

Big Data & Predictive Analytics: The Next Big Leap

Rich Niemiec, Rolta 2013



Innovative Technology for Insightful Impact

Rich's Overview



ORACLE
ACE Director



- Advisor to Rolta International Board
- Former President of TUSC
 - Inc. 500 Company (Fastest Growing 500 Private Companies)
 - 10 Offices in the United States (U.S.); Based in Chicago
 - Oracle Advantage Partner in Tech & Applications
- Former President Rolta TUSC & Former President Rolta EICT International
- Author (3 Oracle Best Sellers – #1 Oracle Tuning Book for over a Decade):
 - Oracle Performing Tips & Techniques (Covers Oracle7 & 8i)
 - Oracle9i Performance Tips & Techniques
 - Oracle Database 10g Performance Tips & Techniques
 - Oracle Database 11g Performance Tips & Techniques
- Former President of the International Oracle Users Group
- Current President of the Midwest Oracle Users Group
- Chicago Entrepreneur Hall of Fame - 1998
- E&Y Entrepreneur of Year & National Hall of Fame - 2001
- IOUG Top Speaker in 1991, 1994, 1997, 2001, 2006, 2007
- MOUG Top Speaker Twelve Times
- National Trio Achiever award - 2006
- Oracle Certified Master & Oracle Ace Director
- Purdue Outstanding Electrical & Computer and Engineer - 2007



Rolta Leaders in Database Technologies and Worldwide Oracle Platinum Partner



- ✓ Oracle Partner-of-the-Year: multiple times
- ✓ Oracle user group leadership
 - Past president of International Oracle User Group
 - Member of Applications & Technology Advisory Councils
 - Current president of Midwest Oracle User Group
 - Service Oriented Architecture (SOA)
 - Fusion
 - Oracle 11g
- ✓ Oracle Magazine Consultants-of-the-Year

Industry Recognitions

- ✓ Nine Times Oracle Titan Award Winners
- ✓ 6 “Oracle Masters” on staff



***One of first few companies worldwide
with highest level of partner certification for Apps & Technology***

Rolta – *Your Partner in Oracle Excellence*



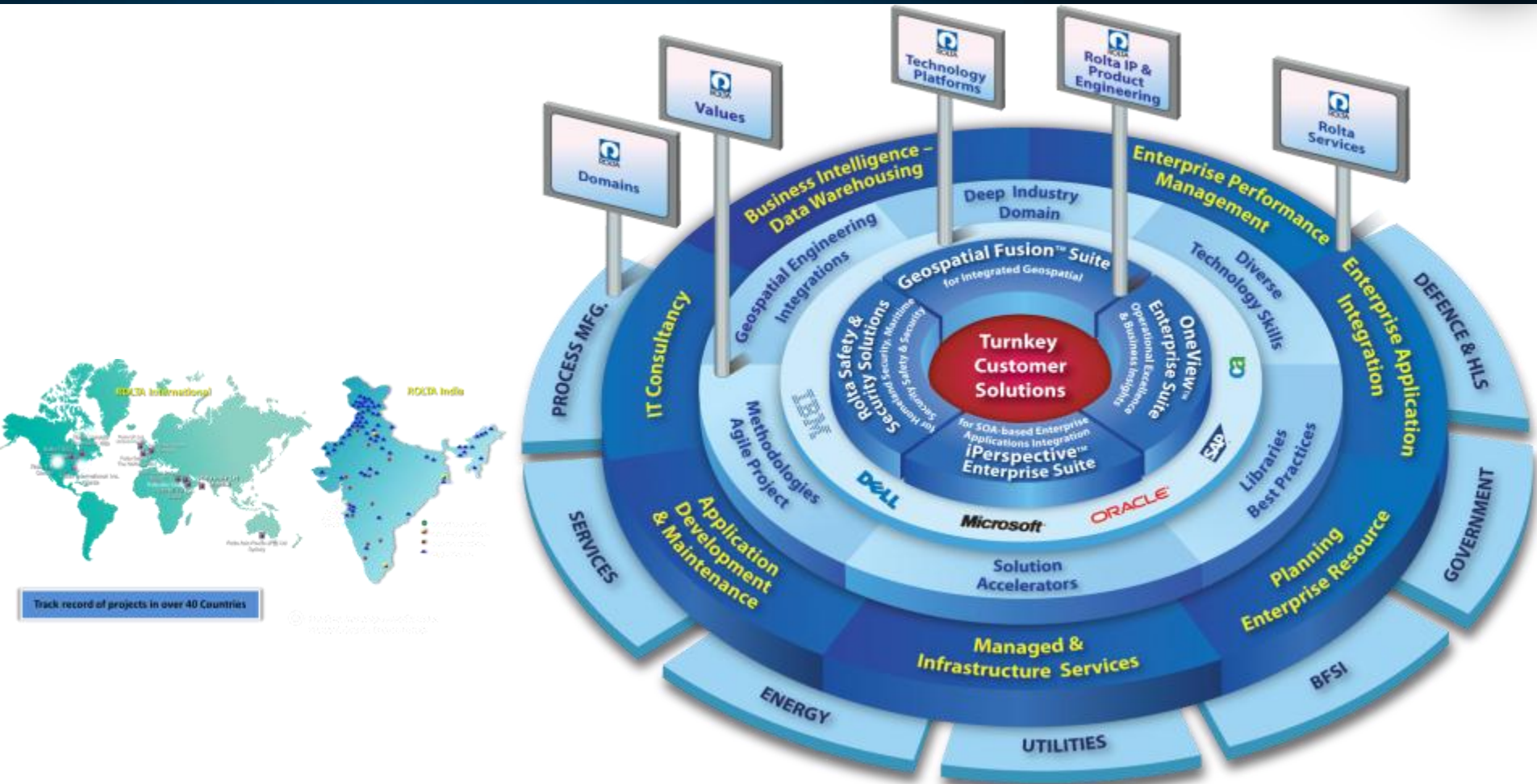
2012 Oracle Excellence Award (BI Application) (9 Partner of the Year / Titans / Excellence Awards)



Prior Years Winner 2002, 2004*, 2007*, 2008, 2010, 2011

***Won 2 Awards**

How to Make a Difference in the World!



- ❑ Oracle Trends
- ❑ Current Scenario
- ❑ Technology Evolution
- ❑ The way forward

Oracle Trends

Know the Oracle



ANNOUNCING ORACLE DATABASE 12^c

A Multitenant Database for the Cloud



Cloning Databases for Test, Development

Fast, flexible copy and snapshot of pluggable databases



ANNOUNCING EXADATA X3

Database In-Memory Machine

The Hardware Foundation of the Oracle Cloud



Exadata X-3: In-Memory Database

4 T DRAM / 22 T Flash Cache



Exadata X3 Database In-Memory Machine



- X3 mass memory hierarchy delivers **extreme performance**
 - Automatically moves all active data from disk to memory
- DRAM memory expanded to 2 or 4 TB for hottest data
 - Up to **40 TB** of compressed user data
- Flash memory expanded **4X** to **22 TB** per rack
 - Up to **220 TB** of compressed user data – **ALL** active data
 - 1.5 Million SQL random read I/Os per second for OLTP
 - Comparable to 15,000 disk drives in 150 array frames
 - 100 GB/sec SQL data scan rate for reporting and warehouses
 - Comparable to 1,000 disk drives in 10 array frames

ORACLE

Oracle Firsts – *Innovation!*



1979 First commercial SQL relational database management system

1983 First 32-bit mode RDBMS

1984 First database with read consistency

1987 First client-server database

1994 First commercial and multilevel secure database evaluations

1995 First 64-bit mode RDBMS

1996 First to break the 30,000 TPC-C barrier

1997 First Web database

1998 First Database - Native **Java** Support; Breaks 100,000 TPC-C

1998 First Commercial RDBMS ported to **Linux**

2000 First database with **XML**

2001 First middle-tier database cache

2001 First RDBMS with **Real Application Clusters**

2004 First **True Grid Database**

2005 First **FREE Oracle Database** (10g Express Edition)

2006 First **Oracle Support for LINUX Offering**

2007 **Oracle 11g Released!**

2008 Oracle Exadata Server Announced (Oracle buys BEA)

2009 Oracle buys Sun – Java; MySQL; Solaris; Hardware; OpenOffice

2010 Oracle announces MySQL Cluster 7.1, Exadata, Exalogic

2011 Oracle X2-2, ODA, Exalytics, SuperCluster, Big Data, Cloud, Social Network

2012 Oracle X3-2, Oracle 12c OEM, Pluggable Databases & X3-8 announced

2013 Oracle12c Released! Oracle X3-8 Exadata, Acquisitions (Acme Packet)!

❑ Recession Over, but Outlook is still challenging?

- Great Recession of 2008-09, biggest blow to economy since 1930
- For some, 2012 was worse than 2008-09: Europe's Debt crisis, General slowing of Global Economy, Domestic Political Troubles in various nations

❑ Many country ratings are reduced, *Most* banking unscathed, Reforms will drive some growth, Big Data will drive more!

❑ World is awash with excess capacities that will find its way into demand centric markets

Worldwide capacity utilization

- Oil Refineries – 84%
- Steel – 76%
- Cement – 75%
- Aluminium – 73%
- Poly-crystalline silicon – 20%

China's production capacity (as % of global capacity)

- Oil Refining ~20% *(excluding OPEC)*
- Steel ~ 50%
- Cement ~ 50%
- Aluminium ~ 40%

Data collated from various Industry Sources

Some of the typical Organizational Challenges and need for Analytics



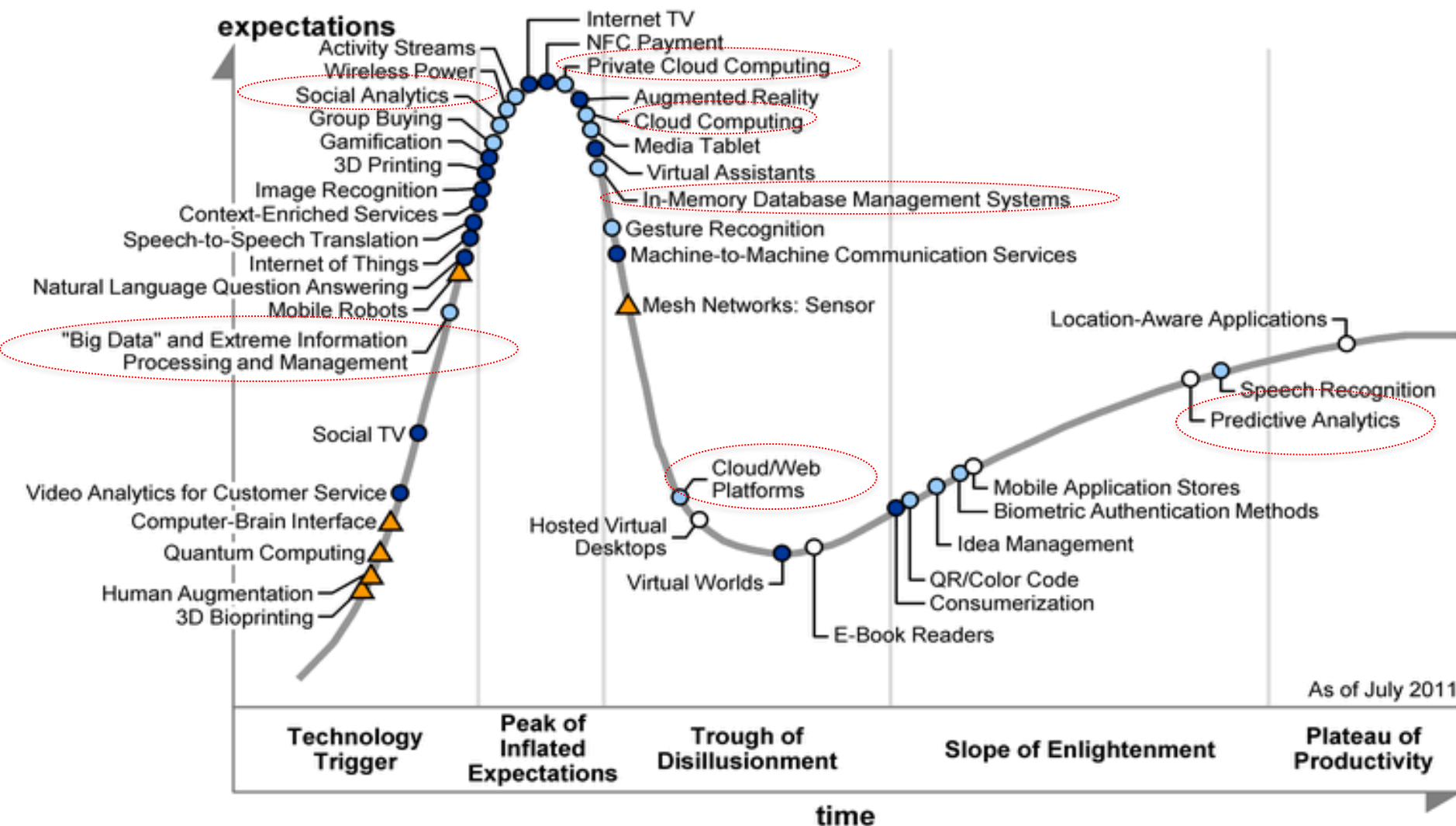
What we
know,
we know

What we
know, we
don't know

What we don't know,
we don't know

- ❑ Cloud Computing, Mobile Computing, Social Media and Big Data Analytics are driving the New Computing Paradigm.
- ❑ This paradigm in-turn sparks-off Business Transformations to improve Efficiency, Compliance with Regulation and overall Business Sustainability based on Customer Centricity.

Technology Trends: Gartner Hype Cycle 2012



Gartner Trends for 2012

Top 10 Strategic Technology Trends for 2012

Human
Experience

1. Media tablets and beyond
2. Mobile-centric applications and interfaces
3. Contextual and social user experience

Business
Experience

4. Internet of things
5. App stores and marketplaces

6. Next-generation analytics

7. Big data

8. In-memory computing

IT Dept.
Experience

9. Extreme low-energy servers

10. Cloud computing

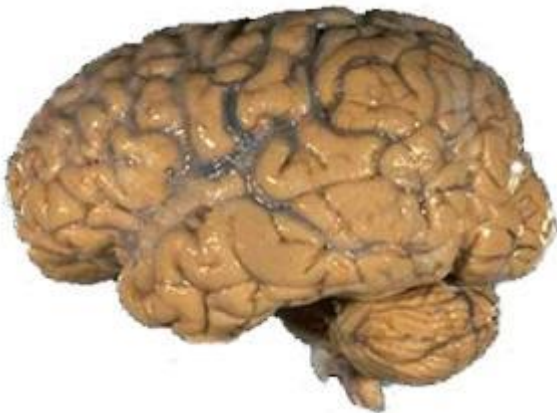
Bigger Data - Data Size Matters...

❖ Worldwide, data is growing rapidly over the years....

- ❑ 2000: 800 Terabytes (10^{12})
- ❑ 2006: 160 Exabytes (10^{18})
- ❑ 2009: 500 Exabytes (just Internet)
- ❑ 2012: 2.7 Zettabytes (10^{21})
- ❑ 2020: 35 Zettabytes ...?

❖ Data generated in ONE day....?

- ❑ Twitter: 7 TB
- ❑ Facebook: > 10 TB



2.8×10^{20} bits of Memory Space – John von Neumann (“Computer and the Brain”, Harvard Lecture Notes, Half Century ago)



Big data: The next frontier for innovation, competition, and productivity

McKinsey Global Institute 2011

Data collated from various online sources

We are drowning in data, but thirsting for Information

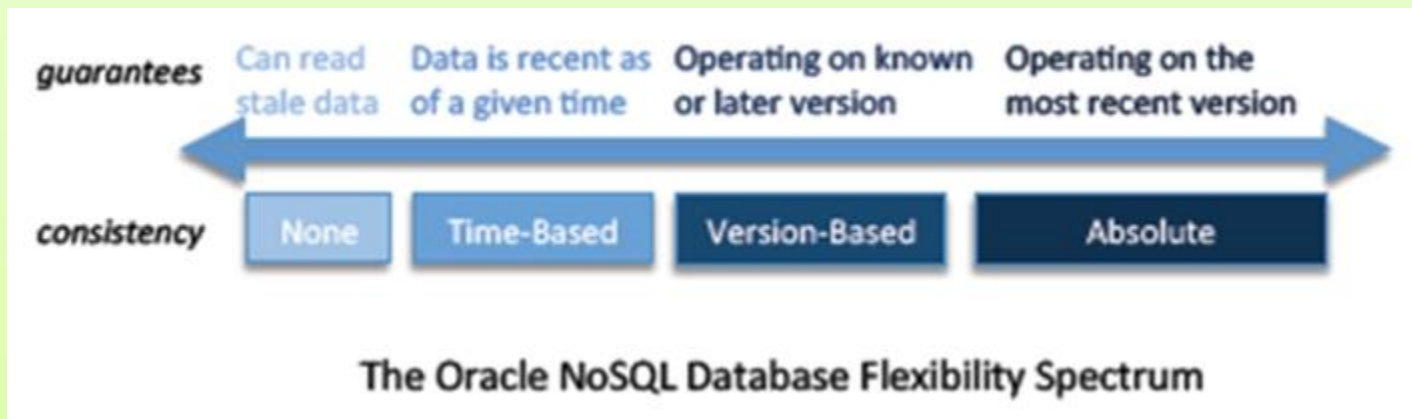
How Much Data ...

- 2004 monthly internet traffic $>1\text{E}$; 2010 it was $21\text{E}/\text{month}$.
- In 2012, **2.5E data created every day** (about $1\text{Z}=1000\text{E}$ per year)
- June 2012 – Facebook has **100P Hadoop cluster**
- Facebook: **500T** processed daily – ($210\text{T}/\text{hr}$ Hive scanned)
- A Single Jet Engine – **$20\text{T}/\text{hour}$** –same rate as Facebook!
- Gmail has **450 million users**
- Wal-Mart – 1 million customer transactions/hour (2.5P DB)
- Large Hadron Collider produced **13P in one year**
- Business data **doubles every 1.2 years**
- 19% of $\$1\text{B}$ companies have **$>1\text{P}$ of data** (31% in 2013)
- 2011 – First **Exabyte tape library** from Oracle
- Decoding Human Genome took 10 yrs; Now takes a week!

- Big Data includes: Social Media, Sensor Data, Biological, Traffic, RFID Data, Environmental, Aerial, Wireless, Security & Video Data, Retail, Medical, Engineering Systems, Search Data, Photographs, Call Records, CRM/ERP data, etc.
- Goal was to **Organize Data without moving it!** – Hadoop HDFS & MapReduce (Cheaper way to access Petabytes)
- **Acquire & Store data** – NoSQL (simple key value storage) – Amazon DynamoDB (hosted), Apache Cassandra, HBase, BigTable, MongoDB, Oracle NoSQL (distributed key value) or just use the original HDFS / GFS & MapReduce (many are **EVENTUALLY** consistent!)
- **Analyze Data** – Google Dremel, Apache Hive Data Warehouse, Oracle Data Warehouse
- 54% of companies doing Big Data say **projects are critical**

NoSQL supports BASE:

- Basically Available
- Soft state
- Eventually consistent



In the Beginning... How did we get here?

- Larry Page & Sergey Brin wrote BigFiles; **GFS** (Google File System) grew out of that & then **MapReduce** which maps problems across cluster a of worker nodes & then collects results & aggregates/reduces result (**used to generate Google's index of WWW**)
- Apache came out with **Hadoop** (*used by Facebook, Yahoo, Amazon EC2 & S3*) which was an Open Source version with **HDFS & MapReduce** – Batch Processing Jobs going after distributed data & processing it near the data (same node) – not super fast (**seconds vs. ms**) & not good for interactive/analytic
- **Google then came out with BigTable** (compressed, high performance data storage) used by Google Maps, Google Reader, Google Earth, YouTube, and Gmail
- Apache adds **NoSQL DB's: Cassandra & HBase**
- The NoSQL onslaught of systems started (over 100 of them) including **Oracle's NoSQL (BerkeleyDB)**.

❑ Hadoop goes Enterprise.

- ❑ Microsoft joins the party (partnership with Yahoo! Spin-off Hortonworks → Hadoop implementation for Windows Server & Azure with connectors to MSSQL)

❑ NoSQL based solutions

- ❑ Security Issues hamper NoSQL
- ❑ Oracle gets in the NoSQL game in a BIGGER way (Big Data Appliance)

“As customers look to manage the huge explosion in data from new and evolving sources, such as the Web, sensors, social networks and mobile applications, Oracle is helping them unlock the value of this data by providing a highly available, reliable and scalable NoSQL database environment.”

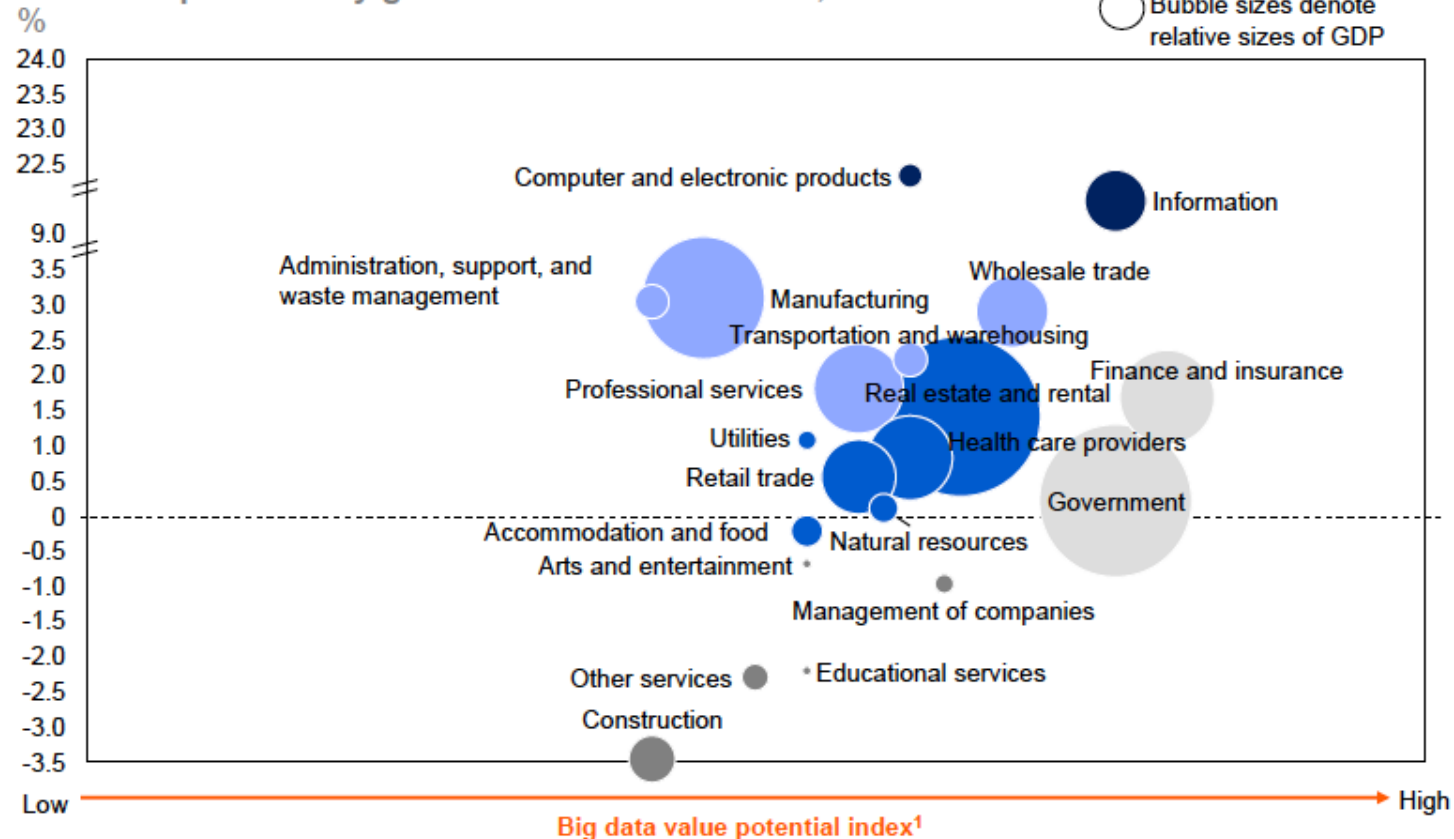
– Oracle SVP, Andrew Mendelsohn

❑ Integration of In-Memory Data Grids and NoSQL, leveraging success stories of Facebook & Twitter

Why Is Big Data Important?

Some sectors are positioned for greater gains from the use of big data

Historical productivity growth in the United States, 2000–08



SOURCE: US Bureau of Labor Statistics; McKinsey Global Institute analysis

Characteristics of Big Data

Volume

Big data comes in one size: large. Enterprises are awash with data, easily amassing terabytes and even petabytes of information.

TB, Records, Transactions, Tables, Files

Velocity

Often time-sensitive, big data must be used as it is streaming in to the enterprise in order to maximize its value to the business.

Batch, Near time, Real time, Streams

Value

Business value of Big Data

Variety

Big data extends beyond structured data, including semi-structured and unstructured data of all varieties: text, audio, video, click streams, log files and more.

Structured, Unstructured, Semistructured



Finance



Telecom



Retail



Life Sciences



Government



Media

Big Data Themes

- HW & SW technologies for large data volumes
- Focus on Web 2.0 technologies
- Database Scale-out
- Relational & Distributed Data Analytics
- Distributed File Systems
- Real Time Analytics

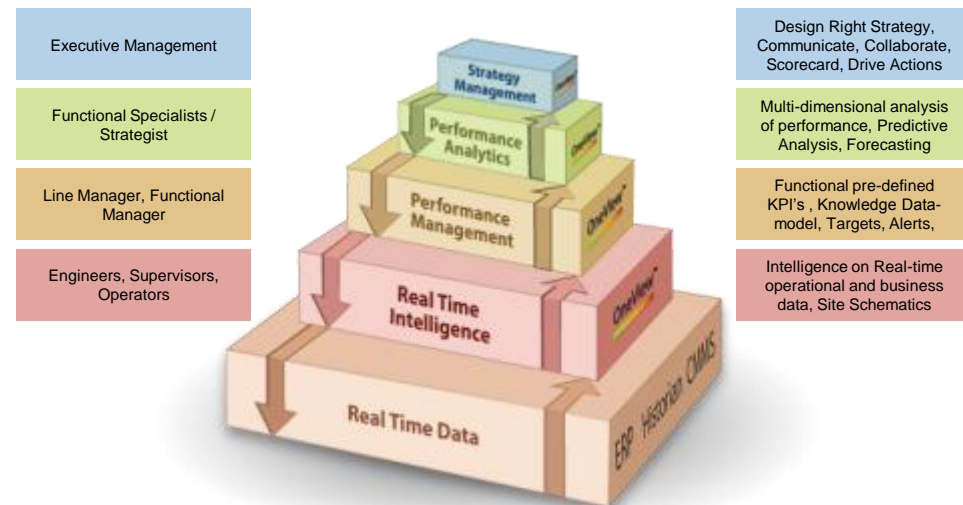
Big Data Domains

- Digital Marketing Optimization
- Data Exploration & Discovery
- Fraud Detection & Prevention
- Social Network & Relationship Analysis
- Machine-generated Data Analytics
- Data Retention

All Data is not similar!

Data Realms	Security	Storage & Retrieval	Modeling	Processing & Integration	Consumption
Master data Transactions Analytical data Metadata	Database, app, & user access	RDBMS / SQL	Pre-defined relational or dimensional modeling	ETL/ELT, CDC, Replication, Message	BI & Statistical Tools, Operational Applications
Reference data	Platform security	XML / xQuery	Flexible & Extensible	ETL/ELT, Message	System-based data consumption
Documents and Content	File system based	File System / Search	Free Form	OS-level file movement	Content Mgmt
Big Data - Weblogs - Sensors - Social Media	File system & database	Distributed FS / noSQL	Flexible (Key Value)	Hadoop, MapReduce, ETL/ELT, Message	BI & Statistical Tools

Analytics



Pre-packaged Analytics are also available...

Oracle BI Applications

Comprehensive, Prebuilt, Best Practice Analytics



Sales	Service & Contact Center	Marketing	Procurement & Spend	Supply Chain & Order Management	Financials	Human Resources
Pipeline Analysis	Service Effectiveness	Campaign Effectiveness	Direct / Indirect Spend	Revenue and Backlog	General Ledger	Employee Productivity
Forecast Accuracy	Customer Satisfaction	Customer Insight	Buyer Productivity	Inventory	Accounts Receivable	Compensation
Sales Team Effectiveness	Resolution Rates	Product Propensity	Off Contract Purchases	Fulfillment Status	Accounts Payable	Compliance Reporting
Up-sell / Cross-sell	Service Rep Efficiency	Loyalty & Attrition	Supplier Performance	Customer Status	Cash Flow	Workforce Profile
Cycle Times	Service Cost	Market Basket Analysis	Purchase Cycle Time	Order Cycle Time	Profitability	Retention Analysis
Lead Conversion	Churn & Service Trends	Campaign ROI	Employee Expenses	BOM Analysis	Expense Management	Return-on Human Capital

Source adapters:

ORACLE **SIEBEL** PeopleSoft.

J D EDWARDS

SAP and Other Operational & Analytic Sources

Oracle BI Suite Enterprise Edition Plus

Oracle Database is loaded with Analytics !!



Analytical Feature	Description
Data Mining	Oracle Data Mining implements complex algorithms to discover patterns, predict probable outcomes, identify key predictors, etc.
Complex data transformations	ETL capabilities and SQL expressions or DBMS_DATA_MINING_TRANSFORM package. for missing values, outlier treatments, binning and normalization.
Statistical functions	SQL statistical functions: hypothesis testing (t-test, F-test), pearson correlation, cross-tab/descriptive statistics (median, mode, etc). DBMS_STAT_FUNCS package adds distribution fitting procedures.
Window / Analytic SQL functions	Computing cumulative, moving, and centered aggregates.
Frequent Itemsets	DBMS_FREQUENT_ITEMSET used as a building block for the Association algorithm used by Oracle Data Mining.
Image feature extraction	Oracle Intermedia supports extraction of color histogram, texture, and positional color.
Linear algebra	UTL_NLA package exposes a subset of the popular BLAS and LAPACK libraries for operations on vectors and matrices.
OLAP	Multidimensional analysis beyond drill-downs and roll-ups, Oracle OLAP also supports time-series analysis, modeling, and forecasting
Spatial analytics	Oracle Spatial's analysis and mining capabilities include binning, pattern detection, spatial correlation, colocation mining, and spatial clustering, topology & NW data model analytics - shortest path, minimum cost spanning tree, nearest-neighbors analysis, traveling salesman problem, etc
Text Mining	Std SQL to index, search, analyze text / documents stored in DB, files, and web with automatic classification and clustering

Predictive in Nature?

Hindsight

What is happening?

- ☐ Historic orientation
- ☐ Typical MIS Reporting or BI
- ☐ Oracle Reports, Hyperion, IBM Cognos, SAP BO, etc



Insight

Why is it happening?

- ☐ Business / Behaviour Analysis, Trends
- ☐ What is currently happening / Why?

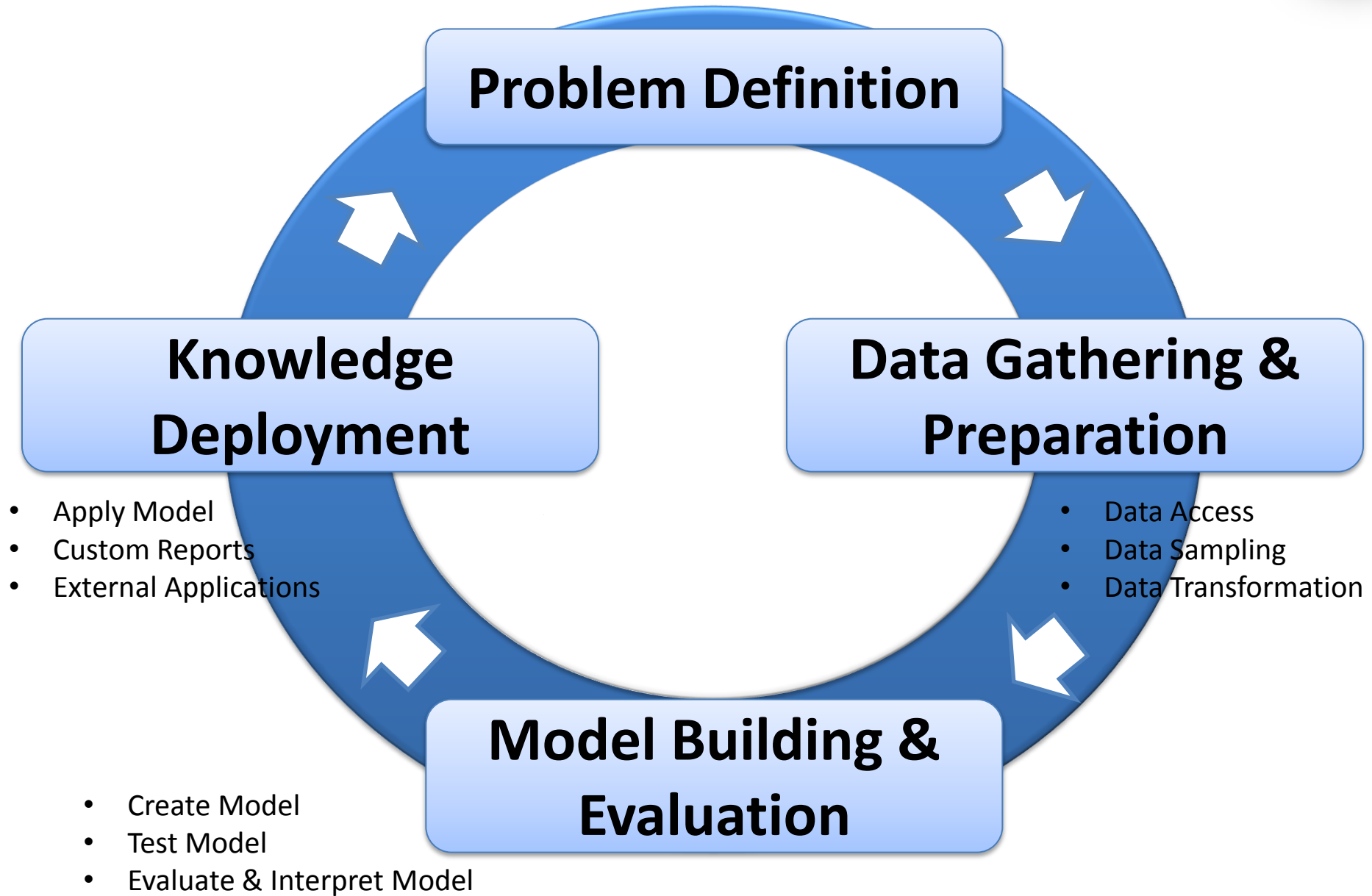


Foresight

What will / should happen?

- ☐ Forecasting
- ☐ Optimization
- ☐ Past behaviour to predict future outcomes

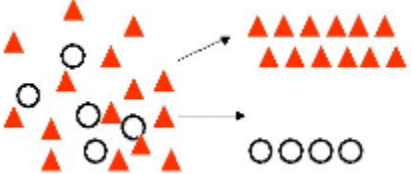

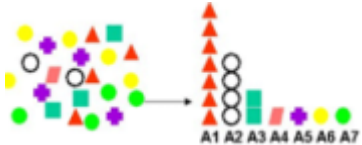

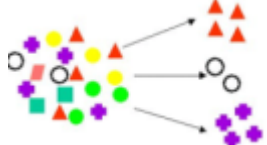
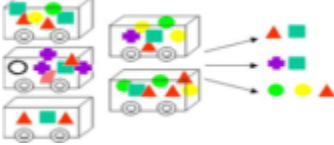
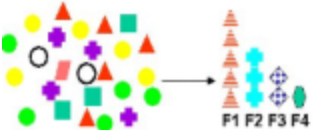




Closely Inter-related Concepts

- ❑ **Artificial Intelligence (AI)** refers to systems that exhibit autonomous intelligence or behaviour
- ❑ **Machine Learning (ML)** are AI techniques that enable devices to learn from their performance and modify their functioning
- ❑ **Data Mining (DM)** applies ML concepts to data
 - Each **Data Mining Function (DMF)** specifies a class of problems that can be modelled and solved.
 - Data mining functions fall generally into two categories: **Supervised and Unsupervised**
 - Each **Data Mining Algorithm (DMA)** is a Mathematical Procedure for solving a class of problems

Popular DMFs & DMAs - supported by Oracle

Functions		Some Examples	Algorithms
Classification		<ul style="list-style-type: none"> High, Medium or Low Value customer Likely Buy / No-Buy 	<ul style="list-style-type: none"> <input type="checkbox"/> Logical Regression <input type="checkbox"/> Naïve Bayes <input type="checkbox"/> Support Vector M/c <input type="checkbox"/> Decision Tree
Regression		<ul style="list-style-type: none"> Customer Lifetime Value Process Yield Rates 	<ul style="list-style-type: none"> <input type="checkbox"/> Multiple Regression <input type="checkbox"/> Support Vector M/c
Attribute Importance		<ul style="list-style-type: none"> Medical diagnosis factors Buyer priorities 	<ul style="list-style-type: none"> <input type="checkbox"/> Minimum Description Length
Anomaly Detection		<ul style="list-style-type: none"> Insurance Frauds Tax compliance 	<ul style="list-style-type: none"> <input type="checkbox"/> One-class Support Vector Machine
Clustering		<ul style="list-style-type: none"> Customer segmentation Life Sciences Discoveries 	<ul style="list-style-type: none"> <input type="checkbox"/> Enhanced K-Means <input type="checkbox"/> Orthogonal Partitioning Clustering
Association		<ul style="list-style-type: none"> Product Bundling Defect Analysis 	<ul style="list-style-type: none"> <input type="checkbox"/> Apriori
Feature Extraction		<ul style="list-style-type: none"> Pattern Recognition Data Projection 	<ul style="list-style-type: none"> <input type="checkbox"/> Non-Negative Matrix Factorization

Oracle Predictive Analytics

EXPLAIN

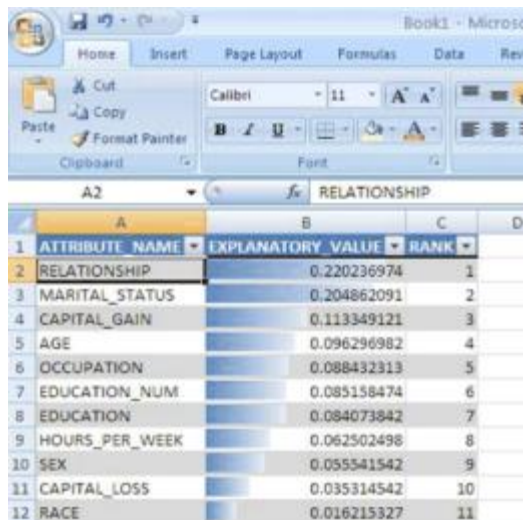
Explains how the individual attributes affect the variation of values in a target column

PREDICT

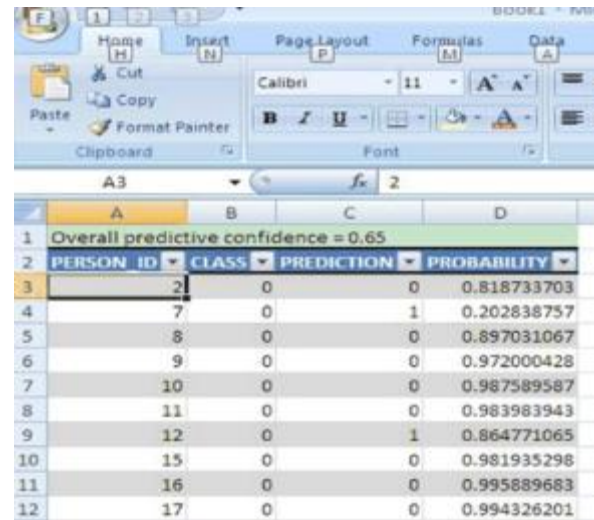
For each case, predicts the values in a target column

PROFILE

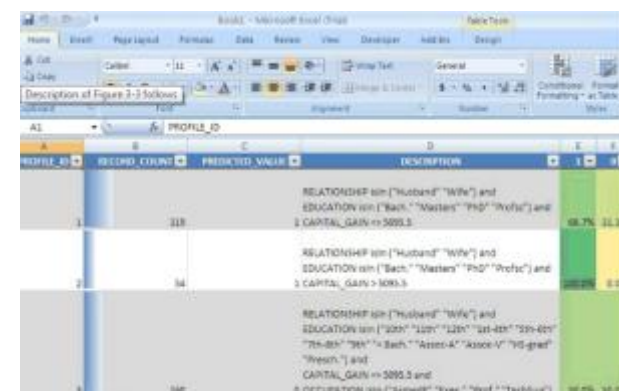
Creates a set of rules for cases that imply the same target value



1	ATTRIBUTE_NAME	EXPLANATORY_VALUE	RANK
2	RELATIONSHIP	0.220236974	1
3	MARITAL_STATUS	0.204862091	2
4	CAPITAL_GAIN	0.113349121	3
5	AGE	0.096296982	4
6	OCCUPATION	0.088432313	5
7	EDUCATION_NUM	0.085158474	6
8	EDUCATION	0.084073842	7
9	HOURS_PER_WEEK	0.062502498	8
10	SEX	0.055541542	9
11	CAPITAL_LOSS	0.035314542	10
12	RACE	0.016215327	11



1	PERSON_ID	CLASS	PREDICTION	PROBABILITY
2	Overall predictive confidence = 0.65			
3	2	0	0	0.818733703
4	7	0	1	0.202838757
5	8	0	0	0.897031067
6	9	0	0	0.972000428
7	10	0	0	0.987589587
8	11	0	0	0.983983943
9	12	0	1	0.864771065
10	15	0	0	0.981935298
11	16	0	0	0.995889683
12	17	0	0	0.994326201

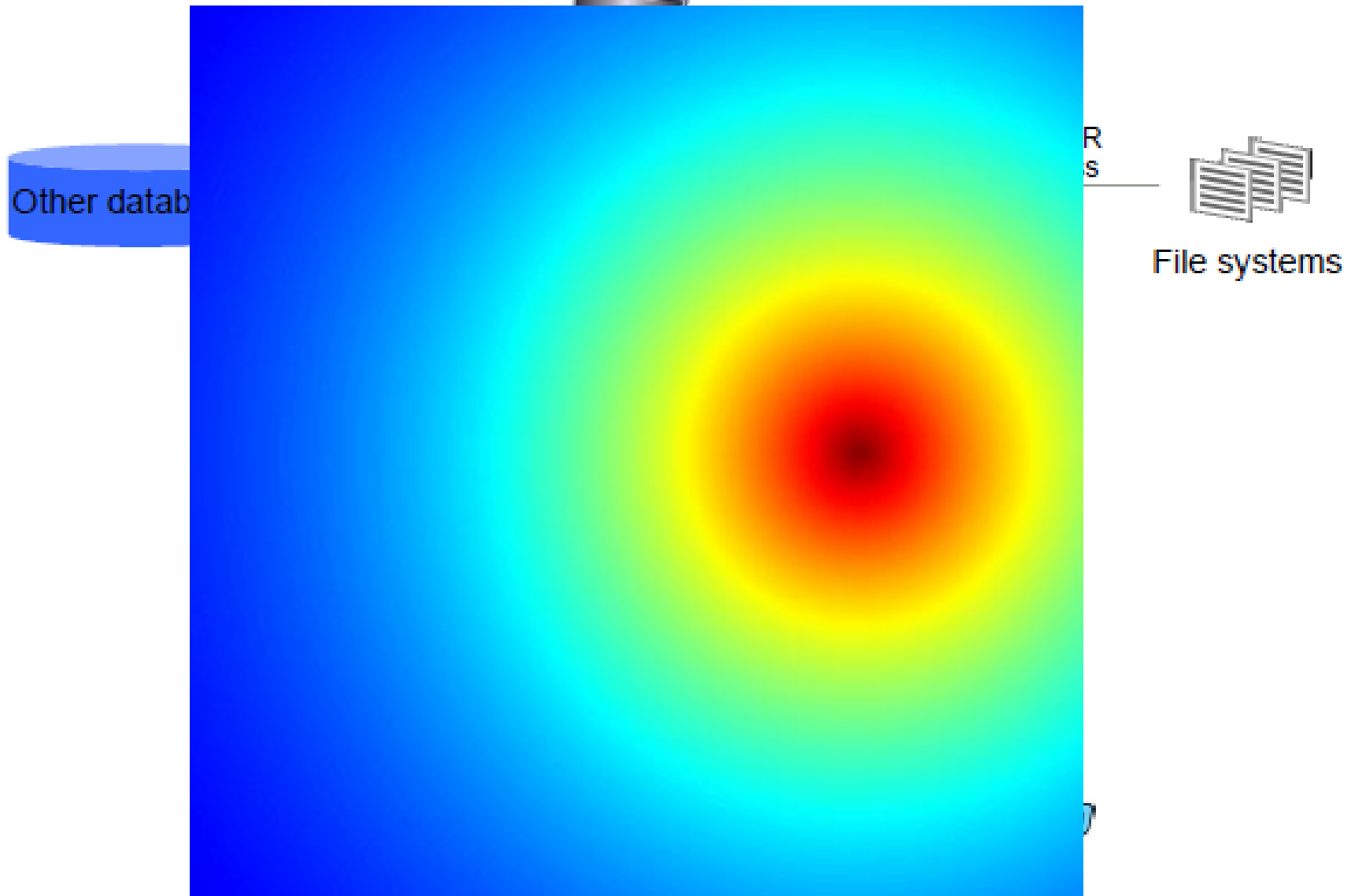


1	PROFILE_ID	RECORD_COUNT	PREDICTED_VALUE	DESCRIPTION
2	118			RELATIONSHIP is ("Husband" "Wife") and EDUCATION is ("Bach." "Masters" "PhD" "Prof") and CAPITAL_GAIN is >3000.5
3	114			RELATIONSHIP is ("Husband" "Wife") and EDUCATION is ("Bach." "Masters" "PhD" "Prof") and CAPITAL_GAIN is >3000.5
4	116			RELATIONSHIP is ("Husband" "Wife") and EDUCATION is ("Bach." "Masters" "PhD" "Prof") and CAPITAL_GAIN is >3000.5 and OCCUPATION is ("Manager" "Exec." "Prof." "Techsup")

Oracle Add-in Spread sheets for Predictive Analytics

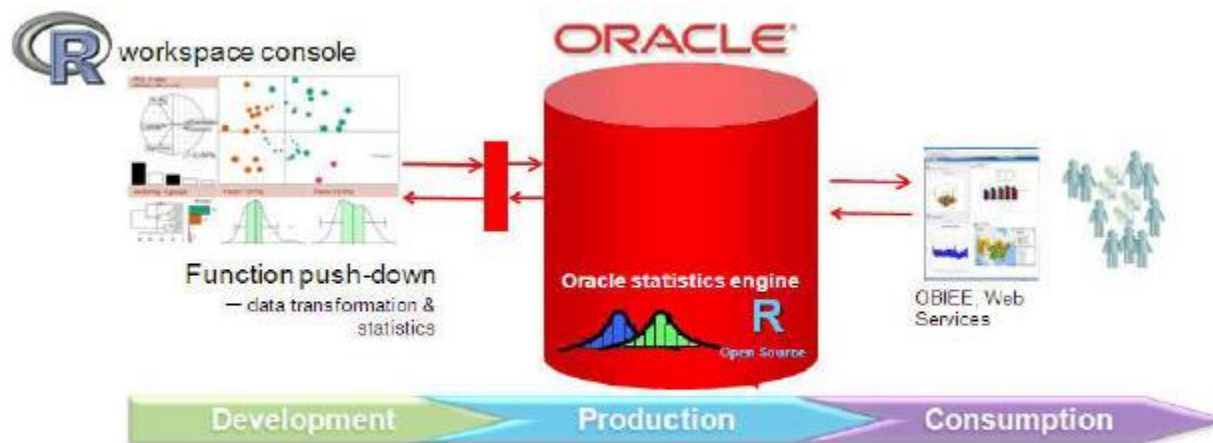
Oracle Data Mining implements Predictive Analytics in PL/SQL & Java APIs

Oracle's "Open" Secret Sauce for Predictive Analytics of Big Data



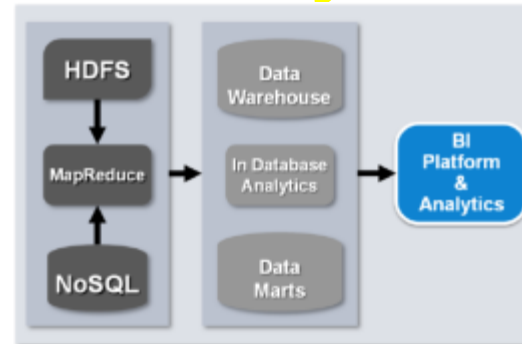
How does it work?

- ❑ Coupled with the power of SQL, Oracle Data Mining eliminates data movement and duplication, maintains security and minimizes latency time from raw data to valuable information
- ❑ Oracle Advanced Analytics Option, makes Oracle tables/views accessible to the R environment as if they are native R objects and transparently leverages the massive scalability of the database for big data analysis problems.



Business Intelligence and Advanced Analytics integrated

Oracle Technologies for – Big Data Predictive Analytics



Oracle Big Data Appliance

Optimized for Hadoop, R, and NoSQL Processing

Oracle Big Data Connectors

Oracle Exadata
“System of Record”
Optimized for DW/OLTP

Oracle Exalytics
Optimized for Analytics Workload



Infiniband

Stream

Acquire

Organize

Analyze and Visualize



My Oracle Big Data Benefits



- **It's actually done and complete** unlike others
- **Full Hadoop integration and loader**
- **Exadata and Exalytics** BI integration & solution
- **Big Data hardware** which includes Hadoop HDFS, MapReduce, R programming language (statistics and regressions...etc.), Oracle NoSQL, ACID compliant, Simple key-value pair data model (hashes keys over many servers - major/minor keys & byte arrays)
- **Based on Oracle's BerkeleyDB** (commercial 8 years!) which integrates with HDFS (Hadoop File System) using external tables if you want,
- Oracle Loader for Hadoop (OLH) takes the analyzed data from MapReduce & **puts into 11g Database as last step** (easier to do)
- **Concurrency is flexible** at any level & it's horizontally scalable
- Oracle knows clustering & HA well (**no single point of failure!**)
- Oracle **Admin tools** are great as are Oracle professionals
- **BerkeleyDB is the worlds most widely used DB toolkit** >200M deployed copies

Rolta Leaders in Big Data Analytics and Worldwide Oracle Platinum Partner



- ✓ Oracle Partner-of-the-Year: multiple times
- ✓ Oracle user group leadership
 - Past president of International Oracle User Group
 - Member of Applications & Technology Advisory Councils
 - Current president of Midwest Oracle User Group
 - Service Oriented Architecture (SOA)
 - Fusion
 - Oracle 11g
- ✓ Oracle Magazine Consultants-of-the-Year

Industry Recognitions

- ✓ Nine Times Oracle Titan Award Winners
- ✓ 6 “Oracle Masters” on staff



***One of first few companies worldwide
with highest level of partner certification for Apps & Technology***

Final Thoughts... Catching your Wave!



*“Things may come to those who wait, but
only the things left by those who hustle.”*

— Abraham Lincoln

Rolta– *Your Partner in Success....*



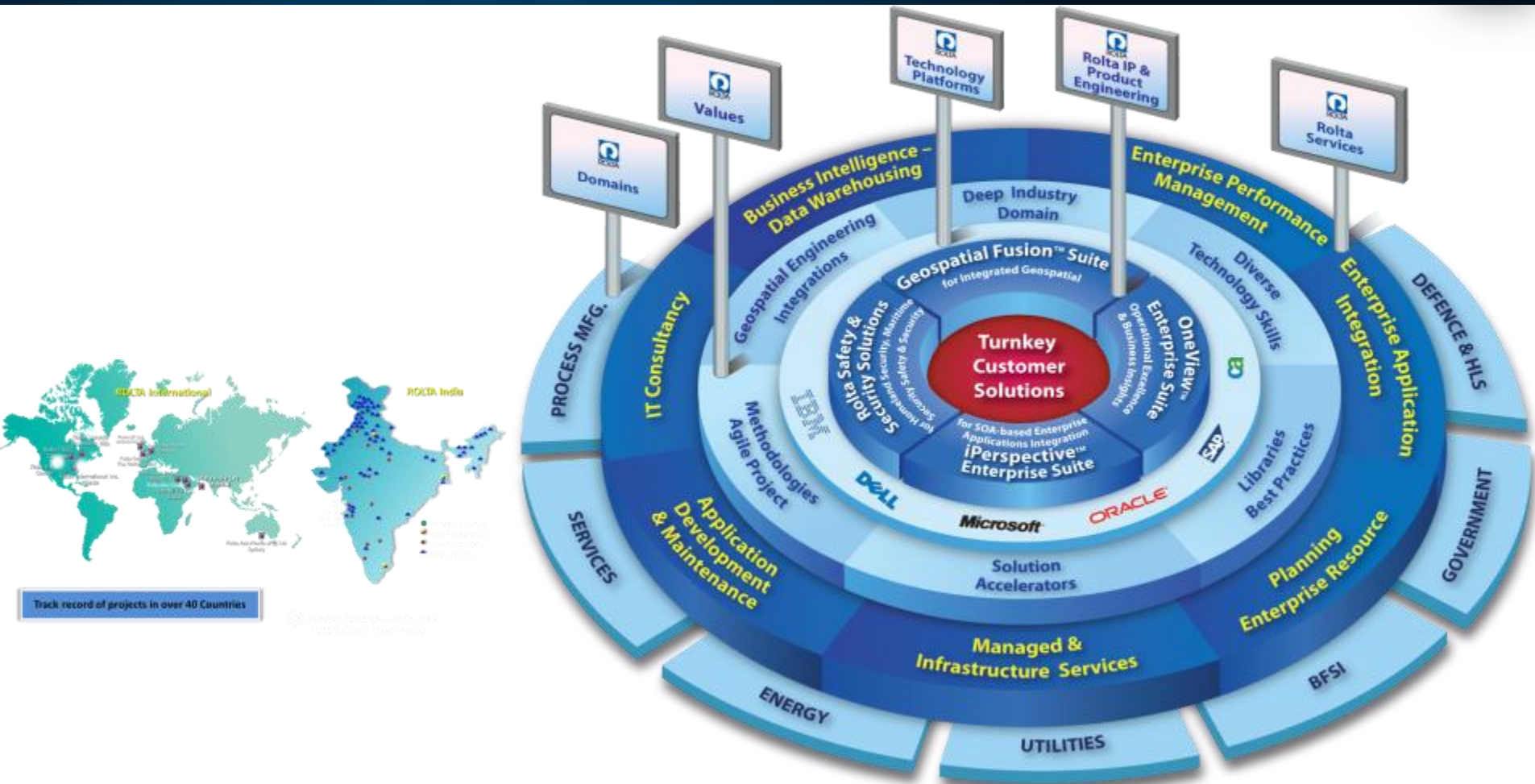
**2012 Oracle Partner of the Year
(9 Titans/Excellence Awards)**



Prior Years: 2002, 2004*, 2007*, 2008, 2010, 2011

***Won 2 Awards**

Thank You & Make a Difference in the World!



Copyright



- Neither Rolta nor the author guarantee this document to be error-free. Please provide comments/questions to rich.niemiec@roltasolutions.com. I am always looking to improve!
- Rich Niemiec/Rolta ©2013. This document cannot be reproduced without expressed written consent from Rich Niemiec or an officer of Rolta, but may be reproduced or copied for presentation/conference use.
- References include Rich Niemiec's Exadata Presentation & Oracle11g Database Performance Tuning Tips & Techniques book, www.oracle.com, en.wikipedia.org, slashgear.com, gifsoup.com, www.amazon.com, Tech Crunch, www.rolta.com, The Matrix movie, Information Week, Gartner, Computerworld, & Oracle OpenWorld

Contact Information

Rich Niemiec: rich.niemiec@roltasolutions.com

www.rolta.com