



Tema 1: Introducción a la Ciencia de Datos

Prof. Jesús Lares

Agenda

Big Data

Científico de Datos

Diagrama de Venn

Arquitectura de Big Data

Apache Hadoop

Escenario: Fuentes de datos

Casos de uso de Big Data

Tendencias en Big Data

Agenda

Big Data

Científico de Datos

Diagrama de Venn

Arquitectura de Big Data

Apache Hadoop

Escenario: Fuentes de datos

Casos de uso de Big Data

Tendencias en Big Data

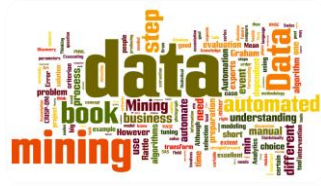
Big Data:

Es una plataforma tecnológica (hardware y software) que permite almacenar (de manera distribuida) y procesar (en forma paralela y distribuida) conjuntos de datos, que por su gran volumen, superan las capacidades de las plataformas TI tradicionales, ya sea porque tomaría demasiado tiempo procesar dichos datos o porque sería muy costoso implementar una arquitectura que soporte tal cantidad de datos.



Ciencia de Datos:

Es la **generación** de **conocimiento** a partir de **grandes volúmenes** de **datos**, aplicando técnicas de procesamiento paralelo y distribuido, para implementar **algoritmos** que permitan **predecir** o **detectar patrones** sobre los datos almacenados. A partir de los resultados obtenidos se podrán construir herramientas que permitan **analizar** los **resultados** y **apoyar** los procesos de **toma de decisiones**.



Unidades de Medidas

IEC prefixes and symbols for binary multiples

Factor	Name	Symbol	Origin	SI Derivation
2^{10}	kibi	Ki	kilobinary (2^{10}) ¹	kilo: $(10^3)^1$
2^{20}	mebi	Mi	megabinary (2^{10}) ²	mega: $(10^3)^2$
2^{30}	gibi	Gi	gigabinary (2^{10}) ³	giga: $(10^3)^3$
2^{40}	tebi	Ti	terabinary (2^{10}) ⁴	tera: $(10^3)^4$
2^{50}	pebi	Pi	petabinary (2^{10}) ⁵	peta: $(10^3)^5$
2^{60}	exbi	Ei	exabinary (2^{10}) ⁶	exa: $(10^3)^6$
2^{70}	zebi	Zi	zettabinary (2^{10}) ⁷	zetta: $(10^3)^7$
2^{80}	yobi	Yi	yottabinary (2^{10}) ⁸	yotta: $(10^3)^8$

Big Data: V's

■ 5 V's Tradicionales

■ Otras V's

Valor

Utilidad
Impacto Social
Impacto Económico

Visualización

Despliegue
Cuadro de Mandos

Variabilidad

Generación
Cambios

Visión

Data-Driven
Cultura del Dato

Volumen

Grandes volúmenes

Velocidad

Generación de datos
Procesamiento

Variedad

Diversos Formatos
Data Stream

Veracidad

Calidad de datos
Precisión

Agenda

Big Data

Científico de Datos

Diagrama de Venn

Arquitectura de Big Data

Apache Hadoop

Escenario: Fuentes de datos

Casos de uso de Big Data

Tendencias en Big Data

Profesional dedicado a analizar e interpretar grandes almacenes o bases de datos

The whiteboard is filled with various hand-drawn diagrams and charts, including:

- Social Media:** A line graph showing 'User Growth' and 'Engagement' over time. It includes labels for 'Email', 'SMS', 'Text', 'Social Network', 'Content Streaming', 'Video', 'Chat', and 'Picture'.
- Business:** A pie chart showing '35%' and 'Market'. It includes labels for 'Idea', 'Yes', 'No', 'Time', 'Business', 'Quality', 'Innovation', and 'Research'.
- Statistics:** A bar chart showing 'Growth' and 'Revenue' over time. It includes labels for 'Statistics', 'Last Year', 'This Year', 'Quality', 'Strategy', 'Vision', 'Success', and 'Innovation'.
- Other Diagrams:** A flowchart showing 'Idea' leading to 'Yes' and 'No', a diagram showing 'Idea' leading to 'Yes' and 'No', a diagram showing 'Idea' leading to 'Yes' and 'No', and a diagram showing 'Idea' leading to 'Yes' and 'No'.



Conocimientos en Big Data



MODERN DATA SCIENTIST

Data Scientist, the sexiest job of 21th century requires a mixture of multidisciplinary skills ranging from an intersection of mathematics, statistics, computer science, communication and business. Finding a data scientist is hard. Finding people who understand who a data scientist is, is equally hard. So here is a little cheat sheet on who the modern data scientist really is.

MATH & STATISTICS

- ☆ Machine learning
- ☆ Statistical modeling
- ☆ Experiment design
- ☆ Bayesian inference
- ☆ Supervised learning: decision trees, random forests, logistic regression
- ☆ Unsupervised learning: clustering, dimensionality reduction
- ☆ Optimization: gradient descent and variants

DOMAIN KNOWLEDGE & SOFT SKILLS

- ☆ Passionate about the business
- ☆ Curious about data
- ☆ Influence without authority
- ☆ Hacker mindset
- ☆ Problem solver
- ☆ Strategic, proactive, creative, innovative and collaborative



PROGRAMMING & DATABASE

- ☆ Computer science fundamentals
- ☆ Scripting language e.g. Python
- ☆ Statistical computing package e.g. R
- ☆ Databases SQL and NoSQL
- ☆ Relational algebra
- ☆ Parallel databases and parallel query processing
- ☆ MapReduce concepts
- ☆ Hadoop and Hive/Pig
- ☆ Custom reducers
- ☆ Experience with xaaS like AWS

COMMUNICATION & VISUALIZATION

- ☆ Able to engage with senior management
- ☆ Story telling skills
- ☆ Translate data-driven insights into decisions and actions
- ☆ Visual art design
- ☆ R packages like ggplot or lattice
- ☆ Knowledge of any of visualization tools e.g. Flare, D3.js, Tableau

Científico de Datos: Trabajo en Equipo



Agenda

Big Data

Científico de Datos

Diagrama de Venn

Arquitectura de Big Data

Apache Hadoop

Escenario: Fuentes de datos

Casos de uso de Big Data

Tendencias en Big Data

Diagrama de Venn

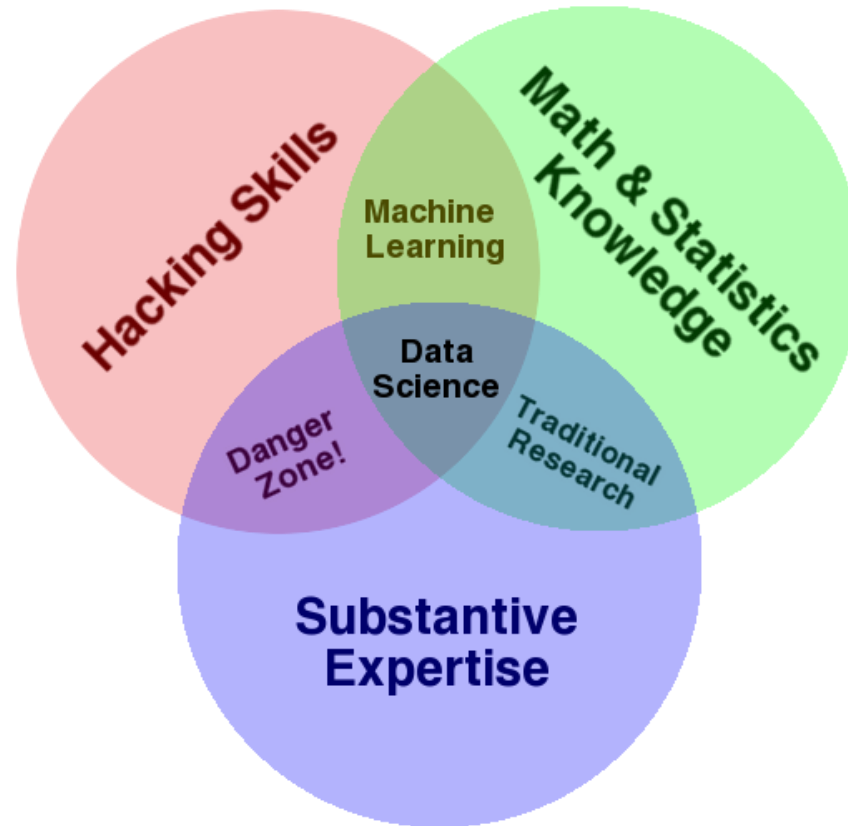
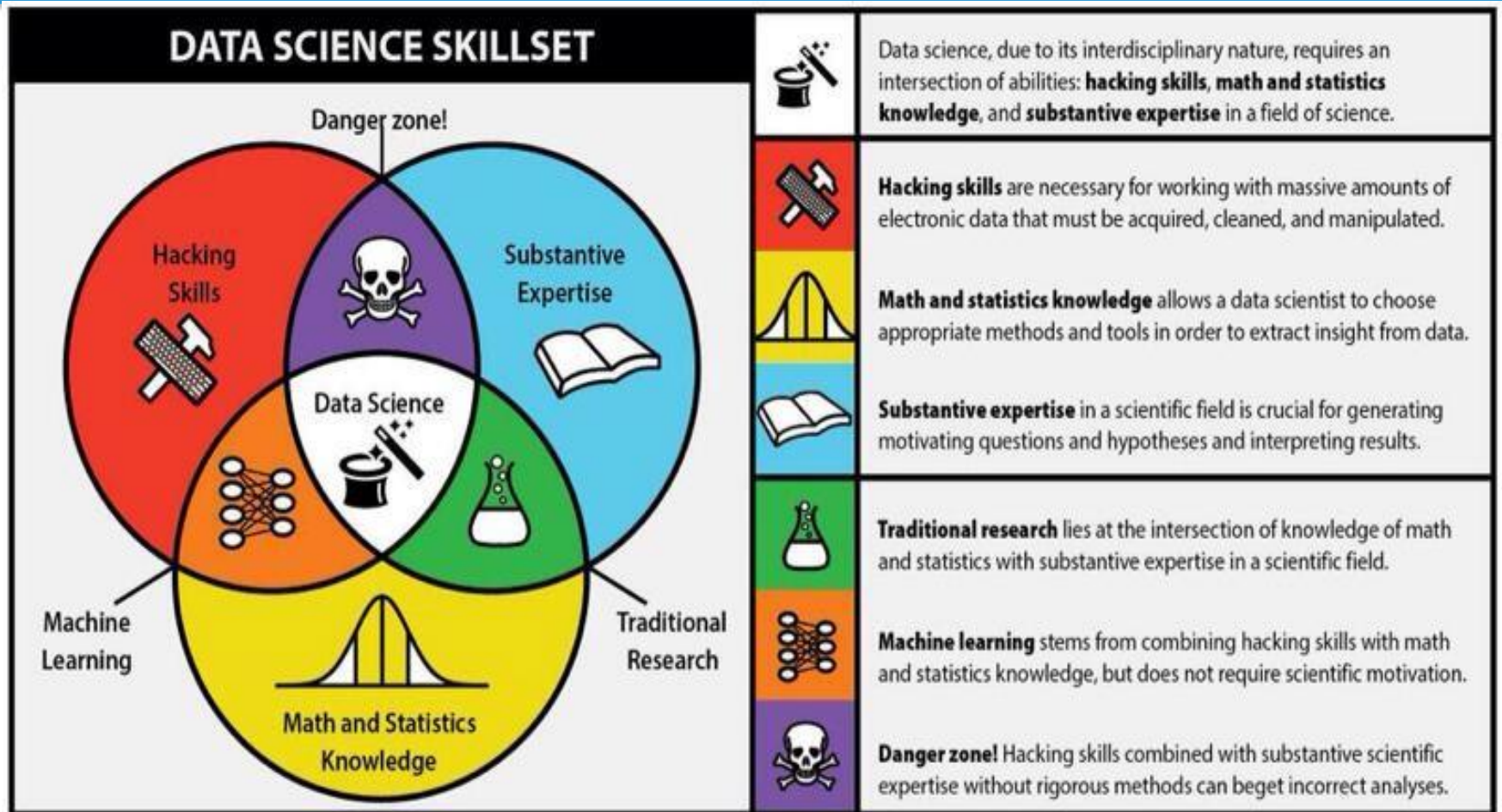


Diagrama de Venn



Agenda

Big Data

Científico de Datos

Diagrama de Venn

Arquitectura de Big Data

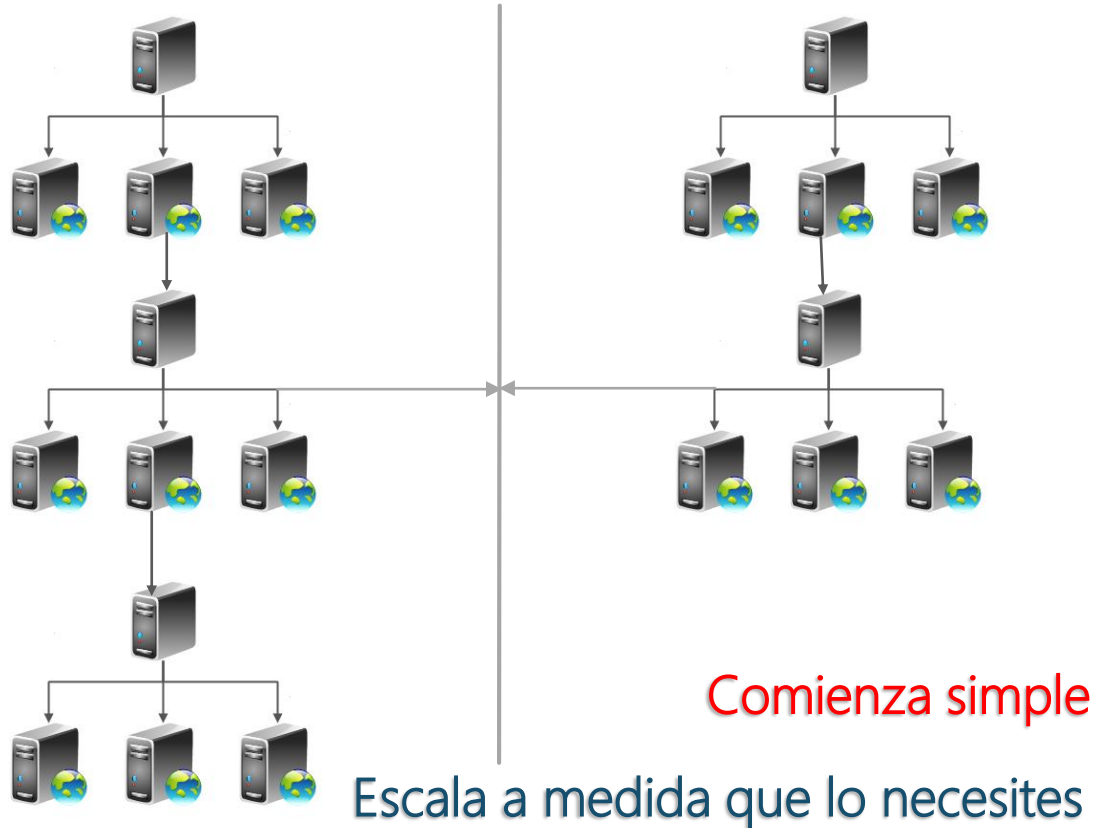
Apache Hadoop

Escenario: Fuentes de datos

Casos de uso de Big Data

Tendencias en Big Data

Hardware: Cluster Maestro-Esclavo



Nodos

- RAM: ≥ 4 GB
- CPU: ≥ 2 Cores
- HDD: ≥ 512 GB

Escalabilidad Horizontal

- Desempeño proporcional al tamaño de la Plataforma
- Compatibilidad con la Nube

Software: Open Source



Linux



Agenda

Big Data

Científico de Datos

Diagrama de Venn

Arquitectura de Big Data

Apache Hadoop

Escenario: Fuentes de datos

Casos de uso de Big Data

Tendencias en Big Data



Antecedentes Apache Hadoop

Publicaciones de Google (Papers):

- Google File System (2003)
- MapReduce (2004)
- Bigtable (2006)

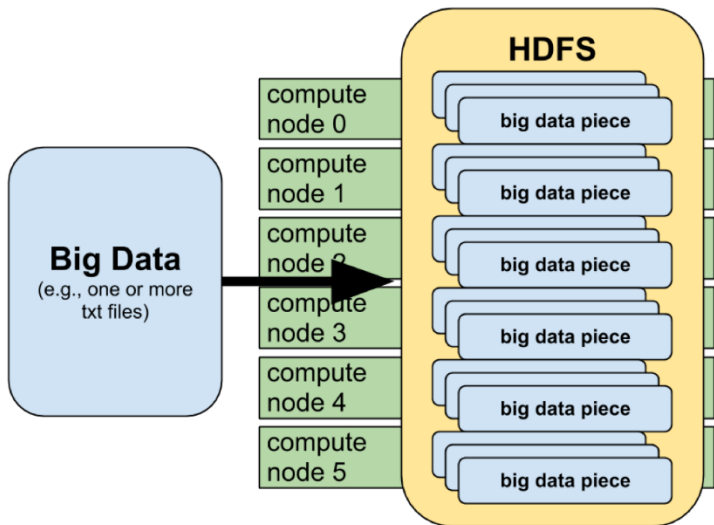


Apache Hadoop

Plataforma de software de código abierto para el almacenamiento distribuido y procesamiento distribuido de grandes volúmenes de datos (Big Data) en clusters de ordenadores contruidos a partir de "Commodity hardware".

Los servicios de Hadoop proporcionan almacenamiento de datos, procesamiento de datos, acceso a datos, la gestión de datos, seguridad y operaciones.

Apache Hadoop: HDFS



HDFS: Hadoop Distributed File System

Almacenamiento Distribuido

Soporta grandes volúmenes de datos

Replicación de datos

Soporta escalabilidad horizontal

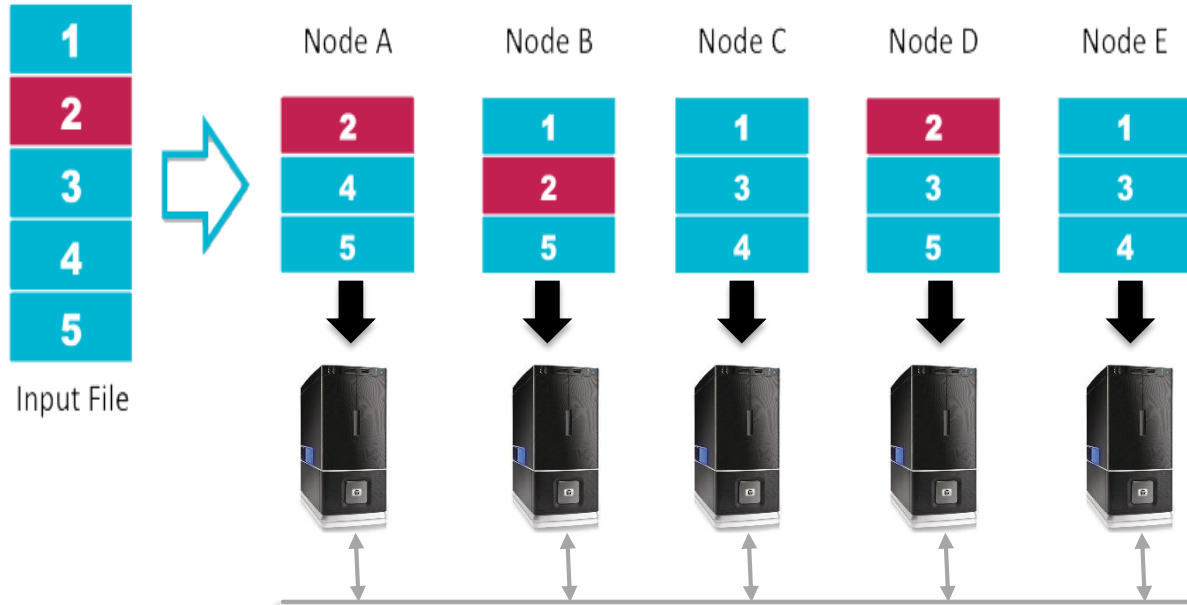
Integración con diversas plataformas

Integración con almacenes de datos

Código abierto (Open Source)



HDFS Data Distribution

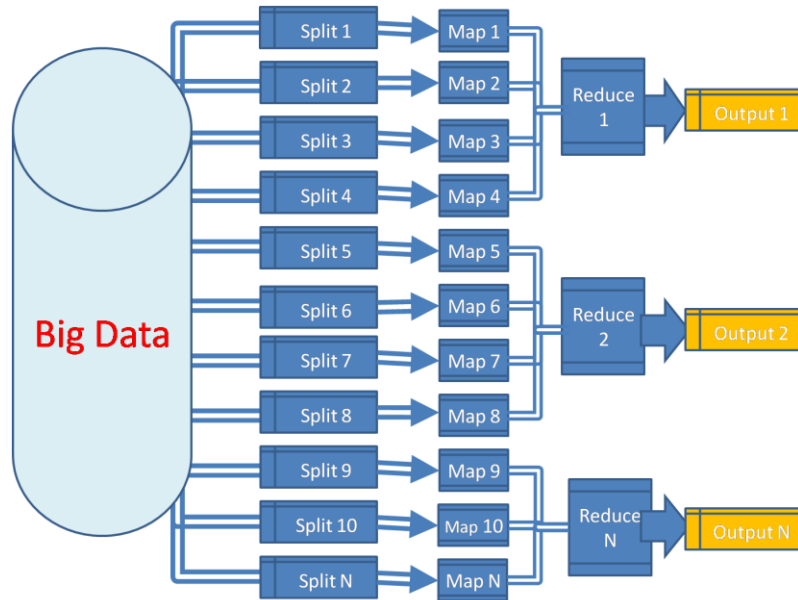


Bloques (64 MB o 128 MB)

Replicación (3 por defecto)

Tolerancia a fallos

Apache Hadoop: MapReduce



Procesamiento paralelo y distribuido

Eficiente sobre grandes volúmenes de datos

Multilenguaje: Java, R, Python

Esquema Divide y Vencerás

Función Mapeo (Map)

Función de Reducción (Reduce)

Cambio de Paradigma de programación



Ejemplo: Conteo de palabras

Archivo de
Entrada



Archivo
dividido



Proceso Map



Proceso
Ordenamiento



Proceso
Reduce



Resultado
Final

Pie Orange Banana

Pizza Orange Pizza

Pie Pie Pear

Pie, 1
Orange, 1
Banana, 1

Pizza, 1
Orange, 1
Pizza, 1

Pie, 1
Pie, 1
Pear, 1

Banana, 1

Orange, 1
Orange, 1

Pear, 1

Pie, 1
Pie, 1
Pie, 1

Pizza, 1

Banana, 1

Orange, 2

Pear, 1

Pie, 3

Pizza, 1

Banana, 1
Orange, 2
Pear, 1
Pie, 3
Pizza, 1

Pie Orange Banana
Pizza Orange Pizza
Pie Pie Pear



Apache Hadoop: YARN



HADOOP 1.0

MapReduce
(cluster resource management
& data processing)

HDFS
(redundant, reliable storage)

HADOOP 2.0

MapReduce
(data processing)

Others
(data processing)

YARN

(cluster resource management)

HDFS

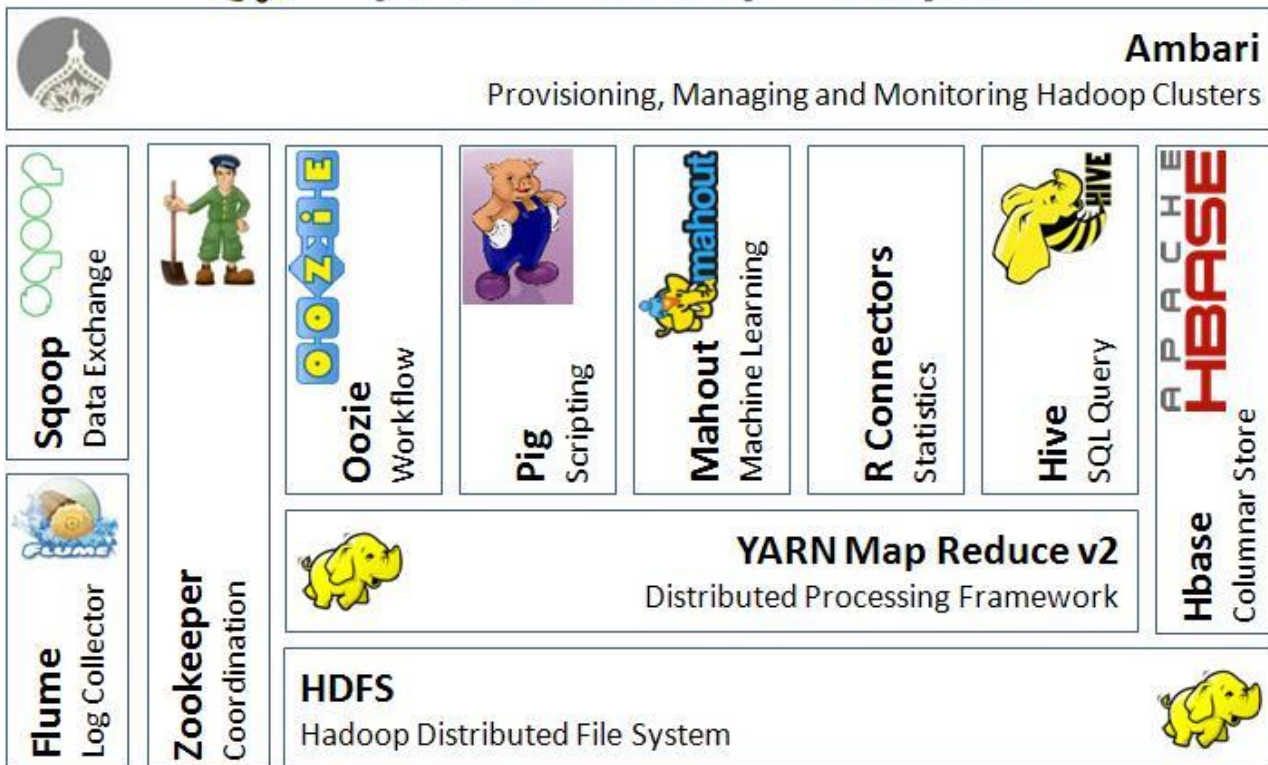
(redundant, reliable storage)

Características	Descripción
Tenencia Múltiple (Multi-tenancy)	Permite que múltiples motores de acceso utilicen Hadoop como el estándar común para los procesamiento en lotes, interactivo y en tiempo real, permitiendo el acceso simultáneamente el mismo conjunto de datos.
Utilización de Cluster	Asigna dinámicamente los recursos del cluster
Escalabilidad	Soporta la escalabilidad o crecimiento del Cluster
Compatibilidad	Mantiene la compatibilidad hacia atrás de los procedimientos MapReduce de las versiones anteriores de Hadoop

YARN: Yet Another Resource Negotiator



Apache Hadoop Ecosystem



Agenda

Big Data

Científico de Datos

Diagrama de Venn

Arquitectura de Big Data

Apache Hadoop

Escenario: Fuentes de datos

Casos de uso de Big Data

Tendencias en Big Data

Integración de fuentes de datos empresariales



Diversas tecnologías: Sistemas, BD

Diversos formatos de datos

Grandes volúmenes de datos

Problemas de integración de datos

Problemas de Calidad de Datos

Complejidad en la generación de reportes

Consultas ineficientes

Esquemas de almacenamiento ineficientes

Escenario: Fuentes de datos

Grandes volúmenes de datos

Esquemas de almacenamiento ineficientes



Sistema de archivos distribuidos

Diversas tecnologías: Sistemas, BD

Problemas de integración de datos

Diversos formatos de datos



Base de Datos NoSQL

Escenario: Fuentes de datos

Consultas ineficientes

Complejidad en la generación de reportes



Procesamiento paralelo y distribuido

Problemas de Calidad de Datos



Analítica Predictiva



Procesamiento paralelo y distribuido

Base de Datos NoSQL

Analítica Predictiva

Sistema de archivos distribuidos

Agenda

Big Data

Científico de Datos

Diagrama de Venn

Arquitectura de Big Data

Apache Hadoop

Escenario: Fuentes de datos

Casos de uso de Big Data

Tendencias en Big Data

Casos de uso de Big Data

Redes Sociales



Internet de las cosas



Investigaciones científicas



Patrones de Compra



Detección de Fraudes



Ataques informáticos



Agenda

Big Data

Científico de Datos

Diagrama de Venn

Arquitectura de Big Data

Apache Hadoop

Escenario: Fuentes de datos

Casos de uso de Big Data

Tendencias en Big Data

Tendencias en Big Data



Mejora del análisis predictivo



Cambio en los modelos de negocio



Sistemas inteligentes basados en el "machine learning"



Integración de los datos en la empresa

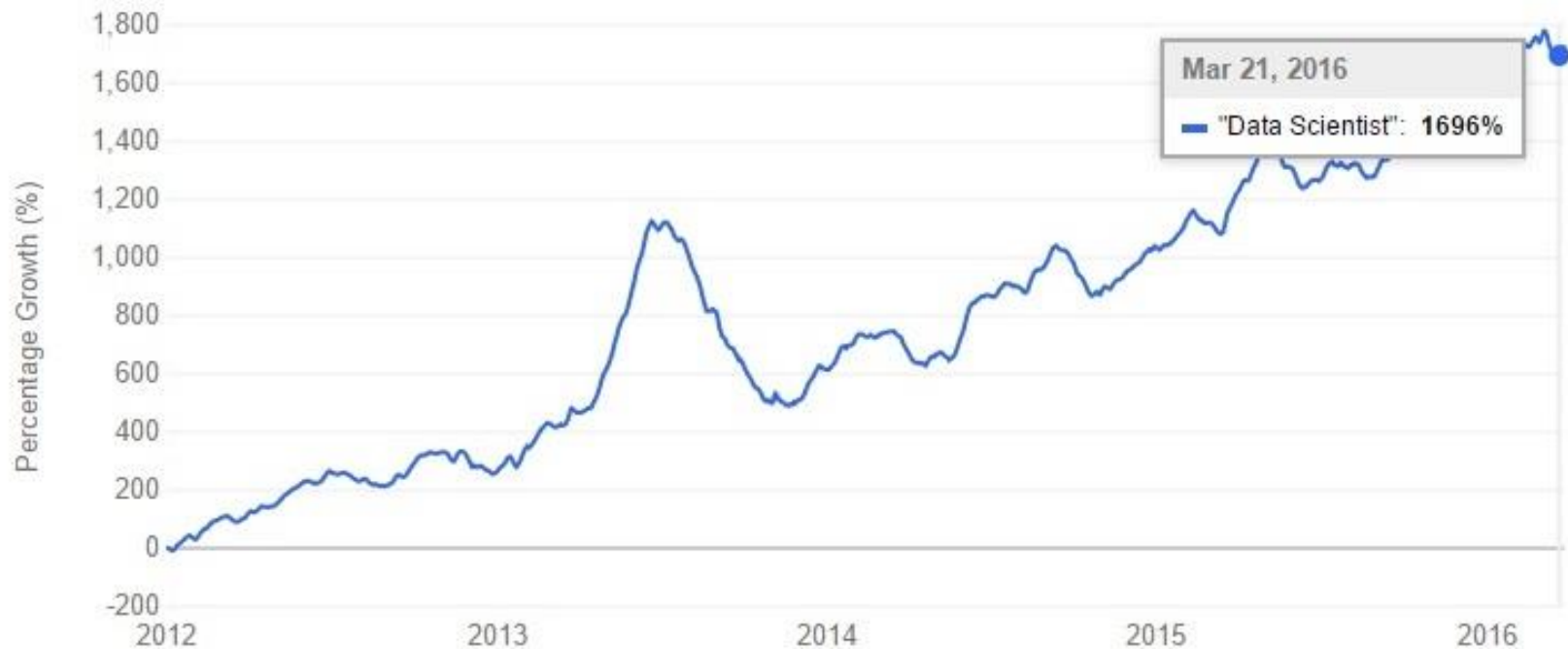


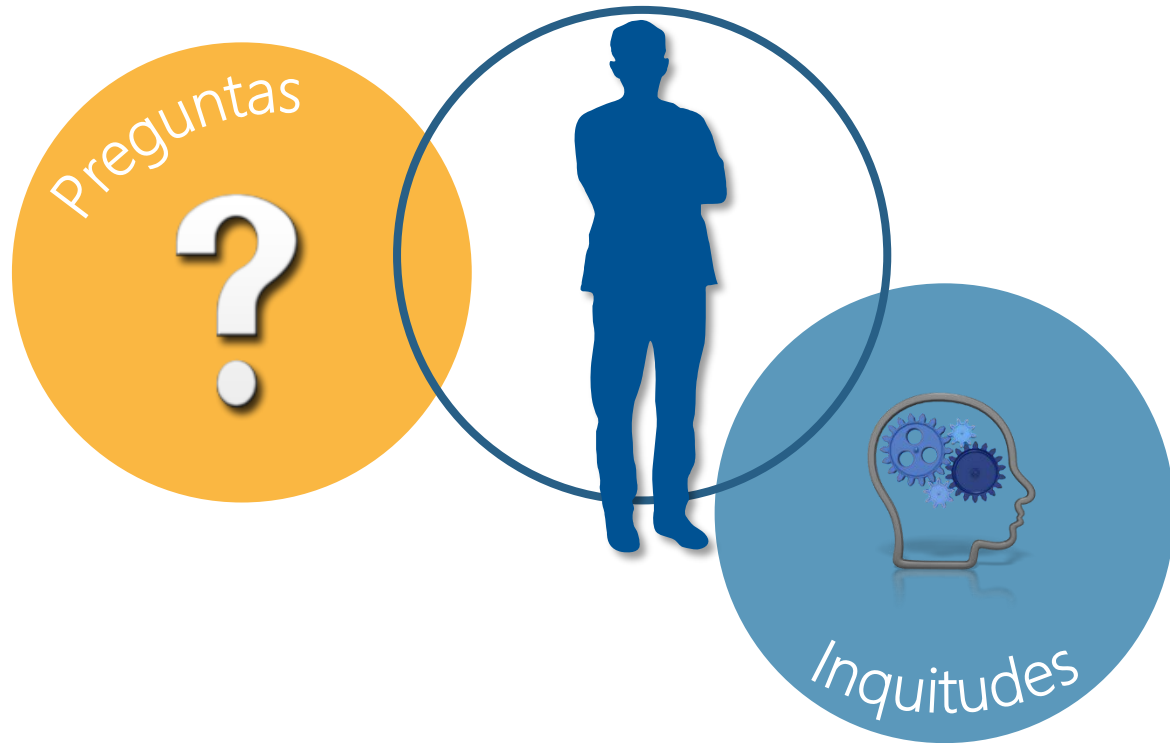
La visualización de datos ofrecerá una panorámica completa



El Big Data conquistará nuevas industrias

"Data Scientist" Job Trends





[illegible]