



ANALISIS DATA

Tingkat Lanjut:

Map Reduce, Hadoop,
dan Analisis Basisdata

Wahyudi
Anasari



eureka
media aksara
Anggota IKAPI
No. 225/UTE/2021

☎ 0858 5343 1992
✉ eurekamediaaksara@gmail.com
📍 Jl. Banjaran RT.20 RW.10
Bojongsari - Purbalingga 53362

ISBN 978-623-120-854-5



9 786231 208545

**ANALISIS DATA TINGKAT LANJUT: MAP
REDUCE, HADOOP, DAN ANALISIS
BASISDATA**

**Wahyudi
Anasari**



eureka
media aksara

PENERBIT CV.EUREKA MEDIA AKSARA

**ANALISIS DATA TINGKAT LANJUT: MAP REDUCE,
HADOOP, DAN ANALISIS BASISDATA**

Penulis : Wahyudi
Anasari

Desain Sampul : Eri Setiawan

Tata Letak : Rizki Rose Mardiana

ISBN : 978-623-120-854-5

Diterbitkan oleh : **EUREKA MEDIA AKSARA, JUNI 2024**
ANGGOTA IKAPI JAWA TENGAH
NO. 225/JTE/2021

Redaksi:

Jalan Banjaran, Desa Banjaran RT 20 RW 10 Kecamatan
Bojongsari Kabupaten Purbalingga Telp. 0858-5343-1992

Surel : eurekamediaaksara@gmail.com

Cetakan Pertama : 2024

All right reserved

Hak Cipta dilindungi undang-undang

Dilarang memperbanyak atau memindahkan sebagian atau seluruh isi buku ini dalam bentuk apapun dan dengan cara apapun, termasuk memfotokopi, merekam, atau dengan teknik perekaman lainnya tanpa seizin tertulis dari penerbit.

KATA PENGANTAR

Puji syukur selalu terucap kepada Allah SWT yang sampai saat ini telah memberikan nikmat sehat, sehingga penulis dapat menyelesaikan buku yang berjudul “Analisis Data Tingkat Lanjut: Map Reduce, Hadoop, dan Analisis Basisdata”. Penulis mengucapkan banyak terima kasih pada semua pihak yang sudah terlibat dalam proses pembuatan buku ini, sehingga buku ini bisa hadir di hadapan pembaca.

Buku yang berada di tangan pembaca ini terdiri dari 3 Bab, yaitu:

Bab 1 Mapreduce dan Hadoop

Bab 2 Analisis Basis Data

Bab 3 Mengoperasionalkan Proyek Analisis Data

Penulis menyadari bahwa buku ini masih jauh dari kesempurnaan, sejatinya kesempurnaan hanya milik Allah yang Maha Kuasa, Oleh karena itu kritik dan saran yang membangun demi penyempurnaan buku ini sangatlah dibutuhkan. Akhir kata penulis mengucapkan terima kasih, semoga buku ini bisa membawa manfaat bagi pengembangan ilmu pengetahuan.

DAFTAR ISI

KATA PENGANTAR	iii
DAFTAR ISI.....	iv
DAFTAR TABEL	v
DAFTAR GAMBAR.....	vi
BAB 1 MAPREDUCE DAN HADOOP	1
A. Analisis Data Tidak Terstruktur.....	2
B. Ekosistem Hadoop.....	24
C. NoSQL	58
BAB 2 ANALISIS BASIS DATA.....	62
A. Dasar-dasar SQL.....	63
B. Analisis Teks dalam Basisdata.....	79
C. SQL Lanjutan	85
BAB 3 MENGOPERASIONALKAN PROYEK	
ANALISIS DATA.....	107
A. Mengomunikasikan dan Mengoperasikan Proyek Analisis.....	108
B. Menciptakan Hasil Akhir.....	112
C. Dasar-dasar Visualisasi Data	144
DAFTAR PUSTAKA.....	177

DAFTAR TABEL

Tabel 1. 1	Fungsi Pig Bawaan.....	27
Tabel 1. 2	Contoh-contoh Penyimpanan Data NoSQL.....	61
Tabel 2. 1	Operator Ekspresi Reguler	80
Tabel 2. 2	Elemen Ekspresi Biasa	80
Tabel 2. 3	Modul-modul MADlib	102
Tabel 3. 1	Perbandingan Materi untuk Presentasi Sponsor dan Analisis.....	117
Tabel 3. 2	Alat Bantu Umum untuk Visualisasi Data.....	144
Tabel 3. 3	Metode Representasi Umum untuk Data dan Grafik	163

DAFTAR GAMBAR

Gambar 1. 1	Contoh Cara Kerja MapReduce.....	7
Gambar 1. 2	File yang disimpan dalam HDFS	15
Gambar 1. 3	Mengacak dan Mengurutkan	17
Gambar 1. 4	Menggunakan Penggabung.....	18
Gambar 1. 5	Menggunakan Partisi Khusus	19
Gambar 1. 6	Memilih Alamat Pengiriman Saat Pembayaran	37
Gambar 1. 7	Komponen Pivotal HD Enterprise	56
Gambar 2. 1	Relasi Diagram	63
Gambar 2. 2	Penjualan Mingguan untuk Pengecer Online	88
Gambar 2. 3	Penjualan Mingguan dengan Rata-Rata Bergerak	90
Gambar 2. 4	Penjualan mingguan dengan rata-rata bergerak dan EWMA.....	97
Gambar 3. 1	Data Analisis Siklus Hidup, Fase 6: mengoperasionalkan.....	108
Gambar 3. 2	Output Utama dari Proyek Analitik yang Sukses	110
Gambar 3. 3	Sinopsis Contoh Studi Kasus Yoyodyne Bank	113
Gambar 3. 4	Rencana Analisis untuk Studi Kasus Yoyodyne Bank	114
Gambar 3. 5	Contoh Tujuan Slide Proyek untuk Studi Kasus YoyoDyne	121
Gambar 3. 6	Contoh Situasi Slide & Tujuan Proyek untuk Studi Kasus Yoyodyne	123

Gambar 3. 7	Contoh Ringkasan Slide Eksekutif untuk Studi Kasus Yoyodyne	124
Gambar 3. 8	Slide Anatomi Ringkasan Eksekutif	126
Gambar 3. 9	Contoh yang Menggambarkan Metodologi Proyek untuk Sponsor Proyek.....	128
Gambar 3. 10	Contoh yang Menggambarkan Metodologi Proyek untuk Analisis dan Ilmuwan Data.....	129
Gambar 3. 11	Contoh Deskripsi Model Untuk Proyek Sains Data.....	131
Gambar 3. 12	Contoh Presentasi Poin-Poin Penting dari Proyek Sains Data yang ditampilkan dalam Bentuk Diagram Batang	133
Gambar 3. 13	Contoh Detail Model yang Menunjukkan Jenis Model dan Variable	135
Gambar 3. 14	Rincian Model Yang Membandingkan Dua Variabel Data.....	137
Gambar 3. 15	Contoh Rekomendasi untuk Proyek Ilmu Data.....	138
Gambar 3. 16	Data Pembukaan Toko Selama Empat Puluh Lima Tahun	147
Gambar 3. 17	Data Pembukaan Toko Selama Tiga Puluh Lima Tahun	148
Gambar 3. 18	Data Pembukaan Toko Selama Empat Puluh Lima Tahun, Ditampilkan dalam Bentuk Peta	149
Gambar 3. 19	Distribusi Frekuensi Skor Pengguna	151
Gambar 3. 20	Distribusi Frekuensi dengan Log Skor Pengguna	152

Gambar 3. 21	Distribusi Frekuensi Skor Pengguna Baru ..	153
Gambar 3. 22	Grafik Analisis Stabilitas untuk Penetapan Harga.....	154
Gambar 3. 23	Grafik yang Membandingkan Harga dalam Dolar AS dengan Skor Loyalitas Pelanggan.....	156
Gambar 3. 24	Grafik yang Membandingkan Harga dalam Dolar AS dengan Skor Loyalitas Pelanggan (dengan Representasi Karpet) ...	157
Gambar 3. 25	Model Penetapan Harga Baru yang Diusulkan dibandingkan dengan Harga dalam Dolar AS dengan Karpet	159
Gambar 3. 26	Evolusi Grafik, Contoh Analis dengan Titik-Titik Pendukung	161
Gambar 3. 27	Evolusi Grafik, Contoh Sponsor	162
Gambar 3. 28	Cara Membersihkan Grafik, Contoh 1 (Sebelumnya).....	165
Gambar 3. 29	Cara Membersihkan Grafik, Contoh 1 (Setelah).....	167
Gambar 3. 30	Cara Membersihkan Grafik, Contoh 1 (Tampilan “Setelah” Alternatif)	168
Gambar 3. 31	Cara Membersihkan Grafik, Contoh 2 (Sebelumnya).....	170
Gambar 3. 32	Cara Membersihkan Grafik, Contoh 2 (Setelah).....	172
Gambar 3. 33	Cara Membersihkan Grafik, Contoh 2 (Tampilan Alternatif “Setelah”)	173
Gambar 3. 34	Diagram Batang Sederhana, dengan Dua Dimensi	175

Gambar 3. 35 Diagram Batang yang menyesatkan,
dengan Tiga Dimensi..... 176



**ANALISIS DATA TINGKAT LANJUT: MAP
REDUCE, HADOOP, DAN ANALISIS
BASISDATA**

**Wahyudi
Anasari**



BAB

1

MAPREDUCE DAN HADOOP

Konsep Utama Ekosistem Hadoop MapReduce NoSQL

Pada buku sebelumnya dengan judul, "Association Rules dan Supervised Learning Menggunakan R," hingga buku yang berjudul, "ANALISIS DATA: TIME SERIES DAN TEXT," membahas beberapa metode analisis yang berguna untuk mengklasifikasikan, memprediksi, dan memeriksa hubungan dalam data. Bab ini dan Bab berikutnya " Analisis dalam Basisdata," membahas beberapa aspek pengumpulan, penyimpanan, dan pemrosesan data tidak terstruktur dan data terstruktur. Bab ini menyajikan beberapa teknologi dan alat utama yang terkait dengan perangkat lunak Apache Hadoop, "kerangka kerja yang memungkinkan pemrosesan dataset besar secara terdistribusi di seluruh kelompok komputer yang menggunakan model pemrograman sederhana".

Bab ini berfokus bagaimana Hadoop menyimpan data dalam sistem terdistribusi dan bagaimana Hadoop mengimplementasikan paradigma pemrograman sederhana yang dikenal sebagai MapReduce. Meskipun bab ini membuat beberapa referensi khusus Java, namun referensi khusus Java, satu-satunya pengetahuan prasyarat yang dimaksudkan adalah pemahaman dasar tentang

BAB 2

ANALISIS BASIS DATA

Konsep-konsep Utama MADlib Ekspresi reguler SQL yang ditentukan pengguna Fungsi-fungsi Window

Analisis basisdata adalah istilah yang luas yang menggambarkan pemrosesan data di dalam repositori. Dalam banyak contoh R sebelumnya, data diekstraksi dari sumber data dan dimuat ke dalam R. Salah satu keuntungan dari **analisis basisdata** adalah bahwa kebutuhan untuk memindahkan data ke dalam alat analitik dihilangkan. Selain itu, dengan melakukan **analisis basisdata**, dimungkinkan untuk mendapatkan hasil yang hampir real-time. Aplikasi analisis basisdata termasuk deteksi penipuan transaksi kartu kredit, rekomendasi produk, dan pemilihan iklan web yang disesuaikan untuk pengguna tertentu.

Basisdata sumber terbuka yang populer adalah PostgreSQL. Nama ini merujuk pada sebuah analitik penting dalam basisdata yang dikenal sebagai **Structured Query Language** (SQL). Bab ini membahas dasar dan juga tingkat lanjut dasar maupun lanjutan dalam SQL. Contoh-contoh kode SQL yang disediakan diuji terhadap basisdata Greenplum, yang didasarkan pada PostgreSQL. Namun, konsep-konsep yang disajikan dapat diterapkan pada lingkungan SQL lainnya.

BAB

3

MENGOPERASIONALKAN PROYEK ANALISIS DATA

Konsep Utama Mengomunikasikan dan mengoperasionalkan proyek analitik Membuat kiriman akhi Menggunakan kumpulan materi inti untuk audiens yang berbeda Membandingkan area fokus utama untuk sponsor dan analisis Memahami prinsip visualisasi data sederhana Membersihkan bagan atau visualisasi

Bab ini berfokus pada fase terakhir dari Siklus Hidup Analisis Data: mengoperasionalkan. Dalam fase ini, tim proyek memberikan laporan akhir, kode, dan dokumentasi teknis. Pada akhir fase ini, tim umumnya mencoba untuk membuat proyek percontohan dan mengimplementasikan model yang dikembangkan dari Fase 4 dalam lingkungan produksi. Seperti yang dinyatakan dalam Bab 2, “Siklus Hidup Analisis Data,” tim dapat melakukan analisis yang akurat secara teknis, tetapi jika mereka tidak dapat menerjemahkan hasilnya ke dalam bahasa yang audiens mereka, orang lain tidak akan melihat nilainya, dan upaya serta sumber daya yang signifikan akan terbuang percuma. Bab ini berfokus pada bagaimana membuat ringkasan narasi yang jelas dan kerangka kerja untuk menyampaikan narasi tersebut kepada para pemangku kepentingan utama.

DAFTAR PUSTAKA

- “apache.org,” [Online]. Available: <http://www.apache.org/>. [Accessed 11 February 2014].
- “BSP Tutorial,” [Online]. Available: http://hama.apache.org/hama_bsp_tutorial.html. [Accessed 20 February 2014].
- “Hadoop Pipes,” [Online]. Available: <http://hadoop.apache.org/docs/r1.2.1/api/org/apache/hadoop/mapred/pipes/package-summary.html>. [Accessed 19 February 2014].
- “Hadoop Wiki Disk Setup,” [Online]. Available: <http://wiki.apache.org/hadoop/DiskSetup>. [Accessed 20 February 2014].
- “Hama,” [Online]. Available: <http://hama.apache.org/>. [Accessed 20 February 2014].
- “HBase Key Value,” [Online]. Available: <http://hbase.apache.org/book/regions.arch.html>. [Accessed 28 February 2014].
- “HBase Regionserver,” [Online]. Available: <http://hbase.apache.org/book/regionserver.arch.html>. [Accessed 3 March 2014].
- “HBase Rowkey,” [Online]. Available: <http://hbase.apache.org/book/rowkey.design.html>. [Accessed 4 March 2014].
- “HDFS Design,” [Online]. Available: http://hadoop.apache.org/docs/stable1/hdfs_design.html. [Accessed 19 February 2014].

- “HDFS High Availability,” [Online]. Available: [http://hadoop.apache.org/docs/ c urrent/hadoop-yarn/hadoop-yarn-site/](http://hadoop.apache.org/docs/current/hadoop-yarn/hadoop-yarn-site/)
- “Mahout,” [Online]. Available: <http://mahout.apache.org/users/basics/algorithms.html>. [Accessed 19 February 2014].
- “Netcraft,” [Online]. Available: <http://news.netcraft.com/archives/2014/02/03/february-2014-web-server-survey.html>. [Accessed 21 February 2014].
- “Number of Column Families,” [Online]. Available: [http://hbase.apache.org/book/ number.of.cfs.html](http://hbase.apache.org/book/number.of.cfs.html).
- “Pig,” [Online]. Available: <http://pig.apache.org/>. [Accessed 11 Feb 2014].
- “pig.apache.org,” [Online]. Available: <http://pig.apache.org/>.
- “Piggybank,” [Online]. Available: <https://cwiki.apache.org/confluence/display/PIG/PiggyBank>. [Accessed 28 February 2014].
- “Pivotal HD,” [Online]. Available: <http://www.gopivotal.com/big-data/pivotal-hd>. [Accessed 8 May 2014].
- “PoweredByYarn,” [Online]. Available: [http://wiki.apache.org/hadoop/ PoweredByYarn](http://wiki.apache.org/hadoop/PoweredByYarn). [Accessed 20 February 2014].
- “wiki.apache.org/hadoop,” [Online]. Available: [http://wiki.apache.org/hadoop/ NameNode](http://wiki.apache.org/hadoop/NameNode). [Accessed 11 February 2014].

“Zookeeper,” [Online]. Available:
<http://hbase.apache.org/book/zookeeper.html>.
[Accessed 21 February 2014].

“Zookeeper,” [Online]. Available:
<http://zookeeper.apache.org/>. [Accessed 11 Feb 2014].

Apache, “Apache Hadoop,” [Online]. Available:
<http://hadoop.apache.org/>. [Accessed 8 May 2014].

Apache, “Hadoop Streaming,” [Online]. Available:
<https://wiki.apache.org/hadoop/HadoopStreaming>. [Accessed 8 May 2014].

B. Minto, *The Minto Pyramid Principle: Logic in Writing, Thinking, and Problem Solving*, Prentice Hall, 2010.

D. Cutting, “Free Search: Rambilings About Lucene, Nutch, Hadoop and Other Stuff,” [Online]. Available:
<http://cutting.wordpress.com>. [Accessed 11 February 2014].

D. Davidian, “IBM.com,” 14 February 2011. [Online]. Available: https://www-304.ibm.com/connections/blogs/davidian/tags/hadoop?lang=en_us. Accessed 11 February 2014].

D. Gottfrid, “Self-Service, Prorated Supercomputing Fun!,” 01 November 2007. [Online]. Available:
<http://open.blogs.nytimes.com/2007/11/01/self-service-prorated-super-computing-fun/>. [Accessed 11 February 2014].

- E. Baldeschwieler, "http://www.slideshare.net," [Online]. Available: <http://www.slideshare.net/ydn/hadoop-yahoo-internet-scale-data-processing>. [Accessed 11 February 2014].
- E. Dumbill, "The Data Lake Dream," *Forbes*, 14 January 2014. [Online]. Available: <http://www.forbes.com/sites/edddumbill/2014/01/14/the-data-lake-dream/>. [Accessed 4 June 2014].
- Eclipse. [Online]. Available: <https://www.eclipse.org/downloads/>. [Accessed 27 February 2014].
- F. Chang, J. Dean, S. Ghemawat, W.C. Hsieh, D.A. Wallach, M. Burrows, T. Chandra, A. Fikes, and R.E. Gruber, "Bigtable: A Distributed Storage System for Structured Data," [Online]. Available: <http://research.google.com/archive/bigtable.html>. [Accessed 11 February 2014].
- G. Reynolds, *Presentation Zen: Simple Ideas on Presentation Design and Delivery*, Berkeley: New Riders, 2011.
- G. Zelazny, *Say It with Charts: The Executive's Guide to Visual Communication*, McGraw-Hill, 2001.
- HDFSHighAvailabilityWithNFS.html. [Accessed 8 May 2014].
- IBM, "IBM.com," [Online]. Available: http://www-03.ibm.com/innovation/us/watson/watson_in_healthcare.shtml. [Accessed 11 February 2014].
- J. Cohen, B. Dolan, M. Dunlap, J. Hellerstein, and C. Welton, "MAD Skills: New Analysis Practices for Big Data," in

Proceedings of the VLDB Endowment Volume 2 Issue 2, August 2009.

- J. Dean and S. Ghemawat, "MapReduce: Simplified Data Processing on Large Clusters," [Online]. Available: <http://research.google.com/archive/mapreduce.html>. [Accessed 11 February 2014].
- K. Muthukkaruppan, "The Underlying Technology of Messages," 15 November 2010. [Online]. Available: <http://www.facebook.com/notes/facebook-engineering/the-underlying-technology-of-messages/454991608919>. [Accessed 11 February 2014].
- LinkedIn, "Hadoop," [Online]. Available: <http://data.linkedin.com/projects/hadoop>. [Accessed 11 February 2014].
- LinkedIn, "LinkedIn," [Online]. Available: <http://www.linkedin.com/about-us>. [Accessed 11 February 2014].
- MADlib, "MADlib Modules" [Online]. Available: <http://doc.madlib.net/latest/modules.html>. [Accessed 10 April 2014]
- MADlib, "MADlib" [Online]. Available: <http://madlib.net/download/>. [Accessed 10 April 2014].
- N. Spiegelberg. [Online]. Available: <http://www.slideshare.net/brizzzdotcom/facebook-messages-hbase>. [Accessed 11 February 2014].
- N. Yau, "flowingdata.com" [Online]. Available: <http://flowingdata.com>.

- N. Yau, *Visualize This*, Indianapolis: Wiley, 2011.
- P. J. Sadalage and M. Fowler, *NoSQL Distilled: A Brief Guide to the Emerging World of Polyglot*, Upper Saddle River, NJ: Addison Wesley, 2013.
- PostgreSQL.org, "Window Functions" [Online]. Available: <http://www.postgresql.org/docs/9.3/static/functions-window.html>. [Accessed 10 April 2014].
- S. Few, *Now You See It: Simple Visualization Techniques for Quantitative Analysis*, Analytics Press, 2009.
- S. Ghemawat, H. Gobioff, and S.-T. Leung, "The Google File System," [Online]. Available: <http://static.googleusercontent.com/media/research.google.com/en/us/archive/gfs-sosp2003.pdf>. [Accessed 11 February 2014].
- S. Singh, "http://developer.yahoo.com/," [Online]. Available: <http://developer.yahoo.com/blogs/hadoop/apache-hbase-yahoo-multi-tenancy-helm-again-171710422.html>. [Accessed 11 February 2014].
- Wikipedia, "IBM Watson," [Online]. Available: http://en.wikipedia.org/wiki/IBM_Watson. [Accessed 11 February 2014].