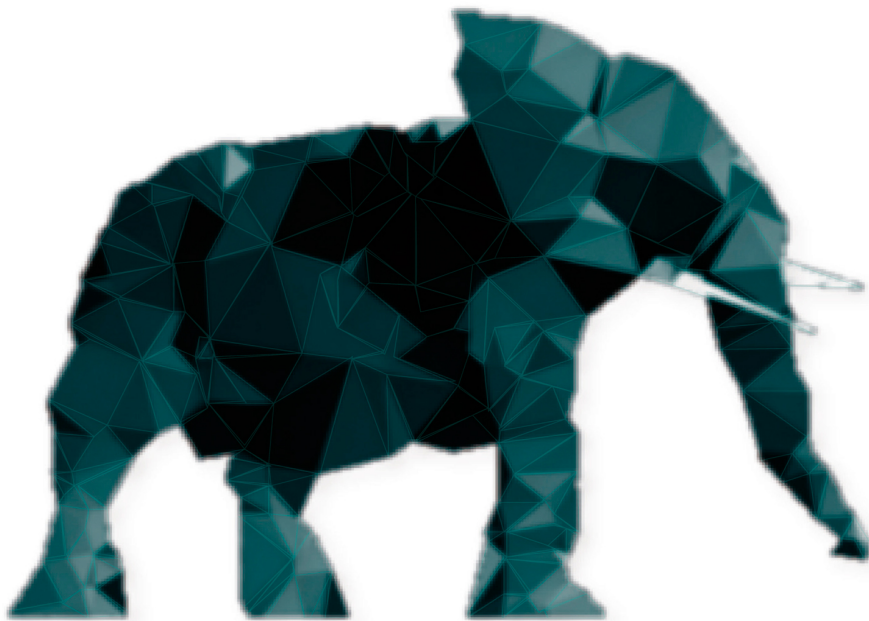
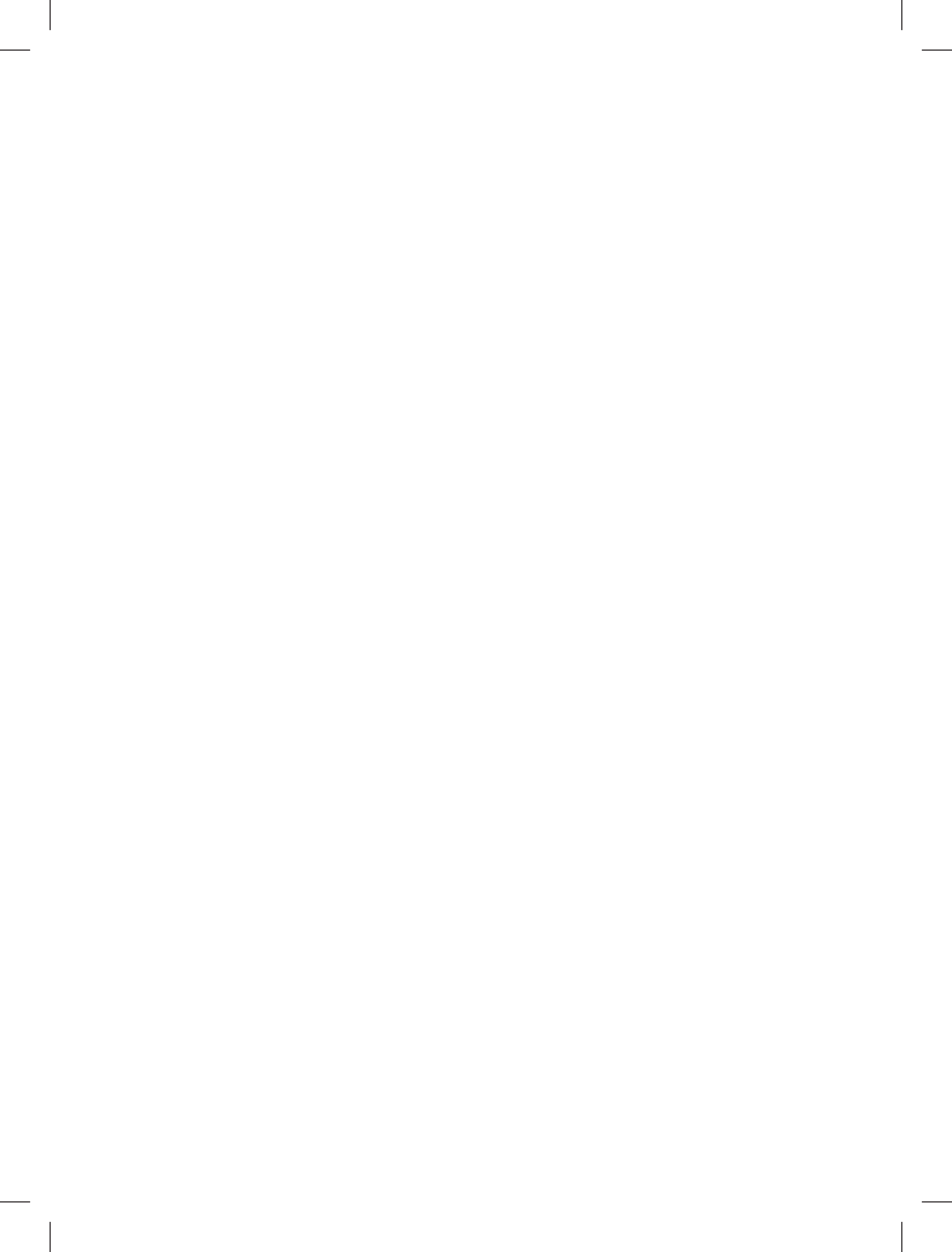


Big Data

ANÁLISIS DE GRANDES VOLÚMENES
DE DATOS EN ORGANIZACIONES

Luis Joyanes Aguilar





Big Data

Análisis de grandes volúmenes de datos en organizaciones

Luis Joyanes Aguilar



Big Data

Análisis de grandes volúmenes de datos en organizaciones

Luis Joyanes Agullar



Big Data. Análisis de grandes volúmenes de datos en organizaciones

Luis Joyanes Aguilar

Derechos reservados © Alfaomega Grupo Editor, S.A. de C.V., México
Primera edición: Alfaomega Grupo Editor, México, julio 2013

Primera edición: MARCOMBO, S.A. 2014

© 2014 MARCOMBO, S.A.
www.marcombo.com

Diseño de cubierta: Iris Biaggini

«Cualquier forma de reproducción, distribución, comunicación pública o transformación de esta obra sólo puede ser realizada con la autorización de sus titulares, salvo excepción prevista por la ley. Diríjase a CEDRO (Centro Español de Derechos Reprográficos, www.cedro.org) si necesita fotocopiar o escanear algún fragmento de esta obra».

ISBN: 978-84-267-2081-8

D.L.: B-20395-2013

Printed in Spain

A Inés y Olivia, que ya viven en la era de la nube y de Big Data disfrutando de sus enormes beneficios y como muestra del inmenso cariño que os tengo.

Luis Joyanes Aguilar



Mensaje del editor

Los conocimientos son esenciales en el desempeño profesional, sin ellos es imposible lograr las habilidades para competir laboralmente. La universidad o las instituciones de formación para el trabajo ofrecen la oportunidad de adquirir conocimientos que serán aprovechados más adelante en beneficio propio y de la sociedad; el avance de la ciencia y de la técnica hace necesario actualizar continuamente esos conocimientos. Cuando se toma la decisión de embarcarse en una vida profesional, se adquiere un compromiso de por vida: mantenerse al día en los conocimientos del área u oficio que se ha decidido desempeñar.

Alfaomega tiene por misión ofrecerles a estudiantes y profesionales conocimientos actualizados dentro de lineamientos pedagógicos que faciliten su utilización y permitan desarrollar las competencias requeridas por una profesión determinada. Alfaomega espera ser su compañera profesional en este viaje de por vida por el mundo del conocimiento.

Alfaomega hace uso de los medios impresos tradicionales en combinación con las tecnologías de la información y las comunicaciones (IT) para facilitar el aprendizaje.

Libros como éste tienen su complemento en una página Web, en donde el alumno y su profesor encontrarán materiales adicionales.

Esta obra contiene numerosos gráficos, cuadros y otros recursos para despertar el interés del estudiante, y facilitarle la comprensión y apropiación del conocimiento. Cada capítulo se desarrolla con argumentos presentados en forma sencilla y estructurada claramente hacia los objetivos y metas propuestas.

Los libros de Alfaomega están diseñados para ser utilizados dentro de los procesos de enseñanza-aprendizaje, y pueden ser usados como textos para diversos cursos o como apoyo para reforzar el desarrollo profesional.

Alfaomega espera contribuir así a la formación y el desarrollo de profesionales exitosos para beneficio de la sociedad.

Acerca del autor

Luis Joyanes Aguilar

Doctor Ingeniero en Informática y Doctor en Sociología, Catedrático de Lenguajes y Sistemas Informáticos de la Universidad Pontificia de Salamanca en el campus de Madrid y profesor invitado en diferentes universidades del mundo. Conferenciante habitual en congresos, seminarios, jornadas y talleres a nivel mundial. Ha escrito numerosos libros y artículos relativos a Tecnologías de la Información.

Patrono de la Fundación de I+D Software Libre de Granada, miembro del Instituto Universitario “Agustín Millares” de la Universidad Carlos III de Madrid y presidente de SISOFT.

Contenido

Parte I. La era de Big Data

CAPÍTULO 1

¿QUÉ ES BIG DATA?	1
Definición de Big Data.....	2
Tipos de datos.....	3
Datos estructurados.....	4
Datos semiestructurados	4
Datos no estructurados.....	5
Integración de los datos: oportunidades de negocio de los Big Data	5
Características de Big Data.....	7
Volumen.....	7
Velocidad	8
Variedad.....	8
Veracidad.....	10
Valor.....	10
El tamaño de los Big Data.....	10
¿Cómo se ha llegado a la explosión de Big Data?	11
El Big Data eclosiona en España (IDC) ...	12
Cómo crear ventajas competitivas a partir de la información: IDC Big Data 2012.....	13
Retos empresariales de Big Data.....	14
El gran negocio de Big Data.....	14
Big Data: <i>the next thing</i> (la siguiente gran tendencia).....	15
La empresa inteligente.....	15
Casos de estudio	16
Una breve reseña histórica de Big Data	18
El origen moderno de Big Data	18
Resumen	20
Notas.....	21

CAPÍTULO 2

FUENTES DE GRANDES VOLÚMENES DE DATOS

Origen de las fuentes de datos	24
Tipos de fuentes de Big Data	25
Los datos de la Web.....	27
El peso de los datos de la Web	29
Los datos de texto	30
Aplicaciones del análisis de texto	31
Otras aplicaciones del análisis de texto	32
Datos de sensores.....	33
Datos de posición y tiempo: geolocalización	34
Datos de RFID y NFC	36
Datos de redes sociales	37
Análisis de redes sociales.....	38
Datos de las operadoras de telecomunicaciones	40
El valor del tráfico de datos	41
Datos de las redes inteligentes de energía (<i>smart grids</i>).....	41
El contador inteligente (<i>smart meter</i>) ..	42
Otros datos de las redes inteligentes....	42
Resumen.....	43
Notas	44

CAPÍTULO 3

EL UNIVERSO DIGITAL DE DATOS. EL ALMACÉN DE BIG DATA

“La era del petabyte” (<i>Wired</i> , 2008)	46
--	----

El universo digital de EMC/IDC (2007-2010)	47
Datos en todas partes (<i>The Economist</i> , 2010)	50
El universo digital de datos: “Extrayendo valor del caos” (2011)	52
La sobrecarga de información cobra forma física	55
El almacenamiento también supera las expectativas	55
La revolución de los datos está cambiando el paisaje de los negocios (<i>The Economist</i> , 2011)	56
La era del exabyte (Cisco, 2012). Hacia la era del zettabyte	57
El universo digital de datos IDC/EMC (diciembre, 2012). El camino a la era del zettabyte	60
Resumen	61
Notas	62

CAPÍTULO 4
SECTORES ESTRATÉGICOS DE BIG DATA Y OPEN DATA

Dominios estratégicos de Big Data	64
Informe McKinsey Global Institute	64
¿Por qué se ha llegado a la explosión de los Big Data?	66
Sectores dominantes en Big Data	67
Sector de la salud	68
El informe “Big Data Healthcare Hype and Hope”	71
Conclusiones del <i>Digital Health Summit</i> , Las Vegas (Enero 2013)	72
Otras consideraciones prácticas	72
Un anticipo a Hadoop	74
Open Data. El movimiento de los datos abiertos	74
Iniciativas Open Data	76
La información pública al servicio del ciudadano	79
La iniciativa de la Unión Europea (enero 2013)	80
Open Data Alliance	81

Open Data Institute (ODI)	81
Resumen	82
Recursos	83
Notas	84

CAPÍTULO 5
BIG DATA EN LA EMPRESA. LA REVOLUCIÓN DE LA GESTIÓN, LA ANALÍTICA Y LOS CIENTÍFICOS DE DATOS

Integración de Big Data en la empresa	86
Presencia del modelo 3 V de Big Data en las empresas	87
Big Data: la revolución de la gestión	89
¿Qué es lo nuevo ahora?	89
Los cinco retos de la gestión	90
Profesionales de análisis de datos: analistas y científicos de datos	92
Ciencia de los datos	94
El científico de datos	96
¿Qué habilidades necesita un científico de datos?	96
Casos de estudio: el ITAM de México DF	99
¿Cómo encontrar los científicos de datos que se necesitan?	99
La inteligencia de negocios en Big Data	100
OLAP	102
Minería de datos	102
Sistemas de apoyo a la decisión (DSS)	103
Herramientas de informes y de visualización	103
Tecnologías de visualización de datos	104
Analítica de Big Data: una necesidad	105
Seguridad y privacidad en Big Data	107
La iniciativa de Cloud Security Alliance (CSA)	108
Privacidad	109
Foursquare. Un caso de estudio en privacidad	109
La seguridad en la Unión Europea	110
Resumen	110
Recursos	111
Notas	112

Parte II. Infraestructura de los Big Data

CAPÍTULO 6
CLOUD COMPUTING, INTERNET DE LAS COSAS Y SOLOMO 113

Origen y evolución de <i>cloud computing</i>	114	
Definición de la nube	115	
Características de <i>cloud computing</i>	117	
Modelos de la nube (<i>cloud</i>).....	120	
Modelos de servicio	121	
Modelos de despliegue de la nube	123	
¿Cómo adaptar la nube en organizaciones y empresas?.....	124	
Consideraciones económicas	124	
Características organizacionales	125	
Acuerdos de nivel de servicio (SLA, Service Level Agreement)	125	
Seguridad	126	
Los centros de datos como soporte de <i>cloud computing</i>	126	
Internet y los centros de datos: una industria pesada	127	
Internet de las cosas	128	
IPv4: El cuello de botella. IPv6: el desarrollo de la Internet de las cosas....	132	
Sensores.....	133	
Bluetooth 3.0/4.0.....	134	
RFID.....	135	
NFC.....	136	
SIM integrada.....	137	
Códigos QR y BIDI	138	
Ciudades inteligentes (<i>smart cities</i>)	139	
¿Qué son los medios sociales (<i>social media</i>)?	139	
El panorama de los medios sociales	141	
Geolocalización	142	
Movilidad	144	
Plataformas móviles.....	145	
Plataformas móviles de código abierto.	147	
Resumen	149	
Recursos.....	150	
Notas.....	152	
 CAPÍTULO 7		
ARQUITECTURA Y GOBIERNO DE BIG DATA		153
La arquitectura de Big Data.....	154	
Fuentes de Big Data	155	
Almacenes de datos (Data Warehouse y Data Marts)	156	
Bases de datos	157	
Hadoop	158	
Plataformas de Hadoop	158	
Integración de Big Data	158	
Analítica de Big Data.....	159	
<i>Reporting, query</i> y visualización.....	159	
Analítica predictiva	160	
Analítica Web	160	
Analítica social y <i>listening social</i>	160	
Analítica M2M	161	
Plataformas de analítica de Big Data	162	
<i>Cloud computing</i>	162	
Gobierno de los Big Data	163	
Gobierno de TI.....	163	
El gobierno de la información.....	165	
Gobierno de Big Data.....	165	
Calidad de los Big Data	166	
Administración de datos maestros	167	
El ciclo de vida de los Big Data.....	168	
Seguridad y privacidad de Big Data.....	168	
Metadatos de Big Data	169	
Arquitectura de Big Data de Oracle	169	
Capacidades de la arquitectura de Big Data	169	
Arquitectura de información de Big Data de Oracle	170	
Plataforma de Big Data de Oracle: productos y soluciones	171	
Arquitectura de Big Data de IBM	173	
Resumen.....	174	
Notas	175	
 CAPÍTULO 8		
BASES DE DATOS ANALÍTICAS: NOSQL Y “EN MEMORIA”		177
Tipos de base de datos actuales	178	
Bases de datos relacionales	178	
Bases de datos heredadas (<i>legacy</i>)	179	
Bases de datos NoSQL	180	
Bases de datos “en memoria”	180	
Sistemas de base de datos MPP	181	
¿Qué es NoSQL?	182	
Bases de datos NoSQL	183	
Diferencias esenciales entre NoSQL y SQL.....	185	
Tipos de base de datos NoSQL.....	185	
Bases de datos clave- valor.....	186	
Bases de datos orientadas a grafos.....	188	
Bases de datos orientadas a BigTable (tabulares/columnares)	189	

Bases de datos orientadas a documentos	191
Bases de datos “en memoria” caché.....	193
Las bases de datos NoSQL en la empresa	193
Breve historia de NoSQL	194
Tendencias para 2013 en bases de datos NoSQL	195
Computación “en memoria”	196
Tecnología “en memoria”	196
Tipos de tecnologías “en memoria”	197
Proveedores de tecnología “en memoria”	198
Analítica “en memoria”	198
Proveedores de computación y bases de datos “en memoria”	199
Bases de datos “en memoria”	200
Uso de la memoria central como almacén de datos	200
Almacenamiento por columnas	202
Paralelismo en sistemas multinúcleo	203
SAP HANA	203
SAP HANA cloud	204
SAP HANA para análisis de sentimientos	205
Oracle.....	205
Microsoft	206
Resumen	206
Recursos.....	207
Notas.....	209

CAPÍTULO 9

EL ECOSISTEMA HADOOP	211
El origen de Hadoop	212
The Google File System	212
MapReduce	213
BigTable	213
¿Qué es Hadoop?	213
Historia de Hadoop	216
El ecosistema Hadoop	218
Componentes de Hadoop	218
MapReduce	220
El enfoque de gestión de MapReduce... ..	221
Hadoop Common Components.....	222
Desarrollo de aplicaciones en Hadoop	222
Hadoop Distributed File Systems (HDFS)	223
Consideraciones teórico-prácticas	224

Mejoras en la programación de Hadoop	225
Pig.....	225
Hive.....	226
Jaql.....	227
Zookeeper.....	227
HBase.....	228
Lucene	228
Oozie.....	228
Avro	228
Cassandra	229
Chukwa	229
Flume.....	229
Plataformas de Hadoop	229
Resumen	231
Recursos	232
Notas	234

Parte III. Analítica de Big Data

CAPÍTULO 10

ANÁLITICA DE DATOS (BIG DATA

ANALYTICS)	237
Una visión global de la analítica de Big Data	238
¿Qué es analítica de datos?	240
Tipos de datos de Big Data	241
Datos estructurados	242
Datos semiestructurados	242
Datos no estructurados	242
Datos en tiempo real	242
Analítica de Big Data.....	243
Tecnologías, herramientas y tendencias en analítica de Big Data	244
Proveedores de analítica de Big Data (distribuciones comerciales)	245
Tecnologías de código abierto de Big Data ..	251
Casos de estudio.....	254
Características de una plataforma de integración de analítica de Big Data	255
Resumen.....	256
Notas	257

CAPÍTULO 11

ANÁLITICA WEB

Analítica Web 2.0.....	260
Breve historia de la analítica Web	261
Enfoques de analítica Web	262
Métricas.....	262

Visitas.....	263
Visitante.....	263
Visitante único.....	264
Tiempo en la página y en el sitio.....	265
Tasa de rebote.....	265
Tasa de salida.....	265
Tasa de conversión.....	266
Compromiso.....	266
Otras métricas.....	267
Indicadores clave de rendimiento (KPI).....	268
Casos prácticos.....	269
Informes (Google Analytics).....	270
Informes estándar.....	270
Informes personalizados.....	271
Informes sociales.....	271
Segmentación.....	271
Herramientas de analítica Web.....	272
Analítica Web móvil (Mobile analytics).....	274
Información de las herramientas de analítica móvil.....	275
Herramientas de analítica móvil.....	275
Caso de estudio: Google Analytics.....	276
Resumen.....	277
Recursos.....	278
Notas.....	279

CAPÍTULO 12

ANÁLITICA SOCIAL	281
El exceso de información: un problema global.....	282
La proliferación de datos sociales.....	283
¿Qué es analítica social?.....	284
Métricas sociales.....	285
Métricas de sitios Web.....	286
Métricas de <i>social media</i>	286
Indicadores clave de rendimiento (KPI).....	288
Diferencias entre métricas y KPI.....	289
Ejemplo práctico simple de métrica versus KPI.....	289
Herramientas de analítica social.....	290
Estadística social.....	291
Herramientas de investigación. Monitorización.....	292
Herramientas globales muy reconocidas.....	293
Herramientas de analítica Web social.....	294
Herramientas de reputación e influencia social.....	295

Herramientas de medida de influencia.....	295
Herramientas de reputación corporativa.....	296
Herramientas de análisis de actividad en redes.....	297
Facebook.....	297
Twitter.....	298
Herramientas de gestión de multiplataforma y multiperfiles.....	299
Análisis de sentimientos.....	300
Herramientas de análisis de sentimientos.....	301
Casos de estudio de analítica social.....	303
BBVA.....	303
Universidad de Alicante.....	303
Social Relationship Management de Oracle.....	303
Otras herramientas.....	304
Resumen.....	304
Notas.....	305

Parte IV. El futuro de la era Big Data

CAPÍTULO 13

LAS NUEVAS TENDENCIAS

TECNOLÓGICAS Y SOCIALES QUE TRAEN

LA NUBE Y LOS BIG DATA..... 307 |

El nexo de la fuerza.....	308
BYOD.....	309
¿Qué es el movimiento BYOD?.....	310
¿Cómo puede el departamento informático gestionar y proteger los dispositivos móviles de los empleados?.....	310
Ventajas y riesgos.....	311
Los hábitos del trabajo.....	311
El impulso debe venir de las compañías.....	312
Consumerización de TI.....	313
El meteórico ascenso de los dispositivos móviles personales.....	315
¿Cómo puede beneficiarse su empresa de la consumerización?.....	315
El informe de ENISA sobre la consumerización en las empresas.....	316
Crowdsourcing.....	317
Casos de estudio.....	318
Crowdfunding.....	319
Características del crowdfunding.....	320
Casos de estudio de crowdfunding.....	320

Reseña histórica del <i>crowdfunding</i>	322
<i>Gamificación /Ludificación</i>	322
¿Dónde utilizar la ludificación?	323
Ventajas de la <i>gamificación</i>	323
Resumen	324
Recursos.....	324
Notas.....	325

CAPÍTULO 14

BIG DATA EN 2020	327
Los retos del futuro	328
Los dominios de Big Data sin explorar... 328	
Necesidad incumplida de proteger los	
datos	329
El protagonismo de los países emergentes	
.....	329
La tercera plataforma.....	330
Analítica M2M: ¿El próximo reto para el Big	
Data?.....	331
M2M: Oportunidad de Big	
Data para operadores móviles	332
Internet de las cosas (<i>the Internet of the</i>	
<i>things</i>)	333
Analítica predictiva	333
Análisis de sentimientos	333
¿Cómo va a cambiar la vida por Big Data en el	
año 2013?	334
¿Cómo Big Data y <i>cloud computing</i> van a cambiar	
el entretenimiento en el año 2013?.....	335
¿Cómo va a cambiar la salud por Big Data? .	336
¿Cómo pueden afectar los Big Data a la actividad	
física y al deporte?	336
La cara humana de Big Data.....	337
Big Data y las tendencias tecnológicas en 2013	
(Gartner)	340
El mercado futuro de Big Data	341
Las cinco grandes predicciones “muy	
profesionales” de Big Data para 2013.....	341
Emergencia de una arquitectura de Big	
Data.....	342
Hadoop no será la única oferta	
profesional	342
Plataformas de Big Data “llave en mano”	
.....	342
El centro de atención será el gobierno	
de datos	342

Emergencia de soluciones de analítica	
“extremo a extremo” (<i>end-to-end</i>)	343
El futuro seguirá sin ser lo que era	343
Notas	344

APÉNDICE A

EL PANORAMA DE BIG DATA (THE BIG	
DATA LANDSCAPE)	347

APÉNDICE B

PLATAFORMAS DE BIG DATA (DOUG	
HENSCHEN)	351

APÉNDICE C

PLATAFORMAS DE HADOOP (DOUG	
HENSCHEN)	361

APÉNDICE D

GLOSARIO	373
-----------------------	-----

APÉNDICE E

BIBLIOGRAFÍA Y RECURSOS WEB ...	393
--	-----

Prólogo

Big Data (*grandes datos* o **macrodatos** según la Fundación Fundéu BBVA) supone la confluencia de una multitud de tendencias tecnológicas que venían madurando desde la primera década del siglo XXI y se han consolidado durante los años 2011 y 2012 cuando han explotado e irrumpido con gran fuerza en organizaciones y empresas, en particular, y en la sociedad, en general: movilidad, redes sociales, aumento de la banda ancha y reducción de su coste de conexión a Internet, medios sociales –en particular, las redes sociales–, Internet de las cosas, geolocalización y, de modo muy significativo, la computación en la nube (*cloud computing*).

Los *grandes volúmenes de datos* han ido creciendo de modo espectacular. Durante 2012, se crearon 2,8 zettabytes (ZB) de datos (1 ZB = 1 billón de gigabytes) según datos de la consultora IDC en el estudio “El Universo Digital de Datos 2012” publicado en diciembre de 2012 y esta cifra se dobla cada dos años. Un dato significativo, Walmart, la gran cadena de almacenes de Estados Unidos, posee bases de datos con una capacidad de 2,5 petabytes y procesa más de un millón de transacciones cada hora. Los *Big Data* están brotando por todas partes y utilizándolos adecuadamente proporcionarán una gran ventaja competitiva a las organizaciones y empresas. La ignorancia de los *Big Data* producirá grandes riesgos en las organizaciones y no las hará competitivas. Para ser competitivas en el siglo actual, como señala Franks (2012)¹, “es imperativo que las organizaciones persigan agresivamente la captura y análisis de estas nuevas fuentes de datos para alcanzar los conocimientos y oportunidades que ellas ofrecen”.

Big Data ya es una realidad consolidada. La consultora Gartner ha cuantificado el gasto en *Big Data* en 2012 en 28.000 millones de dólares y prevé para el año 2013, la cantidad de 34.000 millones de dólares. A su vez, la auditora Deloitte estima que a finales de 2012 más del 90% de las empresas del índice *Fortune 500* podrán poner en marcha iniciativas de *Big Data*. Por estas razones, los profesionales de *Big Data*, los *analistas de datos* y *científicos de datos*, tienen mucho trabajo por delante y será una de las profesiones más demandadas en 2013 y años sucesivos.

En el libro, además de introducir al lector en los fundamentos de los volúmenes masivos de datos, las tecnologías y herramientas de *Big Data*, se estudiarán las diferentes formas en las que una organización puede hacer uso de esos grandes datos para sacar mayor rendimiento en su toma de decisiones y trataremos de las

oportunidades que traerán consigo su adopción y los riesgos de su no adopción, dado el gran cambio social que se prevé producirá el enorme volumen de datos que se irán creando y difundiendo.

El diluvio de datos

La avalancha o aluvión de datos que cada día genera, captura, almacena y analiza las organizaciones y empresas y, por ende, los particulares, ha dado lugar a la nueva tendencia *Big Data*. Situémonos en un día cualquiera del año; imaginemos los millones de usuarios que visitan Facebook, los millones de *tuits* (*tweets*) que se publican a diario, los millones de mensajes y conversaciones que se realizan a través de WhatsApp, Joyn o Line, los millones de correos electrónicos que envían y reciben millones de personas de todo el mundo, los miles de llamadas telefónicas y videoconferencias a través de Skype. Sumemos a toda ese ingente volumen de información, las páginas que visitan dichos usuarios, las noticias que leen, las ofertas de anuncios, ventas, alquileres, etcétera; las visitas a sitios de turismo, de ocio, de cultura... Multiplique esa información personal de más de 2.800 millones de internautas en el mundo. En esencia, hablamos de datos de la Web y de los medios sociales (Social Media).

Por si ese volumen de información no fuera significativo, añadamos ahora los datos que se transfieren entre sí los miles de millones de objetos o cosas, que se comunican entre sí, a través de sensores, chips NFC chips/etiquetas de RFID, etcétera, es decir, la interconexión de datos entre máquinas (M2M) origen del conocido Internet de las cosas o también de los objetos.

Sigamos sumando, datos médicos de los millones de hospitales, hoy prácticamente digitalizados en su gran mayoría; datos de las administraciones públicas, prácticamente todos en línea –al menos en la mayoría de los países del mundo–; datos de posición y geolocalización, sistemas de información geográfica (SIG/GIS) a través de sistemas GPS y teléfonos inteligentes (*smartphones*), además de miles de satélites de comunicaciones, etcétera.

En resumen, la explosión de los grandes volúmenes de datos (*Big Data*) no para de crecer y parece que de modo exponencial. Eric Schmidt, el presidente ejecutivo de Google, ya advertía hace unos años que: “entre el origen de la Tierra y el año 2003 se crearon 5 exabytes de información. Hoy día creamos la misma cantidad cada dos días”. El estudio citado de IDC/ EMC ya confirmaba también que las cifras del Presidente de Google: alrededor de 2,5 exabytes de datos se creaban cada día en el año 2012 y –es más– el número se doblaba cada 40 meses aproximadamente. El estudio de IDC pronostica que el año 2020 se alcanzará en la Tierra, los 40 zettabytes (40 ZB), se han creado 2,8 ZB de datos durante el año 2012, lo que significa que se generarán 5,247 gigabytes (GB) por cada persona existente en el mundo en ese año. Pero lo sorprendente no sólo es este inmenso caudal de datos, sino que en la presentación del informe se revela que menos del

1% de los datos del mundo se analizan para aprovecharse de esa gran ventaja y valor añadido que suponen los *Big Data* y también que menos del 20% de los datos no están protegidos. El informe advierte de las grandes oportunidades que se ofrecen a las empresas para la protección y extracción del valor que suponen este inmenso volumen de datos.

La revolución de la gestión

Andrew McAfee² y Erik Brynjolfsson³ profesores del MIT publicaron un artículo significativo, en el número de octubre de 2012 de la prestigiosa revista *Harvard Business Review*, “Big Data: The Management Revolution”⁴, Las conclusiones fundamentales de su estudio son claras: “La explotación de los nuevos y espectaculares flujos de información pueden mejorar radicalmente el desempeño (rendimiento) de su empresa, Sin embargo, será necesario cambiar su cultura de toma de decisiones”. La propuesta final de su artículo es concluyente: “La evidencia es clara: las decisiones controladas por los datos tienden a ser mejores decisiones: los líderes empresariales o bien adoptan esta situación o serán remplazados por otros que lo hagan”.

Es decir, las organizaciones se debaten entre adoptar o no *Big Data*, al igual que estos primeros años de la actual década, el debate empresarial era la adopción o no de la computación en la nube⁵ con las consiguientes conclusiones para su adopción y migración a la misma de modo radical o gradual, dependiendo de las estrategias de cada organización. La adopción de *Big Data* parece que es un hecho que tarde o temprano deben realizar las organizaciones; los retos y oportunidades que ofrecen compensarán los gastos económicos y de talento que se requerirán y serán compensados con la ventaja competitiva que supondrá dicha adopción y el análisis de esos grandes volúmenes de datos implicarán una gran mejora en la toma de decisiones.

La investigación realizada por los profesores McAfee y Brynjolfsson, no deja lugar a dudas, “las empresas que inyectan *Big Data* y analítica de *Big Data* en sus operaciones muestran tasas de productividad y rentabilidad que son del orden del 5 al 6% más altas que aquellas de la competencia o compañías homólogas”⁶.

El científico de datos (*data scientist*): la nueva profesión sexy del siglo XXI

Así titulan su artículo en el citado número de HBR, los prestigiosos analistas Thomas Davenport⁷ y D. J. Patil⁸. Evidentemente, los *Big Data*, como ya sucedió con la Web 2.0 y el advenimiento de los medios sociales (*social media*), ha traído nuevos roles el mundo del trabajo así como nuevas profesiones.

La Web 2.0 y los medios sociales han traído: Analista Web, especialista SEO y SEM, Community Manager (Gestor de comunidades), Social Media Manager (Director de medios sociales), analistas sociales... Estas nuevas profesiones están dando paso a los analistas de *Big Data* y, de modo muy especial, al científico de datos (*data scientist*) que convive con el analista de datos y el analista de negocios tradicionales.

Las profesiones de la Web y Social Media ya han llegado a las organizaciones y empresas y su formación académica ya ha entrado en la Universidad, en las Escuelas de Negocios y en los departamentos de formación de las grandes empresas así como en las universidades corporativas. Ahora comienzan a llegar los analistas de *Big Data* y de manera muy significativa, por ser el más demandado y más escaso, el **científico de datos**.

¿Qué es un científico de datos?

Aunque el término, como casi siempre sucede con las ideas de impacto, no está totalmente definido, en cuanto a las personas, sí parece claro que nació en las grandes compañías clásicas de Internet controladas por datos, desde el punto de vista de uso en las empresas y en la industria, tales como Google, Amazon, Facebook, Twitter o LinkedIn, y algunas otras también del mundo de negocios de Internet como eBay, PayPal o de gran éxito en *retailing* o en ventas al por menor (grandes almacenes) como es el caso de Walmart.

Sin embargo, hay cierta unanimidad en fijar a D. J. Partil y Jeft Hammerbadier entonces líderes respectivos de análisis de datos en LinkedIn y Facebook, en el año 2008, como acuñadores del término⁹, aunque en 2009, Troy Sadkovsky creó un grupo de investigación en LinkedIn y usó el término para definir una nueva profesión (la suya, por otra parte).

Entonces, ¿qué es un científico de datos? Precisamente, Davenport y Partil, centran todo su artículo de HBR para tratar de definir el nuevo rol del científico de datos ¿Qué tipo de persona debe ser? ¿Qué formación académica le debe respaldar? ¿Qué capacidades y competencia ha de poseer? En el citado artículo, se decanta por un híbrido de hacker, analista, comunicador de datos y asesor de confianza. La combinación es extremadamente potente y rara, confiesan los autores.

Las grandes empresas informática y de Internet parece se decantan por definiciones y roles diversos:

1. “Un buen *data scientist* ha de tener diferentes capacidades: saber matemáticas, tener capacidad analítica y formación en estadística, pero ha de saber contar una historia y tener curiosidad porque se trata de crear significado y valor sobre los datos” (Sonderegger, director senior de *Analytics* en Oracle).
2. “Evolución del analista de datos o de negocios en el contexto de *Big Data*: se considera mitad analista, mitad artista” (Gonzalo Smith, responsables de “*smart analytics*” de IBM GBS España).
3. El científico de datos será aquel que tenga el trabajo más sexy del siglo XXI (Davenport y Partil).
4. Es una persona con habilidades serias en diversas disciplinas técnicas, como ciencias de la computación (informática), analítica, matemáticas, generación de modelos y estadística. Además, debe ser un buen comunicador que sea capaz de entender un problema de negocios, transformar ese problema en un plano analítico, ejecutar el plan y luego dar ¹⁰una solución d negocios (Ani Kaul, CEO de AbsoluData, empresa de analítica e investigación de Alameda, California):

En síntesis, el científico de datos es una profesión emergente. Existen muchos científicos de datos en Google, Amazon, Facebook, LinkedIn, Twitter, y... Todavía existen pocas ofertas de formación en el mundo académico, no sólo iberoamericanas sino de Europa y de Estados Unidos, pero, sin lugar a dudas, estas ofertas irán creciendo poco a poco. En el capítulo 5, se analiza el rol del científico de datos; se citan algunas iniciativas importantes como el caso del prestigioso ITAM de México DF que puso en marcha en 2012, una maestría en Ciencia de Datos.

El análisis de los grandes volúmenes de datos

El análisis de datos y de negocios, son disciplinas antiguas que han experimentado notable crecimiento en todos los campos del saber y, en particular, en organizaciones y empresas, por la necesidad de disponer de herramientas que analicen datos y que éstos sirvan para toma de decisiones eficaces y eficientes.

El análisis de datos ha ido evolucionando a medida que los grandes volúmenes de datos crecían. Las herramientas de inteligencia de negocios han ido recogiendo las tecnologías de OLAP (procesamiento analítico en línea), de informes y consultas (*reporting and query*), de visualización y, especialmente de minería de datos con sus

ya asentadas categorías de minería Web y minería de texto, y las innovadoras minería social en el análisis de datos en medios sociales, que se ha apoyado en técnicas de análisis de sentimiento y de opinión, o minería de opinión y minería de sentimiento como también se la conoce.

La analítica de *Big Data* está emergiendo a la vez que la avalancha de los grandes volúmenes de datos sigue creciendo. “La era de los grandes datos (Big Data) está evolucionando rápidamente y toda nuestra experiencia sugiere que la mayoría de las compañías deben actuar ahora [...]. A medida que las compañías aprendan las destrezas fundamentales (*core skills*) para utilizar los *Big Data*, la construcción de capacidades superiores se pueden convertir pronto en un activo competitivo decisivo” (Barton, Court 2012)¹¹.

En el libro, trataremos las categorías de Analítica de Datos que consideramos fundamentales para el estudio de *Big Data*: analítica Web, analítica Social, analítica de sentimientos, analítica M2M y, en general, analítica de *Big Data*.

Arquitectura de Big Data

Una arquitectura de Big Data debe considerar la integración de las nuevas tecnologías y herramientas de los grandes volúmenes de datos y su integración con los datos tradicionales (bases de datos relacionales y heredadas “*legacy*”) así como la integración con la infraestructura existente de las organizaciones y empresas.

Así pues, en el libro, se ha tratado de considerar los conceptos y componentes básicos que las compañías deberán considerar en la gestión y explotación de sus *Big Data* y que enumeramos a continuación:

- *Fuentes de Big Data*. Los datos proceden de la Web, de los Social Media, de interconexión de objetos M2M mediante sensores conocida como “Internet de las cosas”, de la movilidad, biometría, datos procedentes de las propias personas, etcétera.
- Los *tipos de datos* se clasifican en tres grandes categorías: estructurados (datos transaccionales de las bases de datos relacionales), no estructurados (audio, vídeo, fotografía, textos...) y semiestructurados (procedentes, fundamentalmente de archivos HTML, XML).
- Almacenes de datos empresariales (**EDW**, *Enterprise data warehouse*).
- Bases de datos no relacionales (**NoSQL**) que no siguen, normalmente el estándar SQL.
- Bases de datos analíticas “**en memoria**” y **MPP** (procesamiento masivo paralelo).

- Hadoop: el marco de trabajo por excelencia (*framework*) para procesar y analizar los grandes volúmenes de datos, especialmente los datos no estructurados y semiestructurados.
- Analítica de Big Data. Herramientas de analítica, de informes (*reporting*), de consultas (*query*) y de visualización (*dashboard*) así como los cuadros de mando integral (*balanced scorecard*) que conducirán a analítica de datos, en su sentido general, y analítica Web, analítica social, analítica de sentimientos, analítica M2M, etcétera.

En la actualidad, las tecnologías y herramientas de *Big Data* se deben centrar en la integración de datos estructurados y datos no estructurados o semiestructurados, así como la integración de los datos tradicionales en las bases de datos relacionales con los datos no estructurados en las bases de datos analíticas y NoSQL. Otro aspecto fundamental que se considerará en esta obra será el tema de la seguridad y la privacidad de los grandes volúmenes de datos.

Toda esta arquitectura de *Big Data* requerirá de plataformas que gestionen estos grandes datos para que las organizaciones y empresas puedan obtener el máximo rendimiento. Para ello, se requieren de proveedores de soluciones que hoy día son muy numerosos y que pueden ser agrupados en proveedores de código propietario o de código abierto (*open source*). Entre los cuales podemos destacar: *Soluciones de Big Data propietarias* (SAP, Oracle, IBM, EMC, HP, SAS...), *Bases de datos NoSQL* (Cassandra, MongoDB, CouchDB...), *bases de datos "en memoria"* donde sobresalen SAP HANA, aunque Oracle, IBM, Teradata, EMC ofrecen soluciones similares, integración de todas estas herramientas en el *marco de trabajo Hadoop* con plataformas eficientes e innovadoras como Cloudera, Hortonworks o MapR entre otras.

Innovaciones tecnológicas que han acelerado los Big Data

Las innovaciones que han acelerado la explosión y avalancha de los grandes volúmenes de datos son muchas, pero cuatro son los grandes pilares sobre los que se sustentan las tecnologías de *Big Data*: Los medios sociales "**Social media**", la **Movilidad** (teléfonos inteligentes, tabletas... y aplicaciones (*apps*)), **Cloud Computing** (Computación en la Nube) e **Internet de las cosas** (M2M, sensores de todo tipo, chips **NFC, RFID**...).

Sin embargo, las tendencias se segmentan y las tecnologías, dispositivos y aplicaciones Web más innovadoras crecen casi exponencialmente y son fuentes y origen de datos de todo tipo. Así se pueden considerar las tendencias tecnológicas que vendrán controladas por las multitudes inteligentes (*crowds*) y que influirán en la explosión de los grandes volúmenes de datos: *crowdsourcing*, *crowdfunding*,

BYOD (Bring Your Own Device), *consumerización* y *gamificación*, fundamentalmente. El gran volumen de datos que se irán generando se verán notablemente influenciadas por las tecnologías anteriores y las tendencias que se presuponen tendrán gran impacto en organizaciones y empresas.

Organización del libro

El libro se ha estructurado en cuatro partes que pretenden abarcar las partes fundamentales de las tecnologías y las estrategias de *Big Data*, así como la arquitectura nuclear de los *Big Data*, unidas a la verdadera razón de ser de los *Big Data*, la teoría y herramienta de análisis de los grandes volúmenes de datos con el objetivo esencial de ayudar en la toma de decisiones de organizaciones y empresas. El libro se complementa con varios anexos prácticos donde se describen las plataformas, proveedores y soluciones comerciales más implantadas en la actualidad a nivel mundial, precedidas de un informe sobre el panorama o paisaje actual de las *Big Data* donde se mostrarán de un modo integrado las diferentes plataformas, proveedores y herramientas que componen la oferta comercial a disposición de organizaciones y empresas. Asimismo, se incluye un glosario de términos de *Big Data* que faciliten la comprensión por parte del lector de los numerosos y variados conceptos relacionados con estas tecnologías y tendencias.

PARTE I. LA ERA DE LOS BIG DATA

Capítulo 1. ¿Qué es Big Data?

Capítulo 2. Fuentes de grandes volúmenes de datos

Capítulo 3. El Universo Digital de Datos: El almacén de Big Data

Capítulo 4. Sectores estratégicos de Big Data y Open Data

Capítulo 5. Big Data en la empresa: La revolución de la gestión, la analítica y los científicos de datos

PARTE II. INFRAESTRUCTURA DE LOS BIG DATA

Capítulo 6. Cloud Computing, Internet de las cosas y SoLoMo

Capítulo 7. Arquitectura y gobierno de Big Data

Capítulo 8. Bases de datos analíticas: NoSQL y “en memoria (*in-memory*)”

Capítulo 9. El ecosistema Hadoop

PARTE III. ANALÍTICA DE BIG DATA

Capítulo 10. Analítica de datos

Capítulo 11. Analítica Web

Capítulo 12. Analítica social

Parte IV. EL FUTURO DE BIG DATA

Capítulo 13. Las nuevas tendencias tecnológicas y sociales que traen la Nube y los Big Data

Capítulo 14. Big Data en el horizonte 2020

APÉNDICES

- A. El panorama de Big Data
- B. Plataformas de Big Data
- C. Plataformas de Hadoop
- D. Glosario
- E. Bibliografía y Recursos

Agradecimientos

En primer lugar, deseo expresar mi agradecimiento a mis alumnos de Ingeniería Informática e Ingeniería de Organización Industrial de la Universidad Pontificia de Salamanca en las asignaturas de Sistemas Informáticos, Sistemas de Información, Gestión del Conocimiento e Inteligencia de Negocios, que me han permitido experimentar las tecnologías, herramientas y aplicaciones de *Big Data* tanto en las clases teóricas como en los numerosos talleres y trabajos que ellos han realizado como actividades académicas. A mis alumnos de máster y doctorado de la Universidad Pontificia de Salamanca (campus Salamanca), Universidad Carlos III de Madrid (Documentación) y Universidad Nebrija de Madrid (Empresa) donde he podido implantar todos los conocimientos que he ido adquiriendo en la ciencia y arte de los grandes volúmenes de datos. También quiero agradecer a mis estudiantes de doctorado de la Pontificia de Salamanca del campus de Madrid que investigan conmigo en *cloud computing*, *Big Data*, movilidad y *social media* así como en áreas aplicadas como negocios digitales, educación, gestión del conocimiento,

inteligencia de negocios y sistemas de información geográfica, entre otras, a los que ya han leído su tesis doctoral y a los que espero que lean muy pronto.

Asimismo, deseo agradecer a los alumnos y profesores de las universidades latinoamericanas donde he impartido conferencias, cursos y talleres en los últimos dos años y donde he tratado los temas de **Cloud Computing** y **Big Data: México** (TEC de Monterrey, campus Cuernavaca, Universidad del Valle, Instituto Politécnico de México, Instituto Tecnológico de Tijuana, Instituto Tecnológico Superior de Coahuila (ITESCO), Universidad Autónoma de Baja California –sedes de Tijuana y de Ensenada), **Perú** (Universidad San Martín de Porres, Universidad Tecnológica de Perú y Universidad Garcilaso de la Vega), **Panamá** (Universidad Tecnológica de Panamá), **República Dominicana** (ITLA “Instituto Tecnológico de las Américas”, Universidad Unibe, Universidad APEC y la Fundación Funglode).

Por último, quiero agradecer al director editorial Alberto Umaña y al gerente editorial Marcelo Grillo por su apoyo a la colección de libros **NTICS** que tengo el honor de dirigir y que espero que continúen bien pronto con los siguientes números y en esta ocasión particular a mi editor y, no obstante, amigo, Damián Fernández que desde Buenos Aires me ha acompañado día a día en esta ardua tarea que ha sido la producción de esta obra de *Big Data* y que ha colaborado estrechamente conmigo no sólo como editor sino como un gran amigo que me ha ayudado cuando ha sido menester en esta larga y prolífica tarea de lanzar esta obra sobre un tema tan innovador y de futuro como es *Big Data*. Gracias al equipo editorial.

En **Carchelejo (Sierra Mágina)**, Andalucía (España) y en **México DF** (México), a diecisiete de Mayo de 2013, *Día Mundial de Internet*.

Luis Joyanes Aguilar

Catedrático de Lenguajes y Sistemas Informáticos de la Universidad Pontificia de Salamanca

NOTAS

¹ Bill Franks (2011). *Taming the Big Data Tidal Wave*, New Jersey: Wiley, p.3.

² Investigador del MIT’s Center for Digital Business y autor de *Enterprise 2.0* (Harvard Business School Press, 2009).

³ Director del Center for Digital Business, del MIT's Sloan School of Management. Con Andrew McAfee, son autores de *Race Against the Machine Digital Frontier*, 2012.

⁴ Andrew McAfee⁴ y Erik Brynjolfsson, "Big Data: The Management Resolutor", *Harvard Business Review*, 2012.

⁵ Luis Joyanes. *Computación en la nube. Estrategias de cloud computing en las empresas*. México: Alfaomega, 2012.

⁶ Op. cit.p. 78.

⁷ Profesor visitante de Harvard Business School y autor de varios libros de Gestión del Conocimiento y de Analítica. El autor de esta obra leyó varias de sus obras mientras se formaba en el área de Gestión del Conocimiento.

⁸ Científico de datos y antiguo director de productos de datos de la red social LinkedIn.

⁹ Desde el punto de vista científico, el origen se remonta a los años sesenta donde comenzó a utilizarse en artículos científicos.

¹⁰ Entrevista realizada en la revista *Information Week*, edición de México, publicada el 22 de noviembre de 2012, en su edición digital: www.informationweek.com.mx.

¹¹ Dominic Barton y David Court. *Making Advanced Analytics Work For You*. HBR, Octubre 2012, p. 88.



CAPÍTULO 1

¿QUÉ ES BIG DATA?

Big Data (grandes datos, grandes volúmenes de datos o *macrodatos* como recomienda utilizar la Fundación Fundéu BBVA “Fundación del español urgente”) supone la confluencia de una multitud de tendencias tecnológicas que venían madurando desde la primera década del siglo XXI, y que se han consolidado durante los años 2011 a 2013, cuando han explotado e irrumpido con gran fuerza en organizaciones y empresas, en particular, y en la sociedad, en general: movilidad, redes sociales, aumento de la banda ancha y reducción de su coste de conexión a Internet, medios sociales (en particular las redes sociales), Internet de las cosas, geolocalización, y de modo muy significativo la computación en la nube (*cloud computing*).

Los grandes datos o grandes volúmenes de datos han ido creciendo de modo espectacular. Durante 2011, se crearon 1,8 zettabytes de datos (1 billón de gigabytes) según la consultora IDC, y esta cifra se dobla cada dos años. Un dato significativo, Walmart, la gran cadena de almacenes de los Estados Unidos, posee bases de datos con una capacidad de 2,5 petabytes, y procesa más de un millón de transacciones cada hora. Los Big Data están brotando por todas partes y utilizándolos adecuadamente proporcionarán una gran ventaja competitiva a las organizaciones y empresas. En cambio, su ignorancia producirá grandes riesgos en las organizaciones y no las hará competitivas. Para ser competitivas en el siglo actual, como señala Franks (2012): “Es imperativo que las organizaciones persigan agresivamente la captura y análisis de estas nuevas fuentes de datos para alcanzar los conocimientos y oportunidades que ellas ofrecen”.

Los profesionales del análisis de datos, los analistas de datos y científicos de datos, tienen mucho trabajo por delante y serán una de las profesiones más demandadas en el 2013 y años sucesivos.

En este capítulo, introduciremos al lector en el concepto de Big Data, y en las diferentes formas en que una organización puede hacer uso de ellos para sacar mayor rendimiento en su toma de decisiones. No solo en su concepto con las definiciones más aceptadas, sino que estudiaremos las oportunidades que traerá consigo su adopción, y los riesgos de su no adopción, dado el gran cambio social que se prevé producirá el enorme volumen de datos que se irán creando y difundiendo.

DEFINICIÓN DE BIG DATA

No existe unanimidad en la definición de Big Data, aunque sí un cierto consenso en la fuerza disruptiva que suponen los grandes volúmenes de datos y la necesidad de su captura, almacenamiento y análisis. Han sido numerosos los artículos (*white papers*), informes y estudios relativos al tema aparecidos en los últimos dos años, y serán también numerosos los que aparecerán en los siguientes meses y años; por esta razón, hemos seleccionado aquellas definiciones realizadas por instituciones relevantes y con mayor impacto mediático y profesional. En general, existen diferentes aspectos donde casi todas las definiciones están de acuerdo y con conceptos consistentes para capturar la esencia de Big Data: crecimiento exponencial de la creación de grandes volúmenes de datos, origen o fuentes de datos y la necesidad de su captura, almacenamiento y análisis para conseguir el mayor beneficio para organizaciones y empresas junto con las oportunidades que ofrecen y los riesgos de su no adopción.

La primera definición que daremos es la de Adrian Merv, vicepresidente de la consultora Gartner, que en la revista *Teradata Magazine*, del primer trimestre de 2011, define este término como: “Big Data excede el alcance de los entornos de *hardware* de uso común y herramientas de *software* para capturar, gestionar y procesar los datos dentro de un tiempo transcurrido tolerable para su población de usuarios”¹.

Otra definición muy significativa es del McKinsey Global Institute² que en un informe muy reconocido y referenciado, de mayo de 2011, define el término del siguiente modo: “Big Data se refiere a los conjuntos de datos cuyo tamaño está más allá de las capacidades de las herramientas típicas de software de bases de datos para capturar, almacenar, gestionar y analizar”. Esta definición es, según McKinsey, intencionadamente subjetiva e incorpora una definición cambiante, “en movimiento” de cómo “de grande” necesita ser un conjunto de datos para ser considerado Big Data: es decir, no se lo define en términos de ser mayor que un número dado de terabytes (en cualquier forma, es frecuente asociar el término Big Data a terabytes y petabytes). Suponemos, dice McKinsey, que a medida que la tecnología avanza en el tiempo, el tamaño de los conjuntos de datos que se definen con esta expresión también crecerá. De igual modo, McKinsey destaca que la definición puede variar para cada sector, dependiendo de cuáles sean los tipos de herramientas de software normalmente disponibles; y cuáles, los tamaños típicos de los conjuntos de datos en ese sector o industria. Teniendo presente estas consideraciones, como ya hemos comentado, los Big Data en muchos sectores hoy día, variarán desde decenas de terabytes a petabytes y ya casi exabytes.