



DMIF, University of Udine

Data Management for Big Data

Introduction

Andrea Brunello

andrea.brunello@uniud.it

May 2021



Andrea Brunello, Postdoctoral Researcher at the University of Udine, and member of the Data Science and Automatic Verification Laboratory <https://datasciencelab.dimi.uniud.it/>

My research interests are in the fields of Explainable AI, Virtual Sensing, Machine Learning for Healthcare and, broadly speaking, Data Science. I have also been working on Business Analytics, Data Warehousing and Data Integration projects.

Our laboratory works with local, national, and international partners on Data Science and Automatic Verification projects, with a particular emphasis on the cross-contamination and integration between the two areas.



We are going to spend 8 hours together, covering these topics:

- Relational Databases and SQL
- NoSQL Databases
- Data Warehousing
- Big Data and Business Analytics
- Data Mining
- Business Applications at Gap Srlu company

A Brief History of Information Systems



The Importance of Data

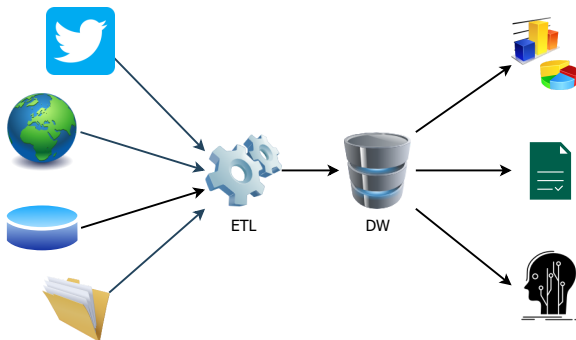
Regardless of the domain, data is driving the future and a massive number of technologies across multiple industries heavily depend on it to thrive.

Data can be defined as a collection of raw, unanalyzed, unorganized material.

Information is data that has been processed, aggregated and organized into a human-friendly format that provides more context (data visualizations, reports, dashboards, ...).

Knowledge derives from a combination of information, experience and intuition. It allows one to draw inferences and develop insights and thus it can assist in decision making.

The Data Analytics Workflow





From Bookkeeping to Magnetic Tapes

Before electronic data processing, companies used to manage their customers, purchases and inventory using traditional bookkeeping methods.

Early electronic data processing came about in the 1950s, initial systems were based on punch cards, naturally exposed to damage and loss of data.

From 1960s magnetic tapes provided better data storage techniques, but sequential access required full tape scan even for 5% of data. Many dedicated hardware as many formats available.



Direct Access Storage Device

The 1970s saw the advent of disk storage, also referred to as direct access storage device (DASD).

No need to go through records $1, 2, 3, \dots, n$ to get to record $n + 1$ once the location address of $n + 1$ is known.

The time to locate a record on a DASD could be measured in milliseconds.

New, more complex data structures were developed, such as lists and trees to be stored on disk. Files + ad-hoc applications.

Along with the DASD, it came a new type of system software known as a database management system (DBMS). Network and hierarchical data models became of widespread use.



From DASD to OLTP

With DBMS it was easier to store and access data on a DASD, moreover, the DBMS took care of tasks such as storing data on a DASD, indexing data, managing access rights, and so forth.

By the mid-1970s, online transaction processing (OLTP) made even faster access to data possible. Applications like bank teller systems and manufacturing control systems became possible.

At the same time, Edgar. F. Codd published his paper on the relational model of systems for managing data. Starting from the 1980s, the relational model has reigned supreme among data models.

SQL becomes industrial standard.



Wal-Mart and the Birth of the DW

Around 1990 Wal-Mart began to achieve wide acclaim for its mastery of supply chain management.

Behind this success was Wal-Mart's data warehouse, and a new way of interacting with data, called OLAP.

Data is collected by its point-of-sales systems to achieve unprecedented insight into the purchasing habits of its 100 million customers and the logistics guiding its 25,000 suppliers.

Wal-Mart's data warehouse was the first commercial Enterprise Data Warehouse to reach 1 terabyte of data in 1992.



Large, Rapid and Heterogeneous Data

Around 2000s, the types of data stored in database systems evolved rapidly, pushed by an ever increasing usage of the Internet and multimedia.

The variety of new data-intensive applications led to NoSQL systems, which gave programmers greater flexibility to work with new types of data, but lacked a high level query language such as SQL.

Distributed storage and computing frameworks were developed, such as Hadoop.

To allow for the interchange of information between systems, formats such as XML and JSON became of widespread usage.

Data Mining applications started to emerge.



Nowadays, most of large and medium sized organizations are using iteratively-evolved information systems to implement their business processes.

As time goes by, these organizations produce a lot of data related to their business, but often these data are not integrated, been stored within one or more platforms.

Such data may be possibly *Big Data*, which are characterized by some very specific issues and opportunities.

Thus, they are difficult to exploit for decision-making processes, though they could be a valuable aiding resource for the management.



A. Silberschatz, H.F. Korth, S. Sudarshan *Database system concepts*, 7th Edition, 2020.

W.H. Inmon, *Building the Data Warehouse*, 4th Edition, 2005.