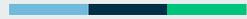


ISA



Data Management and Storage

Lesson 01 - Welcome! Introduction.

Welcome!



Overview of the course:

- Introduction to data management and data management systems
- Concept of relational databases
- Data quality principles, vocabularies and ontologies
- SQL language and application for creating and using databases
- NoSQL systems and their application
- Agro-environment data resources
- Introduction to the concept of large data sets

Assessment



Components:

- Tests (40%) or final exam

Two short tests focused on operational knowledge

Dates: 2024-11-15 (week 9) and 2024-12-20 (week 14)

- Project (40%) - group work

Report on the design and implementation of a data management system

- Participation (20%)

Weekly practice exercises

Project - work group



Goal

Design and implement a data management system integrating data from related domains

Components

- identification of data sources (actual or simulated), but including at least one web data source
- database schema design
- data pre-processing
- database implementation
- data use

Deliverable

- written report including goal, database diagram, documentation and use examples
- database backup dump with README.md instructions on how to implement/recover
- SQL script file with sample queries for data manipulation

Course management



Fenix: <https://fenix.isa.ulisboa.pt/courses/gestadad-283463546571604>

Github: <https://github.com/isa-ulisboa/greends-dms-exercises>

Moodle: <https://elearning.ulisboa.pt/course/view.php?id=9101>

What is data?



Examples of data:

- identify types of data
- identify types of data sources

MIRO Activity

<https://tinyurl.com/greends-dms24-01>



Give examples of types of data

Pick a post-it and reply in few words

Agricultural
crop data -
images, sensor
collected data,
climate data

grafics,
numerics,
precipitation,

Tree
growth
data

sales,
prices,
text,

Dados
meteorológicos
vento, radiação,
temperatura,
humidade do ar

Dates
and
images

precipitation
data

satellite
imagery

sales
reports

wind, UV
index,
topography
data

Number of
rainy days,
colour of
grapes.

Give examples of data sources

Pick a post-it and reply in few words

:~)

university
databases

satellite
sensors

governmental
agencies reports

:~)

private
databases

fertiliser
records

Ine,
copernicus,
ipma

Annual
reports,
scientific
papers

statistic,
social
networks,
researches

LiDAR
sensors,
aeroplane
cameras

What is data?



DATA:

- Facts about things (text, numbers)
 - properties, measurements, descriptions, etc...
- Sounds, images, movies

INFORMATION:

- data with context (data organised and interpretable)

**Algés
is the locality
where I live.**

- data is a collection of facts, information puts them in context
- data not organised, information is
- data points are unrelated, information can provide links
- data does not have meaning *per se*. After being analysed, it may become information
- data does not depend on information. Information depends on data
- data does not support decision making, information can support it.

What is data?



Algés
is the locality
where I live.

- data is a collection of facts, information puts them in context
- data not organised, information is
- data points are unrelated, information can provide links
- data does not have meaning *per se*. After being analysed, it may become information
- data does not depend on information. Information depends on data
- data does not support decision making, information can support it.

Information systems



Data is the new oil!

The evolution from an industrial to an information and knowledge society is represented by the assessment of information as a factor in production. But different of material goods, because:

Representation

Processing

Combination

Age

Original

Vagueness

Medium

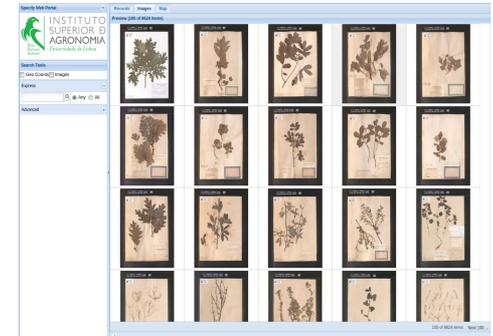
<https://doi.org/10.3390/fi15020071>

Data store

Data store

Repository for persistently storing and managing collections of data. Can be based on

- databases
- files
- file systems
 - Hadoop Distributed File System (HDFS) - big data
- email storage systems
- distributed data stores



Data store



Data store - Can be organised in tree types

Data lake

Storage repository that holds vast amounts of raw, unstructured, semi-structured, and structured data in its **native format**.

- Google drive
- AWS S3
- OneDrive

Data Warehouses

Centralized repository that stores structured (or semi-structured), cleaned, and processed data from various sources.

- relational databases
- cloud databases

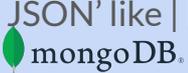
Data Platform

Unified system for efficiently managing and analyzing large datasets. It integrates components like databases, data lakes, and data warehouses to handle structured and / or unstructured data depending on the use cases

- Data lakes
- Data warehouses

Databases

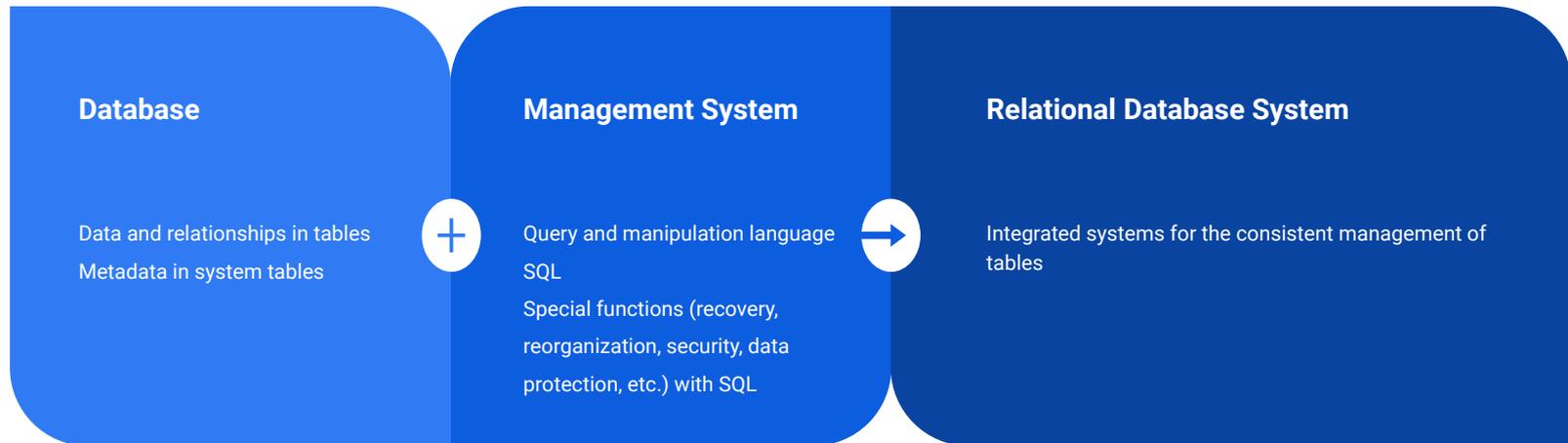
Collection of data stored in a structured manner.

Hierarchical Databases	Relational Databases	Non-Relational Databases
<p>A single object (the “parent”) has one or more objects beneath it (the “child”). Child records can’t have more than one parent.</p> <p>Windows Registry</p>	<p>Also known as RDBMS. Store information in discrete tables, which can be joined together by fields known as foreign keys. Use SQL for operations.</p> <p> MariaDB</p> <p> MySQL</p> <p> ORACLE</p> <p> IBM</p> <p>SQL Server</p> <p>DB2</p>	<p>NoSQL differ from RDBMS in that they can be schema-agnostic, allowing unstructured and semi-structured data to be readily stored and processed.</p> <p>JSON’ like key : value graph</p> <p> mongoDB</p> <p> neo4j</p> <p> cassandra</p>

RDBMS



Relational DataBase Management Systems



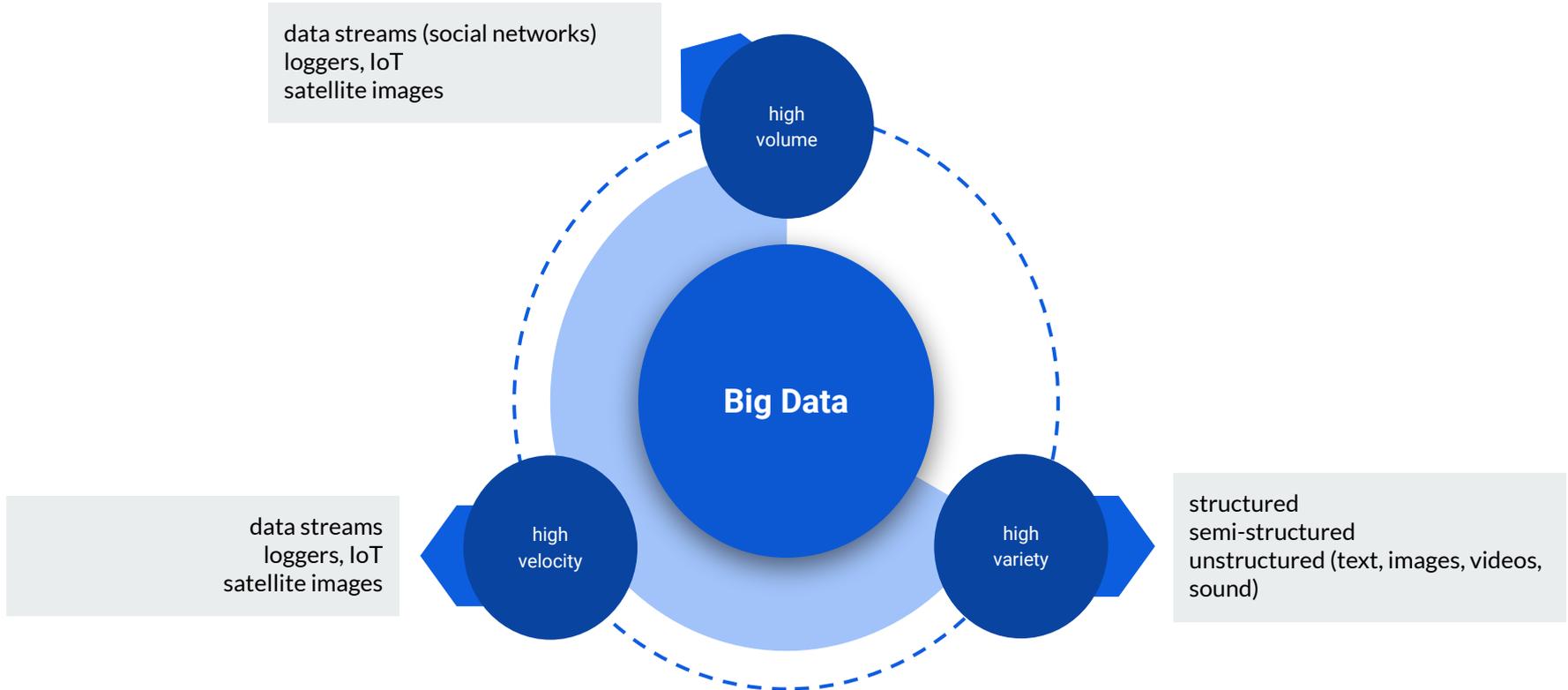
RDBMS



Relational DataBase Management Systems - properties

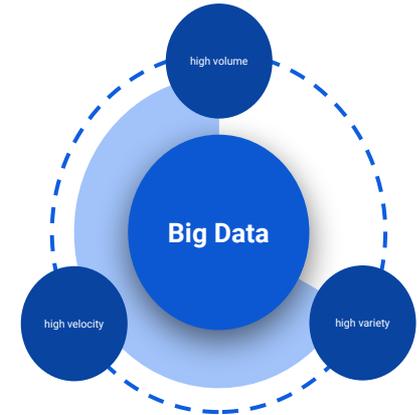
- **Model:** relational model, data represented in tables
- **Schema:** definition of tables and attributes
- **Language:** SQL data definition
- **Architecture:** data independence, data separated from applications
- **Multi-user operation:** several users
- **Consistency assurance:** data integrity
- **Data security and data protection:** protection from loss, destruction, unauthorized access

Big Data



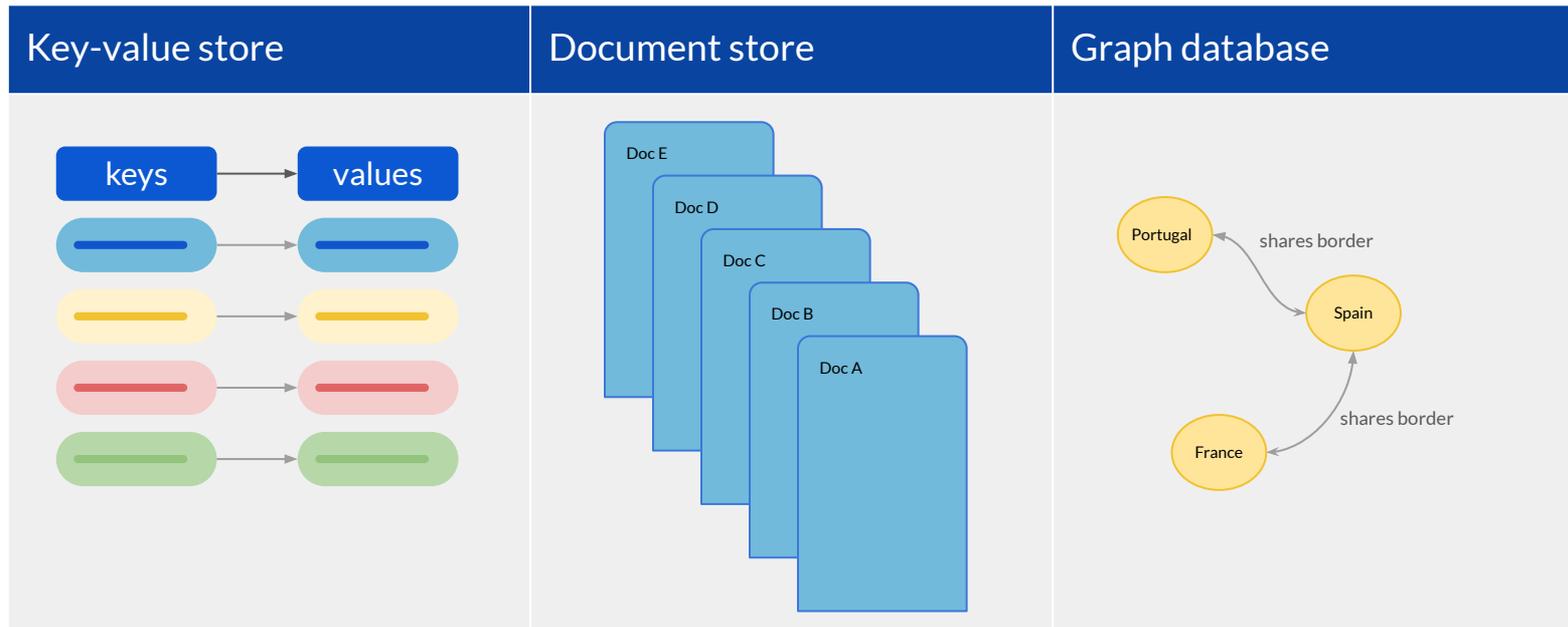
Big Data - examples

01	Text	<ul style="list-style-type: none">• continuous text• structured text• collection of texts, tags, etc
02	Graphics	<ul style="list-style-type: none">• road maps• 3D graphics• technical drawing, etc
03	Image	<ul style="list-style-type: none">• photographs• satellite images• x-ray image, ect
04	Audio	<ul style="list-style-type: none">• language• music• sounds, animal sounds, etc
05	Video	<ul style="list-style-type: none">• film• animation• ads, phone conferences, etc



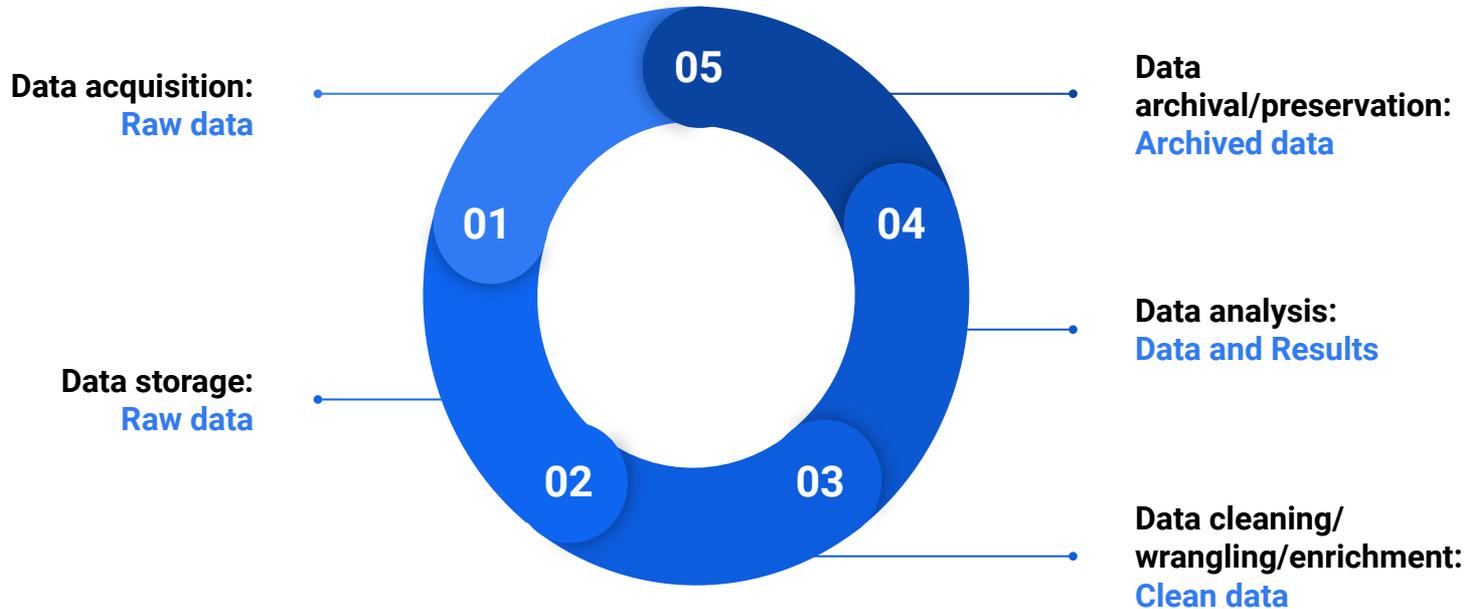
NoSQL Databases

Nonrelational data management approaches (in result of Web 2.0)

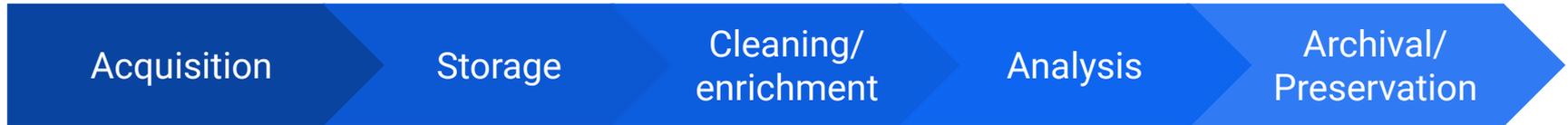


Data Life Cycle

Data life cycle - transformations applied to data and to the states that data goes through as a result of these transformations



Data Life Cycle



- acquire, collect, gather
- sample of the population
- opportunistic (easy to access)
- observational - uncontrolled settings

- extract, transform, load (ETL)
- EDA
- metadata (augmentation)
- visualization techniques

- inconsistencies
- missing, entry errors
- outlier detection
- duplicates
- augmentation
- transformation (scaling, normalization, binarization, etc.)

- dimension reduction
- supervised (regression, classification)
- unsupervised (clustering, kNN)

- deleted
- archived
- published (metadata)

Data management in organisations



Roles	Tasks	Tools
Data architecture	<ul style="list-style-type: none">- Creation and maintenance of data architecture- Definition of data protection rules	<ul style="list-style-type: none">- Data analysis and design methodology- tools for information modeling
Data administration	<ul style="list-style-type: none">- Management of data and methods using standardisation guidelines and international standards- Consultation of developers and end users	<ul style="list-style-type: none">- Data dictionary systems- Tools for cross-reference and usage lists

Data management in organisations



Roles	Tasks	Tools
Data technology	<ul style="list-style-type: none">- Installation, reorganisation and secure data content- Definition of the distribution concept- Disaster prevention and recovery	<ul style="list-style-type: none">- Database system services- Tools for performance optimisation- Tools for recovery/restart
Data utilisation	<ul style="list-style-type: none">- Data analysis and interpretation- Knowledge creation- Pattern detection- Modelling	<ul style="list-style-type: none">- Analysis tools- Report generation- Data mining tools- Visualization methods

Data management in organisations

Data specialists



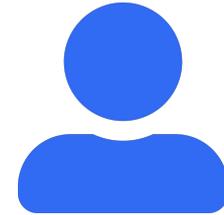
Data architects:

- access adjusted to business model
- distribution
- fragmentation
- replication



Data specialists:

- database technology
- manage infrastructure
- security



Data scientists:

- analysis
- modelling
- interpretation
- prediction

Additional reading



Instituto Nacional de Estatística - Recenseamento Agrícola. Análise dos principais resultados : 2019. Lisboa : INE, 2021.
Disponível em <https://www.ine.pt/xurl/pub/437178558>. ISBN 978-989-25-0562-6