

Data Management for Data Science

*Corso di laurea magistrale in Data Science
Sapienza Università di Roma
2017/2018*

An Introduction to Big Data

Domenico Lembo

Dipartimento di Ingegneria Informatica Automatica e Gestionale A. Ruberti

Availability of Massive Data

- Digital data are nowadays collected at an **unprecedented scale** and in **very many formats** in a variety of domains (e-commerce, social networks, sensor networks, astronomy, genomics, medical records, etc.)
- This has been made possible by the incredible **growth** of the last years **of the capacity of data storage tools**, and of the **computing power of electronic devices**, and as well as the advent of **mobile and pervasive computing**, cloud computing and cloud storage.

Exploitability of Massive Data

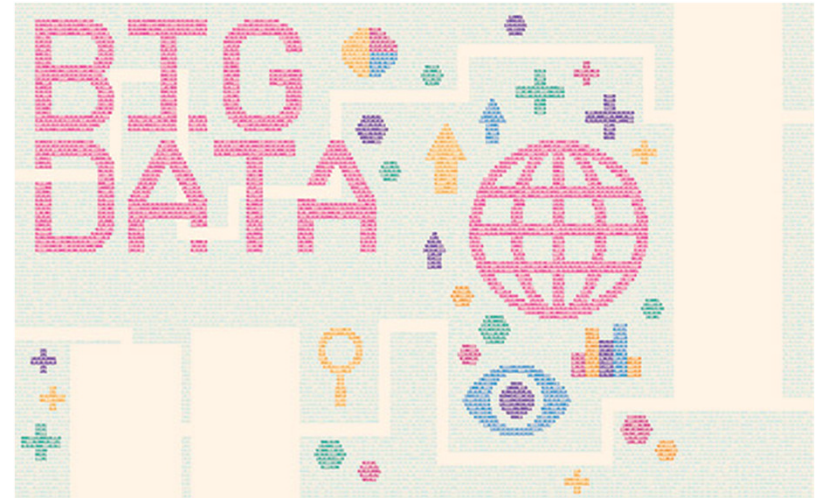
- How to **transform available data into information**, and how to make organizations' business to take advantages of such information are long-standing problems in IT, and in particular in information management and analysis.



- These issues have become more and more challenging and complex in the “**Big Data**” era
- At the same time, facing the challenge may be worthy, since the massive amount of data that is now available may allow for analytical results never achieved before

but be careful!

- “Big data is a vague term for a massive phenomenon that has rapidly become an obsession with entrepreneurs, scientists, governments and the media” (Tim Harford, journalist and economist)*



* <http://www.ft.com/cms/s/2/21a6e7d8-b479-11e3-a09a-00144feabdc0.html#axzz3EvSLWwbu>

Moore's Law for #BigData: The amount of nonsense packed into the term "BigData" doubles approximately every two years (Mike Pluta, Data Architect on Twitter August 2014).

Thinking Big Data*

*"Big Data" has leapt rapidly into one of the most hyped terms in our industry, yet the hype should not blind people to the fact that this is a genuinely important shift about the role of data in the world. The amount, speed, and value of data sources is rapidly increasing. Data management has to change in five broad areas: **extraction** of data from a wider range of sources, changes to the **logistics** of data management with new database and integration approaches, the use of **agile** principles in running analytics projects, an emphasis on techniques for data **interpretation** to separate signal from noise, and the importance of well-designed **visualization** to make that signal more comprehensible. Summing up this means we don't need big analytics projects, instead we want the new data thinking to permeate our regular work."*

Martin Fowler

*<http://martinfowler.com/articles/bigData/>

Thinking Big Data

- Thus, roughly, *Big Data is data that exceeds the processing capacity of conventional database systems*
- But also *Big Data is understood as a capability that allows companies to extract **value** from large volumes of data*
- but, notice, **this does not mean only extremely large, massive databases**
- Besides data dimension, what characterizes Big Data are also the heterogeneity in the way in which information is structured, the dynamicity with which data changes, is the ability of quickly processing it
- This calls for **new computing paradigms or frameworks**, not only advanced data storage mechanisms

The Three Vs

To characterize Big Data, three Vs are used, which are the Vs of

- ***Volume***
- ***Velocity***
- ***Variety***

Volume

- Big data applications are characterized of course by big amounts of data, where big means **extremely large**, e.g., more than a terabyte (TB) or petabyte (PB), or more.
- Some examples:
 - **Walmart**: 1 million transaction per hour (2010)¹
 - **eBay**: data throughput reaches 100 petabytes per day (2013)²
 - **Facebook**: 40 billion photos (2010)¹; 250PB data warehouse with 600TB added to the warehouse every day (2013)³
 - 500 millions of tweet per day (in 2013)
 - And very many other examples, as chatters from social networks, web server logs, traffic flow sensors, satellite imagery, broadcast audio streams, banking transactions, GPS trails, financial market data, biological data, ecc.

¹<http://martinfowler.com/articles/bigData/>

²<http://www.v3.co.uk/v3-uk/news/2302017/ebay-using-big-data-analytics-to-drive-up-price-listings>

³http://www.theregister.co.uk/2013/06/07/hey_presto_facebook_reveals_exabytescale_query_engine/

Volume

- How many data in the world?
 - 800 Terabytes, 2000
 - 160 Exabytes, 2006 (1EB=10¹⁸B)
 - 500 Exabytes, 2009
 - 2.7 Zettabytes, 2012 (1ZB=10²¹B)
 - 35 Zettabytes by 2020



- 90% of world's data generated in the last two years.

Volume

- In a data integration context, the **number of sources providing information can be huge too**, much higher than the number considered in traditional data integration and virtualization systems
- The sheer volume of data is enough to defeat many long-followed approaches to data management
- Traditional centralized database systems cannot handle many of the data volumes, forcing the use of clusters

Velocity

- Data velocity (i.e., **the rate at which data is collected** and made available into an organization) has followed a similar pattern to that of volume
- Many **data sources** accessed by organizations for their business are **extremely dynamic**
- Mobile devices increase the rate of data inflow: data “everywhere”, collected and consumed continuously

Velocity

- **Processing information as soon as it is available**, thus speeding the “feedback loop”, can provide competitive advantages
- As an example, consider online retailers that are able to suggest additional products to a customer at every new information inserted during an on-line purchase
- **Stream processing** is a new challenging computing paradigm, where information is not stored for later batch processing, but is consumed on the fly
- This is particularly useful when **data are too fast to store them entirely** (for example because they need some processing to be stored properly), as in scientific applications, or when the application requires an immediate answer

Variety

- Data is extremely heterogeneous: e.g., in the format in which are represented, but also and in the way they represent information, both at the intensional and extensional level
- E.g., text from social networks, sensor data, logs from web applications, databases, XML documents, RDF data, etc.
- Data format ranges therefore from structured (e.g., relational databases) to semistructured (e.g., XML documents), to unstructured (e.g., text documents)

Variety

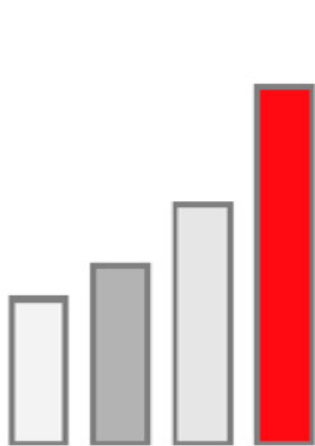
- As for unstructured data, for example, the challenge is to extract ordered meaning for consumption both by humans or machines
- **Entity resolution**, which is the process that resolves (i.e., identifies) entities and detects relationships, then plays an important role
- In fact, these are well-known issues studied since several years in the fields of data integration, data exchange and data quality. In the Big Data scenario, however, they become even more challenging

A fourth V: Veracity*

- Data are of widely different quality
- Traditionally data is thought of as coming from well organized databases with controlled schemas
- Instead, in “Big Data” there is often little or no schema to control its structure
- The result is that there are serious problems with the **quality of the data**
- * The literature often mentions only *three Vs* and does not include veracity. However some authors tend to include veracity as a core characteristic of Big Data (in the other cases, veracity is considered an aspect of variety)

Big Data: V³+ Value

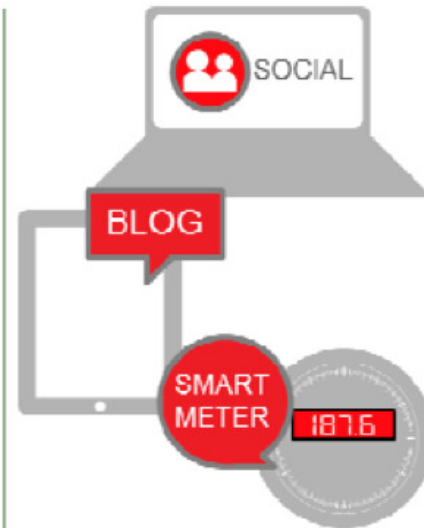
Big Data can generate huge competitive advantages!



VOLUME



VELOCITY



VARIETY



VALUE

The value of Data for organizations

- Although it is difficult to get hard figures on the value of making full use of your data, much of the success of companies such as Amazon and Google is credited to their effective use of data¹
- Thus companies spend large amounts of money to reach this effective use: International Data Corporation (IDC) forecasts that the worldwide **Big Data technology and services market will grow at a 31.7%** compound annual growth rate – about seven times the rate of the overall ICT market – with **revenues reaching \$23.8 billion in 2016**²
- Thus various Big Data solutions are now promoted by all major vendors in data management systems

¹<http://martinfowler.com/articles/bigData/>

²<http://www.idc.com/prodserv/FourPillars/bigData/index.jsp>

Potential value



US health care

- \$300 billion value per year
- ~0.7 percent annual productivity growth



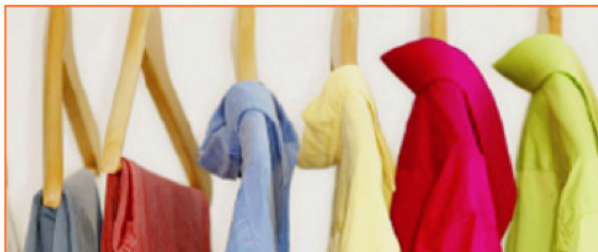
Europe public sector administration

- €250 billion value per year
- ~0.5 percent annual productivity growth



Global personal location data

- \$100 billion+ revenue for service providers
- Up to \$700 billion value to end users



US retail

- 60+% increase in net margin possible
- 0.5–1.0 percent annual productivity growth



Manufacturing

- Up to 50 percent decrease in product development, assembly costs
- Up to 7 percent reduction in working capital

Demand for new data management solutions*

- Despite the popularity and well understood nature of relational databases, it is not the case that they should always be the destination for data
- Depending on the characteristic of data, certain classes of databases are more suited than others for their management
- XML documents are more versatile when stored in dedicated XML store (e.g., MarkLogic)
- Social network relations are graph by nature and graph databases such as Neo4J can make operations on them simpler and more efficient

* From: Edd Dumbill. What is Big data. In Planning for Big Data. O'Reilly Radar Team

Demand for new data management solutions*

- *A disadvantage of the **relational database** is the **static nature of its schema***
- In an agile environment, the results of computation will evolve with the detection and extraction of new information
- Semi-structure **NoSQL databases** meet this need for flexibility: they provide some structure to organize data (enough for certain applications), but do not require the exact schema of the data before storing it

* From: Edd Dumbill. What is Big data. In Planning for Big Data. O'Reilly Radar Team

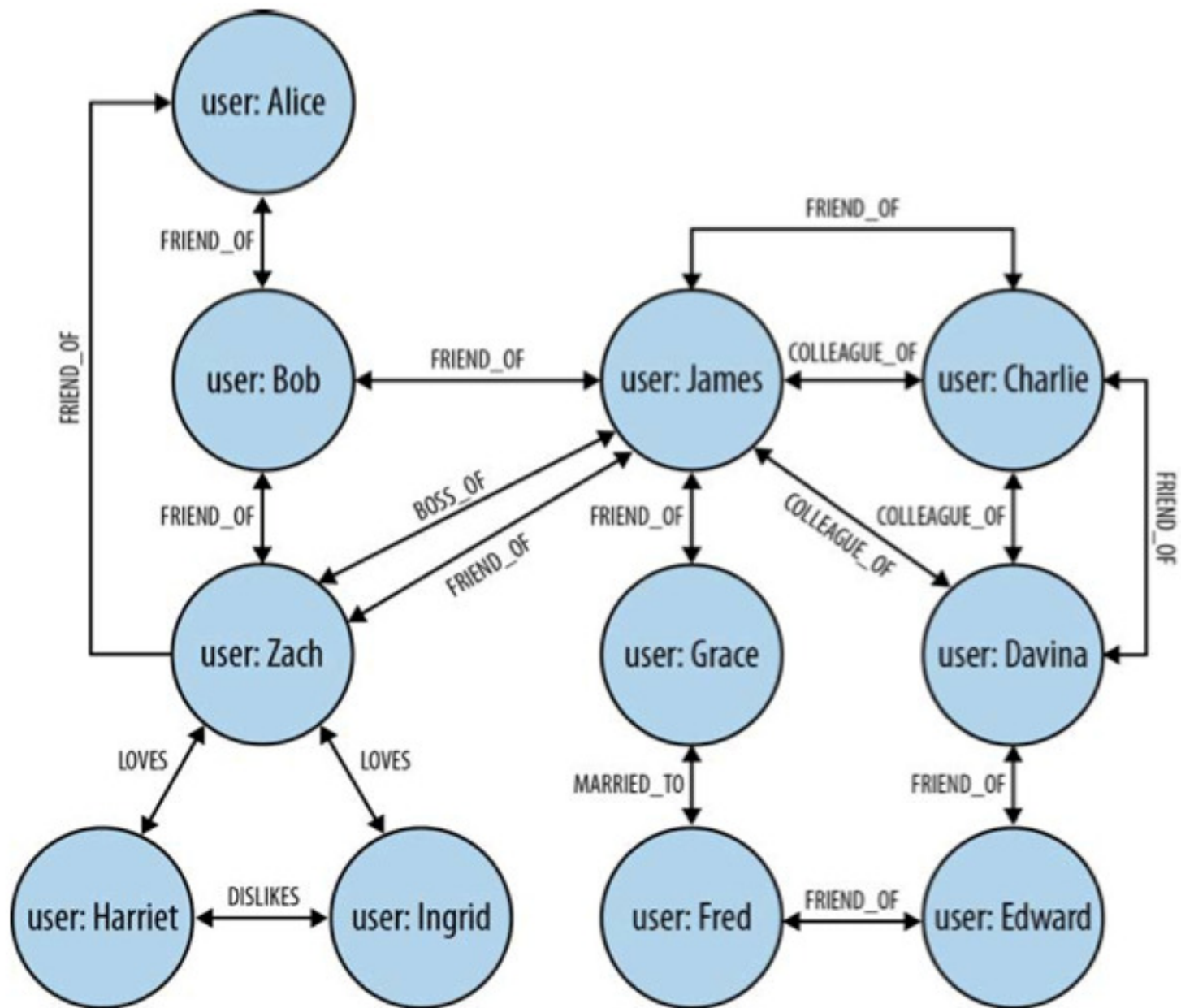
NoSQL databases*

Or better...not only SQL

- The term "NoSQL" is very ill-defined. It's generally applied to a number of recent non-relational databases such as Cassandra, Mongo, Dynamo, Neo4J, Riak, and many others
- They embrace **schemaless data**, **run on clusters**, and have the ability **to trade off traditional consistency** for other useful properties
- Advocates of NoSQL databases claim that they can build systems that are more performant, scale much better, and are easier to program with.

* From: Martin Fowler. NoSQL Distilled. Preface.
(<http://martinfowler.com/books/nosql.html>)

Graph databases



Key-value databases

Key	Value
employee_1	name@Tom-surn@Smith-off@41-buil@A4-tel@45798
employee_2	name@John-surn@Doe-off@42-buil@B7-tel@12349
employee_3	name@Tom-surn@Smith
office_41	buil@A4-tel@45798
office_42	buil@B7-tel@12349

Document databases

Key: "employee_1"



```
{  
  id:" 1" .  
  name:" Tom" .  
  surname:" Smith" .  
  office:{  
    id:" 41" .  
    building:" A4" .  
    telephone:" 45798"  
  }  
}
```

Key: "office_1"



```
{  
  id:" 41" .  
  building:" A4" .  
  telephone:" 45798"  
}
```


Column Family Databases

ColumnFamily: Employees

Key	id	name	surname	office		
employee_1	1	Tom	Smith	id	buil.	tel.
				41	A4	45798

Key	id	name	surname
employee_3	3	Anna	Smith

Key	id	name	surname	office	
employee_2	2	John	Doe	id	buil.
				42	B7

NoSQL databases*

- Is this the first rattle of the death knell for relational databases, or yet another pretender to the throne? Our answer to that is "neither"
- Relational databases are a powerful tool that we expect to be using for many more decades, but we do see a profound change in that relational databases won't be the only databases in use
- Our view is that we are entering a world of Polyglot Persistence where enterprises, and even individual applications, use multiple technologies for data management

* From: Martin Fowler. NoSQL Distilled. Preface.
(<http://martinfowler.com/books/nosql.html>)

Multiple technologies for data management

As an exercise, let us ask google which is the database engine used by Facebook. We get the following tools¹:

- **MySQL** as core database engine (in fact a customized version of MySQL, highly optimized and distributed)
- **Cassandra** (an Apache open source fault tolerant distributed NoSQL DBMS, originally developed at Facebook itself) as database for the Inobx mail search
- **Memcached**, a memory caching system to speed up dynamic database driven websites
- **HayStack**, for storage and management of photos
- **Hive**, an open source, peta-byte scale data warehousing framework based on Hadoop, for analytics, and also **Presto**, a recent exabyte scale datawarehouse²

¹<http://www.techworm.in/2013/05/what-database-actually-facebook-uses.html>

²<http://prestodb.io/>

Data Warehouse

- A data warehouse is a **database used for reporting and data analysis**. It is a central repository of data which is created by integrating data from one or more disparate sources.
- According to Inmon^{*}, a data warehouse is:
 - **Subject-oriented**: The data in the data warehouse is organized so that all the data elements relating to the same real-world event or object are linked together.
 - **Non-volatile**: Data in the data warehouse are never over-written or deleted once committed, the data are static, read-only, and retained for future reporting.
 - **Integrated**: The data warehouse contains data from most or all of an organization's operational systems and these data are made consistent.
 - **Time-variant**: For an operational system, the stored data contains the current value. The data warehouse, however, contains the history of data values.

^{*}Inmon, Bill (1992). *Building the Data Warehouse*. Wiley

Data Warehouse vs. Big Data

- *Are data Warehouses under the hat of Big Data?*
- The concept of data **warehousing dates back to the end of 80s**, and very many data warehouse and business intelligence solutions have been proposed since then.
- BTW, there are many points in common, at least w.r.t. Volume (**data warehouses are large**), Variety (at least in principle, data **warehouses integrate heterogeneous information**), and veracity (data warehouses usually **are equipped with data cleaning solutions**, applied in the so-called extract-transformation-load phase)

Data Warehouse vs. Big Data

- Existing enterprise data warehouses and relational databases excel at processing structured data, and can store massive amounts of data, though at cost.
- However, this requirement for structure imposes an inertia that makes **data warehouses unsuited for agile exploration of massive heterogenous data.**
- The amount of effort required to warehouse data often means that valuable data sources in organizations are never mined.
- Therefore, new computing models and frameworks are needed to make new DW solutions compliant with the Big Data ecosystem.

MapReduce

- MapReduce is a programming framework for parallelizing computation.
- Originally defined at Google.
- Next, there have been various implementations.
- A well-known open source distribution is **Apache Hadoop**.

MapReduce

A MapReduce program is constituted by two components

- **Map()** procedure (*the mapper*) that performs filtering and sorting (it decomposes the problem into parallelizable subproblems)
- **Reduce()** procedure (*the reducer*) devoted to solve subproblems.

The MapReduce Framework manages distributed servers, which execute the various subtasks in parallel.

It both controls communication and data transfers between the various servers, and guarantees fault tolerance and disaster recovery.