

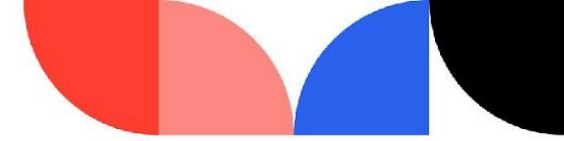
MÀSTER DE FORMACIÓ PROFESSIONAL

INTEL·LIGÈNCIA ARTIFICIAL I BIG DATA

TITULACIÓ OFICIAL FP

M4 R1 C1C
**Fuentes y Tipos
de Datos**

**Clase
2**



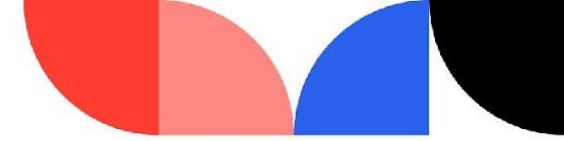
Objetivos de la Clase

- Entender la importancia de combinar diferentes fuentes y tipos de datos en proyectos de análisis de datos.
- **Conceptos Clave:** integración de datos, limpieza de datos, y uniones de datos.

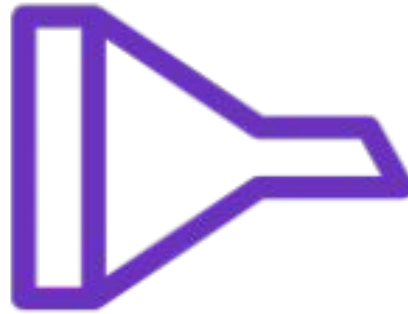
**“Los datos son el
nuevo petróleo”
- Martin Hilbert**



Imagen: The Economist Magazine - 2017



DATOS VS. INFORMACIÓN



Característica o atributo **sin procesamiento**, el cual no informa nada por sí solo.

Unión de **datos procesados**, que se complementan para informar un hecho.

TIPOS DE DATOS



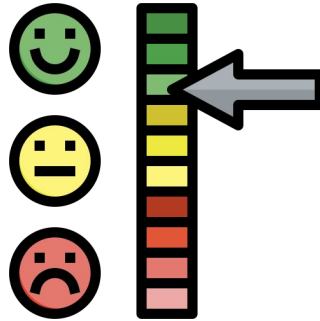
Cualitativos o Categóricos

Nominales



Género, nacionalidad, color
Habilidades, preferencias,
comida o pasiones.

Ordinales



Ratings (bueno, regular,
malo) o Niveles de
educación (primaria,
secundaria, universidad).

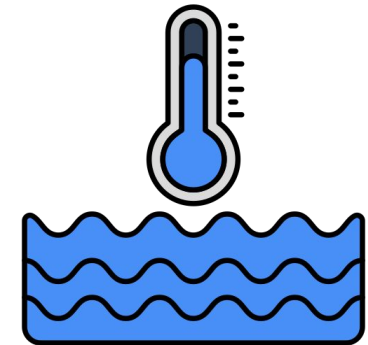
Cuantitativos

Discretos



Pueden tomar ciertos
valores y no pueden ser más
precisos: estudiantes en
clase BI & IA en Stucom

Continuos



Pueden tomar cualquier
valor dentro de un rango
determinado: temperatura,
peso, altura.

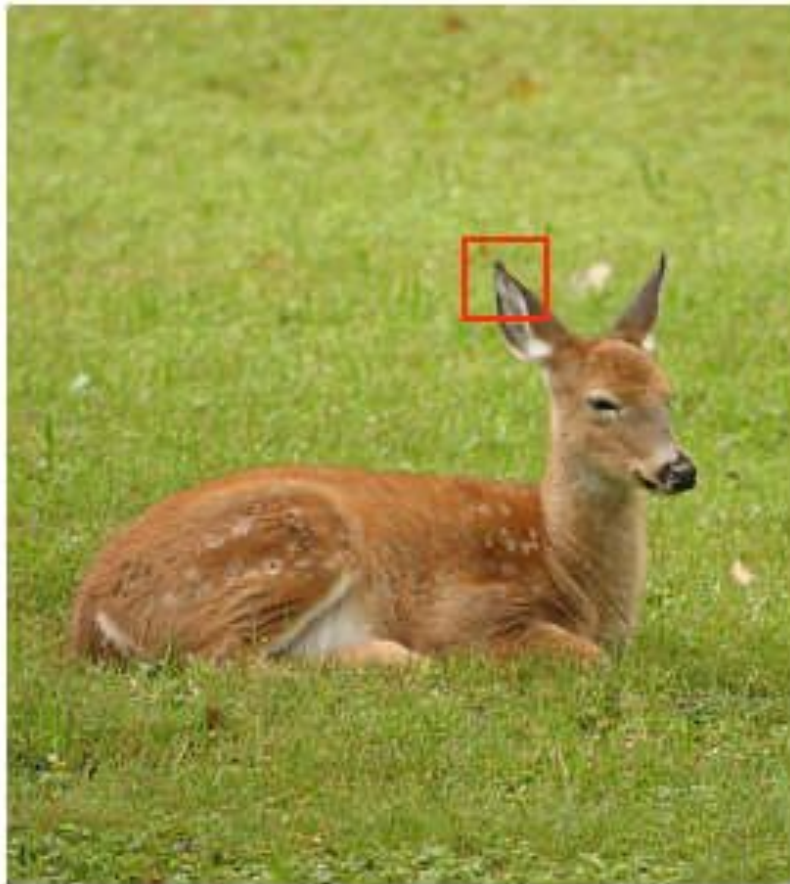
TIPOS DE DATOS	Cualitativos o Categóricos	Cuantitativos
Nominales	Etiquetas sin orden	
Ordinales	Etiquetas con orden	
Discretos		Valores fijos
Continuos		Rangos de valores

Cualitativos o Categóricos	Cuantitativos
Mi mejor amigo pesa 75 kilos	Mi mejor amigo _____
Mi madre ha conducido _____	Mi madre conduce en un coche rojo
El equipo de natación entrenó _____	El abuelo es divertido y sabe escuchar
El terminó la maratón en 3 hr y 45 mins	El tiene unos ojos _____ y un cabello _____



OTROS TIPOS DE DATOS

DATOS DE IMAGEN





A screenshot of a tweet from the official Twitter account (@Twitter). The tweet text is "hello literally everyone". It was posted at 1:27 PM on October 4, 2021, via the Sprinklr platform. The tweet has received 335.6K retweets, 117.8K quote tweets, and 1.6M likes. The interface shows a back arrow, the Twitter logo, and standard engagement icons (reply, retweet, like, share) at the bottom.

Amazon Customer

★★★★★

It works great!

Reviewed in the United States on November 3, 2018

Verified Purchase

We love the lamp! We use it as a night light. It works great. We keep it on red since it slowed me to see the baby and is not bright at all. The white light makes my room too bright and I can't sleep. It has different colors available and it can rotate while you sleep. It doesn't make much noise at all it will let you sleep. (Only makes minimal noise while rotating). It's a great gift. The material does feel cheap but we get what we pay for. I would so buy it again if anything was to happen to this one. Yes the material may be cheap but it works great. Like I said before, we love it!

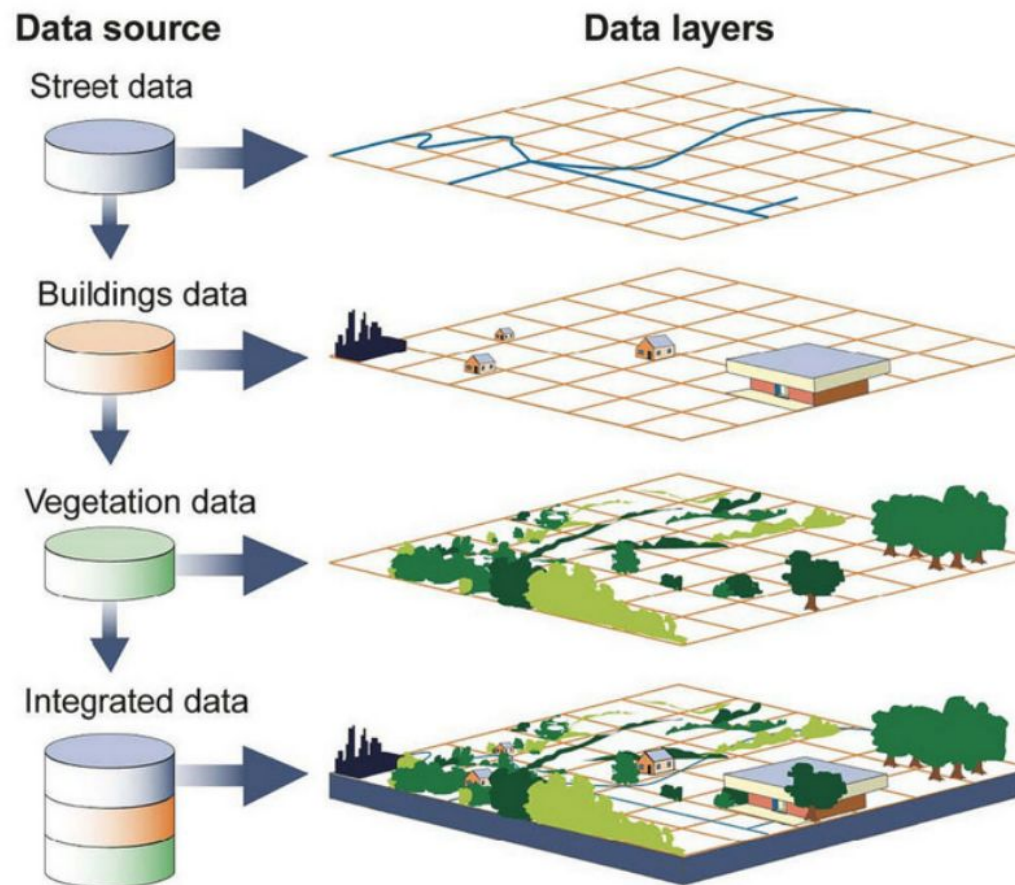
65 people found this helpful

Helpful

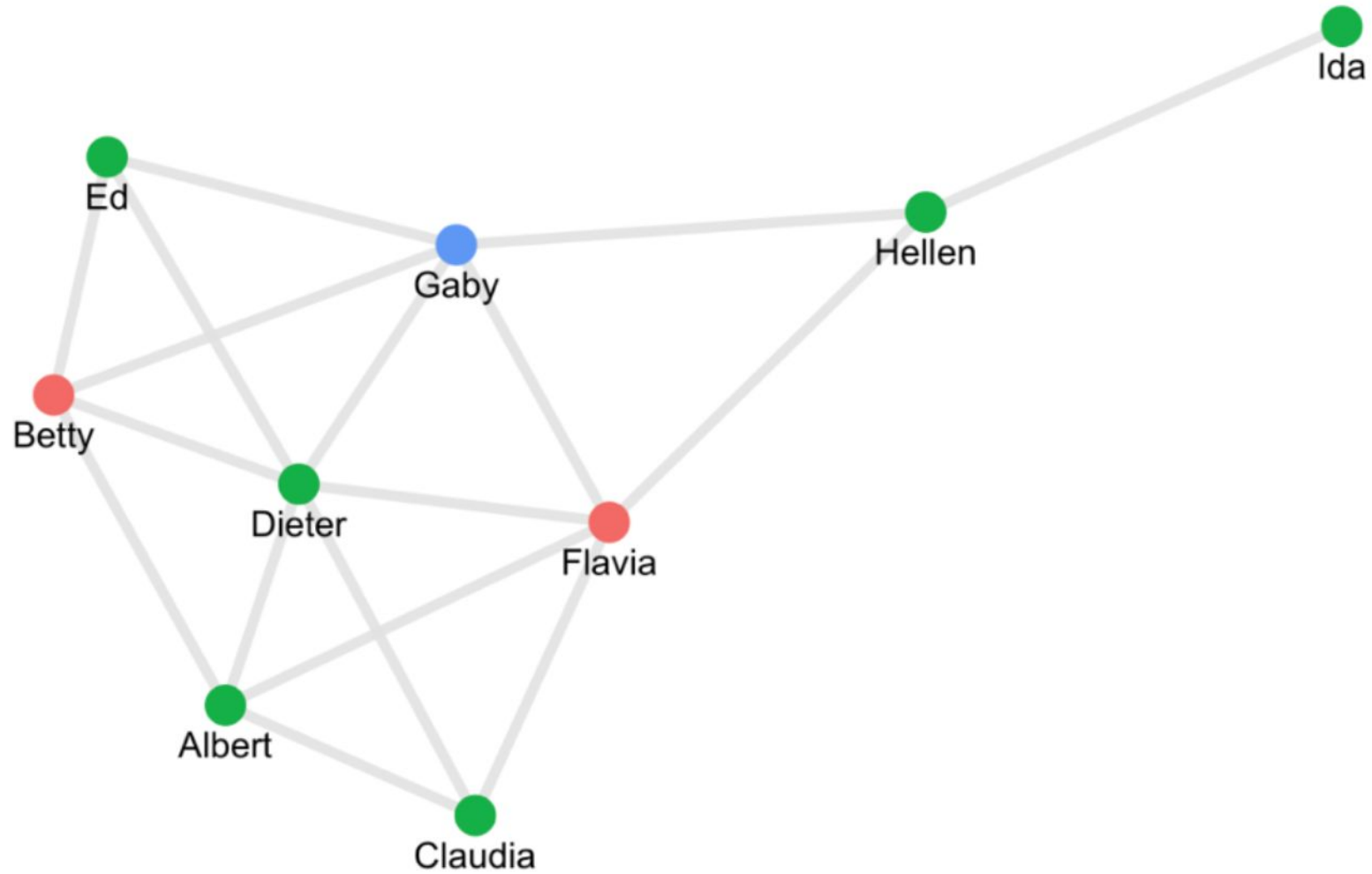
Comment

[Report abuse](#)

DATOS GEOESPACIALES



NETWORK DATA



Net Promoter Score (o NPS) es una métrica común que las empresas utilizan para realizar un seguimiento del éxito de un producto o sitio web. Se mide haciendo una simple pregunta:

¿Qué tan probable es que recomiendes [insertar marca/sitio web/servicio/producto] a un amigo o colega?

Los usuarios responden en una escala de 0 a 10, siendo 0 que no es en absoluto probable que se recomiende y 10 es extremadamente probable que se recomienda.

How likely is it that you would recommend our product to a friend or colleague?

Not at all likely

Extremely likely

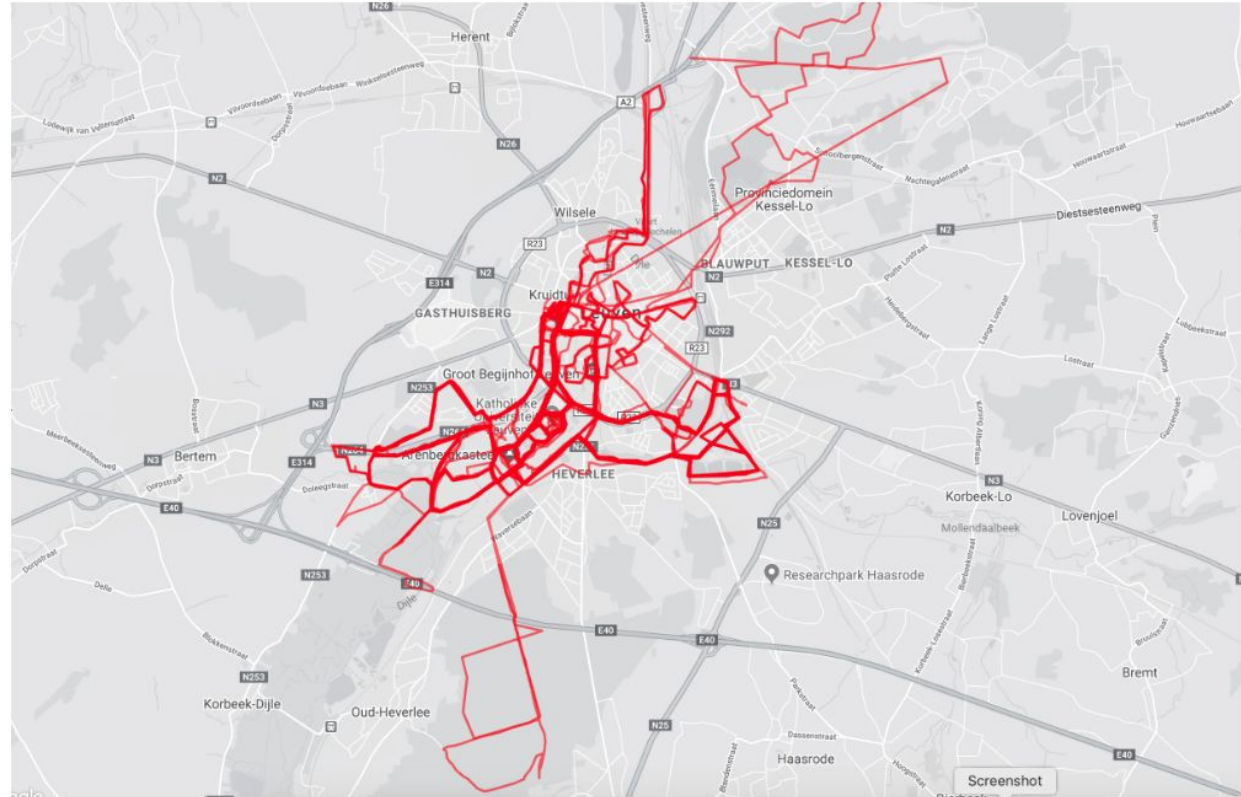
0	1	2	3	4	5	6	7	8	9	10	11
---	---	---	---	---	---	---	---	---	---	----	----

Cualitativo

Cuantitativo

La resolución de Año Nuevo de Jane este año era ponerse en la mejor forma de su vida. Para ayudarla a alcanzar este objetivo, decidió invertir en un rastreador de actividad. Después de algunos meses de rastrear su actividad, hay bastantes datos disponibles.

La empresa que fabricó el rastreador de actividad tiene una API pública que permite el acceso a sus datos personales. Jane está específicamente interesada en los datos GPS de sus carreras porque quiere hacer un mapa de calor que muestre sus rutas de carrera más comunes.



Imagen

Texto

Geoespacial

Network

¡IMPORTANTE!



UBICACIÓN (ALMACENAMIENTO)



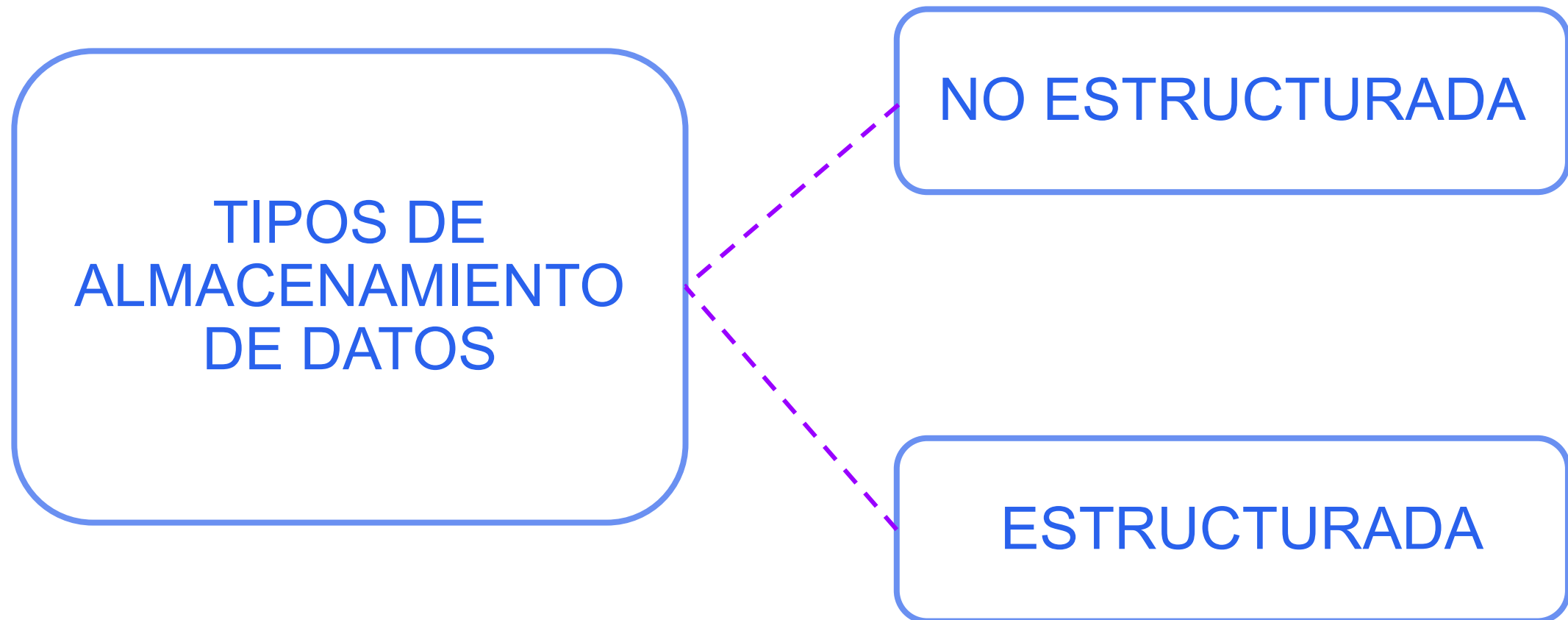
Solución de Almacenamiento Paralelo

- Alto rendimiento local
- Acceso simultáneo y rápido



CLOUD (La Nube)

- Accesibilidad
- Flexibilidad



NO ESTRUCTURADA



Tipos de Datos

Email's

Texto

Páginas Web

Redes Sociales

Almacenamiento

Base de Datos Documental

Recuperación - Lenguaje Query

NoSQL

ESTRUCTURADA

1001 1010	1001 0101	1100 0110
0011 1100	0110 1001	0011 1010
0011 0011	0101 1100	1001 1001

Tipos de Datos

Tabulares

Bases de Datos

Google Sheets

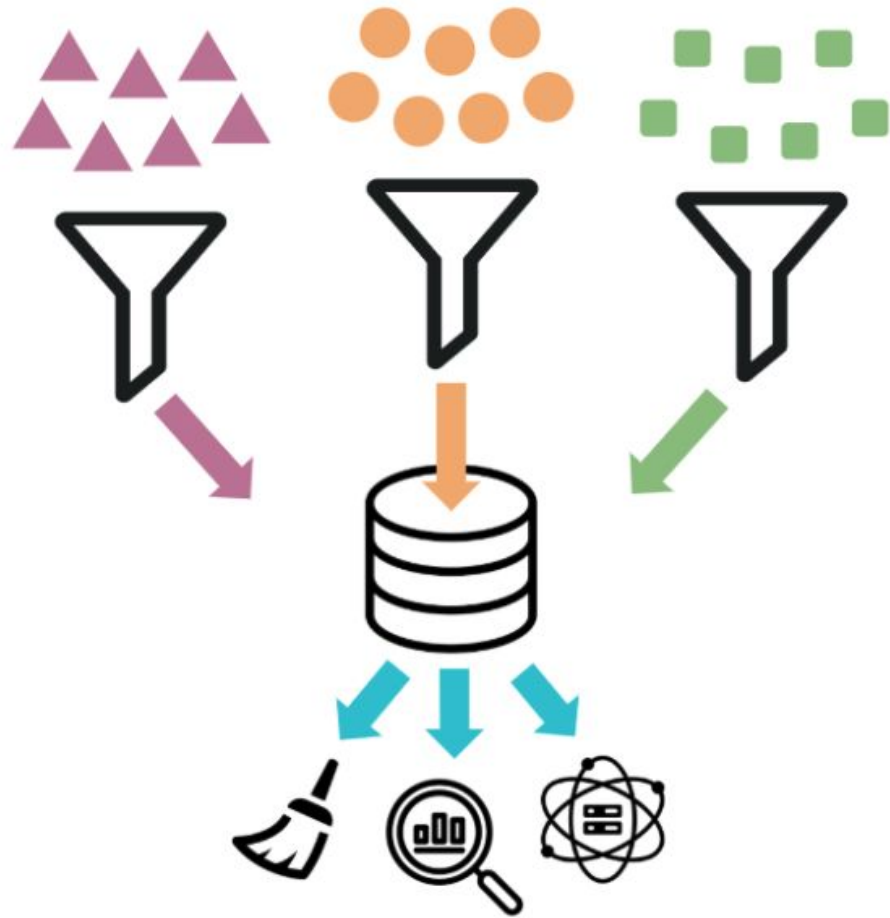
Archivos Excel

Almacenamiento

Base de datos relacionales

Recuperación - Lenguaje Query

SQL (Lenguaje de Consulta Estructurada)



¿Cómo escalamos?

Más de una fuente de datos:

- Bases de Datos
- APIs
- Registros Públicos

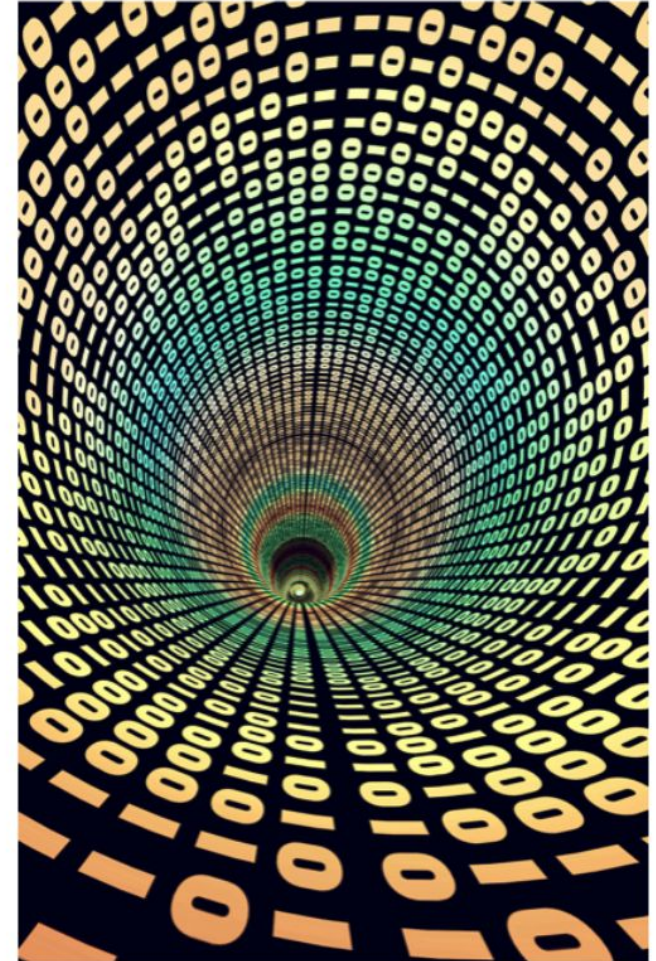
Diferentes tipos de datos:

- Datos estructurados
- Datos Tabulares
- Datos transmisión en tiempo real. Ejemplo: Tweets



Qué es un PIPELINE DE DATOS?

- Movilidad de datos en etapas definidas
- Recolección y almacenamiento de datos automatizado
 - Programados por hora, diariamente, semanalmente, mensualmente, etc
 - Desencadenado por eventos
- Registros Públicos
- Necesarios para grandes proyectos de datos
- Extract Transform Load (ETL)

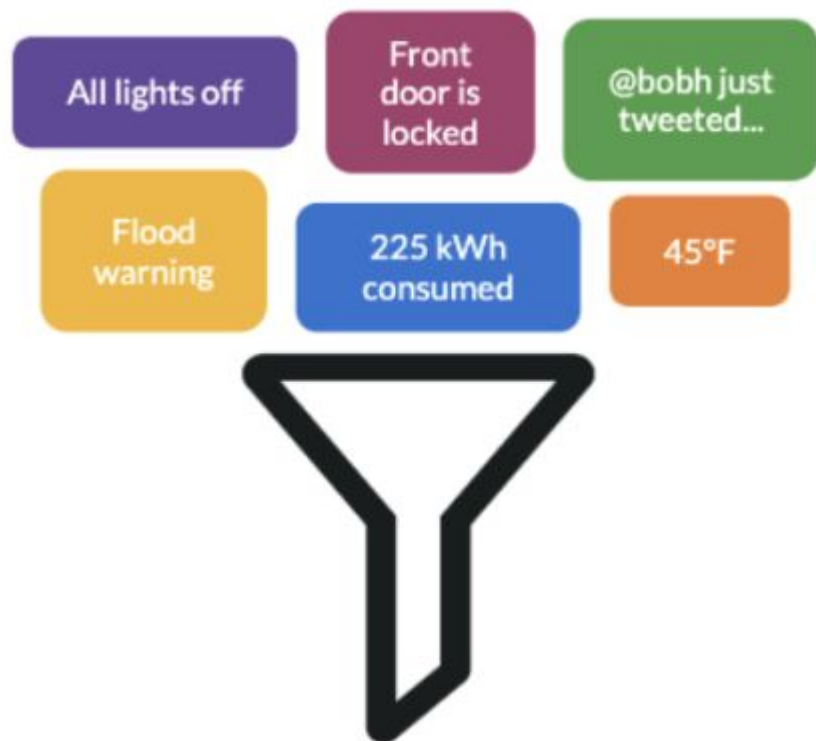


CASO DE ESTUDIO



Dato	Fuente	Frecuencia
Condiciones Ambientales	National Weather Service API	Cada 30 mins
Tweets en tu área	Twitter API	Tiempo-real
Temperatura (Interna)	Smart home thermostat	Cada 5 mins
Estado de luces	Smart light bulbs	Cada 1 min
Estado de bloqueos	Smart door locks	Cada 15 segundos
Consumo de energía	Smart meter	Semanalmente

EXTRACCIÓN

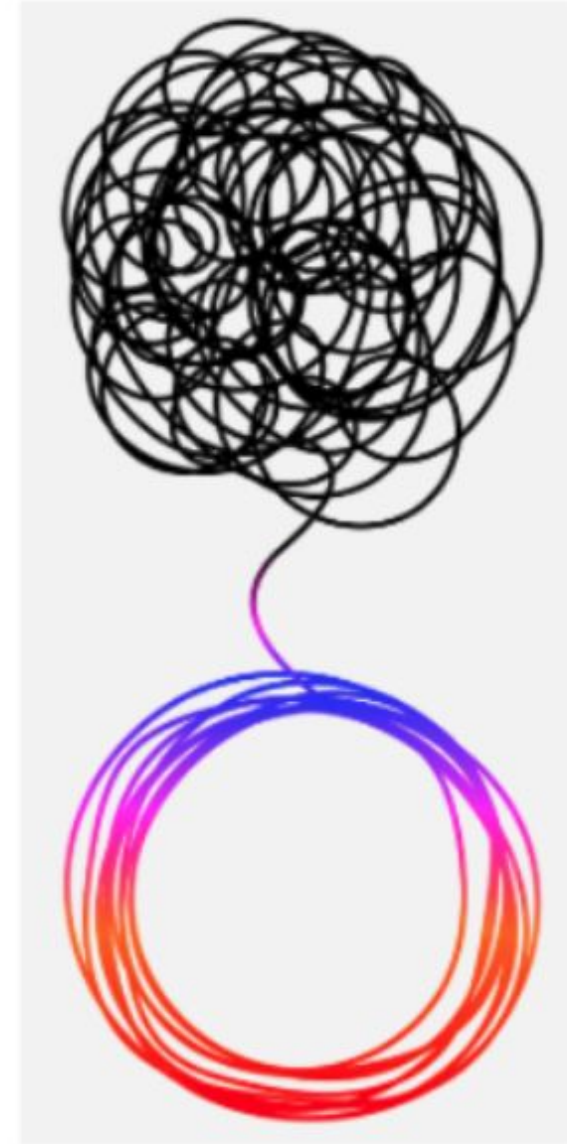


Dato	Fuente	Frecuencia
Condiciones Ambientales	National Weather Service API	Cada 30 mins
Tweets en tu área	Twitter API	Tiempo-real
Temperatura (Interna)	Smart home thermostat	Cada 5 mins
Estado de luces	Smart light bulbs	Cada 1 min
Estado de bloqueos	Smart door locks	Cada 15 segundos
Consumo de energía	Smart meter	Semanalmente

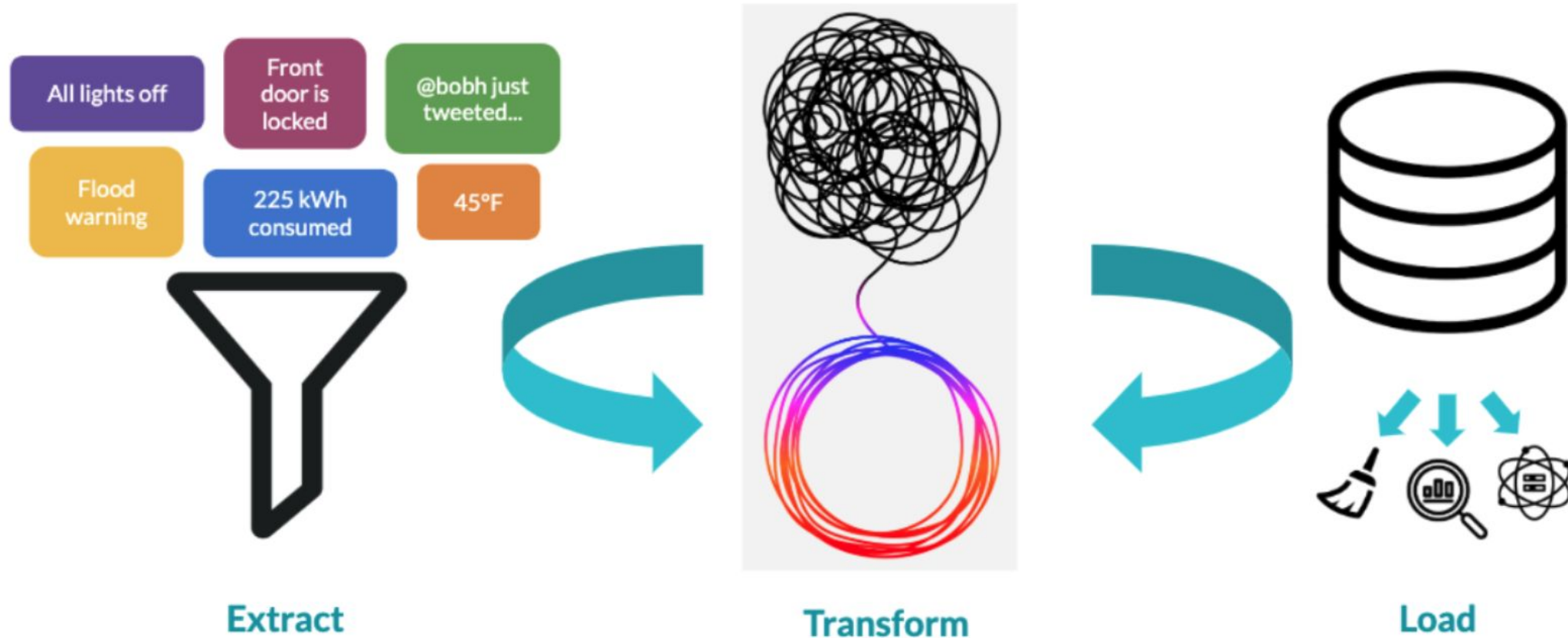


TRANSFORMACIÓN

- Unir diferentes fuentes de datos en un solo conjunto de datos
- Convertir la estructura de datos para que se ajuste al esquema de las bases de datos
 - Ejemplo id_cliente es char y queremos que sea numérico
- Removiendo datos irrelevantes
- La preparación de datos y la exploración de los mismos no ocurren hasta esta fase.



CARGA Y AUTOMATIZACIÓN

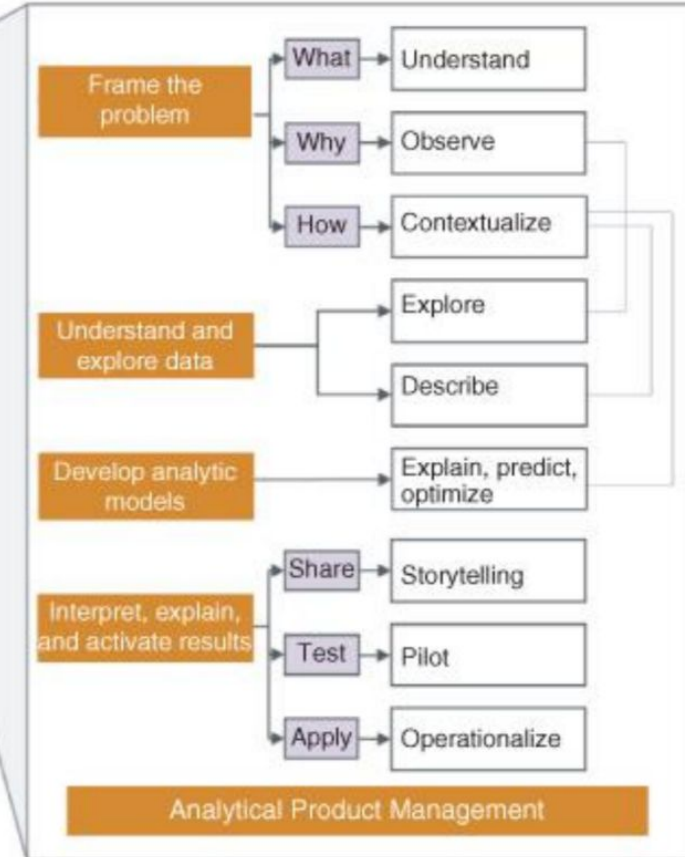


Data Pipeline

Informs requirements for

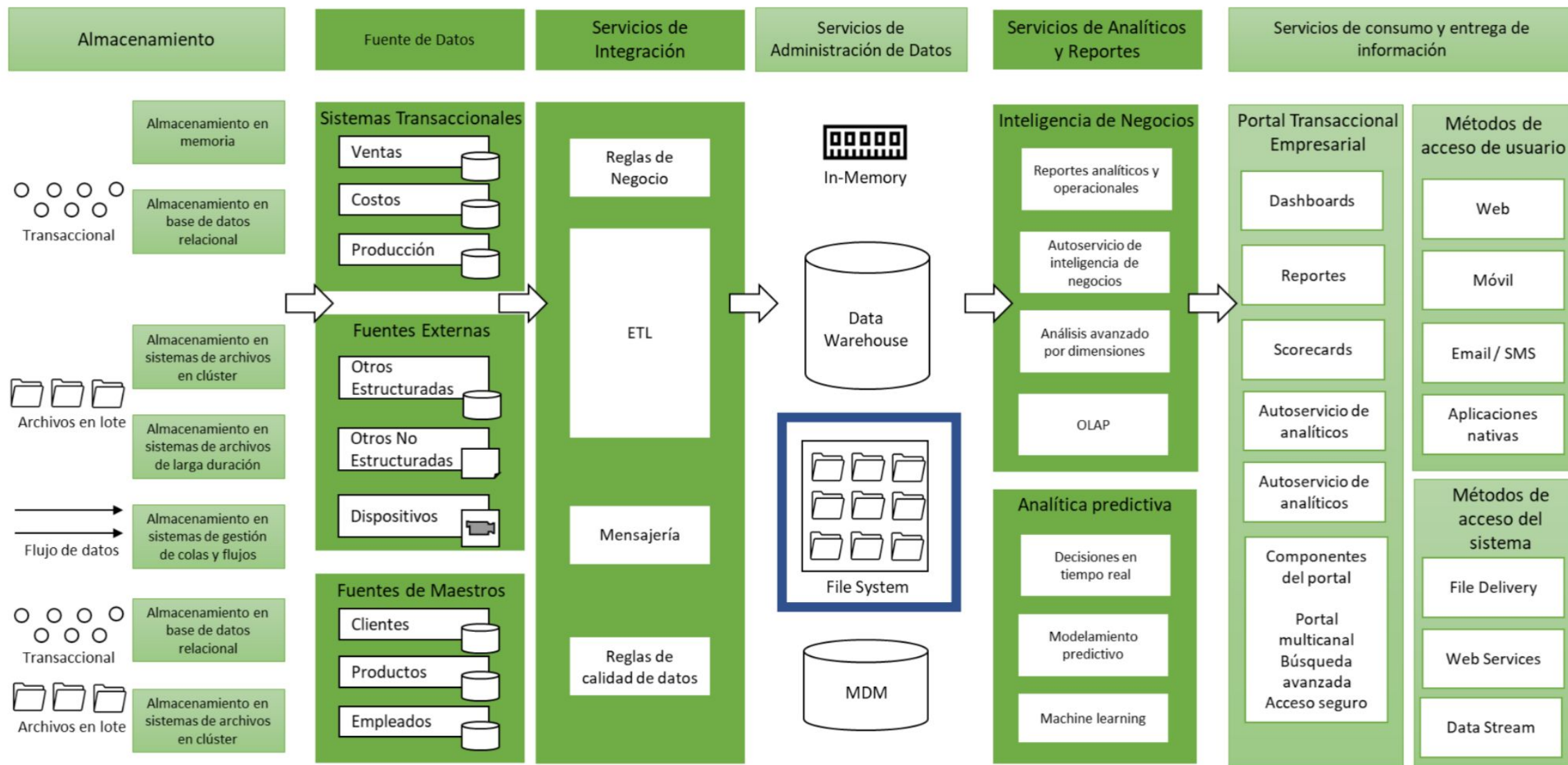
Supports

Analytics Lifecycle



Tomado de: Gregory S. Nelson. The Analytics Lifecycle Toolkit. 2018.







IDE's para el desarrollo de Python

Permite trabajar en un entorno no local y la creación de Notebooks 🚀

Las herramientas que mostramos anteriormente no son las únicas en donde compilar código de Python...



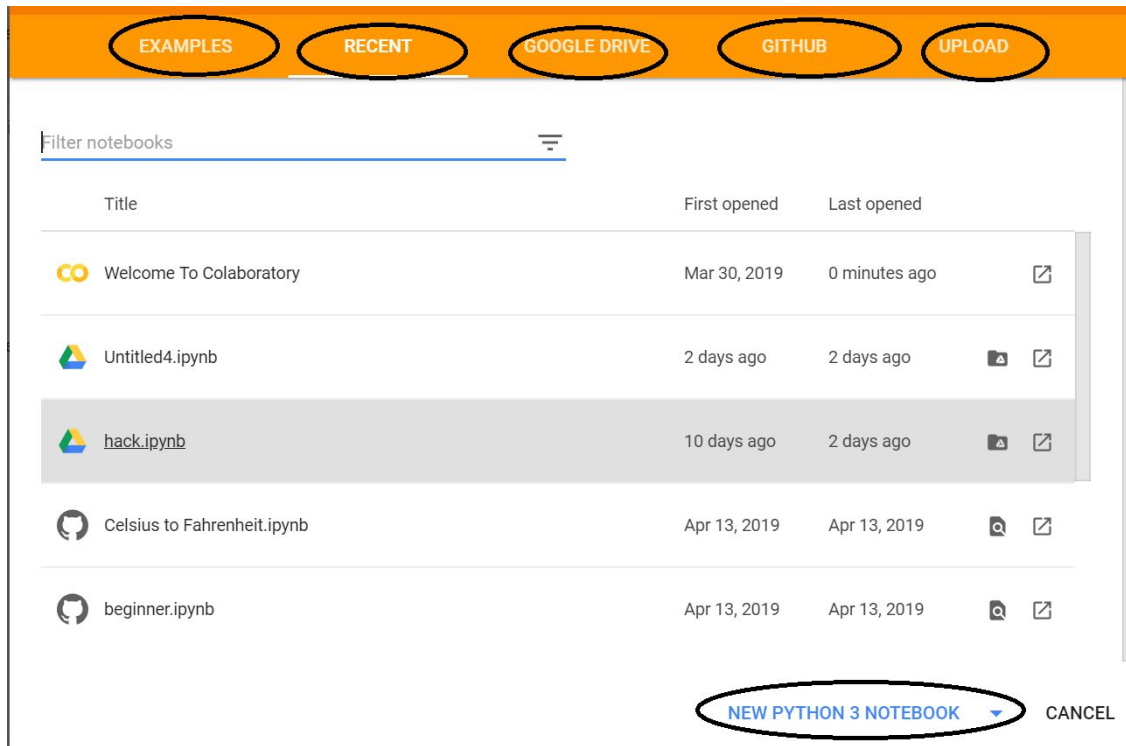
Google Colab

Permite trabajar en un entorno no local y la creación de Notebooks 🚀

- ✓ Es un producto de Google Research. Está especialmente adecuado para tareas de aprendizaje automático, análisis de datos y educación.
- ✓ Jupyter es el proyecto de código abierto en el que se basa Colab.
- ✓ Nos permite compartir notebooks sin la necesidad de descargar ningún software extra.
- ✓ El código se ejecuta en una máquina virtual dedicada a tu cuenta y pueden eliminarse luego de cierto tiempo.



Cómo usar Google Colab



EXAMPLES: Contiene ejemplos de Jupyter notebooks con diversos ejemplos.

RECENT: Jupyter notebooks que has trabajado recientemente.

GOOGLE DRIVE: Jupyter notebooks en tu google drive.

GITHUB: Puedes añadir Jupyter notebooks desde Github pero es necesario conectar Colab con GitHub.

UPLOAD: Si deseas subir un Jupyter notebook desde tu equipo local.

Ir al siguiente enlace: <https://colab.research.google.com>



VAMOS A APLICAR !!!

https://colab.research.google.com/drive/1FC7wngzYTGcO_n_Lv4RB3rHU0r2BVPKm?usp=share_link

<https://www.kaggle.com/datasets/manovirat/groceries-sales-data>





Gracias