

Hands-On Data Science with R

Desde el Procesamiento a la Minería de Datos

Camilo A Herrera R
18 de Junio de 2016

Extracción y Generación de Conocimiento a través de los Datos



*Camilo
Herrera*

Estadístico - Universidad del Valle

Sp. Data Science & Sp. Executive Data Science - Johns Hopkins Bloomberg

Msc candidate - Biometria - Universidad de Buenos Aires

twitter: @hr_mr_zork - web: <http://camilo herrera.co/>

Email: ch@camilo herrera.co

Cronograma

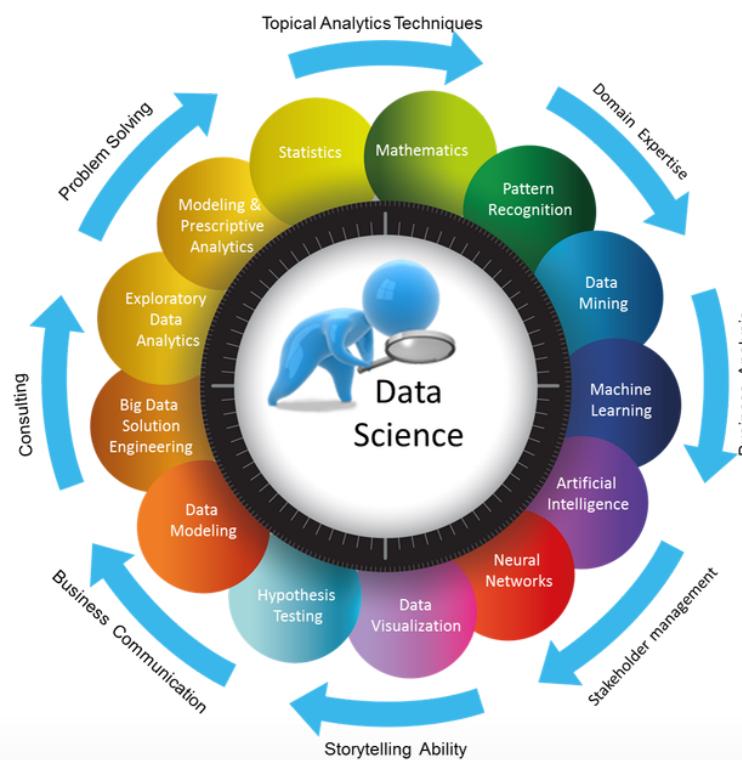
- Apertura del Evento
- Socialización R-Users MeetUp
- Sesión Introductoria Data Science con R
- Conferencias MeetUp:
 - Cleaning Data and Merging Data Sources with R -> Daniel Valencia
 - Visualizando datos con ggplot -> Maria Isabel Arce
- Taller en R (Hands-on)
- Cierre

Usa el Hashtag en twitter #RUsersCali

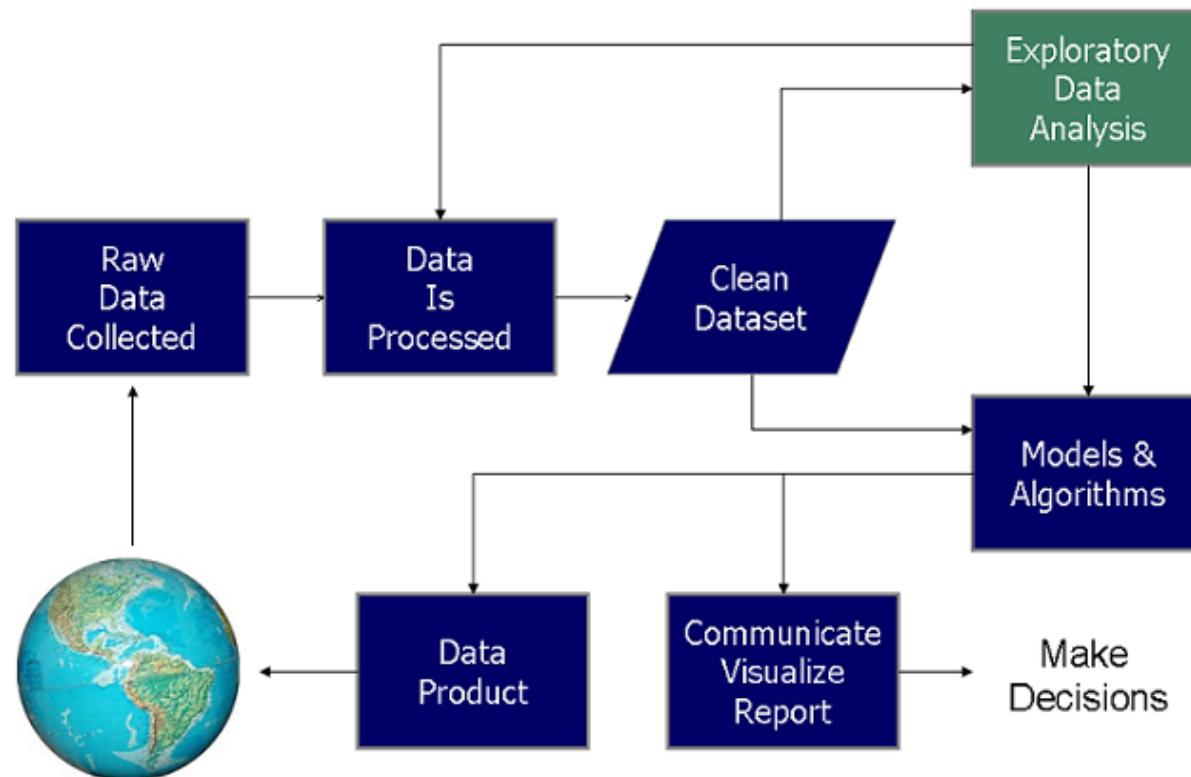
Ciencia de Datos

Que es Ciencia de Datos

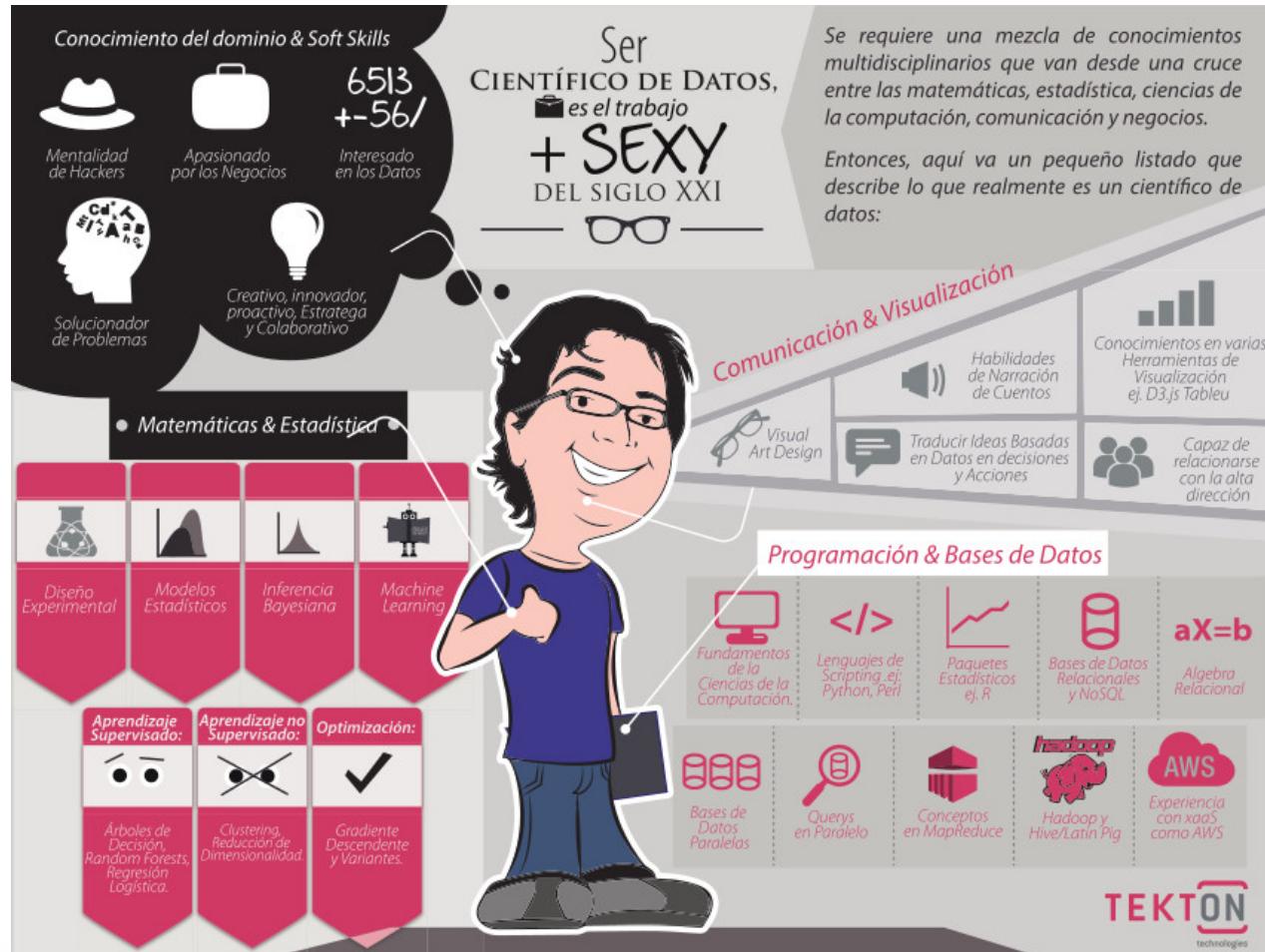
La Ciencia de datos es un campo interdisciplinario que involucra los procesos y sistemas para extraer conocimiento o un mejor entendimiento de grandes volúmenes de datos.



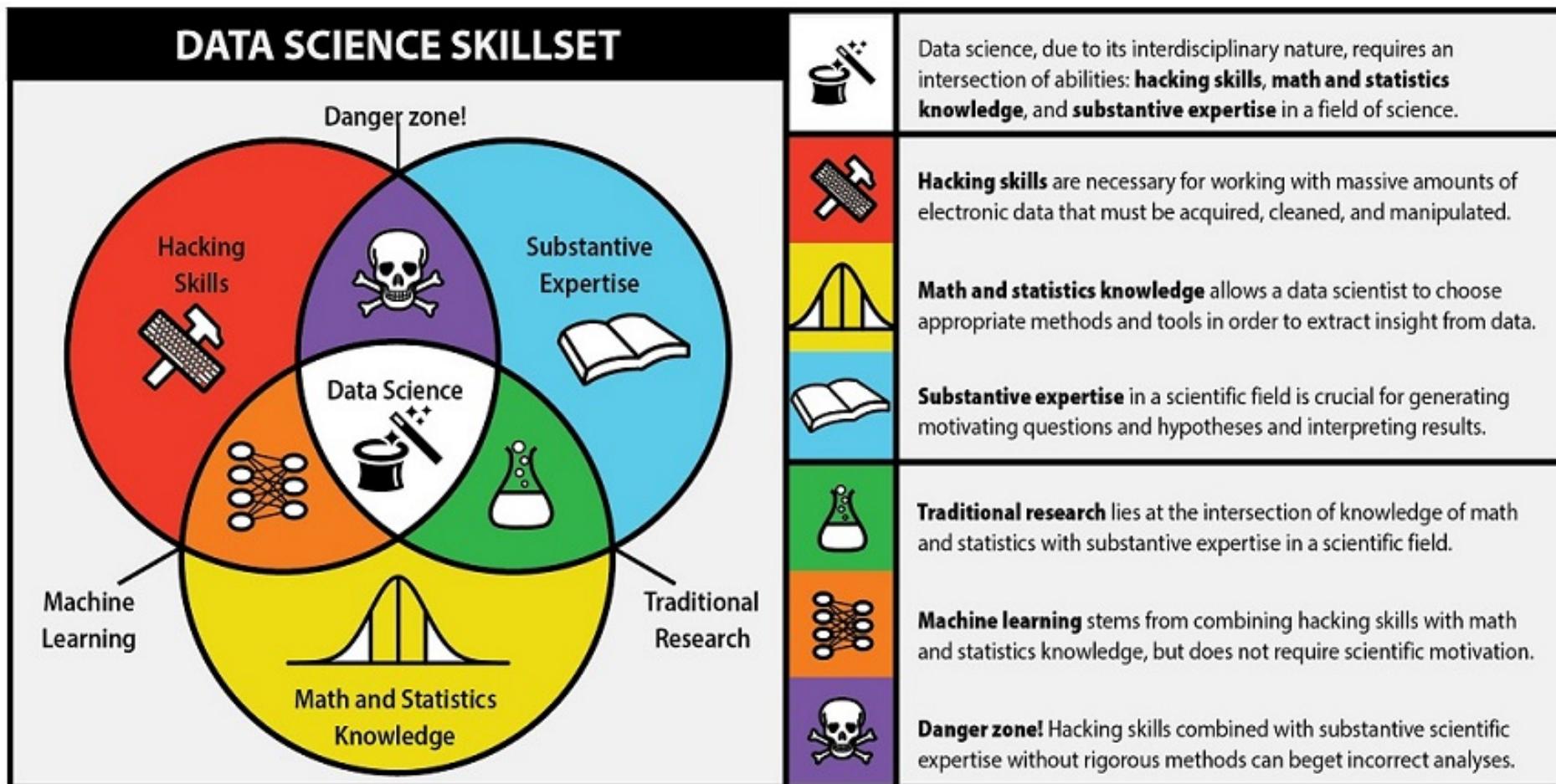
Data Science Process



El Científico de Datos



Habilidades de un Científico de Datos



Una Ruta para ser Científico de Datos



Equipos de Data Science

La ciencia de datos es un proceso que requiere un esfuerzo importante, por lo tanto se necesita de un grupo con un comportamiento equivalente a un equipo deportivo:



Equipos de Data Science

Un equipo de ciencia de datos está compuesto por:

1. Ingenieros de Datos
 - Arquitectura de Datos
 - Infraestructura de Datos
2. Data science
 - Limpieza de Datos
 - Análisis y Comunicación
3. Líder del Equipo de Datos

R

Lenguaje de Programación

¿Por qué R?

- R Es gratis
- Cuenta con un amplio conjunto de paquetes
- Acceso a los datos
- Limpieza de datos
- Análisis
- Generacion de Reportes
- Tiene uno de los mejores entornos de desarrollo - Rstudio
<http://www.rstudio.com/>
- Tiene un increíble ecosistema de desarrolladores
- Los paquetes son fáciles de instalar y "juegan muy bien juntos"

¿Por qué R?

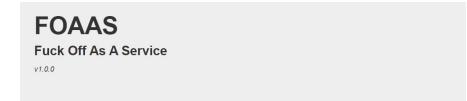


R es considerado la "lingua franca" de la ciencia de datos, por esta razón las empresas se están adaptando rápidamente a R para desarrollar sus programas y productos de "Data Science".

Paquetes en R

Paquetes R en CRAN

Existen paquetes para todo!!



Introduction

FOAAS (Fuck Off As A Service) provides a modern, RESTful, scalable solution to the common problem of telling people to fuck off.

API

Content Negotiation

FOAAS will respond to the following "Accept" values with appropriate content:

- text/plain - Content will be returned as a plain string.
- application/json - Content will be returned as a JSON object ({ "message": "message", "content": "content" }).
- Supports jsonp by enclosing .callback()

```
R> library(rfoaas)
R> greed("R Programming", "Random R Hacker")
The point is, ladies and gentleman, that r programming -- for lack of a
better word -- is good. R Programming is right. R Programming works. R
Programming clarifies, cuts through, and captures the essence of the
evolutionary spirit. R Programming, in all of its forms -- R
Programming for life, for money, for love, knowledge -- has marked the
upward surge of mankind. - Random R Hacker
R> █
```

Visualizador

Plataformas de Analítica y Bigdata que apuestan por R

SQL

SQL Server 2016

Industry leadership in
Mission Critical OLTP

Operational DBMS

Business Intelligence

Data Warehouse

Advanced Analytics

Most secure database
6 years in a row

Year	SQL Server	Oracle	MySQL
2010	35	30	30
2011	30	30	30
2012	30	30	30
2013	30	30	30
2014	30	30	30
2015	30	30	30

Highest performing
data warehouse

System	Performance Rank	Score
SQL Server	1	~100
Oracle	2	~80
SQL Server	3	~60
Oracle	4	~40
SQL Server	5	~20
Oracle	6	~10

End-to-end mobile BI
on any device

System	Cost
Microsoft	\$120
Tableau	\$480
Oracle	\$2,230

In-database
Advanced Analytics

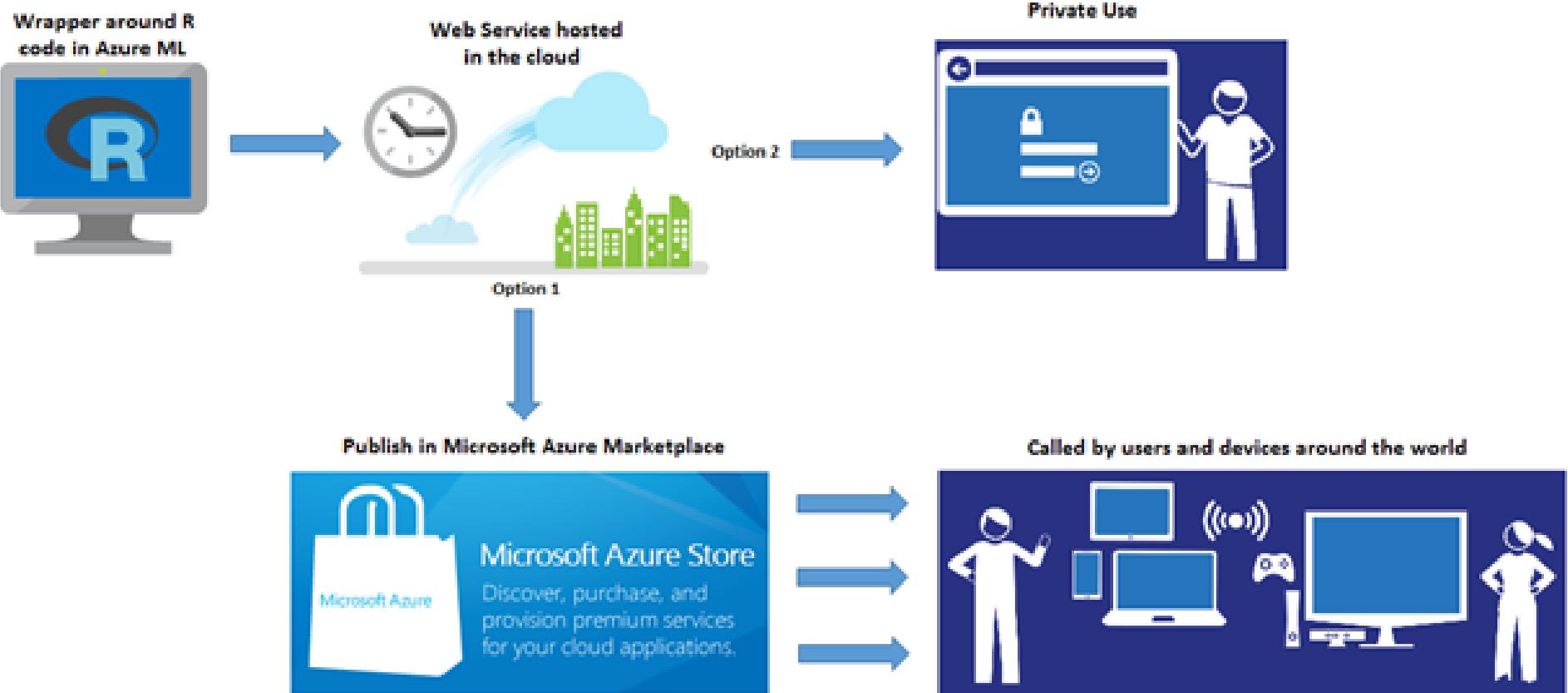
R + in-memory at massive scale

On Premises

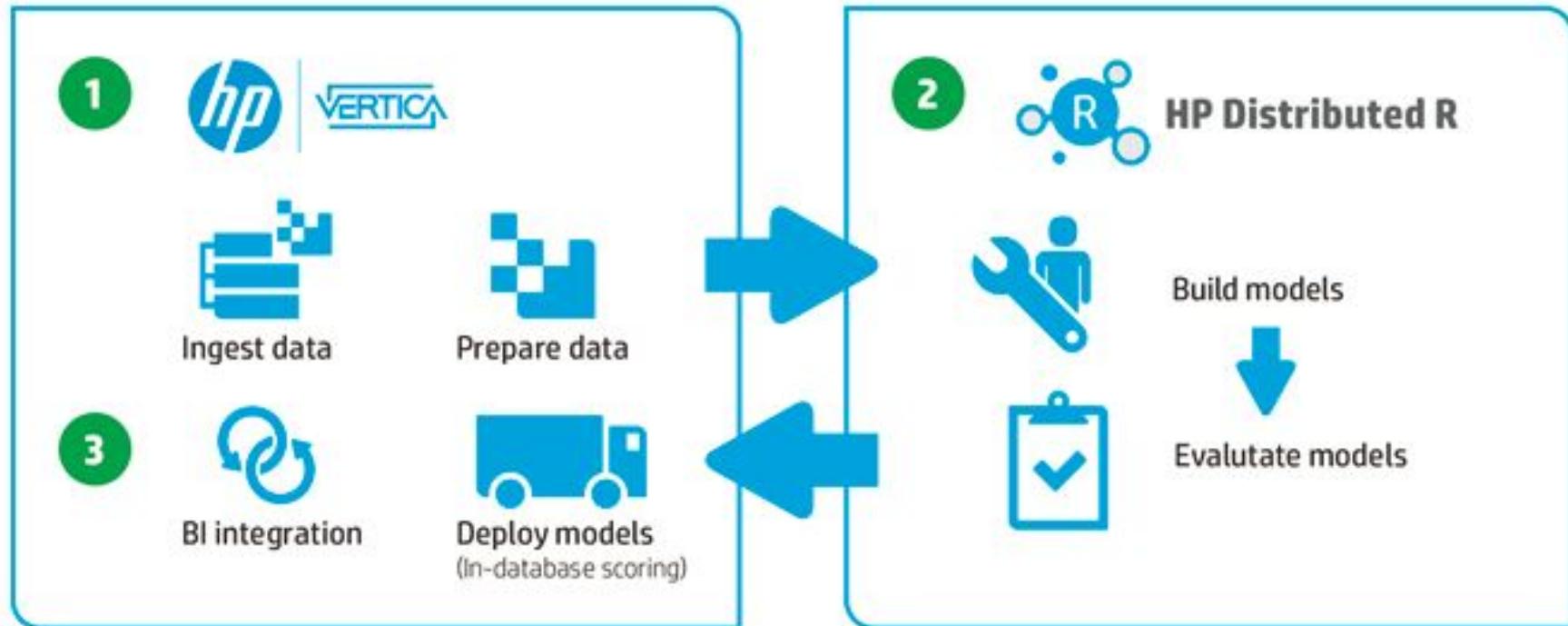
Cloud

18/45

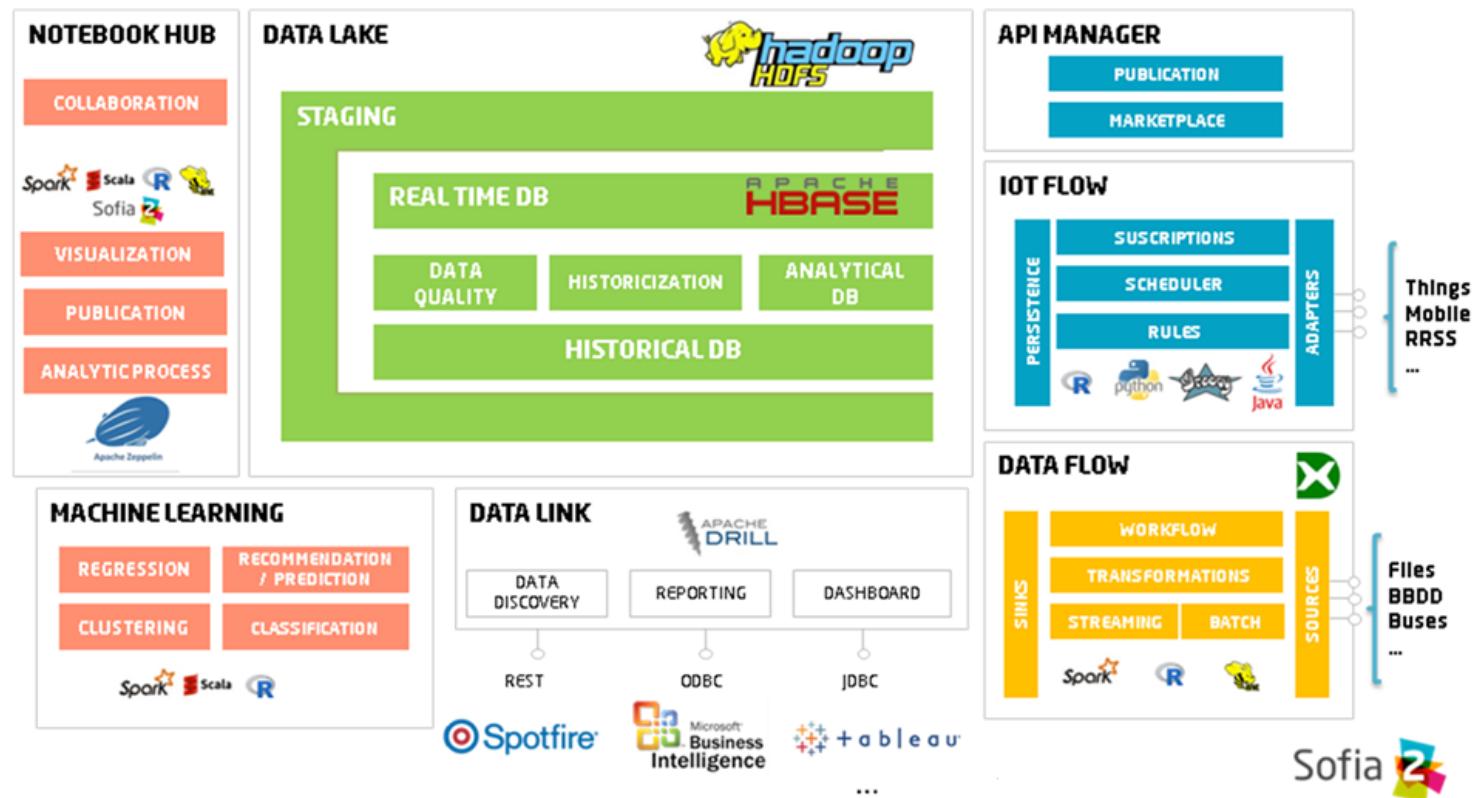
Azure



HP



Sofia2 - Indra



La Apuesta de las Industrias por R

R Consortium

R Foundation Members



Platinum Members



Gold Members



Silver Members



La misión explícita del R Consortium es "avanzar en la promoción mundial y el apoyo para el lenguaje de código abierto R"

Big Data

Guía del Viajero Intergaláctico The Hitchhiker's Guide to the Galaxy

Una de las Historias cuenta, que una raza de **seres hiperinteligentes pandimensionales** construyeron una computadora llamada Pensamiento Profundo («Deep Thought») fabricada con el único objetivo de descifrar la respuesta definitiva.

«el sentido de la vida, el universo y todo lo demás»



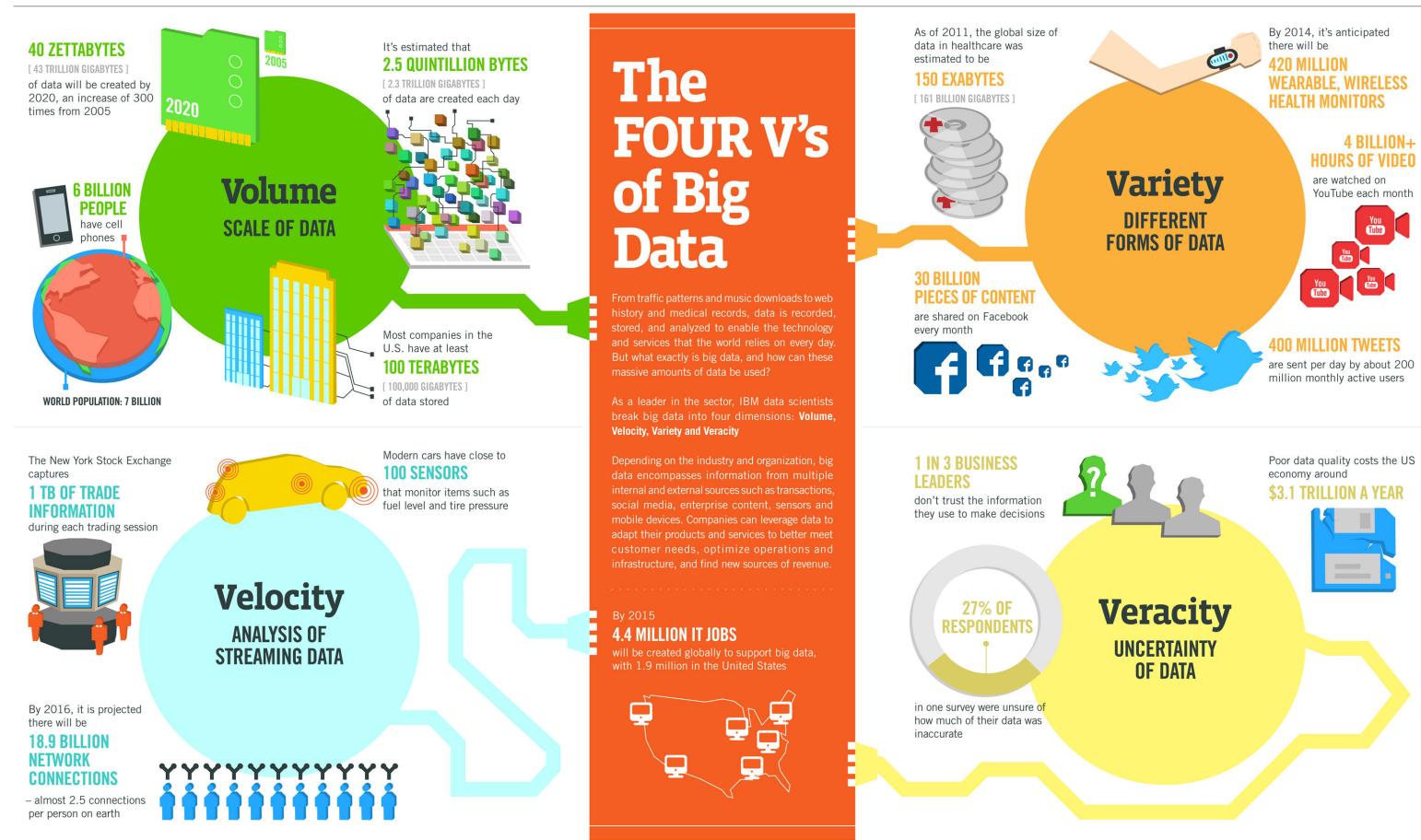
Guía del Viajero Intergaláctico **The Hitchhiker's Guide to the Galaxy**

Pensamiento Profundo se toma **siete millones y medio de años** para dar esa respuesta, la cual, para pesar de muchos, resulta ser **42** sin lugar a dudas.



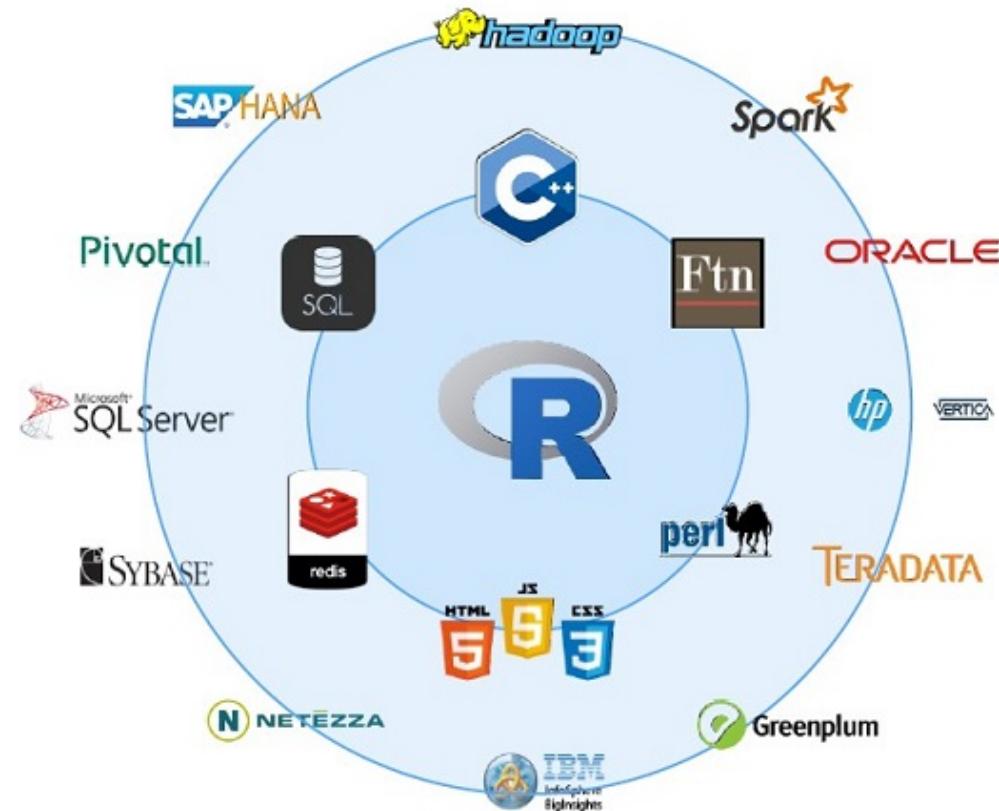
Cuando la respuesta se revela como 42, se ven forzados a construir una computadora aún más poderosa para calcular la «Pregunta máxima», pero sus planes nunca culminan...

Big Data

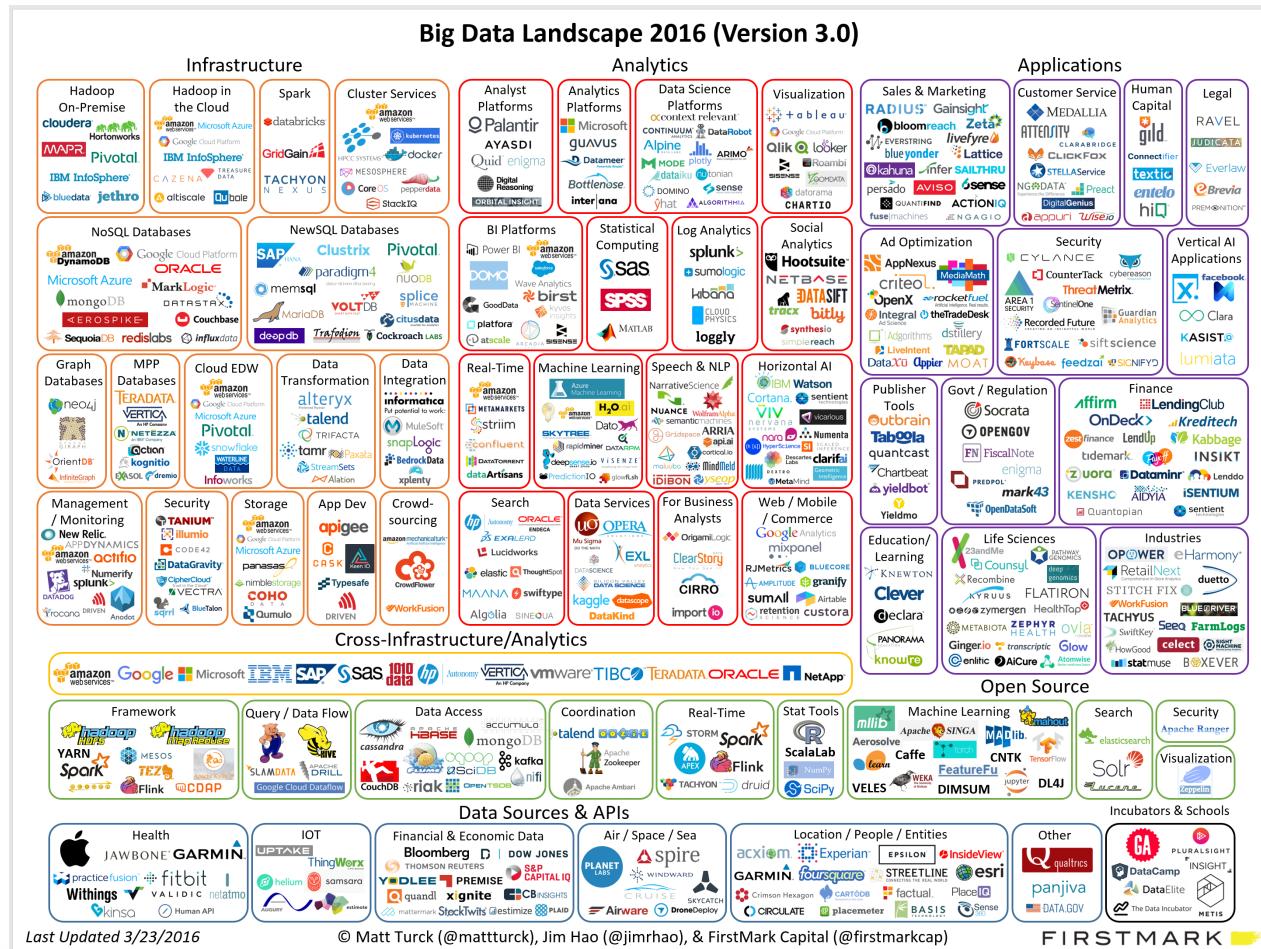


IBM

R y Big Data

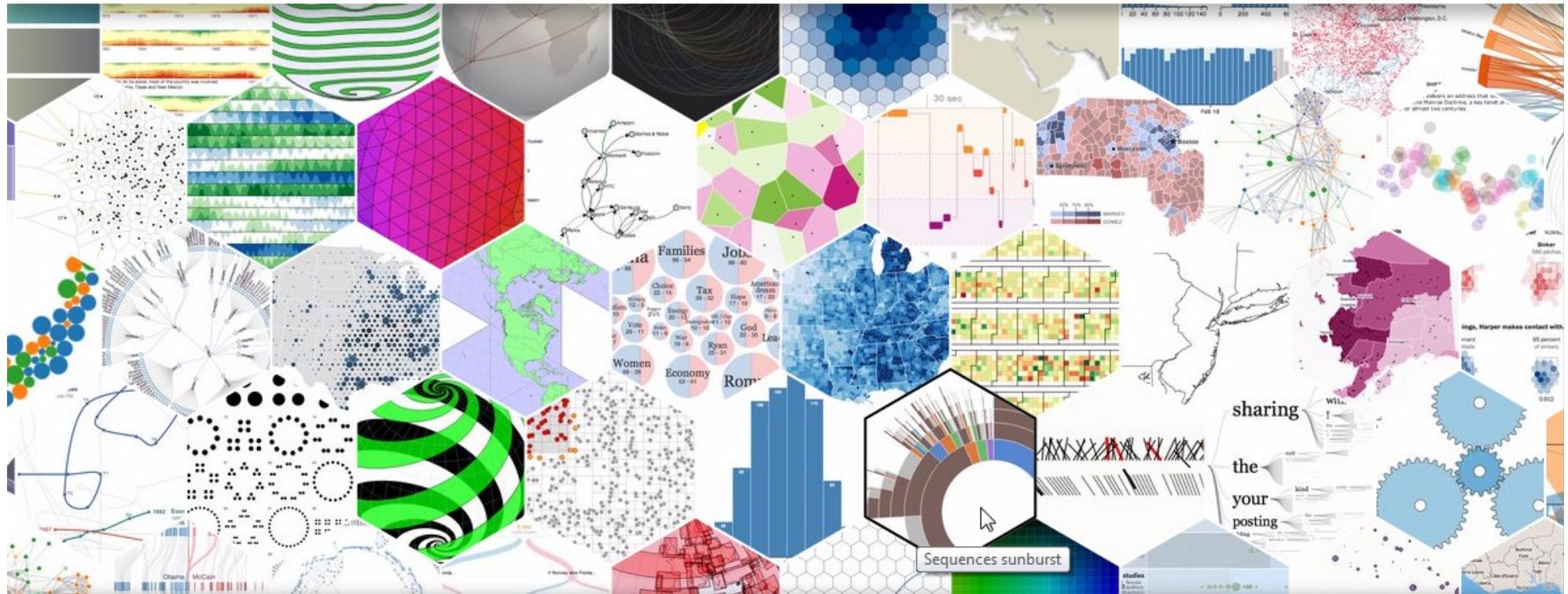


Panorama del Big Data 2016



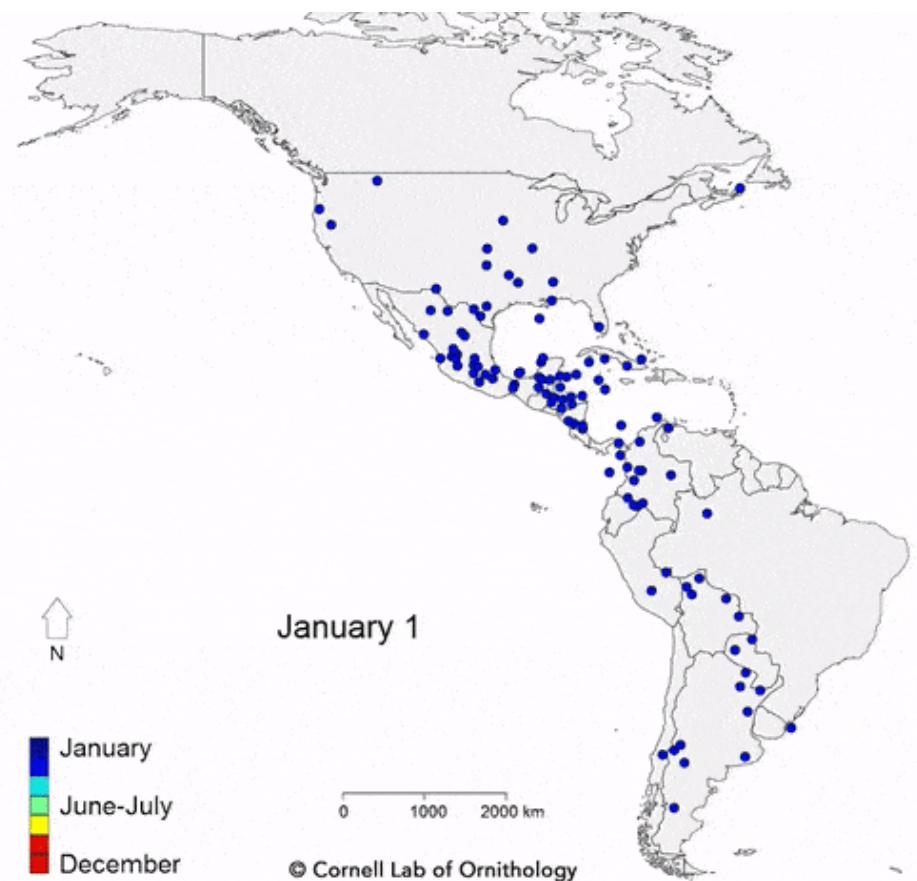
Visualización y Productos de Datos

D3.js



Visualización y Productos de Datos

Generar nuestras propias visualizaciones.



Visualización y Productos de Datos

Generar nuestras propias visualizaciones.

- Gráficos varios dentro de Dashboards
- Gráficos de Redes (nodos)
- Get data y plot (Ejemplo MIO)
- Publica un Tweet con el hashtag: #RUsersCali
- Gráficos y Social tracking (Partido)
- Ejemplo Twitter
- Ejemplo linet (Descargas eléctricas)
- Data Mining - "What if"

Rstudio y el Hadleyverse



Rstudio y el Hadleyverse

Screenshot of Hadley Wickham's GitHub profile page.

Popular repositories:

- ggplot2: An implementation of the Grammar of Graphing... 1,878 stars
- dplyr: Dplyr: A grammar of data manipulation 1,277 stars
- devtools: Tools to make an R developer's life easier 1,276 stars
- adv-r: Advanced R programming: a book 771 stars
- rvest: Simple web scraping for R 584 stars

Repositories contributed to:

- klutometis/roxygen 175 stars
- wesm/feather 862 stars
- davidgohel/gdtools 9 stars
- ggobi/tourr-gui 2 stars
- garrettgman/ggsubplot 61 stars

Contribution activity: 5,766 contributions in the last year

Summary of pull requests, issues opened, and commits. Learn how we count contributions.

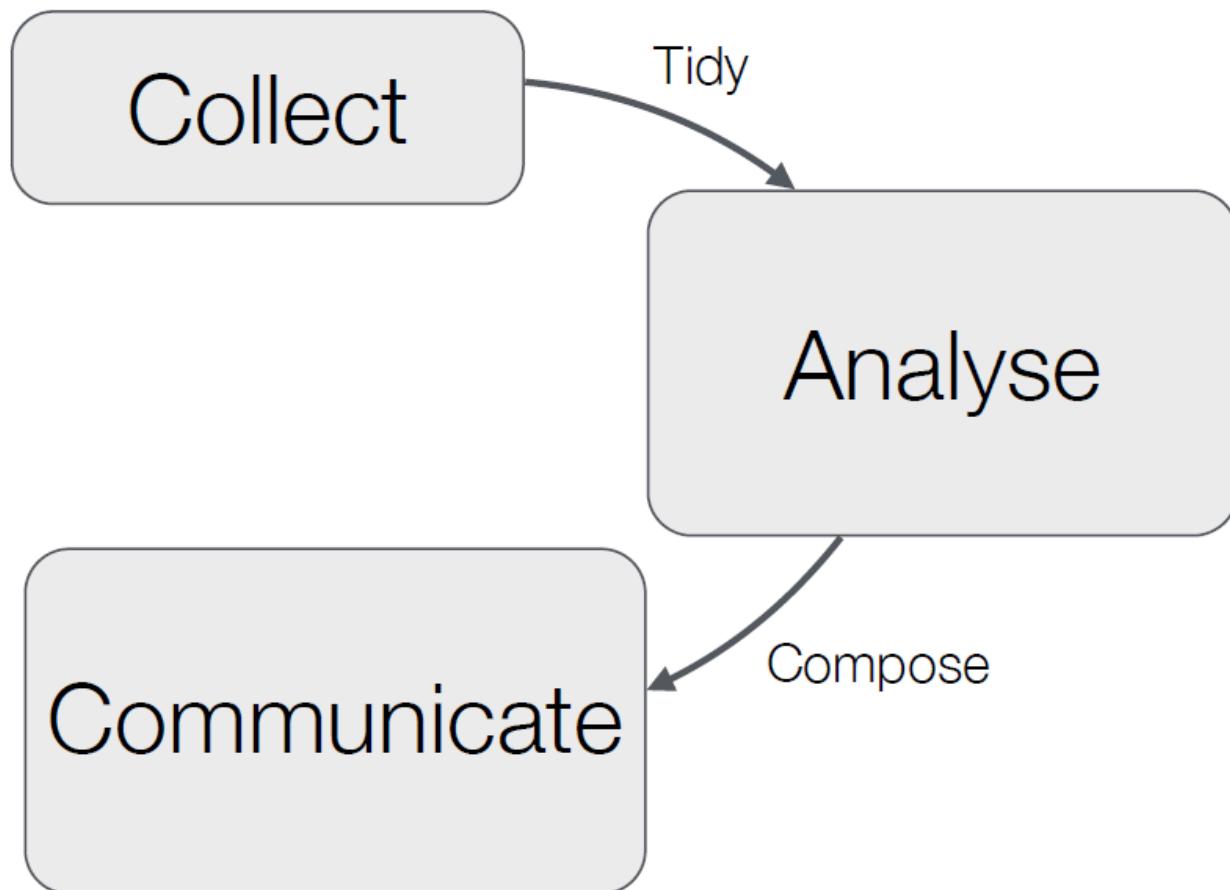
Less More

Period: 1 week

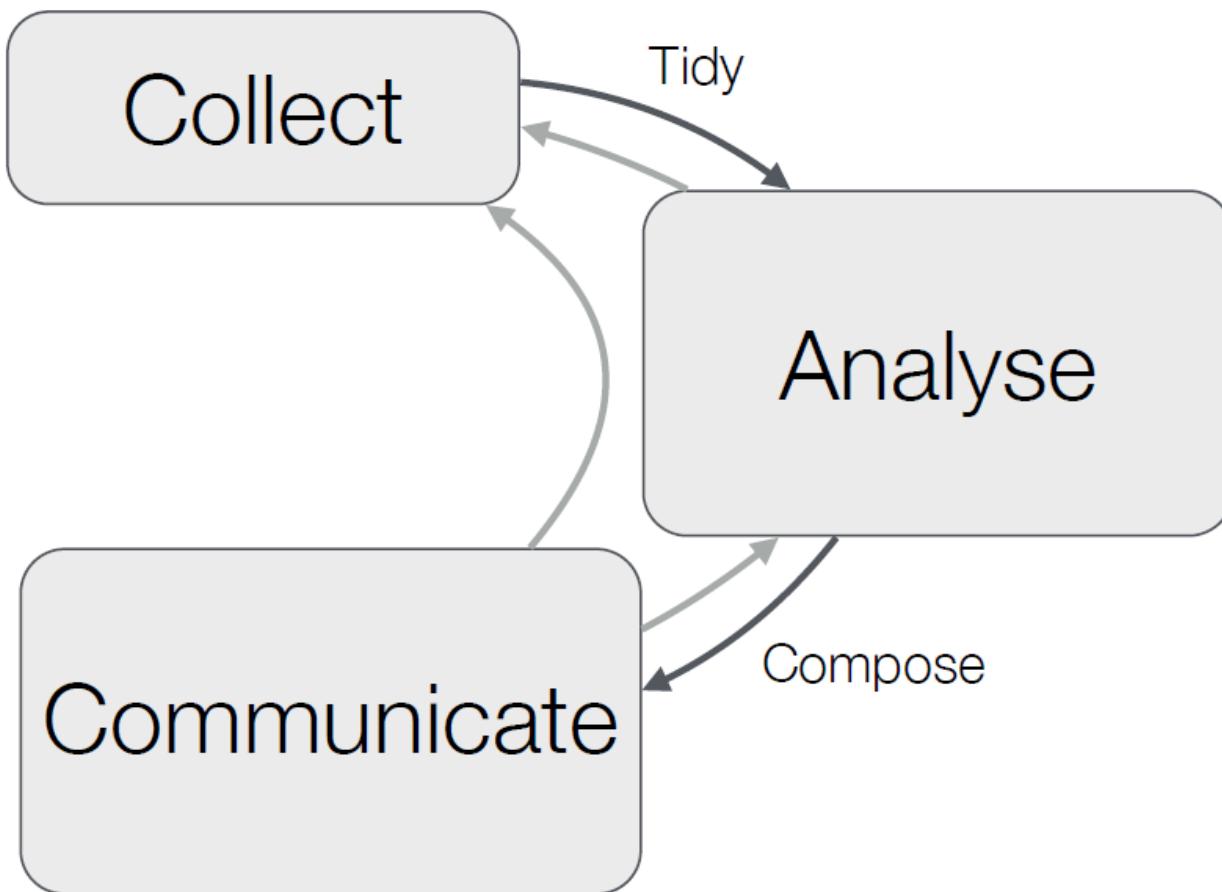
Contribution activity details:

- Pushed 13 commits to hadley/modelr Jun 13 – Jun 15
- Pushed 7 commits to hadley/dplyr Jun 9 – Jun 15
- Pushed 9 commits to hadley/httr Jun 9 – Jun 14

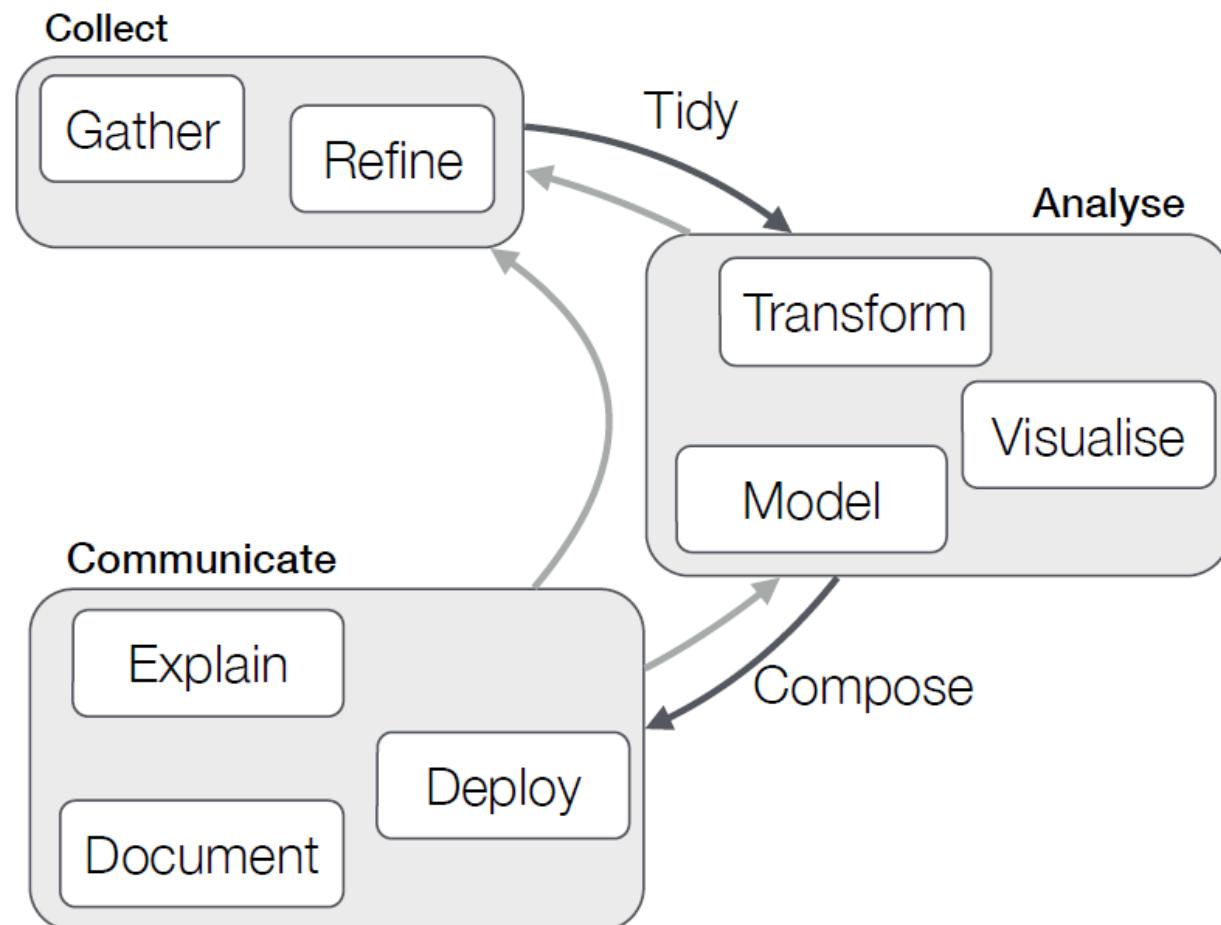
Procesos con Datos



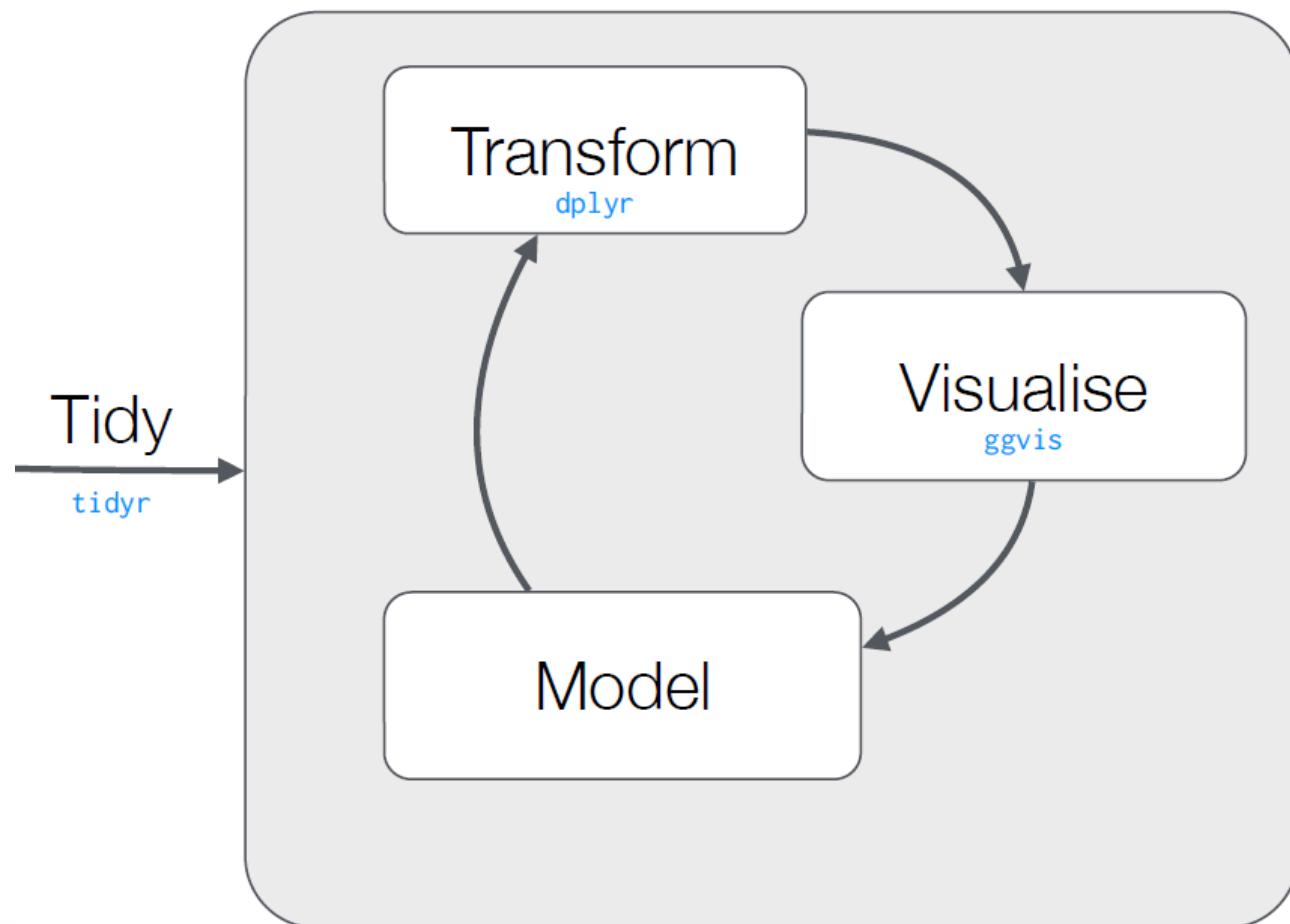
Procesos con Datos



Procesos con Datos



Procesos con Datos



dplyr

- **filter**: keep rows matching criteria
- **select**: pick columns by name
- **arrange**: reorder rows
- **mutate**: add new variables
- **summarise**: reduce variables to values

+ group by

Pipe

%>%

tidyr, dplyr, ggvis, ...

Apps Shiny

Typical Shiny App

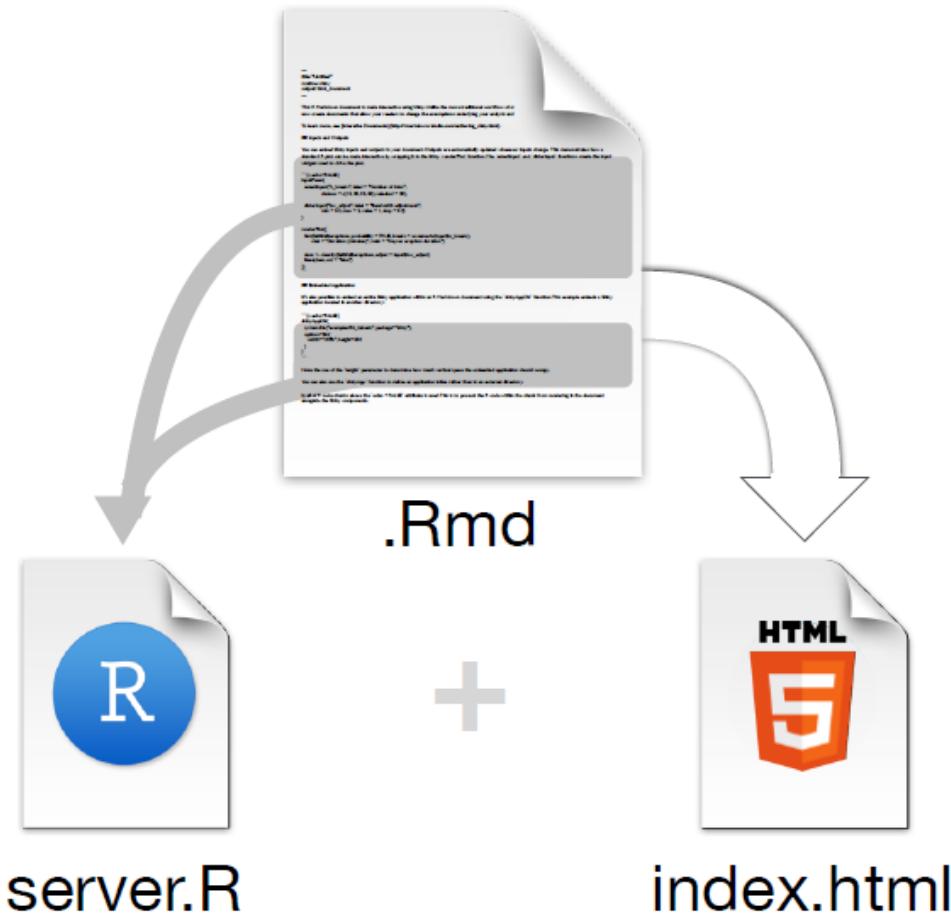


server.R



ui.R

Apps Shiny



Investigación Reproducible

R Markdown



Analytic
Power



LATEX



Microsoft Word



Reveal.js
ioslides, Beamer



Report generation

Fin Sesión Introductoria

