

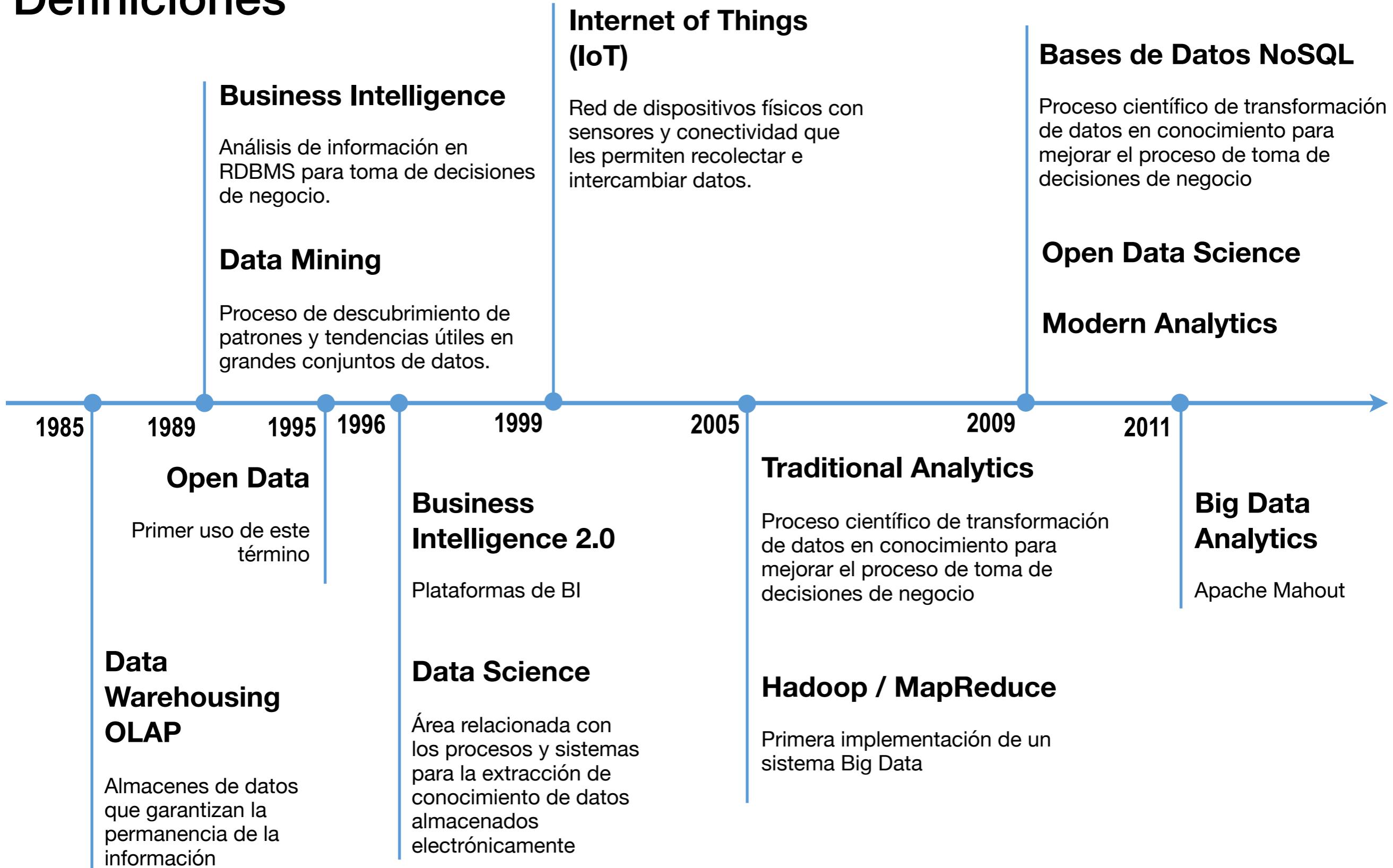
Hacia una visión unificada de Data Science, Analytics y Big Data

(con ejemplos en mercados de energía)

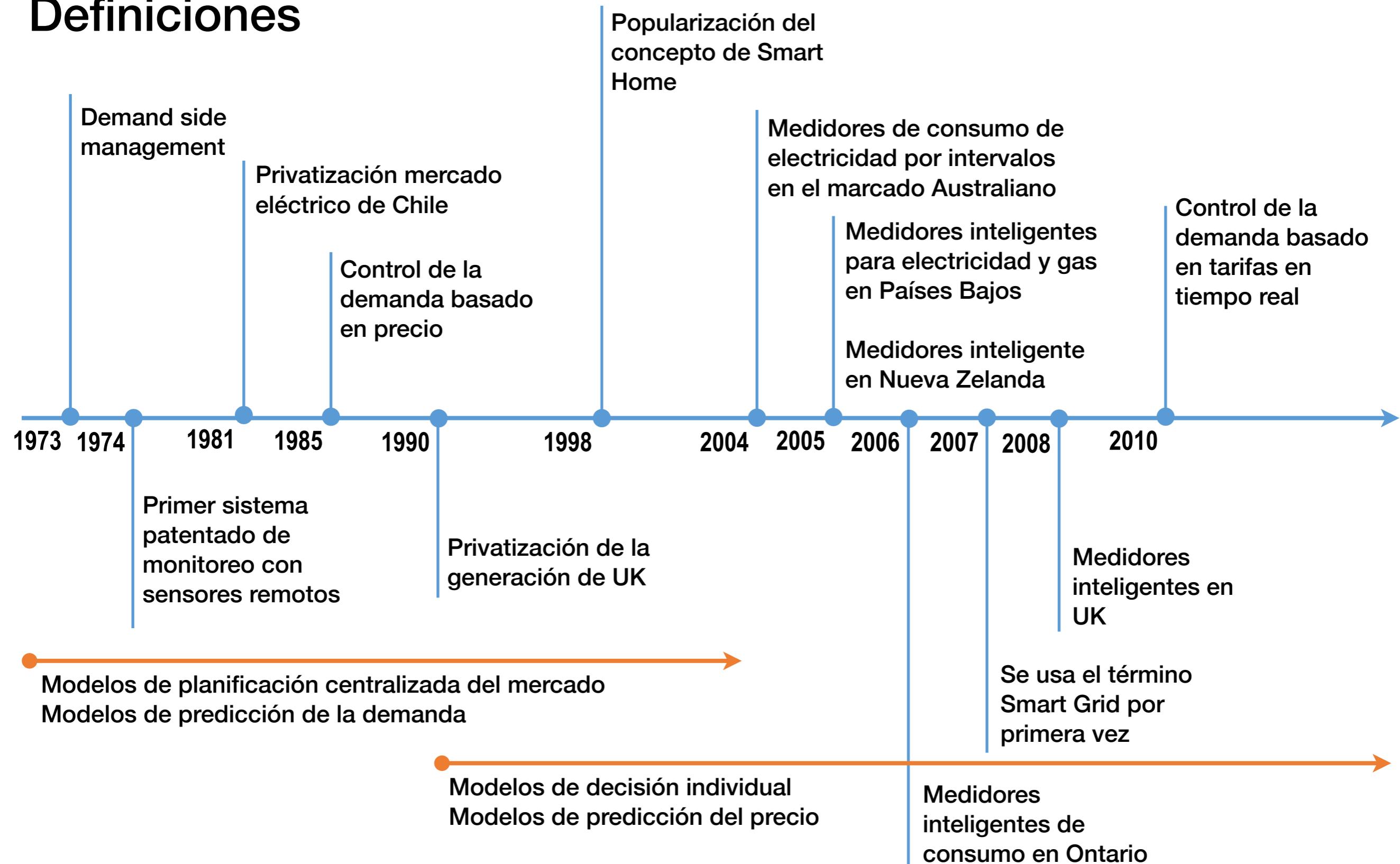
Esta presentación describe, bajo un marco común, los conceptos fundamentales de Data Science, Analytics y Big Data y establece su similitudes y diferencias.

Descargue la última versión de este documento de:
<https://github.com/jdvelasq/data-science-docs/blob/master/ds-analytics-bigdata.pdf>

Definiciones



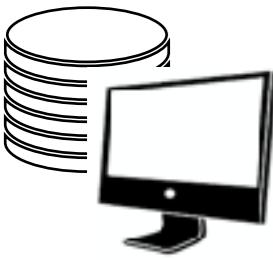
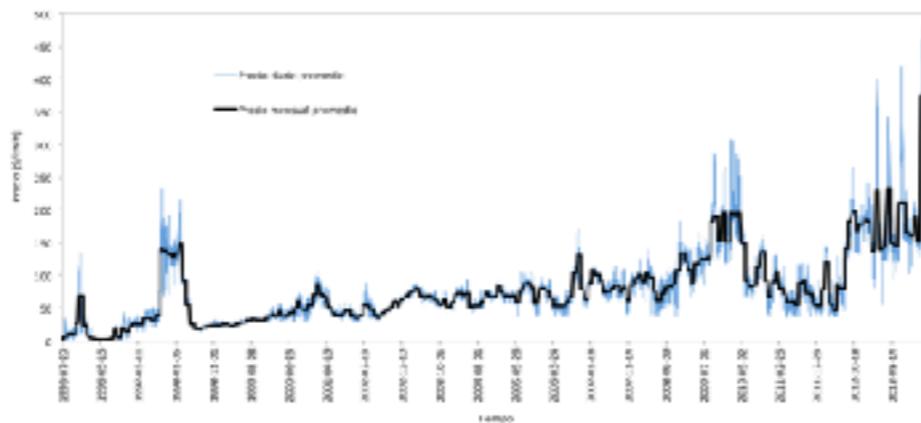
Definiciones



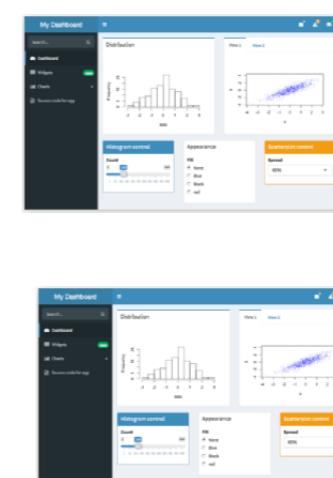
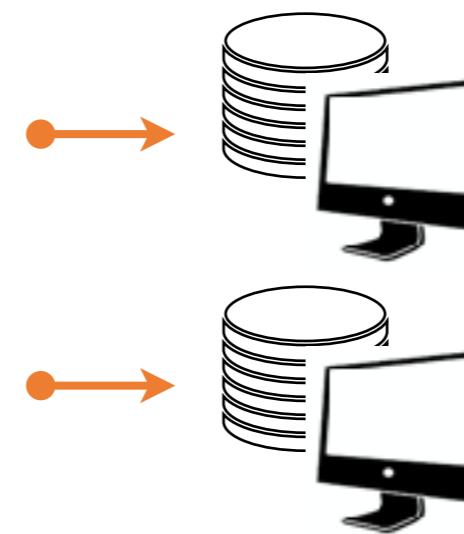
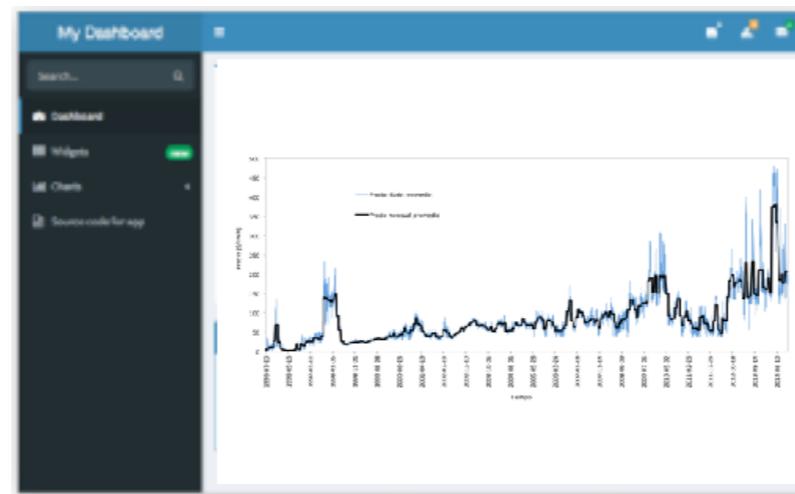
Definiciones



Técnicas de
modelado



Técnicas de modelado
+ Reportes



Producto de Datos

Aplicación que combina datos con algoritmos para inferencia, predicción u optimización para generar más datos e información valiosa.

- Aprendizaje a partir de los datos.
- Auto-adaptación
- Ampliamente aplicable.

Inteligencia de Negocios

Analytics

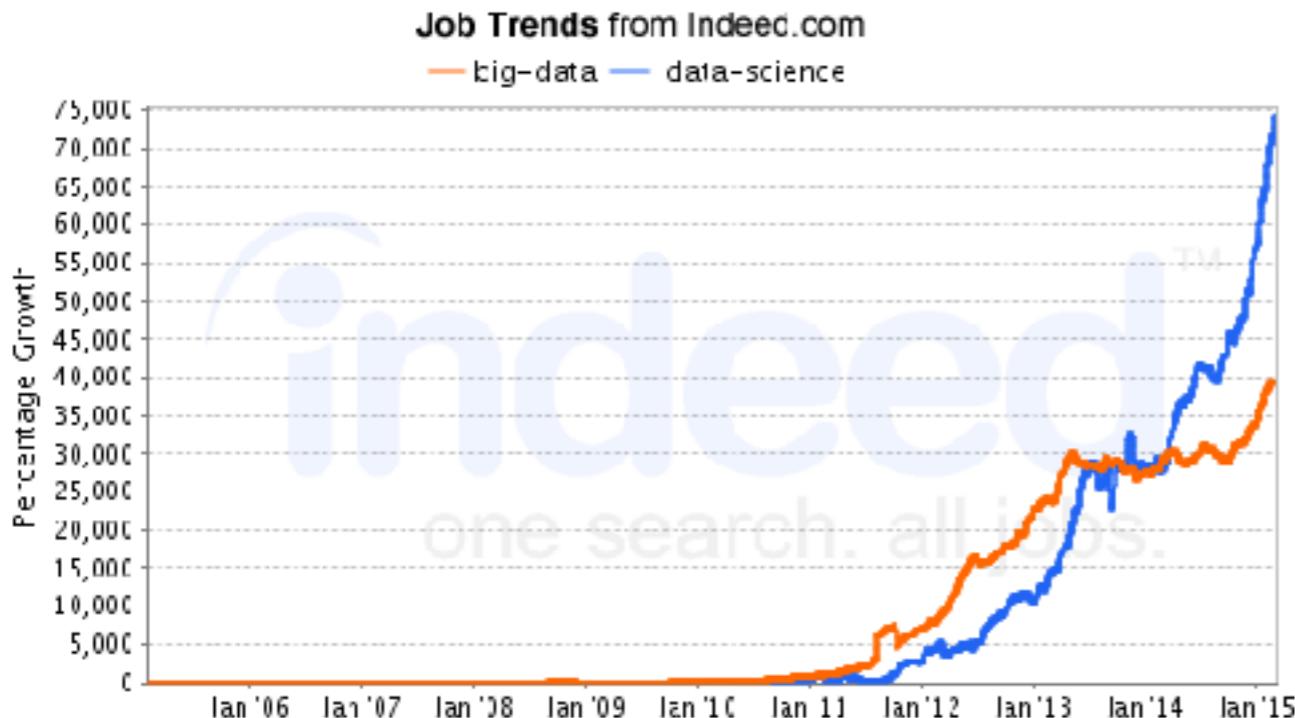
Perspectiva

A recent study by the McKinsey Global Institute concludes, "a shortage of the analytical and managerial talent necessary to make the most of Big Data is a significant and pressing challenge (for the U.S.)." The report estimates that there will be four to five million jobs in the U.S. requiring data analysis skills by 2018, and that large numbers of positions will only be filled through training or retraining. The authors also project a need for 1.5 million more managers and analysts with deep analytical and technical skills "who can ask the right questions and consume the results of analysis of big data effectively."

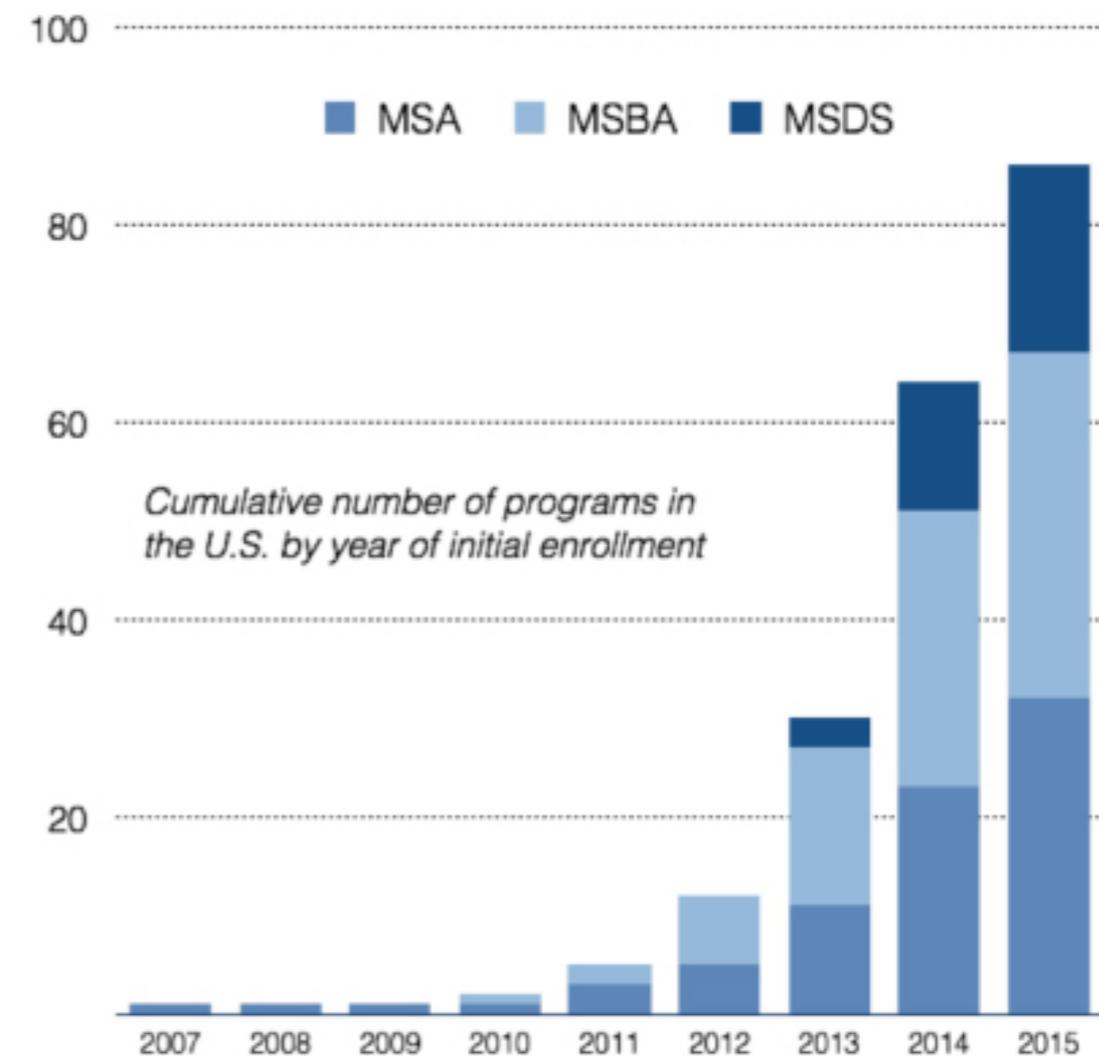
#16	3,433	\$105,395	#1
Highest Paying Job in Demand	Number of Job Openings	Average Base Salary	Best Job in America for 2016

Sources: [25 Best Jobs in America](#) and [25 Highest Paying Jobs in America for 2016](#)

Evolución de empleos/educación en Big Data & Data Science



GROWTH OF MASTER'S DEGREE PROGRAMS IN ANALYTICS AND DATA SCIENCE



Data Science Analytics

Big Data

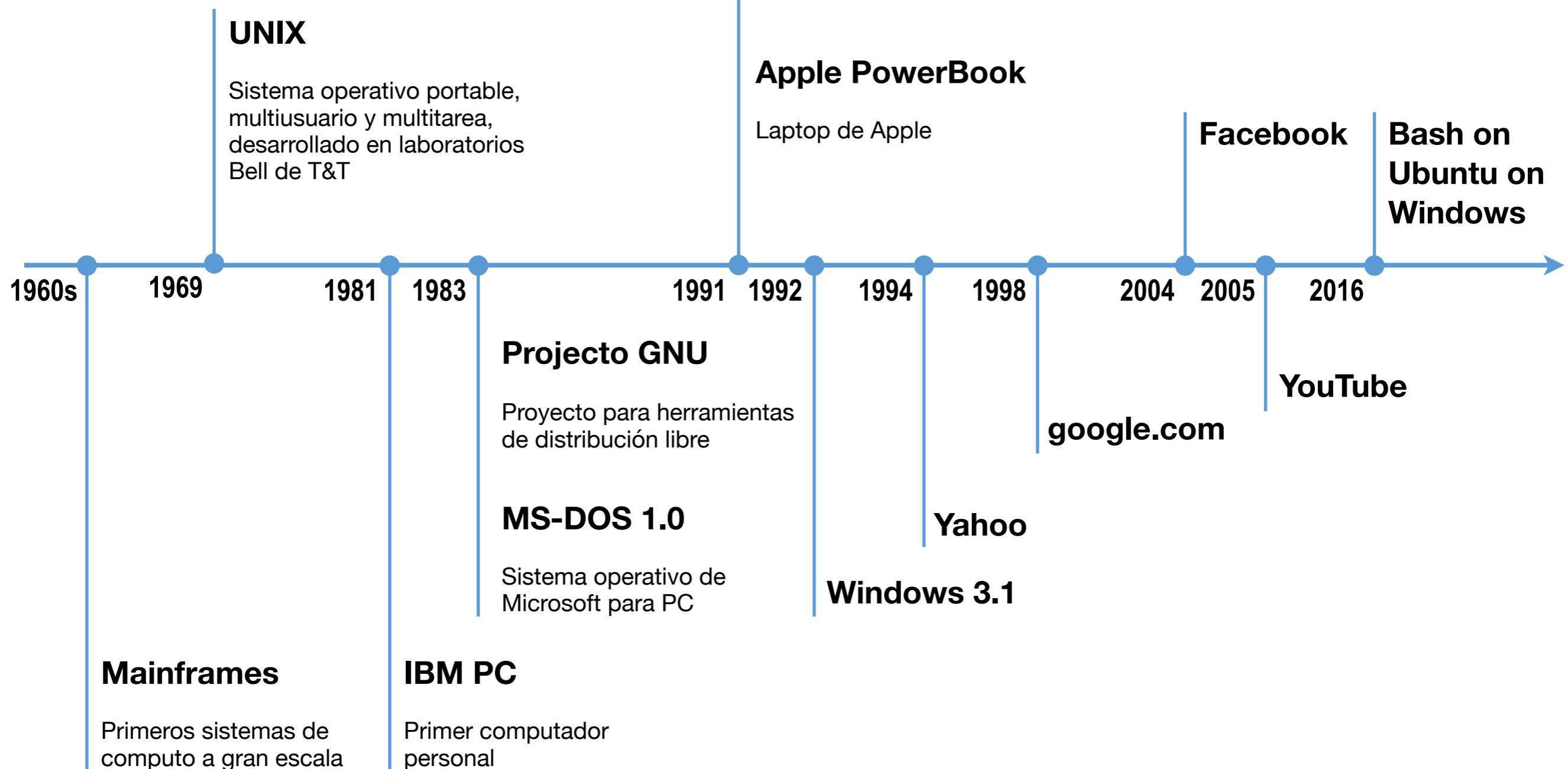
-
- A vertical timeline chart showing the evolution of Data Science and Big Data from 1941 to 2016. The timeline is represented by a blue vertical line with circular markers at each event. The events are listed on the left, and their descriptions are on the right.
- 1941 - Explosión de la información
 - 1966 - Primeros sistemas de cómputo centralizado
 - 1970 - Bases de datos relacionales
 - Concise survey of computer methods - 1974
 - International Association for Statistical Computing - 1977
 - 1st Knowledge discovery in Databases (KDD) workshop - 1989
 - Se usa por primera vez el término Data Science - 1996
 - Journal Data Mining and Knowledge Discovery - 1997
 - Data Science como una disciplina independiente - 2001
 - Data Science Journal - 2002
 - Journal of Data Science - 2003
 - Se usa por primera vez el término Analytics - 2005
 - Se reconoce el término Data Scientist - 2008
 - “... sexy job in the next ten years will be statisticians” - 2009
 - Data Scientist: The Sexiest Job of the 21st Century - 2012
 - Deep Learning y explosión de Machine Learning - 2015
 - 1976 - Lenguaje SQL (Oracle)
 - 1980 - “Data expands to fill the space available”
 - 1985 - Data warehousing
 - 1989 - Inteligencia de negocios (concepto)
 - 1992 - Primer sistema para reporte de bases de datos (Crystal Report)
 - 1995 - Explosión de la Web
 - 1996 - Primeras plataformas de inteligencia de negocios
 - 1997 - 1r uso del término Big Data en un artículo
 - 1998 - Crisis de la BI por el volumen de datos
 - 1999 - Internet of Things (IoT), Predictive Analytics
 - 2000 - Se reconoce la explosión de datos potencialmente relevantes
 - 2001 - Software as a Service (SaaS), Data Volume, Velocity, and Variety
 - 2002 - Web Services
 - 2004 - Google’s MapReduce
 - 2005 - Yahoo desarrolla Hadoop basado en MapReduce de Google
 - 2008 - Google procesa 20 Petabytes en un día.
 - 2009 - Aparece el término NoSQL.
 - 2010 - Cloud
 - 2011 - Aparece el término Data Lakes
 - 2013 - Cloud
 - 2014 - Explosión de IoT
 - 2015 - Smart cities
 - 2016 - Las organizaciones implementan Big Data de forma masiva

Elementos constitutivos

- Infraestructura computacional
- Datos
- Programas (Algoritmos, técnicas y metodologías)

Infraestructura computacional

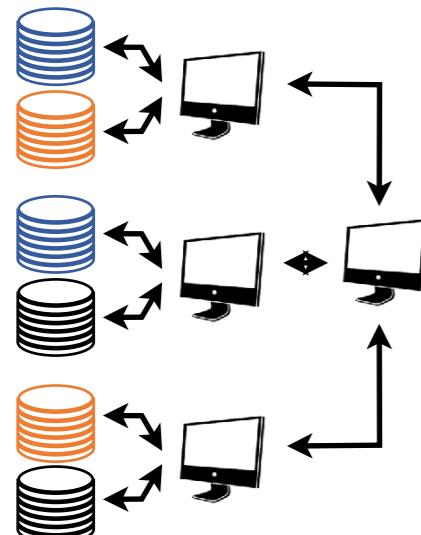
Infraestructura computacional



Infraestructura computacional

Computación local

Servidores + red + clientes



Cloud computing / utility computing

Servidores y almacenamiento en la nube + internet + clientes locales

Software as a Service (SaaS)

Software almacenado en máquinas suministradas por un tercero.

Aplicaciones accesadas vía un cliente o la Web.

Orientado a aplicaciones de usuario final.

Platform as a Service (PaaS)

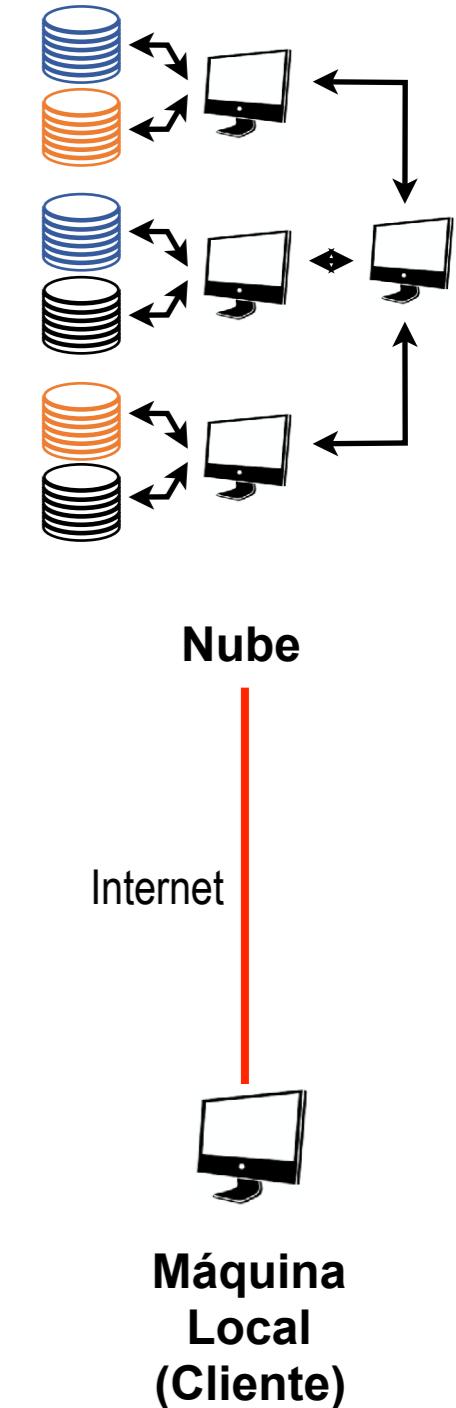
Orientado a desarrolladores.

Ambiente de desarrollo gestionado por un tercero.

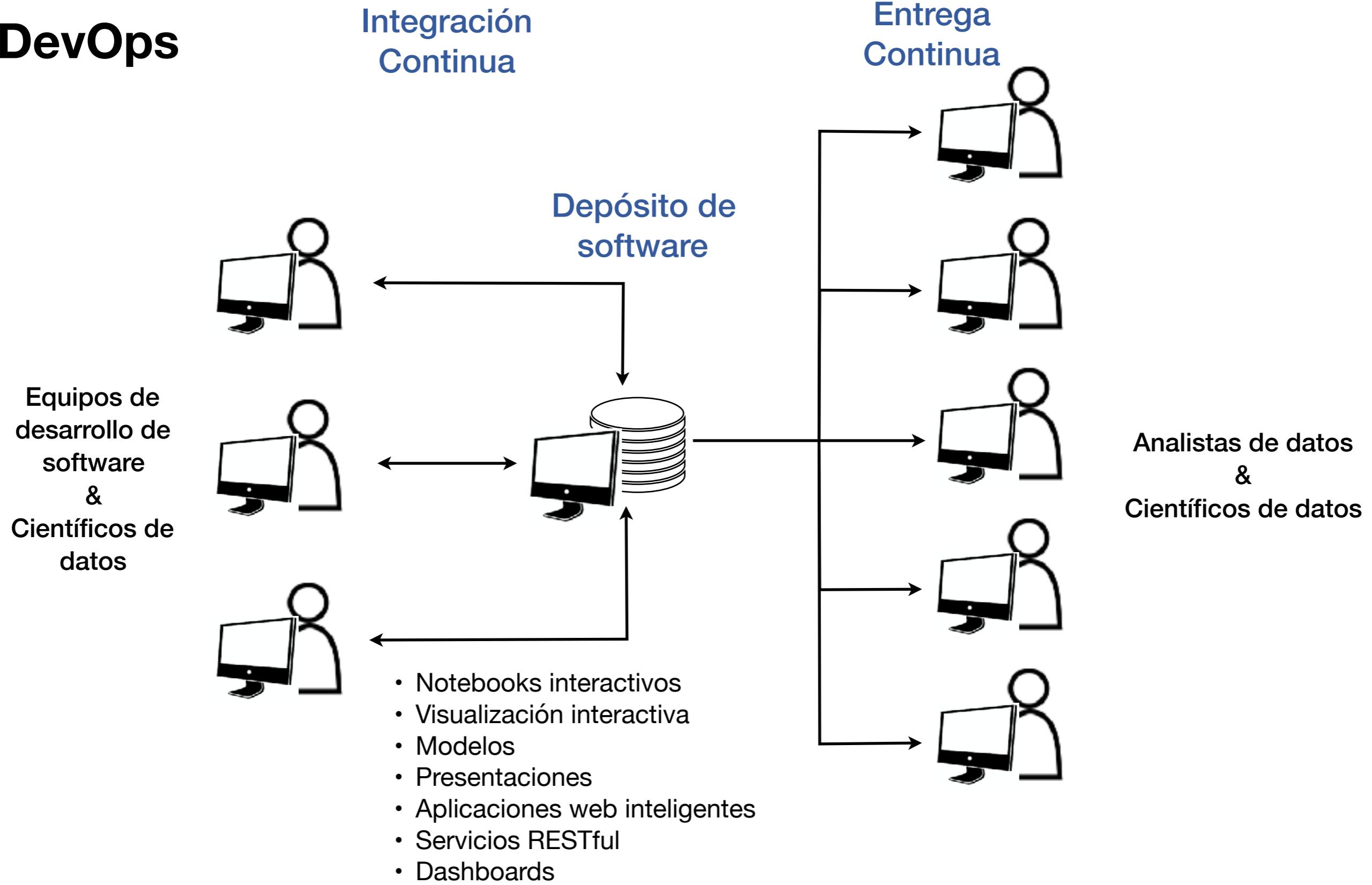
Infrastructure as a Service (IaaS)

Bloques básicos para construcción de ambientes manejados por un tercero

Capacidad de procesamiento, almacenamiento, conectividad, seguridad, etc.



DevOps



Datos

Fuentes de datos

Archivos de datos y Web

- Archivos de texto delimitados
- JSON
- XML
- Archivos de Log
- Archivos específicos de aplicación

Internet of Things

Red de dispositivos físicos con sensores y conectividad que les permiten recolectar e intercambiar datos.

Data warehouse y SQL

- RDBMS
- Cubos de datos

Hogares inteligentes
Ciudades inteligentes
Vehículos eléctricos
Fuentes renovables de energía
Líneas de potencia
Perfil de la demanda
Respuesta de la demanda
Detección de fallos

Hadoop & Spark

Stream de datos

NoSQL

- Almacenes de documentos
- Bases de datos columnar
- Diccionarios (clave, valor)

(Dispositivos usables, ...)

(Data warehousing para gestión del mercado eléctrico)
(Sistemas de bases de datos en organizaciones)

Ejemplo de un reporte del IDEAM (Fuentes de datos)

IDEAM - INSTITUTO DE HIDROLOGIA, METEOROLOGIA Y ESTUDIOS AMBIENTALES																	SISTEMA DE INFORMACION HIDROMETEOROLOGICA -METEORO 3- MAR 1987											
EVALUACION HORARIA DE LA DIRECCION (D) Y VELOCIDAD (V) EN M/S DEL VIENTO EN SUPERFICIE																												
PROCESO	Jul 4-2017 *																ESTACION : 23125060 ALBANIA											
LATITUD	0545 N								DEPARTAMENTO SANTANDER								TIPO EST	CO	SUBZONA HIDR. QDA CANUTILLO									
LONGITUD	7354 W								MUNICIPIO ALBANIA								ENTIDAD	1 IDEAM	ZONA HIDROGR. RIO MIRA									
ELEVACION	1690 m.s.n.m																REGIONAL	8 SANTANDERES-ARA	AREA HIDROGR. M/LENA-CAUCA									
CUADRO NO 1 (1 PARTE)																												
- HORA	00-01	01-02	02-03	03-04	04-05	05-06	06-07	07-08	08-09	09-10	10-11	11-12	12-13	13-14	-	-	-	-	-	-	-							
-	D	V	D	V	D	V	D	V	D	V	D	V	D	V	D	V	D	V	D	V	D							
1	NE	1.6	NE	1.0	NE	2.3	NE	1.3	NE	1.3	NE	1.6	NE	1.6	SW	1.2	SW	1.6	SW	1.3	SW	2.0	NW	1.2				
2	NE	1.3	NE	1.3	SW	1.6	NE	1.6	NE	1.5	NE	2.2	SE	1.2	SE	.6	NW	1.0	SW	1.5	SW	1.8	SW	2.0	SW	2.3		
3	NE	.9	NE	.6	NE	.9	NE	1.2	NE	1.3	NE	1.0	NE	1.2	SW	.9	SW	1.0	SW	1.8	SW	2.2	SW	2.0	SW	2.6		
4	SE	1.0	SE	.6	NE	.9	NE	.9	SE	1.0	NE	.9	NE	.9	NE	.6	SW	.9	SW	1.0	SW	1.5	SW	1.3	SW	2.0		
5	NW	.6	NW	.6	NE	1.3	NE	1.0	SE	.6	SW	.6	NE	.6	SW	1.0	SW	1.0	SW	1.2	SW	2.2	SW	2.2	SW	2.6		
6	NE	1.0	NE	1.3	NE	1.0	NE	1.0	NE	.9	NE	.9	NE	.6	SW	.6	SW	1.2	SW	1.2	SW	2.0	SW	2.0	SW	1.6		
7	NE	.6	SE	.5	VR	.6	NE	.5	NE	.6	SW	.5	SE	.6	SW	.9	SW	1.2	SW	1.5	SW	1.2	SW	2.2	SW	1.2	SW	1.2
8	NE	1.3	NE	1.0	NE	1.2	NE	1.5	NE	1.2	SE	.6	NE	.9	NW	.6	SW	1.3	SW	1.8	NW	.9	SW	1.6	SW	2.4	SW	2.4
9	NE	1.5	NE	1.5	NE	1.6	NE	1.3	NE	.6	NE	.6	SW	.6	SE	.5	SE	.6	SW	.9	SW	1.6	SW	1.5	SW	1.3	SW	1.3
10	SE	.6	SW	.6	SE	.6	NE	.6	NE	.9	NE	1.2	NE	.6	SW	.6	SW	1.3	SW	1.3	SW	1.5	SW	2.0	SW	1.0	SW	2.0
11	NE	1.8	NE	1.5	NE	1.5	NE	1.0	NE	.9	NE	.9	E	.6	NW	.6	SW	1.0	SW	1.6	NW	1.3	SW	2.0	SW	1.8	SW	2.6
12	NE	1.2	NE	1.3	NE	1.5	NE	1.6	NE	1.5	NE	1.2	NE	1.3	SW	.9	SW	1.5	SW	1.8	SW	1.8	SW	2.6	SW	1.6	SW	2.2
13	NE	.9	NE	.9	NE	1.2	NE	1.0	NE	.6	NE	.6	VR	.6	SW	.6	W	1.0	SW	1.5	SW	1.0	SW	1.0	SW	.6	SW	.6
14	NE	.9	NE	.9	NE	.9	NE	.9	NE	.9	NE	.9	NE	.6	SW	.9	SW	1.0	SW	1.5	SW	1.5	SW	2.0	SW	2.2	SW	2.0
15	SE	.5	NE	.5	NE	.6	SE	.6	NE	1.0	NE	.9	NE	.5	SW	.6	SW	1.0	SW	1.0	SW	.9	SW	1.2	SE	1.2	SW	1.0
16	NE	.6	NE	.6	NE	.6	NE	.6	NE	.6	E	.5	NE	.6	NW	.5	SW	.6	SW	.9	SW	1.8	SW	.9	NW	.9	SW	2.0
17	NE	1.3	NE	1.3	NE	1.5	NE	1.0	NE	1.6	NE	1.5	NE	.6	SW	.9	SW	1.2	SW	1.6	SW	2.0	SW	2.6	SW	2.6	SW	2.0
18	NE	1.5	NE	2.0	NE	1.2	NE	1.5	NE	2.2	NE	1.3	NE	.9	SW	.6	SW	1.3	SW	2.0	SW	1.5	SW	2.2	SW	2.4	SW	2.3
19	NE	.9	NE	1.2	NE	.9	NE	1.0	NE	1.0	NE	1.0	NE	1.2	SW	.9	SW	1.5	SW	1.3	SW	1.8	SW	2.0	SW	2.0	SW	2.0
20	NE	1.8	NE	1.6	NE	1.6	NE	2.2	NE	2.3	NE	2.0	NE	1.8	NE	.6	SW	.9	SW	1.2	SW	1.5	SW	2.4	SW	2.3	SW	2.3
21	NE	1.5	NE	1.5	NE	2.2	NE	2.0	NE	1.5	NE	1.5	NE	.6	SW	1.3	SW	1.5	SW	.9	SW	1.2	SW	.9	SW	1.6	SW	1.6

Ejemplo de un reporte de XM

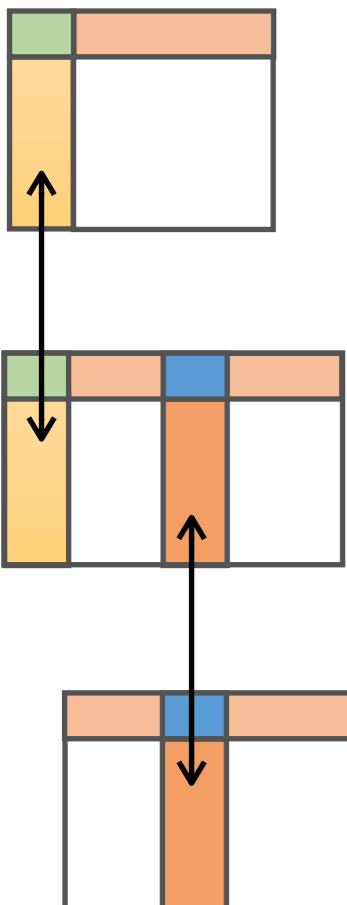
Demanda desagregada a nivel horario.
Registros > 244.000 por año

	A	B	C	D	E	F	G	H	I	J
1	FECHA	CII	Tip	P01	P02	P03	P04	P05	P06	P07
2	2016-01-01	111 R		22,54	21,79	24,04	22,40	22,27	17,57	21,74
3	2016-01-01	111 N		1.499,15	1.378,47	1.398,94	1.353,87	1.299,21	1.261,15	1.165,96
4	2016-01-01	112 R		230,93	228,83	230,83	230,08	227,26	224,94	224,07
5	2016-01-01	112 N		1.052,35	1.066,18	1.067,42	932,78	848,43	821,37	828,92
6	2016-01-01	113 N		829,95	781,96	862,27	765,32	847,33	799,84	972,68
7	2016-01-01	119 N		35,33	37,97	35,04	43,43	46,28	34,90	44,09
8	2016-01-01	121 N		509,30	502,91	505,36	497,51	483,40	674,93	608,67
9	2016-01-01	121 R		13,02	13,12	12,66	12,52	9,90	12,36	11,93
10	2016-01-01	122 R		191,14	191,58	192,88	190,65	181,18	137,20	165,85
11	2016-01-01	122 N		1.902,79	1.867,25	1.837,59	1.812,63	1.789,16	1.765,61	1.699,16
12	2016-01-01	123 N		136,76	136,52	132,77	133,80	131,66	131,97	127,93
13	2016-01-01	123 R		14,38	15,02	14,92	13,82	13,63	14,82	17,69
14	2016-01-01	124 N		2.358,04	2.264,04	2.236,77	2.208,11	2.003,35	1.931,69	1.842,89
15	2016-01-01	124 R		465,83	468,99	415,05	370,19	386,83	474,15	548,28
16	2016-01-01	125 N		13.937,40	14.101,58	13.783,46	13.362,54	12.159,18	10.116,37	9.907,91

Se desea obtener la curva típica de consumo para el día

Almacenamiento de los datos

RDBMS

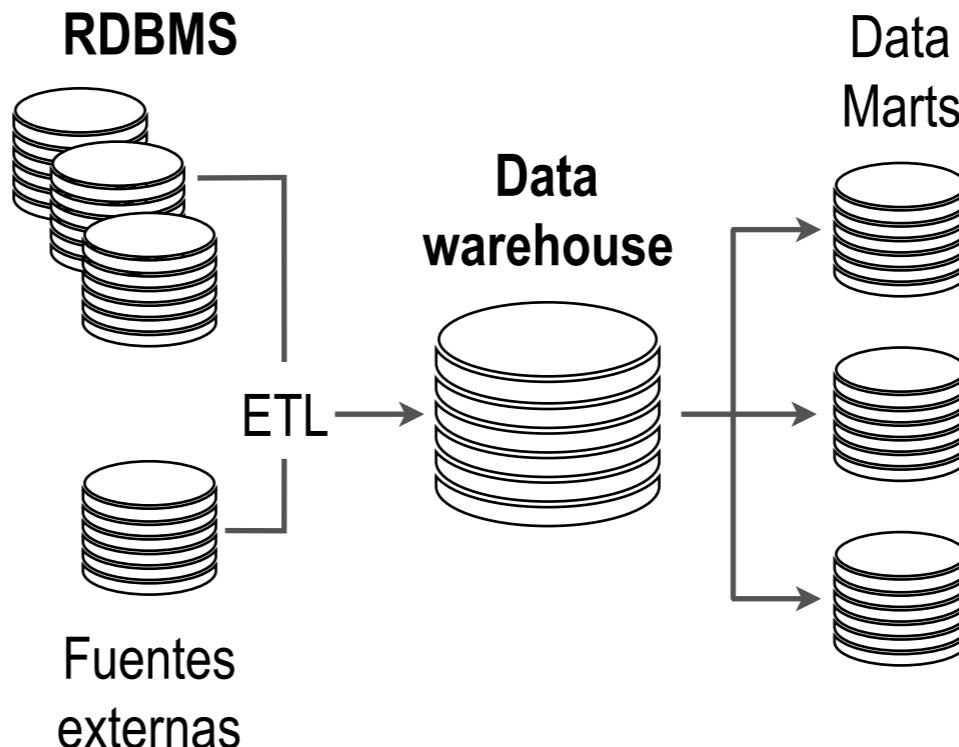


Componentes

- Esquemas
- Tablas
- Consultas
- Reportes
- Vistas
- Otros elementos

Principales RDBMDS

- Oracle
- PostgreSQL
- Microsoft SQL
- Reportes
- Vistas
- Otros elementos



Business Intelligence

Análisis de información en RDBMS para toma de decisiones de negocio.

Data warehouse / OLAP

Bodegas de datos / Procesamiento analítico en línea

- Estructurado
- Integrado
- No volátil (permanencia de la información)
- Variable en el tiempo
- Orientado al análisis y la divulgación de la información
- Cubos multidimensionales

Inteligencia de Negocios

GENSCAPE™

Solutions

Knowledge Center

Events

Blog

News

About

Oil

Power

Natural

Overview

Daily Macro Supply & Demand D

Equity Production Insight

Intrastate Storage Monitoring

Natural Gas Analyst

GENSCAPE™

Solutions

Knowledge Center

Events

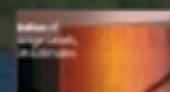
Blog

News

About

The Industry's
Only Proprietary
Natural Gas Data

LEARN MORE ➤



Inteligencia de Negocios

información
INTELIGENTE

Sign In Buscar Metadata

Inicio Demanda Hidrología Intercambios Oferta Transacciones y Precios

Información Inteligente

Indicadores

12,810.78 GWh Volumen Útil Diario Valor Anterior: 12,793.21 GWh ↑14% 2017-10-03 Ver Más	198.01 GWh Aportes Valor Anterior: 210.60 GWh ↓6.36% 2017-10-03 Ver Más	188.14 GWh Demanda Energía SIN Preliminar Valor Anterior: 184.55 GWh ↑1.91% 2017-10-03 Ver Más
155.15 GWh Generación Hidráulica Valor Anterior: 132.90 GWh ↑14.34% 2017-10-02 Ver Más	17.01 GWh Generación Térmica Valor Anterior: 13.45 GWh ↑20.93% 2017-10-02 Ver Más	0.09 GWh Importaciones Preliminar Valor Anterior: 0.12 GWh ↓33.33% 2017-10-03 Ver Más
110.06 \$/kWh Precio Bolsa Promedio Aritmético TX1 Valor Anterior: 100.25 \$/kWh ↑8.91% 2017-10-02 Ver Más	32,521.85 \$M Transacciones en contratos Valor Anterior: 32,521.85 \$M ↔0.00% 2017-09-30 Ver Más	2,918.62 \$M Restricciones sin Alivios Valor Anterior: 2,618.43 \$M ↑10.29% 2017-10-02 Ver Más

Inteligencia de Negocios

Información Inteligente > Transacciones > Histórico Transacciones

Tipo	Nombre	Descripción Contenido	Tamaño de archivo	Modificado
📁	AGC			25/01/2013 9:47
📁	Bolsa			20/11/2012 8:20
📁	Contratos			19/11/2012 16:33
📁	Desviaciones			19/11/2012 16:38
📁	Precios			19/11/2012 16:39
📁	Reconciliaciones			19/11/2012 16:29
📁	Servicios			08/08/2012 16:18
📄	Compras_Bolsa_Internacional_(kWh)_2016		106 KB	01/10/2017 14:42
📄	Compras_Bolsa_Internacional_(kWh)_2017		77 KB	03/10/2017 14:30
📄	Compras_Bolsa_Nacional_(kWh)_2016		6469 KB	01/10/2017 14:47
📄	Compras_Bolsa_Nacional_(kWh)_2017		4756 KB	03/10/2017 14:41
📄	Compras_Bolsa_TIE_(kwh)_2016		76 KB	01/10/2017 14:48
📄	Compras_Bolsa_TIE_(kwh)_2017		53 KB	03/10/2017 14:50
📄	Compras_Contrato_(kWh)_2016		7053 KB	01/10/2017 14:52
📄	Compras_Contrato_(kWh)_2017		5014 KB	03/10/2017 9:46
📄	Costo_Marginal_Despacho_Programado_(SkWh)_2016		123 KB	01/10/2017 14:54
📄	Costo_Marginal_Despacho_Programado_(SkWh)_2017		112 KB	03/10/2017 9:48
📄	Delta_Internacional_Delta_Nacional_(SkWh)_2016		44 KB	01/10/2017 14:57
📄	Delta_Internacional_Delta_Nacional_(SkWh)_2017		39 KB	03/10/2017 9:51
📄	Desviaciones_(kWh)_2016		333 KB	01/10/2017 15:00

Inteligencia de Negocios

energone

in t

- HOME
- ABOUT US
- MARKETS SERVED
- PRODUCTS
- SERVICES
- INVESTORS
- CONTACT US

EnergyDashboard

The EnergyDashboard enables managers of wholesale energy portfolios to see and manage "at a glance" ALL the key features, status and requirements of all their wholesale energy operations, such as:

- Market data, prices and chosen analytics
- Portfolio status (bid compliance in multiple markets)
- Contracted position
- Energy operations (workflow) status and alerts

Furthermore, users can easily switch between the various aspects of the operational functions, seamlessly moving between market operations and bidding and contracts.

CONTACT US

Inteligencia de Negocios

[Login](#)



[Crude Oil](#) [Oil Products](#) [Coal](#) [News & Insight](#) [About Arg](#)
[LPG/NGL](#) [Bioenergy](#) [Power](#) [Argus Consulting](#) [Methodolo](#)
[Emissions](#) [Natural Gas/LNG](#) [Transportation](#) [Services](#) [Reference](#)
[Fertilizer](#) [Petrochemicals](#) [Metals](#) [Events](#) [Contact Us](#)

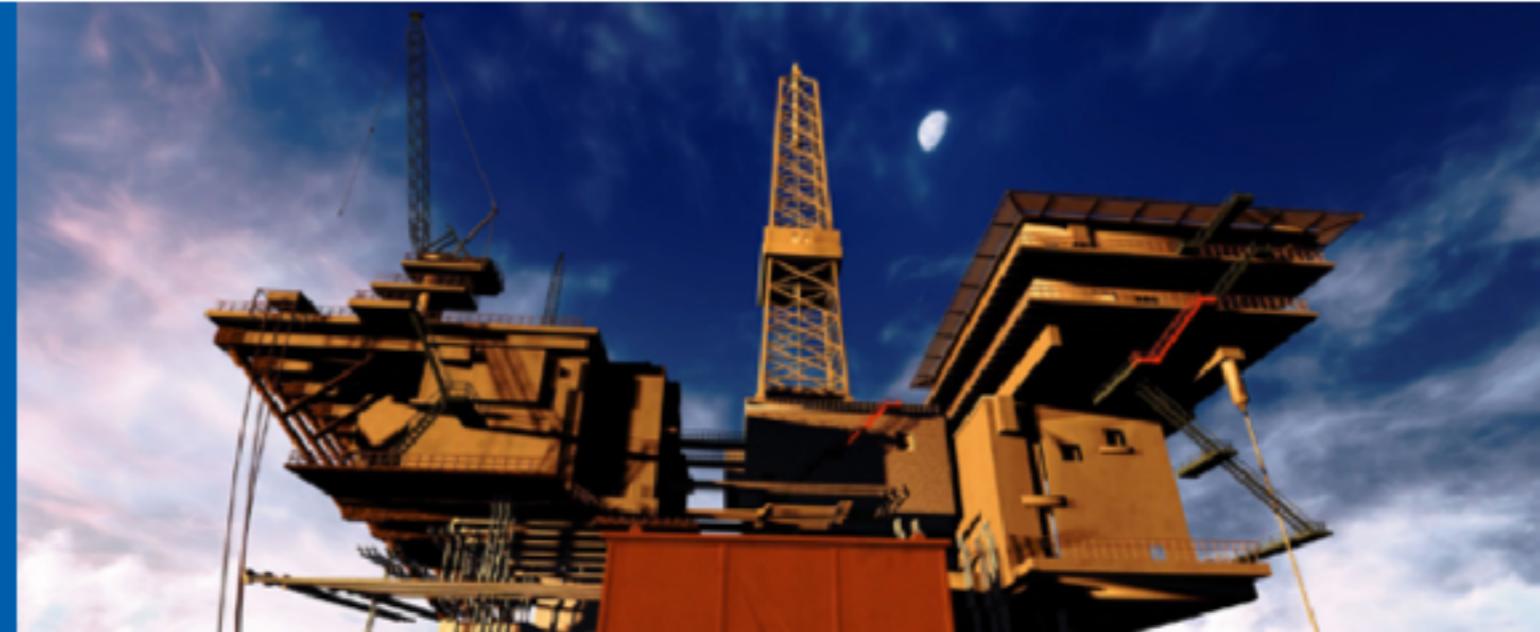
[Home](#) / Crude Oil

[Русский](#) | [Portuguese](#) | [日本](#) | [中国](#) | [Español](#)

Crude Oil

Argus crude services cover global and specialist local markets, offering exclusive prices and insight.

Our product range includes daily international prices and trading information, intelligent analysis and fundamental data on supply and demand, stock levels as well as logistical updates, events and consulting services.



Watch the latest Irma and Harvey updates here: Week 2 Price Effects of Harvey

What will happen to crude in the wake of Irma? What is happening with the US Gulf Coast in Harvey recovery?

[Watch](#)

Almacenamiento de los datos

Datos tabulares

KEY	Fecha	Planta	Generación
001	2017-10-01	Jaguas	100.2
002	2017-10-01	Playas	23.1
003	2017-10-01	Guatape	130.1

Document (JSON/XML)

```
[  
  {  
    Fecha:2017-10-01,  
    Planta:Jaguas,  
    Generación: 100.2  
  },{  
    Fecha:2017-10-01,  
    Planta:Playas,  
    Generación:23.1,  
  },{  
    Fecha:2017-10-01,  
    Planta:Guatapé,  
    Generación:130.1  
  }]  
]
```

Pares <clave, valor>

Tabla001.Fecha=2017-10-01
Tabla001.Planta=Jaguas
Tabla001.Generación=100.2
Tabla002.Fecha=2017-10-01
Tabla002.Planta=Playas
Tabla002.Generación=23.1
Tabla003.Fecha=2017-10-01
Tabla003.Planta=Guatapé
Tabla003.Generación=130.1

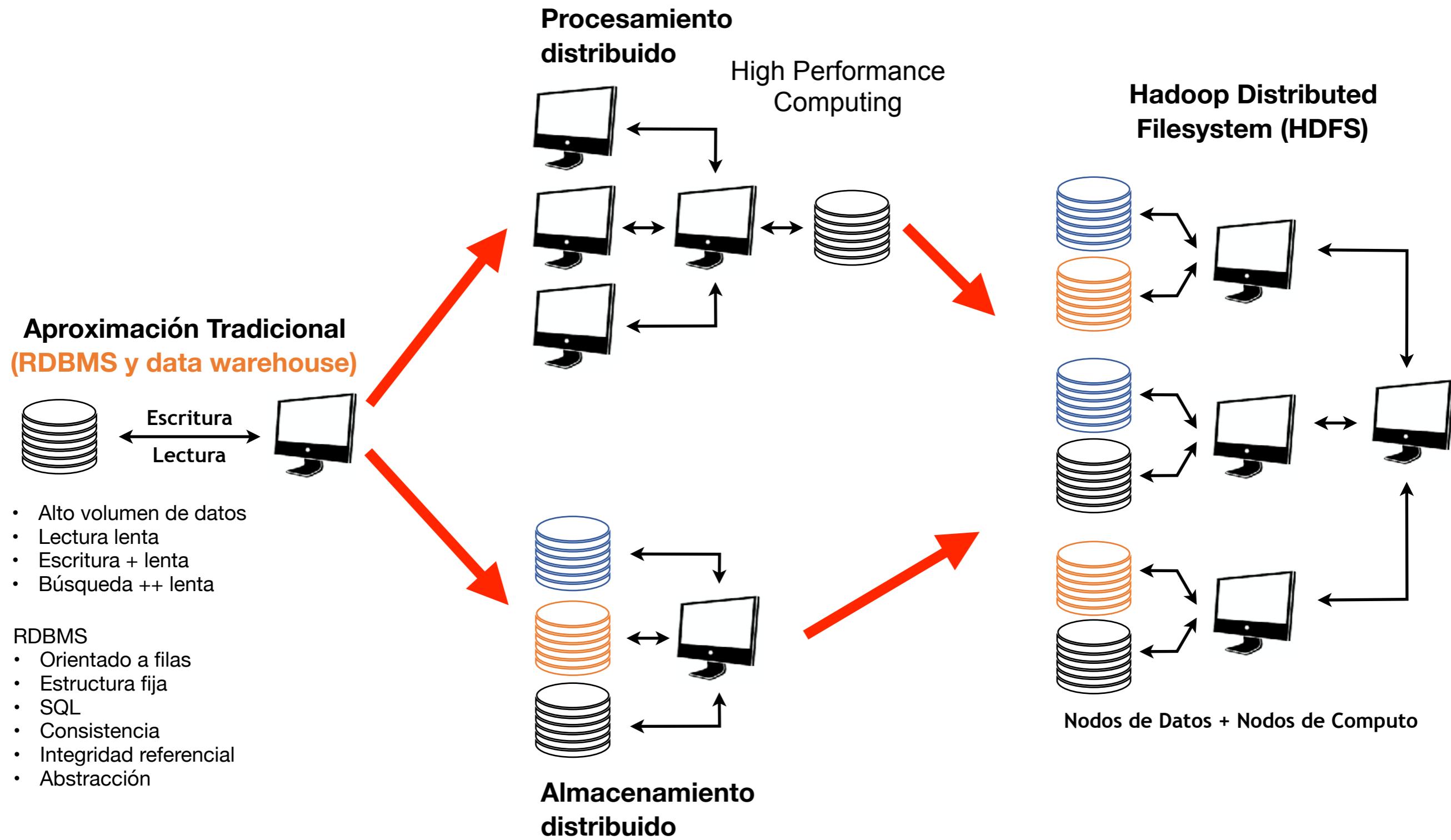
Sistema orientado a filas

001:2017-10-01,Jaguas,100.2
002:2017-10-01,Playas,23.1
003:2017-10-01,Guatape,130.1

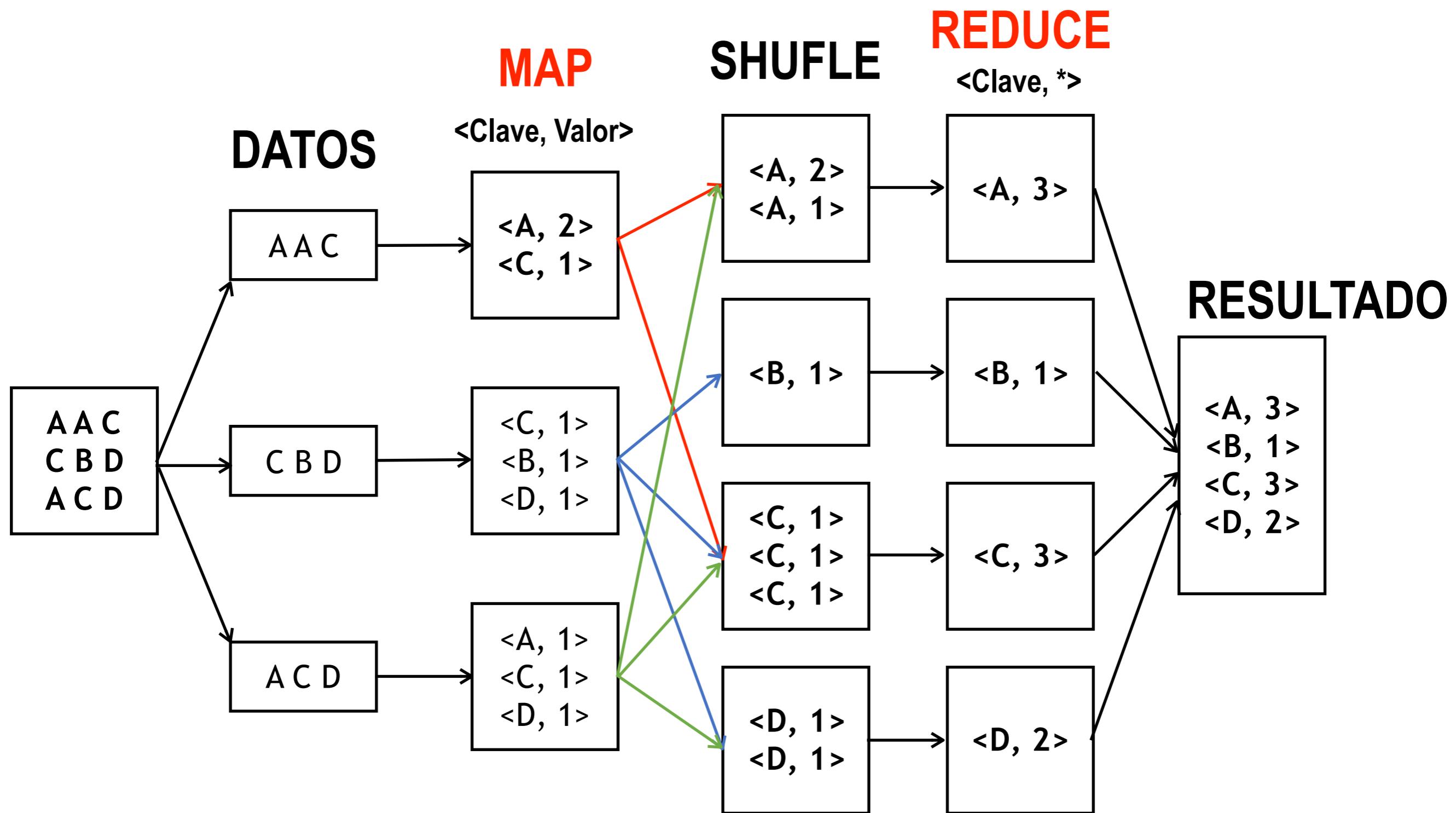
Column family database

001:{Fecha:2017-10-01, Planta:Jaguas, Generación:100.2}
002:{Fecha:2017-10-01, Planta:Playas, Generación:23.1}
003:{Fecha:2017-10-01, Planta:Guatapé, Generación:130.1}

Big Data (Hadoop)



Hadoop / MapReduce



Ecosistema Apache Hadoop

HDFS
Operaciones básicas del sistema de archivos

MapReduce
(Java)

Hive
Lenguaje de consultas similar a SQL.

Mahout
Algoritmos de Machine Learning implementados usando MapReduce (Java)

Pig
Lenguaje para el procesamiento de datos.

HBase
Base de datos orientada a columnas.

Sqoop
RDBMS <--> Hadoop

AVRO
Sistema de serialización de datos

ZooKeeper
Coordinación

YARN
Programación de trabajos y manejo de recursos en clusters

CASSANDRA
Sistema de Base de Datos

Oozie
Flujo y programación de trabajos

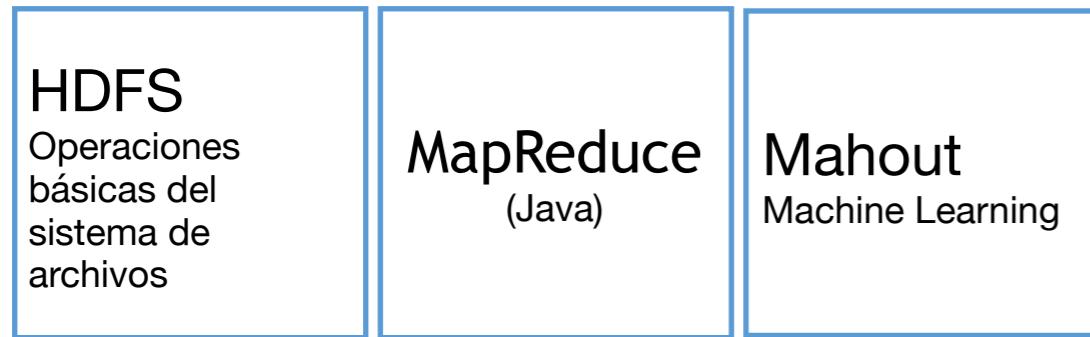
Ejemplo de Pig

```
records = LOAD 'sample.txt' AS (year:chararray, temperature:int, quality:int);
filtered_records = FILTER records BY temperature;
grouped_records = GROUP filtered_records BY year;
max_temp = FOREACH grouped_records GENERATE group, MAX(filtered_records.temperature);
DUMP max_temp;
```

Ejemplo de Hive

```
CREATE TABLE records (year STRING, temperature INT, quality INT)
ROW FORMAT DELIMITED FIELDS TERMINATED BY '\t';
LOAD DATA LOCAL INPATH 'sample.txt' OVERWRITE INTO TABLE records;
SELECT year, MAX(temperature) FROM records GROUP BY year;
```

Mahout, MLib, R y Python

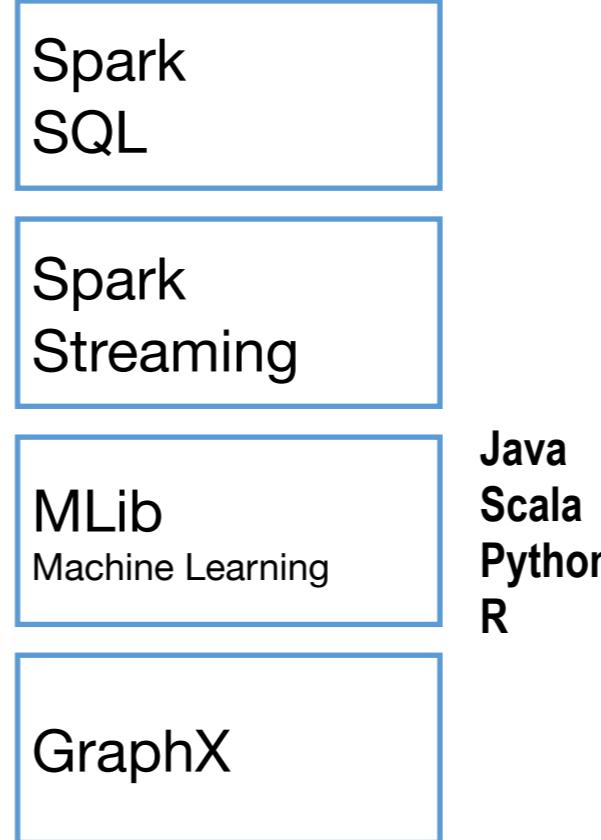
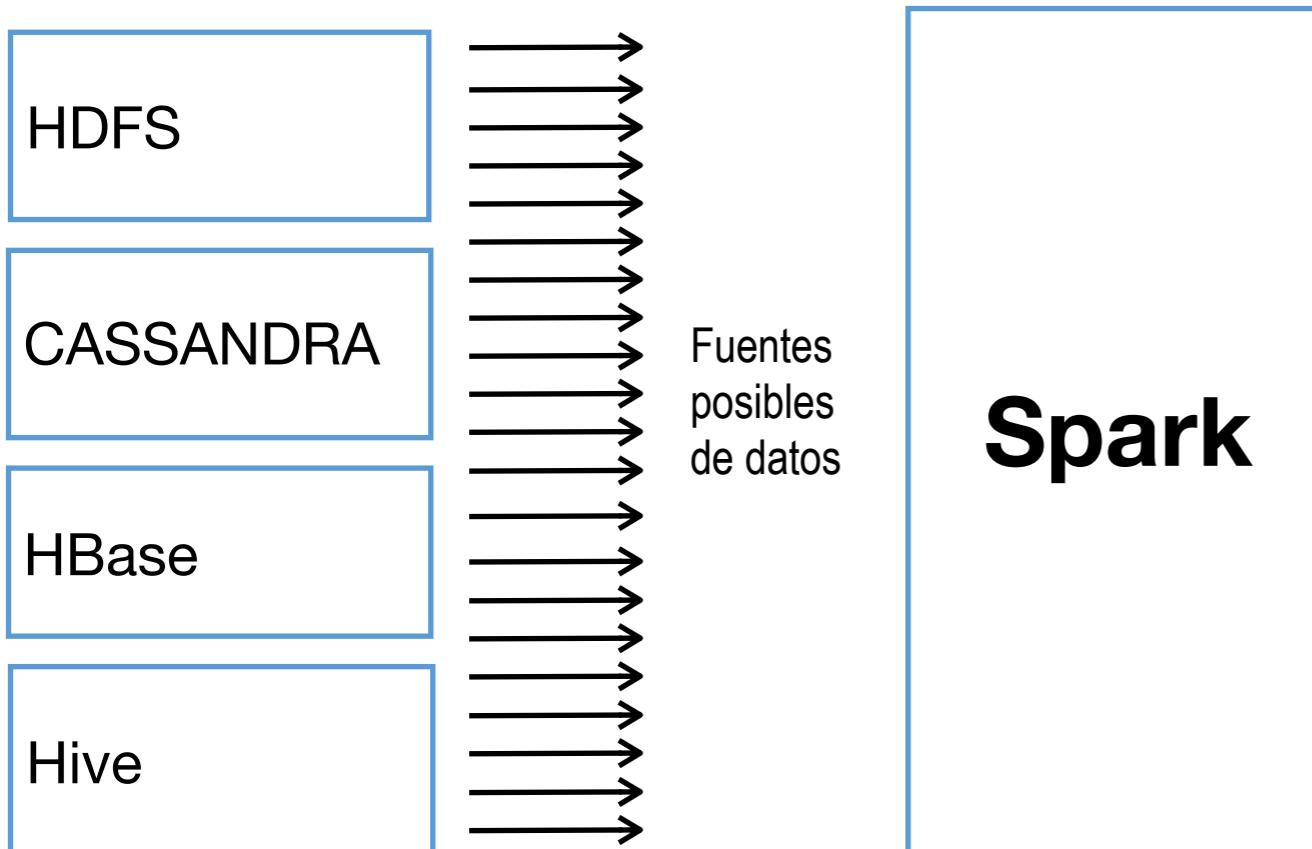


Regresión logística
Regresión lineal
Clustering
Filtrado colaborativo

RHadoop
rdfs
rnr
rhbase

<http://mahout.apache.org/users/basics/algorithms.html>

Hadoop / MapReduce



<https://spark.apache.org/mllib/>

Mahout vs MLib

MLib Machine Learning

Algorithms

MLlib contains many algorithms and utilities.

ML algorithms include:

- Classification: logistic regression, naive Bayes,...
- Regression: generalized linear regression, survival regression,...
- Decision trees, random forests, and gradient-boosted trees
- Recommendation: alternating least squares (ALS)
- Clustering: K-means, Gaussian mixtures (GMMs),...
- Topic modeling: latent Dirichlet allocation (LDA)
- Frequent itemsets, association rules, and sequential pattern mining

ML workflow utilities include:

- Feature transformations: standardization, normalization, hashing,...
- ML Pipeline construction
- Model evaluation and hyper-parameter tuning
- ML persistence: saving and loading models and Pipelines

Other utilities include:

- Distributed linear algebra: SVD, PCA,...
- Statistics: summary statistics, hypothesis testing,...

Collaborative Filtering with CLI drivers

User-Based Collaborative Filtering

Item-Based Collaborative Filtering

Matrix Factorization with ALS

Matrix Factorization with ALS on Implicit Feedback

Weighted Matrix Factorization, SVD++

Classification with CLI drivers

Logistic Regression - trained via SGD

Naive Bayes / Complementary Naive Bayes

Hidden Markov Models

Clustering with CLI drivers

Canopy Clustering

k-Means Clustering

Fuzzy k-Means

Streaming k-Means

Spectral Clustering

Dimensionality Reduction note: most scale reduction algorithms are available through the [Core Library for all engines](#)

Singular Value Decomposition

Lanczos Algorithm

Stochastic SVD

PCA (via Stochastic SVD)

QR Decomposition

Topic Models

Latent Dirichlet Allocation

Miscellaneous

RowSimilarityJob

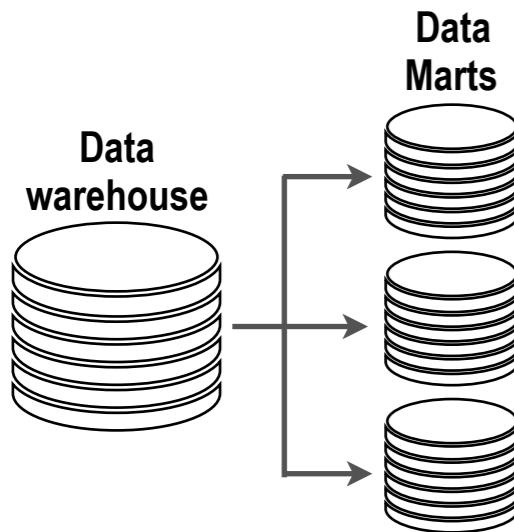
Collocations

Sparse TF-IDF Vectors from Text

XML Parsing

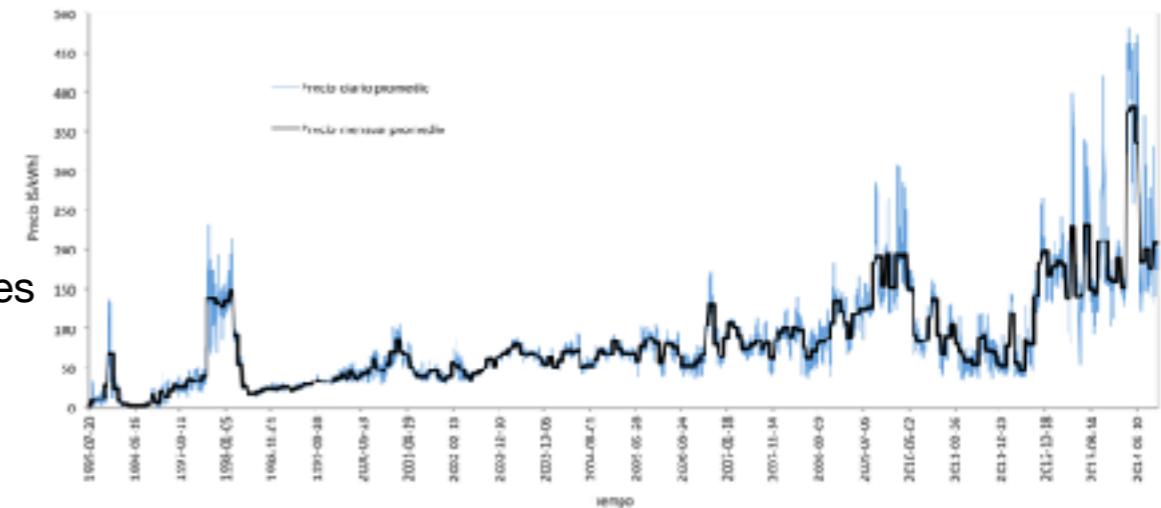
Email Archive Parsing

Evolutionary Processes



Data Mining

Proceso de descubrimiento de patrones y tendencias útiles en grandes conjuntos de datos.



Data Science

Área relacionada con los procesos y sistemas para la extracción de conocimiento de datos almacenados electrónicamente (¿para la toma de decisiones? ¿para probar hipótesis?)

Analytics

Proceso científico de transformación de datos en conocimiento para mejorar el proceso de toma de decisiones [Informs].

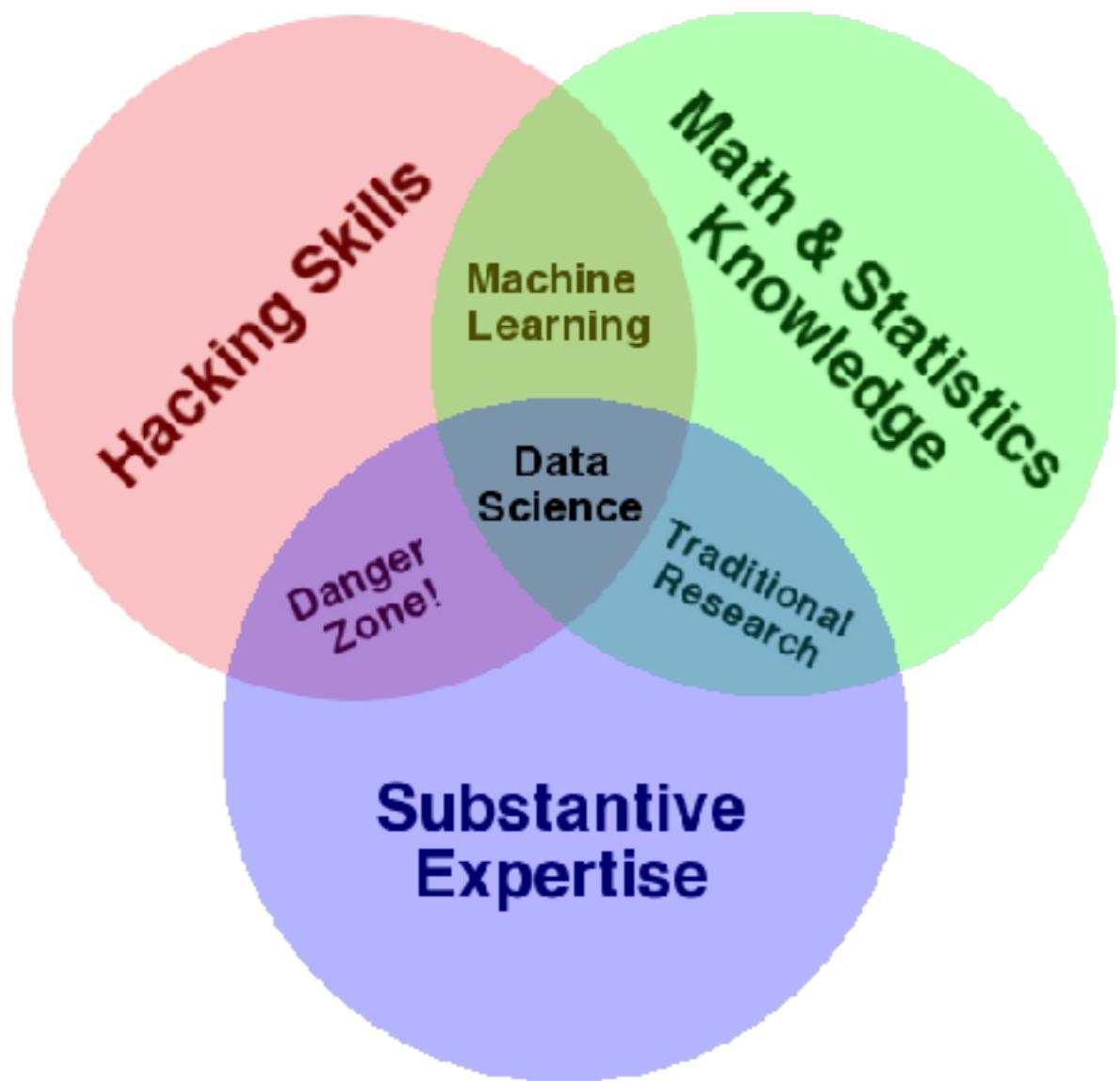
- I. Problema organizacional
- II. Transformación en un problema de analytics
- III. Datos
- IV. Selección de la metodología
- V. Desarrollo del modelo
- VI. Puesta en marcha (deploy)
- VII. Gestión del ciclo de vida del modelo

Big Data Analytics

Analytics sobre grande conjuntos de datos

Algoritmos, técnicas y metodologías

Diagrama de Ven explicando qué es Data Science y habilidades requeridas



Data Science and Data Scientists: What's in a Name?

Saunders, 2013

Data Architect Data Engineer

Diseño y estructura de las bases de datos.

Data Manager

Gestiona la creación y mantenimiento de las bases de datos.

ETL Developer

Gestiona la extracción, transformación y carga de los datos a las bases de datos.

Data Analyst

Fuentes y usos de los datos.

Business Intelligence Practitioner

Combinación de negocios + tecnología con el fin de proveer información a las unidades de negocios para toma de decisiones

Data Scientist

Habilidades en la programación de computadores para manejo de datos y modelado predictivo (estadística, aprendizaje de máquinas, minería de datos, etc.).

Analytics Practitioner

Data Science + Optimización + Simulación

Disciplina	Tecnología	Habilidades	Foco
Inteligencia de Negocios	<ul style="list-style-type: none">ETL/SQLRDBMSReportesVisualización	<ul style="list-style-type: none">ProgramaciónAnálisis de datosModelado de datosDesarrollo de reportesEstadística BásicaAnálisis del negocio & EstrategiaPresentación oral	<ul style="list-style-type: none">Suministro de información y reporteVisualización de datosEstadísticos descriptivosIntegración de datos y consolidación
Análisis de datos	<ul style="list-style-type: none">Software para modelado de datosSoftware para diagramaciónSoftware para documentaciónSQLSoftware para perfilado de datos	<ul style="list-style-type: none">Modelado de datosAnálisis del negocioManipulación de datosEstadística básica	<ul style="list-style-type: none">Reglas de negocioDefinición de datosRelaciones entre datosAtributos de datosEstructuras de datosFuentes y usos de datosCalidad de datos
Ciencia de los Datos (Analytics)	<ul style="list-style-type: none">Software estadísticoDatos columnaresMap-ReduceNoSQLLenguajes de programaciónSoftware para graficaciónSoftware para optimización, simulación, predicción y análisis de decisiones	<ul style="list-style-type: none">Estadística avanzadaProgramaciónAnálisis del negocioArquitecturas y tecnologías modernas para el manejo de datosDesarrollo de productos de datosSimulación de sistemasOptimizaciónPredicción	<ul style="list-style-type: none">Modelado predictivoAnálisis estadístico avanzadoMinería de datosManejo de datos no estructuradosManejo de grandes volúmenes de datosI+DAnálisis de decisiones

Similitudes y Diferencias

DATA SCIENCE

Programación.

Adquisición, limpieza, preprocesamiento y visualización de datos.

Investigación reproducible.

Modelado de Datos (minería de datos)

Inferencia Estadística.

Modelos estadísticos

Aprendizaje de Máquinas

Productos de Datos.

ANALYTICS

Programación

Adquisición, limpieza, preprocesamiento y visualización de datos.

Modelado de Datos (modelado predictivo)

Inferencia Estadística.

Modelos estadísticos

Aprendizaje de Máquinas

Productos de Datos.

Inteligencia de Negocios.

Simulación.

Optimización.

Métodos prescriptivos: modelos predictivos + optimización.

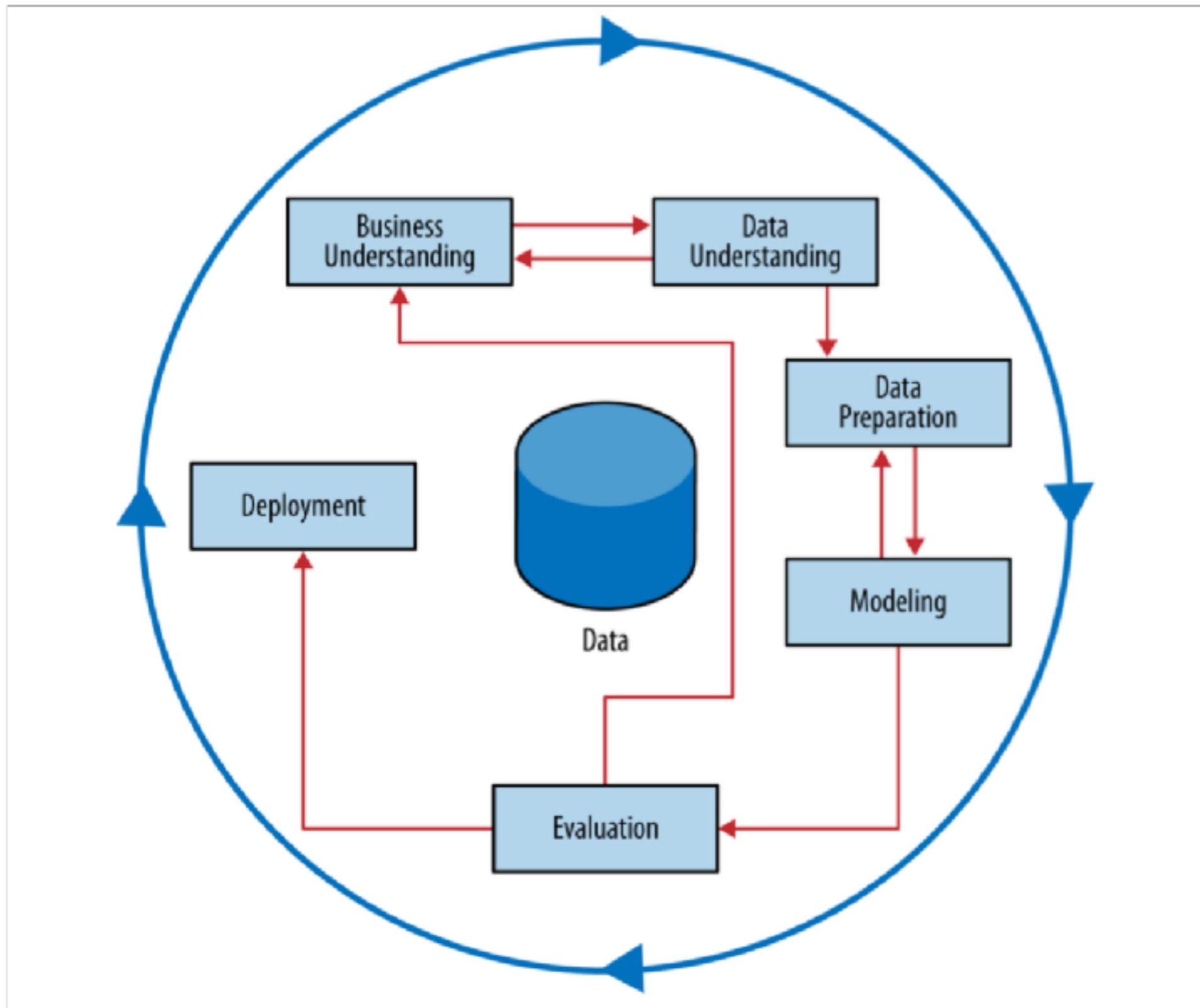
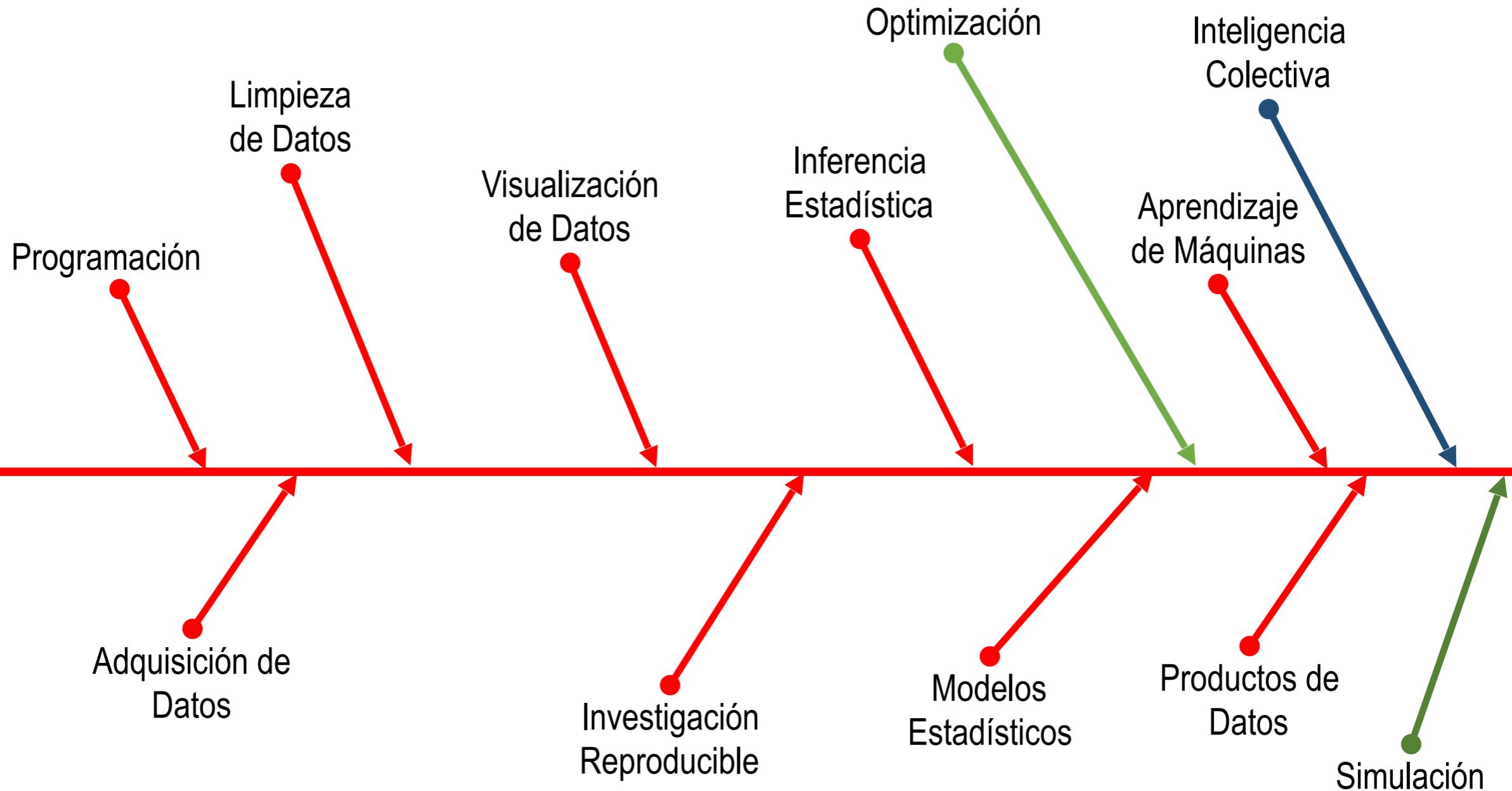


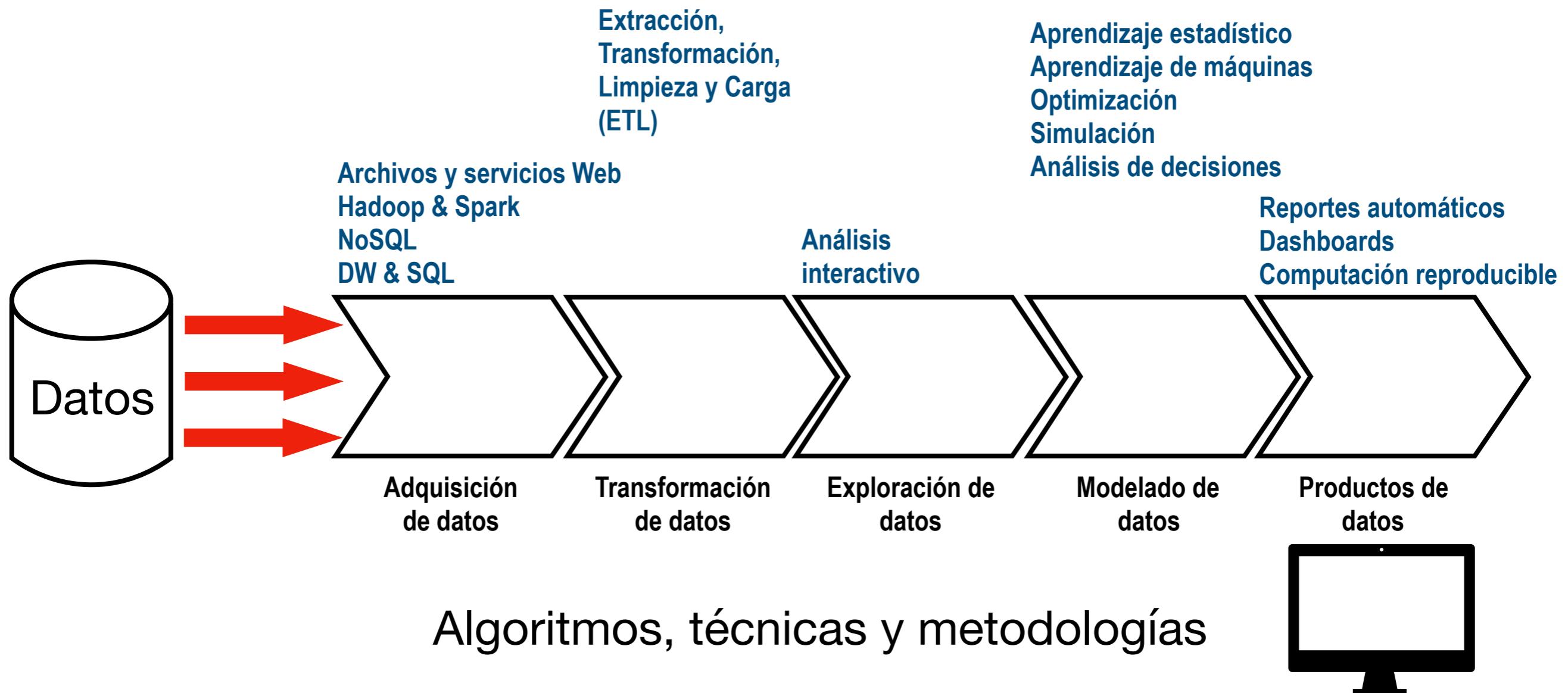
Figure 2-2. The CRISP data mining process.

Componentes de Data Science / Analytics



Data-driven decision making!

Proceso de modelado en Analytics



Infraestructura
computacional

{ Un procesador
Muchos procesadores

{ Computación en máquinas locales
Computación en la nube

Programación -- ¿Usted sabe programar ... / Es capaz de ...?

¿Ordenar un vector de números?

Programación para
ingeniería
Cómputo numérico.

¿Calcular la suma de los primeros 20 números primos?

¿Computar la inversa de una matriz?

Programación para
Computer Sciences

Manipulación de texto.

```

echo "ESTACION;FECHA;ANO;MES;DIA;HORA;HHMMSS;DIRECCION;VELOCIDAD" > datos
tail +2 AQUITANIA.csv >> datos

## Elimina lineas vacias
sed -e '/^$/d' datos > out.1

## borra lineas en blanco
sed -e '/;;;/d' out.1 > datos

## llena las horas vacias
sed -e 's/;;;/00:00:00;/g' datos > out.1

## etcetera ...

## promedio para cada hora
csvsql --query "select ESTACION, FECHA, ANO, MES,
DIA, HORA, DIRECCION, avg(VELOCIDAD) as VELOCIDAD from 'out'
group by ESTACION, FECHA, HORA" out.5 > out.6

```

ESTACION;FECHA;HORA;DIRECCION;VELOCIDAD
AQUITANIA;2005-04-16;11:10:00;135;6,3
AQUITANIA;2005-04-16;11:20:00;135;5,1
AQUITANIA;2005-04-16;11:30:00;135;6,3
AQUITANIA;2005-04-16;11:40:00;113;6,1
AQUITANIA;2005-04-16;11:50:00;135;4,1
AQUITANIA;2005-04-16;12:00:00;135;5,5
AQUITANIA;2005-04-16;12:10:00;135;5,4
AQUITANIA;2005-04-16;12:20:00;135;5,5
AQUITANIA;2005-04-16;12:30:00;90;4,6
AQUITANIA;2005-04-16;12:40:00;90;6,7

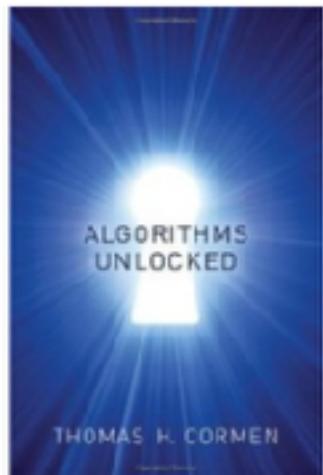
ESTACION,FECHA,ANO,MES,DIA,HORA,DIRECCION,VELOCIDAD
AQUITANIA,2005-04-16,2005,4,16,11,135,5.58
AQUITANIA,2005-04-16,2005,4,16,12,90,5.45
AQUITANIA,2005-04-16,2005,4,16,13,135,4.86666666666667
AQUITANIA,2005-04-16,2005,4,16,14,135,3.666666666666665
AQUITANIA,2005-04-16,2005,4,16,15,135,3.466666666666667
AQUITANIA,2005-04-16,2005,4,16,16,135,3.699999999999993
AQUITANIA,2005-04-16,2005,4,16,17,135,4.83333333333333
AQUITANIA,2005-04-16,2005,4,16,18,135,4.76666666666667
AQUITANIA,2005-04-16,2005,4,16,19,135,4.350000000000005
AQUITANIA,2005-04-16,2005,4,16,20,135,2.68333333333333
AQUITANIA,2005-04-16,2005,4,16,21,135,3.199999999999997

Ejemplo de Pig

```
records = LOAD 'sample.txt' AS (year:chararray, temperature:int, quality:int);
filtered_records = FILTER records BY temperature;
grouped_records = GROUP filtered_records BY year;
max_temp = FOREACH grouped_records GENERATE group, MAX(filtered_records.temperature);
DUMP max_temp;
```

Ejemplo de Hive

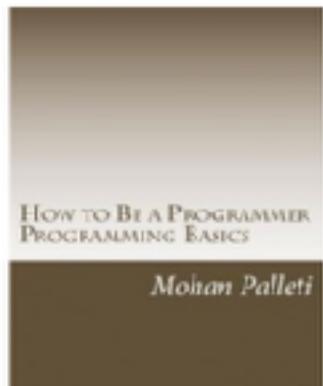
```
CREATE TABLE records (year STRING, temperature INT, quality INT)
ROW FORMAT DELIMITED FIELDS TERMINATED BY '\t';
LOAD DATA LOCAL INPATH 'sample.txt' OVERWRITE INTO TABLE records;
SELECT year, MAX(temperature) FROM records GROUP BY year;
```



Algorithms Unlocked

By: Thomas H. Cormen

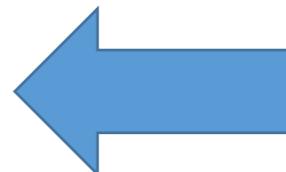
Have you ever wondered how your GPS can find the fastest way to your destination, selecting one route from seemingly countless possibilities in mere seconds? How your credit card account number is protected when you make a purchase over the Internet? The answer is algorithms. And how do...



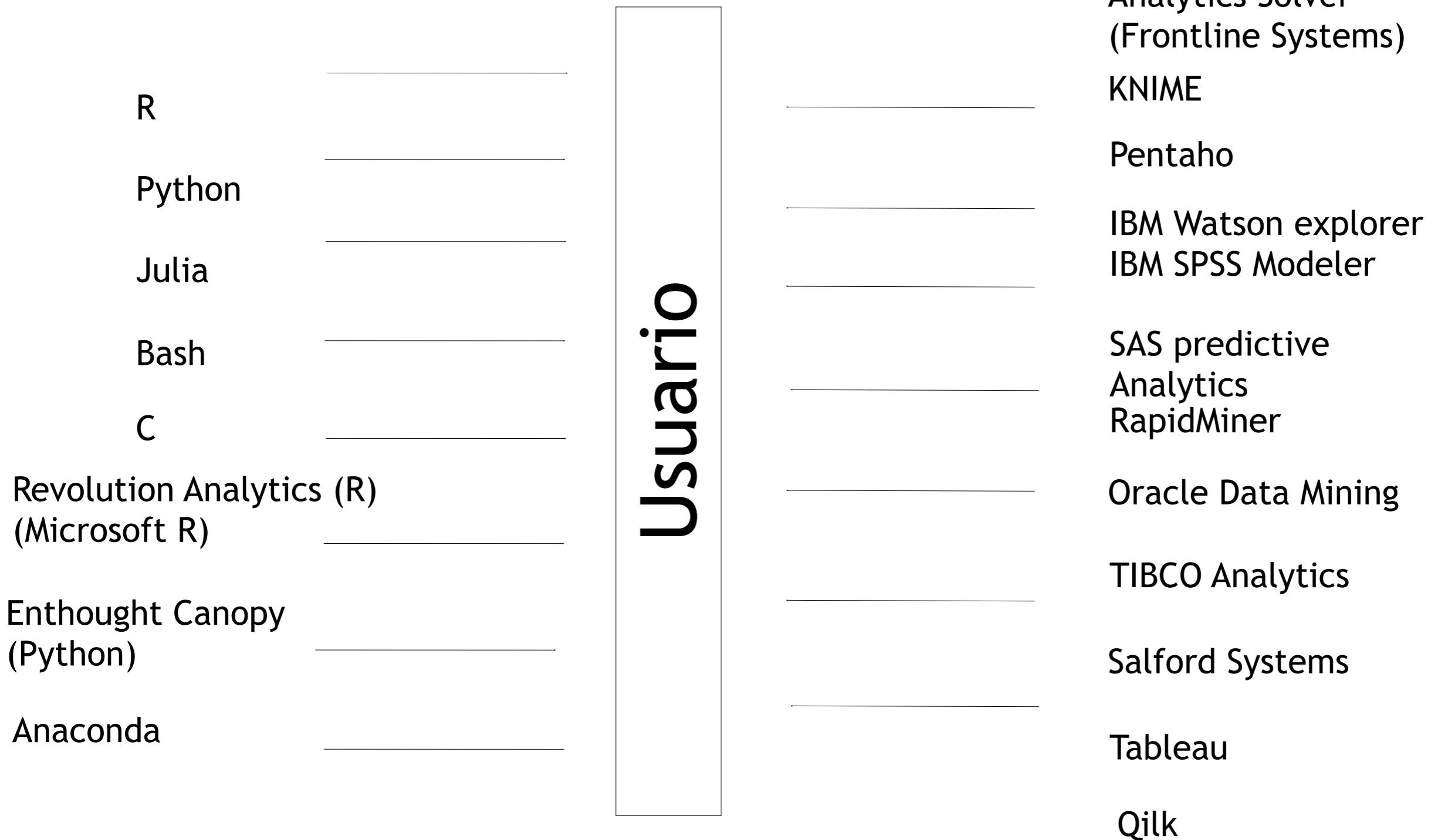
How to Be a Programmer: Programming Basics

By: Mohan Palleti

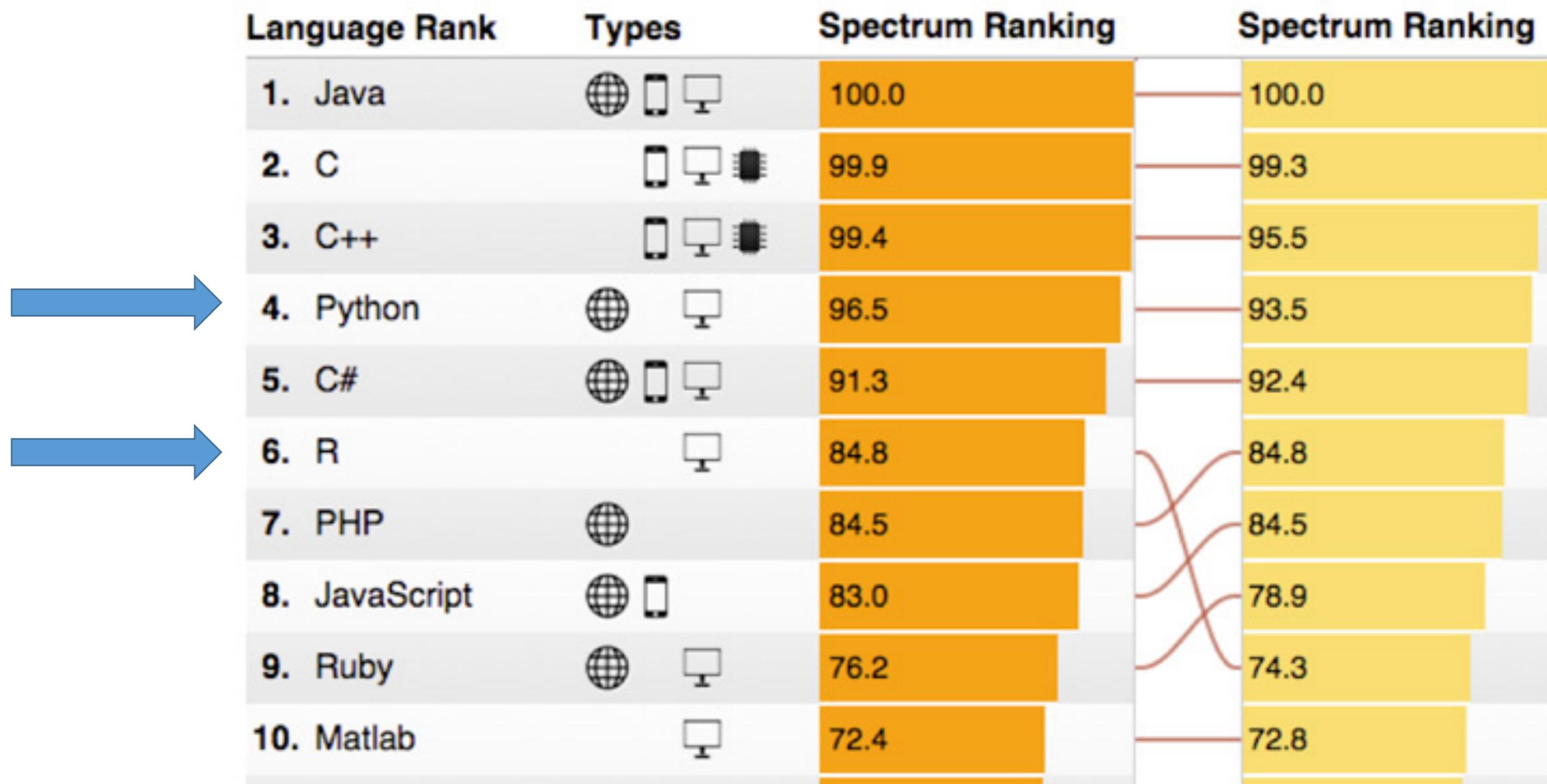
A Self-help 97 pages book to learn the basics of programming using Microsoft Excel's VBA tools. Ideal resource for school teachers and educators wanting to teach programming basics.



Programación vs Aplicaciones de usuario final



The 2015 Top Ten Programming Languages (IEEE Spectrum)



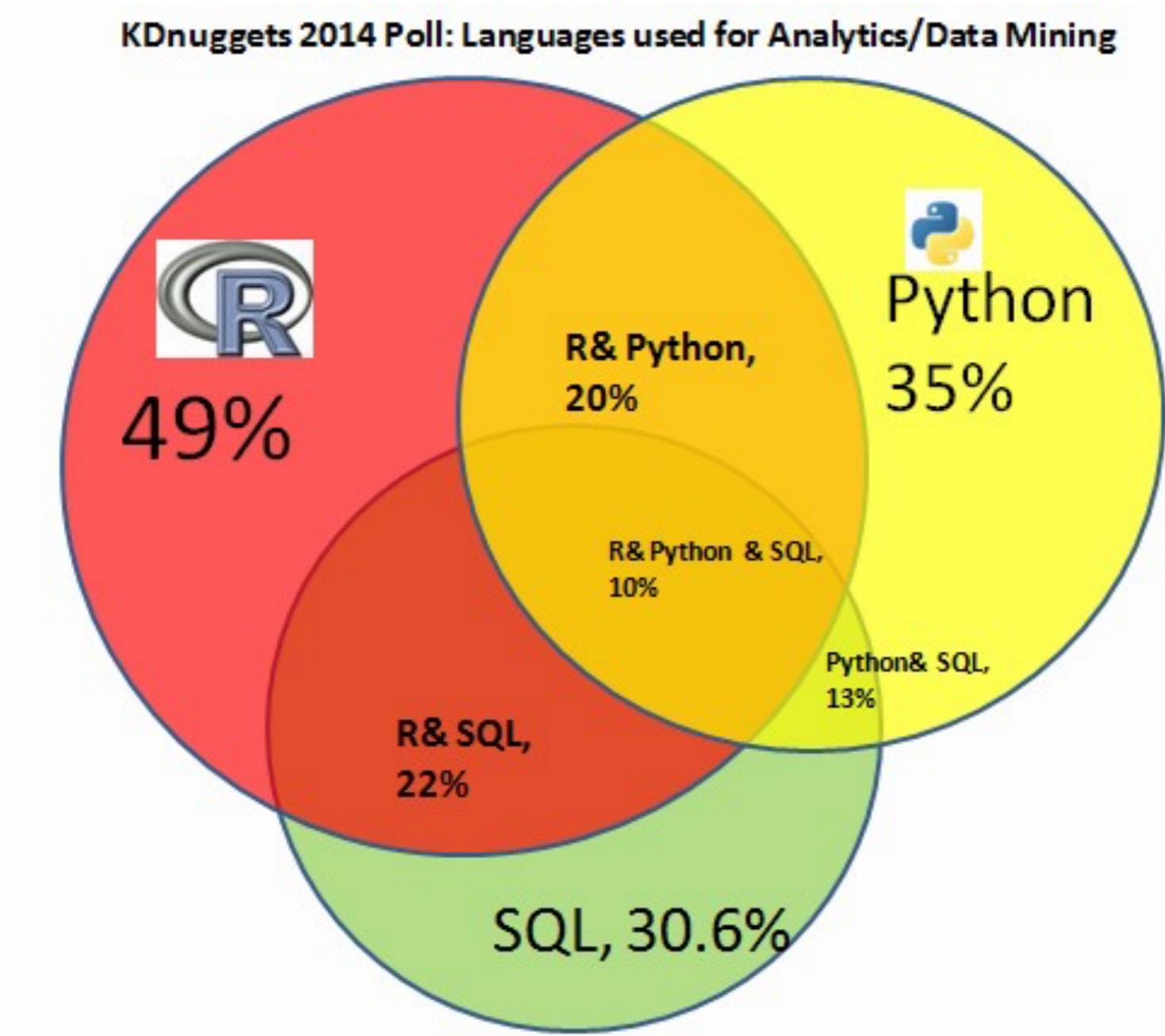
The 2016 Top Ten Programming Languages (IEEE Spectrum)

Language Rank	Types	Spectrum Ranking
1. C		100.0
2. Java		98.1
3. Python		98.0
4. C++		95.9
5. R		87.9
6. C#		86.7
7. PHP		82.8
8. JavaScript		82.2
9. Ruby		74.5
10. Go		71.9

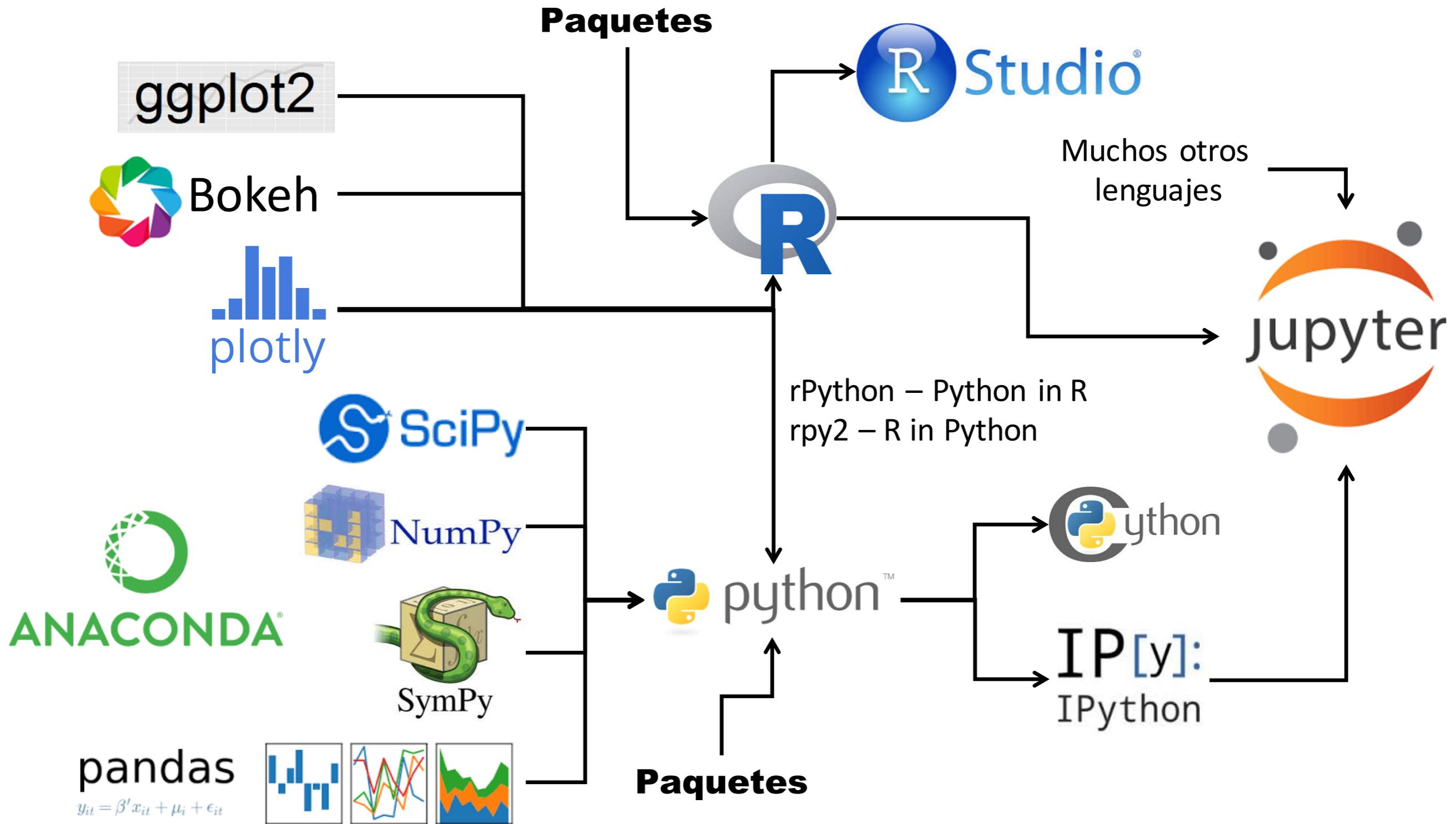
The 2017 Top Ten Programming Languages (IEEE Spectrum)

Language Rank	Types	Spectrum Ranking
1. Python		100.0
2. C		99.7
3. Java		99.5
4. C++		97.1
5. C#		87.7
6. R		87.7
7. JavaScript		85.6
8. PHP		81.2
9. Go		75.1
10. Swift		73.7

Popularidad de los lenguajes



Ecosistema de herramientas y lenguajes de programación



Adquisición y Limpieza de Datos

TXT, Excel, CSV, PDF, *.docx.

Páginas web (HTML) y Google Groups.

Bases de datos relacionales.

Lenguaje Natural.

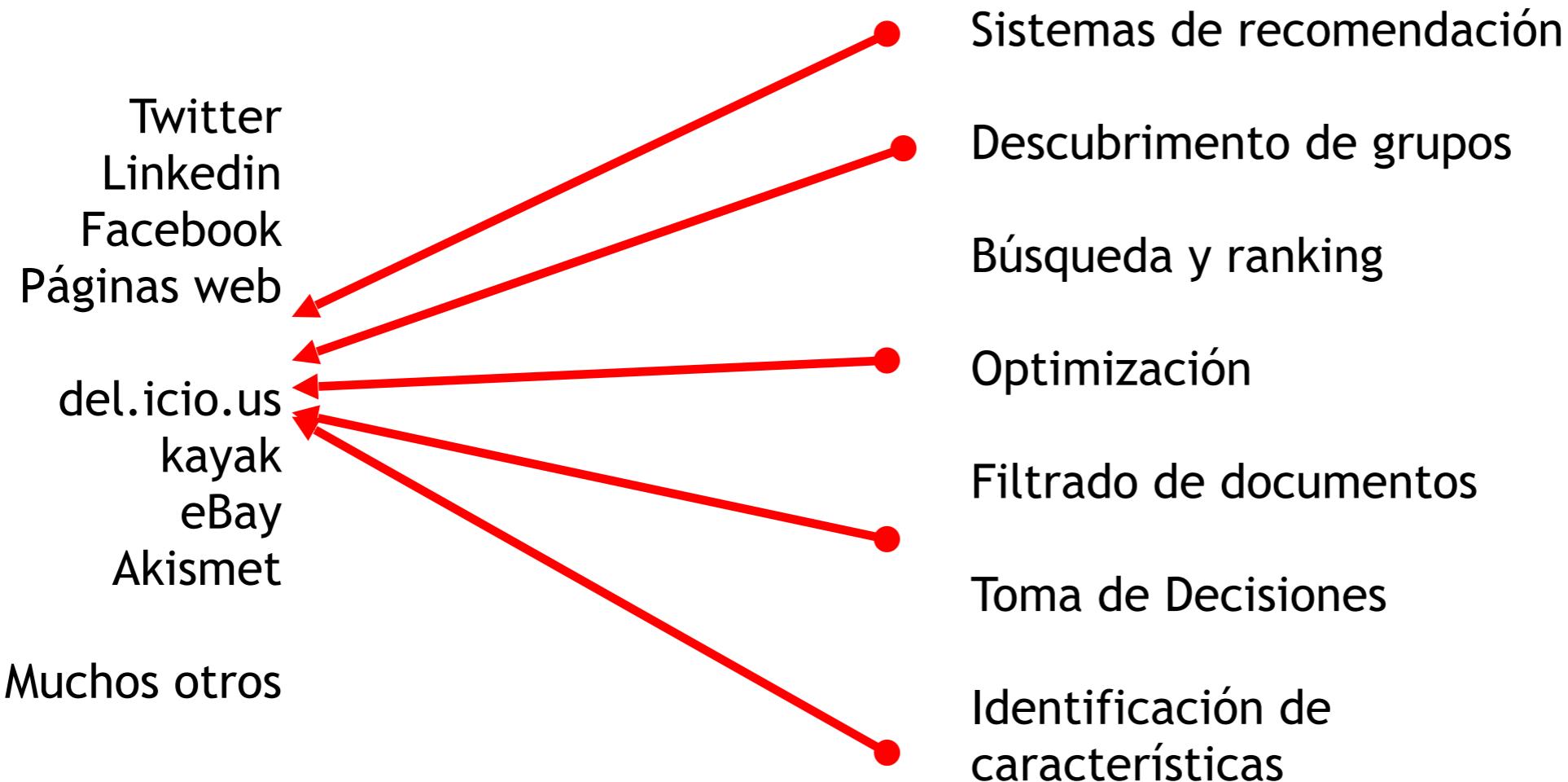
Imágenes (Captcha)

Manipulación de texto

Conversión de un formato a otro

Detección de datos faltantes, datos nulos, datos inconsistentes

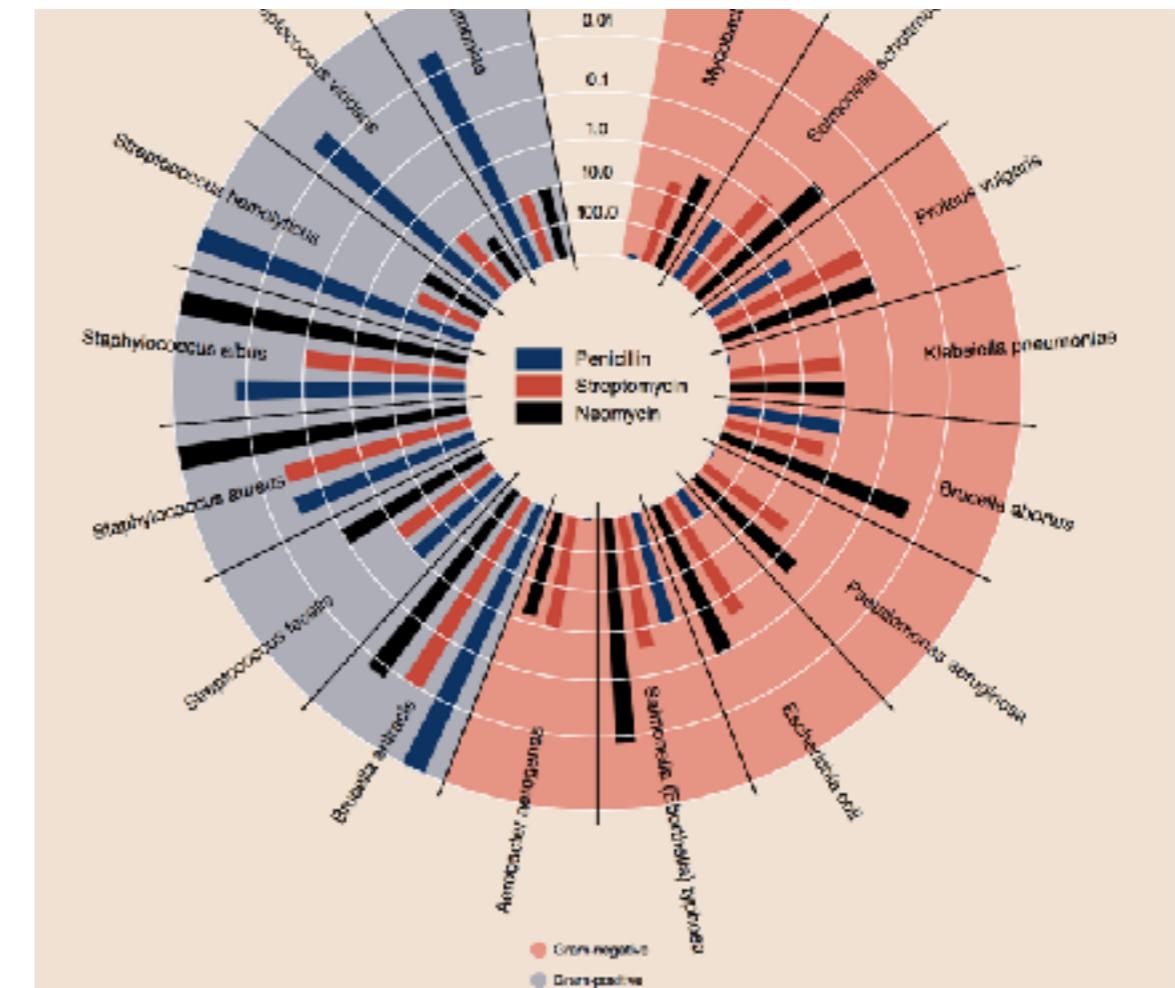
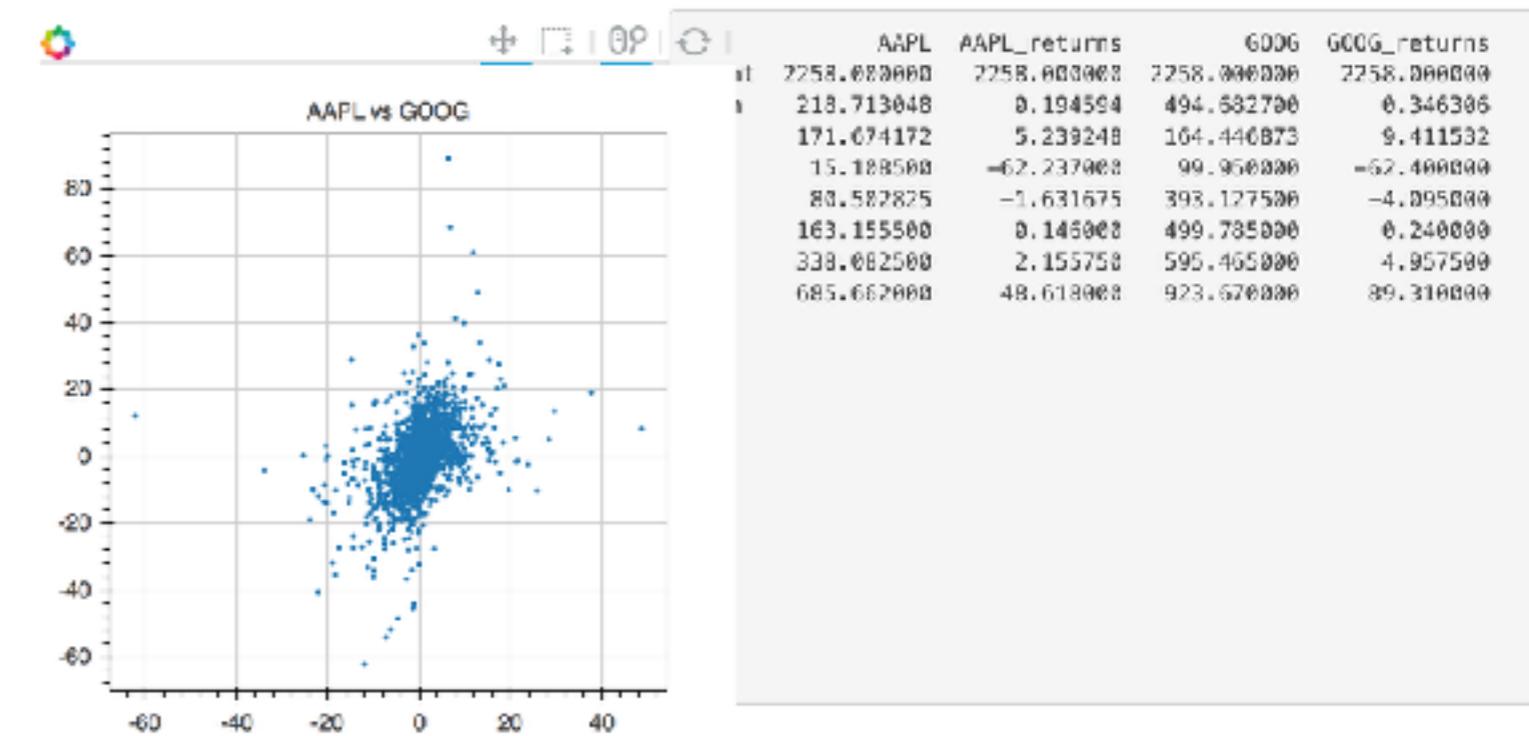
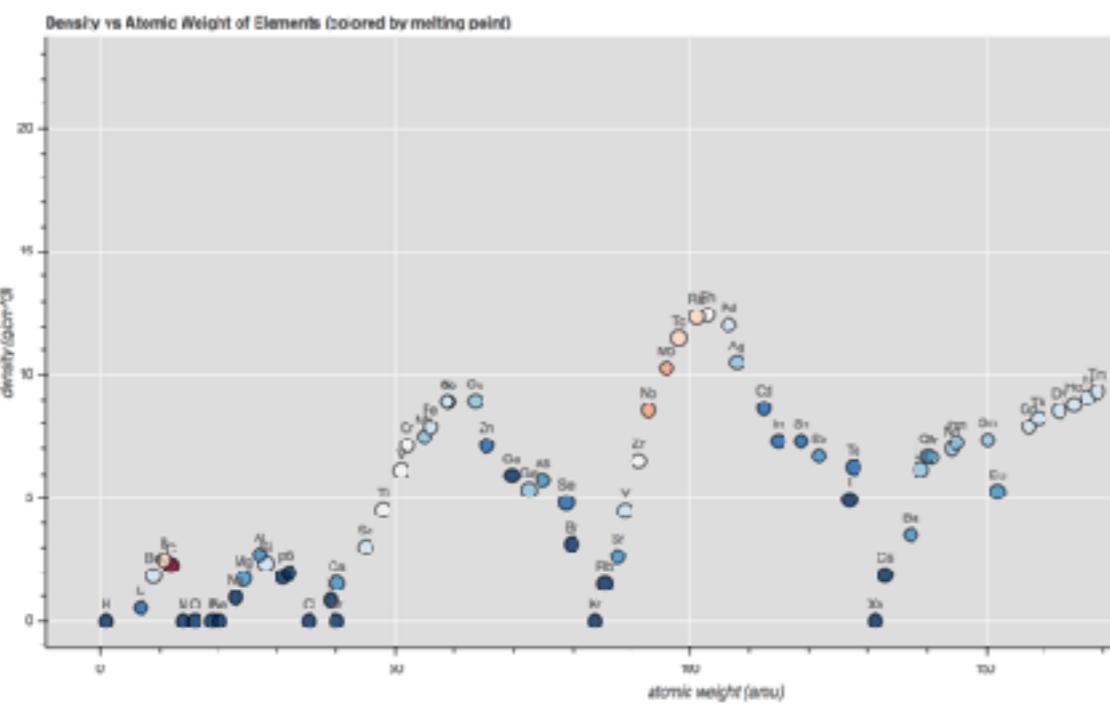
Adquisición de datos -- Inteligencia colectiva



[Link to this](#)

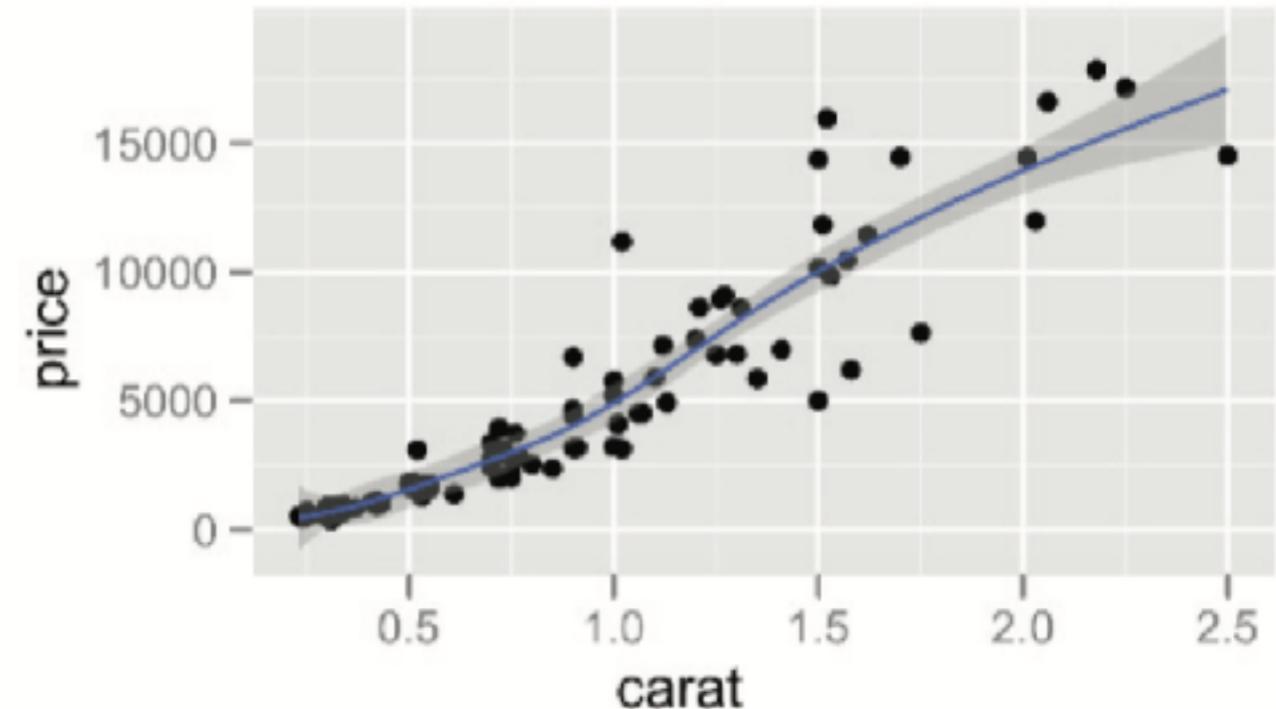
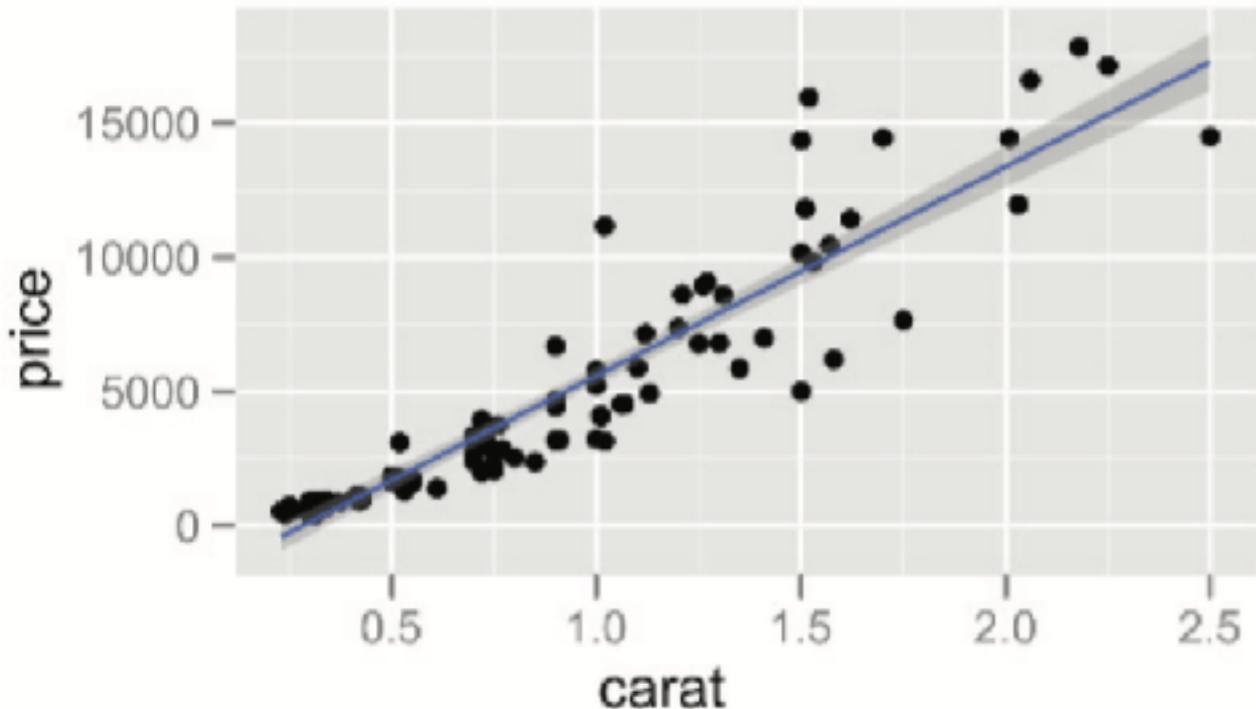
Visualización de Datos

- Bokeh
- Matplotlib
- R
-



Modelado Estadístico y Aprendizaje de Máquinas

Aplicación clásica



¿Y si hay 10 millones de datos?

Computación Reproducible (Markdown)

Markdown Editor

Input

```
Inline link: [destination](<index.html>)
Reference link: [destination][1]
Reference link: [reference link]

[1]: <index.html>
[reference link]: <http://www.infopark.com> "Link title"

Automatic link: <http://daringfireball.net/projects/markdown />

This is a blockquote
(pre + code)

-----
Heading 1
-----
Heading 2
-----
### Heading 3
#####
Heading 4
#####
#####
Heading 5
#####
#####
Heading 6
#####
* List item 1
* List item 2
  * Subitem 2.1
  * Subitem 2.2
```

Preview

```
Inline link: destination
Reference link: destination
Reference link: reference link
Automatic link:

This is a blockquote
(pre + code)
```

Heading 1

Heading 2

Heading 3

Heading 4

Heading 5

Heading 6

- List item 1
- List item 2
 - Subitem 2.1
 - Subitem 2.2

[Help on this page](#)

?

Ok

Cancel

Investigación Reproducible (Markdown + R)

The screenshot shows the RStudio interface with a document titled "knitr-ex1.Rmd". The toolbar includes a "Push here" button, which is highlighted with a blue arrow pointing from the text above. Below the toolbar, the "Knit HTML" button is also highlighted with a blue arrow. The code editor contains the following R Markdown code:

```
1 My First knitr Document
2 -----
3
4 This is some text (i.e. a "text chunk").
5
6 Here is a code chunk
7 ```{r}
8 set.seed(1)
9 x <- rnorm(100)
10 mean(x)
11 ```
```

The resulting HTML output is displayed in a preview pane:

My First knitr Document

This is some text (i.e. a "text chunk").

Here is a code chunk

```
set.seed(1)
x <- rnorm(100)
mean(x)
```

[1] 0.1089

Two blue callout boxes are present: one labeled "Code input" pointing to the code chunk, and another labeled "Numerical output" pointing to the resulting numerical value.

Computación Reproducible (Jupyter Notebook)

The screenshot shows a Jupyter Notebook window with the following content:

```
import scipy
import sys

# make nice plots
import plt_fmt

Populating the interactive namespace from numpy and matplotlib
```

"m" key denotes a markdown cell

```
In [8]: kk = rand(5,2)

(r1,r2) = kk[1]::
print (kk[1]::)
print (r1)
print (r2)

[ 0.20757795  0.01992547]
0.207577947999
0.019925471486
```

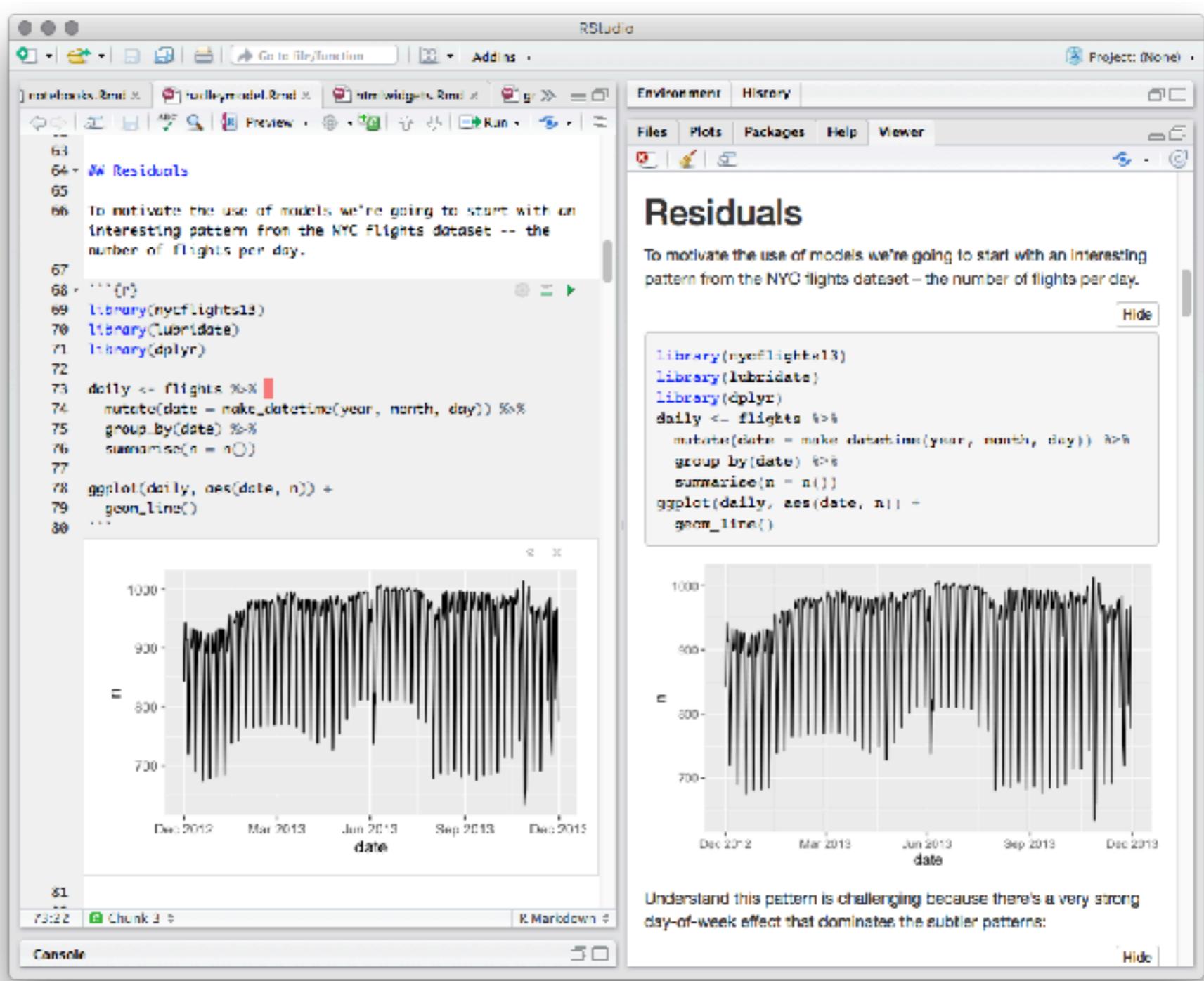
```
In [4]: def vfield(n,time, param):
    """
        param is an Nx2 matrix specifying the parameters for
        the dynamical system
    """

    (r1, r2) = param[0,:]
    (M1, M2) = param[1,:]
```

```
Out[4]: []
```

Productos de datos

- Informes autocalculables.
- Tableros de control (Dashboards)
- Aplicaciones de datos

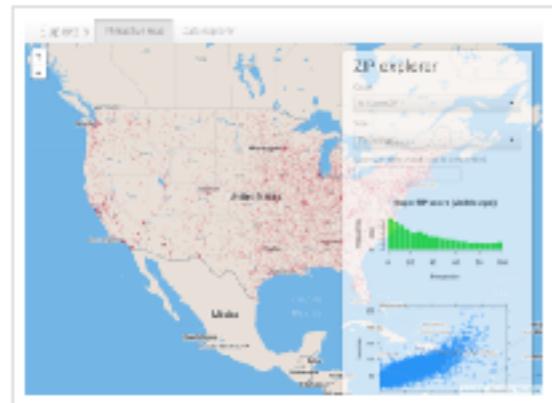


Gallery

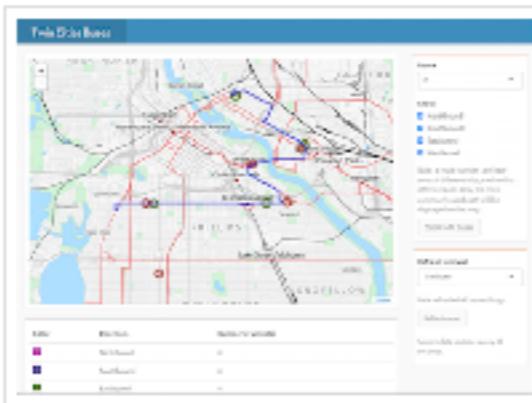
This gallery contains useful examples to learn from. Visit the [Shiny User Showcase](#) to see an inspiring set of sophisticated apps.

Interactive visualizations

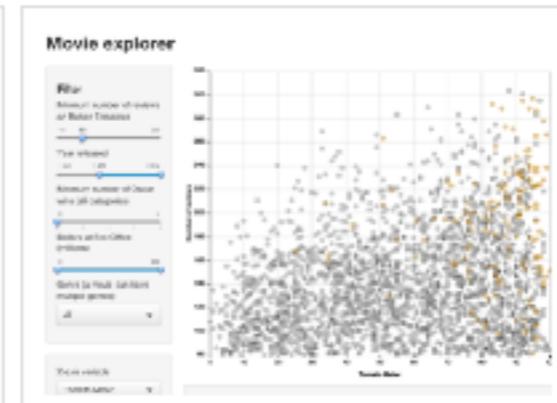
Shiny is designed for fully interactive visualization, using JavaScript libraries like [d3](#), [Leaflet](#), and [Google Charts](#).



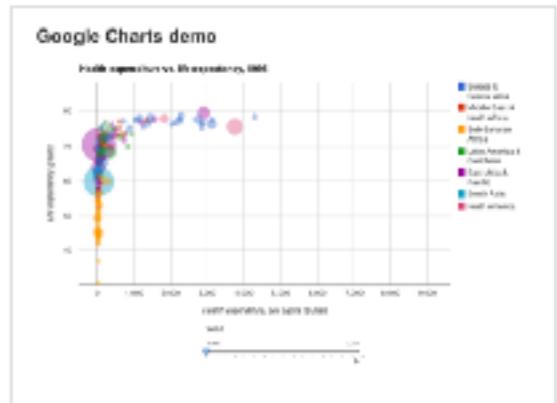
SuperZip example



Bus dashboard



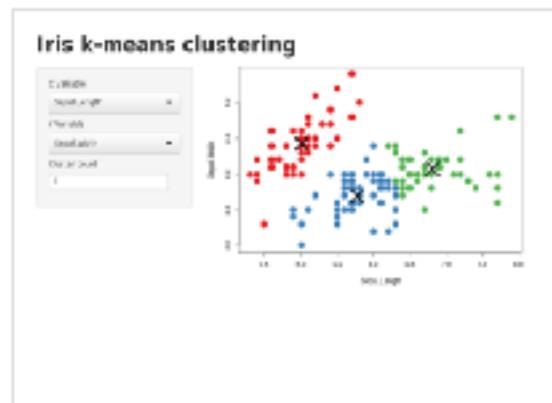
Movie explorer



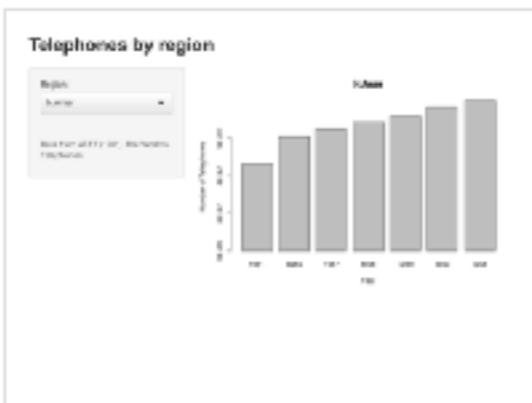
Google Charts

Start simple

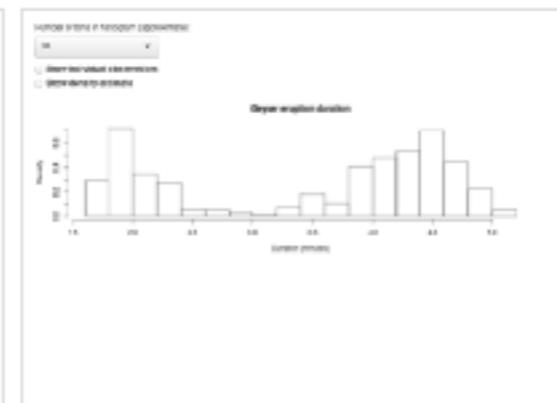
If you're new to Shiny, these simple but complete applications are designed for you to study.



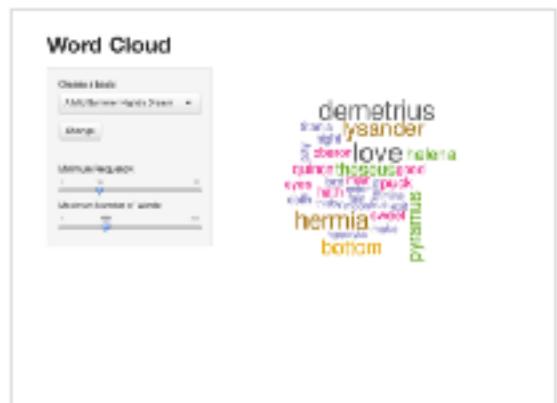
Kmeans example



Telephones by region

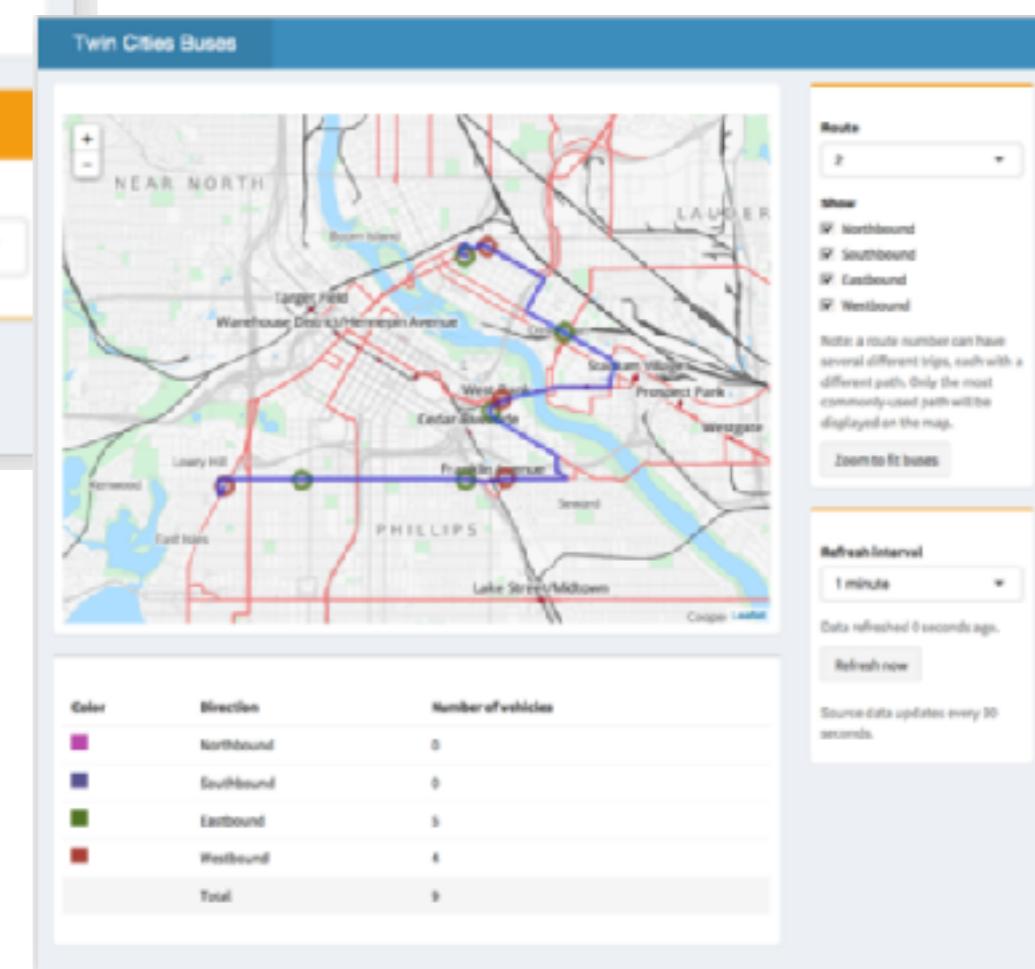
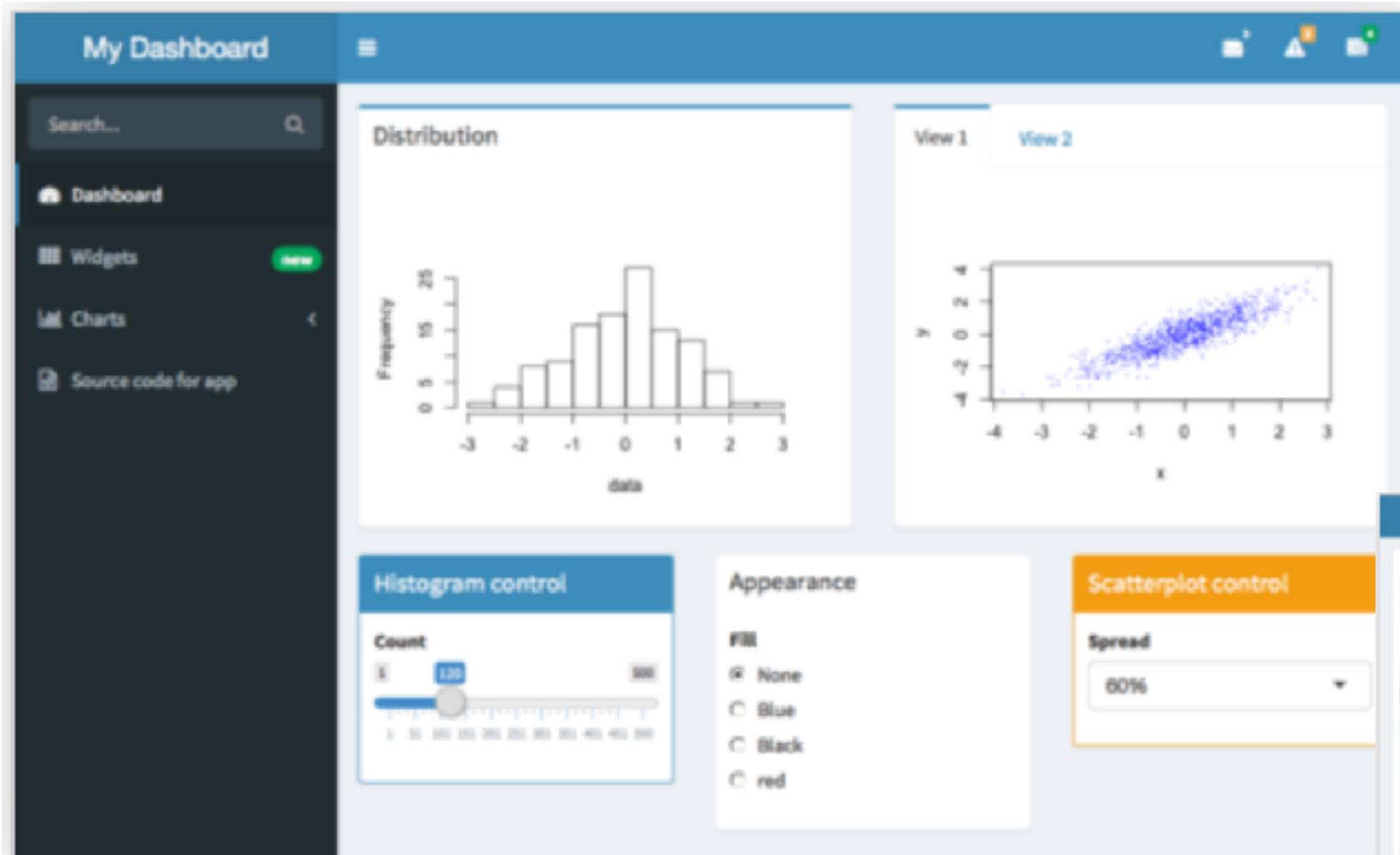


Faithful



Word cloud

R Dashboards



Hacia una visión unificada de Data Science, Analytics y Big Data

(con ejemplos en mercados de energía)

Esta presentación describe, bajo un marco común, los conceptos fundamentales de Data Science, Analytics y Big Data y establece su similitudes y diferencias.

Descargue la última versión de este documento de:
<https://github.com/jdvelasq/data-science-docs/blob/master/ds-analytics-bigdata.pdf>

JUAN DAVID VELÁSQUEZ HENAO, MSc, PhD

Profesor Titular

Departamento de Ciencias de la Computación y la Decisión

Facultad de Minas

Universidad Nacional de Colombia, Sede Medellín

 jdvelasq@unal.edu.co

 @jdvelasquezh

 <https://github.com/jdvelasq>

 <https://goo.gl/prkjAq>

 <https://goo.gl/vXH8jy>