

## Introducción a CS de datos · Conceptos

### I. tipo de datos

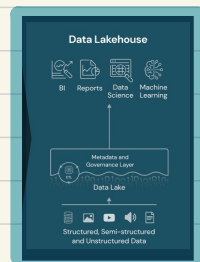
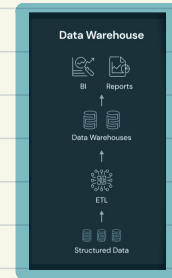
- A. Estructurado : - organizado en un formato predefinido y homogéneo  
(Excel, BDD sql,) - está en tablas con filas y columns
- B. semiestructurados : - no están ordenados con filas-columns  
(JSON, XML, mails) - Están jerarquizados ; utilizan sistema de llaves / etiquetas  
similar a un diccionario en python: tiene "key-values"
- C. no estructurada : - no tiene estructura ni jerarquía  
(logs, posts, RSS, ★ Se pueden almacenar en DATA LAKES!!  
PDF)



OBS: un dataset es un formato tabular en donde los datos (registros) se organizan por atributos.

### II. ¿Cómo guardamos y organizamos los datos?

- a. Data warehouse → enfoque con fin analítico
  - Integra distintas bases de datos
  - Cumple modelo ACID: Atomicidad Consistencia Aislamiento Durabilidad
    - ↳ en resumen garantiza la integridad y fiabilidad en procesos dentro de la base de datos
- b. Data Lakehouse
  - Permite guardar y organizar datos estructurados, semiestructurados y no estructurados en un solo lugar / sistema
  - Flexible y escalable



Fuente: Databricks

### III. ¿De Donde sacamos datos?

#### III.a Web



- Ⓐ → URL (clásico descarga archivos)
- Ⓑ → API ("interfaz de programación de aplicaciones"): permite la comunicación entre 2 componentes de Software a través de reglas y protocolos (funcionan en sentido cliente-servidor)

OBS: para el curso se utilizan las librerías requests y BeautifulSoup

→ métodos en requests: .get(), .post() / peticiones protocolo http

NOTA: apis pueden ser públicas o privadas, las respuestas que se obtienen pueden estar en formato JSON.

- Ⓒ → Webscrapping!! : extraer automatizada los datos de una página web
  - ★ Aquí entra BeautifulSoup → trae datos html o XML !!

★ Concepto: DATA WRANGLING → Proceso de limpiar, transformar y organizar los datos

★ Concepto2: ETL (extraer transformar cargar)

Ambos tienen el objetivo de preparar los datos pero varían en la ejecución y formato

- ETL es más repetitivo, para grandes volúmenes, formal y planificado
- Data Wrangling es más flexible y enfocado en análisis específicos que se requieran. Semi automatizado / manual

