

HOW DATA SCIENCE REALLY LOOKS LIKE



FAVIO
VÁZQUEZ
DATA SCIENTIST



CLOSTER

H₂O.ai

WHO AM I?

- Venezuelan
- Bachelor of Physics and Computer Engineer
- Master's at PCF-UNAM (Cosmology)
- Data Scientist at **H2O.ai**
- Faculty Member and **Emeritus** Instructor
- CEO and Chief Data Scientist at **Closter**
- Creator of **Ciencia y Datos**

Favio Vázquez Following ▾

Data scientist, physicist and computer engineer. Love sharing ideas, thoughts and contributing to Open Source in Machine Learning and Deep Learning ;).



Medium member since November 2018 · Editor of Ciencia y Datos · Top writer in Artificial Intelligence

2 Following 8K Followers ·

- Very active on LinkedIn ;)
- Editor International Journal of Business Analytics and Intelligence
- Writer at Towards Data Science, Heartbeat, Becoming Human and Planeta Chatbot
- Medium Top Writer in Business, AI and Technology

Contents

- Motivation
- Important concepts
- Data Science Timeline
- Creating a data product (theory)
- Learning data science
- **Data science problem and solution (Workshop)**

Motivation: Data science products are complex

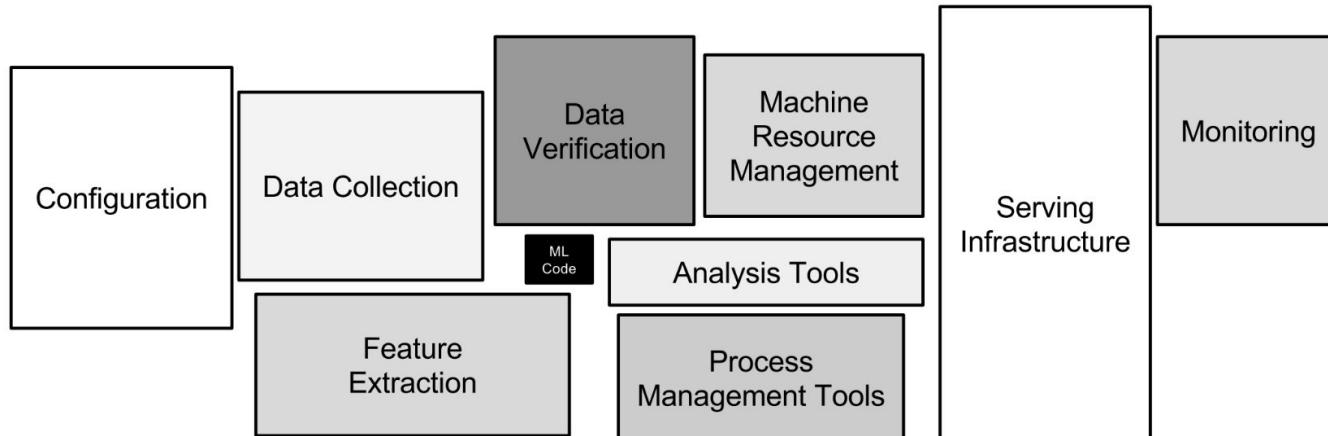
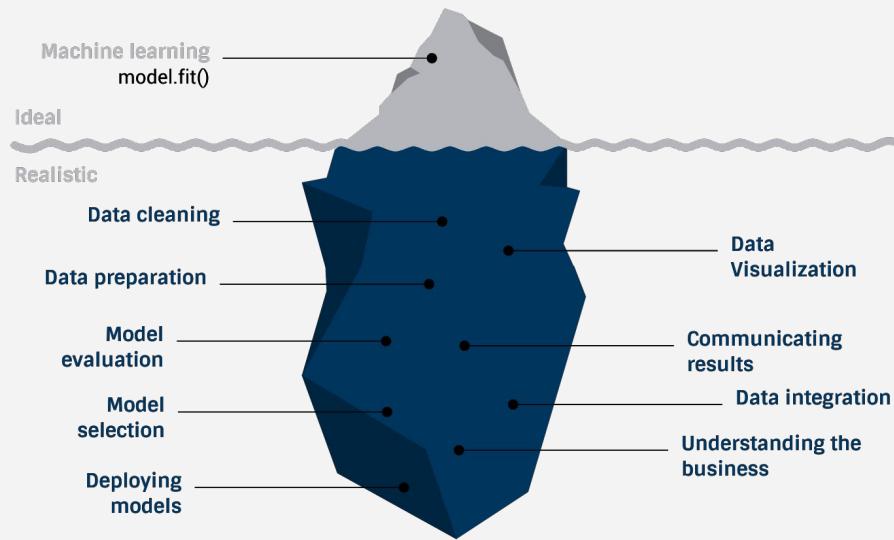


Figure 1: Only a small fraction of real-world ML systems is composed of the ML code, as shown by the small black box in the middle. The required surrounding infrastructure is vast and complex.

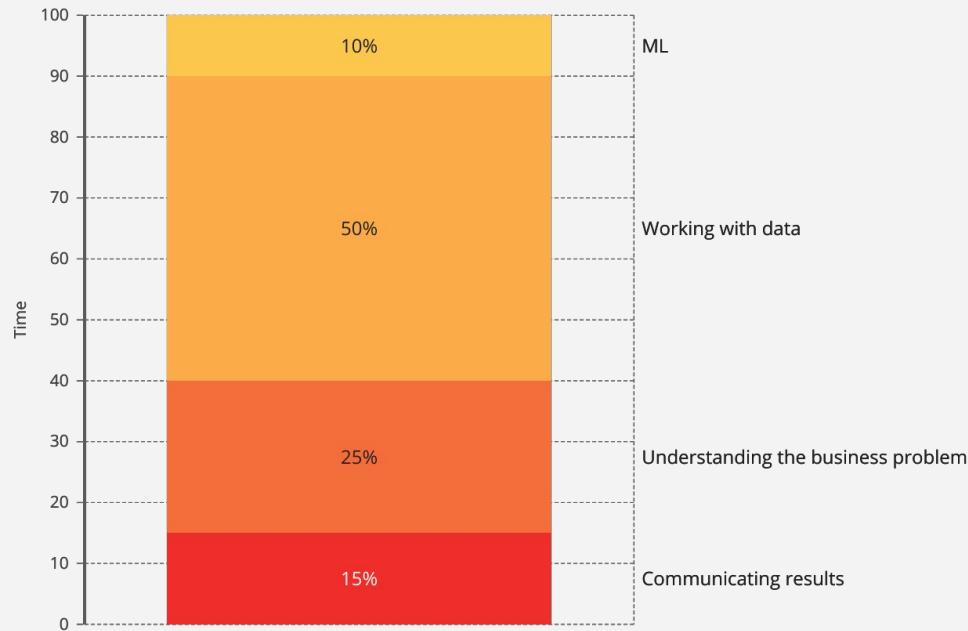
Motivation: Machine learning is not only `model.fit()`

What students think data science is



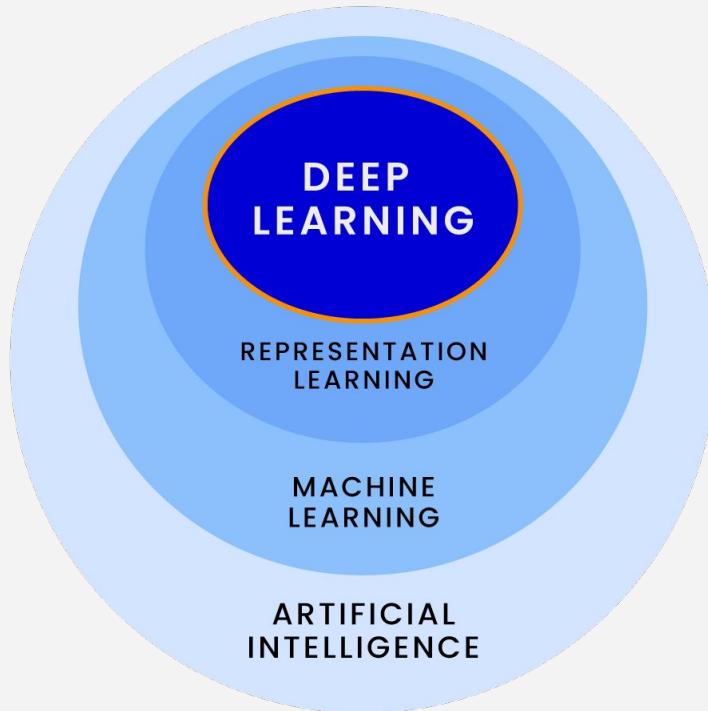
Motivation: How data science really looks like

How a data scientist time is spent



Important concepts

Panorama of AI

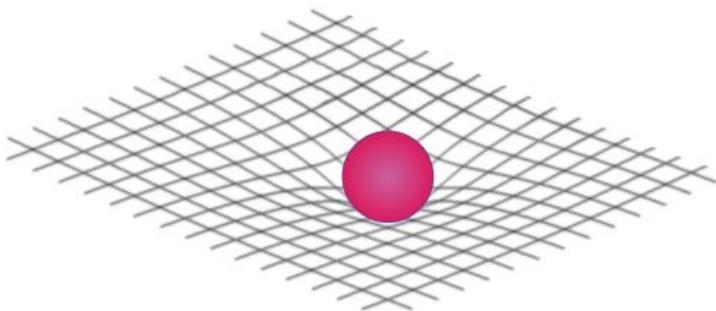


Artificial Intelligence (AI)



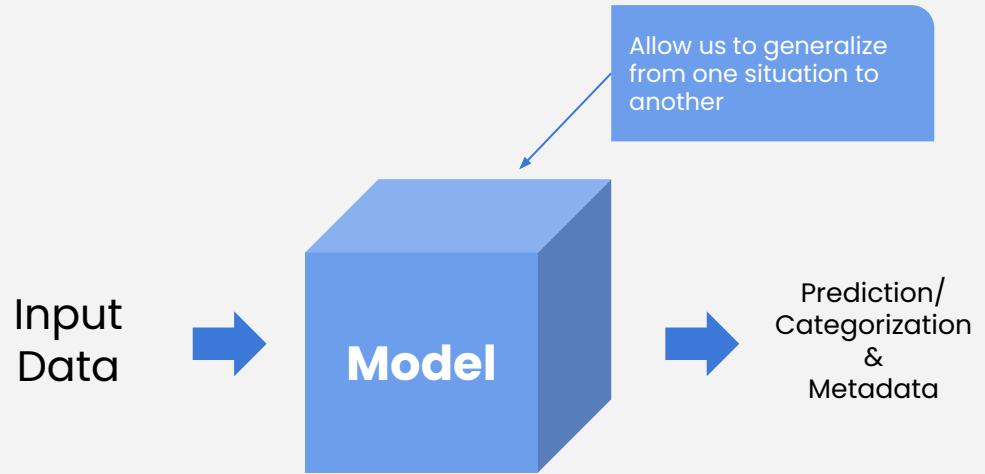
Automated
processes to
emulate human
skills and
capabilities with
computing.

Model

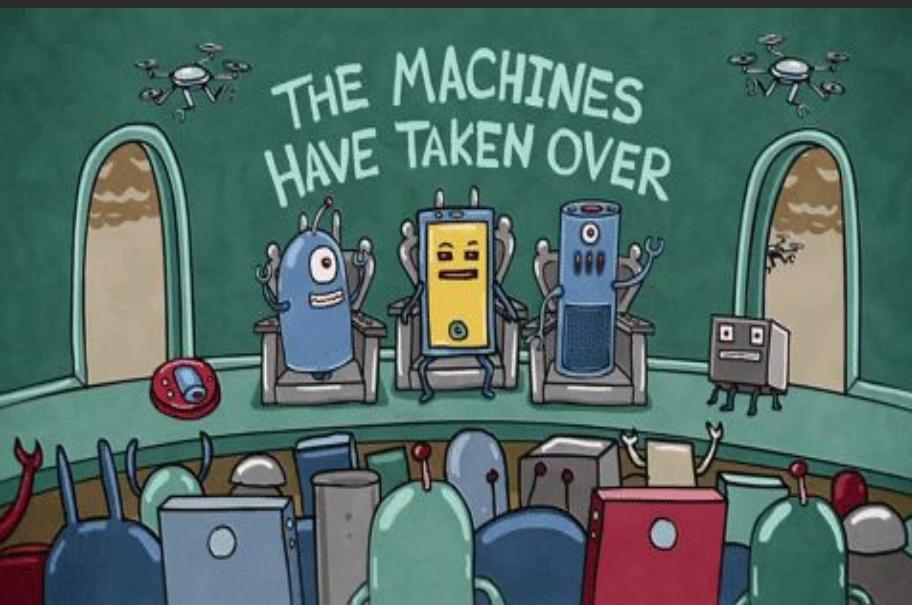


$$G_{\mu\nu} + \Lambda g_{\mu\nu} = \frac{8\pi G}{c^4} T_{\mu\nu}$$

Abstraction of reality to understand a process or phenomenon using mathematical tools.



Machine Learning

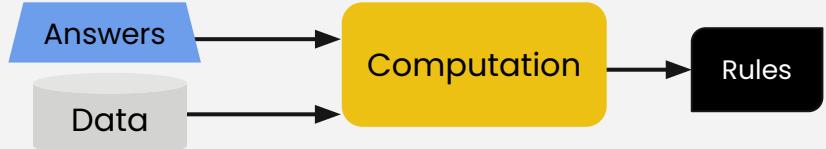


Discovery of patterns from data to predict the operation of a process or phenomenon using algorithms that are not explicitly programmed.

Traditional Programming



Machine Learning



methodology
sciences development effort areas
evaluation **support project**
software center engineering campus
collaboration space **director**
impact **working** career
students reproducibility
studio **science** groups
budget evaluation
institute faculty work tools statistics
university environment **scientists**
activities new domain
education **research**
graduate program



CLOSTER

H₂O.ai

THE ROLE OF DATA SCIENCE

Data science is the mediator on the path of using AI to impact business.



Perhaps the name "data science" is not very comfortable for everyone, but I will try to show that we can take advantage of it for now.



how does **data**
science emerge in
a company?

NATURE



...
...

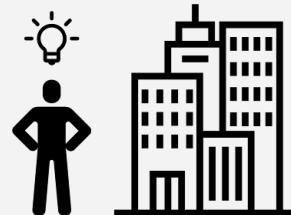


```
10100  
00101  
10100
```



IT PEOPLE

ILLUSTRATIVE



SUPERBOSS

We have a lot of
data, we need to
do something
about it

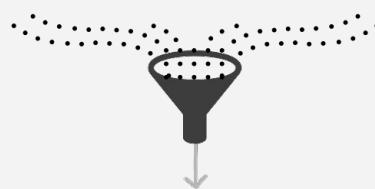


CLOSTER

H₂O.ai

ILLUSTRATIVE

NATURE



10100
00101
10100

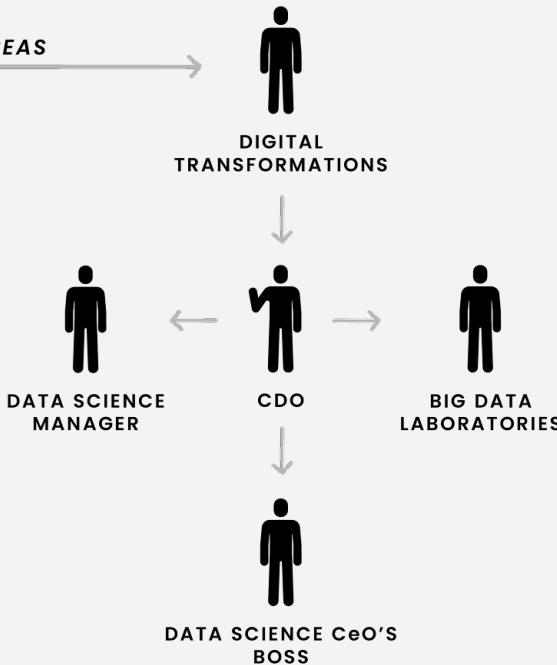


IT PEOPLE



SUPERBOSS

CREATION OF NEW AREAS



ILLUSTRATIVE

NATURE



10100
00101
10100



IT PEOPLE



SUPERBOSS



BUSINESS PROBLEMS



DIGITAL TRANSFORMATIONS



DATA SCIENCE
MANAGER



CDO



BIG DATA
LABORATORIES



TI PEOPLE



DATA SCIENCE CeO'S
BOSS

NATURE



...
...



10100
00101
10100



IT PEOPLE



EMAIL CAMPAIGNS



IMPROVE SATISFACTION



INCREASE PROFITS



CALL-CENTERS



GEOSPATIAL PROBLEMS



SUPPLY CHAINS



IMPROVE LOW METRICS



DATA SCIENCE
MANAGER



DATA SCIENTISTS



CDO



DATA SCIENCE CeO'S
BOSS



BIG DATA
LABORATORIES

DATA
ENGINEERS

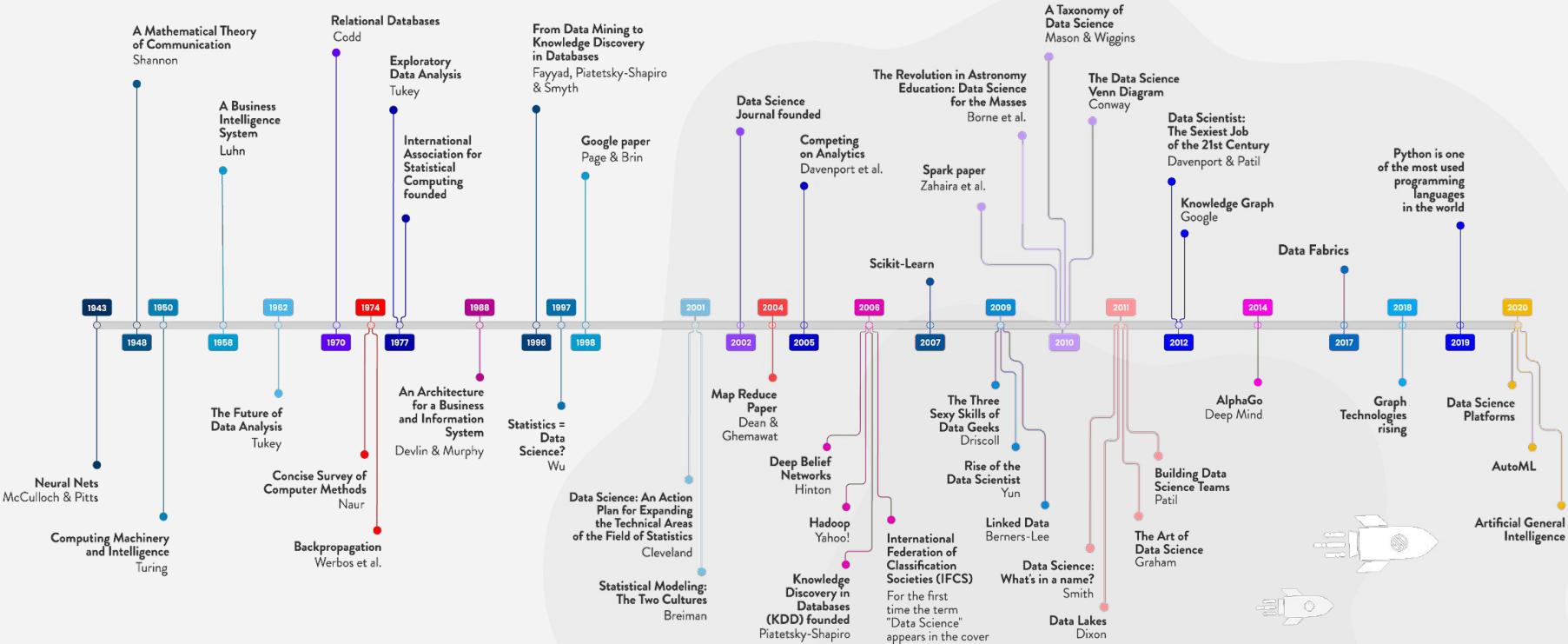


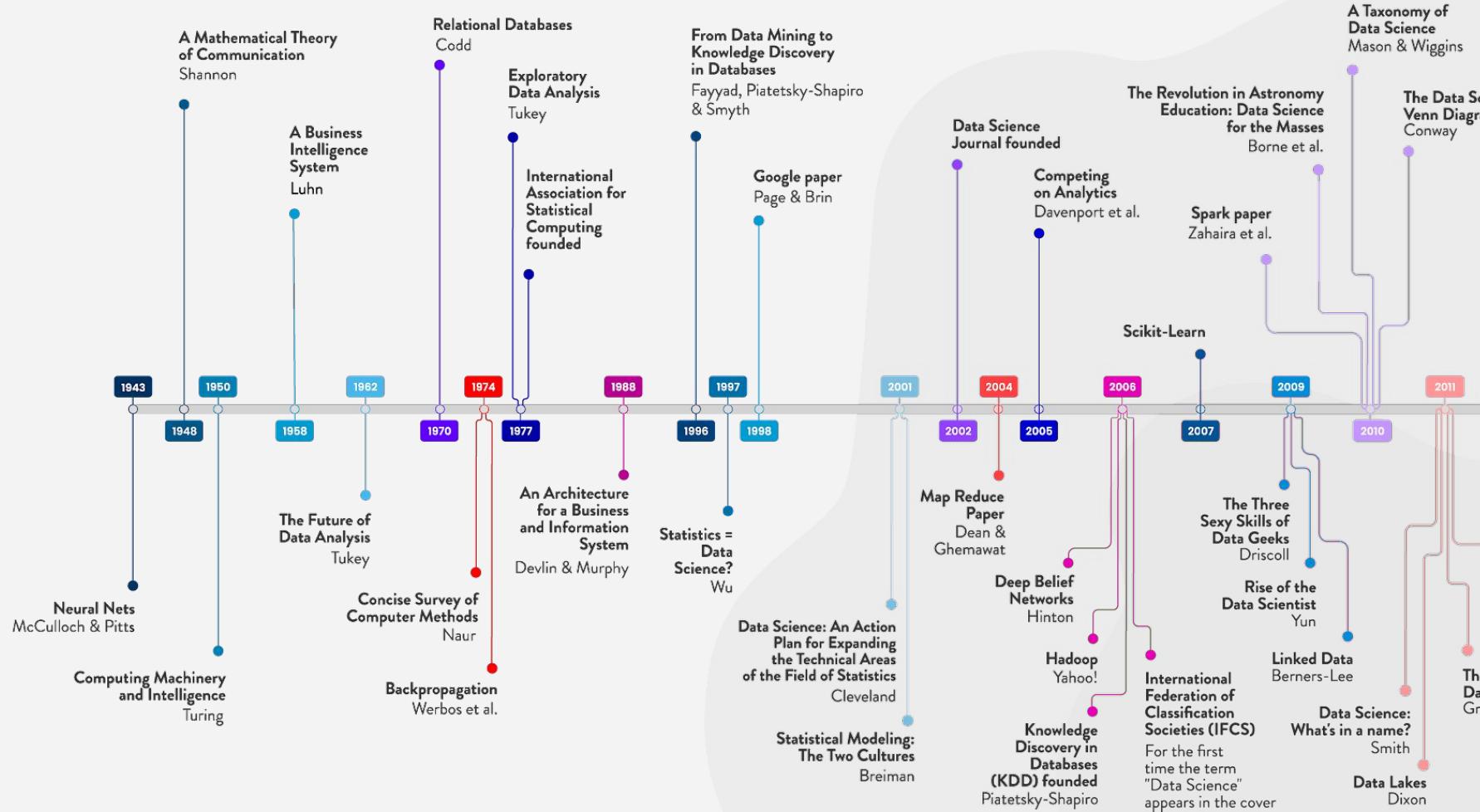
DATA
ARCHITECTS

DATA SCIENTISTS

ILLUSTRATIVE

DATA SCIENCE TIMELINE





From Data Mining to
Knowledge Discovery
in Databases

Fayyad, Piatetsky-Shapiro
& Smyth

Google paper
Page & Brin

Data Science
Journal founded

Competing
on Analytics
Davenport et al.

The Revolution in Astronomy
Education: Data Science
for the Masses
Borne et al.

Spark paper
Zaharia et al.

A Taxonomy of
Data Science
Mason & Wiggins

The Data Science
Venn Diagram
Conway

Data Scientist:
The Sexiest Job
of the 21st Century
Davenport & Patil

Python is one
of the most used
programming
languages
in the world

1996
1997
1998

Statistics =
Data
Science?
Wu

Data Science: An Action
Plan for Expanding
the Technical Areas
of the Field of Statistics
Cleveland

Statistical Modeling:
The Two Cultures
Breiman

Map Reduce
Paper
Dean &
Ghemawat

Deep Belief
Networks
Hinton

Hadoop
Yahoo!

Knowledge
Discovery in
Databases
(KDD) founded
Piatetsky-Shapiro

Scikit-Learn

The Three
Sexy Skills of
Data Geeks
Driscoll

Rise of the
Data Scientist
Yun

Linked Data
Berners-Lee

International
Federation of
Classification
Societies (IFCS)
For the first
time the term
"Data Science"
appears in the cover

Data Lakes
Dixon

Data Science:
What's in a name?
Smith

The Art of
Data Science
Graham

Building Data
Science Teams
Patil

AlphaGo
Deep Mind

Data Fabrics

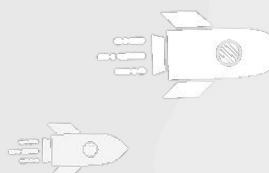
Graph
Technologies
rising

Data Science
Platforms

Artificial General
Intelligence

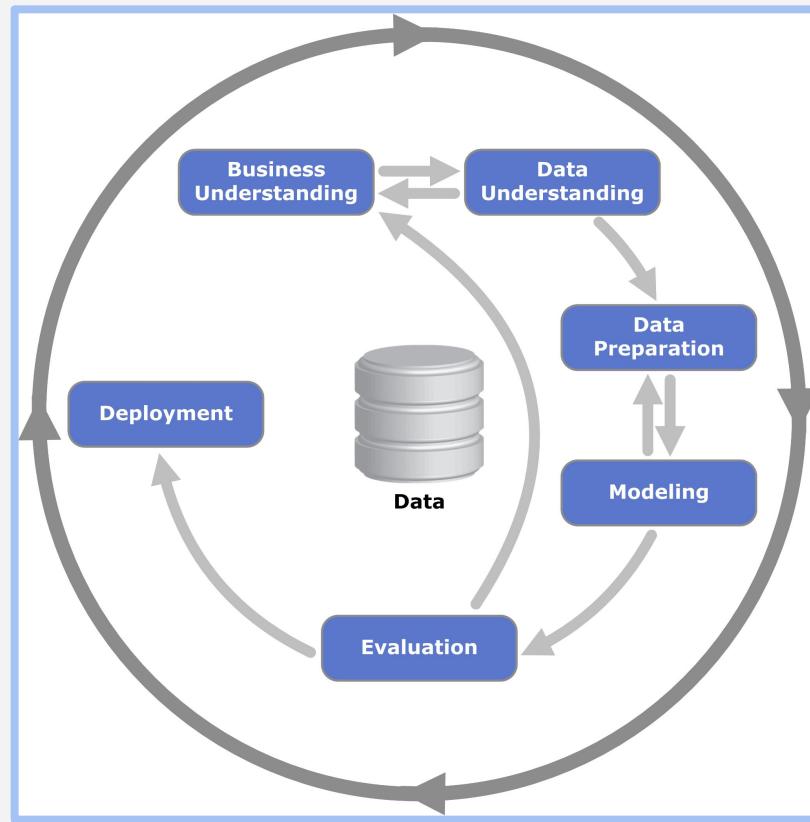
CLOSTER

H₂O.ai

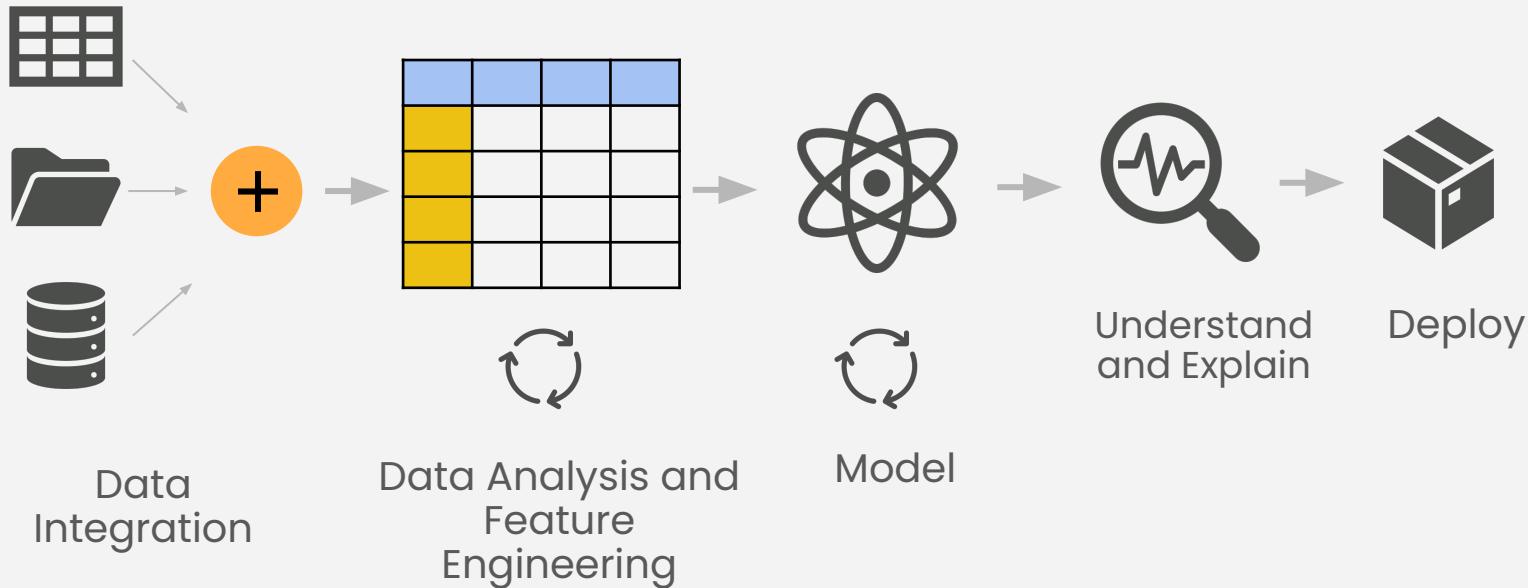


Data Science Products

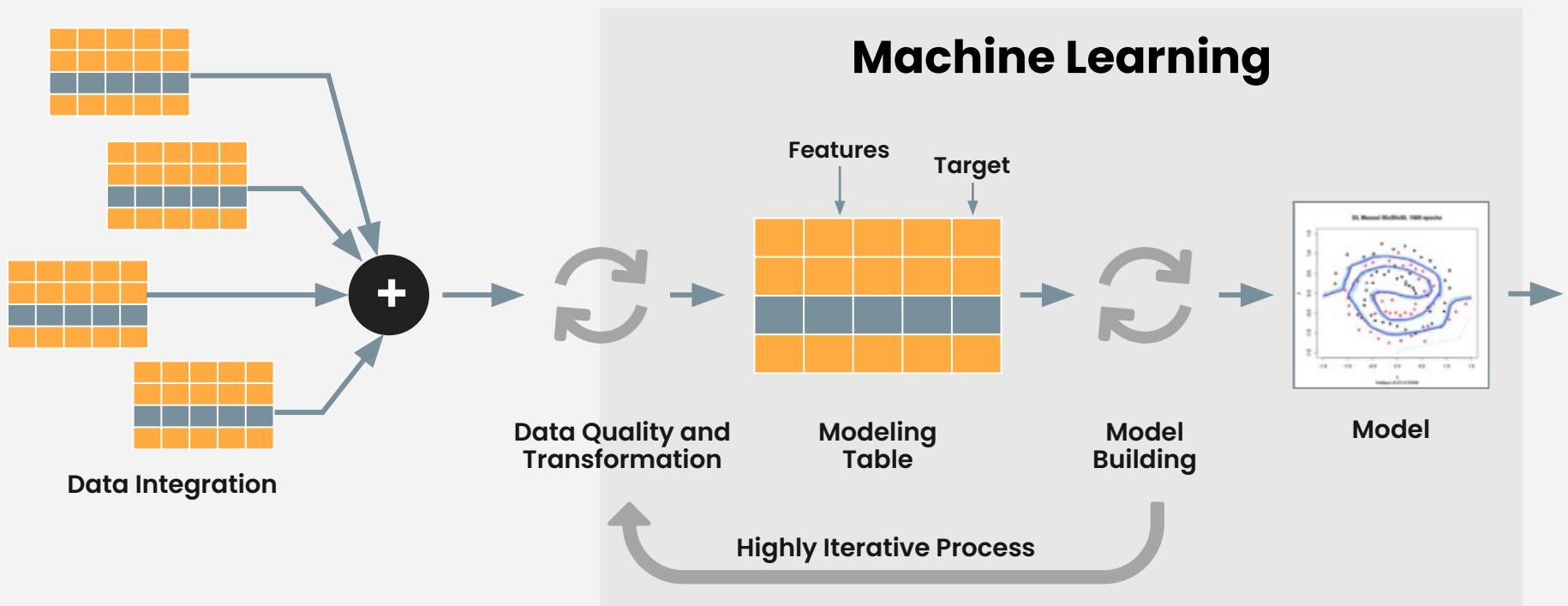
CRISP-DM (One of the fathers of DS)



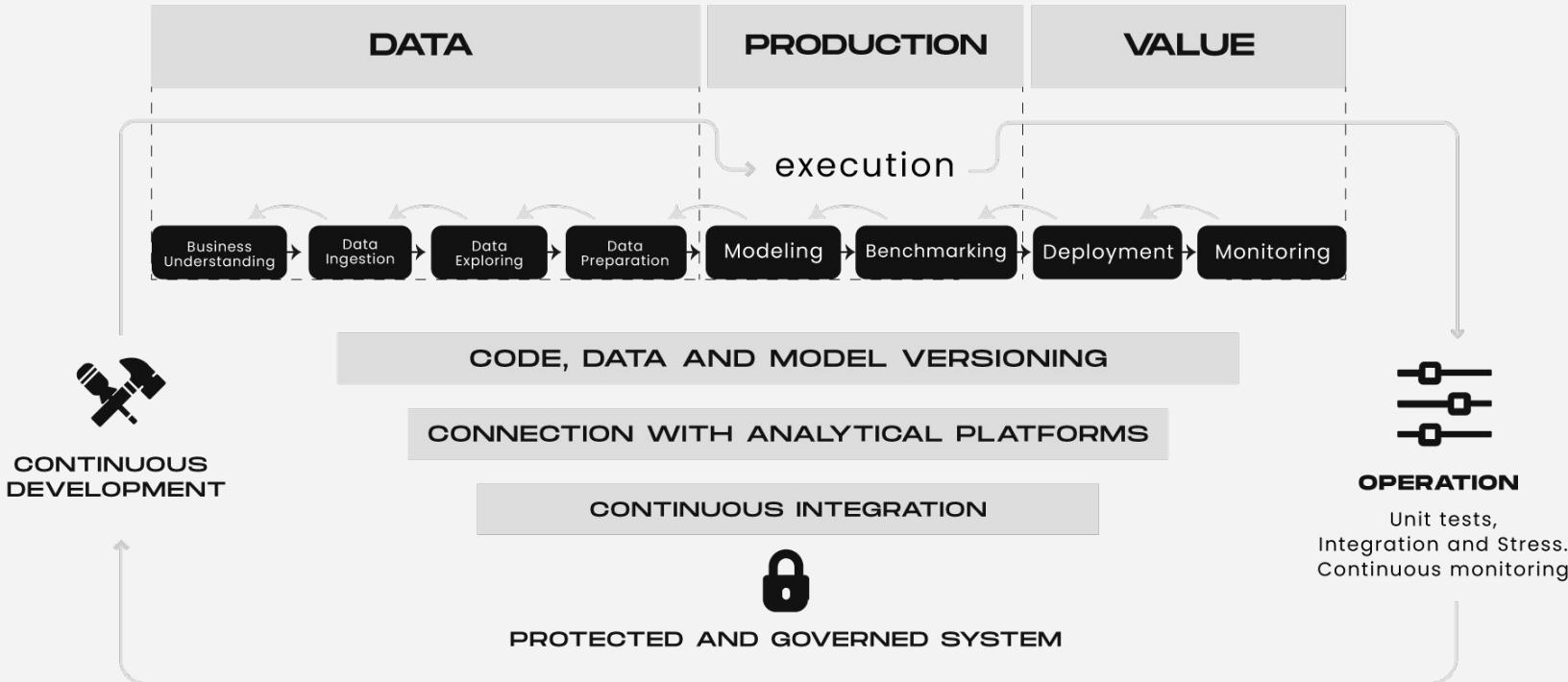
Data Science Workflow



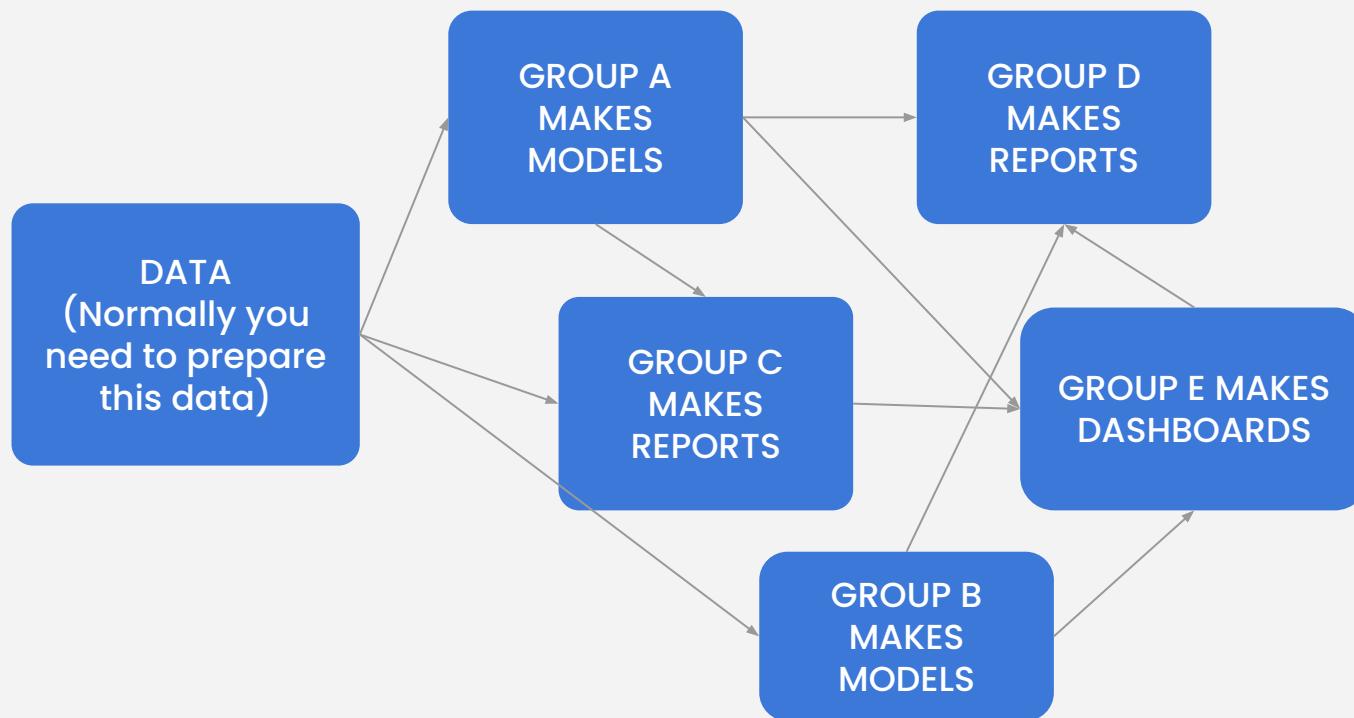
Data Science Workflow



DATA SCIENCE METHODOLOGY



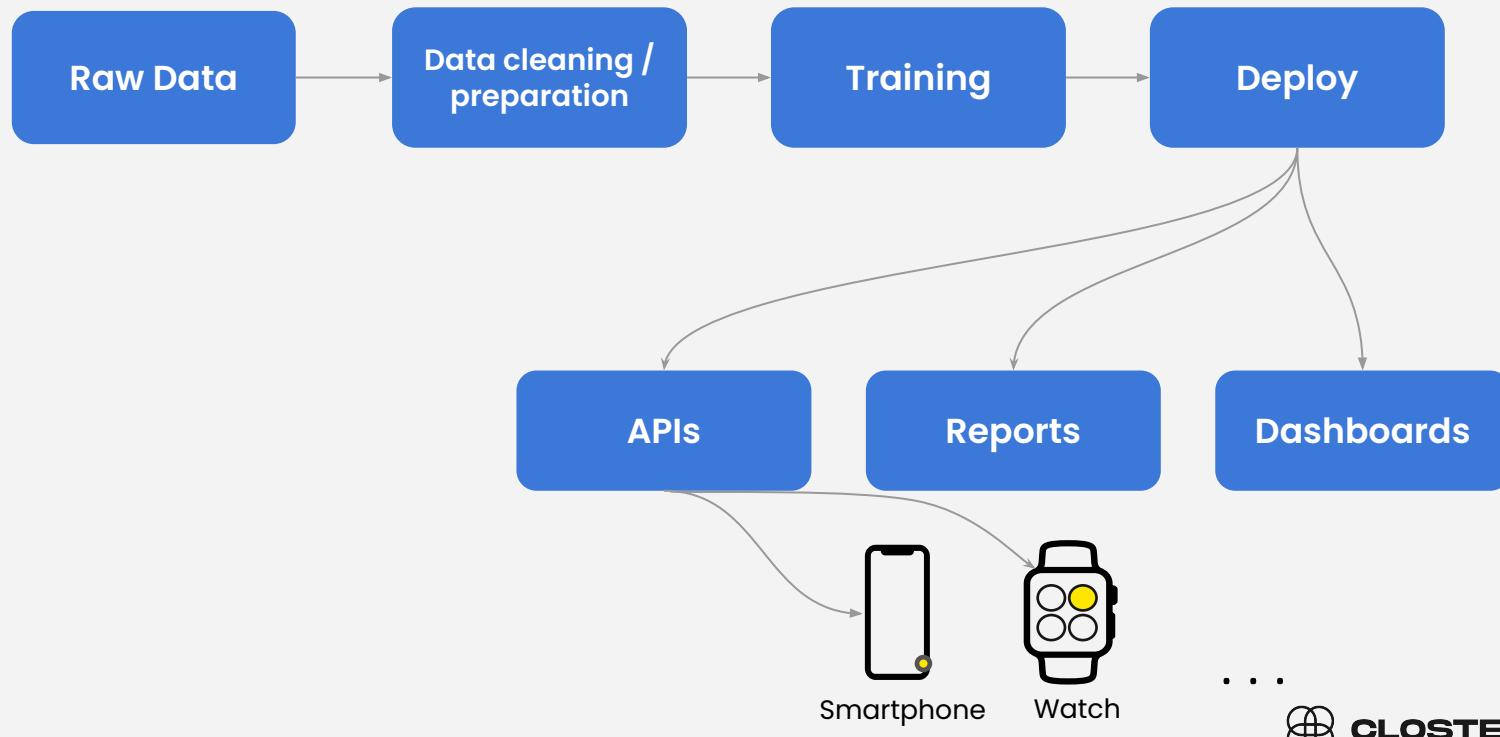
Data science products are complex



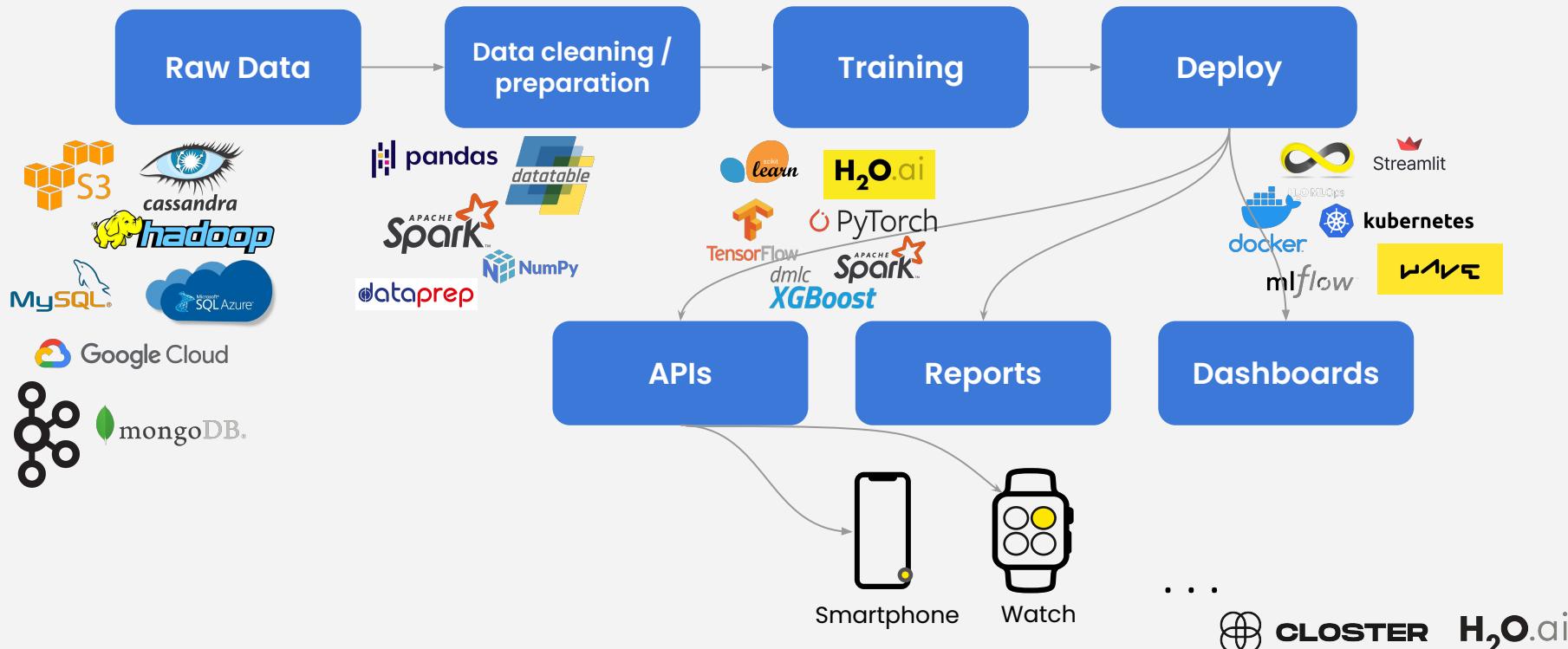
CLOSTER

H₂O.ai

Data science products are complex



Data science products are complex



Data science products are complex



Demo links

<https://github.com/FavioVazquez/dsgo-workshop-2021>

<https://aquarium.h2o.ai/>