

# CENTRO DEL PROFESORADO DE MÁLAGA

## Introducción al Big Data y Business Intelligence

Nombre: José Manuel García Nieto

Fecha: 6 de abril de 2017. 16:30 a 20:30



JOSÉ MANUEL GARCÍA NIETO  
WEB PERSONAL



**MENU**

- Inicio
- Investigación
- Docencia
- Publicaciones
- Software
- Enlaces



**José Manuel García Nieto**

- Doctor Ingeniero en Informática  
PhD in Computer Science
- Ingeniero en Informática  
- Ingeniero Técnico en Informática de Gestión

**Departamento de Lenguajes y Ciencias de la Computación**

- E.T.S.I. Informática, Universidad de Málaga
- Ampliación Campus de Teatinos
- 29071 Málaga-España
- C/ Doctor Ortiz Ramos. Ed. Ada Byron, A.2.1.
- Tfn:+34 951 952 924
- @: [jnieto@lcc.uma.es](mailto:jnieto@lcc.uma.es)
- Grupo GISUM

**ResearcherID**

Click here to see my profile

Ver mi perfil en **LinkedIn**

- Tesis Doctoral. *Emergent Optimization: Design and Applications in Telecommunications and Bioinformatics*. (RD 1393/2007) Ingeniería del Software e Inteligencia Artificial. Universidad de Málaga (Spain). [\[Memoria-PDF\]](#) [\[Presentación-PDF\]](#)
- Máster de Postgrado. Ingeniería del Software e Inteligencia Artificial. Universidad de Málaga (Spain).
- Proyecto de Máster: *"Algoritmos Basados en Inteligencia Colectiva para la Resolución de Problemas de Bioinformática y Telecomunicaciones"*. Aplicación de algoritmos de *Particle Swarm Optimization* a problemas del campo de las Telecomunicaciones (Location Area Management) y del campo de la Bioinformática (Selección y Clasificación de Genes en Microarrays de ADN) [\[Proyecto-PDF\]](#) [\[Memoria-PDF\]](#)
- Ingeniero en Informática. E.T.S. de Ingenieros en Informática. Universidad de Málaga (Spain).
- Proyecto Fin de Carrera (Ing. Informática): *"Algoritmos Basados en Cúmulos de Partículas para la Resolución de Problemas Complejos"*. Aplicación de algoritmos de *Particle Swarm Optimization* a problemas del campo de las Telecomunicaciones (Location Area Management) y del campo de la Bioinformática (Ordenación de genes en Microarrays de ADN) [\[PDF\]](#)
- Ingeniero Técnico en Informática de Gestión. E.T.S. de Ingenieros en Informática. Universidad de Málaga (Spain).

**7 Pageviews**  
Apr 1st - Apr 30th



[Click to See Details](#)



Última Actualización 18/10/2007. Visitor number 11,907



## Esquema de Contenidos:

1. Introducción al Big Dada
2. Trabajando con los datos
3. Almacenamiento y Procesamiento: Técnicas, Herramientas y Plataformas
4. Análisis mediante la Minería de Datos
5. Visualización y Consumo de Datos
6. Seguridad y Gobernanza
7. Aplicaciones Reales de Negocio: Casos de Éxito

# Presentación: Grupo Khaos – Universidad de Málaga

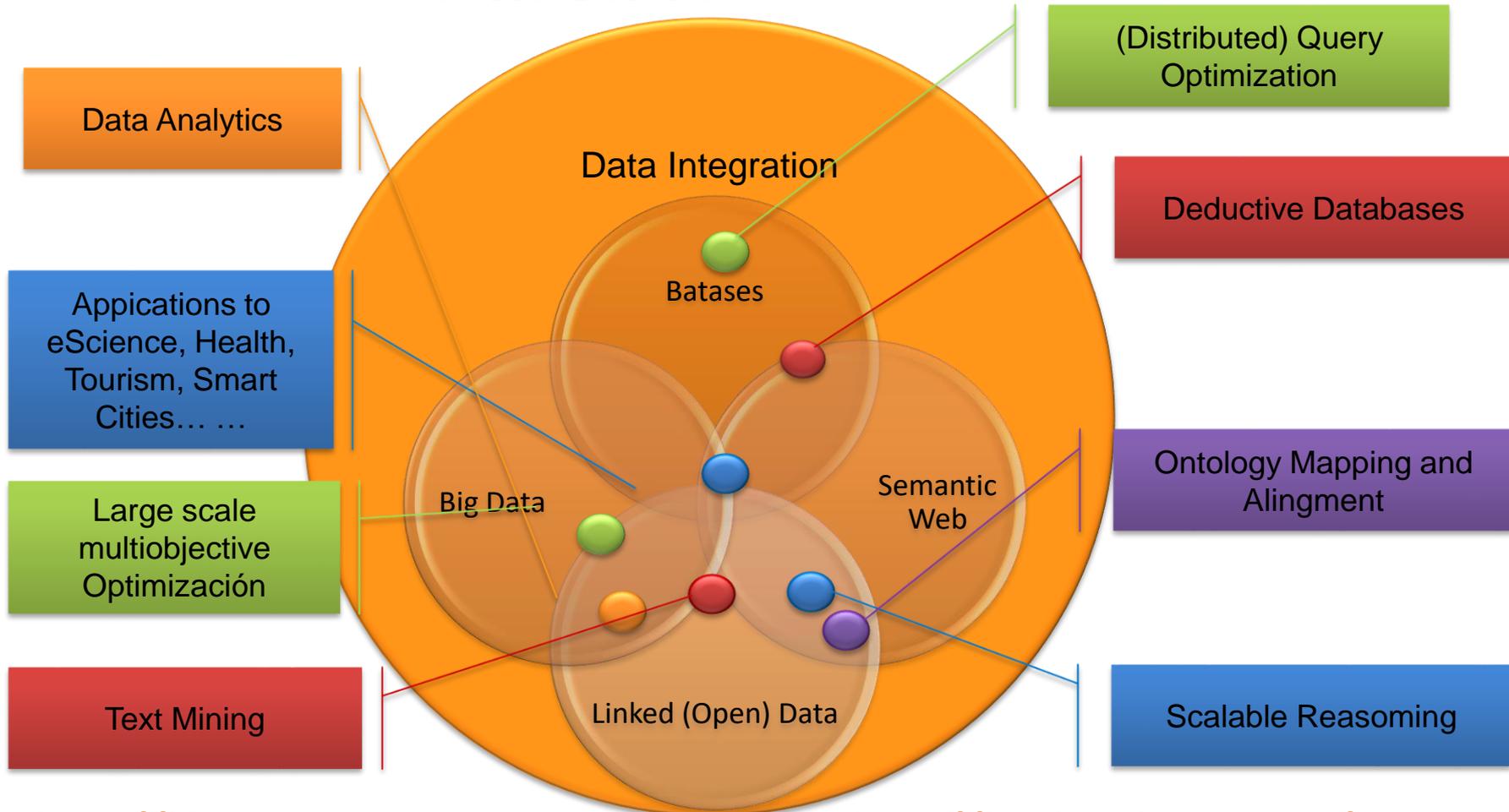


Ada Byron  
Research  
Center

<http://khaos.uma.es>

<http://bigdata.lcc.uma.es/>

# Presentación: Grupo Khaos – Universidad de Málaga



<http://khaos.uma.es>

<http://bigdata.lcc.uma.es/>

• Líneas de trabajo:

- **Data management and Integration** (Relational Databases, NoSQL databases, Linked Data, Open Data)
- **Data Analytics**: Data minig, Text mining, (NLP&Semantics, social networks analytics, short texts,...), Information retrieval, Big Data Analytics. Semantic annotation and enrichment. Integrated-Data Analytics.
- **Semantic Technologies** (Scalable Reasoning, Schema Mapping and Alingment, Semantic Assigation to Mobile Objects, Semantic base Content Recommendation, Semantics for eScience, Smart Data for Big Data interpretation)
- **Metaheuristics**. Mono and multi-Objective Optimization (scalable, high throughput, Globus, MapReduce, SPARK, Storm). Metaheuristics and new semantic aware Machine Learning algorithms (based on Mahout)

• Application domains:

- Health: **Bioinformatics, Translational Medicine, System Medicine. Also Electronic Health Records and Medical Decision Support Systems.**
- **eCommerce and Retailing**
- **Cultural Heritage and Tourism**
- **Smart Cities**

# Máster Propio Universitario en Advanced Analytics on **Big Data**



Con aplicaciones reales en Smart Cities y movilidad urbana, Smart Grids y eficiencia energética y eHealth y medicina de precisión

## Contenidos

Hadoop, Spark, SQL, MongoDB, Excel, Hive, Open Data, D3JS, análisis predictivo, OLAP, Visualización, Smart Grids, análisis de textos, Internet of Things, Cloud Computing, Seguridad

## Características

- o 90 créditos ECTS
- o Semipresencial
- o Precio: 6000€
- o Docentes universidad y empresa
- o Prácticas remuneradas en empresas (6 meses, 1000 euros/mes)

**Inicio: octubre 2016**

## Información y contacto

Grupo de Investigación Khaos  
Teléfono: 951 952 918  
Email: [bigdata@lcc.uma.es](mailto:bigdata@lcc.uma.es)  
<http://bigdata.lcc.uma.es>

<http://bigdata.lcc.uma.es/>

## Empresas colaboradoras en el programa de prácticas



## Organizan



## Entidades colaboradoras





## LA COLECCIÓN BIG DATA Y LOS NEGOCIOS

- Libro 1: Introducción al Big Data
- Libro 2: Big Data: gestión inteligente de los datos
- Libro 3: Introducción al trabajo con datos
- Libro 4: Modelado en Big Data
- Libro 5: Almacenamiento de Big Data.
- Libro 6: Procesamiento y Análisis Inteligente de Big Data
- Libro 7: Visualización y consumo de Big Data
- Libro 8: Big Dada: seguridad y gobernanza.





# Introducción al Big Data

## ¿Por qué formarse en Big Data?

### OFERTAS DE EMPLEO

#### RELACIONADAS CON BIG DATA

¿Para qué es esto?

¿eres una empresa?  
PUBLICA UNA OFERTA GRATIS

¿por qué publicar una oferta?

356 visitas

#### Desarrolladores Big Data (Hadoop – Spark)

STRATIO

En Stratio buscamos desarrolladores Big Data para incorporarse en proyectos punteros e interesantes. Con experiencia como desarrollador Java y/o Sacala. Conocimiento de las principales

Localización: [Madrid, España](#)

Publicadas: 19/01/2015

154 visitas

#### Arquitecto JEE Big Data

STRATIO

En Stratio buscamos arquitectos JEE Big Data con al menos 4 años de experiencia (Spring, Struts, Hibernate) y con conocimientos en Tecnologías Big Data (Cassandra, Hadoop, Spark, Hive, Pig...)

Localización: [Madrid, España](#)

Publicadas: 19/01/2015

305 visitas

#### Formador Big Data

CULTURE LAB TS S.L.

Culture Lab es un centro de formación referente en el entorno de las tecnologías TIC para empresas y particulares. Actualmente estamos interesados en aumentar nuestra plantilla de especialistas en el área

Localización: [Madrid, España](#)

Publicadas: 11/01/2015

279 visitas

#### Arquitectos y Desarrolladores Big Data

ANYHELP INT.

Buscamos un Arquitecto para liderar el diseño, implementación y mantenimiento de un ecosistema Big Data en un apasionante proyecto en el sector banca. Será necesaria experiencia previa en

Localización: [Madrid, España](#)

Publicadas: 08/01/2015

322 visitas

#### Big Data Analyst

ALLMINDS EXECUTIVE SEARCH.

Buscamos para empresa del sector Comunicación Digital, ubicada en Barcelona, un Analista Big Data. Se requiere experiencia en análisis digital, con una experiencia de 2-3 años mínimo. Conocimientos en

Localización: [Barcelona, España](#)

Publicadas: 23/12/2014

1152 visitas

#### BIG DATA MEANS NOTHING WITHOUT YOU

STRATIO

Big data architects, Data analysts and people who love to get involved in a project where learning is the first. We are willing to hear about you. Currently in Stratio (an IT company specialized in

Localización: [Madrid, España](#)

Publicadas: 17/12/2014

178 visitas

#### Arquitecto/a Big Data

RAY HUMAN CAPITAL

Proyectos con soluciones Big. Buscamos incorporar en Madrid a un Arquitecto Big. Data, especialmente en soluciones open source y el ecosistema Hadoop (Map....

Localización: [Madrid, España](#)

Publicadas: 10/12/2014

200 visitas

#### Ingenieros/as Big Data

RAY HUMAN CAPITAL

Se valorarán conocimientos en tecnologías de Big Data y entornos distribuidos, conocimientos en entornos y tecnologías de Cloud Computing, experiencia en...

Localización: [Madrid, España](#)

Publicadas: 10/12/2014

## La Explosión de Datos

La creación de datos se está produciendo a un ritmo record:

hasta el 2010 se generaron en el mundo 1ZB de datos

y

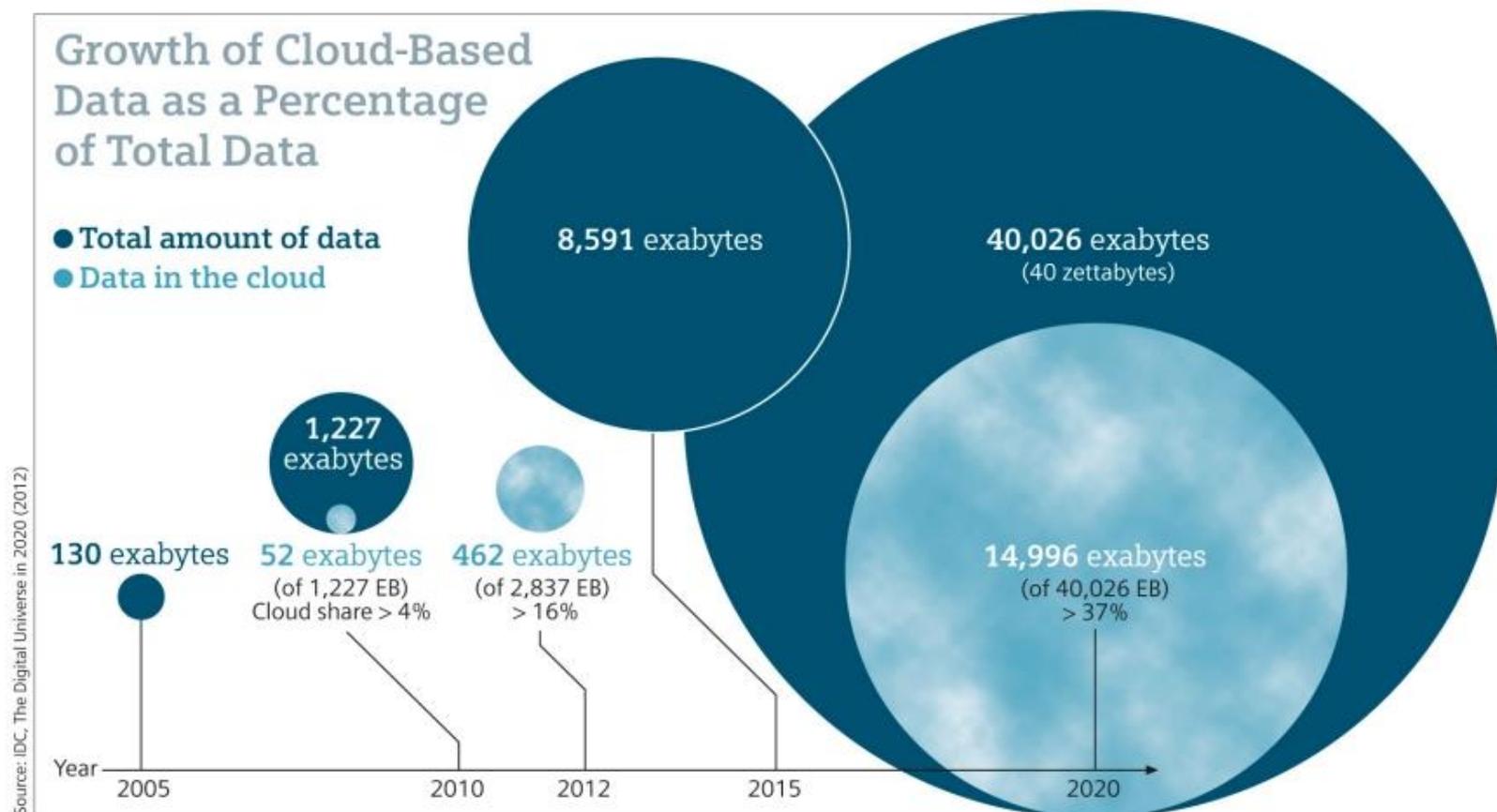
en 2014, se generaron más de 7ZB en un año

### Unidades básicas de información

#### Prefijos del Sistema Internacional

Múltiplo - (Símbolo)	<u>SI</u>	<u>Binario</u>
<u>kilobyte</u> (kB)	$10^3$	$2^{10}$
<u>megabyte</u> (MB)	$10^6$	$2^{20}$
<u>gigabyte</u> (GB)	$10^9$	$2^{30}$
<u>terabyte</u> (TB)	$10^{12}$	$2^{40}$
<u>petabyte</u> (PB)	$10^{15}$	$2^{50}$
<u>exabyte</u> (EB)	$10^{18}$	$2^{60}$
<b>zettabyte</b> (ZB)	$10^{21}$	$2^{70}$
<u>yottabyte</u> (YB)	$10^{24}$	$2^{80}$

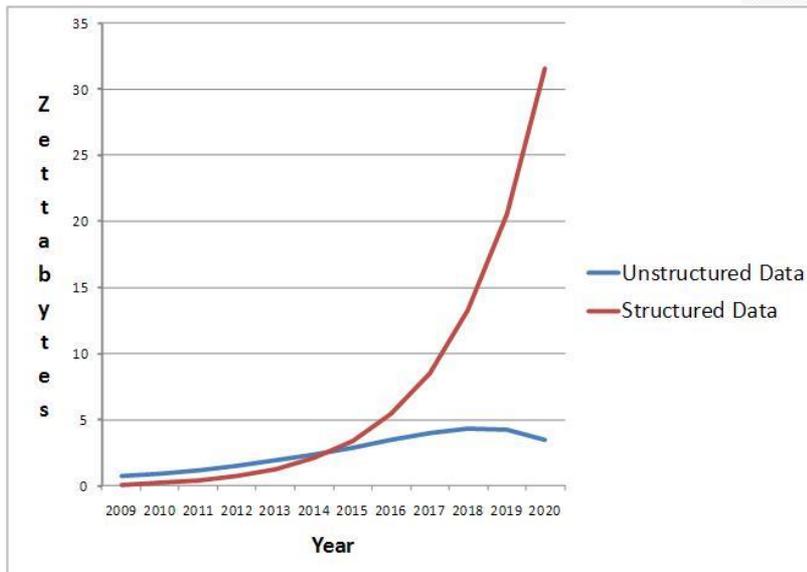
## La Explosión de Datos



## La Explosión de Datos



## ESTAMOS EN LA ERA BIG DATA



### Oracle and Big Data

#### Big Data for the Enterprise

The term big data draws a lot of attention, but behind the hype there's a simple story. For decades, companies have been making business decisions based on transactional data stored in relational databases. Beyond that critical data, however, is a potential treasure trove of less structured data: weblogs, social media, email, sensors, and photographs that can be mined for useful information.

que a los clientes les resulte fácil y rentable procesar y extraer información de volúmenes masivos de datos.

**BigQuery** enables businesses and developers to gain real-time business insights from massive amounts of data without any upfront hardware or software investments. Imagine a big pharmaceutical company optimizing daily marketing spend using worldwide sales and advertisement data. Or think of a small online retailer that makes product recommendations based on user clicks. Today, we are making BigQuery publicly available, an important milestone in our effort to bring Big Data analytics to all businesses via the cloud.

enrich your information by connecting to the world's data and services.

#### Key Capabilities

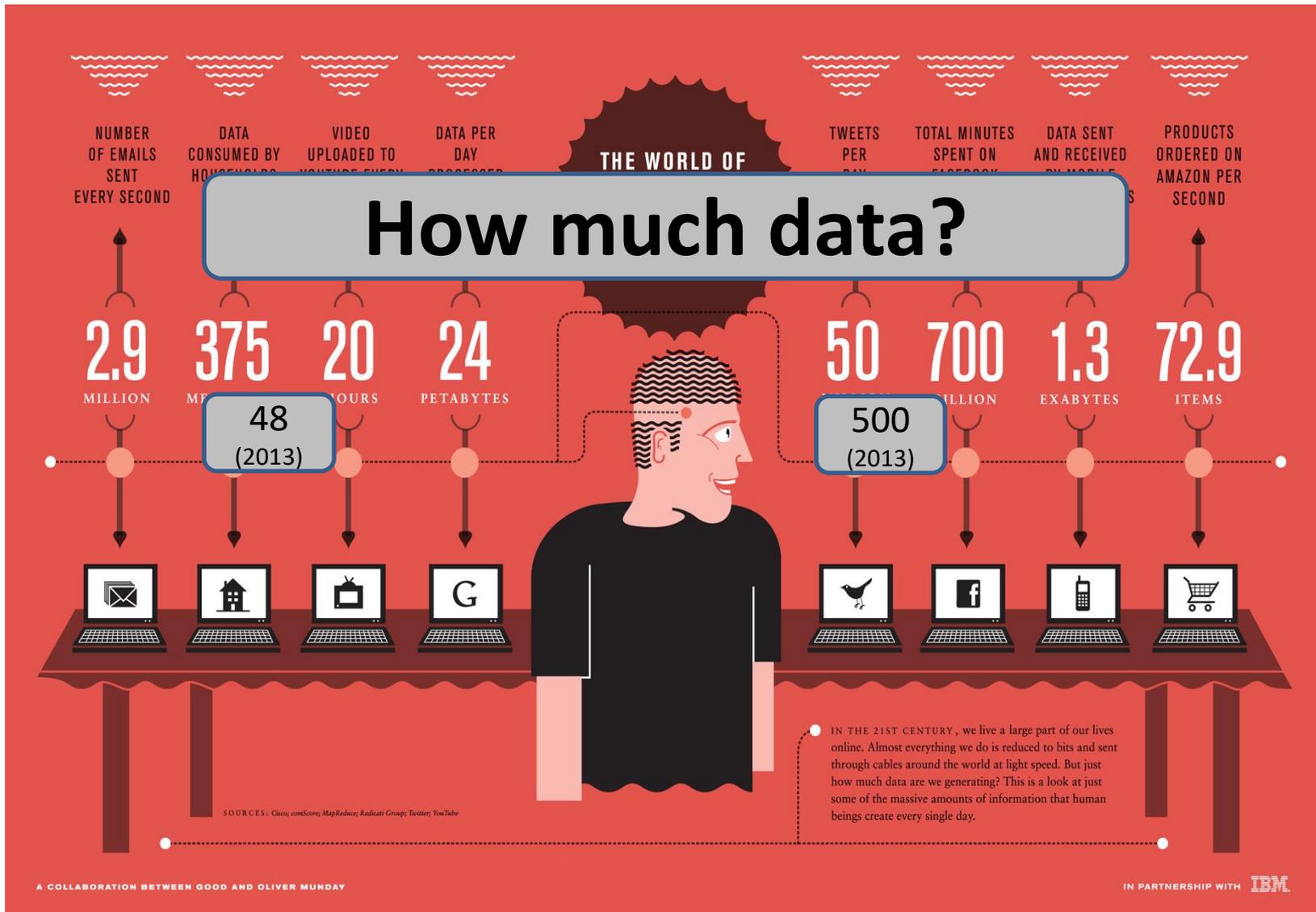
analyze Big Data and gain insights with familiar tools, such as Excel and SharePoint.  
leverage PowerPivot and Power View for Big Data using Excel.  
connect to the world's data and services with Windows Azure Marketplace.  
scale on demand in the private and public cloud.

Every day, we create 2.5 quintillion bytes of data — so much that 90% of the data in the world today has been created in the last two years alone. This data comes from everywhere: sensors used to gather climate information, posts to social media sites, digital pictures and videos, purchase transaction records, and cell phone GPS signals to name a few. This data is **data**.

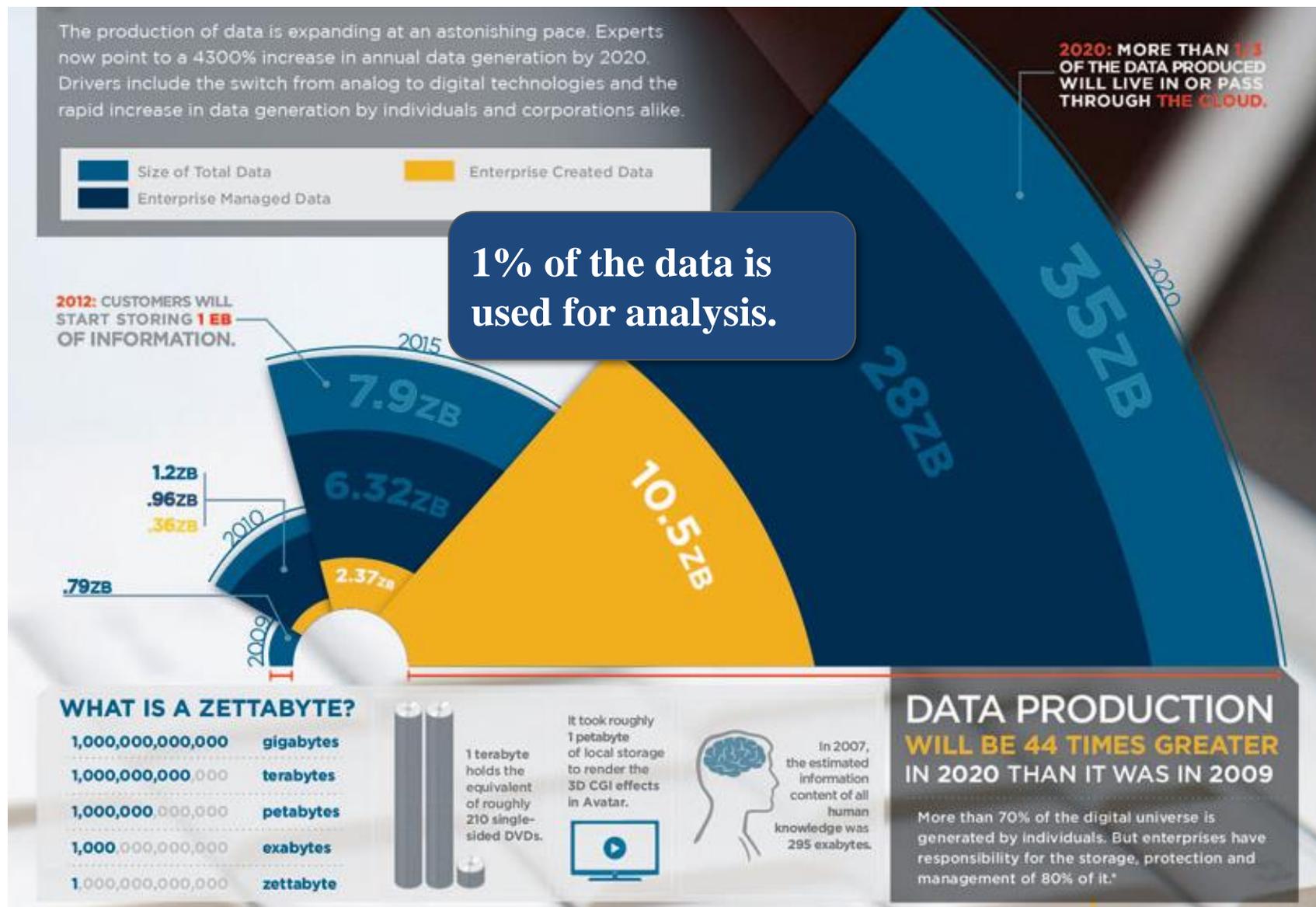
- > Analysis
- > Complex Event Processing
- > Office Based Tools
- > Predictive Analytics
- > Reporting



**data spans three dimensions: Volume, Velocity and Variety.**



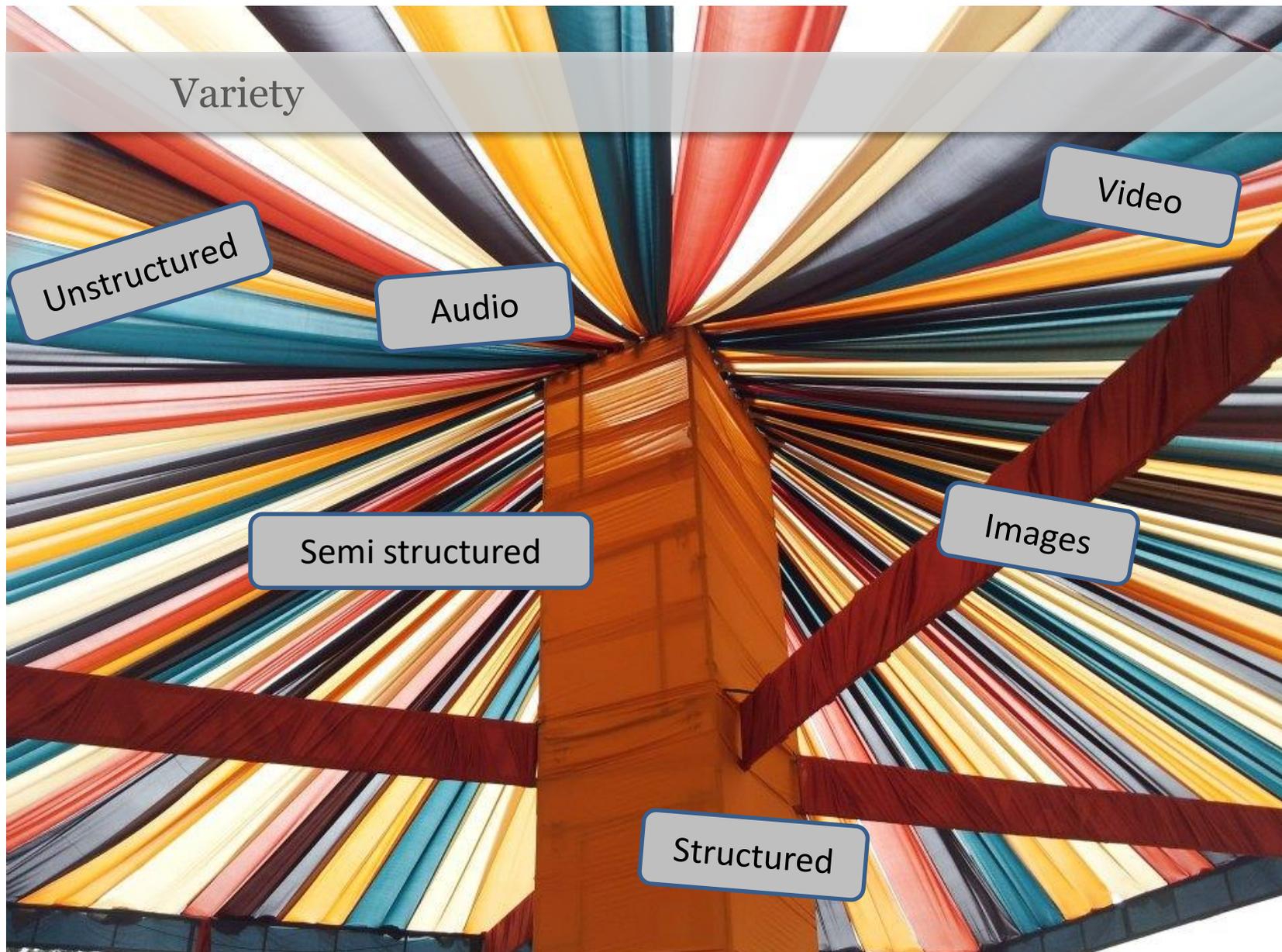
# Introducción al Big Data



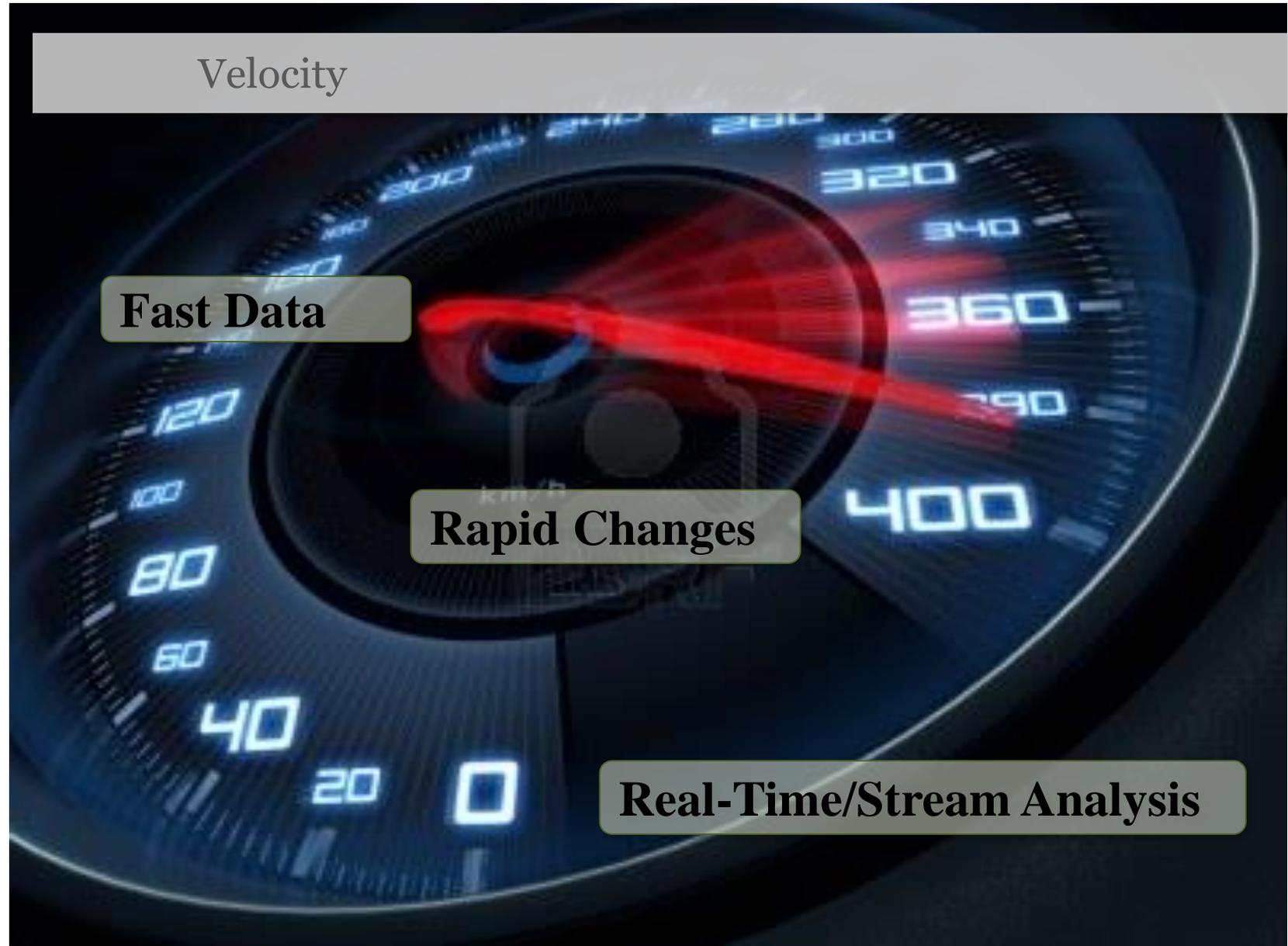
[http://www.csc.com/insights/flxwd/78931-big\\_data\\_growth\\_just\\_beginning\\_to\\_explode](http://www.csc.com/insights/flxwd/78931-big_data_growth_just_beginning_to_explode)

<http://www.guardian.co.uk/news/datablog/2012/dec/19/big-data-study-digital-universe-global-volume>

# Introducción al Big Data



# Introducción al Big Data

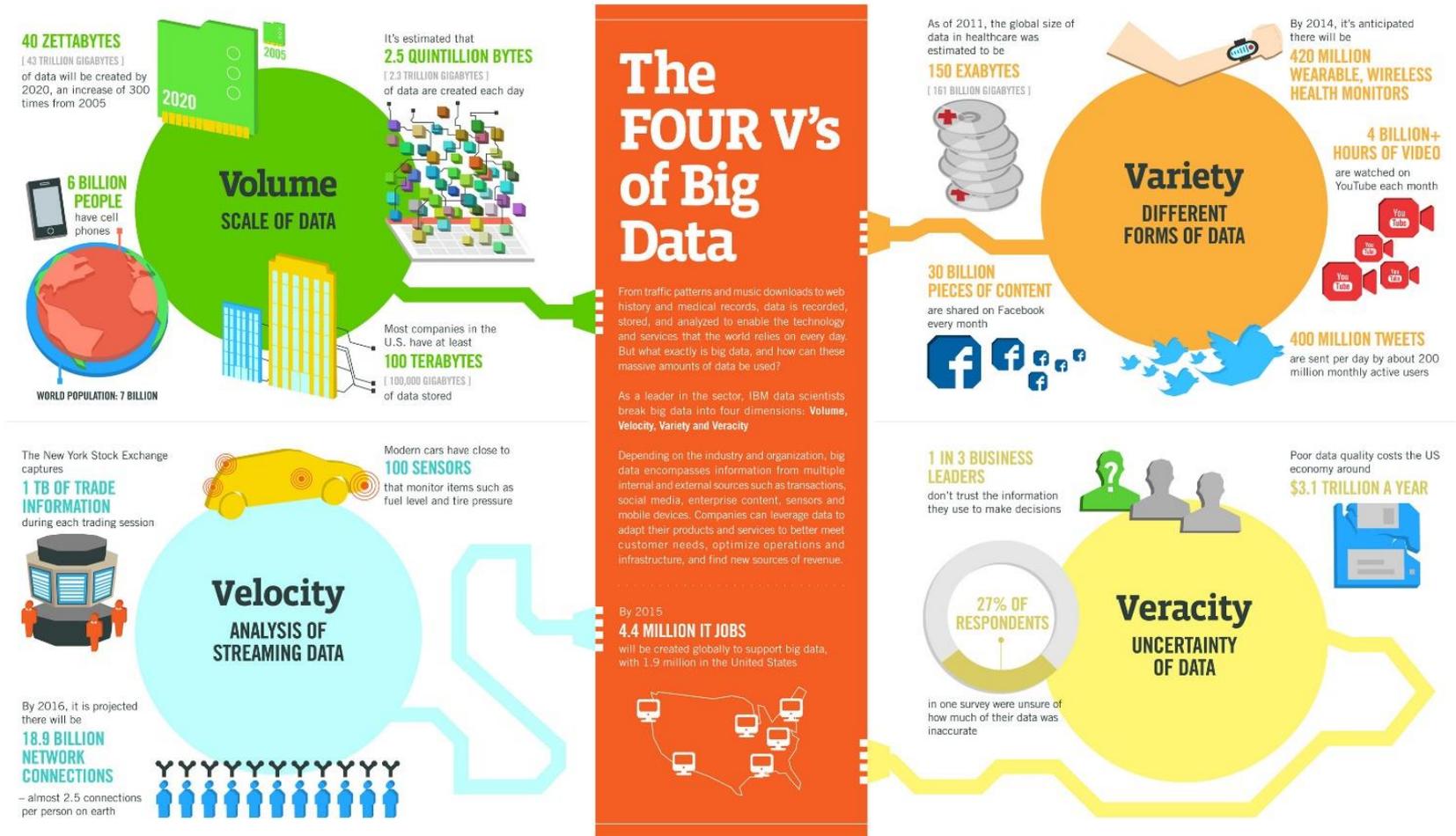


## No sólo cantidad

- El Big Data no es sólo una cuestión de **gestión de grandes cantidades de datos...**
  - ...es también una oportunidad para alcanzar nuevos descubrimientos y conocimiento emergente a partir del **análisis de los datos**

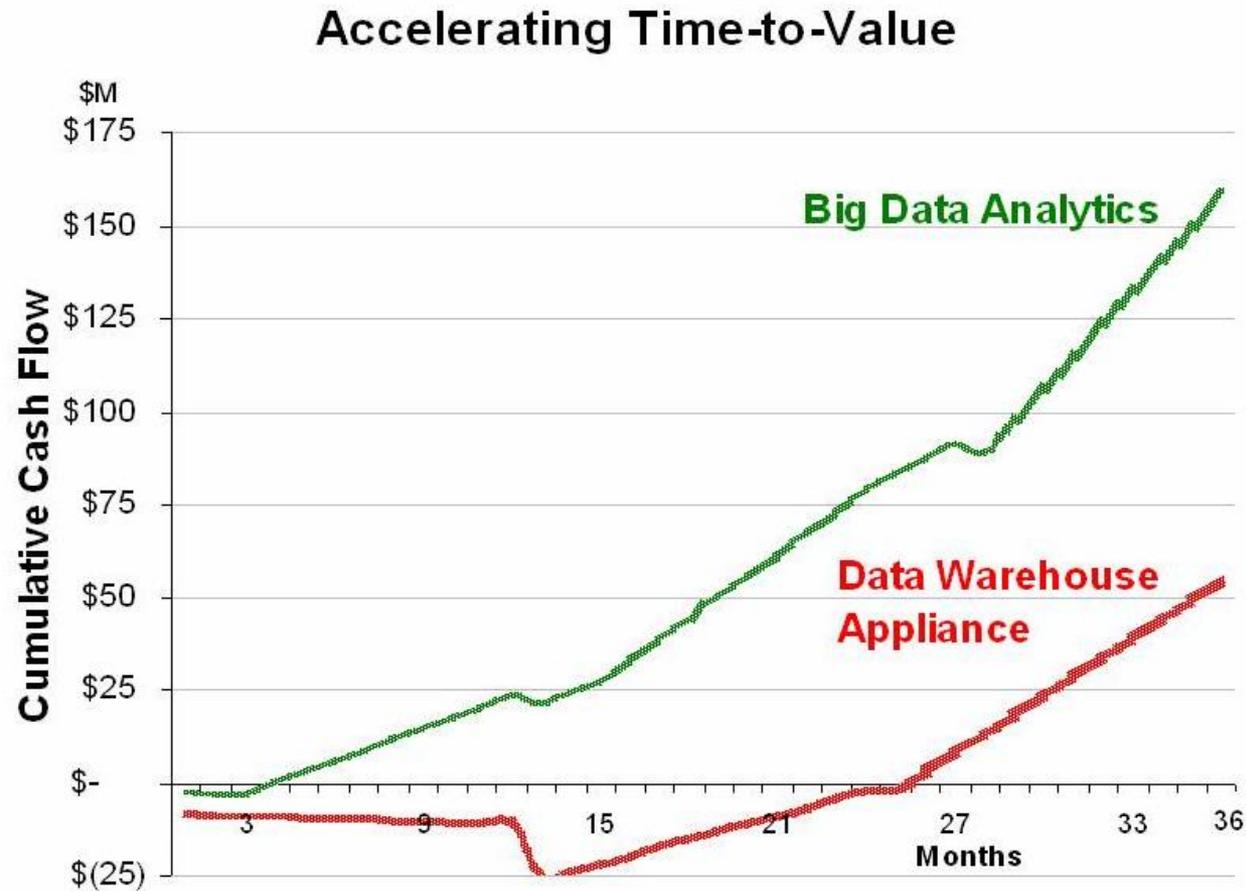
**volumen, velocidad,  
variedad y veracidad**

# Introducción al Big Data

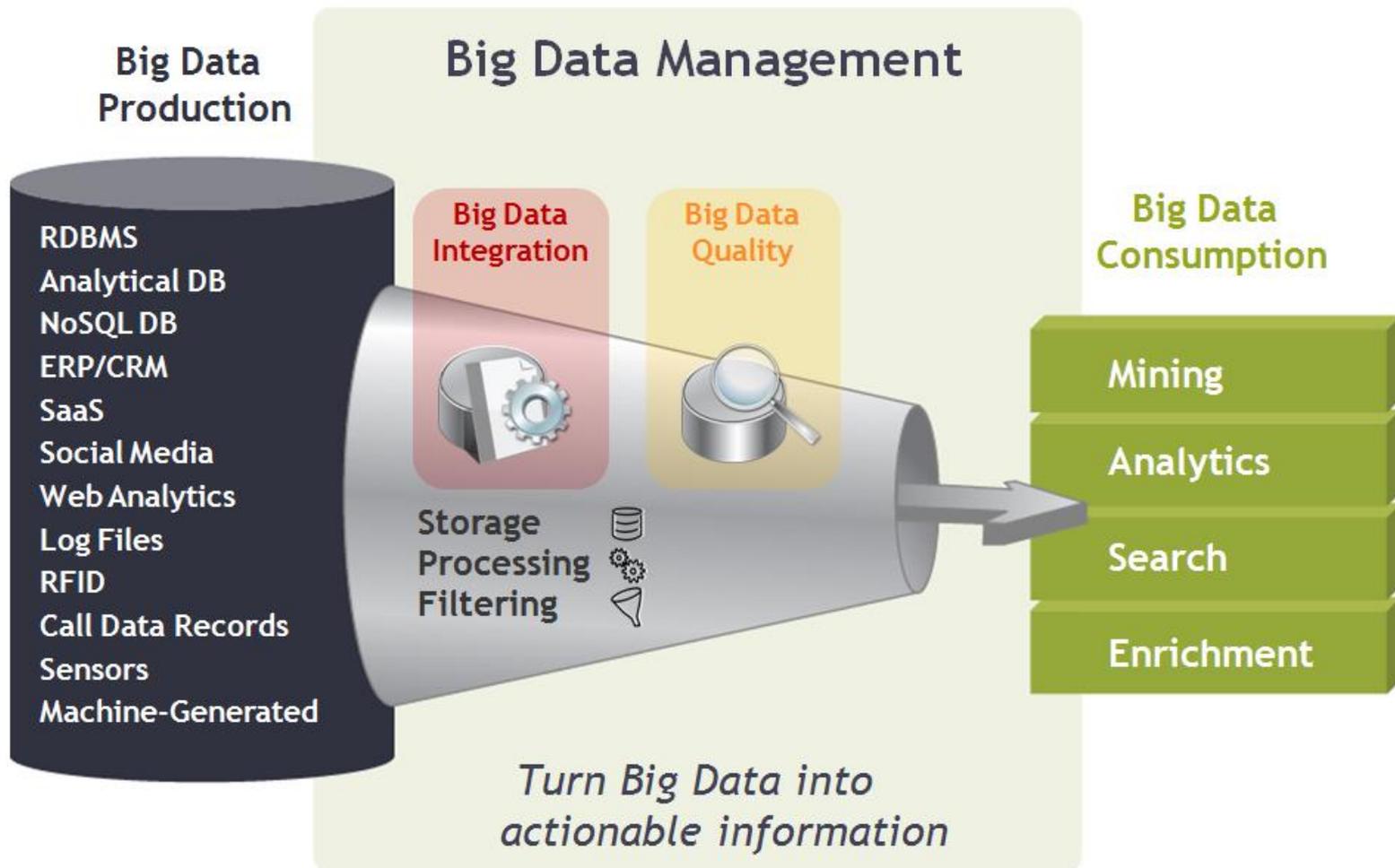


Sources: McKinsey Global Institute, Twitter, Cisco, Gartner, EMC, SAS, IBM, MEPTEC, QAS

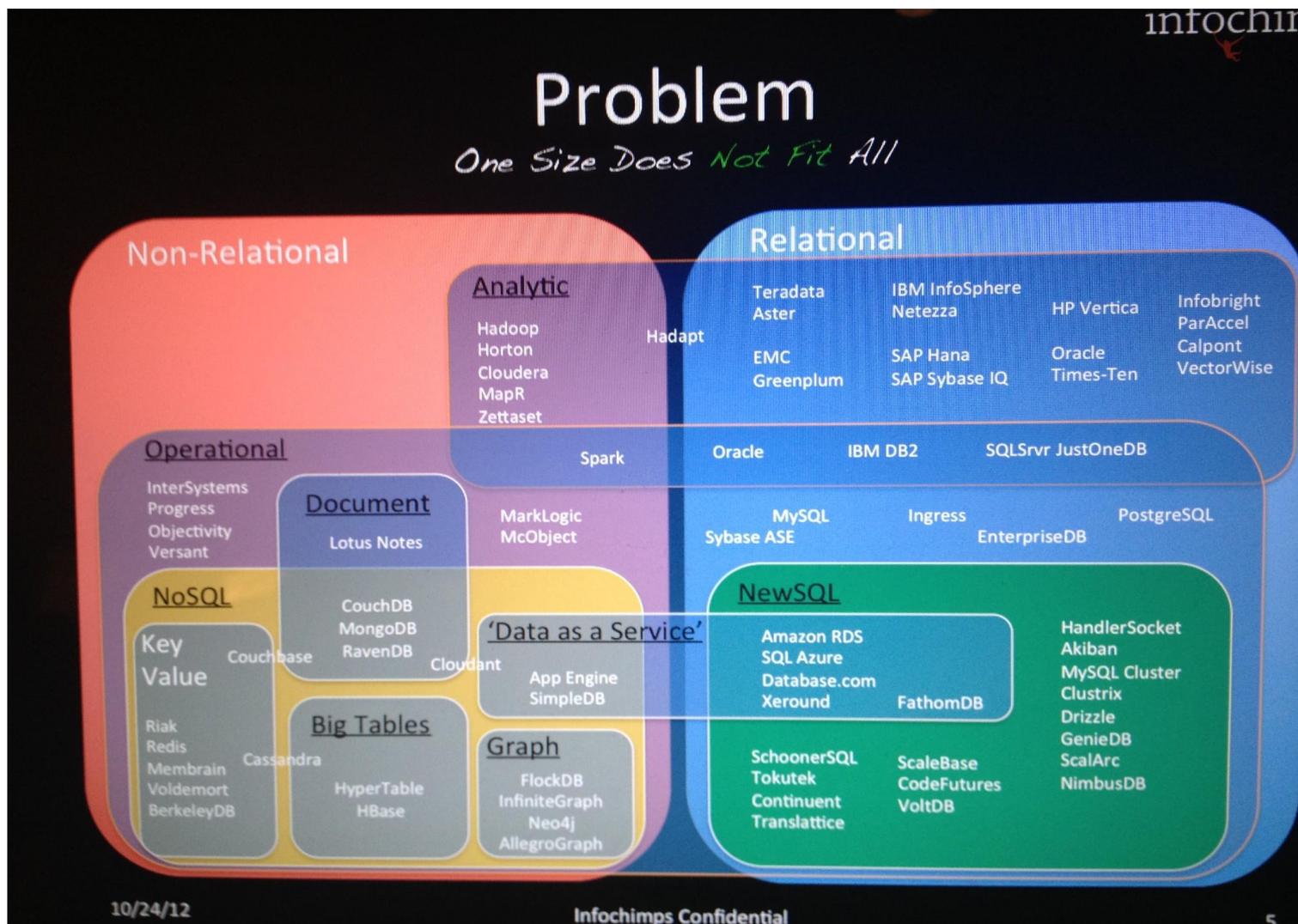
# Introducción al Big Data



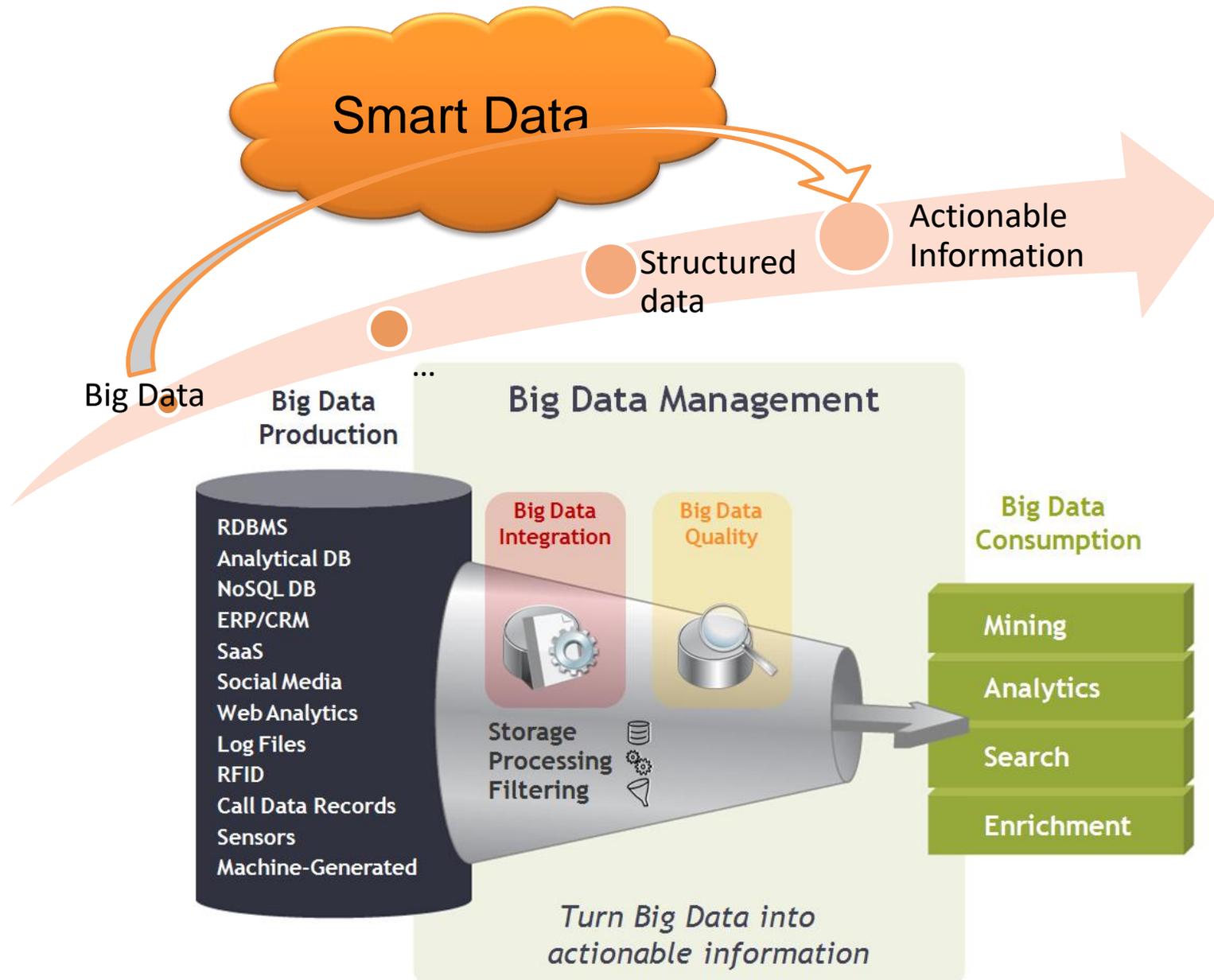
## También Gestión y Análisis



## También Gestión y Análisis

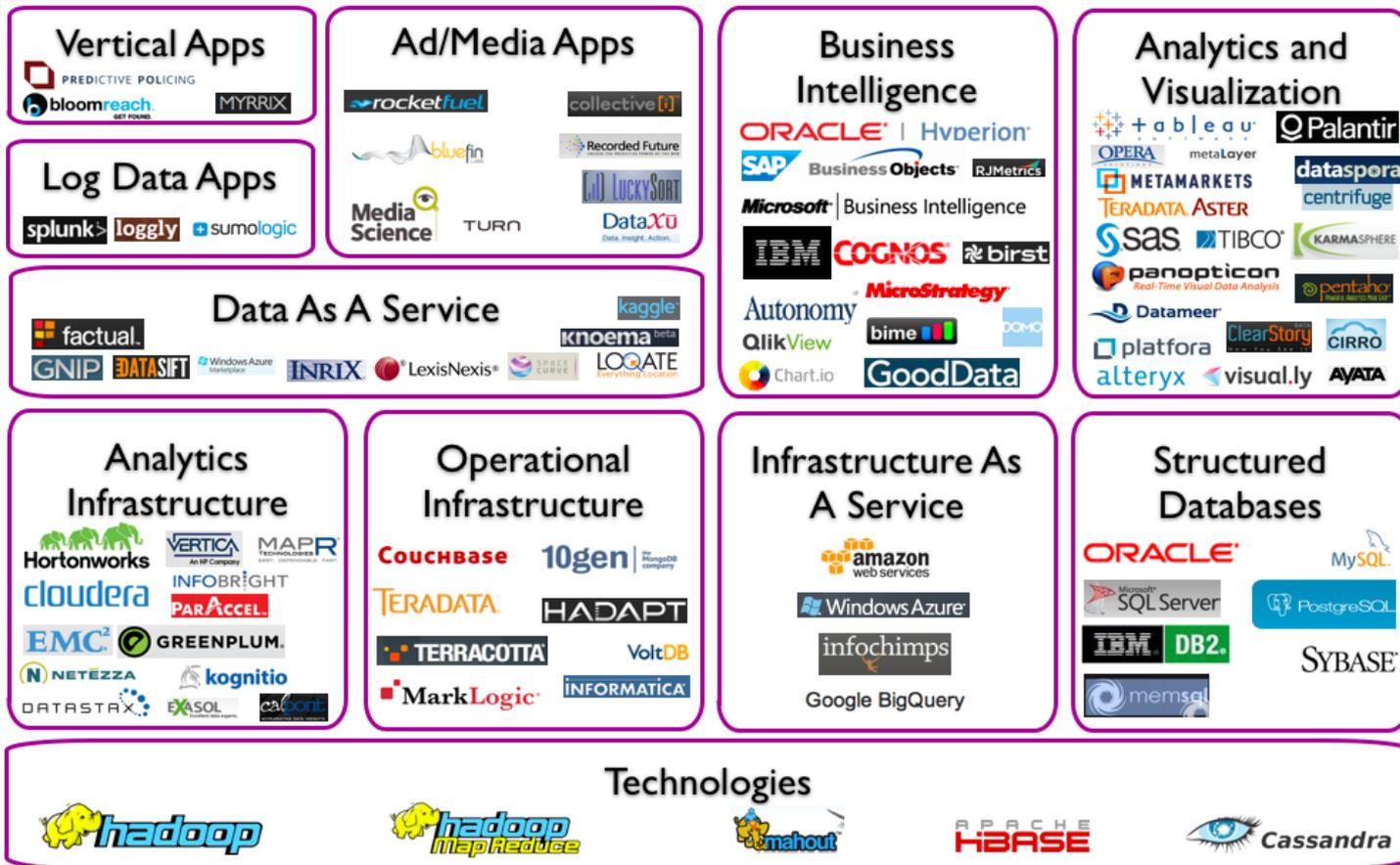


# Introducción al Big Data



# Introducción al Big Data

## Big Data Landscape



# Introducción al Big Data

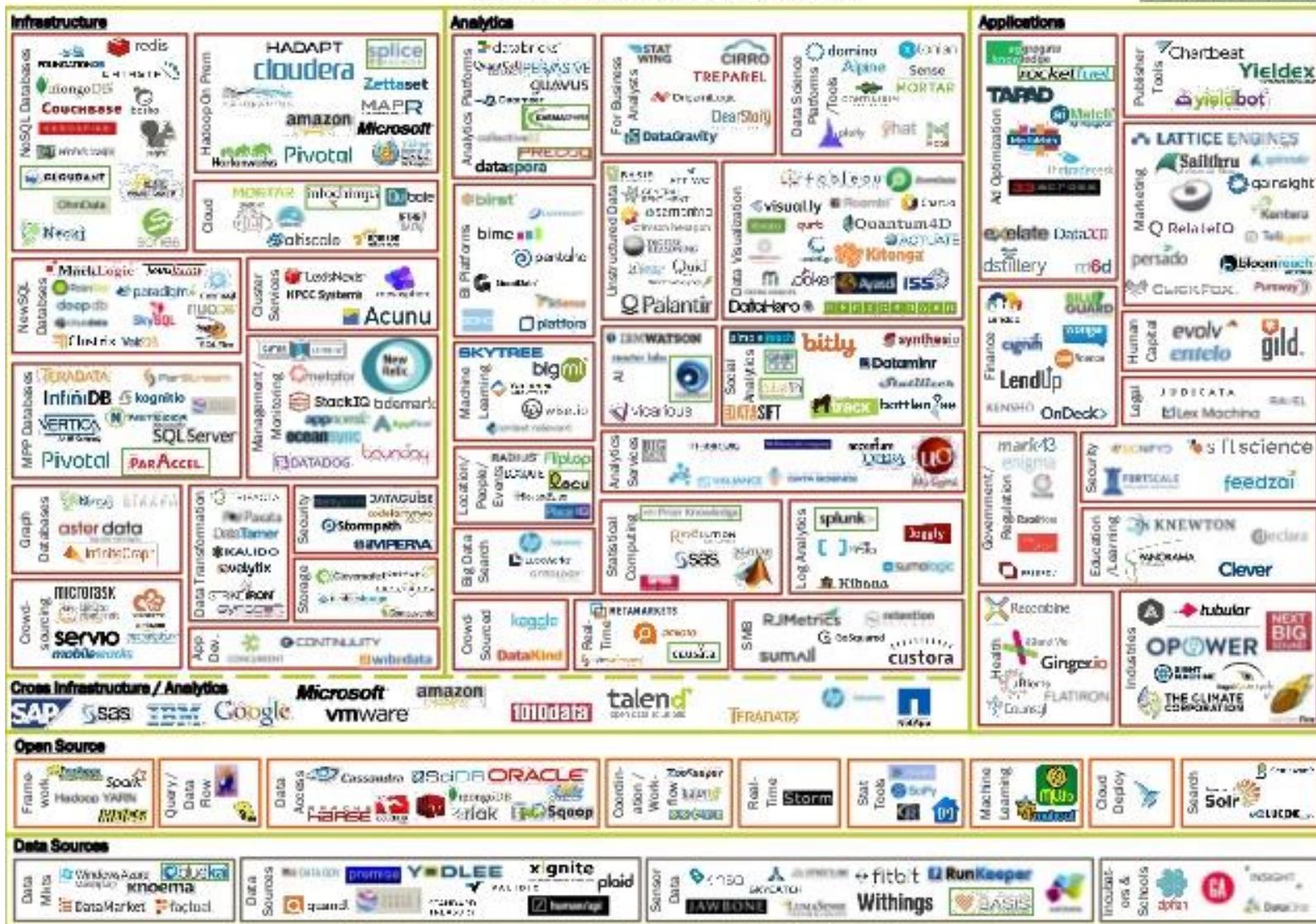
## Big Data Landscape (Version 2.0)



# Introducción al Big Data

## BIG DATA LANDSCAPE, VERSION 3.0

Exit: Acquisition or IPO



## Visión Global

Hardware Local

Plataformas y Servicios de Cloud

Gestión de Consumo Energético

Gestión de Redes de Comunicaciones

## Visión Global

Hadoop

Hadoop as a Service

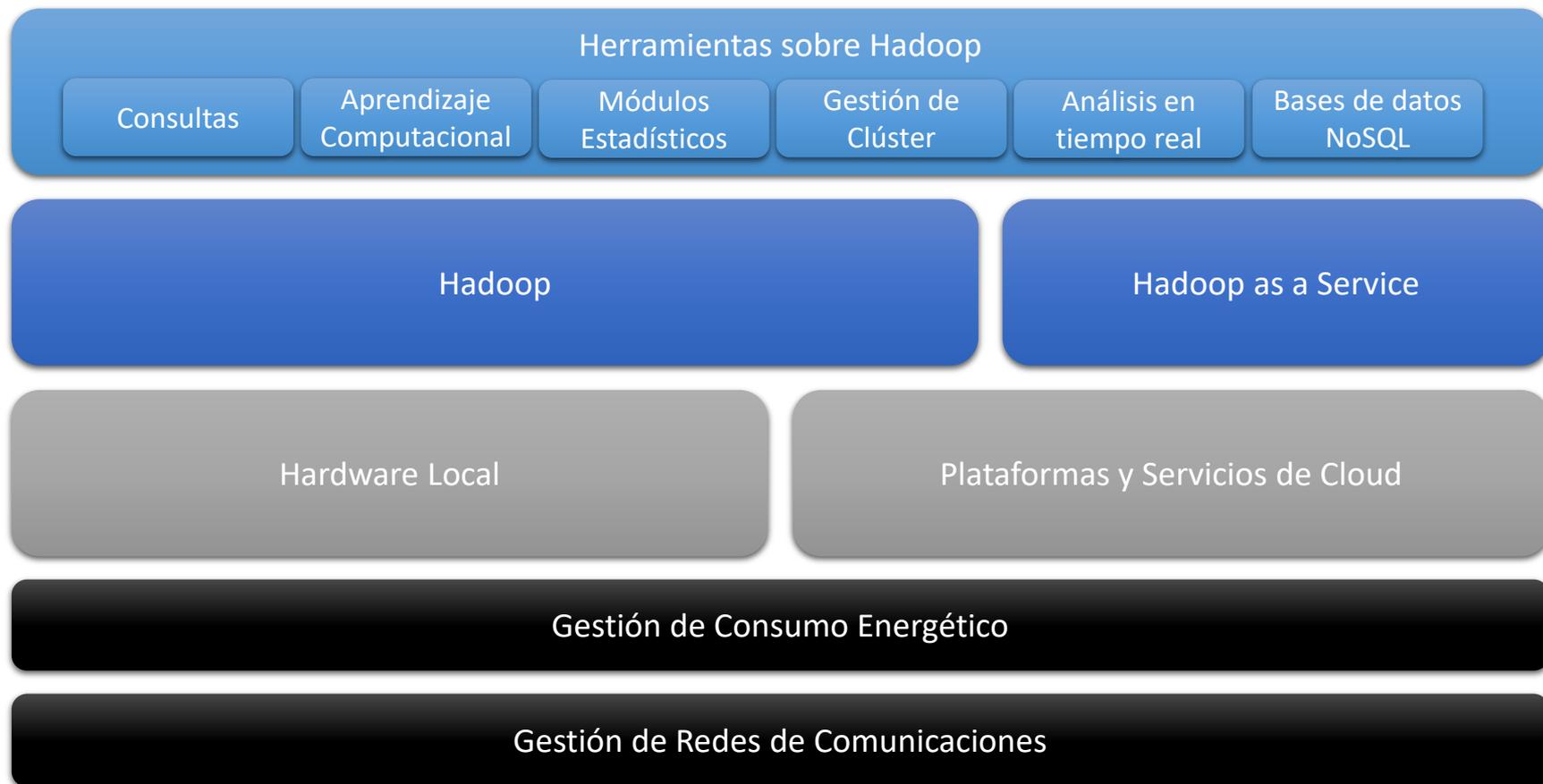
Hardware Local

Plataformas y Servicios de Cloud

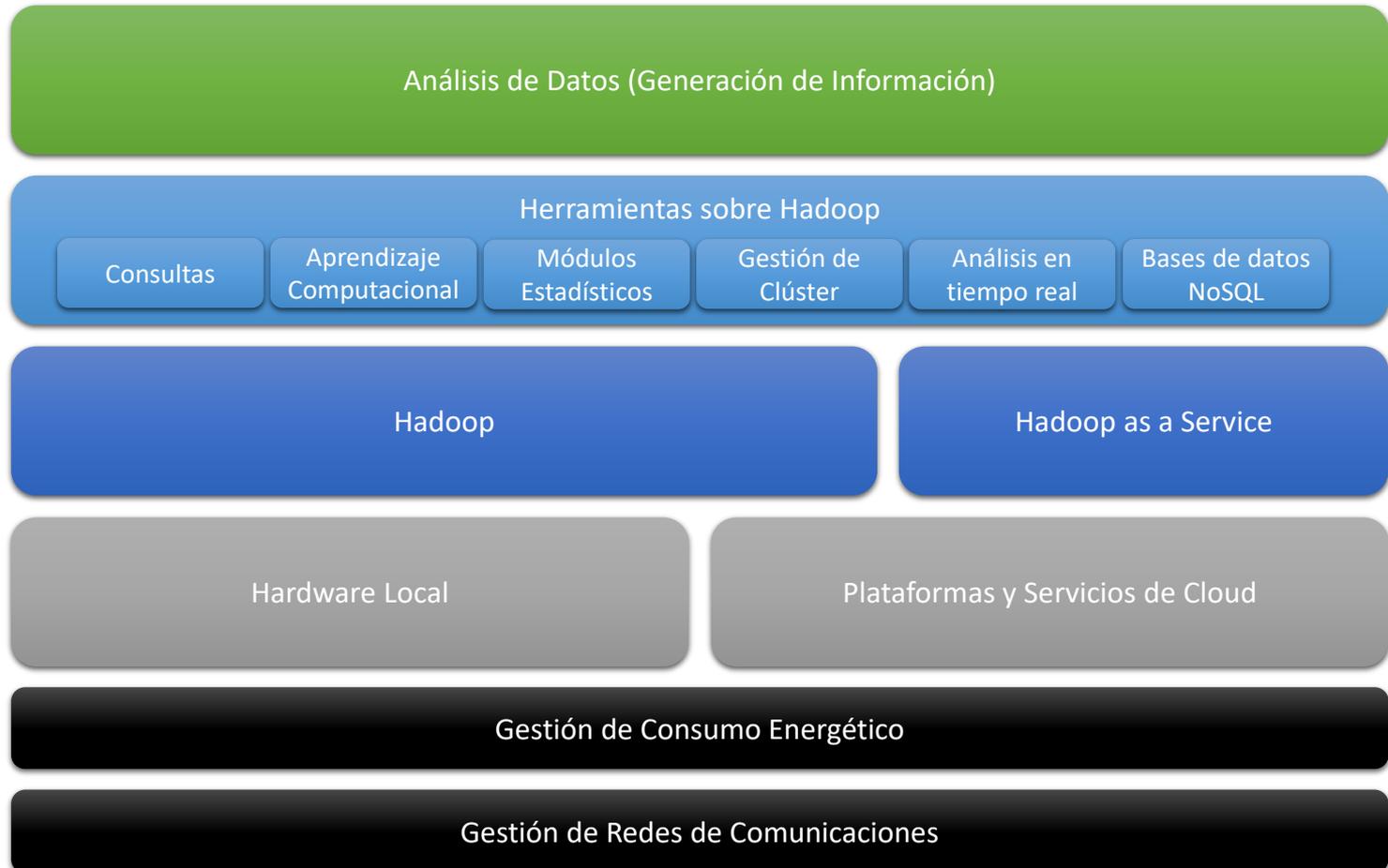
Gestión de Consumo Energético

Gestión de Redes de Comunicaciones

## Visión Global



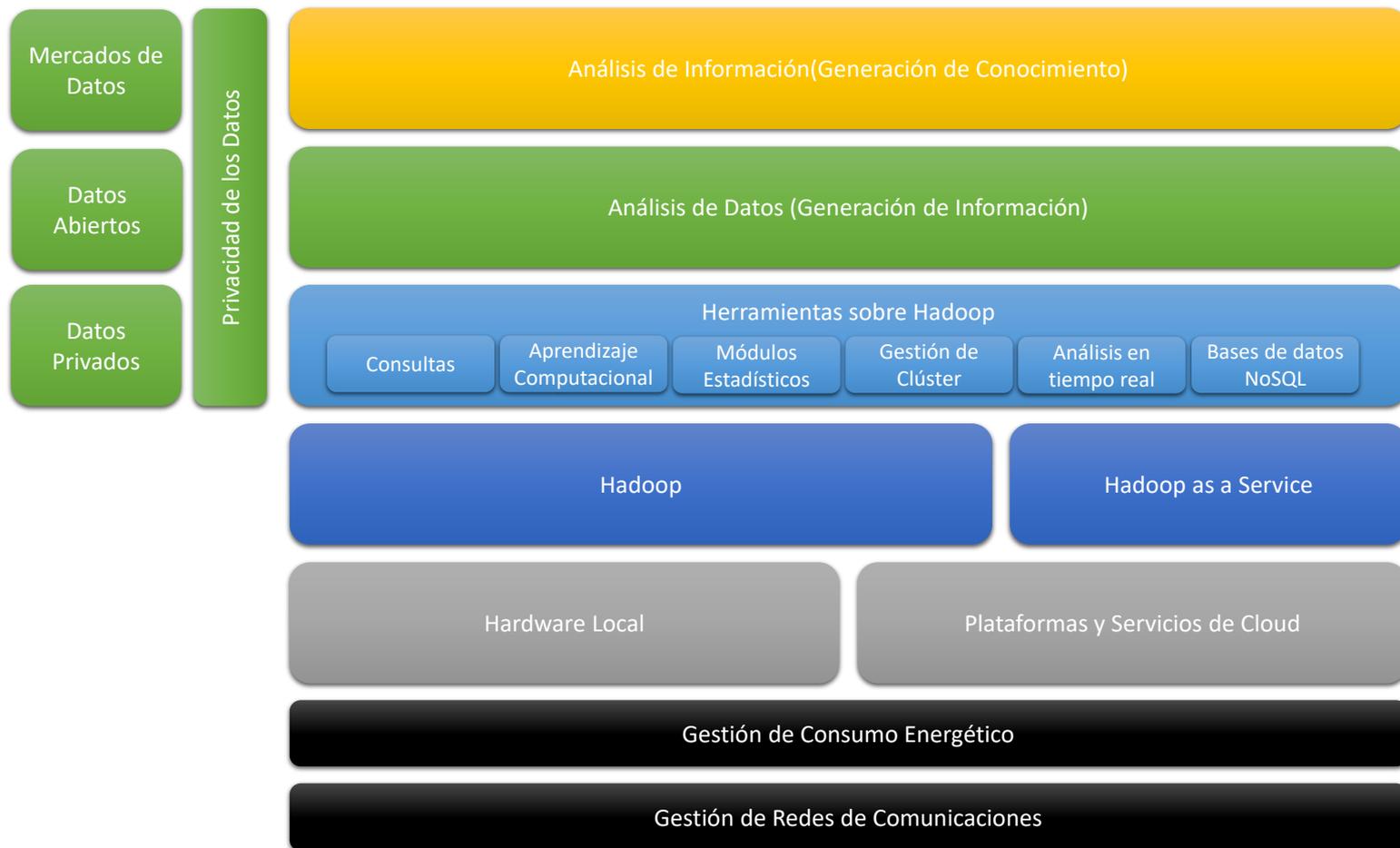
## Visión Global



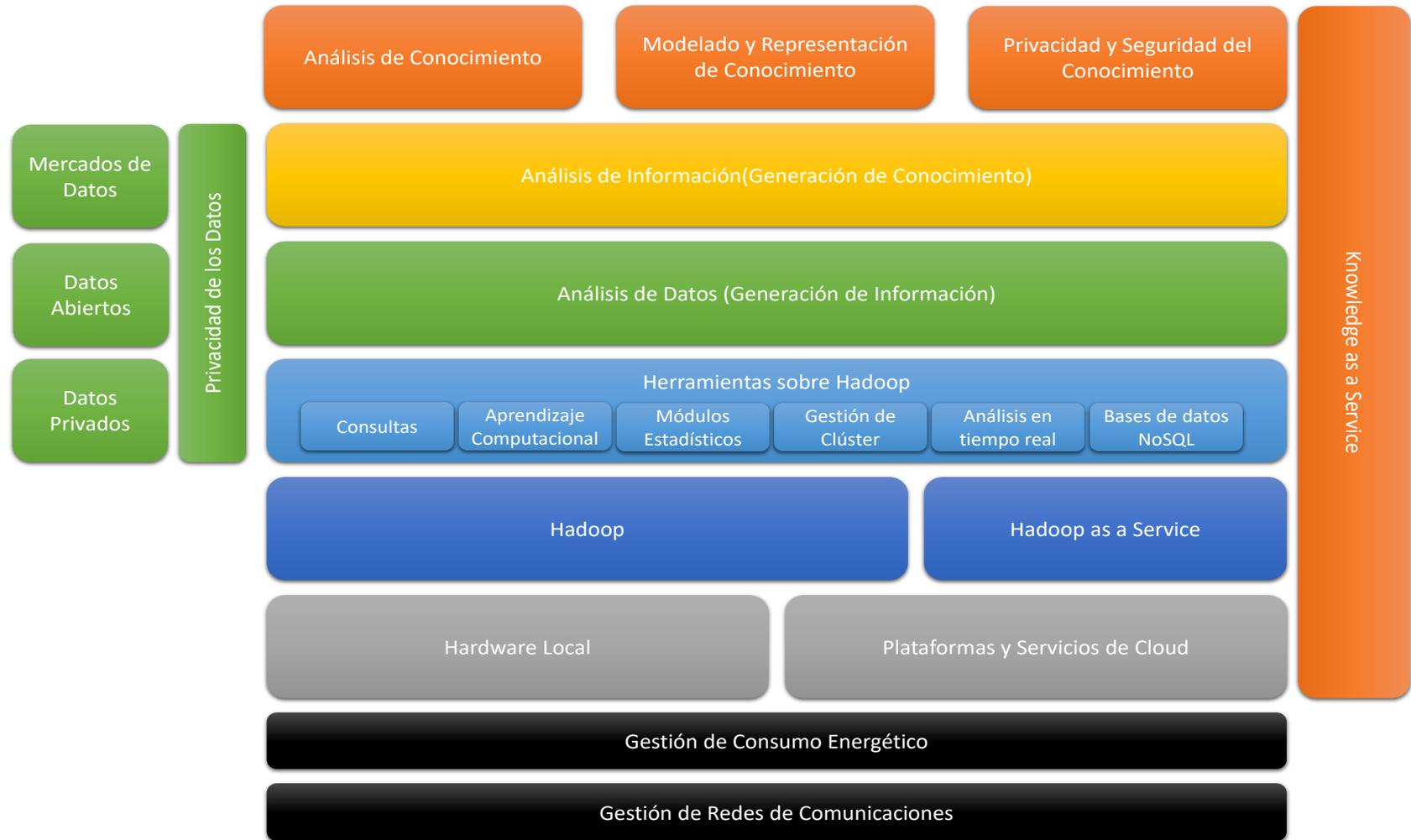
## Visión Global



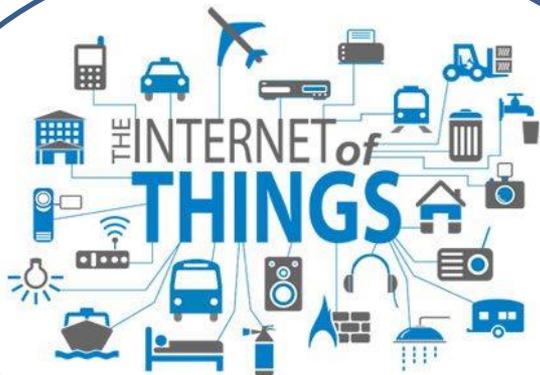
## Visión Global



## Visión Global



## Conceptos alrededor del Big Data



Big Data





## Conceptos alrededor del Big Data



## Conceptos alrededor del Big Data



Download from  
**Dreamstime.com**  
This watermark-free image is for previewing purposes only.

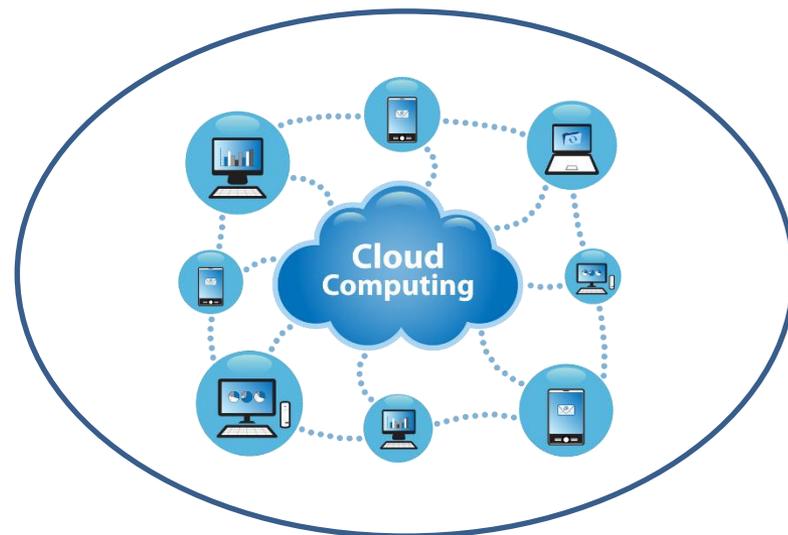


Download from  
**Dreamstime.com**  
This watermark-free image is for previewing purposes only.



## Conceptos alrededor del Big Data

La **computación en la nube**, conocida también como **servicios en la nube**, **informática en la nube**, **nube de cómputo** o **nube de conceptos** (del [inglés](#) *cloud computing*), es un [paradigma](#) que permite ofrecer [servicios](#) de computación a través de una red, que usualmente es [Internet](#).



**WIKIPEDIA**  
La enciclopedia libre

# Conceptos alrededor del Big Data

## 5 Reasons Businesses Use the Cloud

Every year, more and more businesses are adopting cloud. While traditionally thought of as an IT decision, cloud is increasingly being considered a business decision to enable company functions. Take a look at five reasons why more businesses are adding the cloud to their technology arsenals.

**1 It offers better insight and visibility**  
Businesses are using cloud technology to support their analytics efforts. Of leading organizations:

- 54% use analytics extensively to derive insights from big data
- 59% use cloud to share data seamlessly across applications
- 59% intend to use cloud to access and manage big data in the future

**2 It makes collaboration easy**  
Cloud allows work to be accessed from anywhere on multiple devices, making cross-functional collaboration much easier. Here's what leading organizations—those that are gaining competitive advantage through cloud—cite as popular uses:

- 58% collaborate across the organization and ecosystem
- 59% improve integration between development and operations

**3 It can support a variety of business needs**  
Companies are forging a tighter link between technology and business outcomes. Take a look at the business functions companies have migrated to the cloud.

- 18% messaging
- 15% storage
- 13% office/productivity suites

**4 It allows for rapid development of new products and services**  
The cloud offers businesses valuable capabilities. Here's what leading organizations say it enables them to do:

- 52% use it to innovate products & services rapidly
- 24% are able to offer additional products & services

**5 The results are proven**  
From business growth to increased efficiency, businesses using cloud are realizing benefits across the company.

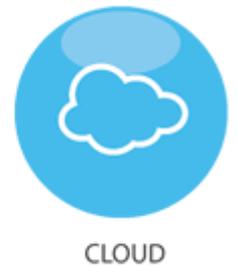
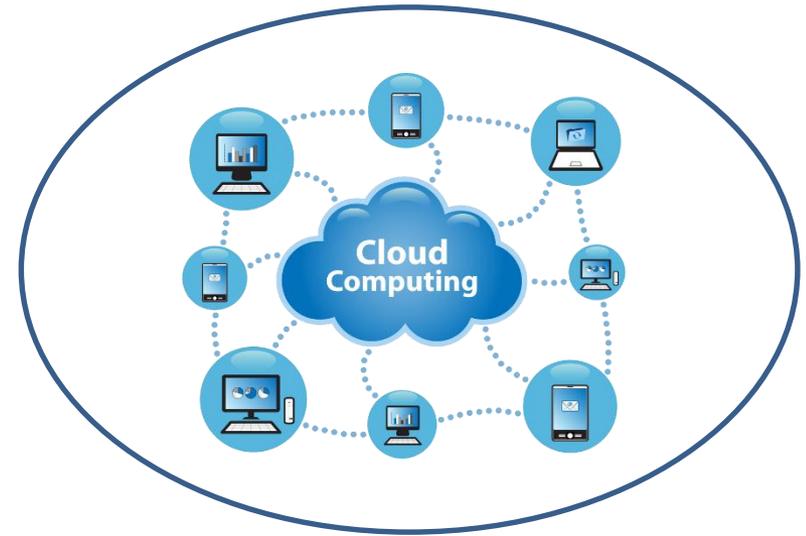
- 25% of businesses saw a reduction in IT costs
- 55% saw an increase in efficiency
- 48% saw improvement in employee mobility



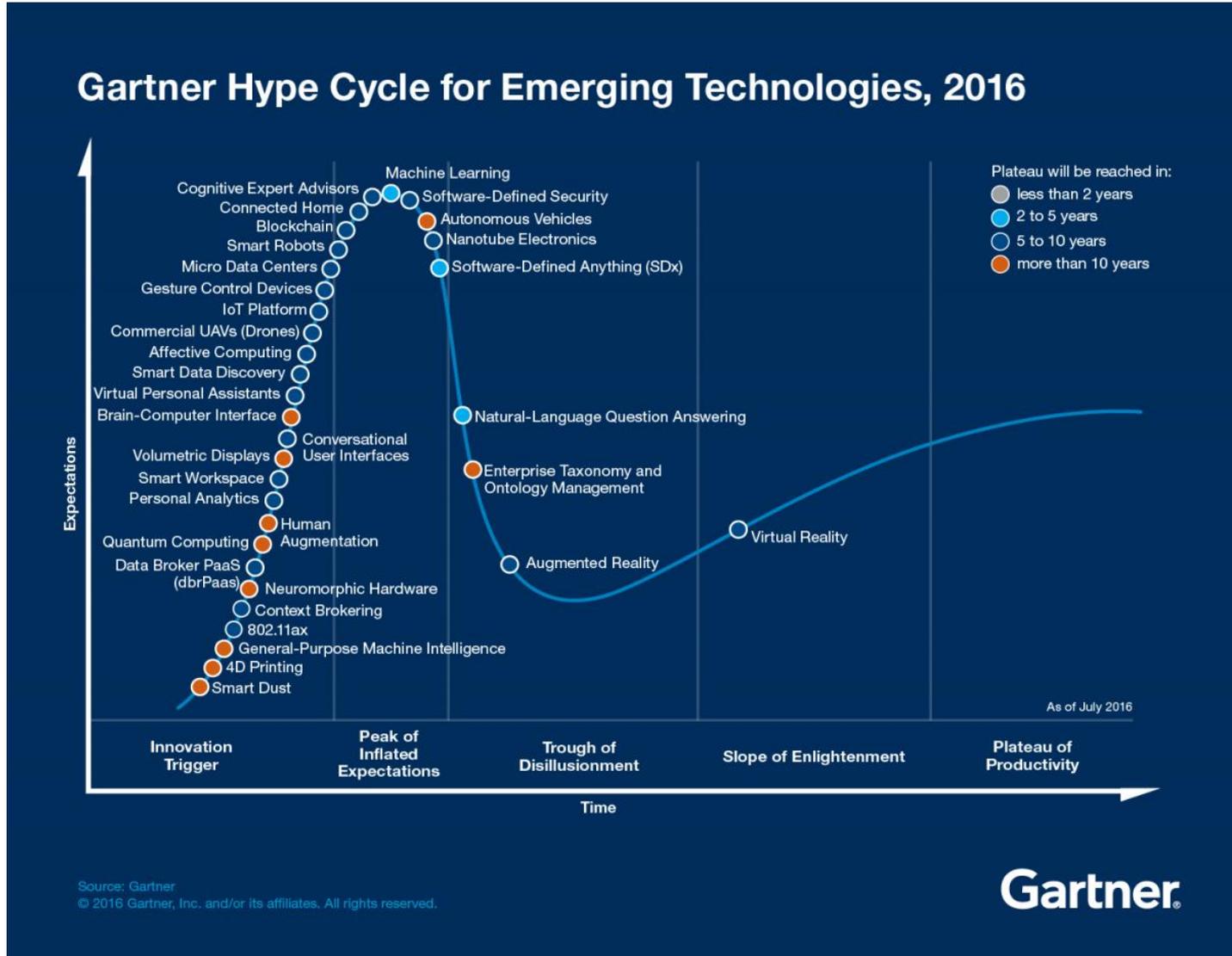
Sources: CDW, IBM Center for Applied Insights



## Conceptos alrededor del Big Data



## Conceptos alrededor del Big Data



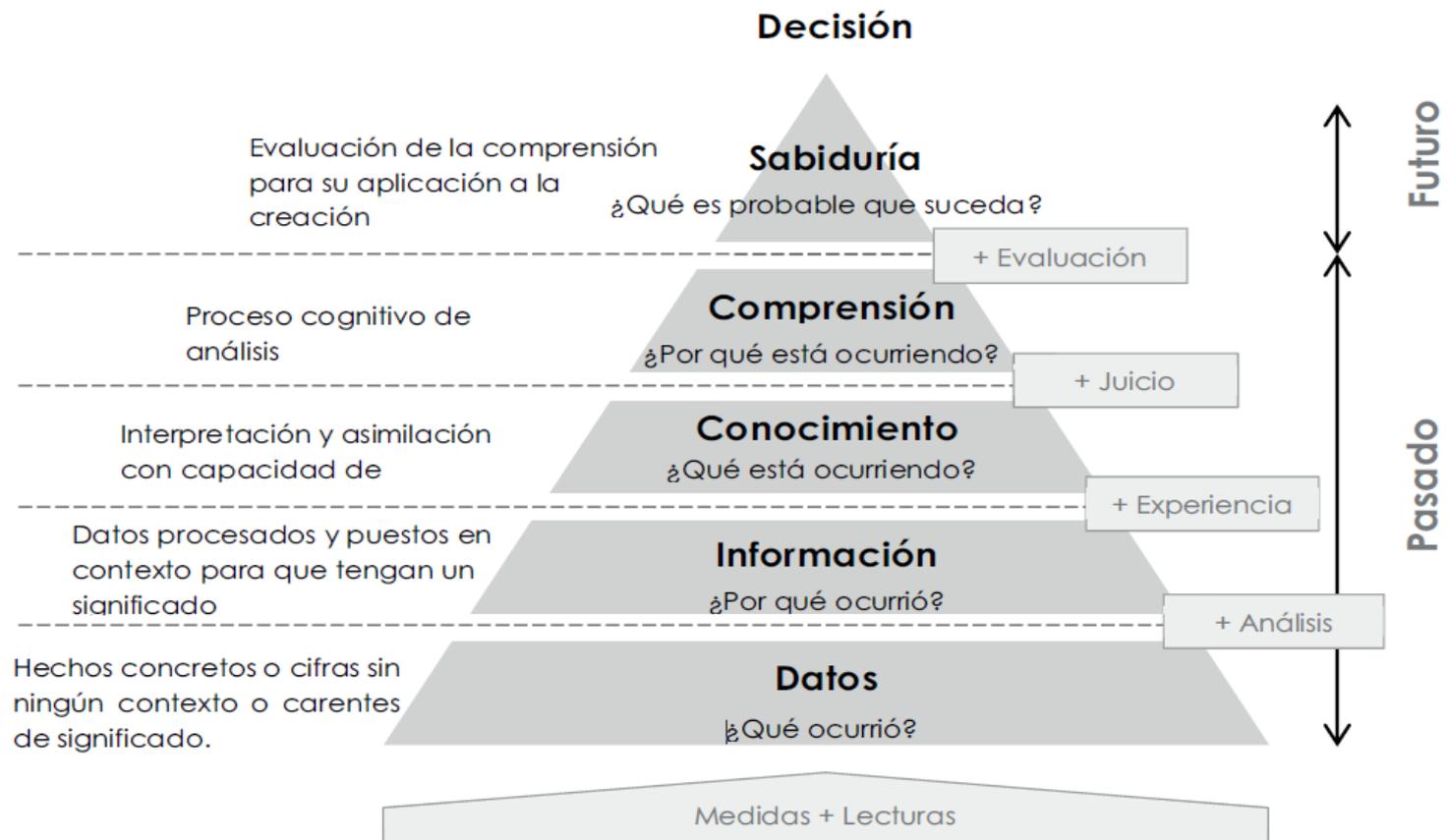
## Esquema de Contenidos:

1. Introducción al Big Dada
2. Trabajando con los datos
3. Almacenamiento y Procesamiento: Técnicas, Herramientas y Plataformas
4. Análisis mediante la Minería de Datos
5. Visualización y Consumo de Datos
6. Seguridad y Gobernanza
7. Aplicaciones Reales de Negocio: Casos de Éxito

# Trabajando con los datos

## Fundamentos del Trabajo con Datos

- Descubrimiento de Conocimiento en Bases de Datos (KDD)



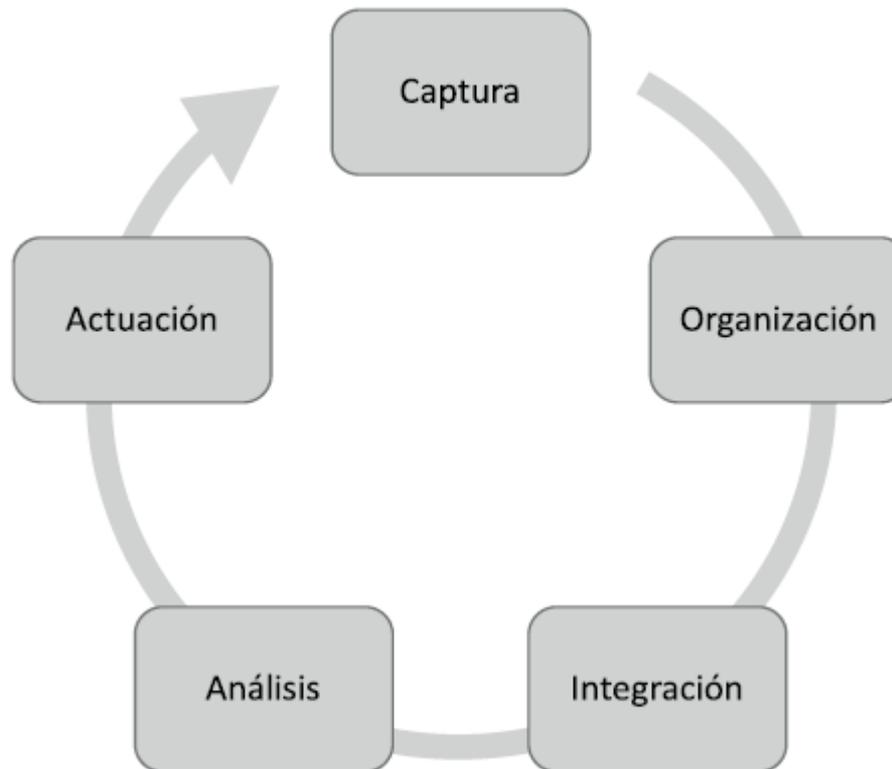
## *Fundamentos del Trabajo con Datos*

- ¿Qué es la Inteligencia de Negocio (Business Intelligence)?
- Necesidad Creciente de Información Estratégica
  - La crisis de la información
- Sistemas Operacionales frente a la Toma de Decisiones
- Almacenes de Datos

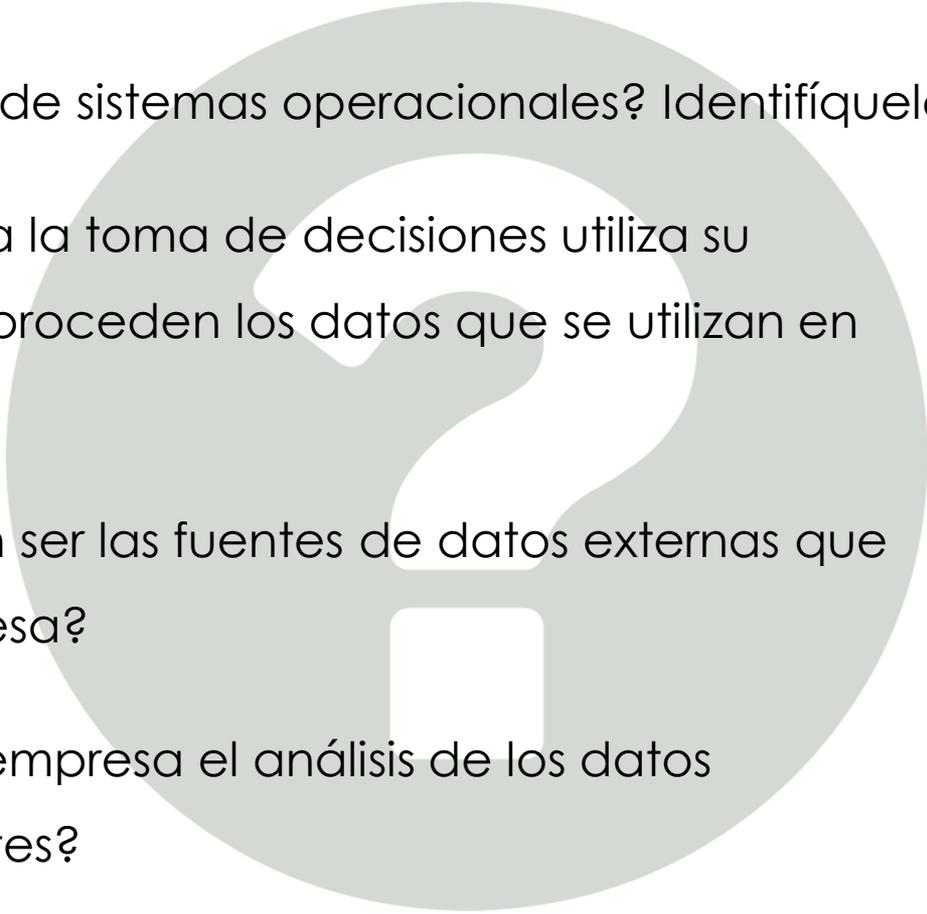
# Trabajando con los datos

## *El Ciclo de Vida de los Datos*

- Gestión del Ciclo de Vida de los Datos

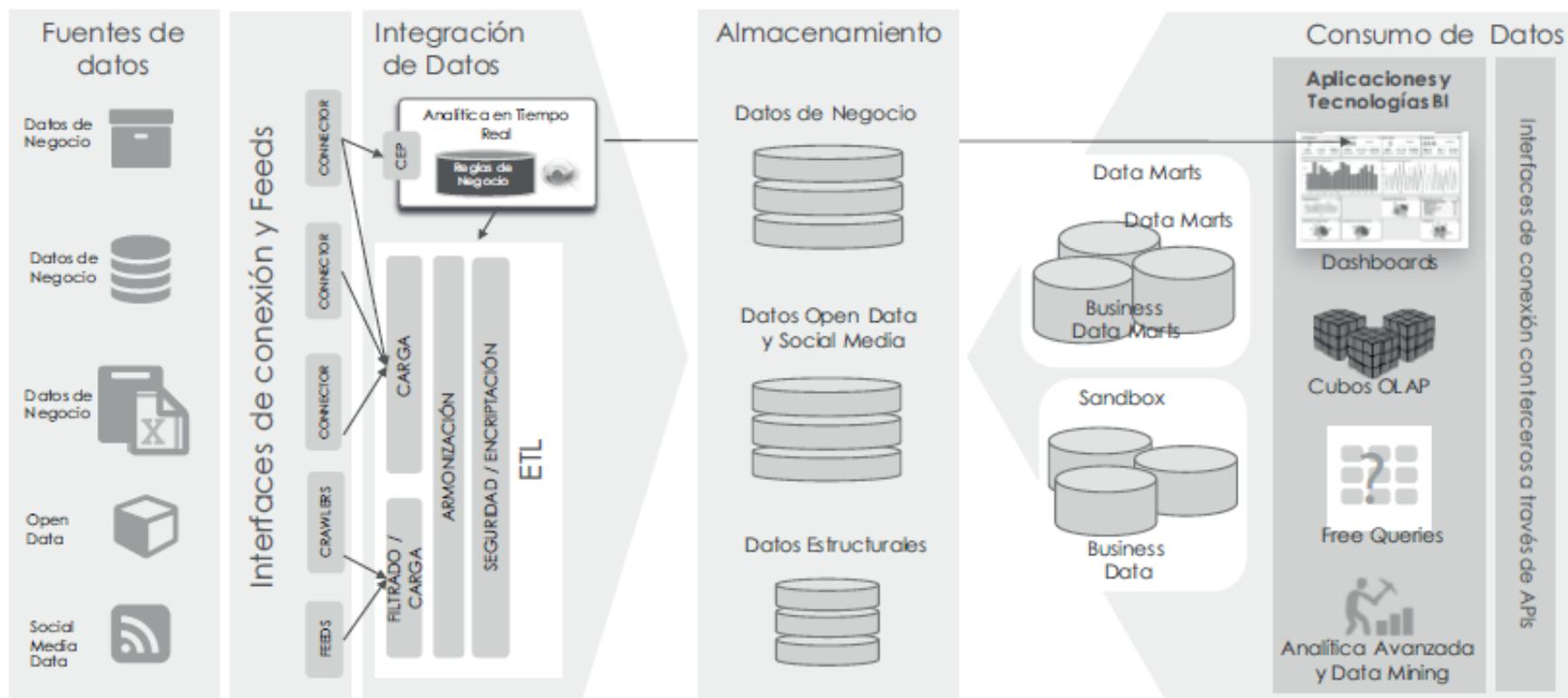


## *La Gestión de los Datos*

- ¿Dispone su organización de sistemas operacionales? Identifíquelos
  - ¿Qué sistemas de apoyo a la toma de decisiones utiliza su organización? ¿De dónde proceden los datos que se utilizan en estos sistemas?
  - ¿Cuáles cree que pueden ser las fuentes de datos externas que aportarían valor a su empresa?
  - ¿Cómo beneficiaría a su empresa el análisis de los datos procedentes de estas fuentes?
- 

# Trabajando con los datos

## Componentes de un Sistema Big Data para Business Intelligence

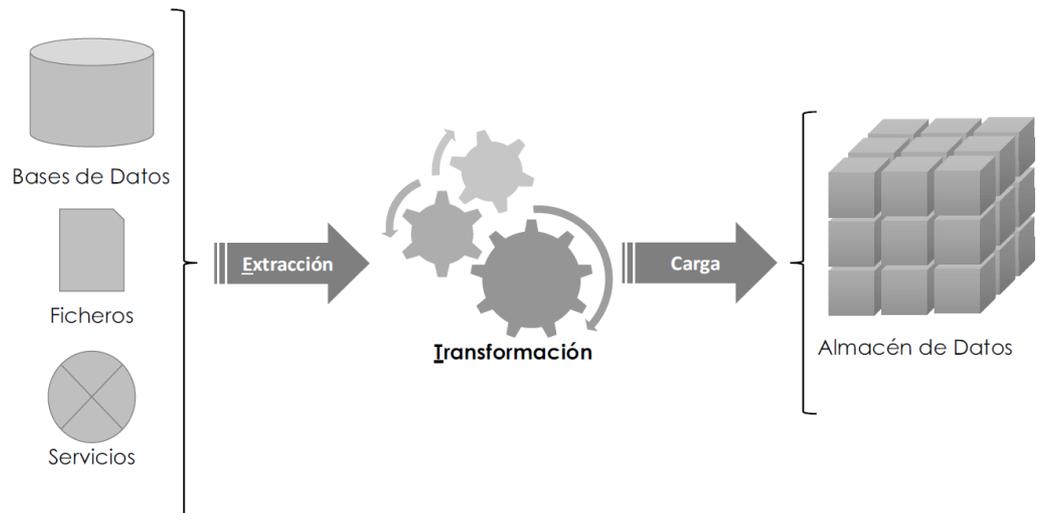


## *El Dilema del Directivo: Intuición vs. Datos*

- Machine learning
- La intuición como forma de tomar decisiones.
- La gestión basada en datos (Data Driven Management).
- Proceso de transición a la toma de decisiones basada en datos.
- Ser capaz de hacer las preguntas correctas.

# Trabajando con los datos

- **Extracción, transformación y carga**
  - La extracción implica “limpieza” de los datos
  - La transformación normalmente conlleva un “cambio” de modelo
  - La carga está unida a “validación” de la calidad de los datos



# Trabajando con los datos

## ▪ La calidad de los datos

- Las organizaciones actúan bajo la premisa de que la información disponible es precisa y válida. Si la información no es válida, entonces no pueden responder de las decisiones basadas en ella.
- ✓ El “profiling” de los datos es el análisis estadístico de los valores de los datos para evaluar su consistencia, originalidad y lógica
- ✓ El “sampling” (muestreo) estima la calidad de los datos mediante el análisis de un subconjunto de éstos
- ✓ Existen herramientas como Talend Open Studio que nos ayudan a gestionar la calidad de los datos



Resources | Community | Blog | Contact | Log In |

Products Downloads Support & Services Customers Company

### Watch and Learn Tutorials

Dimension	Description	Monitoring
Accuracy	Does the data accurately represent reality in a real-world context?	Accuracy
Completeness	Do broken links exist between data that should be linked?	Completeness
Consistency	Is there a single representation of data?	Consistency
Validity	Is the data valid against a set of rules or constraints?	Validity
Uniqueness	Is the data unique (no duplicates) and well-organized?	Uniqueness
Timeliness	Is the data up-to-date and available when needed?	Timeliness
Accessibility	Is the data easily accessible, understandable, and user-friendly?	Accessibility
Privacy	Is the data stored and processed in a secure manner?	Privacy
Performance	Is the information update frequency adequate?	Performance

Talend Open Studio for Data Quality for Dummies

Advanced Data Profiling

Fuzzy Matching Strategies

Embed Real-Time Data Quality into your Business Processes and Applications

<https://www.talend.com/products/talend-open-studio>

# Trabajando con los datos

- **La preparación de los datos**

- Manipulación y transformación de los datos sin refinar, para que la información contenida en el conjunto de datos pueda ser descubierta o estar accesible de forma más fácil.
  - ✓ Operaciones: normalización, restablecimiento de incompletos, limpieza, conversión, agrupación, eliminación de redundancias, simplificación, conversión, etc.

- **Transformación de los datos**

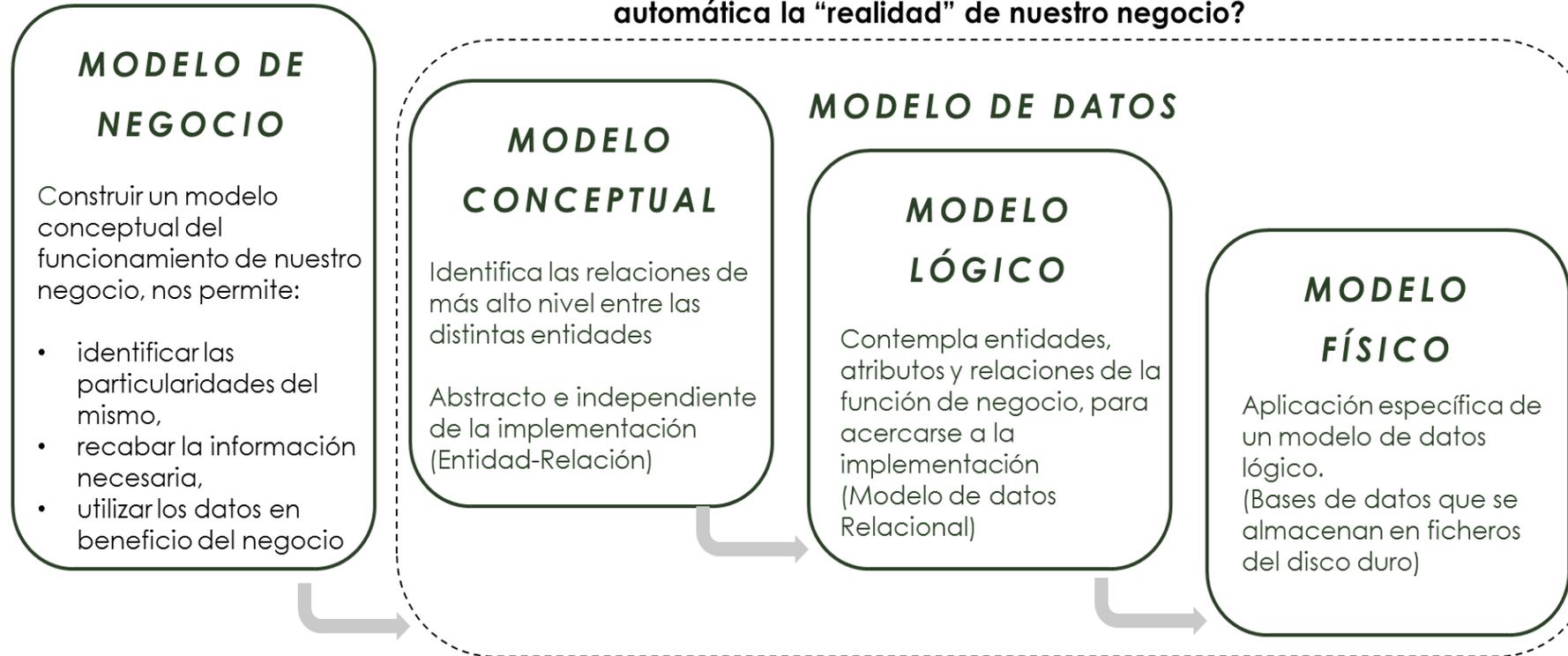
- Aplicación de una serie de funciones o reglas de negocio sobre los datos extraídos para convertirlos en datos que, a continuación, serán cargados en la nueva fuente.

- Ejemplo: Datos numéricos a categóricos
  - ✓ 18-30 Años: Jóvenes.
  - ✓ 31-65 Años: Adultos en edad laboral
  - ✓ 66 Años en adelante: Clientes en edad de jubilación.

# Trabajando con los datos

## Recolección de datos

- ¿Cómo almacenamos de manera automática la “realidad” de nuestro negocio?

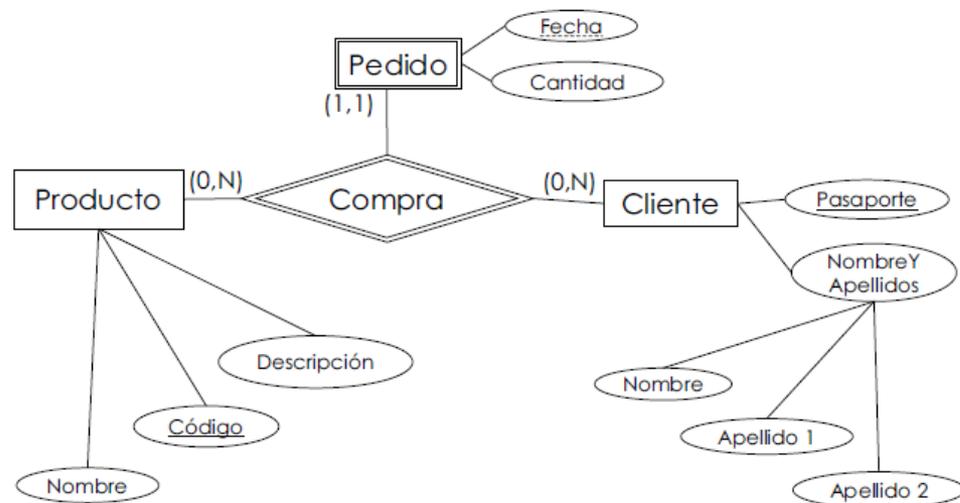


# Trabajando con los datos

## Recolección de datos

### El Modelo Conceptual

- Trasladar el modelo de nuestro negocio –algo que utiliza un lenguaje y conceptos propios del mundo empresarial– al lenguaje y los conceptos propios del mundo de la tecnología
- Hacer entender a los técnicos o tecnólogos nuestro negocio y sus particularidades e interrelaciones

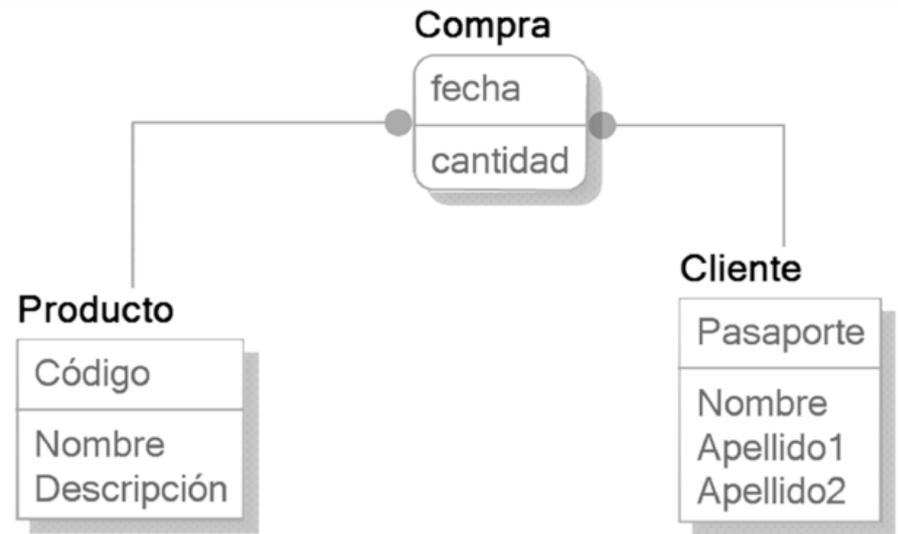


# Trabajando con los datos

## Recolección de datos

### ▪ El Modelo Lógico

- Está comprometido con un modelo de implementación concreto. Por ejemplo, el modelo de datos relacional
- Se adecua su estructura y contenido para facilitar su implementación en un sistema gestor de bases de datos concreto

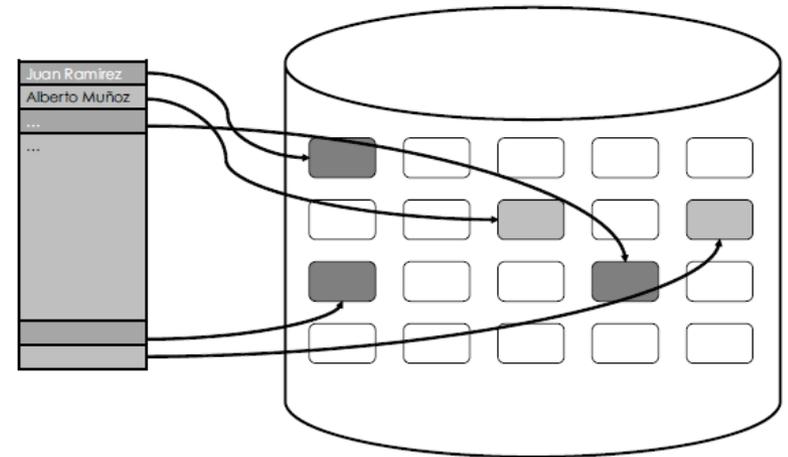


# Trabajando con los datos

## Recolección de datos

### ▪ El Modelo Físico

- El modelo físico describe, cómo los datos que siguen el modelo lógico, son almacenados en el disco y cómo acceder eficientemente a estos datos
- Las bases de datos en las que se apoyan la mayoría de las aplicaciones de gestión son el modelo físico que surge de la aplicación del modelo lógico definido.
- El modelo relacional utiliza un lenguaje específico para hacer consultas a los datos almacenados conocido como lenguaje «**SQL**», acrónimo de «Structured Query Language», o lenguaje de consultas estructurado.



## Esquema de Contenidos:

1. Introducción al Big Dada
2. Trabajando con los datos: Procesos ETL
3. Almacenamiento y Procesamiento: Técnicas, Herramientas y Plataformas
4. Análisis mediante la Minería de Datos
5. Visualización y Consumo de Datos
6. Seguridad y Gobernanza
7. Aplicaciones Reales de Negocio: Casos de Éxito

## *Almacenamiento Masivo de Datos*

- En un PC, almacenamos datos en nuestro disco duro.
- Los datos pueden ser de cualquier tipo:
  - Tablas, ficheros de texto, páginas web...
  - ...imágenes, vídeos, música...
- BIG DATA: ¿Qué pasa cuando los datos no caben en un disco duro?

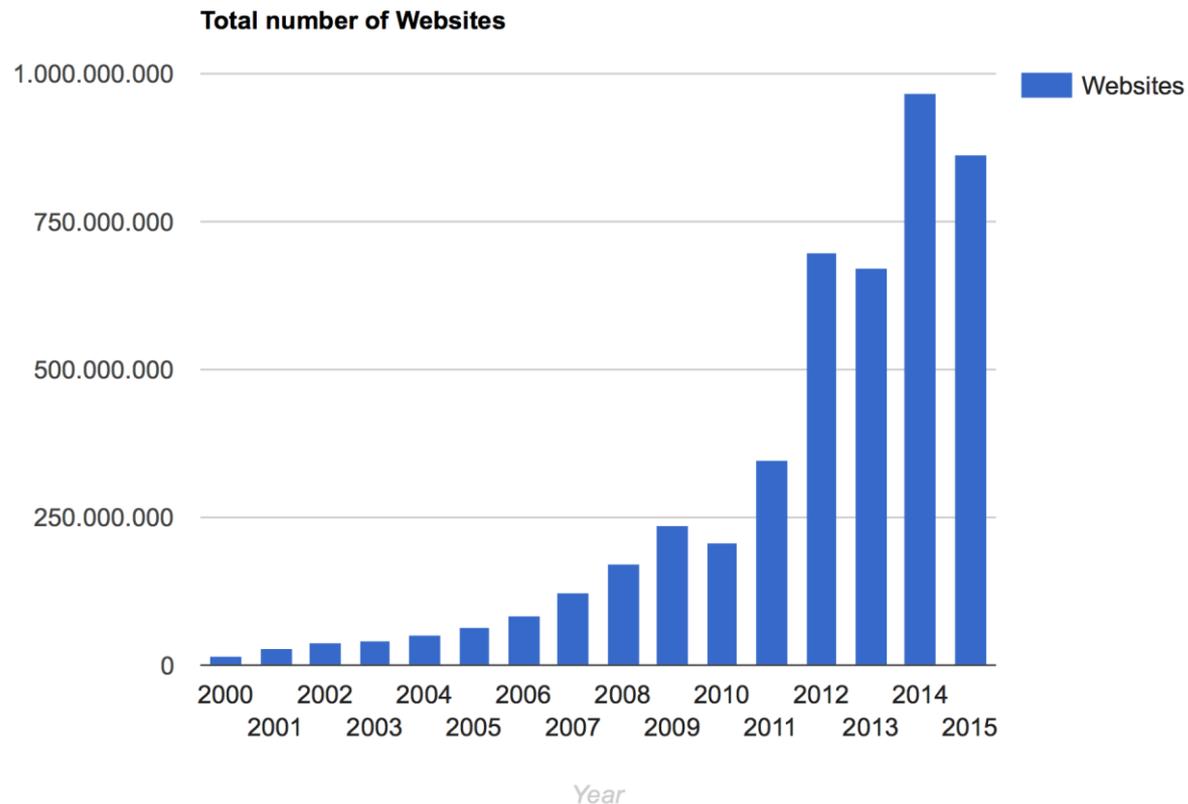


## *El Enfoque de Google*

- Google es uno de los máximos exponentes del Big Data.
- El buscador web tiene que indexar el contenido de todas las webs posibles.
- Nace en 1998. La web ha crecido mucho desde entonces...

## El Enfoque de Google

- Evolución del número de sitios web (no de páginas web)



Fuente: [www.internetlivestats.com/total-number-of-websites/](http://www.internetlivestats.com/total-number-of-websites/)

## *El Enfoque de Google*

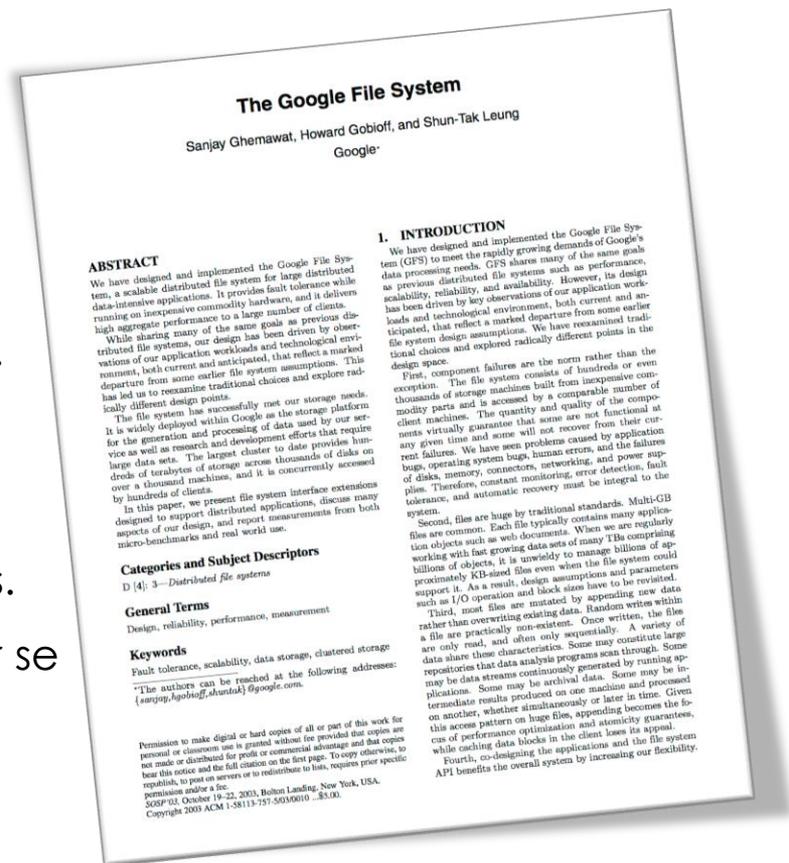
- Google no es solo un motor de búsqueda.
- En los primeros 5 años de vida (1998-2003), consigue lanzar estos productos (y más)



# Almacenamiento y Procesamiento

## Google File System

- En 2003, Google anuncia Google File System (GFS)
- Es un sistema de ficheros distribuido. Permite:
  - Almacenar grandes cantidades de datos.
  - Distribuir (repartir) los datos entre varios ordenadores.
  - Permitir lecturas de ficheros más eficientes.
  - Evitar la pérdida de datos si un ordenador se estropea.



## *Hadoop Distributed File System*

- En 2004, D. Cutting y M. Cafarella liberan una versión de código abierto de GFS.
- En 2005, este sistema se denomina **HDFS** y nace el proyecto HADOOP.



# Almacenamiento y Procesamiento

## *HDFS – Distribución*

- Cada fichero se divide en bloques de 64 o 128 MB.
- Los bloques se almacenan divididos entre varios ordenadores.
  - Cada ordenador se denomina **nodo**.
  - El conjunto de ordenadores se denomina **cluster**.
- VENTAJA 1: Si se duplica el número de nodos, caben el doble de datos.
- VENTAJA 2: Se pueden almacenar ficheros que no quepan en un ordenador.

- VENTAJA 3: Las lecturas son más rápidas:

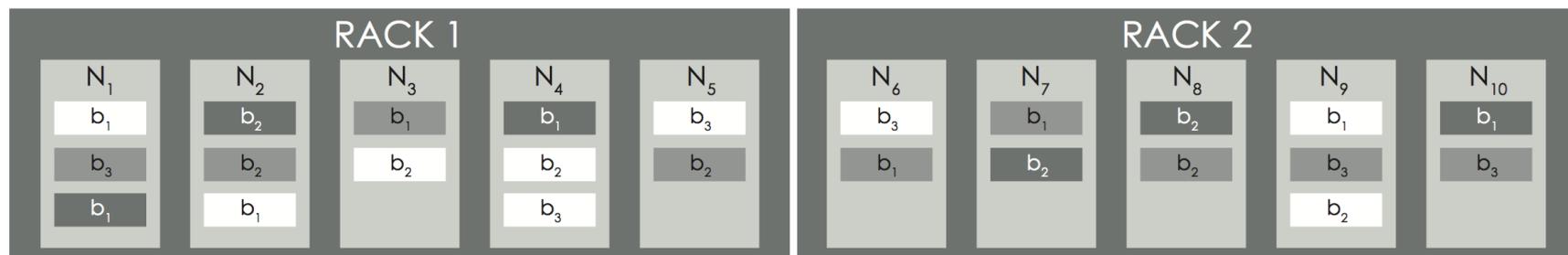
se puede leer un fichero de varios

ordenadores a la vez.



## HDFS – Replicación

- PROBLEMA: Si hay muchos nodos, la probabilidad de que uno falle es alta.
- SOLUCIÓN: Cada bloque se almacena en varios nodos, por defecto 3.
- Cada bloque almacenado en un nodo se denomina **réplica**.
- Las réplicas tienen en cuenta la topología: mejor almacenarlas en *racks*\* o centros de datos diferentes.



\* Un rack o armario es un contenedor donde se almacenan varios nodos que comparten sistema de alimentación y de red.

## *Limitaciones HDFS*

- HDFS es un sistema de ficheros
- Permite almacenar datos sin estructura (cualquier tipo de ficheros).
- Los datos no están indexados.
- Es ideal cuando queremos recuperar y procesar datos de forma masiva.
- INCONVENIENTE: No es útil si queremos recuperar un conjunto pequeño de los datos.

## *SOLUCIÓN: Bases de Datos*

- Una base de datos almacena datos con una cierta estructura.
  - Esta estructura se denomina **modelo de datos**.
- Normalmente se pueden definir **índices**, que permiten recuperar datos más eficientemente.
- Operaciones habituales (CRUD)
  - Recuperación.
  - Inserción / Actualización.
  - Borrado.

## *Bases de Datos NoSQL*

- Desde hace años se han utilizado sobre todo **bases de datos relacionales**.
  - También denominadas **SQL**, por el lenguaje de consulta que implementan.
- En los últimos años se han extendido las BBDD **NoSQL**.
- Características principales (no todas aplican siempre)
  - Almacenan datos semiestructurados.
  - Son escalables (permiten almacenar Big Data).
- Cada tipo de base de datos NoSQL tiene sus propias ventajas e inconvenientes.

## **NoSQL: BBDD Clave-Valor**

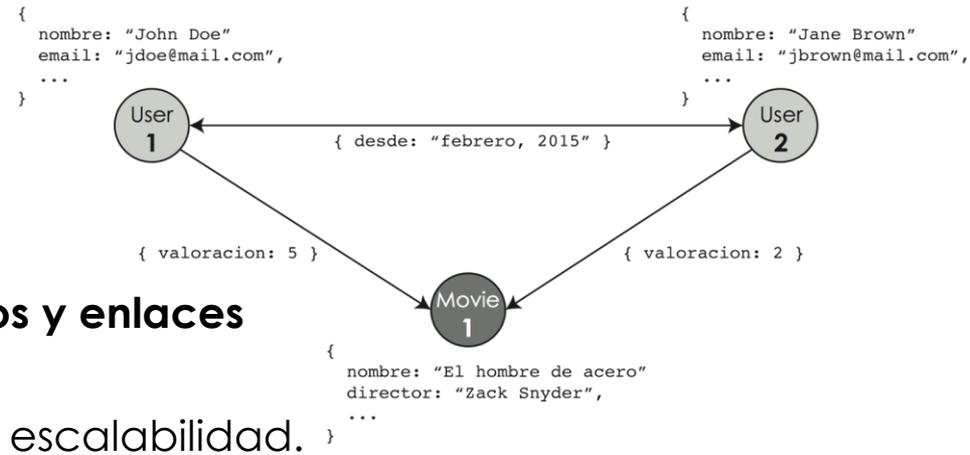
- Son las bases de datos con el modelo de datos más simple.
  - Almacenan **valores** (de cualquier tipo) asociados a una **clave**.
- ESCALABILIDAD: Son muy escalables, basta con distribuir los datos.
- RENDIMIENTO: Son muy rápidas, porque las claves están indexadas.
- OPERACIONES: Su operativa es muy limitada:
  - Recuperar un valor a partir de su clave.
  - Insertar un dato (par clave-valor).
  - Borrar un valor dada su clave.

## **NoSQL: BBDD Documentales**

- Almacenan documentos **semiestructurados**. Cada documento puede tener campos distintos.
- ESCALABILIDAD: Soportan **sharding**, es decir, distribución de los documentos.
- RENDIMIENTO: Son rápidas, se pueden definir índices en cualquier campo.
- OPERACIONES: Permiten una amplia operativa, pero no soportan JOIN.
  - Recuperar documentos por cualquier campo, con consultas avanzadas.
  - Insertar nuevos documentos y actualizar los existentes.
  - Borrar documentos.
  - Realizar operaciones de agregación de datos.
- La más extendida es **MongoDB**.

# Almacenamiento y Procesamiento

Figura 37 – Ejemplo de grafo representando una red social



## NoSQL: BBDD Orientadas a Grafos

- Almacenan datos en forma de **nodos y enlaces**
- ESCALABILIDAD: No destacan por su escalabilidad.
- RENDIMIENTO: Son rápidas, se pueden definir índices en cualquier campo.
- OPERACIONES: Permiten una amplia operativa, soportando JOIN.
  - Recuperar documentos por cualquier campo, pudiendo recuperar sus enlaces.
  - Insertar nuevos nodos o enlaces, y actualizar los existentes.
  - Borrar nodos o enlaces.
  - Realizar operaciones de agregación de datos.
- La más extendida es **Neo4j**.

## **NoSQL: BBDD Orientadas a Columnas**

- Almacenan datos en tablas, donde cada fila es un registro y existen familias de columnas.
- Cada registro tiene una **clave** (*rowkey*), y cada familia de columnas almacena pares clave-valor.
- ESCALABILIDAD: Son muy escalables, distribuyen los datos en varios nodos.
- RENDIMIENTO: Son eficientes para recuperar datos por su *rowkey*, que está indexada.
- OPERACIONES: Soportan operaciones básicas, pero no soportan JOIN.
  - Recuperar registros, pudiendo obtener solo algunas familias de columnas o valores.
  - Insertar nuevos registros, y actualizar los existentes.
  - Borrar registros.

## Comparación de BBDD SQL y NoSQL

Figura 38 – Comparación entre tecnologías de bases de datos

	Relacionales <i>MySQL</i>	Clave-Valor <i>Voldemort</i>	Documentos <i>MongoDB</i>	Grafos <i>Neo4j</i>	Columnas <i>HBase</i>
<b>Modelo de Datos</b>	Tablas y relaciones	Pares clave-valor	Documentos BSON/JSON	Nodos y enlaces	Tablas con columnas no fijas
<b>Relaciones (JOIN)</b>	✓	✗	✗	✓	✗
<b>Transacciones ACID</b>	✓	✗	✗	✓	✗
<b>Índices</b>	✓	Solo en la clave	✓	✓	Solo en la clave de fila (rowkey)
<b>Datos semiestructurados</b>	✗	✓	✓	✓	✓
<b>Escalabilidad horizontal</b>	✗	✓	✓	✗	✓
<b>Replicación de datos</b>	✓	✓	✓	✓	✓

## ***Almacenamiento de Big Data en la Nube***

- Los sistemas de Cloud Storage nos permiten almacenar grandes cantidades de datos:
  - Ficheros
  - Datos estructurados en BBDD relacionales
  - Datos semiestructurados en BBDD no relacionales
- No tenemos que adquirir equipos: ahorramos en discos duros, ordenadores, etc.
- No tenemos que preocuparnos de gestionar *backups*, la empresa lo gestiona por nosotros.
- Pagamos por la cantidad de datos que almacenamos, o bien por las veces que accedemos a ellos.

## Almacenamiento de Ficheros

- Existe el concepto de “*bucket*” (cubo).
  - Un cubo es un espacio donde almacenamos ficheros y directorios.
- Dentro de cada *bucket*, encontramos una jerarquía de directorios y ficheros, como es habitual.
- El servicio de *cloud* gestiona por nosotros la disponibilidad, *backups*, etc.
- Principales competidores:
  - Amazon S3
  - Google Cloud Storage



Google Cloud Storage

## *Bases de Datos NoSQL*

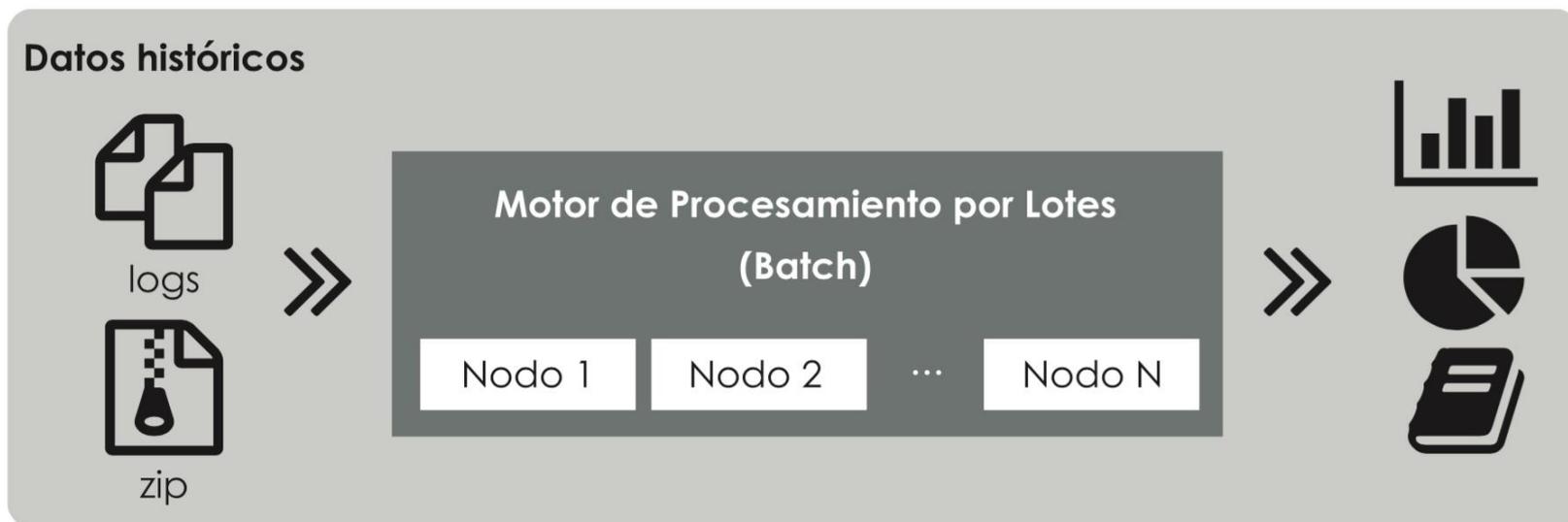
- Principal solución: Amazon DynamoDB.
- El modelo de datos está compuesto por tablas.
- Cada tabla contiene objetos, que a su vez contienen atributos.
- Los objetos deben tener una clave primaria. El resto de atributos no está definido a priori.
- Escala horizontalmente para almacenar Big Data.



## ***Procesamiento de Big Data***

- Dos paradigmas principales:
  - Por lotes.
  - En tiempo real.
- Procesamiento por lotes:
  - Se analizan grandes volúmenes de datos históricos.
  - El procesamiento puede llevar minutos, horas, días...
- Procesamiento en tiempo real:
  - Se analizan grandes volúmenes de datos que llegan a gran velocidad.
  - El procesamiento debe ser inmediato.

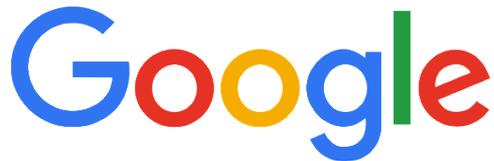
## Procesamiento de Big Data por Lotes (Batch)



# Almacenamiento y Procesamiento

## El (Nuevo) Enfoque de Google

- En 2004, Google sigue lanzando nuevos productos...



## La Solución de Google

- En 2004, Google anuncia MapReduce
- Es un *framework* para el procesamiento de datos. Permite:
  - Procesar datos almacenados sobre GFS.
  - Hacerlo de manera distribuida y paralela.
  - Evitar la pérdida de datos si un ordenador se estropea.



## *Hadoop MapReduce*

- En 2005, una versión de código abierto de **MapReduce** se incluye en HADOOP.

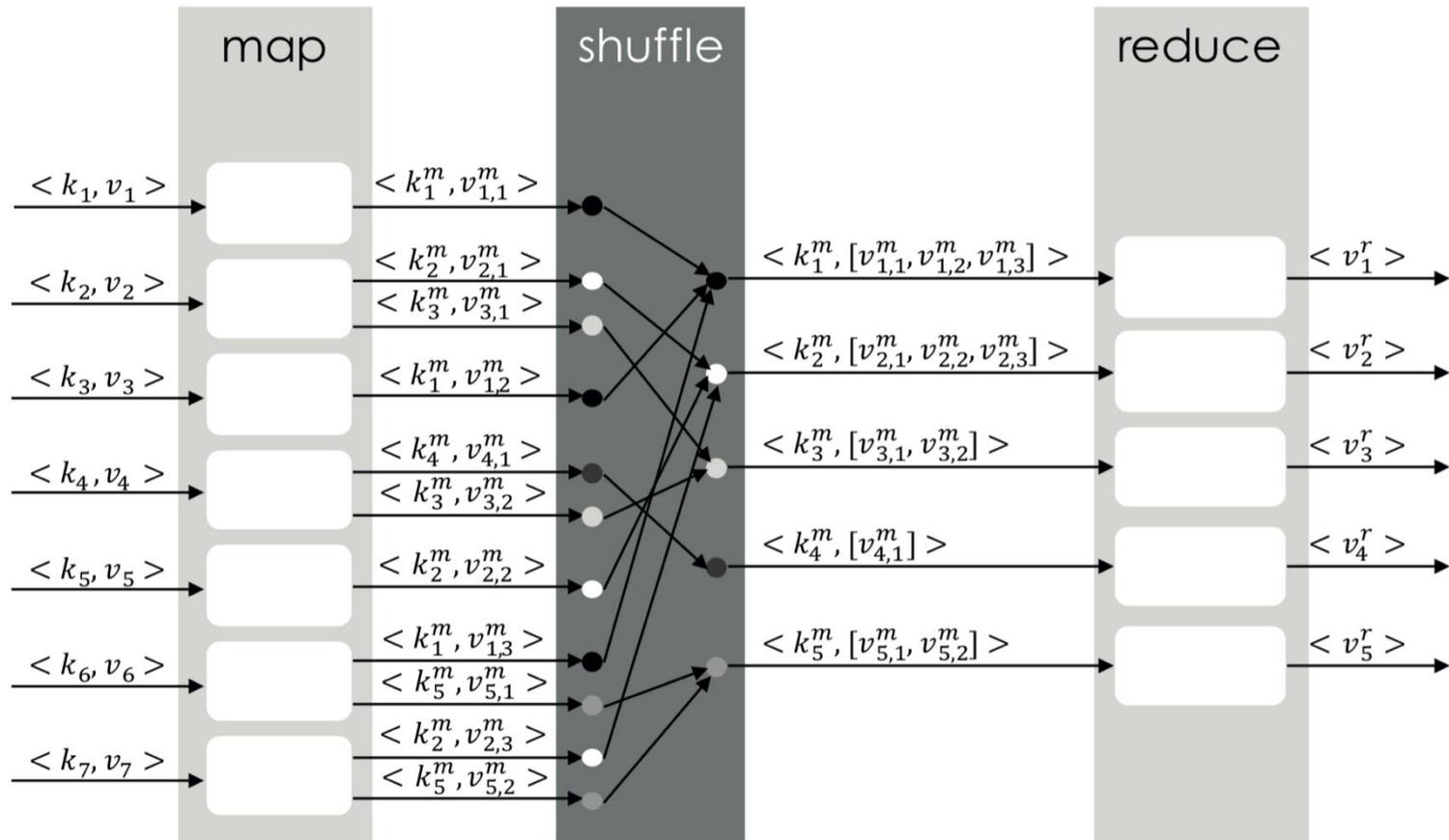


## *El Paradigma MapReduce*

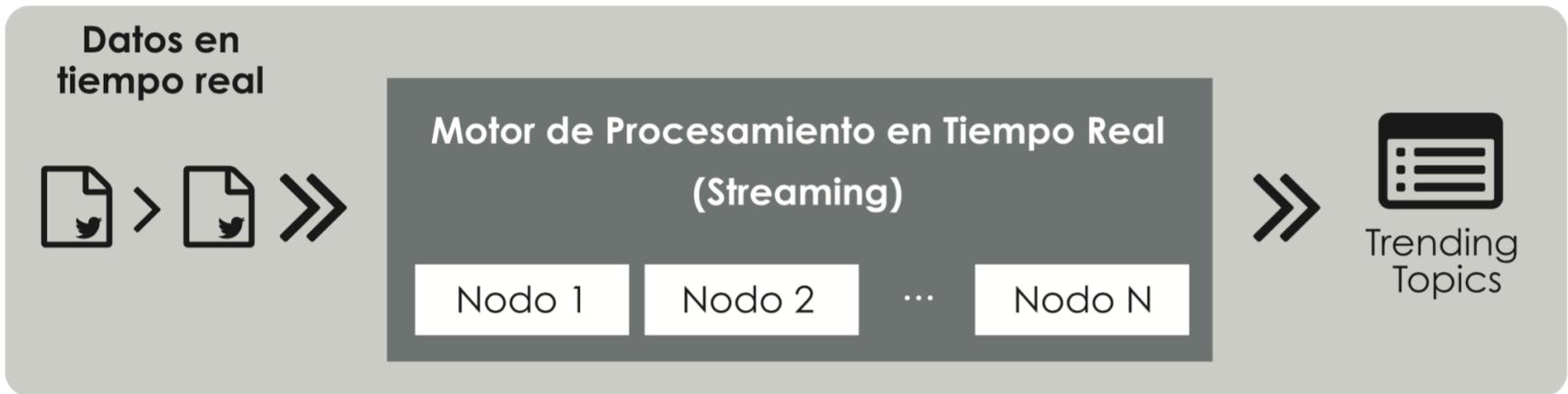
- Consta de dos rutinas principales.
- La rutina **map** realiza una transformación de los datos.
  - ENTRADA: una tupla clave-valor.
  - SALIDA: una o varias tuplas clave-valor.
- Una fase intermedia denominada **shuffle** se encarga de agrupar la salida del map por claves.
- La rutina **reduce** realiza una agrupación o agregación de los datos:
  - ENTRADA: una clave, junto con todos sus valores asociados.
  - SALIDA: un valor.

# Almacenamiento y Procesamiento

## El Paradigma MapReduce



## Procesamiento de Big Data en Tiempo Real (Streaming)



## *La Revolución del Procesamiento de Datos: Apache Spark*

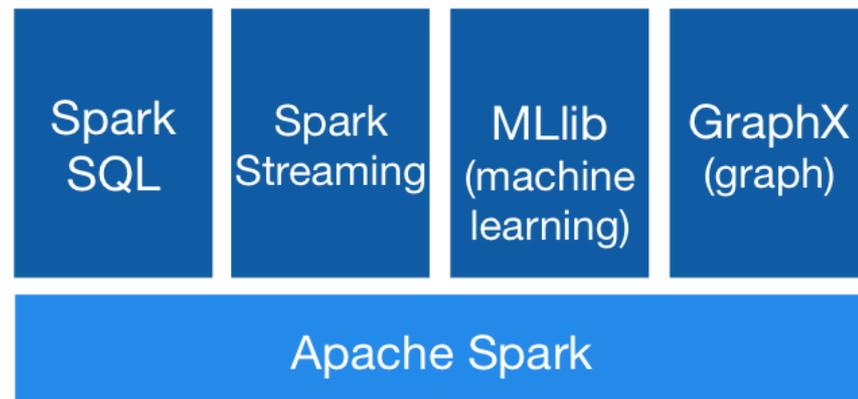
- Apache Spark nace en 2010, y es creado por Matei Zaharia.
- Se le considera el sucesor de Hadoop, y es muy utilizado por la comunidad del Big Data.
- En 2014 Spark bate un record: ordena 100 TB de datos en solo 23 minutos.
- Su clave es el concepto de “RDD”: *Resilient Distributed Dataset*.
  - Es una abstracción de un conjunto de datos.
  - El programador trabaja con él como lo haría normalmente.
  - Pero por detrás, el conjunto de datos está distribuido.



# Almacenamiento y Procesamiento

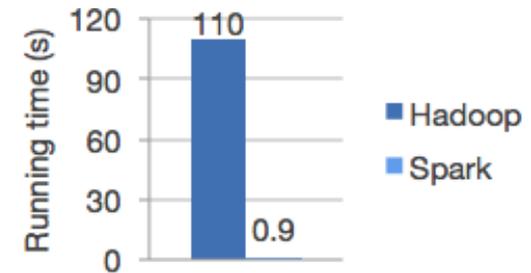
## *Spark: una Solución Polivalente*

- Spark está compuesto por varios módulos que permiten abordar muchos problemas:
  - **Spark Core:** procesamiento por lotes de datos, incluyendo operaciones MapReduce.
  - **Spark Streaming:** procesamiento en tiempo real de Big Data.
  - **Spark MLlib:** aprendizaje automático sobre Big Data.
  - **Spark SQL:** consultas SQL sobre Big Data.
  - **Spark GraphX:** procesamiento y análisis de grandes grafos.



## Spark vs. Hadoop

- Spark se anuncia como el competidor de Hadoop, al ser más rápido.
  - Esto se debe a que realiza optimizaciones en memoria.
- Spark es cómodo para los desarrolladores.
  - Se integra con Java, Scala y Python.
- Spark y Hadoop se complementan:
  - Spark puede instalarse en un cluster de Hadoop.
  - Puede leer datos de HDFS-



```
text_file = spark.textFile("hdfs://...")

text_file.flatMap(lambda line: line.split())
           .map(lambda word: (word, 1))
           .reduceByKey(lambda a, b: a+b)
```

## *Procesamiento por Lotes en la Nube*

- Los servicios de *Cloud Processing* nos permiten:
  - Procesar grandes cantidades de datos.
  - Sin necesidad de invertir en comprar máquinas potentes.
  - Apagar las máquinas (interrumpiendo el gasto) cuando haya finalizado el procesamiento.
  - Escalabilidad: crear más máquinas cuando haya más datos que procesar.
  - Elasticidad: mejorar las capacidades de nuestras máquinas cuando sea necesario.
- Se denomina **instancia** a una máquina virtual que procesa datos en la nube.

## ***Amazon Elastic MapReduce***

- Amazon Elastic MapReduce permite iniciar un *cluster* de Hadoop en minutos.
- Podemos definir el número de máquinas y su rendimiento.
- Además, podemos escoger las herramientas que queremos instalar.k
- Nos permite un procesamiento por lotes rápido y económico



## *Google BigQuery*

- Permite realizar consultas sobre terabytes de datos en pocos minutos.
- Se emplea el lenguaje SQL, lo que simplifica la consulta sobre los datos.
- Es menos flexible que Hadoop MapReduce: los datos deben ser estructurados.
- Permite un procesamiento por lotes cómodo, sencillo y económico.
- Google cobra por el volumen de datos leídos.



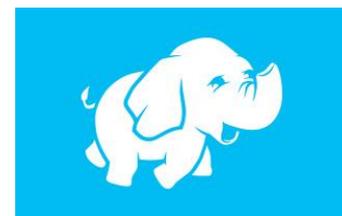
Google BigQuery

## Azure Spark

- Azure es el grupo de servicios en la nube de Microsoft.
- Azure HDInsights proporciona servicios de Big Data: Hadoop, Storm, Spark...
- Azure Spark permite desplegar un *cluster* de Spark en cuestión de minutos.
- Una vez creado, podemos utilizar un *notebook* para desarrollar aplicaciones con Spark.



Microsoft  
Azure



## Esquema de Contenidos:

1. Introducción al Big Dada
2. Trabajando con los datos
3. Almacenamiento y Procesamiento: Técnicas, Herramientas y Plataformas
4. Análisis mediante la Minería de Datos
5. Visualización y Consumo de Datos
6. Seguridad y Gobernanza
7. Aplicaciones Reales de Negocio: Casos de Éxito

## *Análisis de Big Data*

- Analizar datos supone sacar un valor de los mismos:
  - Algo que no se ve a simple vista en los datos.
  - Algo que proporciona un valor de negocio.
  - Algo que nos ayuda a verificar o desmentir una hipótesis.
- Dos tipos de análisis importantes son el **análisis predictivo** y el **análisis de patrones**.
- Se enmarcan dentro del ámbito del **aprendizaje automático** (*machine learning*).

## ***Ejemplos de Problemas de Aprendizaje Automático***

- Detectar si un correo electrónico es *spam* o no lo es.
- Predecir si un gasto con tarjeta de crédito es legítimo o fraudulento.
- Predecir si un cliente va a tener problemas para pagar un crédito en función de su historial bancario.
- Estimar el gasto que va a realizar un cliente en nuestro comercio en función de su perfil demográfico.
- Segmentar a nuestros usuarios o clientes en función de su perfil, para personalizar nuestras ofertas.
- Recomendar un producto a un usuario en función de sus potenciales intereses.

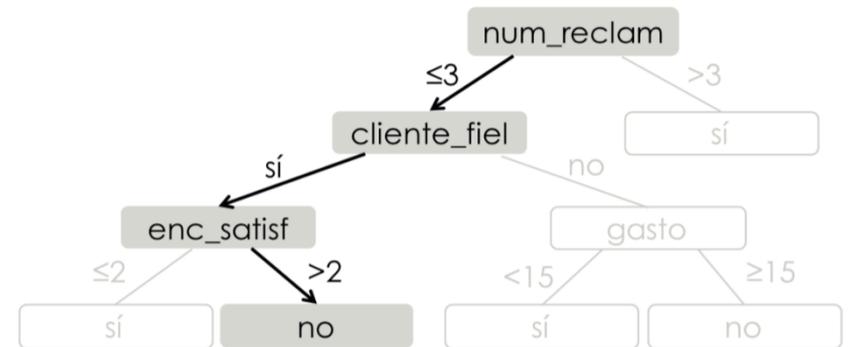
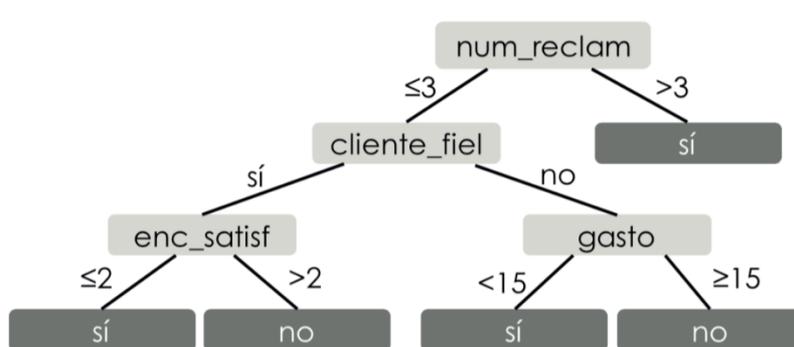
## *Análisis Predictivo*

- El propósito es aprender un modelo a partir de los datos existentes.
- Este modelo servirá para realizar predicciones con nuevos datos.
- Algo de vocabulario:
  - Una **instancia** es un dato.
  - Un **atributo** es una característica sobre nuestros datos.
  - La **clase** (o **salida**) es el valor que queremos predecir o estimar.
- Dos tipos de problemas similares, pero con técnicas distintas:
  - **Clasificación:** la salida es categórica, es decir, una etiqueta de entre varias posibles.
  - **Regresión:** la salida es numérica.

# Análisis mediante la Minería de Datos

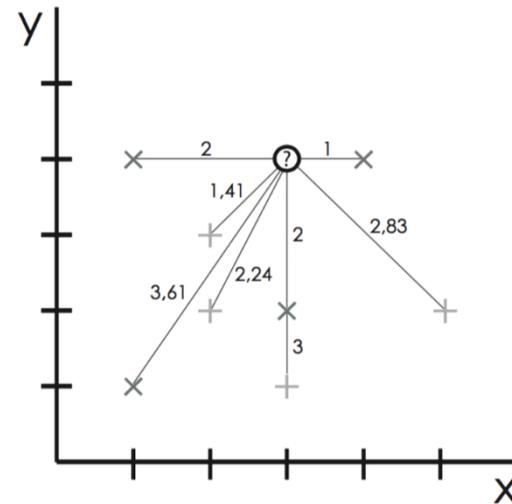
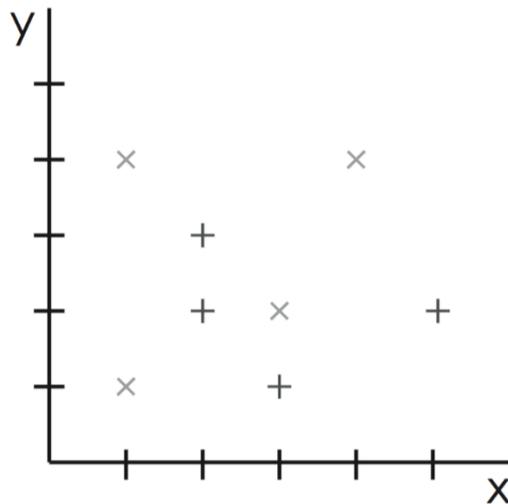
## Análisis Predictivo: Árboles de Decisión / Regresión

- Un árbol de decisión permite ir “haciendo preguntas” a los datos para ir bajando por ramas del árbol.
- Cada nodo del árbol es una pregunta a los datos.
- Cuando llegamos a una hoja (no hay más ramas), entonces tenemos nuestra salida.



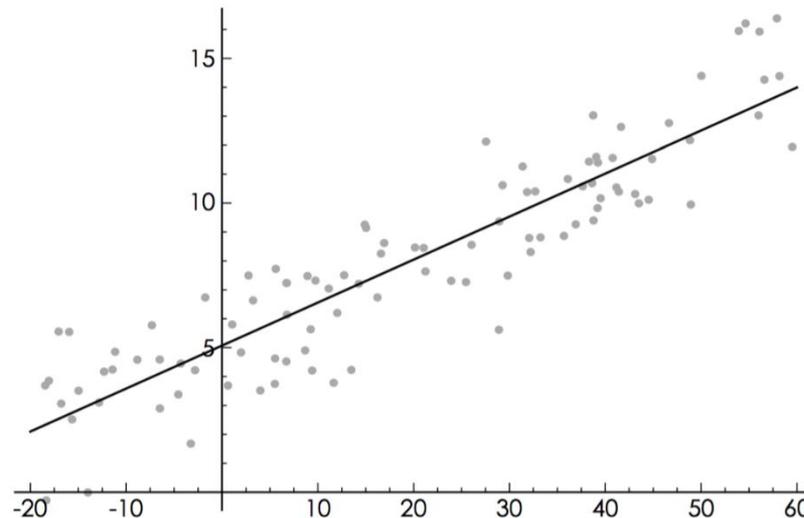
## Análisis Predictivo: Modelos Geométricos

- Se basan en funciones de distancia para estimar la clase o salida a predecir.
- En resumen: ¿cuáles son los datos más parecidos a este que quiero predecir?
- La técnica más empleada es *k-Nearest Neighbors* (k-NN).



## *Análisis Predictivo: Regresión Lineal*

- Dado un conjunto de puntos en el espacio, queremos aprender un plano (o recta) que los ajuste.
- El objetivo es minimizar el error, es decir, la distancia entre el plano y los puntos.
- La ecuación de este plano (o recta) nos servirá para estimar la salida de los nuevos datos.



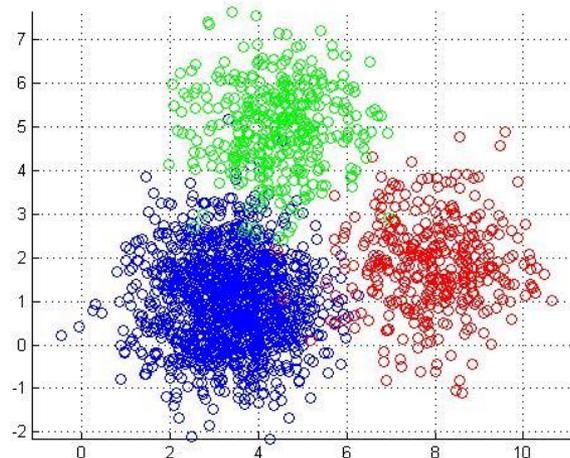
## *Análisis de Patrones*

- No se busca predecir una clase o salida.
- En su lugar, se busca agrupar datos que son parecidos.
- Ejemplo: ¿qué usuarios son parecidos entre sí?
- También se denomina **clustering**, agrupación o segmentación.

# Análisis mediante la Minería de Datos

## *Análisis de Patrones*

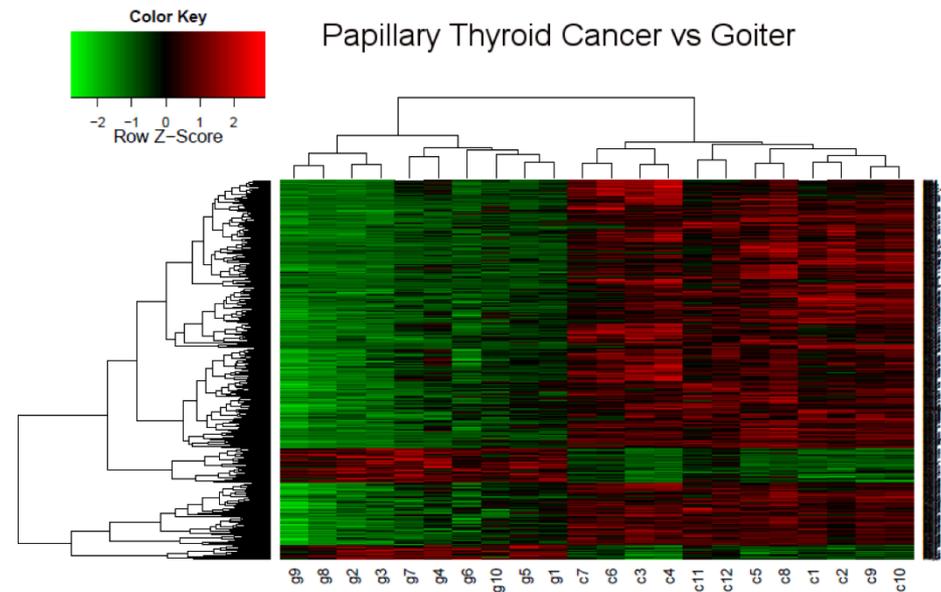
- No se busca predecir una clase o salida.
- En su lugar, se busca agrupar datos que son parecidos.
- Ejemplo: ¿qué usuarios son parecidos entre sí?
- También se denomina **clustering**, agrupación o segmentación.
- Se tiene que definir una medida de distancia (métrica) entre los datos.



# Análisis mediante la Minería de Datos

## Análisis de Patrones: técnicas jerárquicas

- Pueden ser **aglomerativas** o **divisivas**.
- En las aglomerativas, se van uniendo los datos más cercanos de dos en dos en cada paso.
- En las divisivas, todos los datos comienzan en un único **cluster** y se van separando en cada paso.



## *Recomendación*

- La recomendación combina análisis de patrones y análisis predictivo.
- Por un lado, queremos segmentar usuarios o clientes parecidos entre sí.
- Por otro lado, queremos segmentar productos parecidos entre sí.
- A continuación, realizamos predicción: ¿qué producto le gustará a un usuario?
- Para ello, miramos los usuarios parecidos a él, y los productos que han adquirido.
- Existen dos enfoques: **colaborativo** o **basado en contenidos**.

# Análisis mediante la Minería de Datos

## Análisis Escalable de Big Data

- El análisis de Big Data introduce nuevos desafíos.
- El software clásico de aprendizaje automático puede dar problemas cuando:
  - Los conjuntos de datos tienen muchas instancias.
  - Los datos tienen muchos atributos.
- Existe software específico para realizar aprendizaje automático con Big Data:
  - Apache Mahout
  - Apache Spark Mlib
  - Apache Flink
  - H2O



H<sub>2</sub>O – The Killer-App for Spark

MLlib	H <sub>2</sub> O	SQL	<b>In-Memory</b>	Big Data, Columnar
<b>H<sub>2</sub>ORDD</b>			<b>ML</b>	100x faster Algos
<b>HDFS=DATA</b>			<b>R</b>	CRAN, API, fast engine
			<b>API</b>	Spark API, Java MM
			<b>Community</b>	Devs, Data Science



## *Análisis de Big Data en la Nube*

- Otra opción es emplear herramientas para realizar aprendizaje automático en la nube.
- Los grandes competidores ofrecen herramientas de este tipo:
  - Amazon ML.
  - Google Cloud Prediction.
  - Microsoft Azure ML.
- Una opción más intuitiva es BigML.

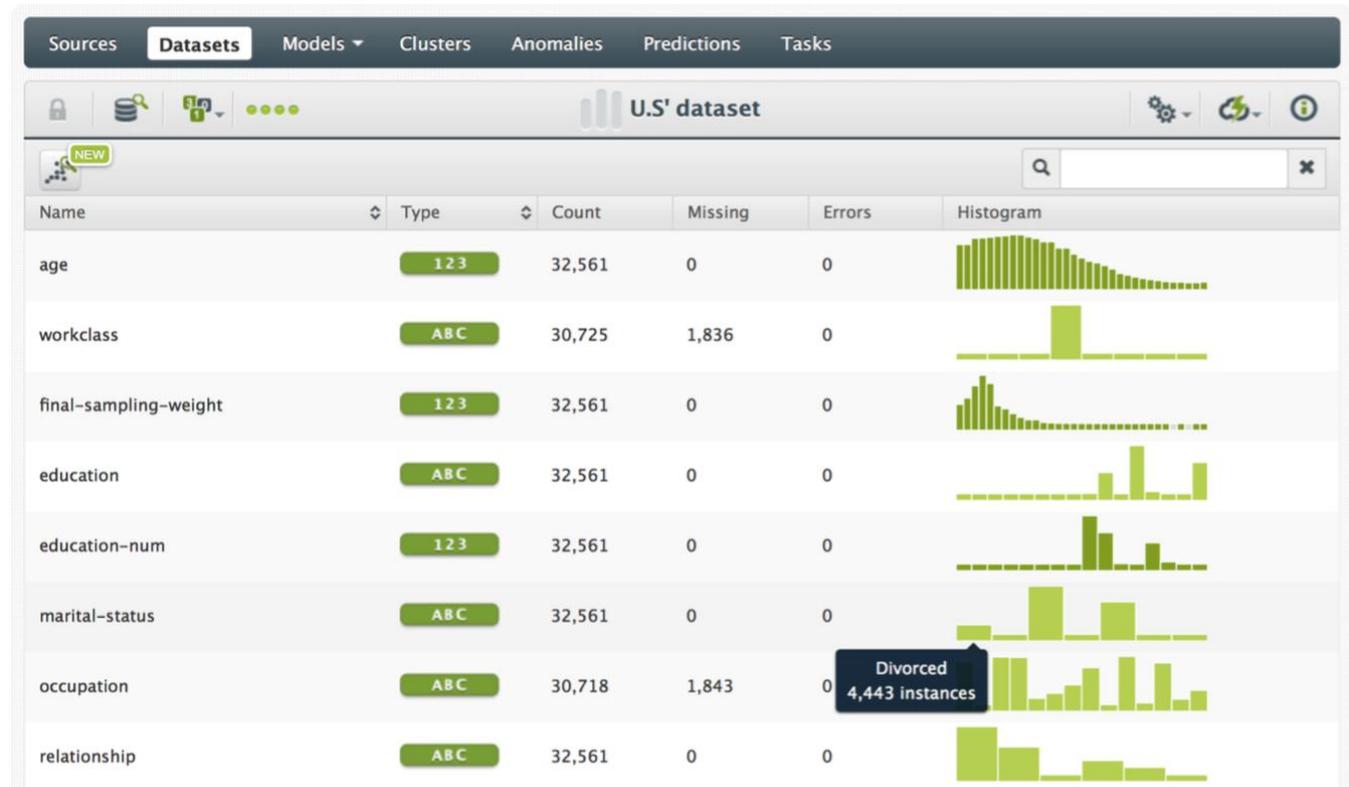
<https://bigml.com/>



# Análisis mediante la Minería de Datos

## BigML: Gestión de Datos

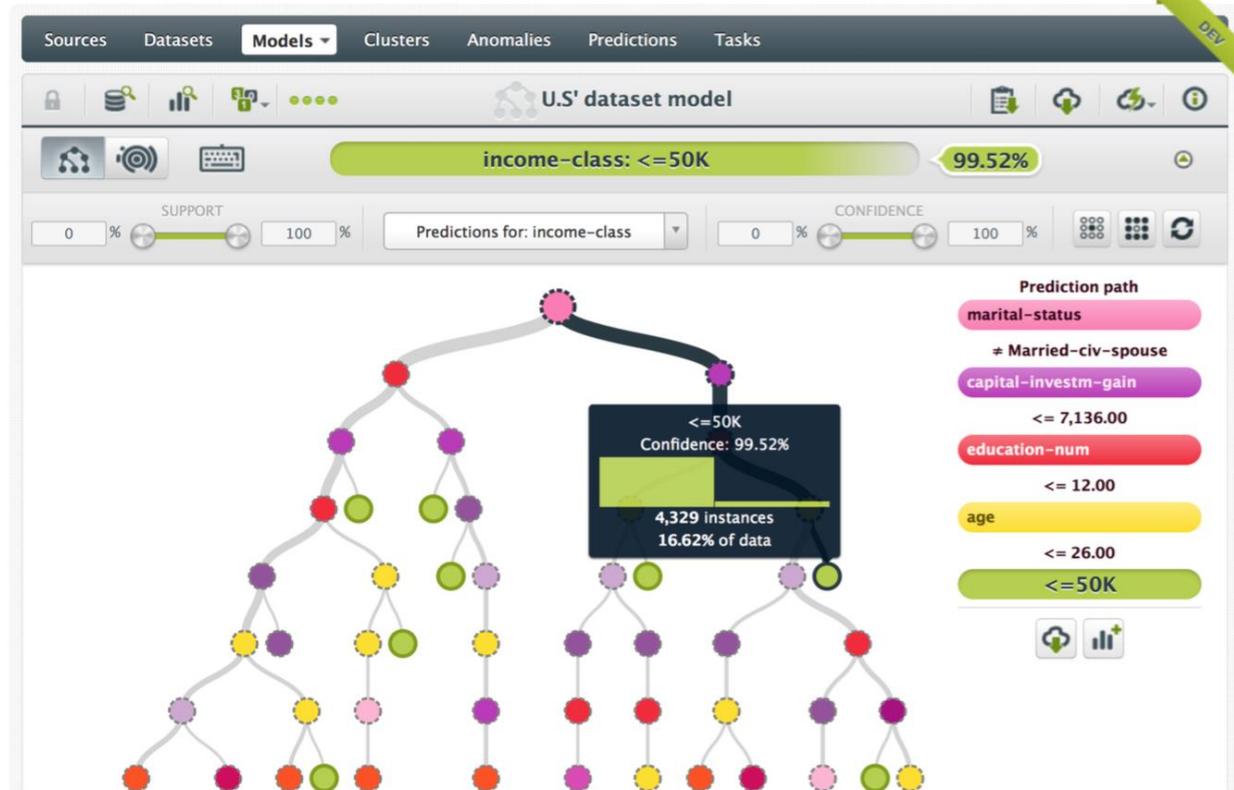
- BigML permite gestionar los datos, mostrando las distribuciones de valores de cada atributo.



# Análisis mediante la Minería de Datos

## BigML: Análisis Predictivo

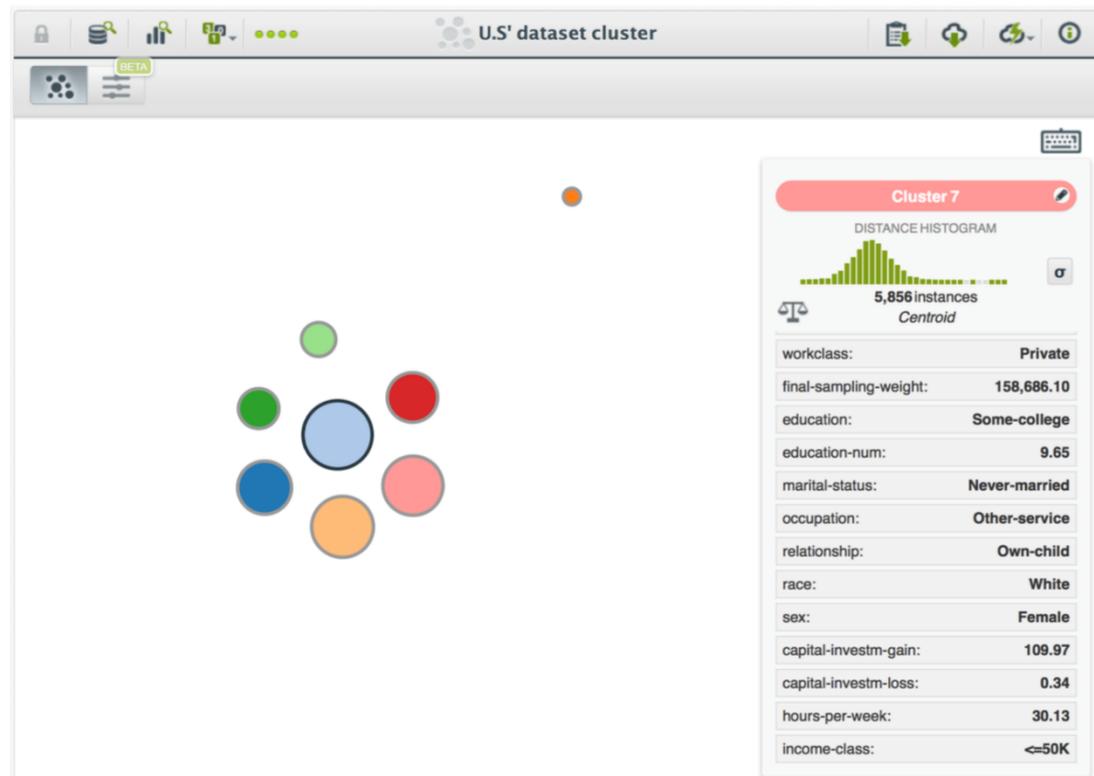
- BigML permite aprender automáticamente árboles de decisión a partir de datos y visualizarlos.



# Análisis mediante la Minería de Datos

## BigML: Análisis Predictivo

- BigML permite realizar tareas de agrupación de datos para encontrar *clusters* similares.



## Esquema de Contenidos:

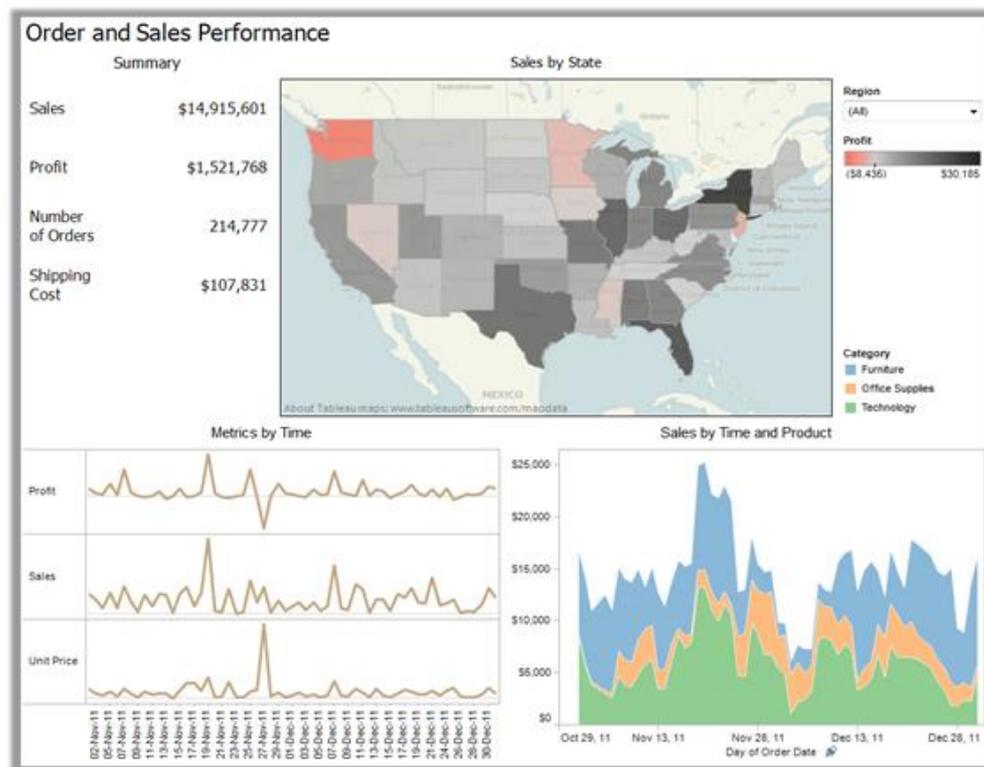
1. Introducción al Big Dada
2. Trabajando con los datos
3. Almacenamiento y Procesamiento: Técnicas, Herramientas y Plataformas
4. Análisis mediante la Minería de Datos
5. Visualización y Consumo de Datos
6. Seguridad y Gobernanza
7. Aplicaciones Reales de Negocio: Casos de Éxito

## *Objetivos de la visualización de datos*

- Aprovechar la habilidad humana de extraer patrones a partir de imágenes.
- Ayudar al usuario a comprender más rápidamente patrones descubiertos automáticamente por un sistema automático de extracción del conocimiento o KDD (Knowledge Discovery in Databases)

# Visualización y Consumo de Datos

- Ejemplo de panel de monitorización de ventas mediante la herramienta Tableau.
- Se muestran diferentes medidas (ventas, beneficios, pedidos y costes de envío) sobre un mapa, medidas (beneficios, ventas, precio unitario) a lo largo del tiempo y medida de ventas a lo largo de tiempo por categoría de productos

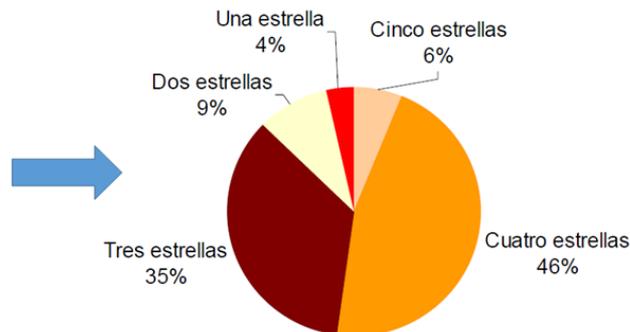
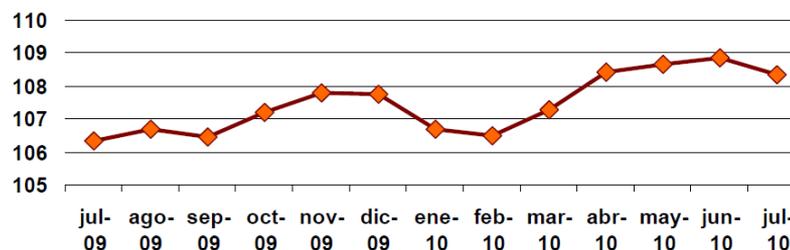
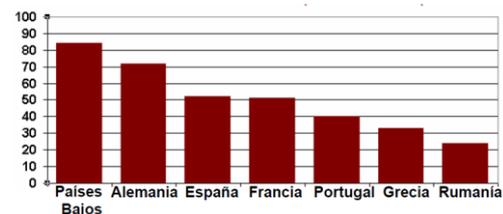


# Visualización y Consumo de Datos

## Tipos Fundamentales de Gráficos Estadísticos

- Gráfico de Barras
- Gráfico de Líneas
- Gráfico de Sectores
- Pictograma
- Gráfico de Dispersión
- Cartograma

Orientación vertical y orden por frecuencias



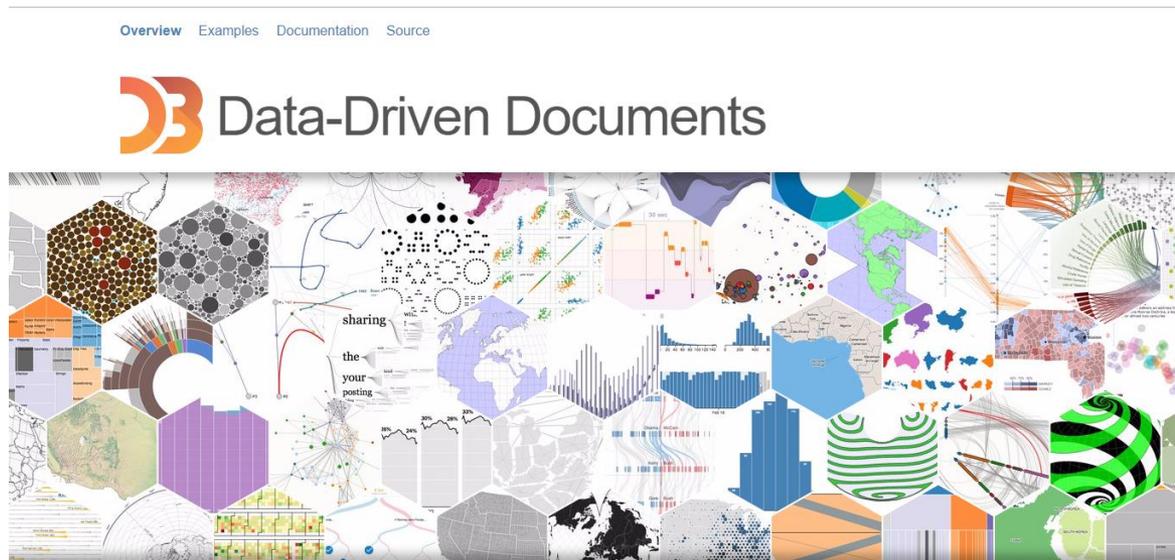
## *Nuevo paradigma*

- El uso de los nuevos métodos de visualización se amplía, concediendo a la interfaz gráfica un rol protagonista en ámbitos aparentemente tan alejados de los negocios como el del arte y la cultura.
- Tanto la interfaz como las nuevas herramientas de visualización constituyen ya mucho más que meros elementos (aunque con funciones destacadas) en el ámbito cultural, conformando un nuevo paradigma.

# Visualización y Consumo de Datos

## Librerías y APIs para la Visualización de Datos

- Google Chart Tools (<http://chart.apis.google.com>)
- JavaScript InfoVis Toolkit
- D3.js (<http://d3js.org/>)



The screenshot shows the D3.js website. At the top, there are navigation links: Overview, Examples, Documentation, and Source. Below this is the D3.js logo, which consists of a stylized orange 'D' and '3', followed by the text 'Data-Driven Documents'. The main content area is a collage of various data visualizations, including bar charts, pie charts, maps, and network diagrams. Below the collage, there is a paragraph of text describing D3.js as a JavaScript library for manipulating documents based on data. To the right of this text is a link that says 'See more examples.'. Below the text, there is a link to download the latest version (4.7.4) and a list of download options, including 'd3.zip'. At the bottom, there is a snippet of code to link directly to the latest release.

Overview Examples Documentation Source

## D3 Data-Driven Documents

D3.js is a JavaScript library for manipulating documents based on data. D3 helps you bring data to life using HTML, SVG, and CSS. D3's emphasis on web standards gives you the full capabilities of modern browsers without tying yourself to a proprietary framework, combining powerful visualization components and a data-driven approach to DOM manipulation.

See more examples.

Download the latest version (4.7.4) here:

- [d3.zip](#)

To link directly to the latest release, copy this snippet:

# Visualización y Consumo de Datos

## Webmapping: Visualización de Datos en Mapas Web

- Cada vez son más las aplicaciones que consumen grandes volúmenes de datos referentes a la “geolocalización”. Herramientas
  - OpenLayers (<http://openlayers.org/>)
  - Leaflet
  - CartoBD

estacionamientos m...

PUBLIC Map not published yet

LAYERS WIDGETS

ADD

sedeswifi

0 ANALYSES 2 WIDGETS

sedeswifi

trafico\_rodado\_veh...

0 ANALYSES 0 WIDGETS

trafico\_rodado\_vehiculos\_id

dismuni\_t

0 ANALYSES 0 WIDGETS

dismuni\_t

poi\_carriles

0 ANALYSES 0 WIDGETS



## Esquema de Contenidos:

1. Introducción al Big Dada
2. Trabajando con los datos
3. Almacenamiento y Procesamiento: Técnicas, Herramientas y Plataformas
4. Análisis mediante la Minería de Datos
5. Visualización y Consumo de Datos
6. Seguridad y Gobernanza
7. Aplicaciones Reales de Negocio: Casos de Éxito

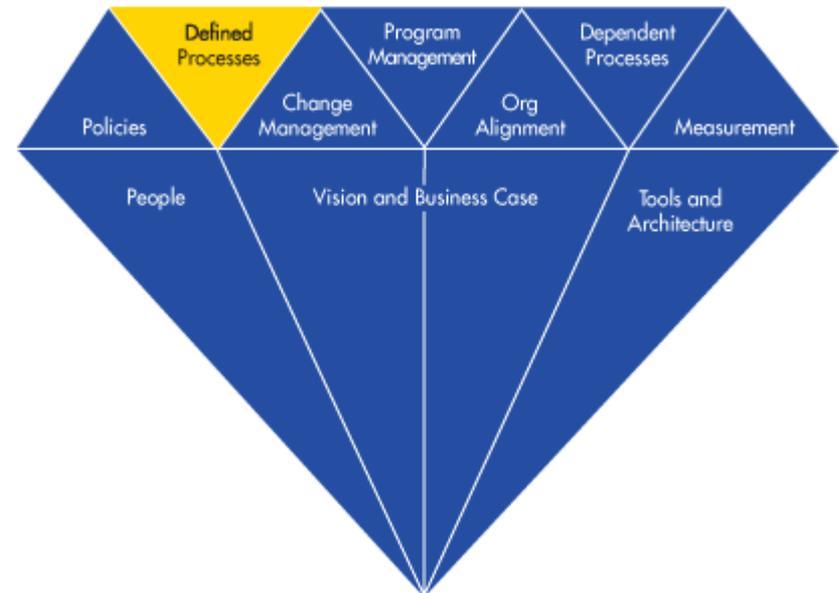
# Seguridad y Gobernanza

Según el estudio 'Transforming the Data Center', elaborado por Cloudera e Intel, la seguridad y la gobernanza de los datos se han convertido en los grandes desafíos de Big Data. Así lo piensan el 60% de los participantes en esta investigación.



## Gobernanza de datos integral

- La gobernanza de datos integral significa gestionar los datos por completo en todos los silos organizativos, arquitectónicos y políticos de la empresa. Requiere adaptar las personas, el proceso, **las políticas y la tecnología para garantizar la entrega de datos fiables y seguros**, de modo que se pueda **cumplir las normativas del sector, reducir el coste** que supone hacer negocios y expandir el negocio.



## Esquema de Contenidos:

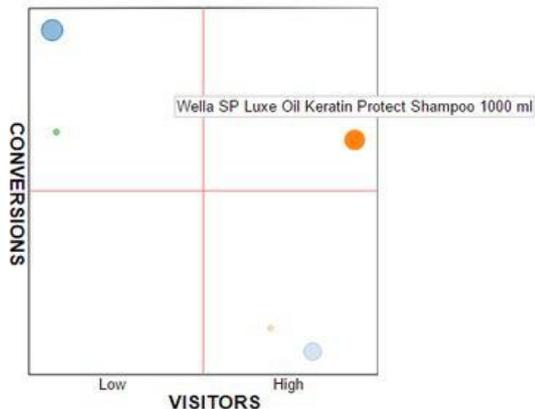
1. Introducción al Big Dada
2. Trabajando con los datos
3. Almacenamiento y Procesamiento: Técnicas, Herramientas y Plataformas
4. Análisis mediante la Minería de Datos
5. Visualización y Consumo de Datos
6. Seguridad y Gobernanza
7. Aplicaciones Reales de Negocio: Casos de Éxito

# Use Case 1: Clickstream

## Analizando la Huella Digital en Comercio Electrónico

📅 26 marzo, 2016 👤 José Manuel Garcia Nieto

Dentro del contexto del comercio electrónico, cuando hablamos de la “huella digital” nos referimos a la secuencia de pasos, en forma de “clicks”, y acciones que realiza un cliente a lo largo de su visita a una tienda virtual. Imagine que entra en su sitio de comercio electrónico favorito con la intención de comprar un determinado producto de su catálogo. Por ejemplo, estamos interesados en comprar un ordenador portátil.



<https://elblogdelainnovaciondigital.wordpress.com/2016/03/26/analizando-la-huella-digital-en-comercio-electronico/>

# Use Case 1: Clickstream



GET STARTED

Descargue sandbox

TUTORIAL SERIES

REAL WORLD EXAMPLES

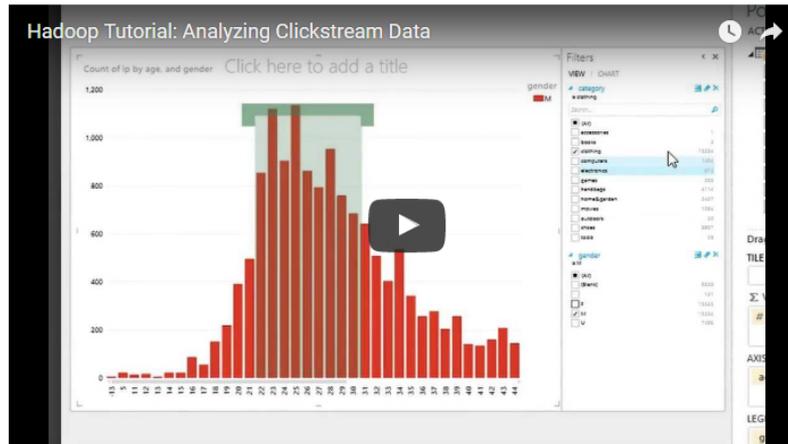
1. Indexing and Searching text within images with Apache Solr
2. Incremental Backup of Data from HDP to Azure using Falcon for Disaster Recovery and Burst capacity
3. Realtime Event Processing in Hadoop with NIFI, Kafka and Storm
4. **Visualize Website Clickstream Data**
  - Introduction
  - Lab 1: Perform Web Log Analysis with Hive
  - Lab 2: Visualize Clickstream Logs with Excel
5. Refine and Visualize Server Log Data
6. Analyzing Social Media and Customer Sentiment With Apache NIFI and HDP Search
7. Analyze HVAC Machine and Sensor Data
8. Natural Language Processing and Sentiment Analysis for Retailers using HDP and ITC Infotech Radar

## VISUALIZE WEBSITE CLICKSTREAM DATA

### INTRODUCTION

Your home page looks great. But how do you move customers on to bigger things—like submitting a form or completing a purchase? Get more granular with customer segmentation. Hadoop makes it easier to analyze, visualize and ultimately change how visitors behave on your website.

In this demo, we demonstrate how an online retailer can optimize buying paths to reduce bounce rate and improve conversion.

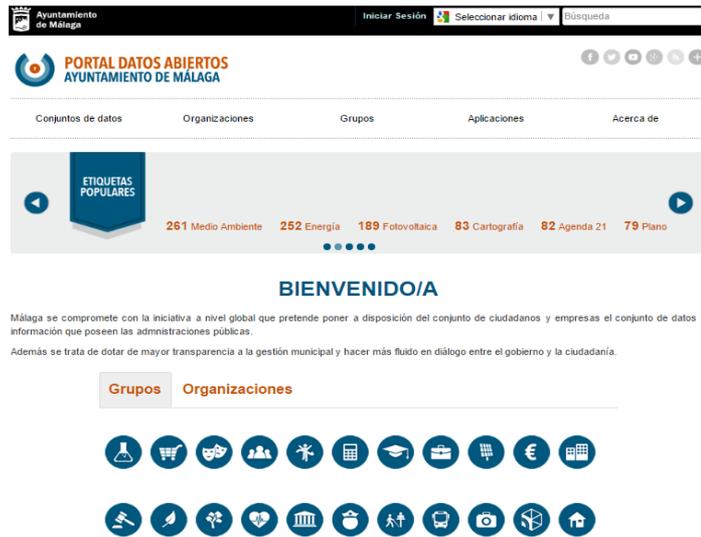


[https://github.com/hortonworks/hadoop-tutorials/blob/master/Community/T01\\_RHadoop\\_visitors\\_prediction.md](https://github.com/hortonworks/hadoop-tutorials/blob/master/Community/T01_RHadoop_visitors_prediction.md)

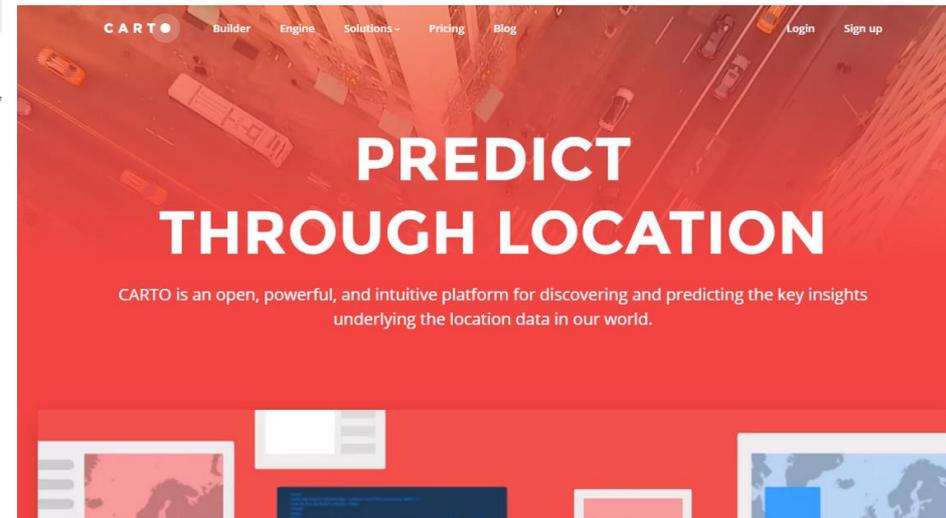
# Use Case 2: Open Data Visualization with Carto

Open Data from

<http://datosabiertos.malaga.eu/>



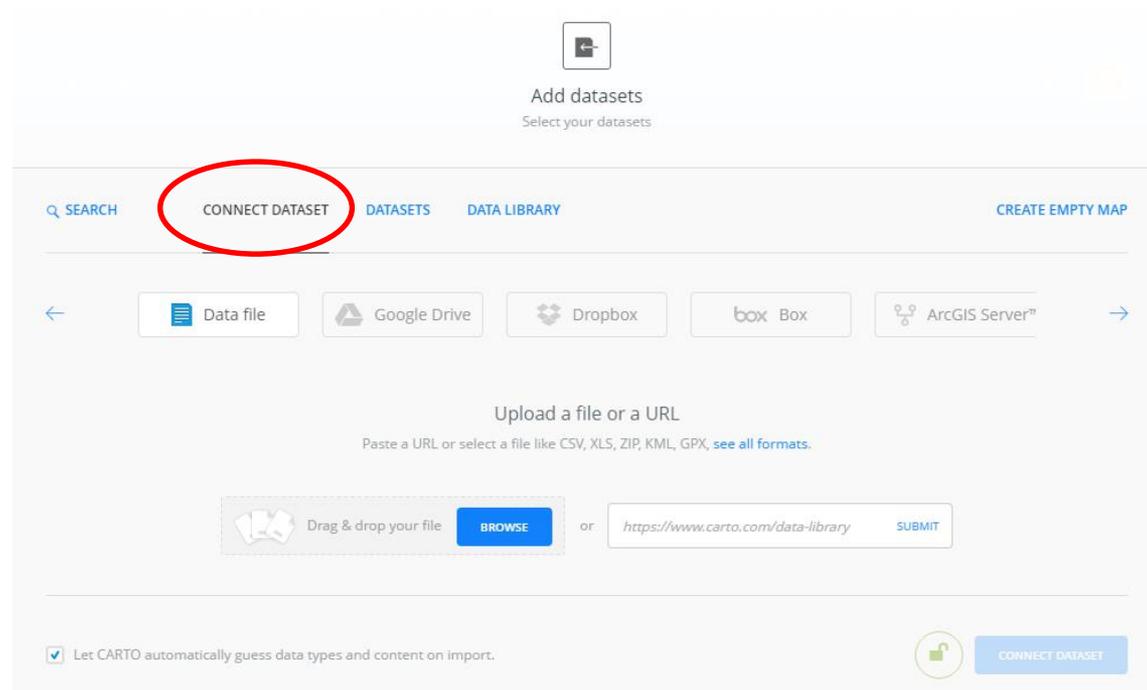
Registration on Carto <https://carto.com/>



David Bueno (CEMI) : <https://www.youtube.com/watch?v=DQjNZ4baSsw>

# Use Case 2: Open Data Visualization with Carto

- Connecting Dataset



# Use Case 2: Open Data Visualization with Carto

- Connecting Dataset: Málaga Bici

The screenshot shows the 'PORTAL DATOS ABIERTOS AYUNTAMIENTO DE MÁLAGA' website. The search bar contains 'Málaga bici' and shows '4 conjuntos de datos encontrados para "Málaga bici"'. The results list 'Málaga bici' as the first item, with a description: 'Información relacionada con el sistema de bike sharing de la Empresa Malagueña de Transportes (EMT), que ofrece la posibilidad de poner a disposición de los ciudadanos que...'. A red oval highlights this search result. The left sidebar shows navigation options like 'Organizaciones', 'Grupos', and 'Etiquetas'.

The screenshot shows the 'Málaga bici' dataset page. The page title is 'Málaga bici' and it shows '2 Seguidores'. The 'Datos y Recursos' section includes a 'CSV' icon and the text 'Estacionamientos CSV' and 'Estado de los estacionamientos en formato CSV'. A red oval highlights the 'Explorar' button in the top right corner. Below this, there are tags for 'aparcamiento', 'bici', 'bicicleta', and 'estacionamiento'. The 'Información Adicional' section contains a table with the following data:

Campo	Valor
Autor	Responsable del Área de Accesibilidad y Movilidad
Descargas recientes	36
Descargas totales	1854
Fecha última actualización del fichero	06 Enero 2017, 17:02
Frecuencia	1 minuto

# Use Case 2: Open Data Visualization with Carto

- Connecting Dataset: Málaga Bici

**Estacionamientos CSV** [Ir al recurso](#) [API de datos](#)

URL: <http://datosabiertos.malaga.eu/recursos/transporte/EMT/estacionamientos/Estacionamientos.csv>

Estado de los estacionamientos en formato CSV.

Data Explorer | Mapa | Tabla | Gráfico

[Incrustar](#)

Add Filter

« 1 – 23 » 23 records

Search data ... Go »

_id	ID	NOMBRE	NOMBR...	DIRECC...	ID_EST...	NOMBR...	NUM_D...	NUM_LI...	NUM_O...	LAT	I
19	13	19-Arroy...	MALAGA	Arroyo d...	5	Con inci...	255	7	14	36.731595	
20	4	20-Barb...	MALAGA	Avenida ...	5	Con inci...	255	12	12	36.711781	
3	5	03-Pz. ...	MALAGA	Plaza de...	5	Con inci...	255	10	11	36.723308	
11	9	11-Cdad...	MALAGA	Pollideo...	5	Con inci...	255	15	8	36.749820	
12	15	12-Rect...	MALAGA	Av. Cerva...	5	Con inci...	255	13	10	36.719897	
13	26	13-Av. A...	MALAGA	Avda An...	5	Con inci...	255	7	16	36.717187	
14	2	14-Est. ...	MALAGA	Estación...	5	Con inci...	255	10	6	36.71276	
2	18	02-Pz. T...	MALAGA	Plaza de...	5	Con inci...	255	20	9	36.720893	
4	11	04-C.A.C.	MALAGA	Av. Com...	5	Con inci...	255	15	6	36.71543	

Q SEARCH CONNECT DATASET DATASETS DATA LIBRARY CREATE EMPTY MAP

Data file Google Drive Dropbox box Box ArcGIS Server

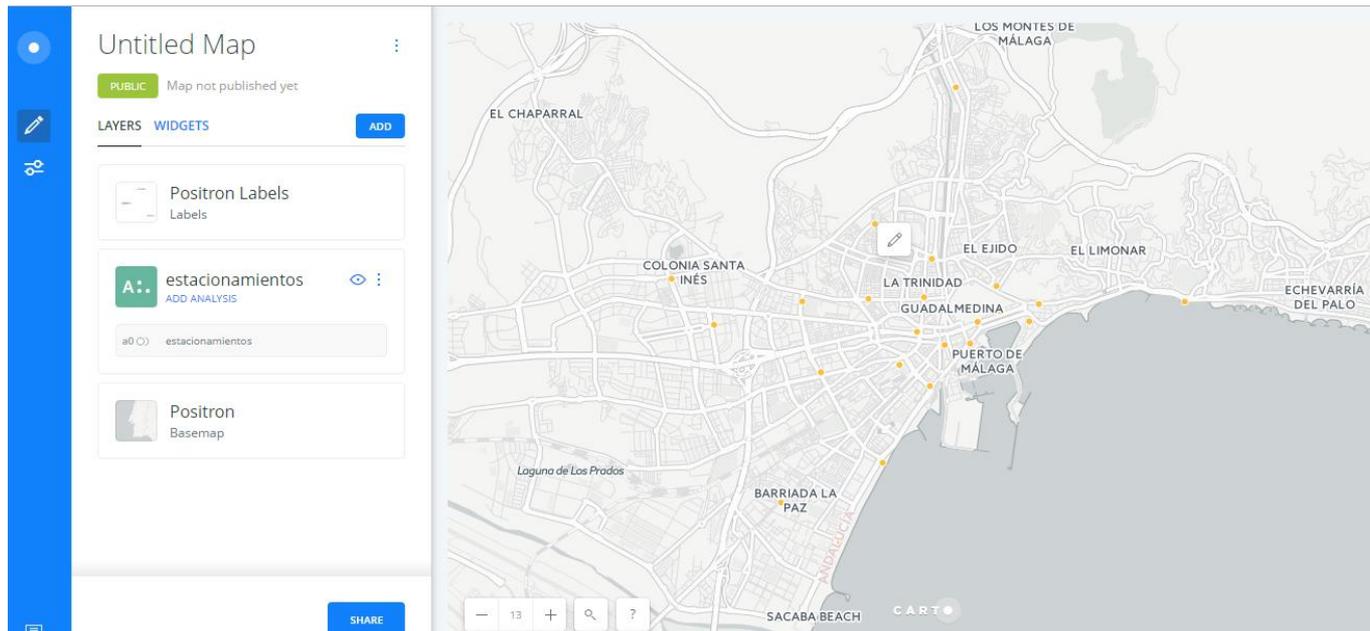
Upload a file or a URL  
Paste a URL or select a file like CSV, XLS, ZIP, KMZ, GPX, see all formats.

Drag & drop your file BROWSE <http://datosabiertos.malaga.eu/recursos> SUBMIT

Let CARTO automatically guess data types and content on import. [CONNECT DATASET](#)

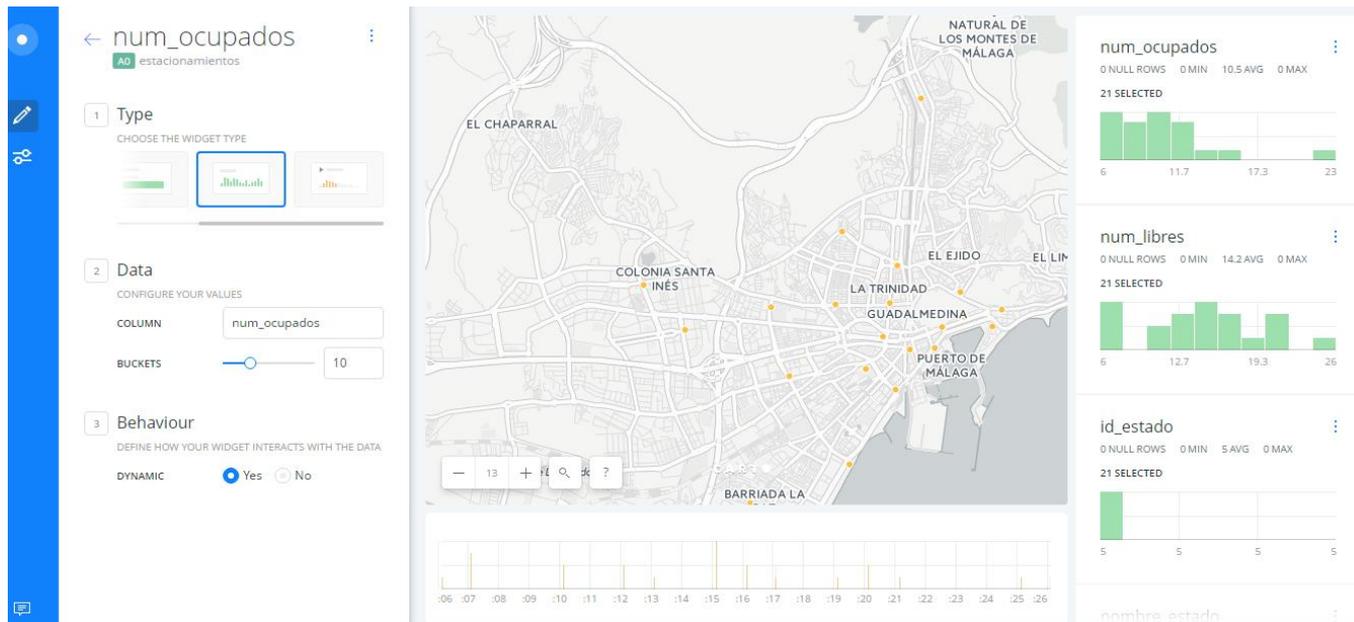
# Use Case 2: Open Data Visualization with Carto

- Connecting Dataset: Málaga Bici



# Use Case 2: Open Data Visualization with Carto

- Adding Widgets: Málaga Bici



# Use Case 2: Open Data Visualization with Carto

- Adding new layer: Málaga Carriles Bici

PORTAL DATOS ABIERTOS  
AYUNTAMIENTO DE MÁLAGA

Conjuntos de datos | Organizaciones | Grupos | Aplicaciones | Acerca de

/ Conjuntos de datos

Organizaciones

- PRESIDENCIA (4)
- ACCESIBILIDAD Y MOV... (1)

Grupos

- Urbanismo e infraes... (4)
- Transporte (1)

Etiquetas

- territorio (4)
- agenda 21 (4)

carriles

5 conjuntos de datos encontrados para "carriles" Ordenar por: Relevancia

**Carriles bici**  
Carriles bici en Málaga  
GeoJSON KML CSV

**Agenda 21 - Carriles Bici 2013**  
Territorio y configuración de la ciudad. Localización de los carriles bici existentes en la

Carriles bici

Seguidores 0

Organización

ACCESIBILIDAD Y MOVILIDAD  
Accesibilidad y movilidad leer más

Social

Conjunto de datos | Grupos | Flujo de Actividad

**Carriles bici**

Datos y Recursos

- Carriles bici GeoJSON**  
Ubicación de los carriles bici en formato GeoJSON
- Carriles bici KML**  
Ubicación de los carriles bici en formato KML
- Carriles bici CSV**  
Ubicación de los carriles bici en formato CSV

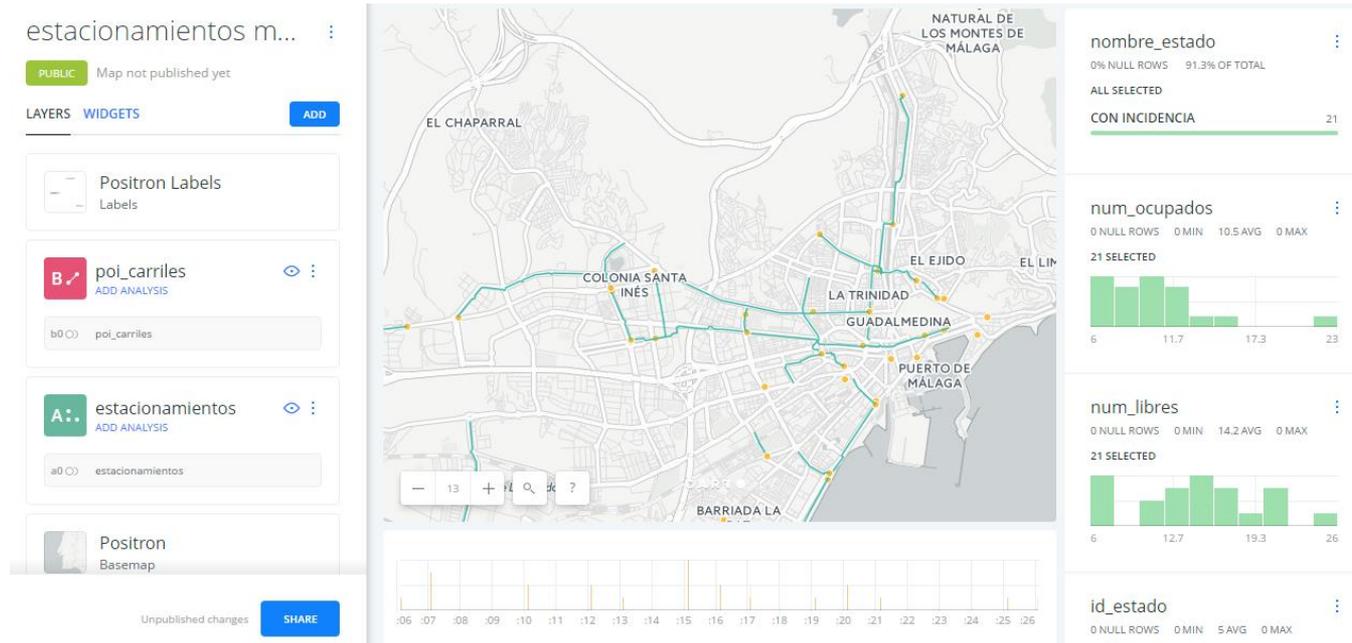
bici carriles movilidad transporte

Información Adicional

Campo	Valor
Fuente	<a href="http://movilidad.malaga.eu/export/sites/default/movilidad/trafico/portal/menu/seccion_0003/documentos/carriles_bici.pdf">http://movilidad.malaga.eu/export/sites/default/movilidad/trafico/portal/menu/seccion_0003/documentos/carriles_bici.pdf</a>
Autor	Responsable de movilidad

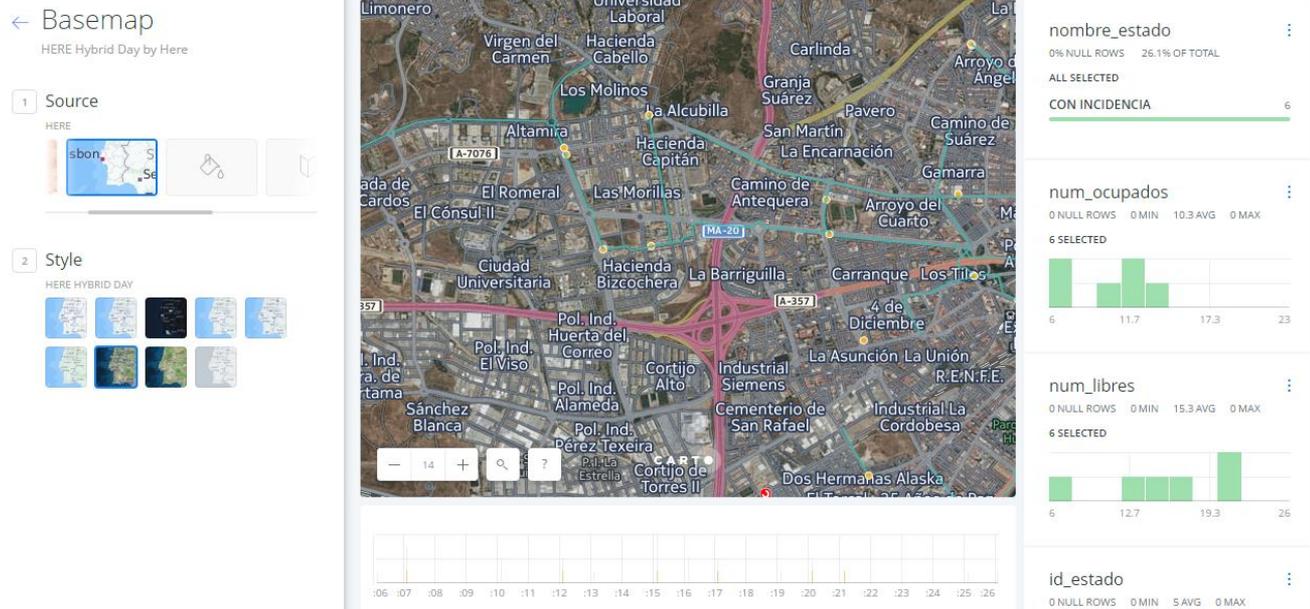
# Use Case 2: Open Data Visualization with Carto

- Adding new layer: Málaga Carriles Bici



# Use Case 2: Open Data Visualization with Carto

- Changing Basemap: Málaga Carriles Bici



# Use Case 2: Open Data Visualization with Carto

- Adding new layer: Málaga Traffic Noise

The screenshot shows the Carto interface for a specific dataset. On the left, there is a sidebar with the following information:

- Información SIG Mapa Estratégico de Ruido de Málaga - Trafico rodado vehiculos Índice Ldía**
- Seguidores: **0**
- Organización: MEDIO AMBIENTE Y SOSTENIBILIDAD (with a leaf logo)

The main content area displays the dataset details:

- Navigation tabs: Conjunto de datos, Grupos, Flujo de Actividad
- Dataset title: **Información SIG Mapa Estratégico de Ruido de Málaga - Trafico rodado vehiculos Índice Ldía**
- Section: **Datos y Recursos**
- Dataset name: **Trafico rodado vehiculos Índice Ldía shp** (with a data icon)
- Description: Trafico rodado vehiculos Índice Ldía en formato shp
- Buttons: Explorar, Más información, Ir al recurso
- Tags: ambiente, ruido, tráfico, vehiculos

Below the dataset information, there is an **Información Adicional** section with the following table:

Campo	Valor
Autor	Responsable del Área de Medio Ambiente y Sostenibilidad
Descargas recientes	1
Descargas totales	32

# Use Case 2: Open Data Visualization with Carto

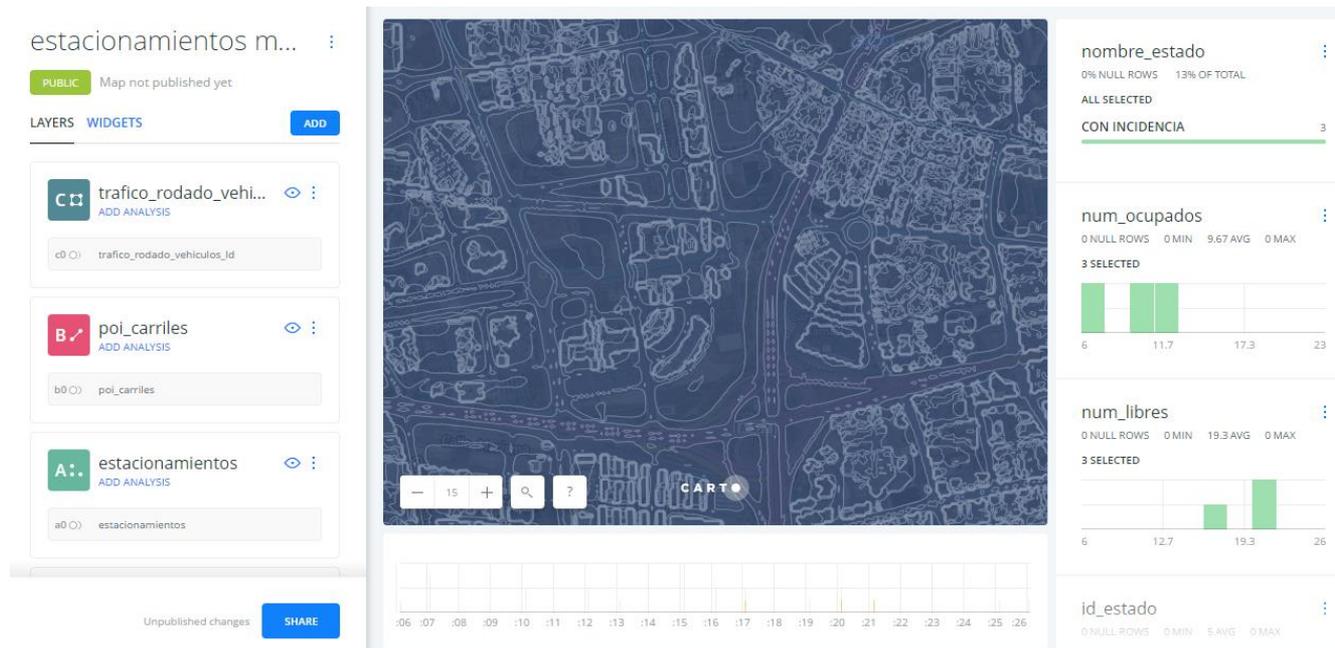
- Adding new layer: Málaga Traffic Noise

The screenshot shows the 'Portal Datos Abiertos Ayuntamiento de Málaga' website. The main content area displays the dataset 'Trafico rodado vehiculos Índice Ldía shp'. A blue button labeled 'Ir al recurso' is visible. Below the title, the URL is provided: <http://datosabiertos.malaga.eu/storage/f/2013-10-31T06%3A40%3A14.968Z/trafico-rodado-vehiculos-ld.zip>. The description states 'Trafico rodado vehiculos Índice Ldía en formato shp' and notes 'Todavía no existen vistas creadas para este recurso.' On the left, there are tabs for 'Recursos', 'Social', and 'Google+'. The 'Recursos' tab is active, showing a table with the title 'Información adicional'.

Campo	Valor
Última actualización	Junio 12, 2014

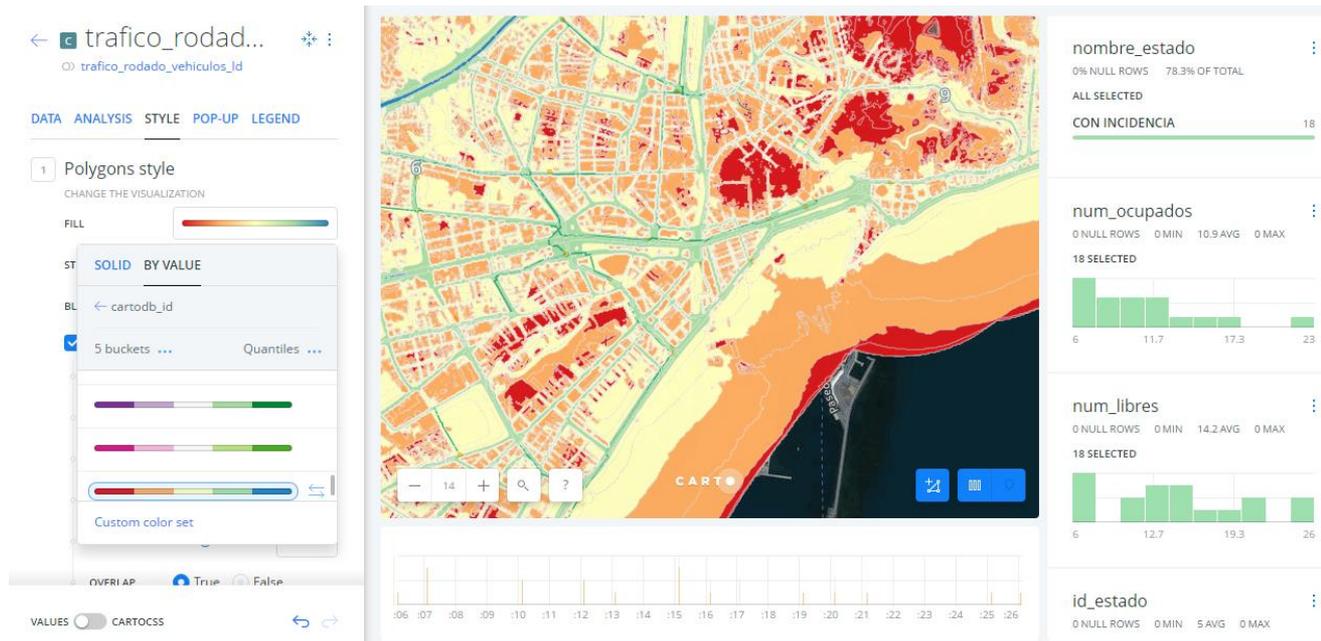
# Use Case 2: Open Data Visualization with Carto

- Adding new layer: Málaga Traffic Noise



# Use Case 2: Open Data Visualization with Carto

- Adding new layer: Málaga Traffic Noise



# Use Case 2: Open Data Visualization with Carto

- Adding new layer: Malaga Council Districts

**Sistema de Información Cartográfica ED50 - Distrito Municipal**

Seguidores: 0

Organización: **ORDENACIÓN DEL TERRITORIO Y VIVIENDA**

Conjunto de datos | Grupos | Flujo de Actividad

### Sistema de Información Cartográfica ED50 - Distrito Municipal

Datos y Recursos

**Distrito Municipal shp**  
Descripción: Contiene la delimitación administrativa y territorial del...

Explorar

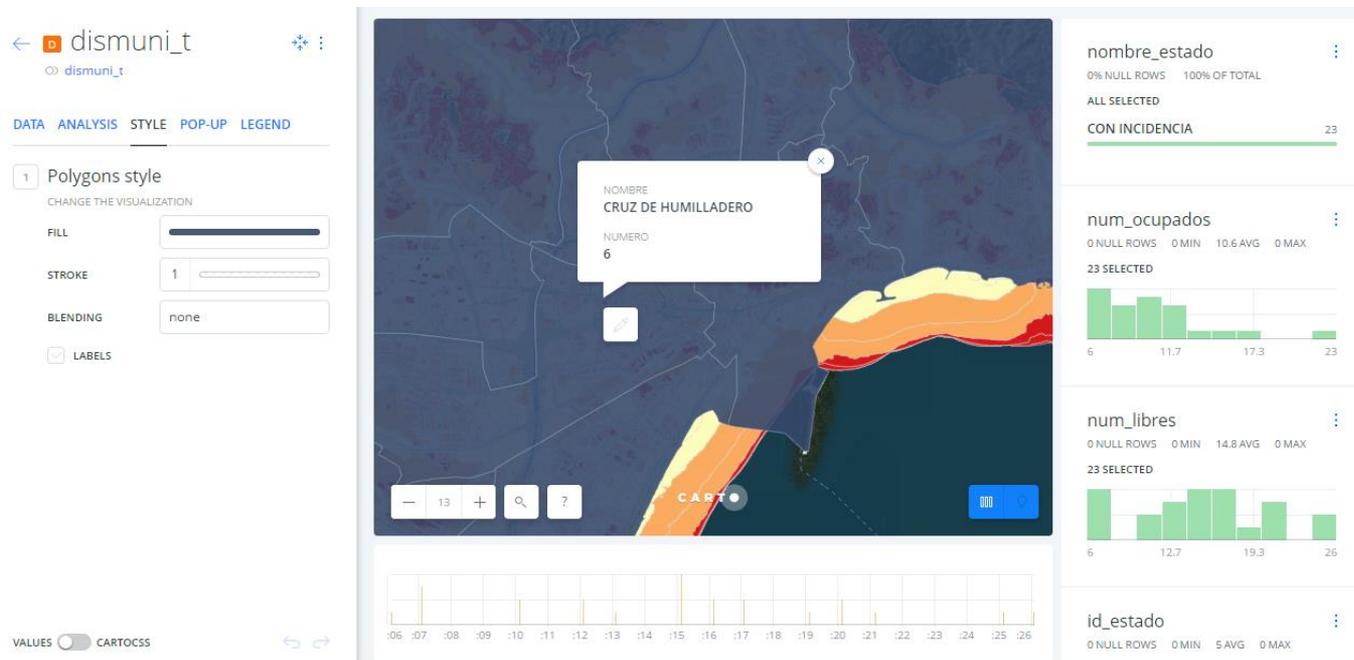
ED50 | callejero | cartografía | manzana | parcela | plano

Información Adicional

Campo	Valor
Autor	Responsable de la información de ordenación del territorio y vivienda
Descargas recientes	11
Descargas totales	360

# Use Case 2: Open Data Visualization with Carto

- Adding new layer: Malaga Council Districts



# Use Case 2: Open Data Visualization with Carto

- Adding new layer: Malaga Council Wifi Points

The screenshot displays the 'Portal Datos Abiertos Ayuntamiento de Málaga' website. The main navigation bar includes 'Conjuntos de datos', 'Organizaciones', 'Grupos', 'Aplicaciones', and 'Acerca de'. The breadcrumb trail shows the path: 'Organizaciones / ORDENACIÓN DEL TERRITORIO Y ... / Sedes Wifi Municipales'. The page title is 'Sedes Wifi Municipales'. On the left sidebar, there is a section for 'Seguidores' with a count of '0' and an 'Organización' section with a house icon. The main content area features a 'Conjunto de datos' tab selected, with sub-tabs for 'Grupos' and 'Flujo de Actividad'. Below this, the dataset name 'Sedes Wifi Municipales' is displayed with an RDF icon. A 'Datos y Recursos' section lists three data formats: CSV, JSON, and SHP, each with a description and an 'Explorar' button. At the bottom, there are filter tags for 'equipamiento', 'sede', and 'wifi', and a section for 'Información Adicional'.

# Use Case 2: Open Data Visualization with Carto

- Adding new layer: Malaga Council Wifi Points

