiencia en Spark y la computación paralela en general por sus investigaciones y desarrollos en Barcelona. Sin duda, han elaborado un libro completo y fácil de seguir, con muchos ejemplos, y no solo cubren los fundamentos de Spark, sino también las bibliotecas más utilizadas del ecosistema que conforma Apache Spark.

Espero que este libro sea solo una introducción a tu viaje al procesado paralelo de datos en el mundo *big data*. Las ideas aquí tratadas representan algunos de los mejores métodos ideados

para trabajar con datos hoy en día. El procesado avanzado de datos sigue siendo una de las áreas de investigación más activas dentro de las ciencias de la computación, y estoy seguro de que están por llegar muchas nuevas ideas de otros campos de la informática para abrirse paso dentro de este campo. Espero que Spark siga aportándote alguna de estas ideas y que este libro te permita empezar a aprender sobre esta nueva y emocionante área.

Matei Zaharia, CTO en Databricks y vicepresidente de Apache Spark

Prefacio

En el marco de los cambios tecnológicos que estamos viviendo en pleno siglo XXI, el crecimiento exponencial de la información disponible representará una transformación completa de la actual forma de vivir, trabajar y pensar. Y todo ello nos lleva a la siguiente gran tendencia que dominará la industria de las TIC en los próximos años, conocida como *big data*.

En este escenario, nos encontramos en un momento en el que Apache Spark aparece con un vigor sin precedentes en el universo de los programadores del *big data*, porque ya se vislumbra que va a ser una de las plataformas de gran impacto en los próximos años.

Este libro proporciona al lector una oportunidad para empezar a programar y manejar datos a través del ecosistema Apache Spark. Spark es actualmente uno de los paquetes de código abierto más importantes en el espacio del *big data* y por el que importantes empresas como IBM, SAP, Oracle o Amazon han apostado, siendo asimismo grandes contribuidoras.

Este libro de carácter introductorio, que puede utilizarse como texto de autoestudio o de soporte a cursos que requieran una introducción a Apache Spark, presenta una descripción del ecosistema de paquetes que conforman Apache Spark, así como sus características básicas, e incluye ejemplos de código en Python para que el lector pueda tener una comprensión de primera mano de algunas de las posibilidades de Apache Spark, probando estos ejemplos en su propio PC si así lo desea.

Pero no debemos olvidar que el proyecto Apache Spark está evolucionando muy rápidamente. Cuando empezamos a escribir

este libro, en junio de 2015, la versión 1.4.0 aún estaba en fase beta. En el momento de imprimirse, octubre de 2015, la versión 1.5 de Spark acaba de ver la luz.

En general, Apache Spark mantiene cuidadosamente la compatibilidad binaria y de fuente para interfaces de programación estables como las que se presentan en este libro introductorio. Por tanto, esperamos que los ejemplos básicos usados en el libro sigan funcionando en futuras versiones Apache Spark 1.x y si fuera necesario los iremos actualizando en la página web del libro.

Los autores esperamos que esta obra sirva para ayudar a que esta nueva tecnología, que sin duda va a marcar el futuro inmediato del *big data*, llegue a todos aquellos estudiantes (universitarios o de bachillerato), científicos y técnicos que quieran dar un paso más en la programación del *big data*. Y al igual que nos dice Matei en su prólogo, nosotros también deseamos que este libro sea solo el principio de un largo camino en este fascinante mundo que ofrece el análisis del *big data*.

Esperamos que el lector disfrute aprendiendo con este libro de la misma manera que nosotros, los autores, lo hemos hecho escribiéndolo.