

Especialización en Inteligencia de Datos Aplicada



Facultad de Informática Universidad Nacional del Comahue

EXTRACCIÓN, PREPARACIÓN Y ALMACENAMIENTO DE LOS DATOS

Segundo Cuatrimestre Unidad I

EXTRACCIÓN, PREPARACIÓN Y ALMACENAMIENTO DE LOS DATOS

CONTENIDOS MÍNIMOS

Captura de datos. Fuente de datos. Captura de datos en la Web. Almacenamiento en bases de datos relacionales y no relacionales. Repositorios NoSQL. Depósitos de datos. Limpieza de datos. Lagos de datos.

UNIDAD I

IDENTIFICACIÓN Y RECOLECCIÓN DE DATOS

Ingesta de datos. Tipos de datos fuente. Extracción de datos de la Web (Web Scraping, Web Crawling). Patrones de ingesta de datos en el contexto de Big Data.

UNIDAD II

PREPARACIÓN DE DATOS

Calidad de Datos. Limpieza de Datos. Ruido y detección de anomalías en los datos. El proceso ETL y ELT. Integración de los datos.

UNIDAD III

ALMACENAMIENTO DE DATOS

Tipos de almacenamiento.
ACID-CAP-BASE. NoSQL,
Distribuidas. Depósitos de Datos vs
Lago de Datos. Conceptos del
almacenamiento de grandes
volúmenes de datos.



Especialización en Inteligencia de Datos Aplicada



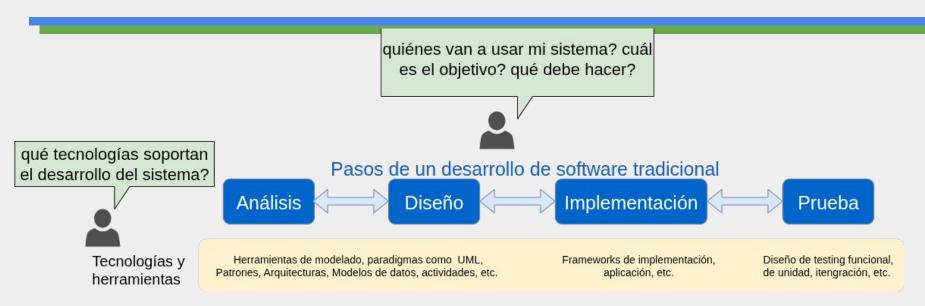
Facultad de Informática Universidad Nacional del Comahue

pero,

cómo comenzamos?



Cómo se construye un sistema?



Qué tipo de sistema voy a construir?

Tipos de Sistemas

- Podemos distinguir entre dos tipos fundamentales:
 - Sistemas Transaccionales: gestión de las transacciones diarias,
 automatización de los procesos empresariales.
 - Sistemas para la Toma de Decisiones: proporcionar una visión amplia, integral y profunda de los datos de una organización, de forma de tomar decisiones informadas.

Tipos de Sistemas



LA TOMA DE DECISIONES

SISTEMAS
TRANSACCIONALES

Soporte según tipos de sistemas

OLTP (OnLine Transaction

Processing): Procesamiento en línea para el manejo transaccional



bancos, aerolineas, universidades, seguros, etc. OLAP (OnLine Analytical Processing): Procesamiento en línea para la toma de decisiones

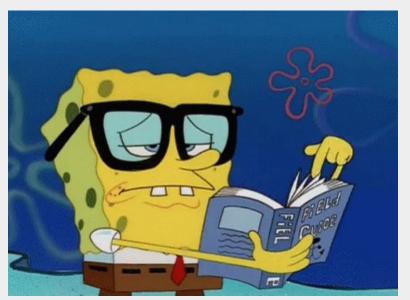




lagos de datos

información agregada, gerencial, toma de desiciones

Ahora sí.....





Sistemas para la Toma de Decisiones

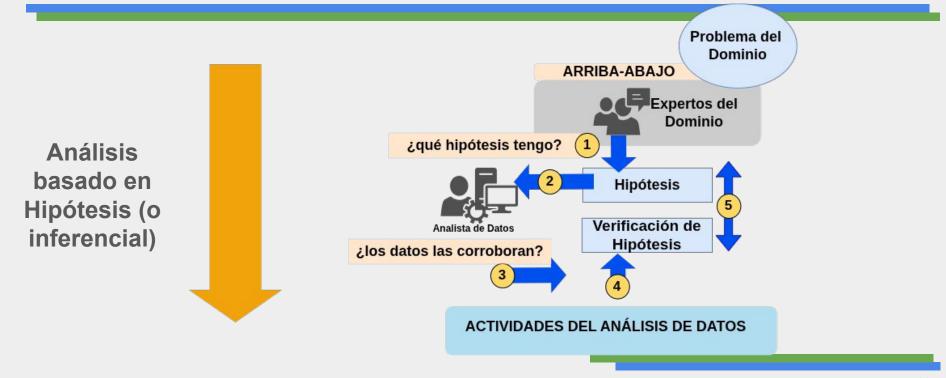
- Se deben seguir una serie de pasos para su construcción
- Al conjunto de estos pasos lo podemos llamar

PROCESO DE ANÁLISIS DE DATOS

Basado en Hipótesis

Exploratorio

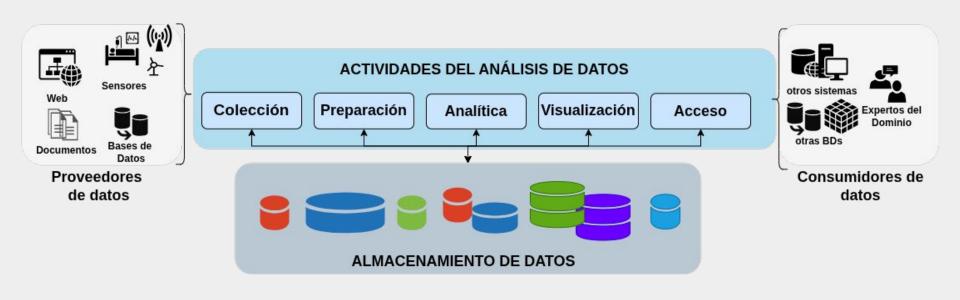
PROCESO DE ANÁLISIS DE DATOS de arriba a abajo



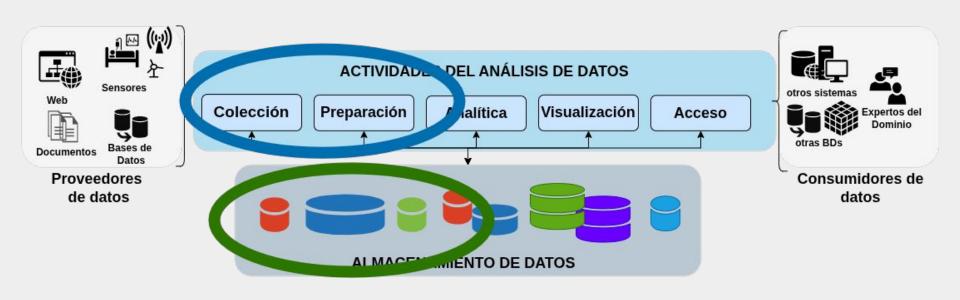
PROCESO DE ANÁLISIS DE DATOS de abajo a arriba



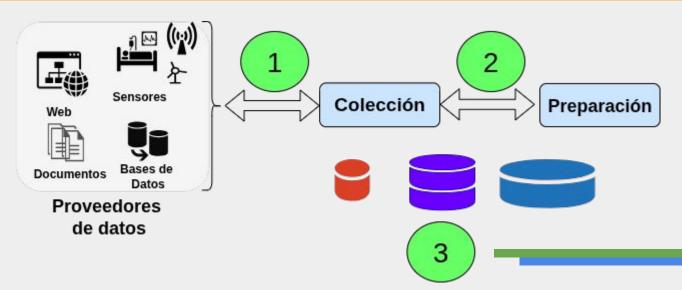
PROCESO DE ANÁLISIS DE DATOS Actividades

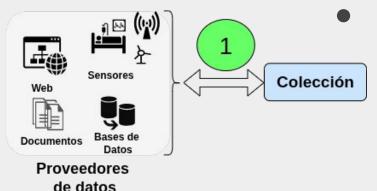


PROCESO DE ANÁLISIS DE DATOS Actividades



Es el proceso de coleccionar datos de diferentes fuentes y colocarlos en un repositorio destino





Cómo deben ser identificados y recolectados los datos?

- Muchas fuentes de datos involucradas
- Diferentes formatos de datos
- y además....

- Qué aspectos debemos considerar?
 - Evolución de las fuentes
 - Las fuentes de datos pueden cambiar, por ejemplo las Web, estructuras de las BDs, nuevas fuentes a agregar.



- Qué aspectos debemos considerar?
 - Seguridad y Privacidad
 - De dónde obtengo los datos? pueden los usuarios conocer esa información?



- Qué aspectos debemos considerar?
 - Periodicidad de recolección
 - Con qué frecuencia las fuentes deben ser recuperadas/coleccionadas? De qué depende?



PROCESO DE ANÁLISIS DE DATOS Ingestión de Datos - Periodicidad

- La periodicidad en que se necesitan los datos se enmarca en tres modos específicos:
 - Por lotes (batch)
 - Tiempo real (stream)
 - Micro-lotes (micro-batch)



- Se coleccionan grandes volúmenes de datos cada cierta cantidad de tiempo como horas, días, meses, etc.
- No requieren una acción inmediata.
- Los lotes pueden estar definidos por tamaño o por tiempo.
- Usos en Análisis: Análisis de datos históricos, predicciones, optimización, clasificación.



PROCESO DE ANÁLISIS DE DATOS Ingestión de Datos - Tiempo real

 Se coleccionan los datos ni bien son (o casi) por las fuentes.



- Requieren una acción inmediata.
- Complejidad en cuanto al envío de la información que se genera (fallas de red, consistencia, envíos duplicados).
- *Usos en Análisis*: detección de fraudes, logística (transporte), alertas, etc.

PROCESO DE ANÁLISIS DE DATOS Ingestión de Datos - Micro-Lotes

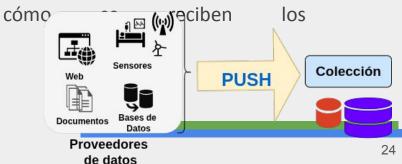
- Se coleccionan lotes más chicos y con más frecuencia.
- *Usos en Análisis*: comportamientos de usuarios, minería de la web, etc.



- Otros tipos de patrones para la colección de los datos son:
 - Pull (extracción): el destino lee datos directamente de las fuentes
 - Push (inserción): las fuentes envían datos al destino.



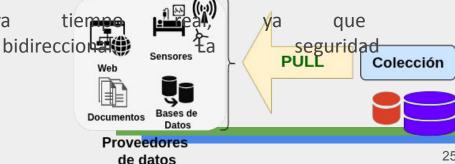
- Los datos se envían desde el origen (fuentes) al destino (sistema) en cuanto están disponibles.
- El origen puede generar muchos/pocos datos en corto/largo tiempo.
- **Ventajas**: *Tiempo real* ya que nuevos datos se envían al destino, *eficiente* cuando las fuentes producen datos constantemente y *seguridad* ya que se hace en destino.
- Desventajas: Es difícil controlar datos en cuanto a su frecuencia, tamaño, etc.



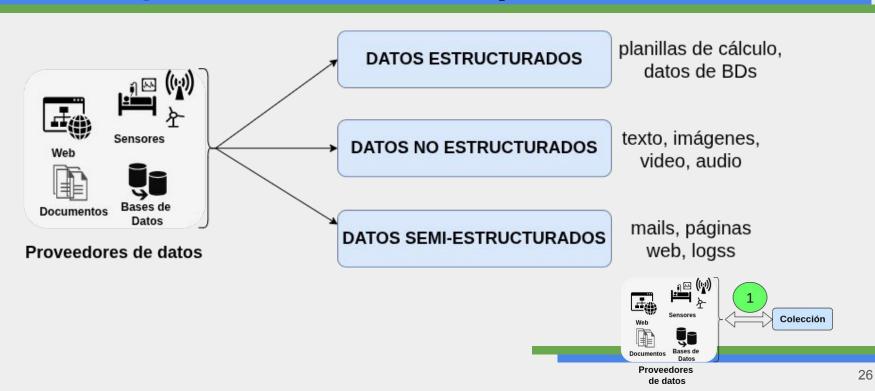
- El destino envía una solicitud al origen y obtiene una respuesta con o sin datos.
- Permite al consumidor obtener datos de forma controlada. El consumidor puede retrasarse y recuperarse cuando lo desee.
- Ventajas: escalabilidad ya que se puede modificar según se desee y a preferencia del consumidor.

para

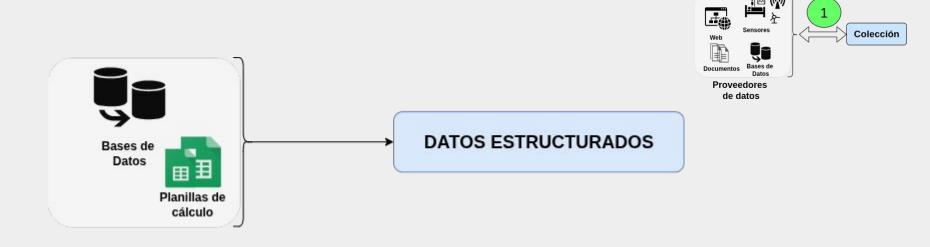
 Desventajas: No propicio la comunicación es está en las fuentes.



PROCESO DE ANÁLISIS DE DATOS Ingestión de Datos -Tipos de Datos



PROCESO DE ANÁLISIS DE DATOS Ingestión de Datos - Estructurados



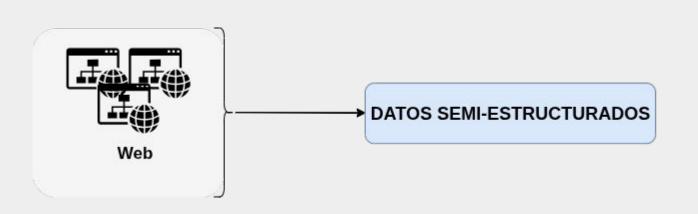
PROCESO DE ANÁLISIS DE DATOS Ingestión de Datos - Estructurados



PROCESO DE ANÁLISIS DE DATOS A TRABAJAR!!



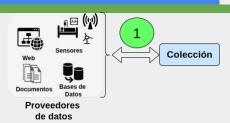
PROCESO DE ANÁLISIS DE DATOS Ingestión de Datos -SemiEstructurados





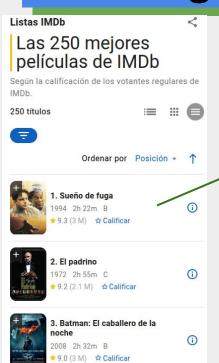
Colección de datos desde la Web

- basados en tags de HTML, XML y contenido
- Recursos:
 - web crawling
 - web scraping



- Web crawling (o rastreadores Web) es el proceso de explorar la Web (mediante programas automáticos) para recuperar datos.
 - Se navega las páginas usando los enlaces (hyperlinks) y recuperando las páginas siguientes.
 - Los datos coleccionados se utilizan para indexación (de los motores de búsqueda) o procesamiento posterior.
- Web scraping (o raspado Web)

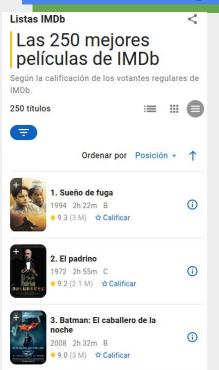




```
<ur><ur>class= ibc-meradara-fist-summary-frem
  <span class="ipc-metadata-list-summary-item t" aria-disabled="false"></span>
▼ <div class="sc-e2db8066-1 QxXCO cli-parent li-compact"> flex
  ▼ <div class="sc-e2db8066-0 iXaKBg"> flex
    ▼ <div class="sc-ee514ad1-0 kYZRWL cli-poster-container"> flex
      <div class="ipc-poster ipc-poster--base ipc-poster--media-radius ipc-poster--dynamic-width ipc</p>
     </div>
    ▼<div class="sc-f30335b4-0 eefKuM cli-children"> flex
       wedin class="inc title inc title hase inc title title inc title link no icon inc title-on-te
        <<a href="/title/tt0111161/?ref =chttp t 1" class="ipc-title-link-wrapper" tabi dex="0">
          </a>
       </div>
              class="sc-f30335b4-6 kGhnhC cli-title-metadata"> flex
                class="sc-f30335b4-7 jhjEEd cli-title-metadata-item">1994</span>
                             30335b4-7 ihiEEd cli-title-metadata-item">2h 22m</span>
6 VIDEOS
                ■ 99+ FOTOS
           Period Drama
                                      ref="/whats-on-tv/?ref =nv tv ontv" tabindex="-1"
                               one
                                      ref="/chart/toptv/?ref =nv tvv 250"
                                                                            tabindex="-1"
                               one"
     por matar a su esposa y amante. Tras una
                                     nref="<u>/chart/tymeter/?ref =nv tvv mp v</u>" ta<del>bi</del>ndex="-1" <u>aria-</u>disabled="false"
                               one"
     dura adaptación, intenta mejorar las
                                     ref="/feature/genre/" tabindex="-l" aria-disabled="false"> ... </a> (flex)
     condiciones de la prisión y dar esperanza
                               one
                                     nier<u>= /news/cv/?ier =nv nw cv capind</u>ex="-1" aria-disabled="false">... </a>
```

- Web crawling (o rastreadores Web)
- Web scraping (o raspado Web) se centra en extraer datos específicos de los sitios y exportarlos para un almacenamiento en formato XML, Excel, JSON o SQL. Su automatización permite recopilar datos para la minería.



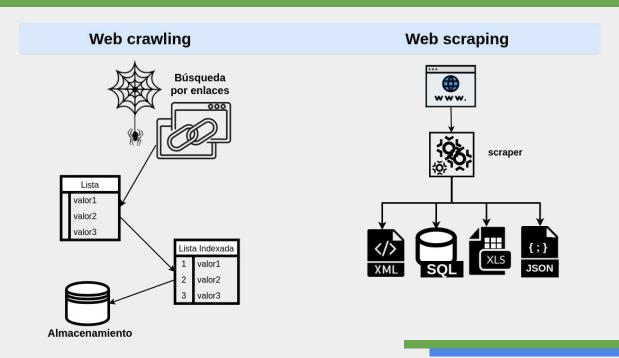


titulo de pelicula	año	duración	categoría	calificación
Sueño de Fuga	1994	2hs	В	9.3
El padrino	1972	3hs	С	9.2
Batman	2008	2.30hs	В	9

scraping

.CSV

titulo de pelicula,año,duración,categoría,calificación Sueño de Fuga,1994,2hs,B,9.3 El padrino,1972,3hs,C,9.2 Batman,2008,2.30hs,B,9



PROCESO DE ANÁLISIS DE DATOS A TRABAJAR!!



PROCESO DE ANÁLISIS DE DATOS con qué seguimos?

