

# **Analítica Avanzada: Un Enfoque Metodológico usando KNIME y CRISP-DM**

victor.toledo@coredevx.com

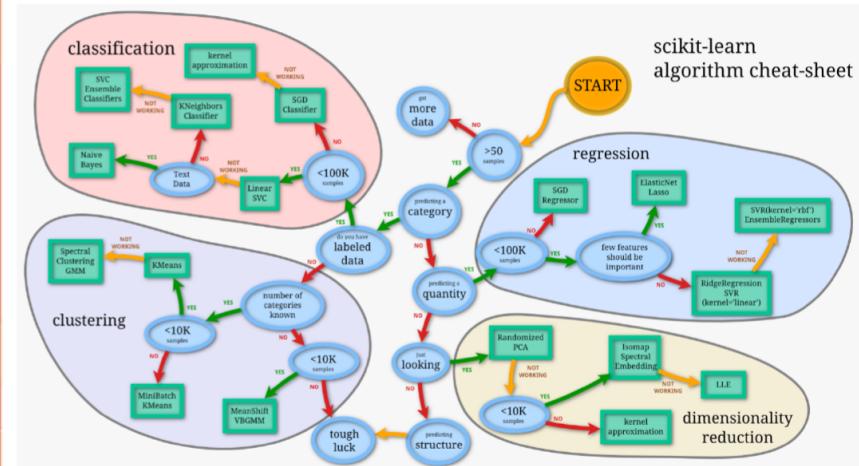
R+D Leader

# Outline

- Qué (creemos) es Analítica Avanzada (AA)
- Que Desafíos / Dimensiones existen en AA (Talento, Implementación, Infraestructura)
- Cuál es el Retorno con AA: Ecuación de ROI
- Aproximaciones a Desafíos en cada Dimensión (KNIME, CRISP-DM, DS-CI/CD)

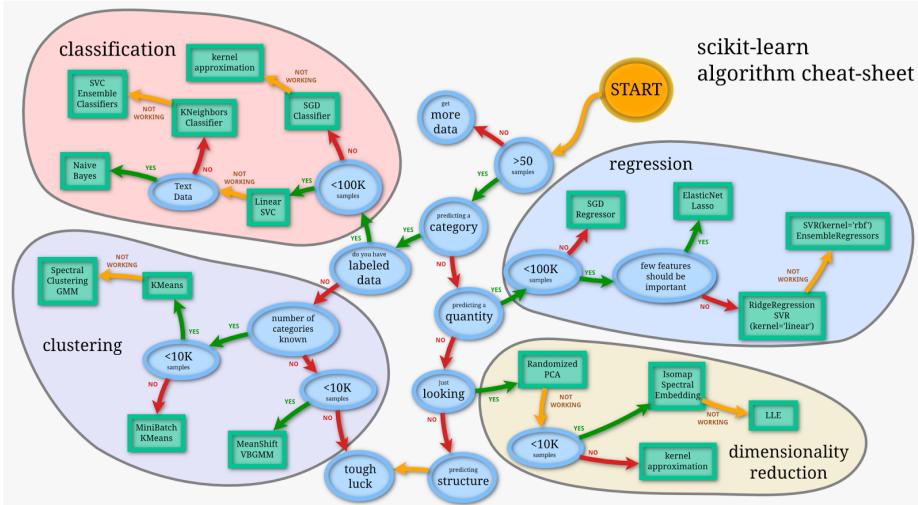
# Analítica Avanzada (AA)

	Business Intelligence	Advanced Analytics
Orientation	Rearview	Future
Types of questions	What happened When, who, how many	What will happen? What will happen if we change this one thing? What's next?
Methods	Reporting (KPIs, metrics) Automated Monitoring/Alerting (thresholds) Dashboards Scorecards OLAP (Cubes, Slice & Dice, Drilling) Ad hoc query	Predictive Modeling Data Mining Text Mining Multimedia Mining Descriptive Modeling Statistical / Quantitative Analysis Simulation & Optimization
Big Data	Yes	Yes
Data types	Structured, some unstructured	Structured and Unstructured
Knowledge Generation	Manual	Automatic
Users	Business Users	Data scientists, Business analysts, IT, Business Users
Business Initiatives	Reactive	Proactive

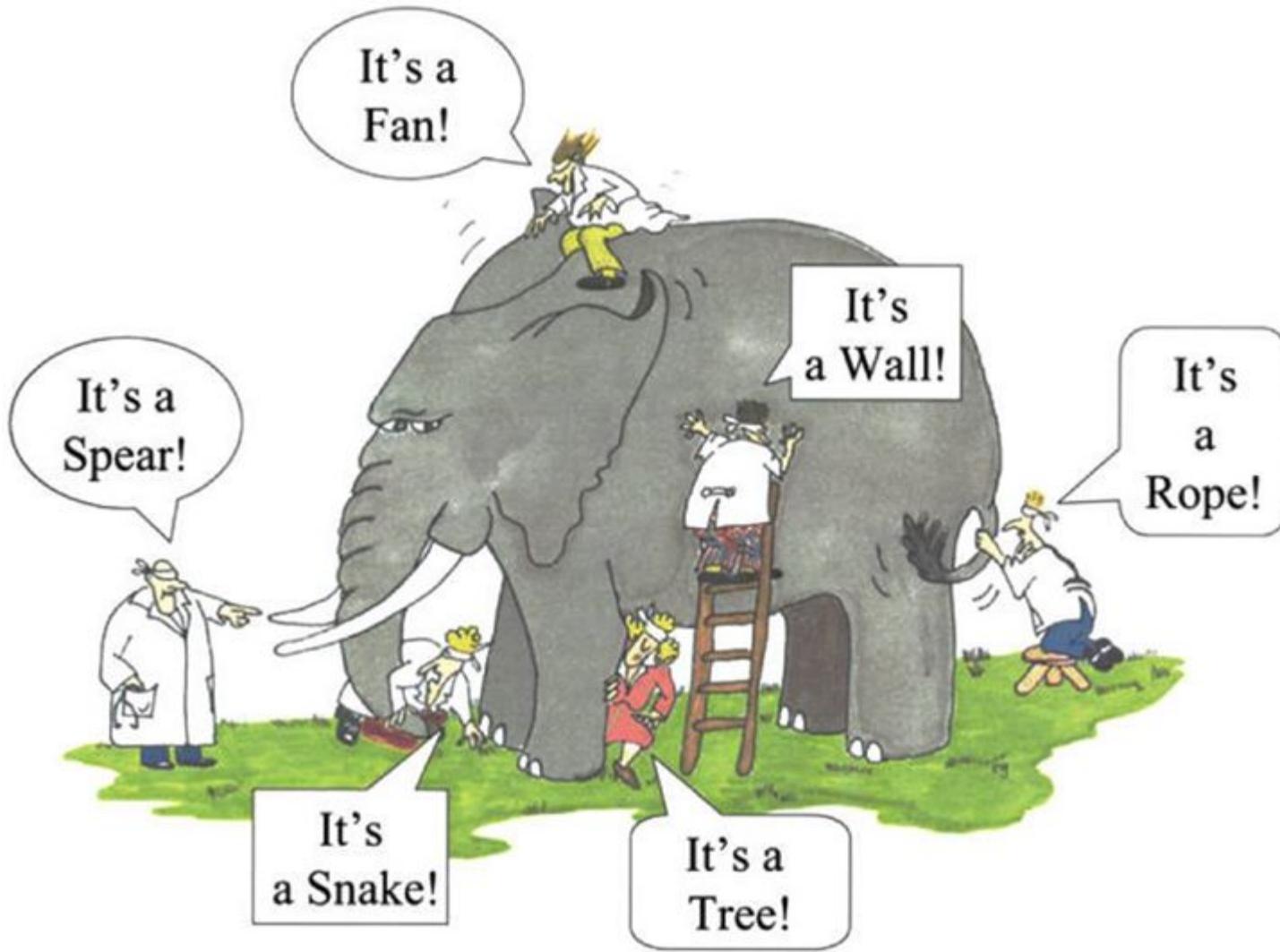




# AA: Landscape



**Hay un Elefante en la  
Sala!**

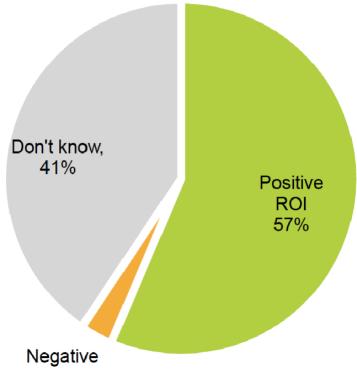


**El Desafío:  
Entregar valor a la  
Industria usando AA en  
tiempo y recursos  
definidos**

# 41% of Organizations Don't Know if Big Data ROI Will Be Positive or Negative

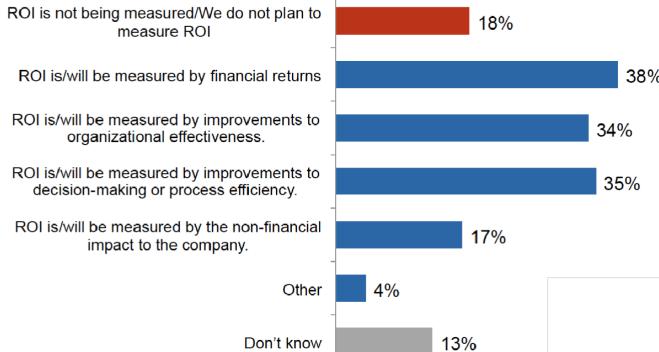
Please indicate whether your organization's big data investment has /is expected to have - a positive or negative ROI?

Positive/negative ROI

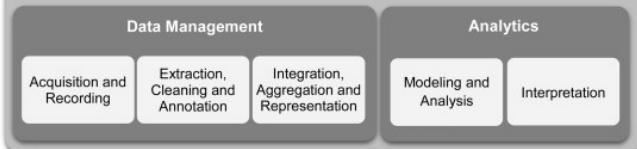


How is ROI being measured/How will ROI be measured for your organization's big data investment?

ROI Measurement



Big data Processes

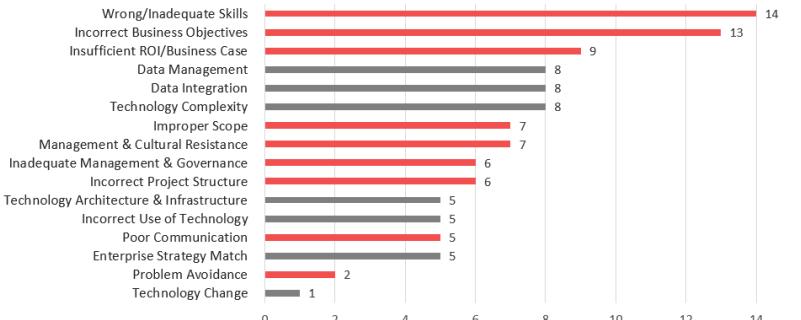


## Big Data ROI



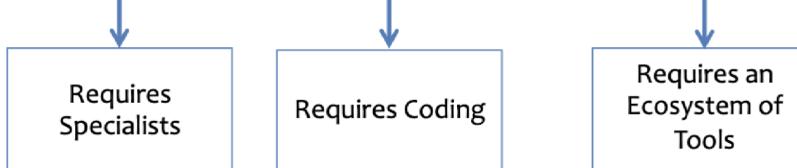
## Becker: Leading Causes in Big Data Project Failures

(red color added to highlight project management and organizational issues)



# AA: ROI Equation

$$[\text{Cost of Talent}] \times [\text{Time Spent}] + [\text{Overall Infrastructure}]$$



First, identify the talents you need to have access to:

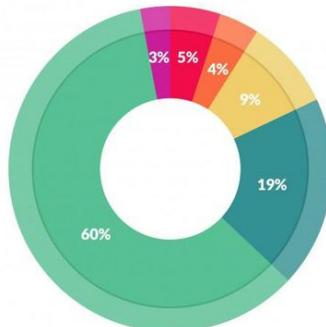
Project management	Subject expertise
Data wrangling	Design
Data analysis	Storytelling

Next, map talents to team members:

Person	Talent	Person	Talent
Anand		Roberto	■■■■■
Cameron	■■■■■	Stephani	■■■■■
Emily	■■■■■	Susan	■■■■■
Kevin	■■■■■	Xia-Li	■■■■■

Finally, assess how much depth you have for each type of talent:

Talent	Depth
Project management	■■
Data wrangling	■■■■■
Data analysis	■■■■■
Subject expertise	■■■■■
Design	■■■■■
Storytelling	■■■■■



## What data scientists spend the most time doing

- Building training sets: 3%
- Cleaning and organizing data: 60%
- Collecting data sets: 19%
- Mining data for patterns: 9%
- Refining algorithms: 4%
- Other: 5%



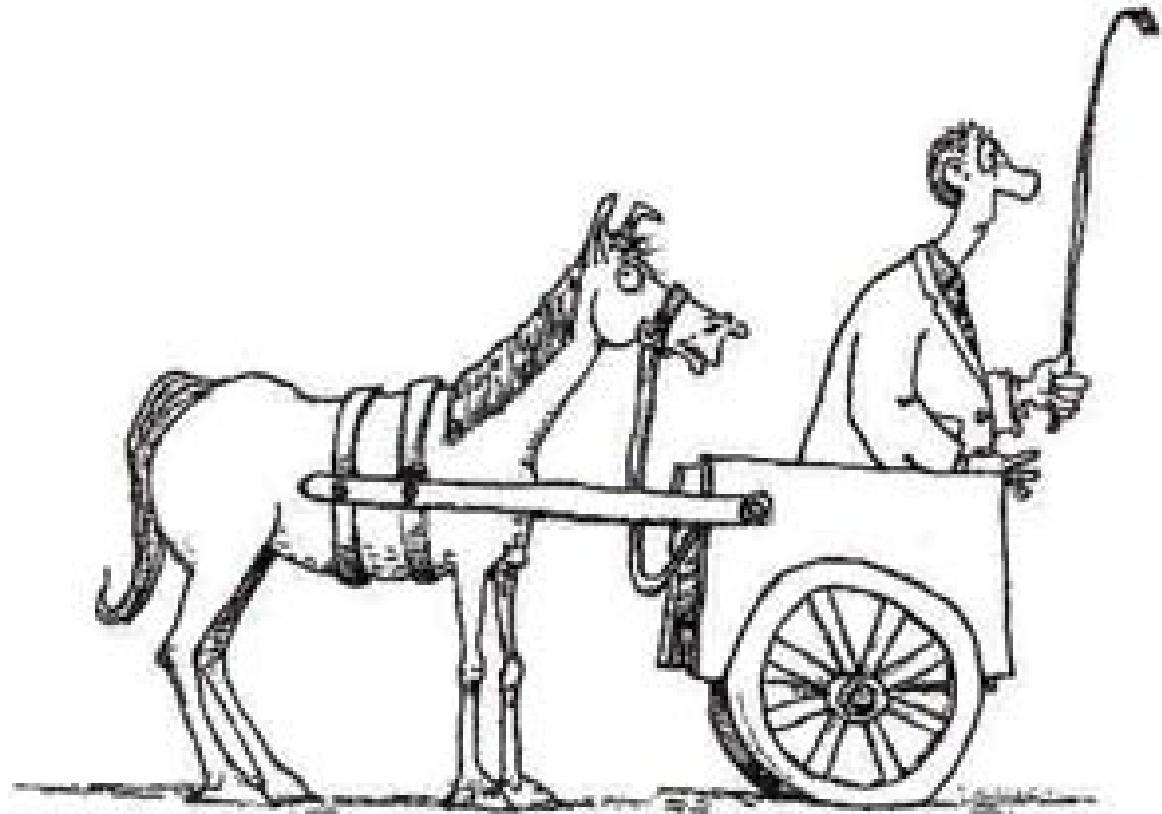
From: "Data Science and the Art of Persuasion,"  
by Scott Berinato, January–February 2019

HBR



SKETCH2X.COM

**Una Respuesta:  
Definir un Framework  
Robusto, y actuar de  
acuerdo a éste**



What main methodology are you using for your analytics, data mining, or data science projects ? [200 votes total]

**votes total]**

2014 poll 2007 poll

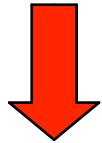
A horizontal bar chart comparing the percentage of respondents using various machine learning tools across different methodology categories. The categories are listed on the left, and the tools are represented by colored bars (red, green, blue) on the right. The legend indicates that red represents Python, green represents Anaconda, blue represents scikit-learn, and orange represents R.

Methodology Category	Tool 1	Tool 2	Tool 3	Tool 4
RISP-DM (86)	43%	42%	0%	0%
My own (55)	27.5%	19%	0%	0%
SEMMA (17)	8.5%	13%	0%	0%
Other, not domain-specific (16)	8%	4%	0%	0%
KDD Process (15)	7.5%	7.3%	0%	0%
My organizations' (7)	3.5%	5.3%	0%	0%
A domain-specific methodology (4)	2%	1:Python	46%	0%
None (0)	0%	6:Anaconda	46%	0%
		9:scikit-learn	49%	109%

KDnuggets 2018 Data Science, Machine Learning Software Poll:  
Top Tools Associations



# Framework



Business Understanding	Data Understanding	Data Preparation	Modeling	Evaluation	Deployment
<b>Determine Business Objectives</b> <i>Background</i> <i>Business Objectives</i> <i>Business Success Criteria</i>	<b>Collect Initial Data</b> <i>Initial Data Collection Report</i>	<b>Select Data</b> <i>Rationale for Inclusion/Exclusion</i>	<b>Select Modeling Techniques</b> <i>Modeling Technique</i> <i>Modeling Assumptions</i>	<b>Evaluate Results</b> <i>Assessment of Data Mining Results w.r.t. Business Success Criteria</i> <i>Approved Models</i>	<b>Plan Deployment</b> <i>Deployment Plan</i>
<b>Assess Situation</b> <i>Inventory of Resources Requirements, Assumptions, and Constraints</i> <i>Risks and Contingencies</i> <i>Terminology</i> <i>Costs and Benefits</i>	<b>Describe Data</b> <i>Data Description Report</i>	<b>Clean Data</b> <i>Data Cleaning Report</i>	<b>Construct Data</b> <i>Derived Attributes</i> <i>Generated Records</i>	<b>Review Process</b> <i>Review of Process</i>	<b>Plan Monitoring and Maintenance</b> <i>Monitoring and Maintenance Plan</i>
<b>Explore Data</b> <i>Data Exploration Report</i>	<b>Explore Data</b> <i>Data Exploration Report</i>	<b>Build Model</b> <i>Parameter Settings</i> <i>Models</i>	<b>Integrate Data</b> <i>Merged Data</i>	<b>Determine Next Steps</b> <i>List of Possible Actions</i> <i>Decision</i>	<b>Produce Final Report</b> <i>Final Report</i> <i>Final Presentation</i>
<b>Verify Data Quality</b> <i>Data Quality Report</i>	<b>Format Data</b> <i>Reformatted Data</i>	<b>Assess Model</b> <i>Model Assessment</i> <i>Revised Parameter Settings</i>	<b>Dataset</b> <i>Dataset Description</i>	<b>Review Project</b> <i>Experience Documentation</i>	<b>Review Project</b> <i>Experience Documentation</i>

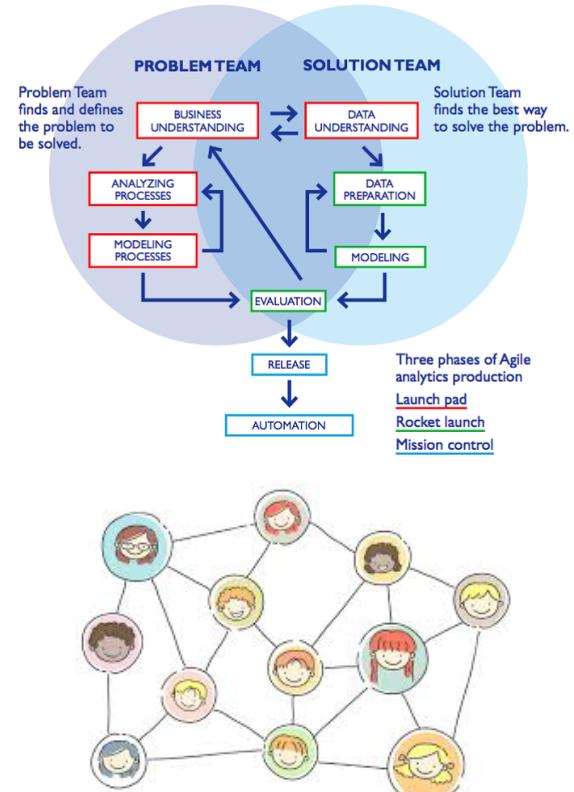
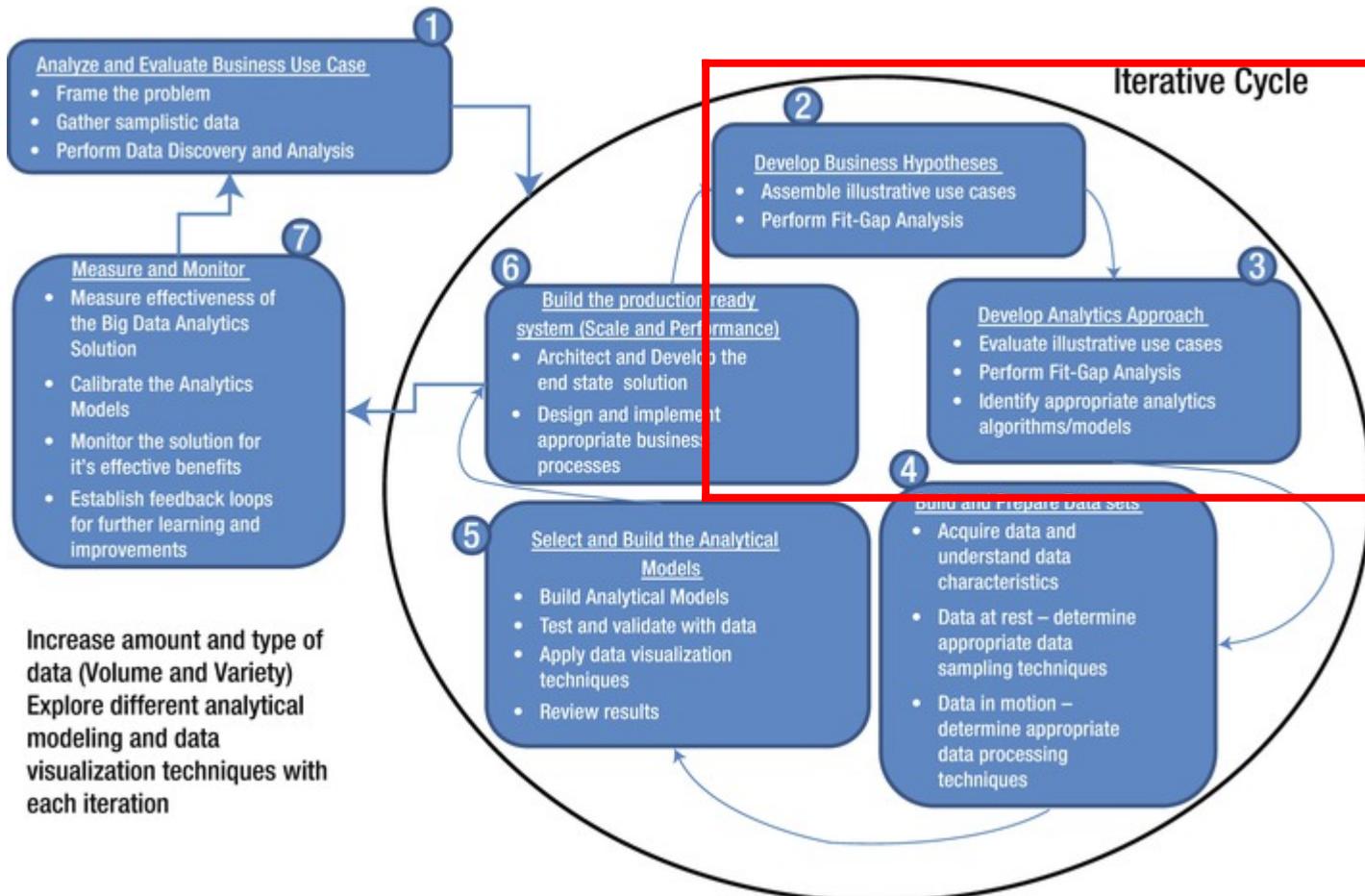


Figure 3: Generic tasks (bold) and outputs (italic) of the CRISP-DM reference model

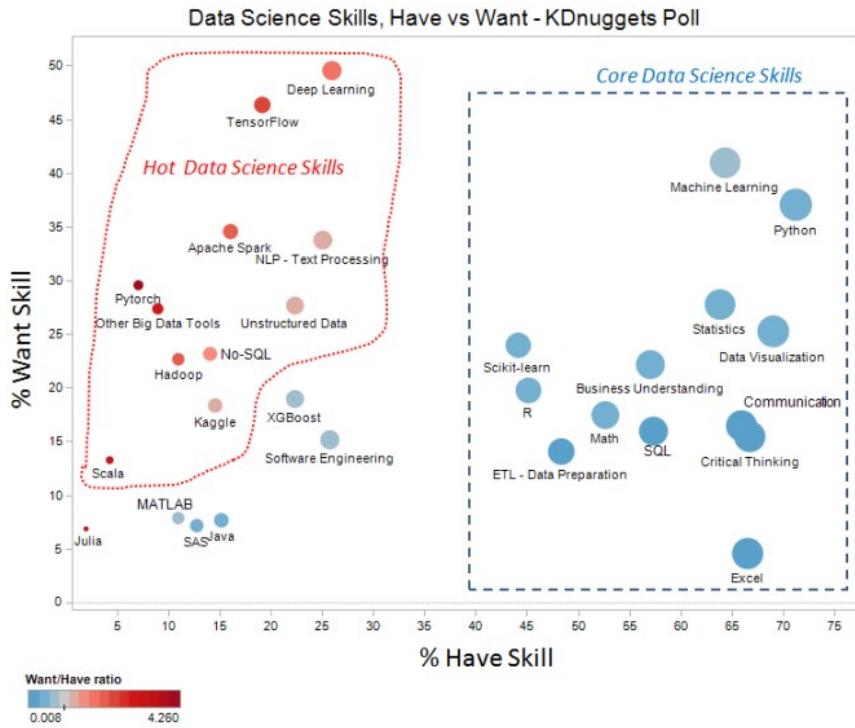
# Metodología: CRISP-DM

- Define pasos generales (As-Appropriate)
- No específica en Entregables.
- No específica en Artefactos
- No específica en Roles y Responsabilidades

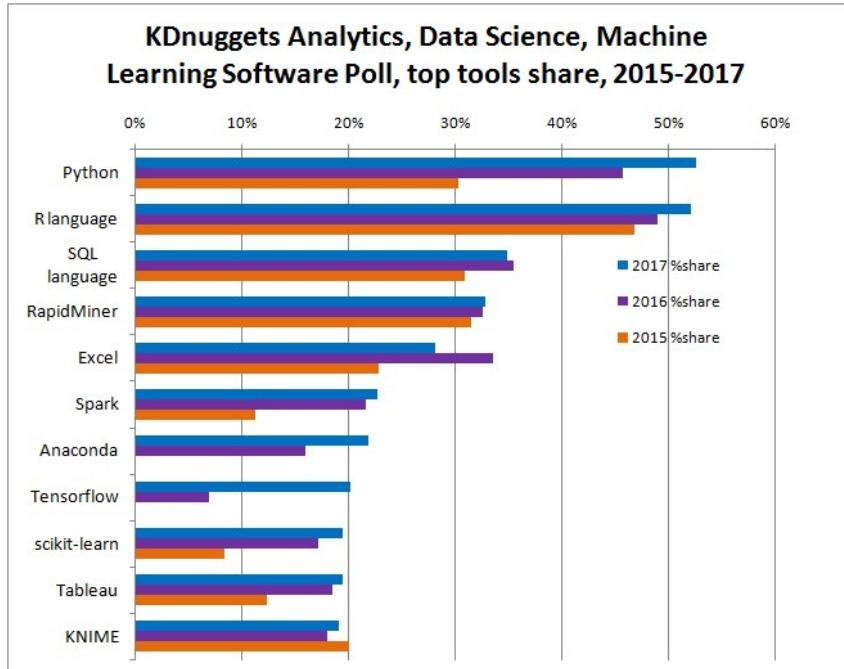
Requerimos especificarlo para asegurar tiempos y recursos involucrados



# AA: Skills+Tools

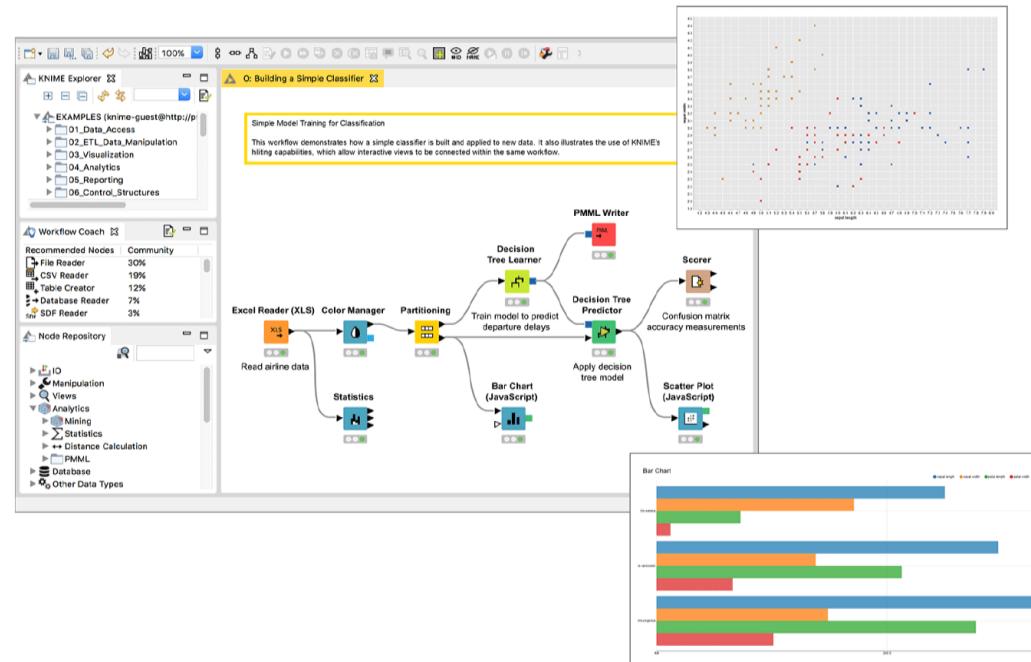


KDnuggets Analytics, Data Science, Machine Learning Software Poll, top tools share, 2015-2017

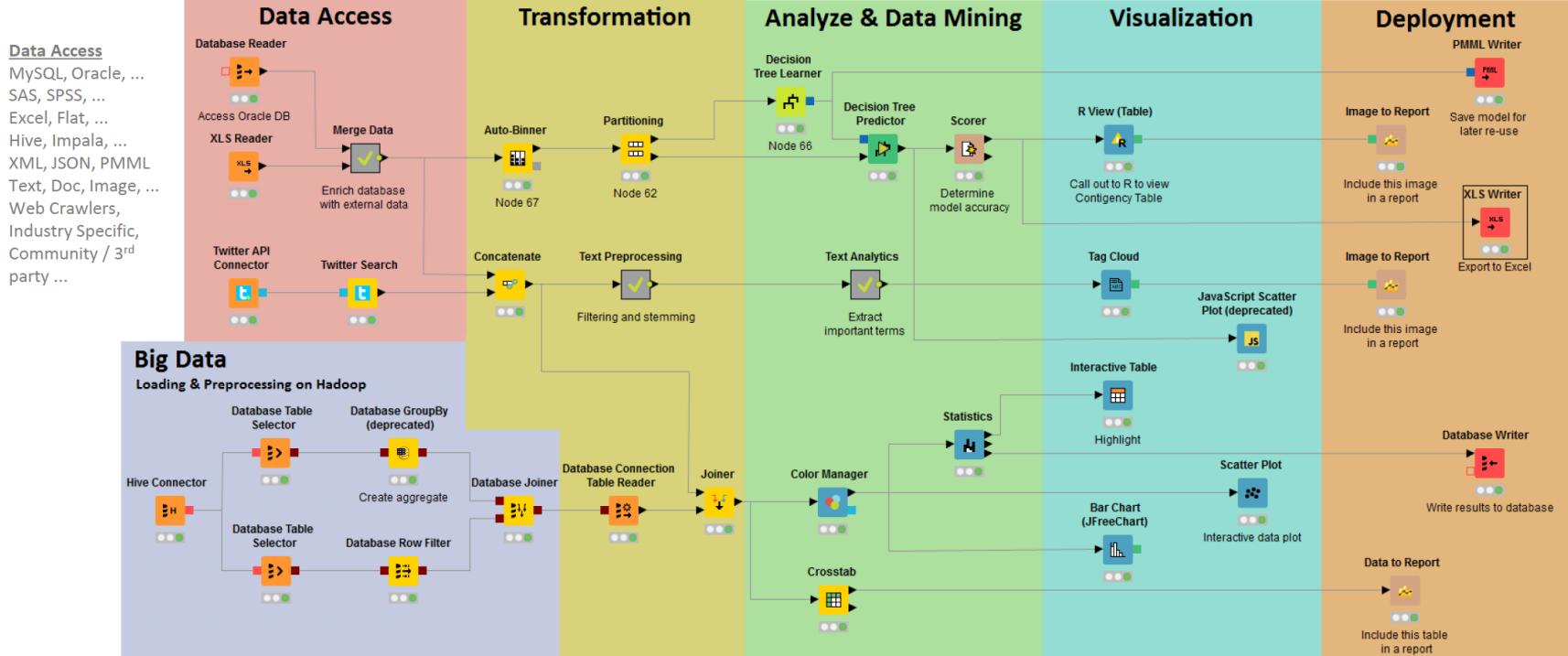


# What is KNIME Analytics Platform?

- A tool for data analysis, manipulation, visualization, and reporting
- Based on the graphical programming paradigm
- Provides a diverse array of extensions:
  - Text Mining
  - Network Mining
  - Cheminformatics
  - Big Data
  - Many integrations, such as Java, R, Python, Weka, H2O, etc.

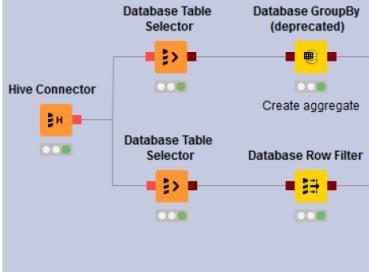


# Over 2000 native and embedded nodes included:



## Big Data

### Loading & Preprocessing on Hadoop



## Big Data

Hive, Impala, HDFS Vertical, Teradata/Aster, Spark, MLlib, Community / 3<sup>rd</sup> party, ...

## Transformation

Row, Column, Matrix, Text, Image, Networks, Time Series, Java, Python, Community / 3<sup>rd</sup> party, ...

## Analysis & Mining

Statistics, Machine Learning, Data Mining, Web Analytics, Text Mining, Network Analysis, Social Media Analysis, R, Weka, Python, Community / 3<sup>rd</sup> party, ...

## Visualization

R, Python, JFreeChart, JavaScript, Community / 3<sup>rd</sup> party, ...

## Deployment

via BIRT, PMML, XML, JSON, Databases, Excel, Flat, etc., Text, Doc, Image, Industry Specific, Community / 3<sup>rd</sup> party, ...

# Implementación: Reducir Codificación

KNIME Explorer

EXAMPLES (knime@hub.knime.com)

- 00\_Components
- 01\_Data\_Access
- 02\_ETL\_Data\_Manipulation
- 03\_Visualization
- 04\_Analytics
- 05\_Reporting
- 06\_Control\_Structures
- 07\_Scripting
- 08\_Other\_Analytics\_Types
  - 01\_Text\_Processing
  - 02\_Chemistry\_and\_Life\_Sciences

Node Repository

File Reader

Linear Correlation

Histogram

Pie/Donut Chart

Number To String

Partitioning

Variable Numérica a Segmentar Train / Test Predecir

CommentDataAnalysis

XRAY

Description

File Reader

This node can be used to read data from an ASCII file or URL location. It can be configured to read various formats.

When you open the node's configuration dialog and provide a filename, it tries to guess the reader's settings by analyzing the content of the file. Check the results of these settings in the preview table. If the data shown is not correct or an error is reported, you can adjust the settings manually (see below).

The file analysis runs in the background and can be cut short by clicking the "Quick scan", which shows if the analysis takes longer. In this case the file is not analyzed completely, but only the first fifty lines are taken into account. It could happen then, that the preview appears looking fine, but the execution of the File Reader fails, when it reads the lines it didn't analyze. Thus it is recommended you check the settings, when you cut an analysis short.

Dialog Options

ASCII file location

Enter a valid file name or URL. When you press ENTER, the file is analyzed and the settings pre-set. You can also choose a previously read file from the dropdown list, or select a file from the "Browse..." dialog.

Preserve user settings

If checked, the checkmarks and column names/types you explicitly entered are preserved even if you select a new file. By default, the analyzer starts with fresh default settings for each new file location.

Rescan

If clicked, the file content is analyzed again. All settings are reset (unless the

KNIME Hub Search

Search workflows, nodes, and more...

```
graph LR; FR[File Reader] --> NC[Linear Correlation]; FR --> H[Histogram]; FR --> PD[Pie/Donut Chart]; FR --> N1[Node 1]; FR --> N2[Node 2]; FR --> N12[Node 12]; FR --> N13[Node 13]; FR --> NTS[Number To String]; FR --> P[Partitioning]; NTS --> V[V]; V --> N9[Node 9]; N9 --> DTL[Decision Tree Learner]; DTL --> N10[Node 10]; N10 --> SP[Scorer]; N10 --> EV[Evaluador]; P --> SRL[Random Forest Learner]; SRL --> N3[Node 3]; N3 --> RFL[Random Forest Predictor]; RFL --> N4[Node 4]; N4 --> SP; N4 --> EV; RFL --> N11[Node 11]; N11 --> SP; N11 --> EV; N12 --> L[Decision Tree Predictor]; L --> N5[Node 5]; N5 --> SP; N5 --> EV; N13 --> L; L --> N6[Node 6]; N6 --> SP; N6 --> EV; N2 --> SP; N2 --> EV; SP --> E[Scorer]; SP --> EV; E --> EV;
```

# Implementación: Code Oriented

## Deep Learning

### Assignment 1

The objective of this assignment is to learn about simple data curation practices, and familiarize you with some of the data we'll be reusing later.

This notebook uses the [notMNIST](#) dataset to be used with python experiments. This dataset is designed to look like the classic [MNIST](#) dataset, while looking a little more like real data: it's a harder task, and the data is a lot less 'clean' than MNIST.

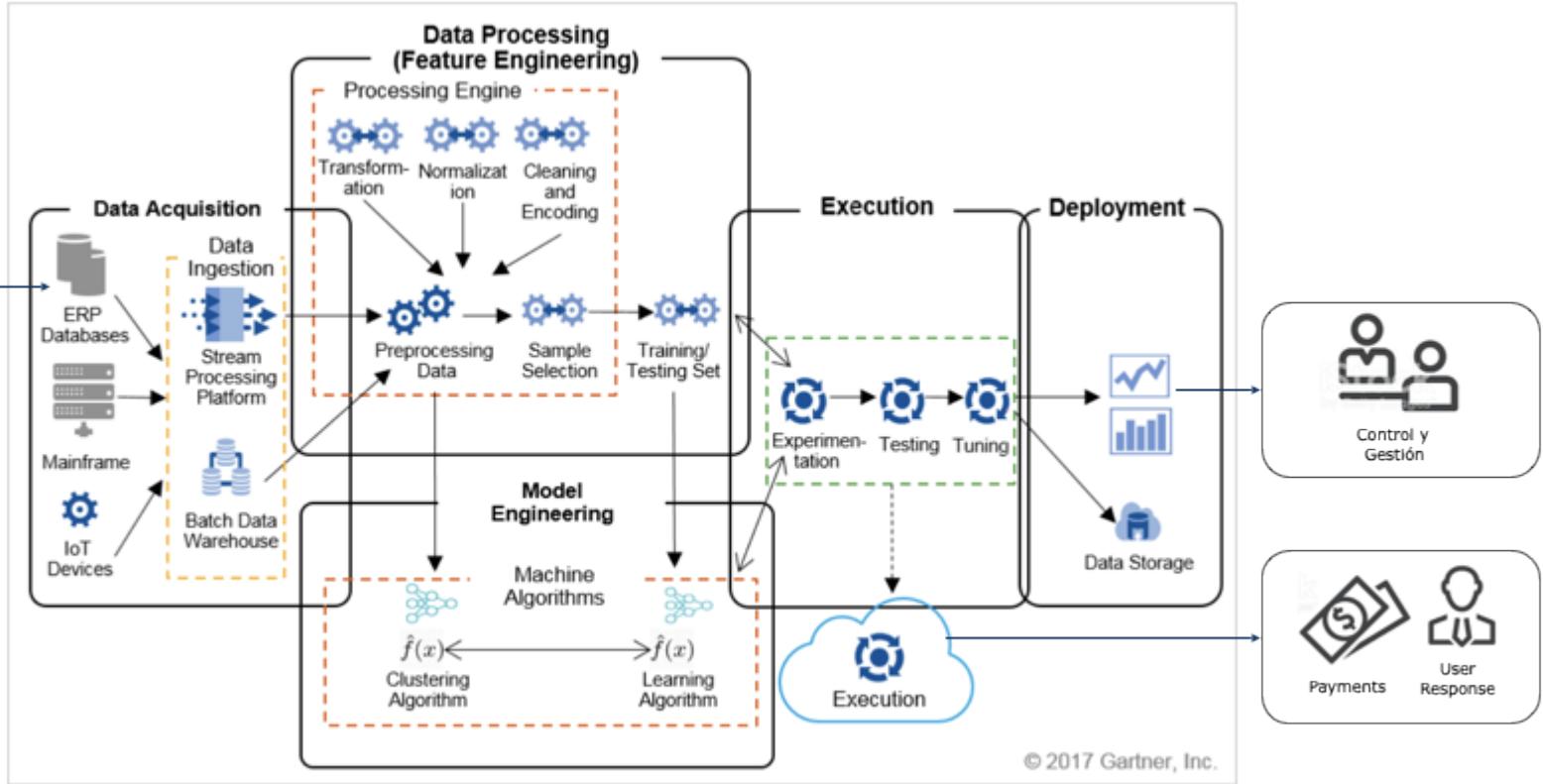
```
In [0]: # These are all the modules we'll be using later. Make sure you can import them
# before proceeding further.
import matplotlib.pyplot as plt
import numpy as np
import os
import sys
import tarfile
from IPython.display import display, Image
from scipy import ndimage
from sklearn.linear_model import LogisticRegression
from six.moves urllib.request import urlretrieve
from six.moves import cPickle as pickle
```

First, we'll download the dataset to our local machine. The data consists of characters rendered in a variety of fonts on a 28x28 image. The labels are limited to 'A' through 'J' (10 classes). The training set has about 500k and the testset 19000 labelled examples. Given these sizes, it should be possible to train models quickly on any machine.

```
In [0]: url = 'http://yaroslavvb.com/upload/notMNIST/'

def maybe_download(filename, expected_bytes):
    """Download a file if not present, and make sure it's the right size."""
    if not os.path.exists(filename):
        filename, _ = urlretrieve(url + filename, filename)
        statinfo = os.stat(filename)
        if statinfo.st_size == expected_bytes:
```

# AA: Infrastructure

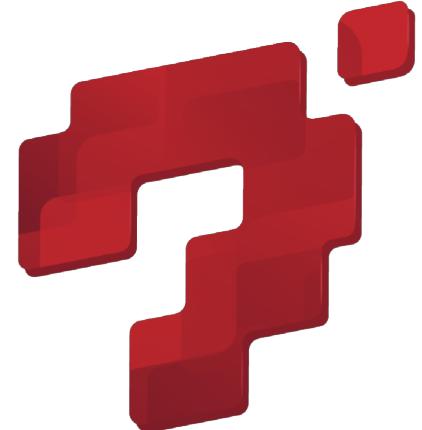


KNIME PLATFORM

# Conclusiones

- En que patrón de AA estamos? Orientado a la Programación? Orientado al Modelamiento?
- Necesitamos una **metodología** clara desde la iniciación de una idea de AA hasta la puesta en producción.
- Necesitamos medir cuantitativamente los éxitos del proyecto : ROI + Metodología.
- Necesitamos definir una **comunidad** y **prácticas locales** (No Deus Ex-Machina).

**¡Gracias!  
Preguntas?**



victor.toledo@coredevx.com

