

DATA ANALYTICS Y DATA SCIENCE

Una introducción corta

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>

Data Science vs Analytics

Data Science (¿Data Analytics?):

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

Data Mining

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

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

Perspectiva

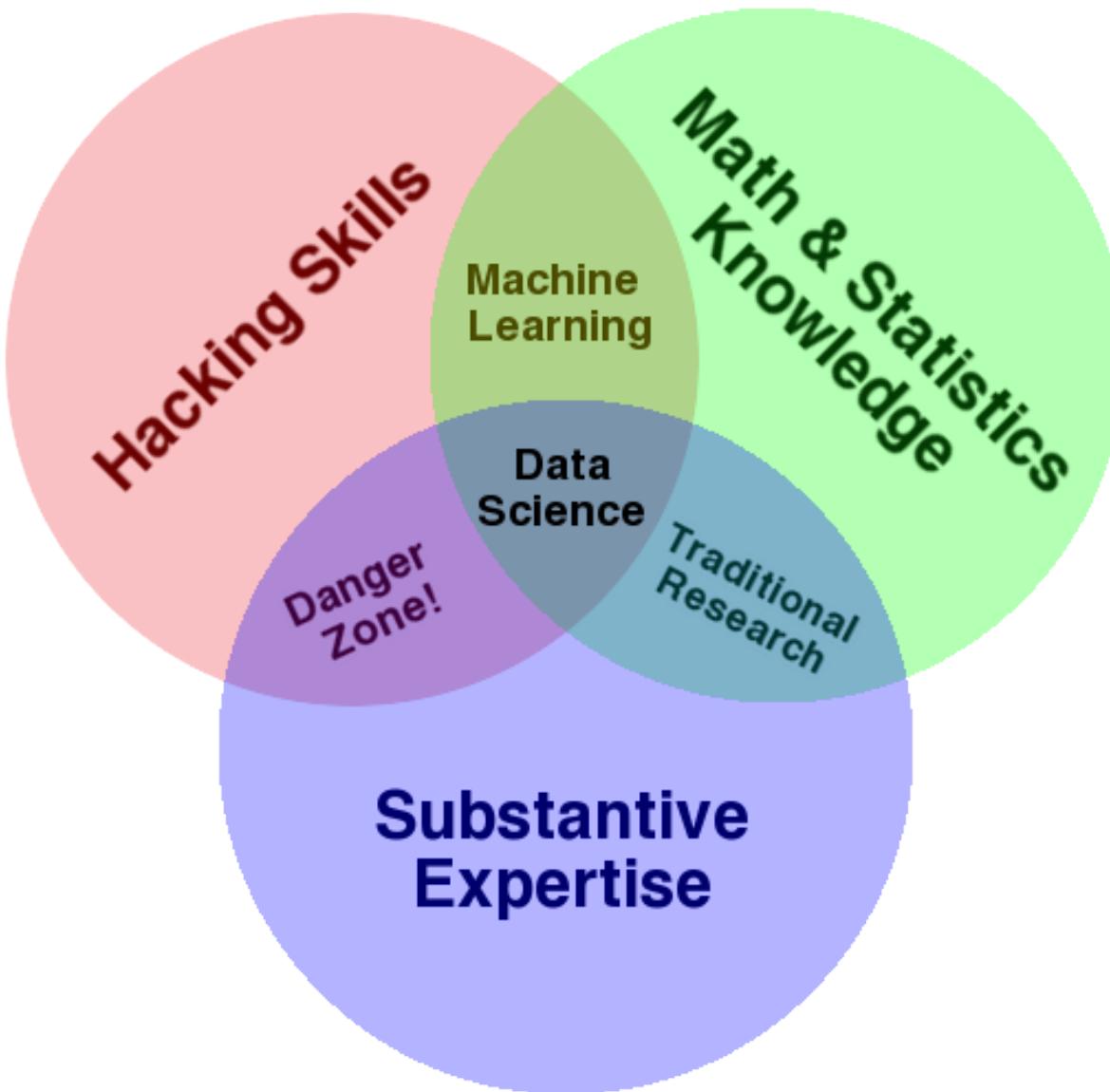
A recent study by the McKinsey Global Institute concludes, "a shortage of the analytical and managerial talent necessary to make the most of Big Data is a significant and pressing challenge (for the U.S.)." The report estimates that there will be four to five million jobs in the U.S. requiring data analysis 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."

#15	3,433	\$105,395	#1
-----	-------	-----------	----

Highest Paying Job in Demand	Number of Job Openings	Average Base Salary	Best Job in America for 2016
------------------------------	------------------------	---------------------	------------------------------

Sources: <http://www.glassdoor.com/blog/jobs-america/> and <http://www.glassdoor.com/blog/highest-paying-jobs-demand/>

Diagrama de Ven explicando qué es Data Science



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

Saunders, 2013

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

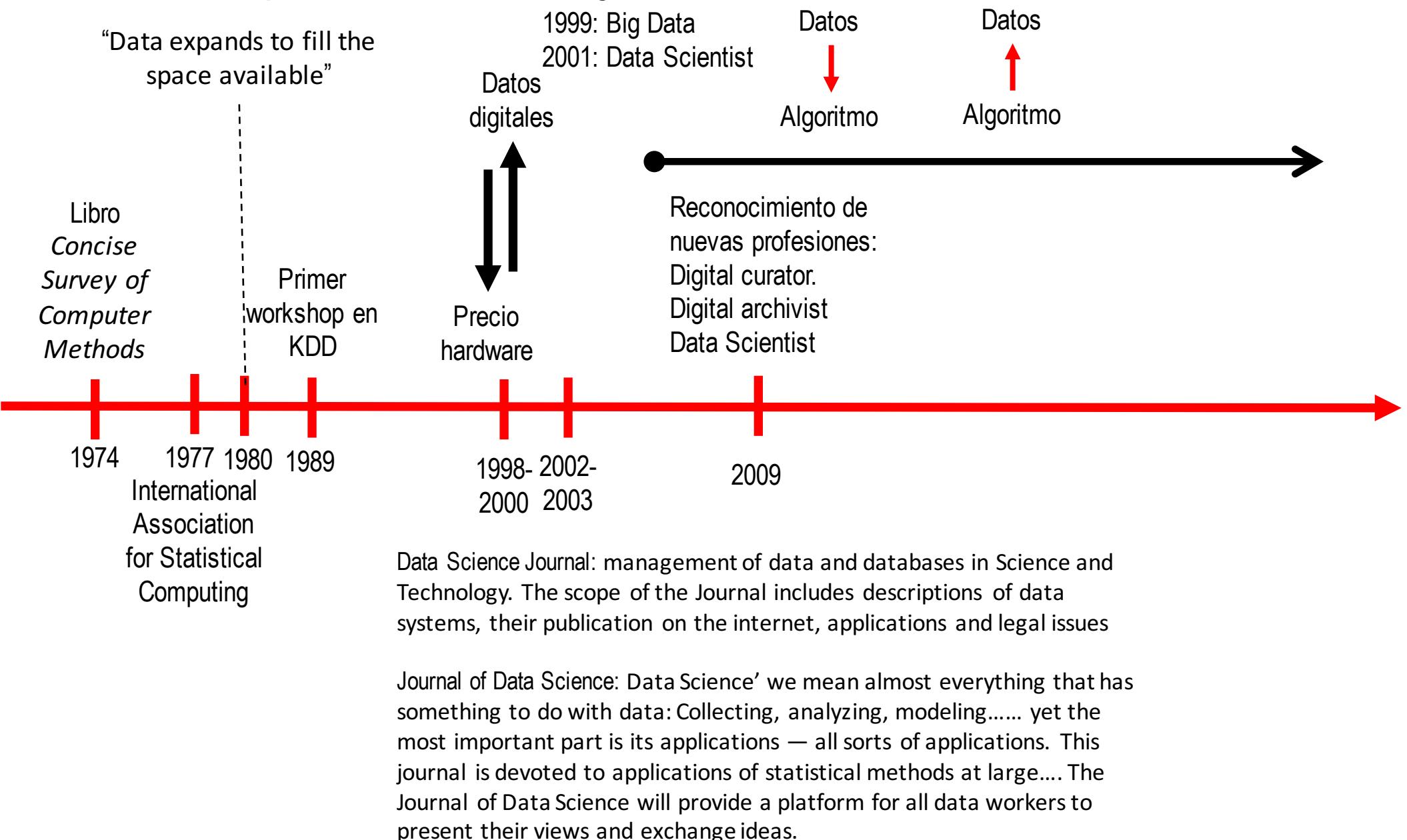
Habilidades en la programación de computadores para manejo de datos y modelado predictivo (estadística, aprendizaje de máquinas, minería de datos, etc.).

Analytics

Data Science + Optimización + Simulación

DISCIPLINE	TECHNOLOGIES	SKILLS	FOCUS
BUSINESS INTELLIGENCE	<ul style="list-style-type: none"> • ETL Tools / SQL • RDBMS • Reporting • Visualization 	<ul style="list-style-type: none"> • Programming • Data Analysis • Data Modeling • Report Development • Basic Statistics • Technical Architecture • Business Analysis & Strategy • Presentation 	<ul style="list-style-type: none"> • Information Delivery and Reporting • Data Visualization • Descriptive Statistics • Data Integration and Consolidation
DATA ANALYSIS	<ul style="list-style-type: none"> • Data Modeling Software • Diagramming Software • Documentation Software • SQL • Data Profiling Software 	<ul style="list-style-type: none"> • Data Modeling • Business Analysis • Data Manipulation • Basic statistics 	<ul style="list-style-type: none"> • Business Rules • Data Definitions and Lineage • Data Entity Relationships • Data Attributes • Data Structures • Sources and Targets of Data • Data Quality
DATA SCIENCE	<ul style="list-style-type: none"> • Statistics Software • Columnar Data • Map-Reduce • NoSQL • Programming Languages • Graphing/Charting Software 	<ul style="list-style-type: none"> • Advanced Statistics • Programming • Business Analysis • Modern Data Management Technologies and Architectures 	<ul style="list-style-type: none"> • Predictive Modeling • Advanced Statistical Analysis • Data Mining • Unstructured Data Management • Large data volumes • Research

Línea de tiempo. Data Science / Big Data



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.

Simulacion.

Optimización.

Métodos prescriptivos: modelos predictivos + optimización.

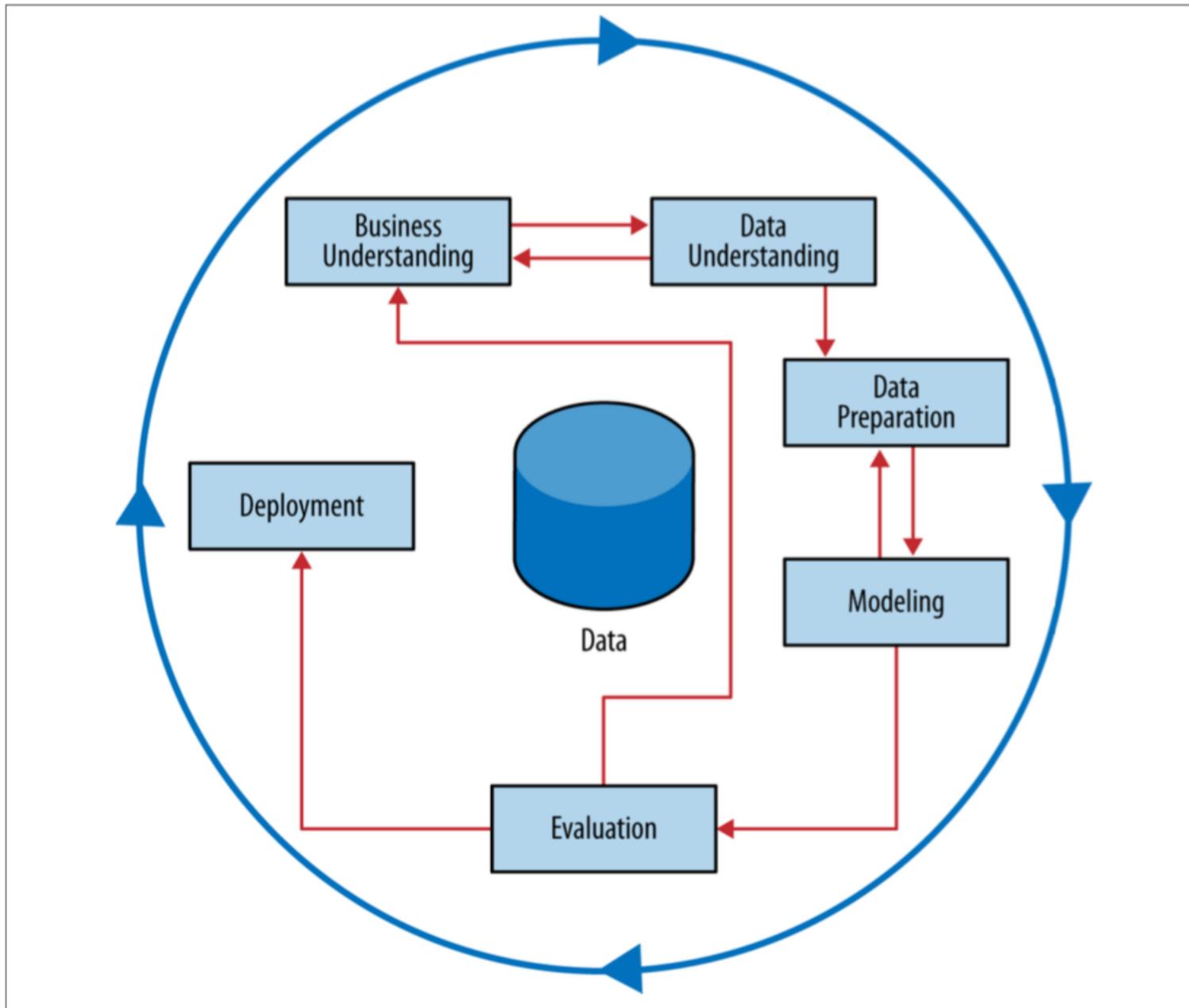
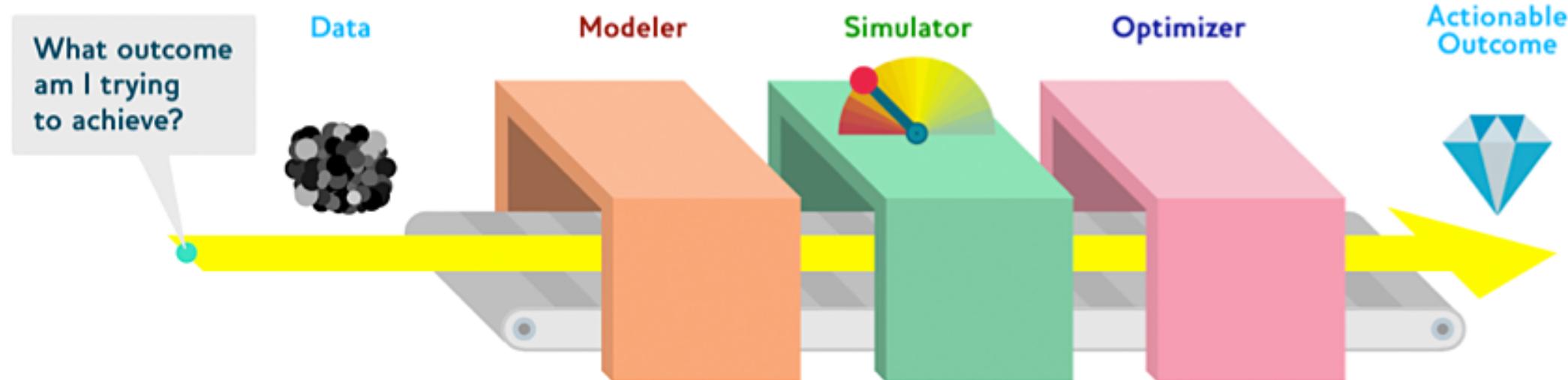


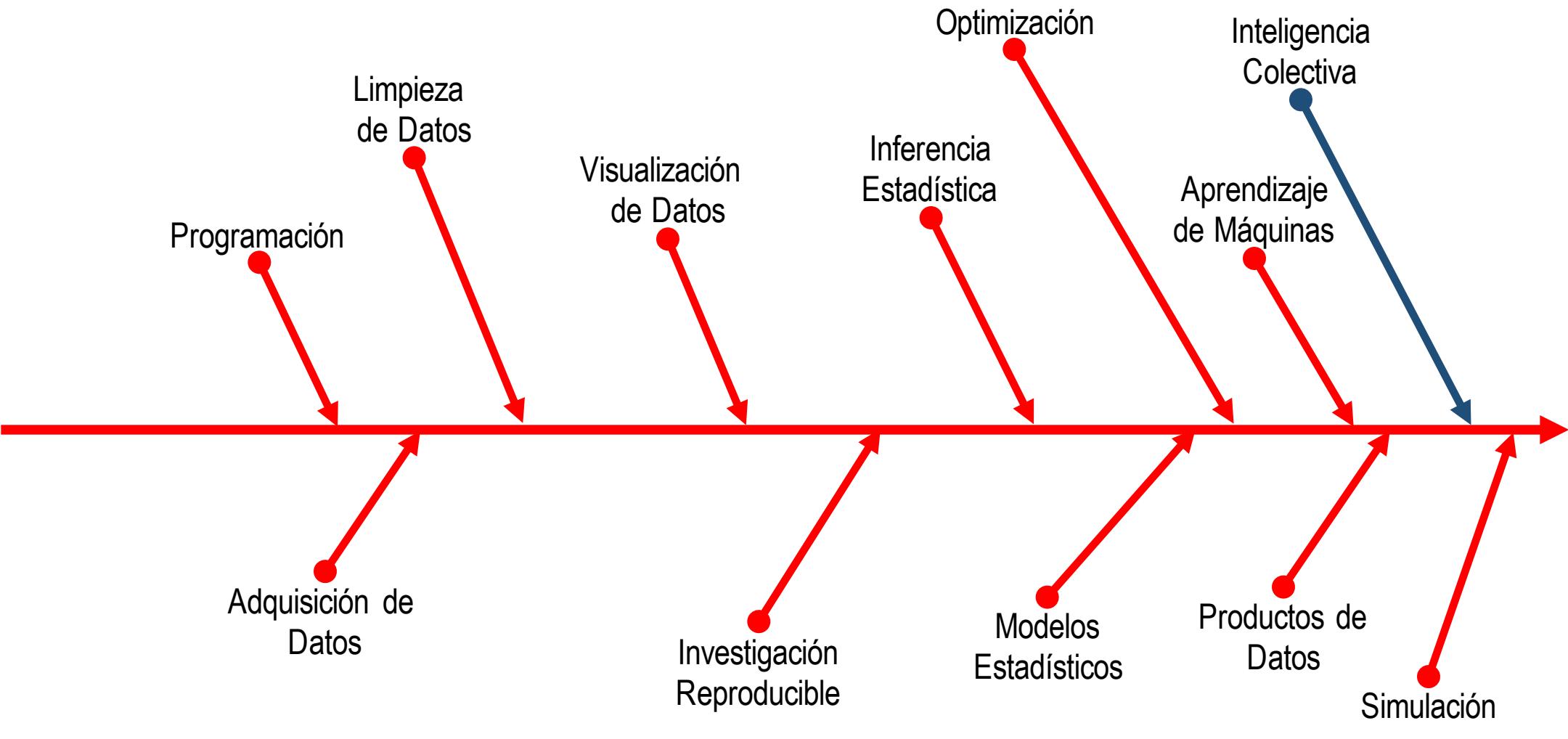
Figure 2-2. The CRISP data mining process.

Designing great data products

<http://radar.oreilly.com/2012/03/drivetrain-approach-data-products.html>



Data Science / Analytics



Data-driven decision making!

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

¿Ordenar un vector de números?

¿Extraer la tercera línea de texto de un conjunto de archivos?

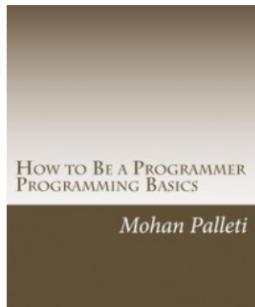
¿Calcular la suma de los primeros 20 números primos?



Algorithms Unlocked

By: Thomas H. Cormen

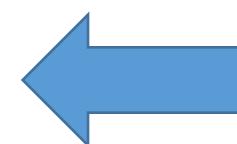
Have you ever wondered how your GPS can find the fastest way to your destination, selecting one route from seemingly countless possibilities in mere seconds? How your credit card account number is protected when you make a purchase over the Internet? The answer is algorithms. And how do...



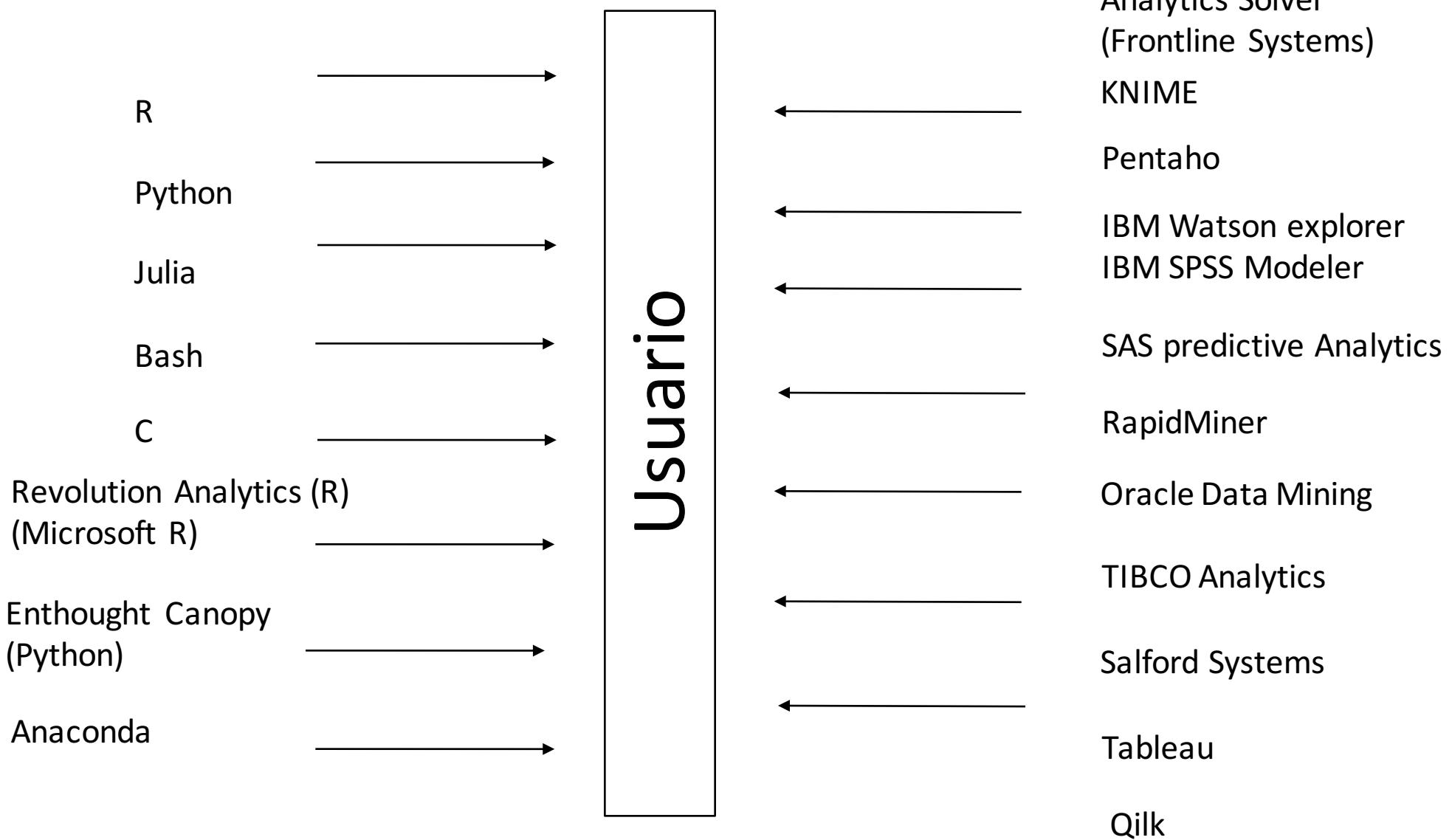
How to Be a Programmer: Programming Basics

By: Mohan Palleti

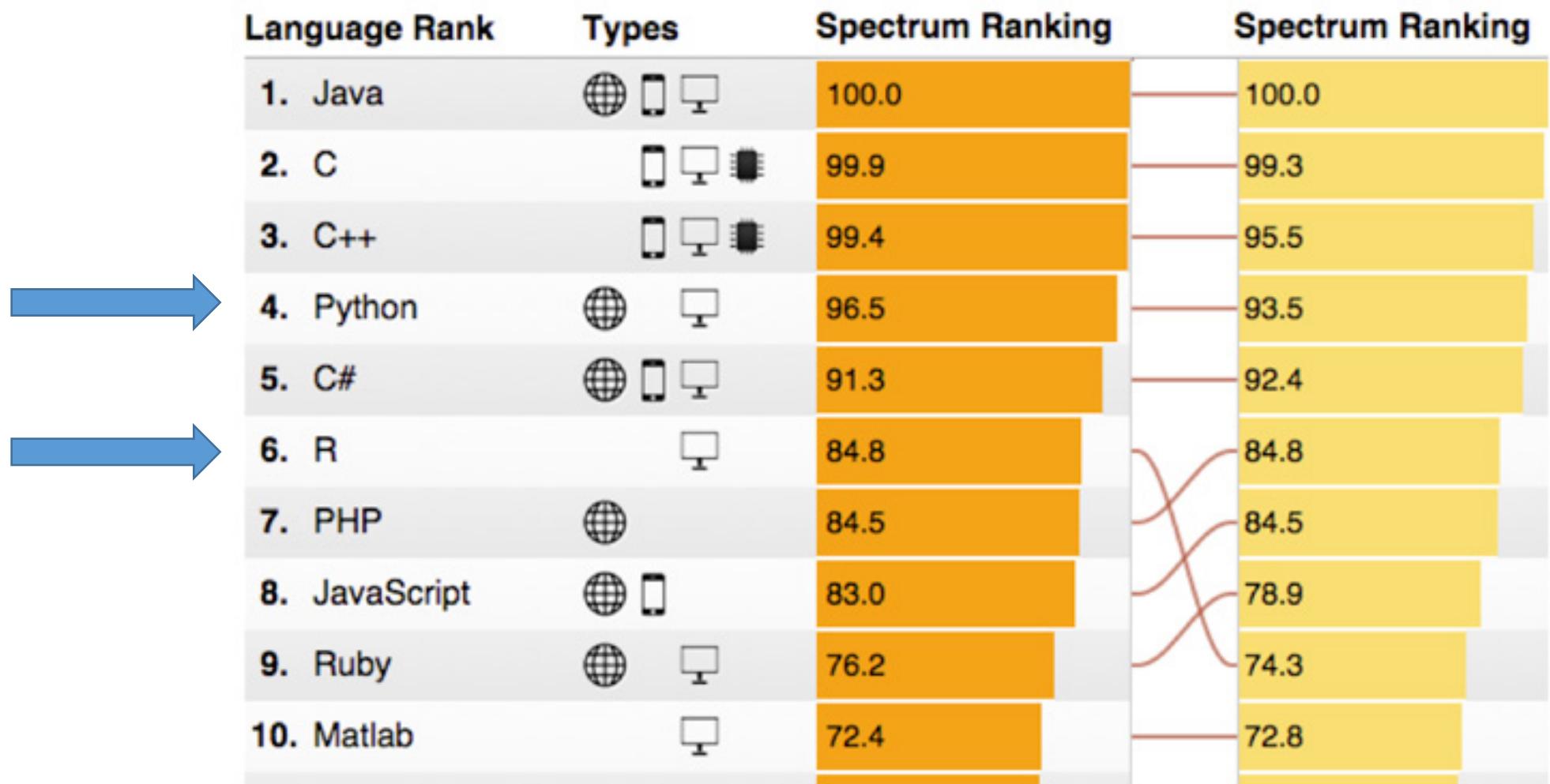
A Self-help 97 pages book to learn the basics of programming using Microsoft Excel's VBA tools. Ideal resource for school teachers and educators wanting to teach programming basics.



Programación vs Aplicaciones de usuario final



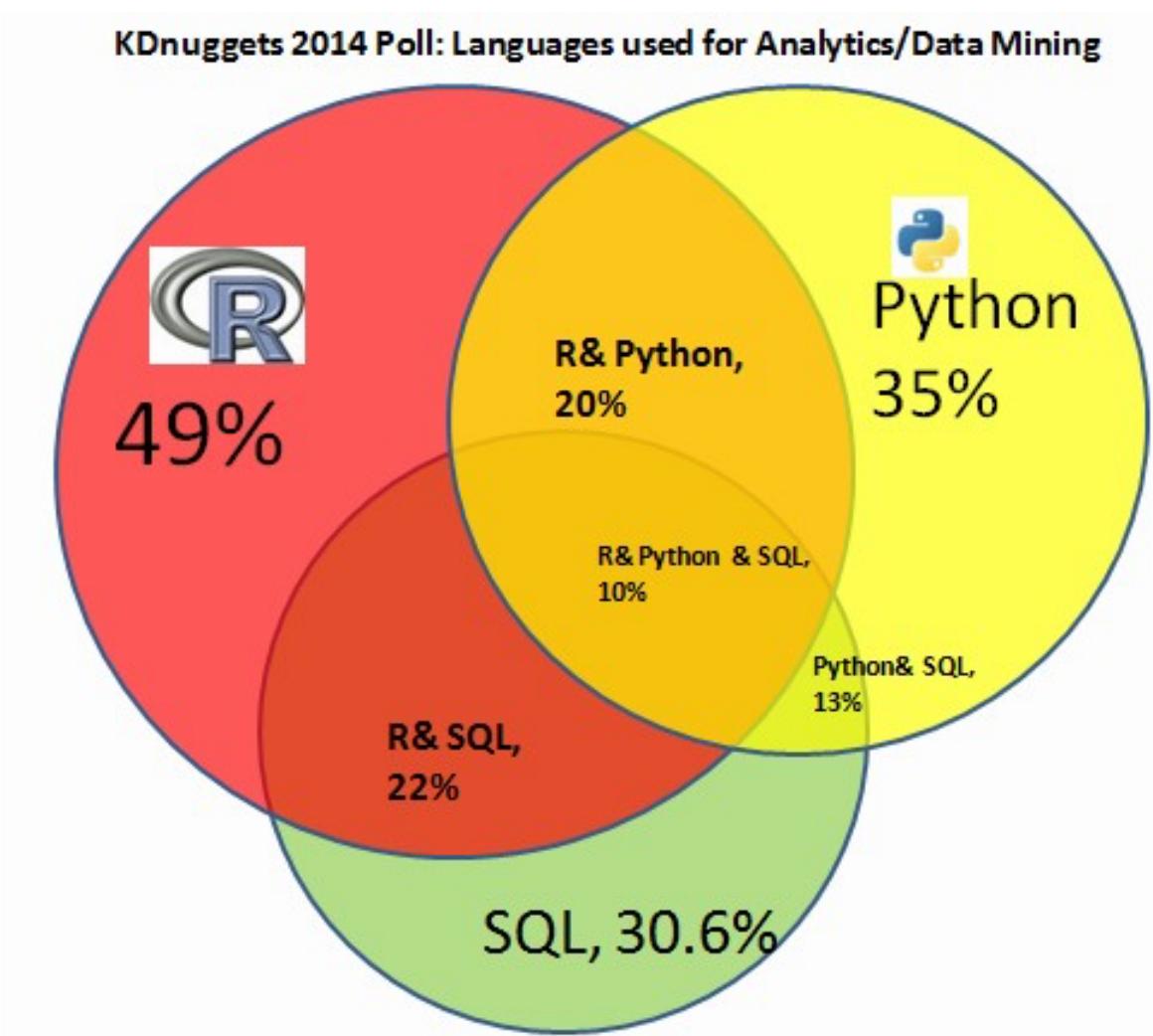
The 2015 Top Ten Programming Languages (IEEE Spectrum)



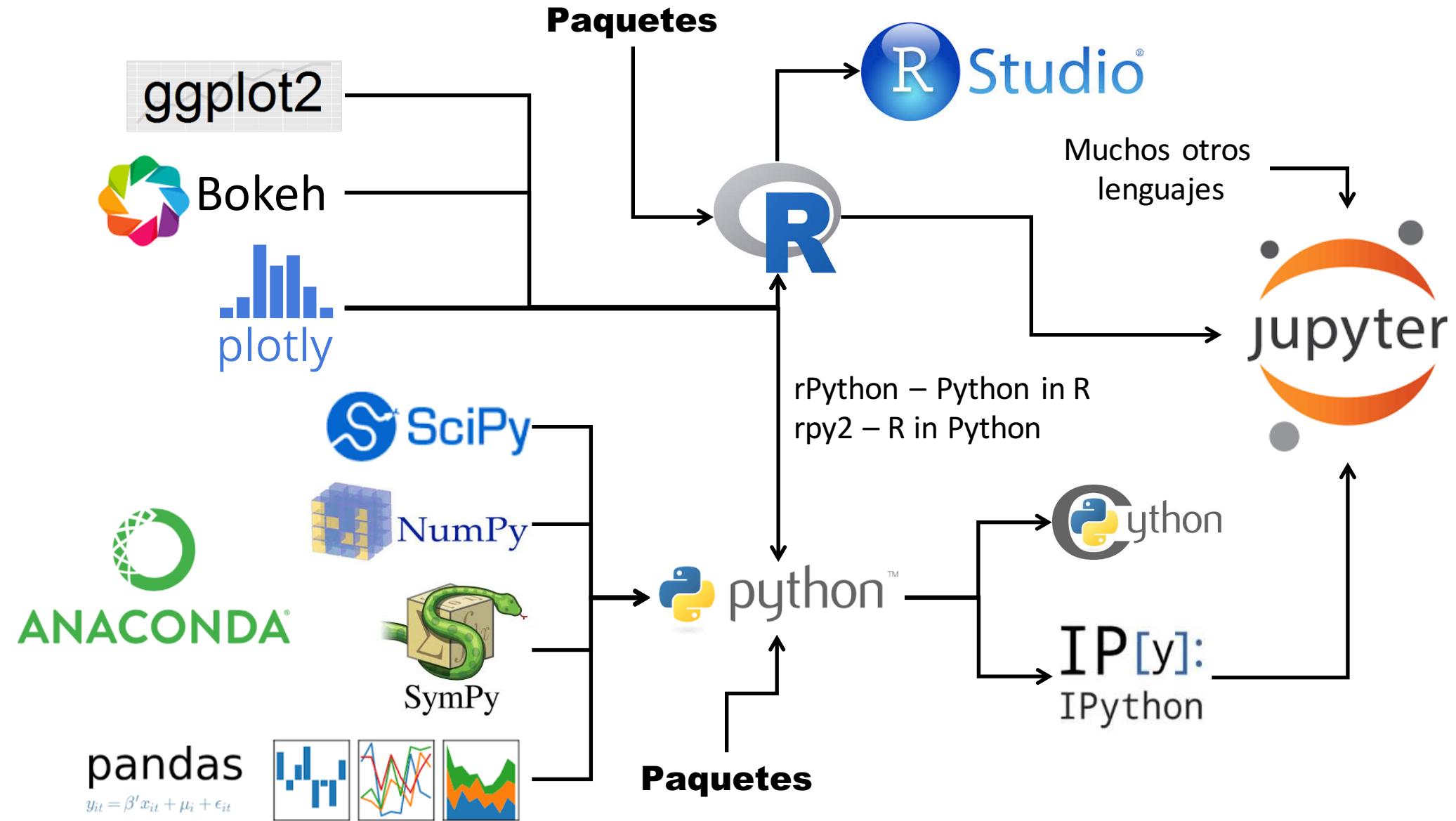
The **2016** Top Ten Programming Languages (IEEE Spectrum)

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

Popularidad de los lenguajes



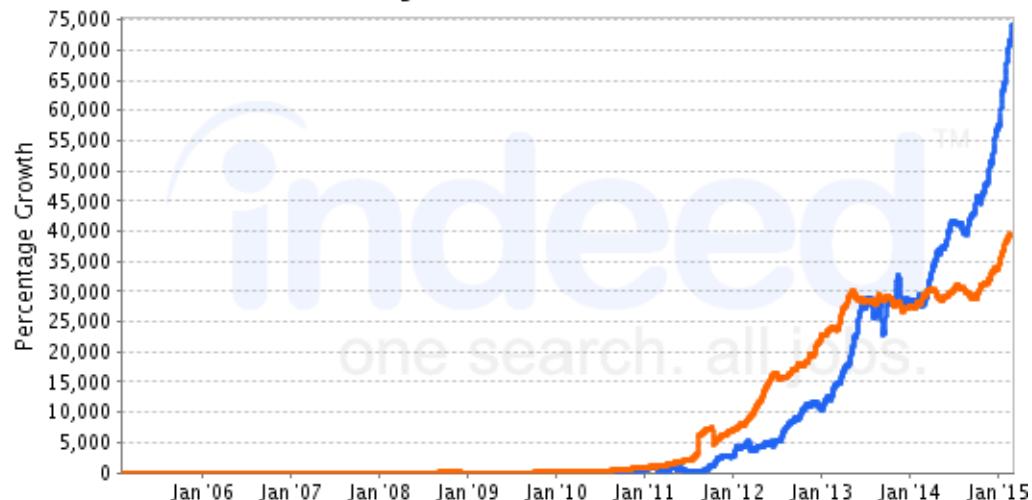
Ecosistema de computación científica: Python y R



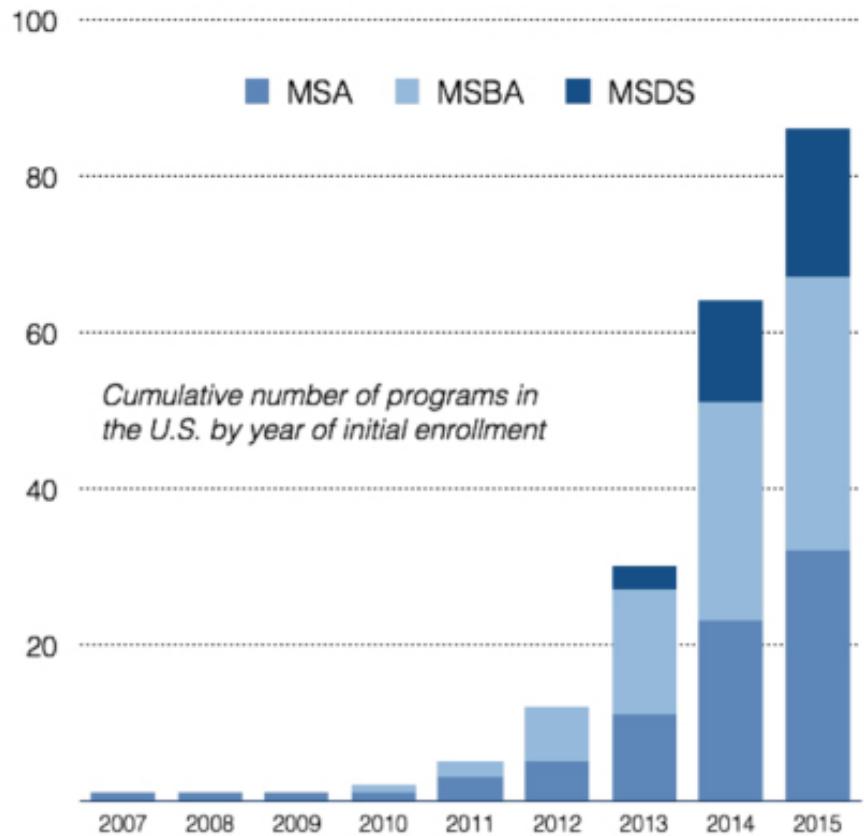
Evolución de empleos/educación en Big Data & Data Science

Job Trends from Indeed.com

— big-data — data-science

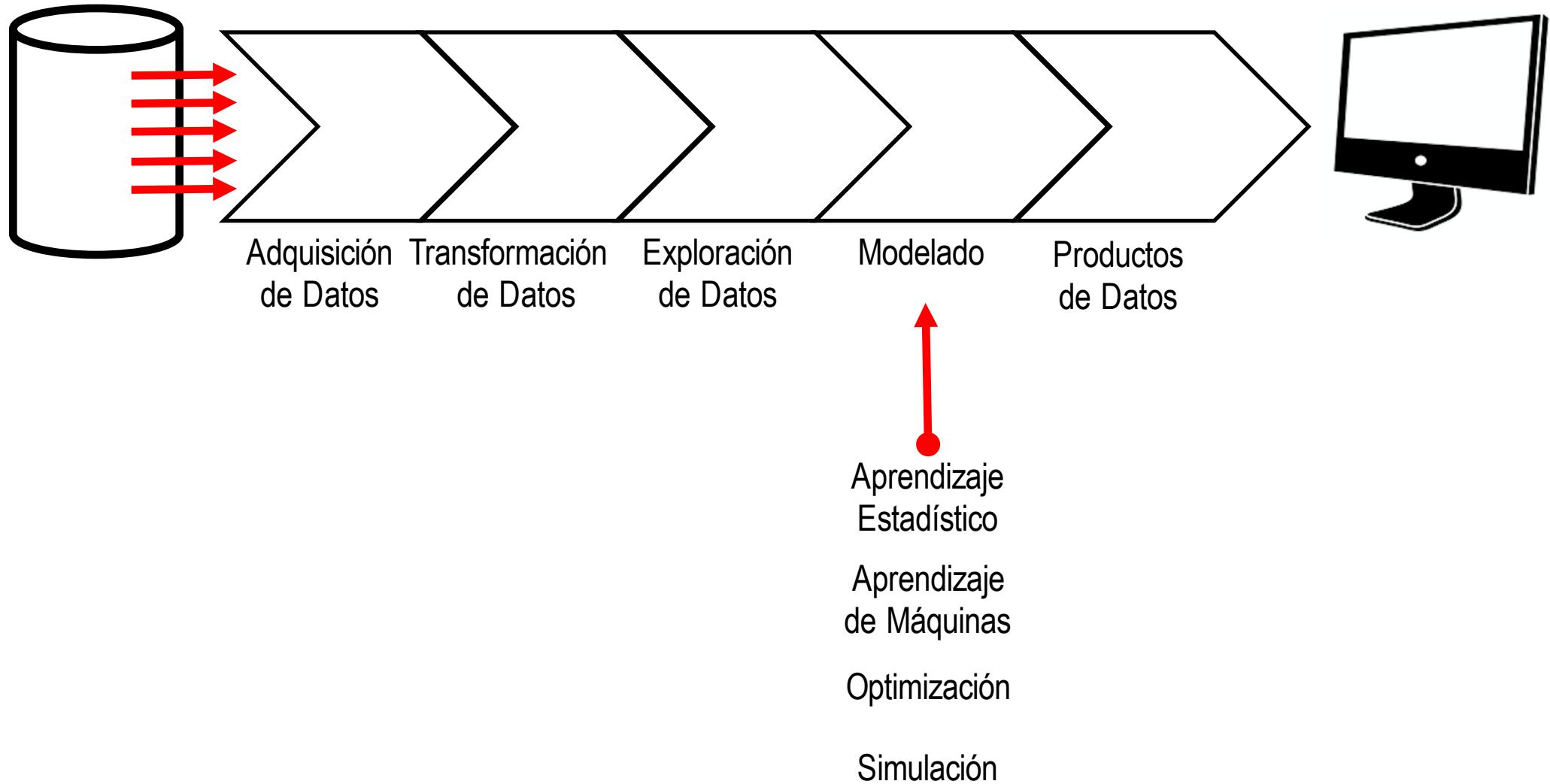


GROWTH OF MASTER'S DEGREE PROGRAMS IN ANALYTICS AND DATA SCIENCE

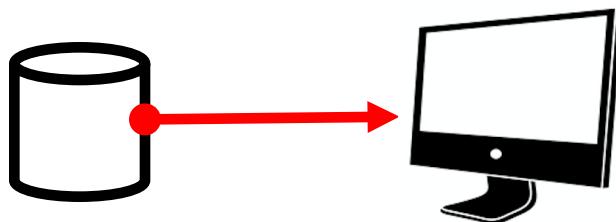


http://analytics.ncsu.edu/?page_id=4184

Fases en Data Science / Analytics

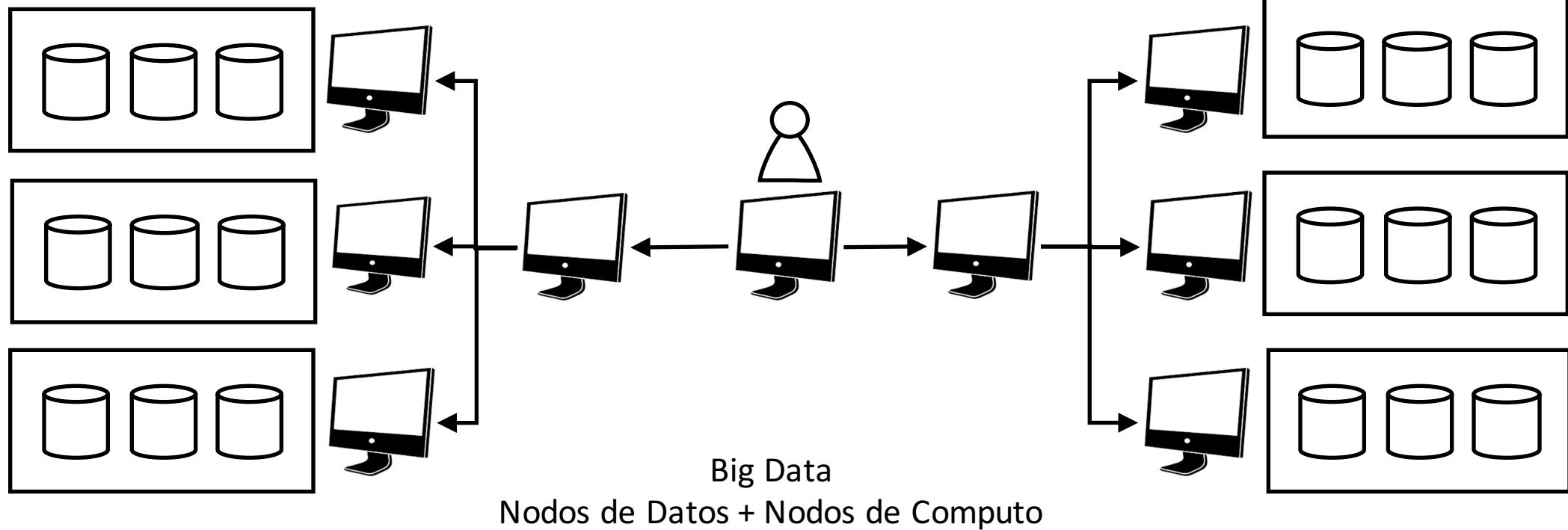


Aproximación tradicional vs Big Data

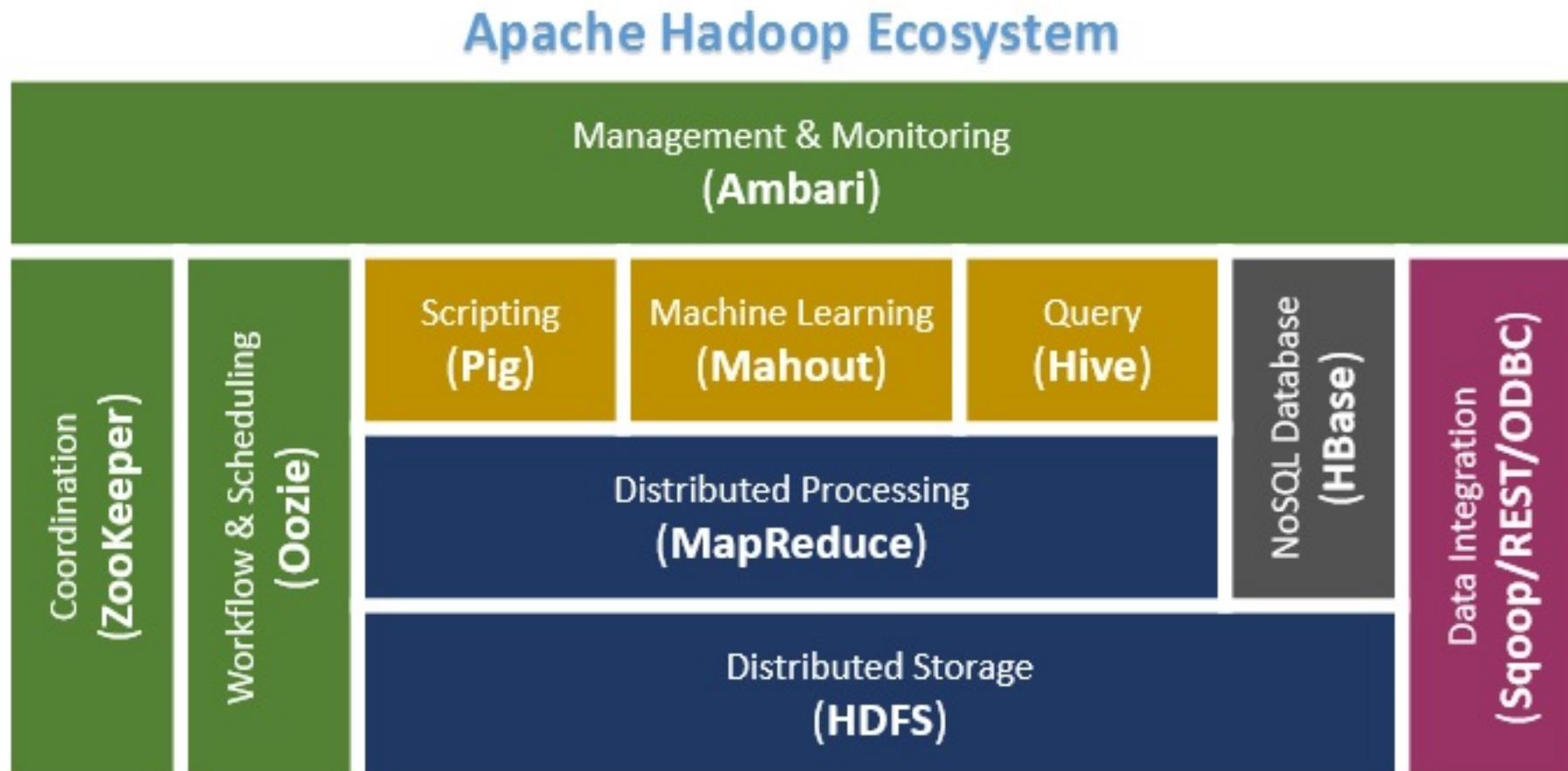


Programación Tradicional

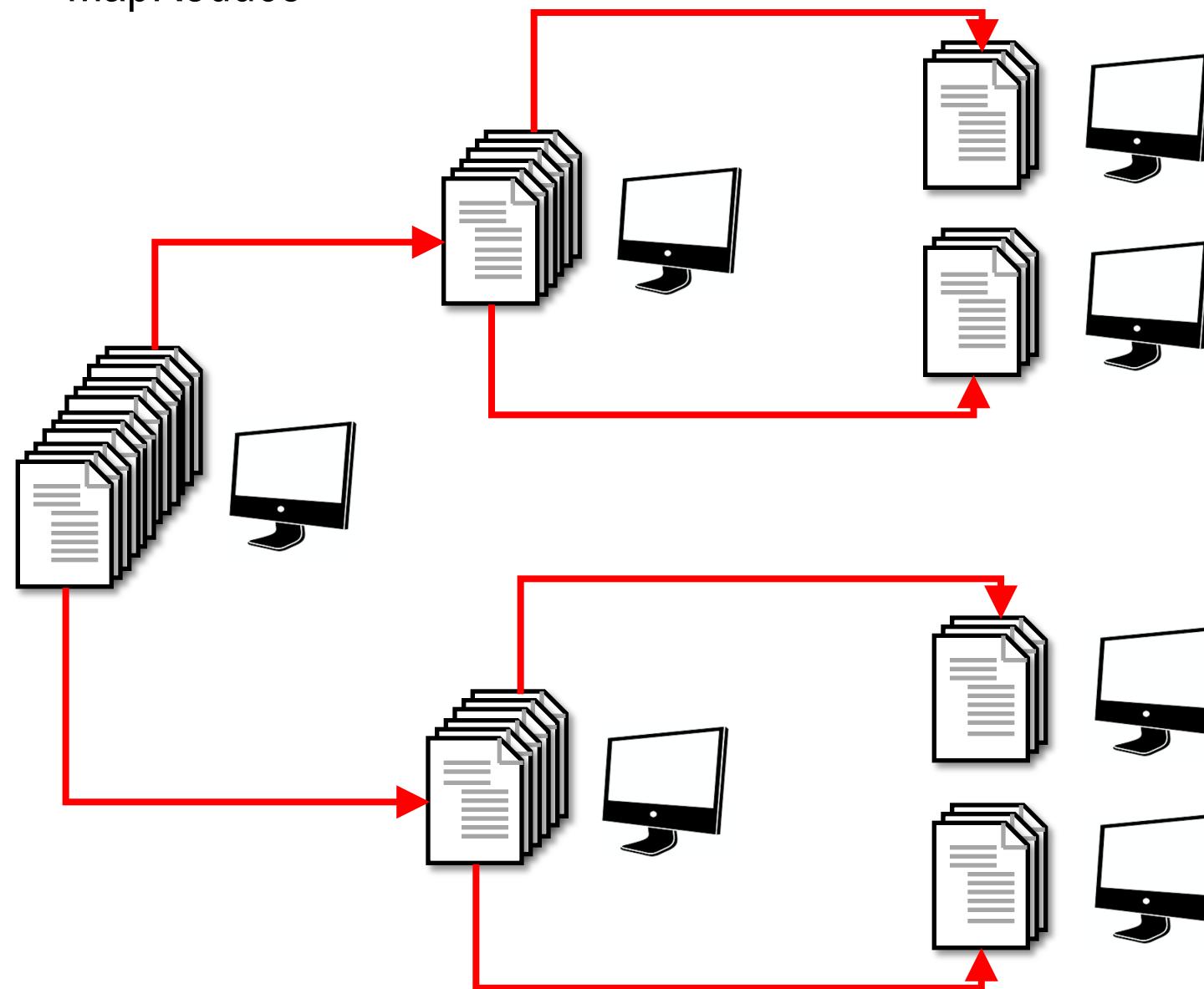
¿Cómo ejecutar algoritmos
tradicionales que son voraces en
recursos computacionales?



Big Data – Apache Hadoop



MapReduce



Adquisición y Limpieza de Datos

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

Páginas web (HTML) y Google Groups.

Bases de datos relacionales.

Lenguaje Natural.

Imágenes (Captcha)

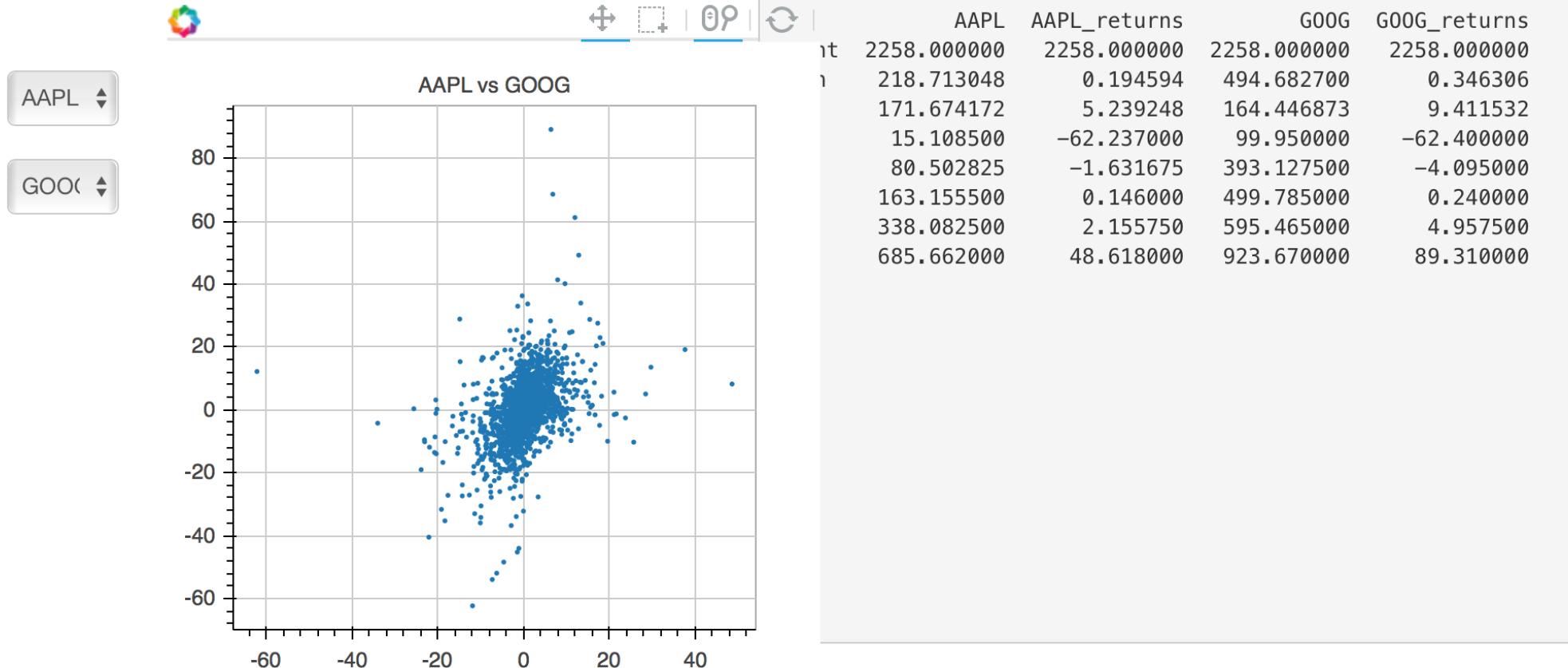
Manipulación de texto

Conversión de un formato a otro

Detección de datos faltantes, datos nulos, datos inconsistentes

Visualización de Datos

[link to this](#)



Investigación 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

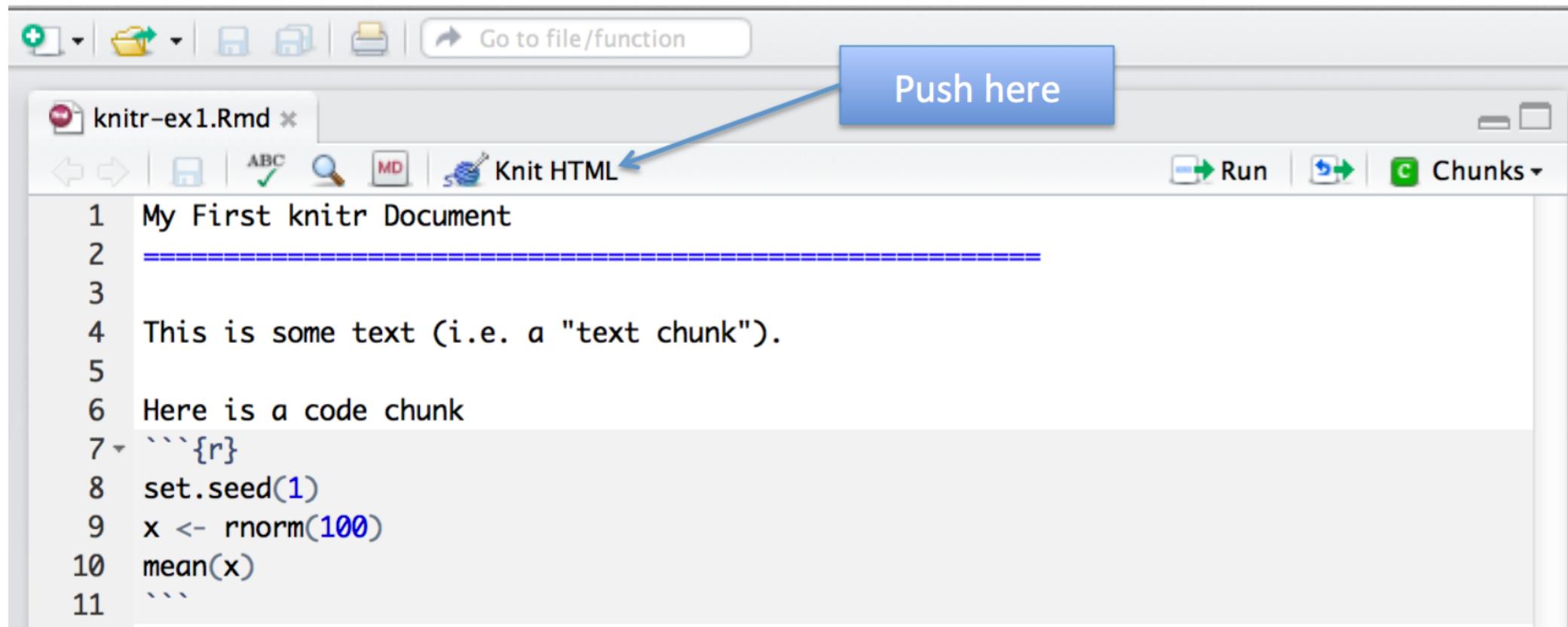
- List item 1
- List item 2
 - Subitem 2.1
 - Subitem 2.2

[Help on this page](#)

?

Ok Cancel

Investigación Reproducible (Markdown + R)



The screenshot shows the RStudio interface with a knitr document open. The document contains the following text:

```
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 ```
```

A blue arrow points from a blue box containing the text "Push here" to the "Knit HTML" button in the toolbar.

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)
```

Code input

```
## [1] 0.1089
```

Numerical output

Investigación Reproducible (Jupyter Notebook)

The screenshot shows a Jupyter Notebook interface 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,:]
```

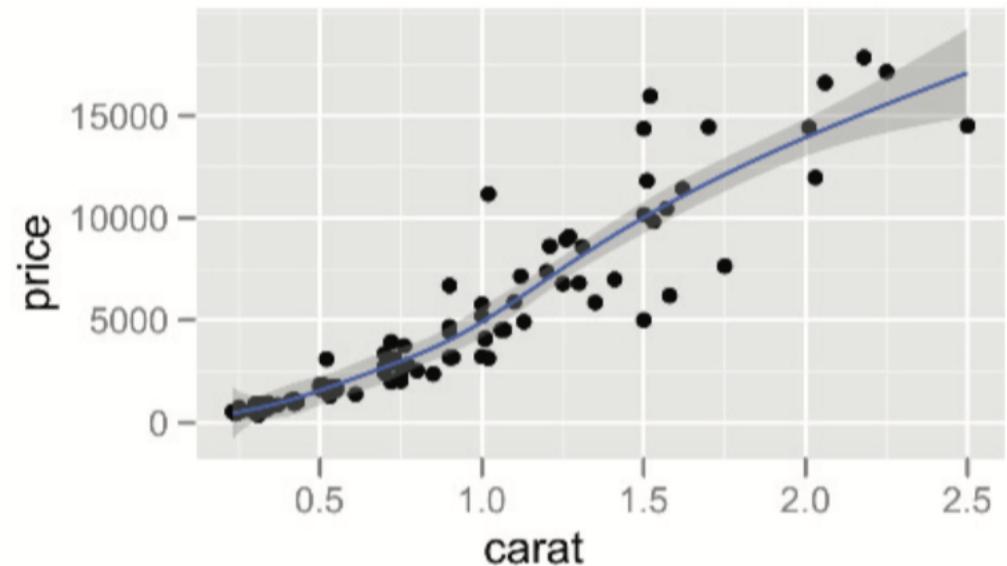
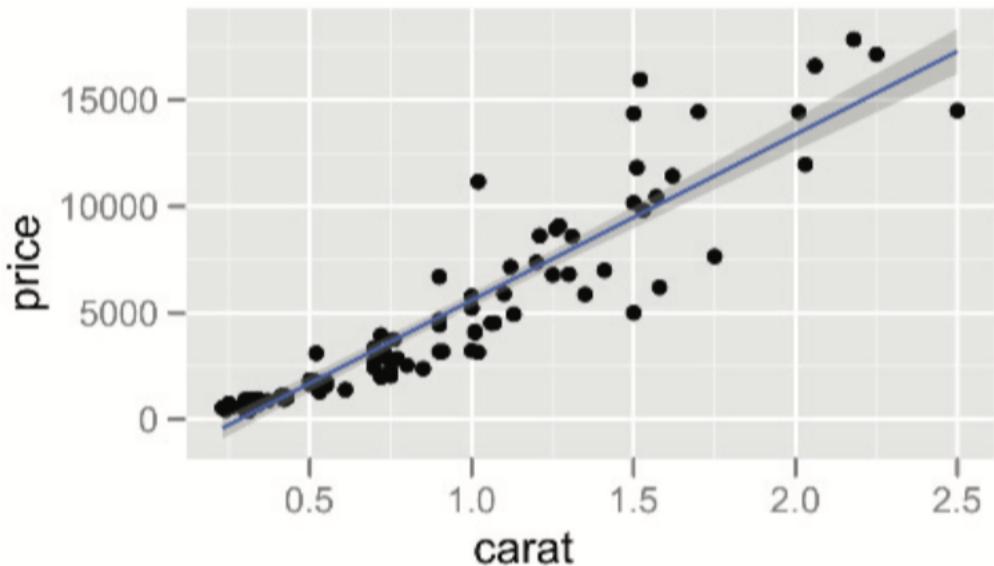
```
Out[4]: [

A plot window is visible at the bottom left, showing a blue line segment starting near the bottom right corner of a white square frame.


```

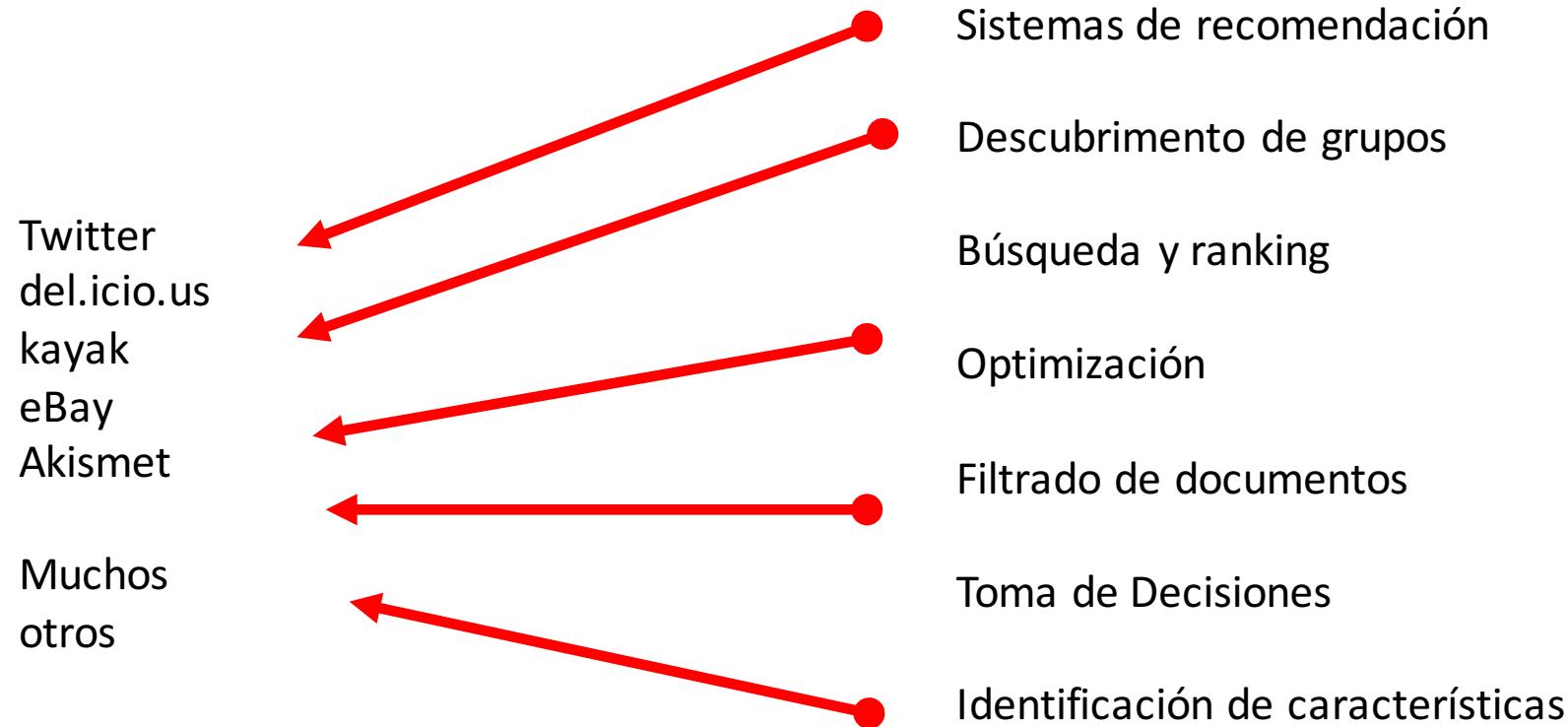
Modelos Estadísticos

Aplicación clásica



¿Y si hay 10 millones de datos?

Inteligencia Colectiva



Especialización en Analítica

Ciencia de los datos aplicada.

Decisiones bajo incertidumbre en las organizaciones.

Sistemas de bases de datos masivos.

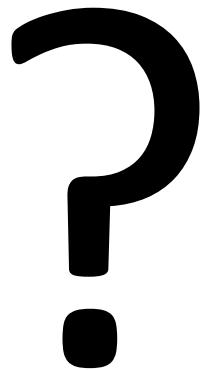
Aprendizaje de máquinas para datos masivos.

Optimización y simulación.

Modelado predictivo y series de tiempo.

Inteligencia de negocios.

En resumen....



Gracias por su atención

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>