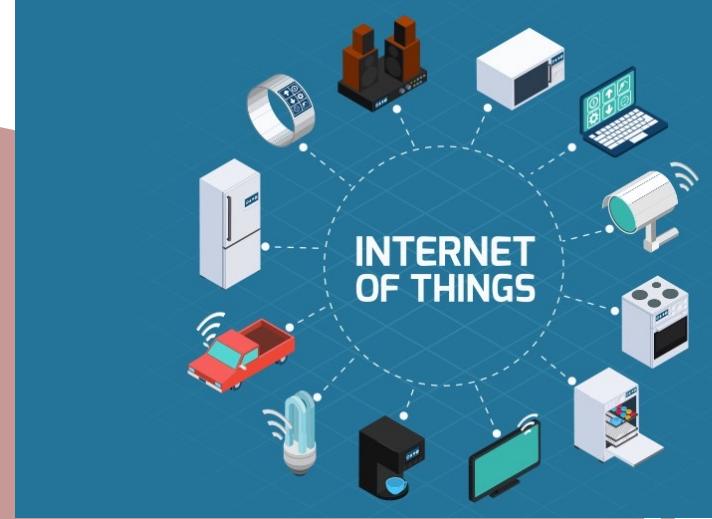


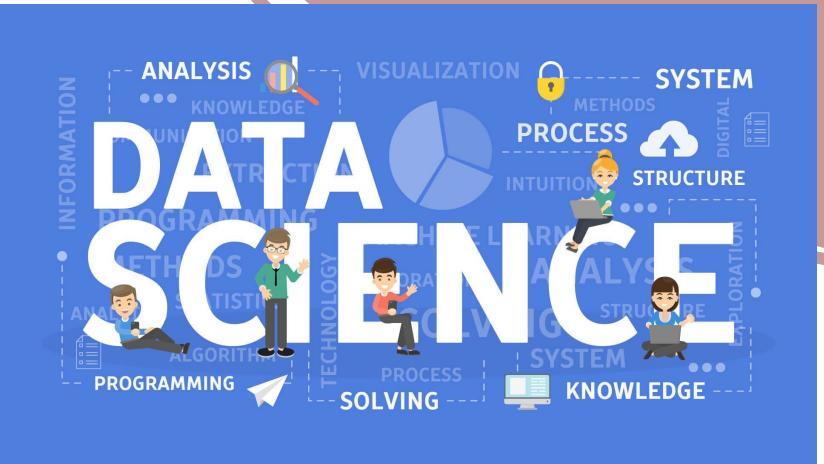
Data Science

André Maletzke

2021



Data + Science



Data

- 2.5 quintillion bytes of data (every day)
2,500,000,000,000,000
- This amount of data created will continue to grow
- 3.7 billion humans use the internet every day (src: Forbes.com)
- Google processes more than 40,000 searches every second
 - 3.5 billion searches each day



Data is the new oil!

"A NEW commodity spawns a lucrative, fast-growing industry, prompting antitrust regulators to step in to restrain those who control its flow. A century ago, the resource in question was oil. Now similar concerns are being raised by the giants that deal in data, the oil of the digital era."

Source: theeconomist.com

Leaders

May 6th 2017 edition >

Regulating the internet giants

The world's most valuable resource is no longer oil, but data

The data economy demands a new approach to antitrust rules



David Parkins

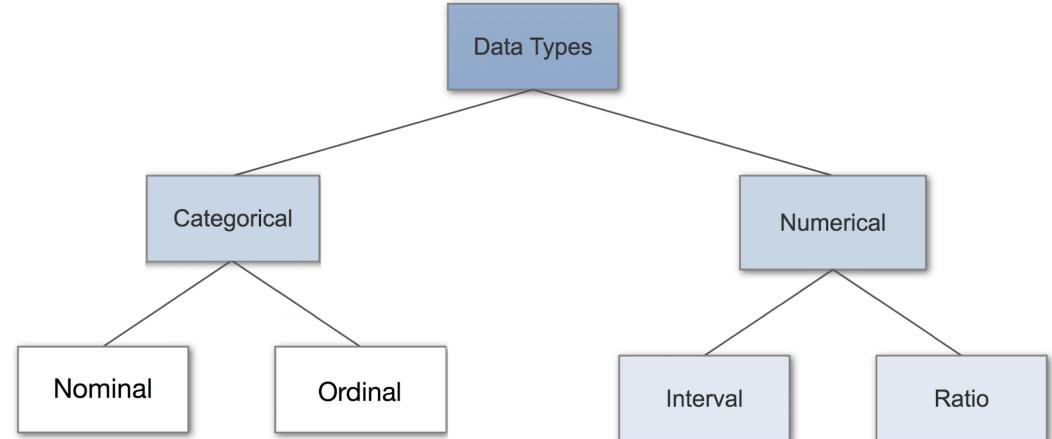
What is data?



WIKIPEDIA
The Free Encyclopedia

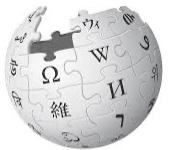
Data are characteristics or information, usually numerical, that are collected through observation.^[1] In a more technical sense, data are a set of values of qualitative or quantitative variables about one or more persons or objects

Although the terms "data" and "information" are often used interchangeably, these terms have distinct meanings. In some popular publications, data are sometimes said to be transformed into information when they are viewed in context or in post-analysis.



Sciencia

Science (from the [Latin](#) word *scientia*, meaning "knowledge")^[1] is a systematic enterprise that [builds](#) and [organizes](#) [knowledge](#) in the form of [testable explanations](#) and [predictions](#) about the [universe](#).



WIKIPEDIA
The Free Encyclopedia

Scientific method

Em sentido estrito, ciência refere-se ao [sistema](#) de adquirir conhecimento baseado no [método científico](#) bem como ao corpo organizado de conhecimento conseguido através de tais [pesquisas](#).

Scientific method

Scientific
method is the
base of science!

Every baby
knows the
scientific method!



2 Form a hypothesis.



1

Make an observation.



3 Perform the experiment.



4

Analyze the data.



5

Report your findings.



6

Invite others to reproduce the results.

Data driven decisions... Historical notes

Dr. Ignaz Semmelweis



Puerperal fever was common in mid-19th-century hospitals and often fatal

1846 – He was appointed assistant to Professor Johann Klein (First Obstetrical Clinic of the Vienna General Hospital)

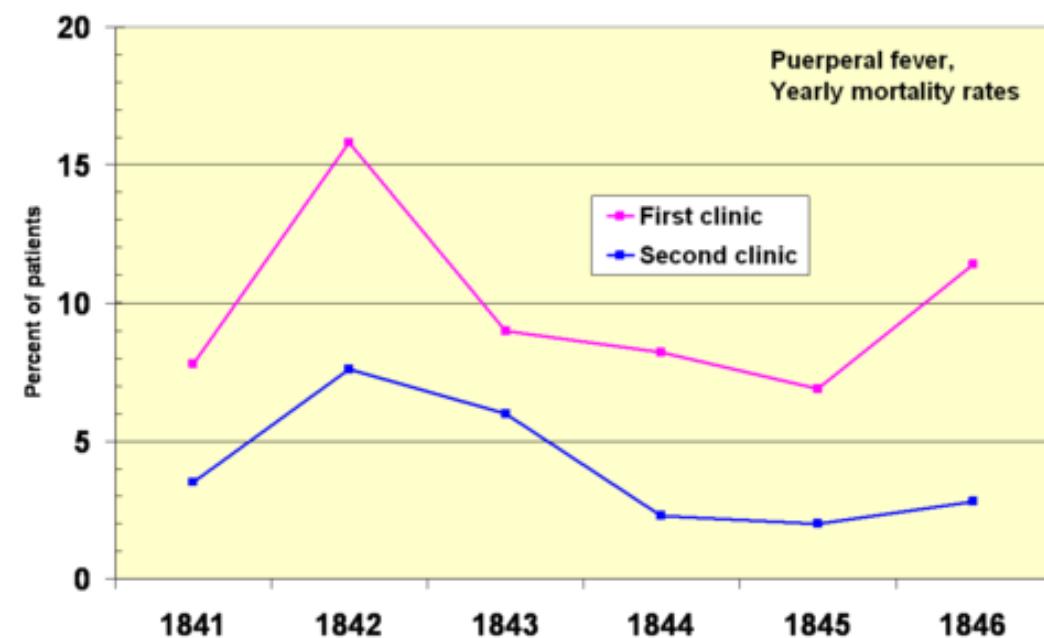
One of his duties was check the medical **records**

Observation:

The First Clinic had an average maternal mortality rate of about **10%** due to puerperal fever. The Second Clinic's rate was considerably lower, averaging less than **4%**.

Data driven decisions... Historical notes

Why?



The two clinics used almost the same techniques

He excluded "overcrowding" as a cause (Second Clinic was always more crowded and yet the mortality was lower)

He eliminated climate as a cause because the climate was the same.

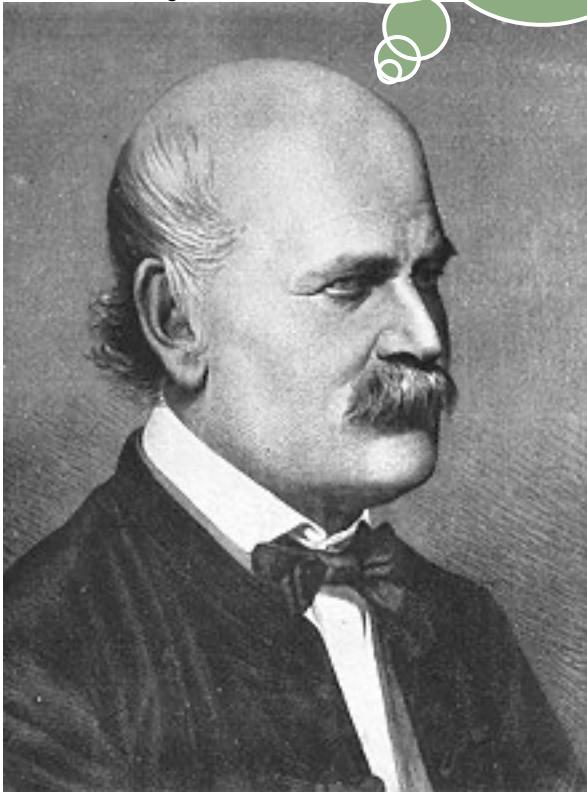
The only major difference was the individuals who worked there:

- (1) The First Clinic was the teaching service for medical students
- (2) The Second Clinic had been selected in 1841 for the instruction of midwives only

Data decisions... Historical notes

Theory of
cadaverous
poisoning

Dr. Ignaz Semmelweis



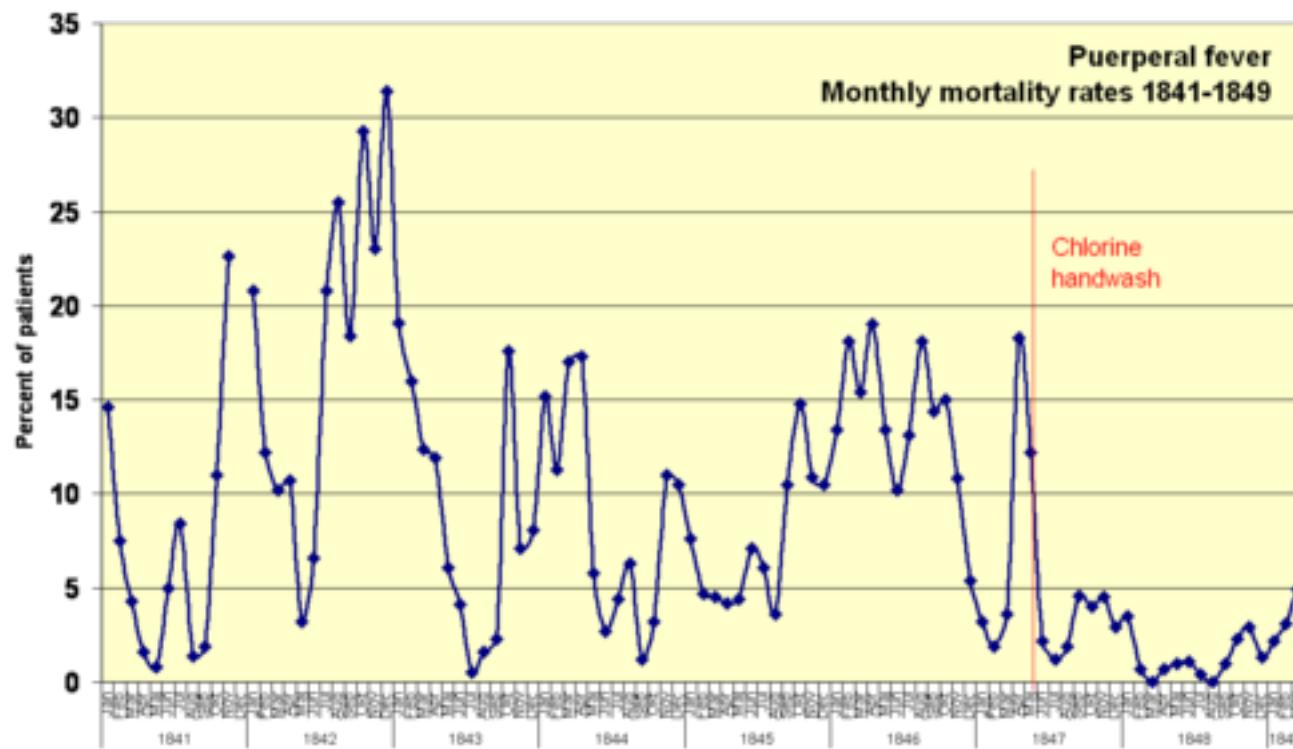
1847 - The death of Semmelweis' good friend (Jakob Kolletschka) - who had been accidentally poked with a student's scalpel while performing a *post mortem* examination

Semmelweis immediately proposed a connection between cadaveric contamination and puerperal fever

Dr. Semmelweis proposed that the medical students carried "cadaverous particles" on their hands

He instituted a policy of using a solution of chlorinated lime for washing hands between autopsy work and the examination of patients.

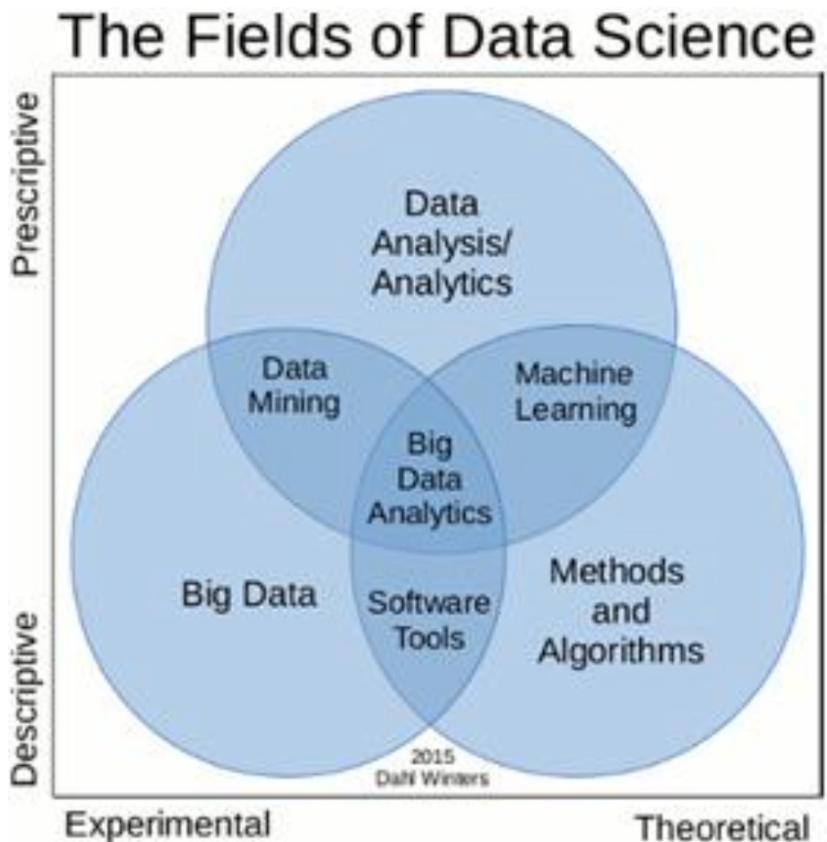
Data driven decisions... Historical notes



The result was the mortality rate in the First Clinic declined 90%

Source: en.wikipedia.org

What is Data Science?



The purpose of computing is insight, not numbers.

– Richard W. Hamming

Like any emerging field, it hasn't been completely defined yet ...

Steven S. Skiena

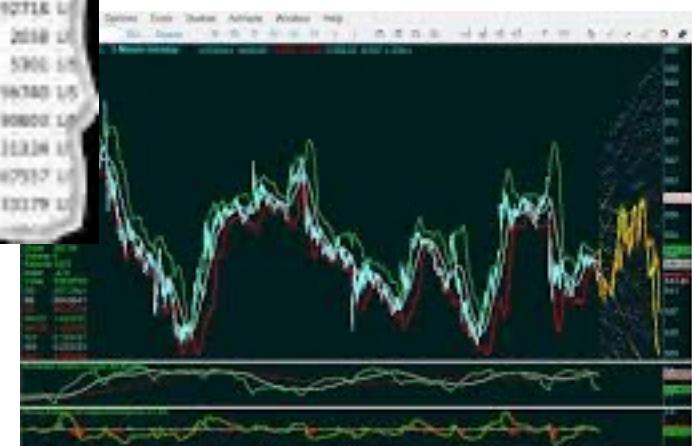
Data Science



id	name	sex	age	height	weight	date	status	city	state	zipCode
1	John	Male	25	5'10"	180	2023-01-01	Active	New York	NY	10001-0001
2	Mary	Female	28	5'5"	150	2023-01-01	Active	Los Angeles	CA	90001-0001
3	Sarah	Female	30	5'7"	165	2023-01-01	Active	Chicago	IL	60601-0001
4	William	Male	32	5'9"	190	2023-01-01	Active	Boston	MA	02101-0001
5	Emily	Female	27	5'6"	170	2023-01-01	Active	Houston	TX	77001-0001
6	David	Male	29	5'11"	195	2023-01-01	Active	Seattle	WA	98101-0001
7	Olivia	Female	26	5'4"	145	2023-01-01	Active	Phoenix	AZ	85001-0001
8	Isabella	Female	24	5'3"	135	2023-01-01	Active	Philadelphia	PA	19101-0001
9	Lucas	Male	22	5'7"	175	2023-01-01	Active	Dallas	TX	75201-0001
10	Charlotte	Female	21	5'2"	125	2023-01-01	Active	Minneapolis	MN	55401-0001
11	Matthew	Male	23	5'10"	185	2023-01-01	Active	Portland	OR	97201-0001
12	Ava	Female	20	5'5"	155	2023-01-01	Active	San Antonio	TX	78201-0001
13	Noah	Male	19	5'7"	170	2023-01-01	Active	Austin	TX	78701-0001
14	Madison	Female	18	5'4"	140	2023-01-01	Active	San Diego	CA	92101-0001
15	Benjamin	Male	17	5'6"	160	2023-01-01	Active	Seattle	WA	98101-0001
16	Harper	Female	16	5'3"	130	2023-01-01	Active	Phoenix	AZ	85001-0001
17	Elijah	Male	15	5'2"	120	2023-01-01	Active	Philadelphia	PA	19101-0001
18	Charlotte	Female	14	5'0"	105	2023-01-01	Active	Minneapolis	MN	55401-0001
19	Levi	Male	13	5'1"	115	2023-01-01	Active	Portland	OR	97201-0001
20	Amelia	Female	12	5'3"	130	2023-01-01	Active	San Antonio	TX	78201-0001
21	Lucas	Male	11	5'5"	150	2023-01-01	Active	Austin	TX	78701-0001
22	Isabella	Female	10	5'2"	120	2023-01-01	Active	San Diego	CA	92101-0001
23	Noah	Male	9	5'7"	170	2023-01-01	Active	Seattle	WA	98101-0001
24	Madison	Female	8	5'4"	140	2023-01-01	Active	Phoenix	AZ	85001-0001
25	Benjamin	Male	7	5'6"	160	2023-01-01	Active	Philadelphia	PA	19101-0001
26	Harper	Female	6	5'3"	130	2023-01-01	Active	Minneapolis	MN	55401-0001
27	Elijah	Male	5	5'2"	120	2023-01-01	Active	Portland	OR	97201-0001
28	Charlotte	Female	4	5'0"	105	2023-01-01	Active	San Antonio	TX	78201-0001
29	Levi	Male	3	5'1"	115	2023-01-01	Active	Austin	TX	78701-0001
30	Amelia	Female	2	5'3"	130	2023-01-01	Active	San Diego	CA	92101-0001
31	Noah	Male	1	5'5"	150	2023-01-01	Active	Seattle	WA	98101-0001



For the basis of the grande statue de Zouave, a Olympic, Phidias and represented by Dorian Diogenes. Based in Seville (Spain) and in Los Angeles (USA).
Dies ist ein Blindtext. An then last text.
vie
set
On
Gra
kau
les
Ma
fr
for
p
the
environments that are capab
contribution of this paper is to propose a new method for each object using weight of each pixel. In the pre
processing step, the algorithm calculates the weight of each pixel X, Y, Z and RGB value of each data point and then it
P each data point's normal vector using the point's neighbor. After
preprocessing, our algorithm calculates weights for each data point; each weight indicates membership. Rounding is done
objects of the scene.



Data Science



- Missing values
- Inconsistencies
- Redundancies
- Data from different sources
- Different formats

“...people usually spend about 60–70% of their time just on gathering and cleaning the data.”

Data Science



Questions about the problem

- Does the past represent the future?
- What do I want to model?
- How will the model be used?
- What data do I need? Or What data do I have?
- How hard to get the data?

Data Science problems

Classification

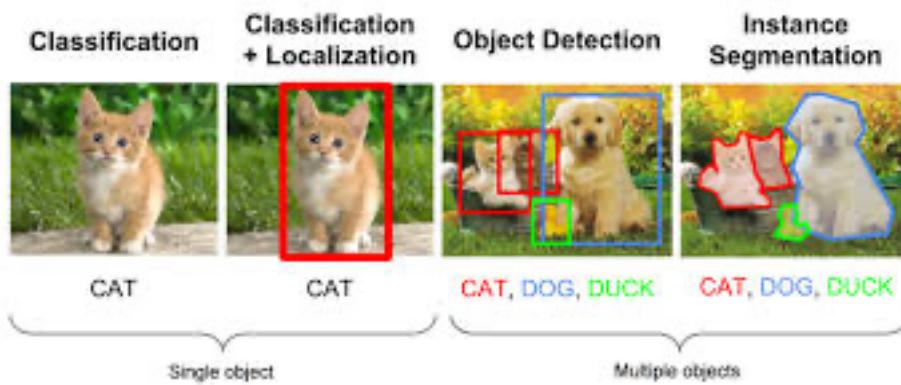


Image classification

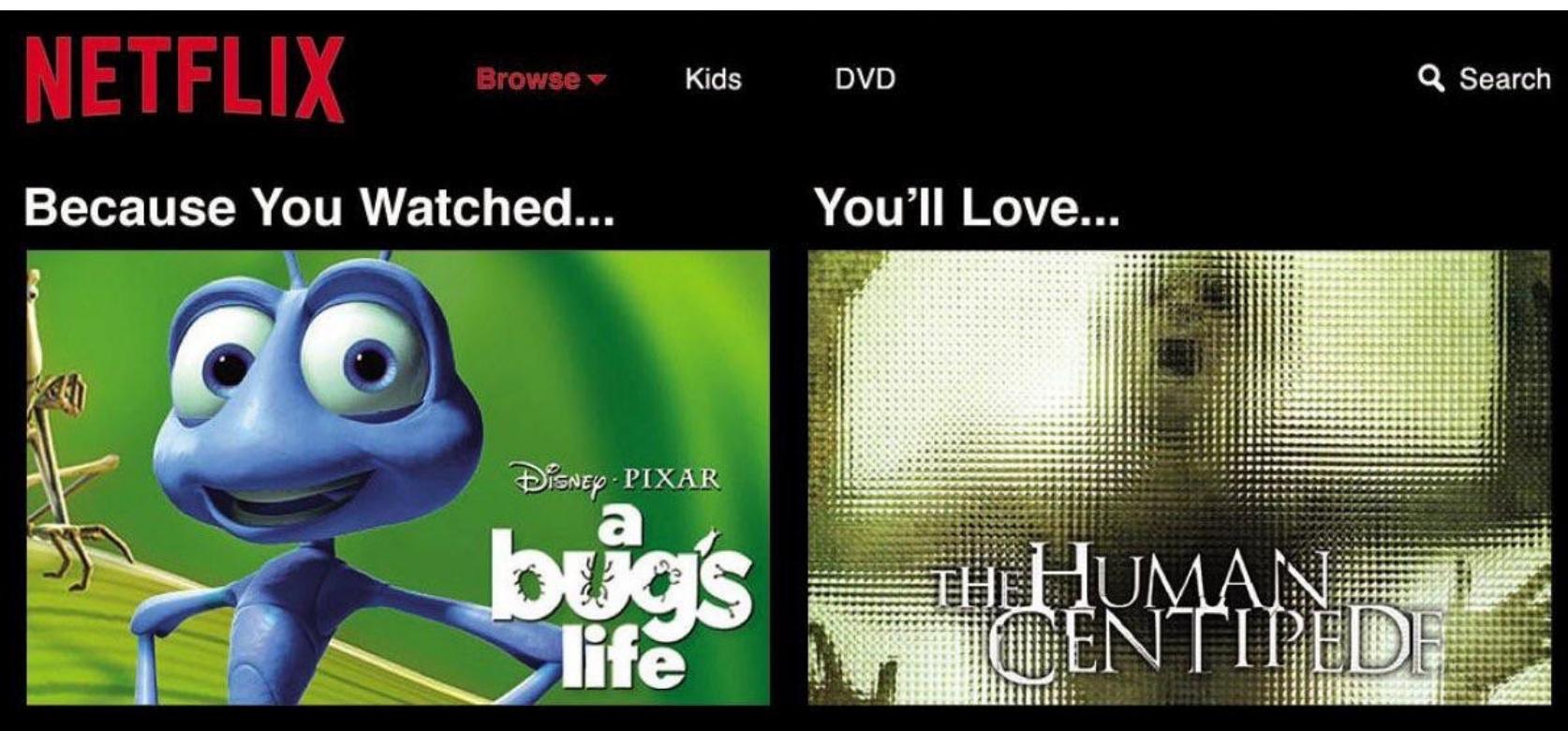
Regression



Stock price prediction

Data Science problems

Clustering or ranking



How (and when) do we best present that movie recommendation to the user in a way that maximizes viewership and monthly subscriber loyalty?

Data Science problems



Counting cars from aerial images.

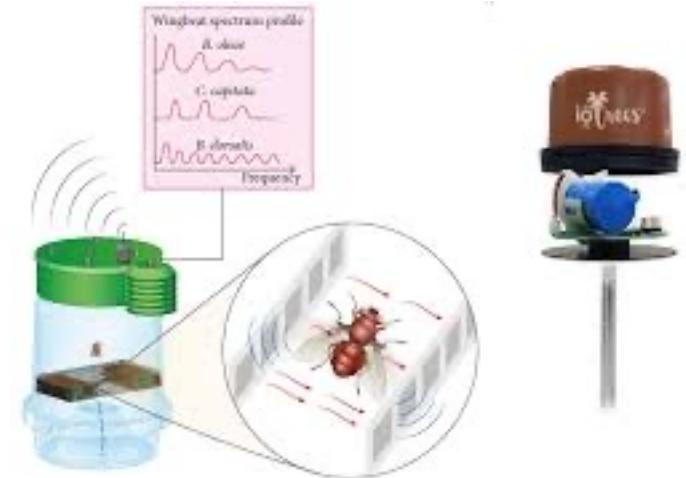
Source: Amar et al., (2019)

Counting

**Regression
Classification
Clustering
Counting**



Precision agriculture technologies



Insect pests

**Classification
Counting**

Data Science problems

Counting / Quantification

Positive ratings



Counting consumers' opinion (Gao and Sebastiani, 2015)

Where does the data come from?



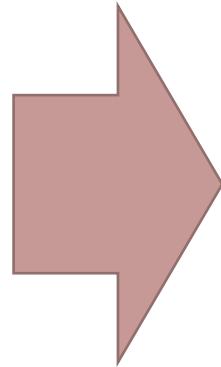
How does the data look like?

- Rating
 - Reviews
 - Others

Data Science problems



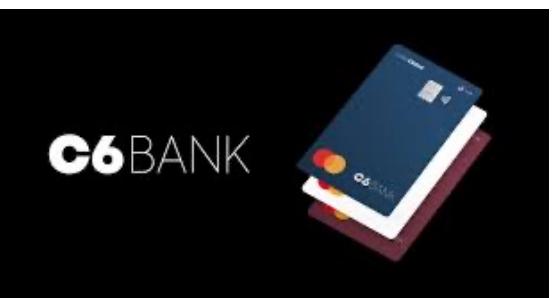
Estimate the prevalence of support for different political candidates (Hopkins and King, 2010)



Data Science problems



Credit score prediction



Who will get the credit card?

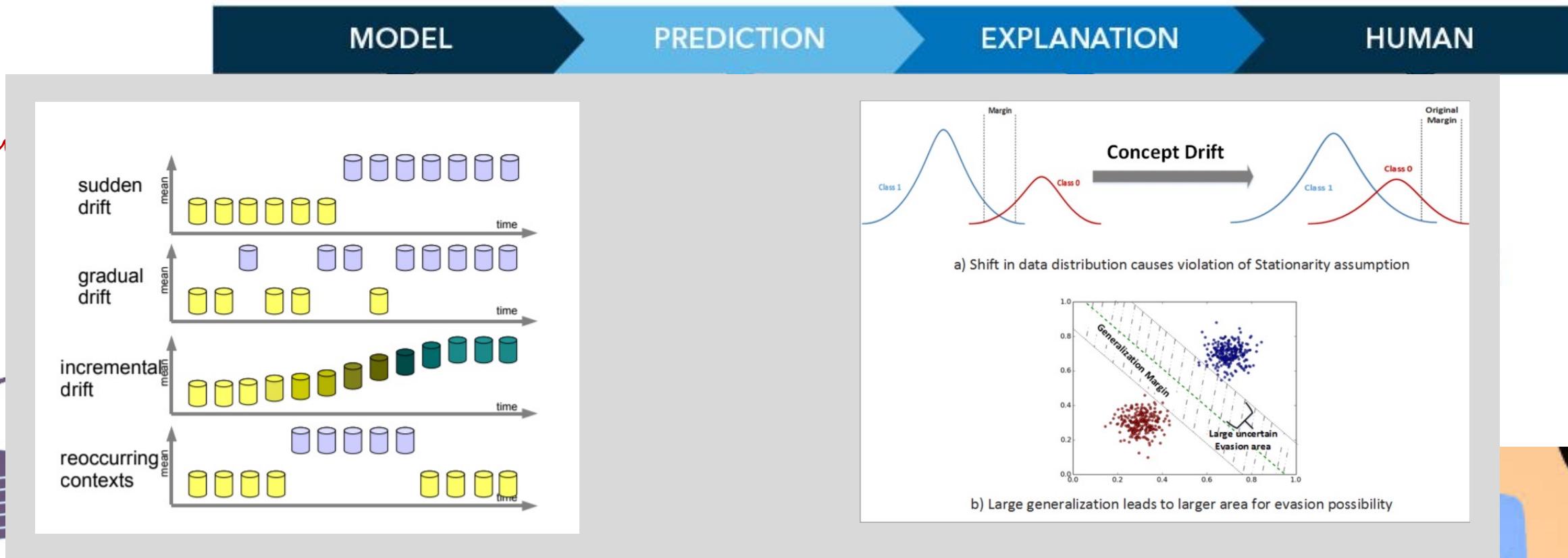
Who will get a bank loan?

Etc.



Data Science problems

Sec

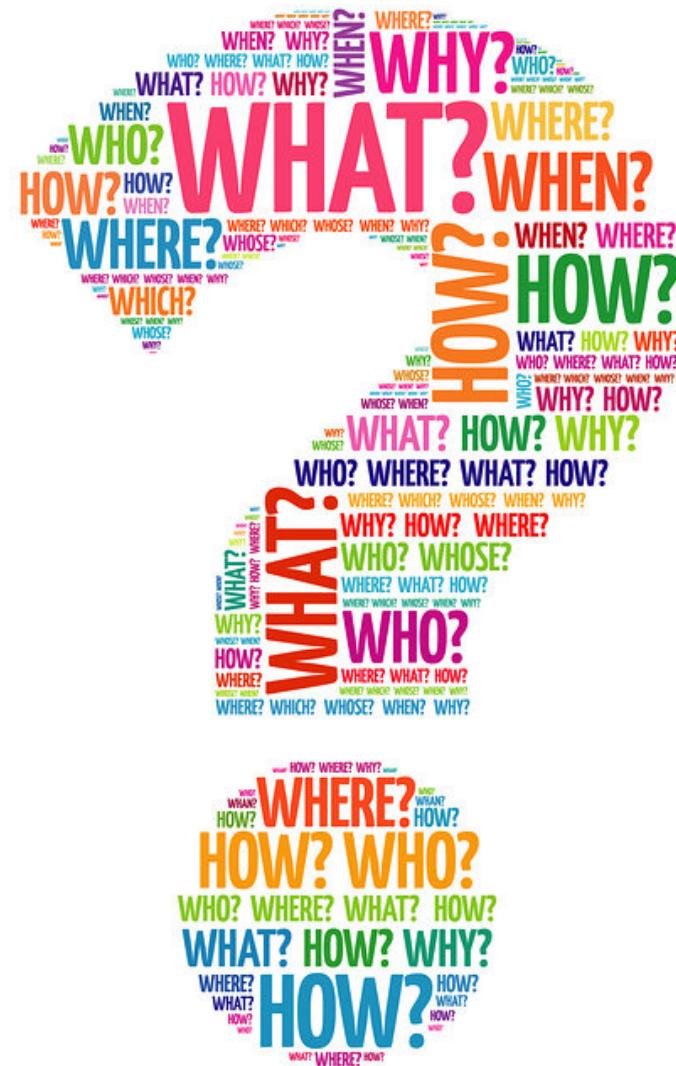


Data is a
natural
resource.



It's everywhere!

How to ask the right questions to get right answers



Startups based on Data Science

BigML

Headquartered in Oregon, [BigML](#) was founded in 2011 with the goal of designing and building a platform that would make machine learning accessible to everyone.

DataRobot

[DataRobot](#) is a Boston-based business that has created AI technology and return-on-investment enablement services for businesses competing in what it describes as today's "intelligence revolution".

Skin Analytics

[Skin Analytics](#) is a UK start-up that was founded by Neil Daly in 2012. The start-up has developed and tested an AI platform that aims to detect skin cancer at an early stage, to help those who need medical help get it sooner while giving peace of mind to those who don't need medical assistance.

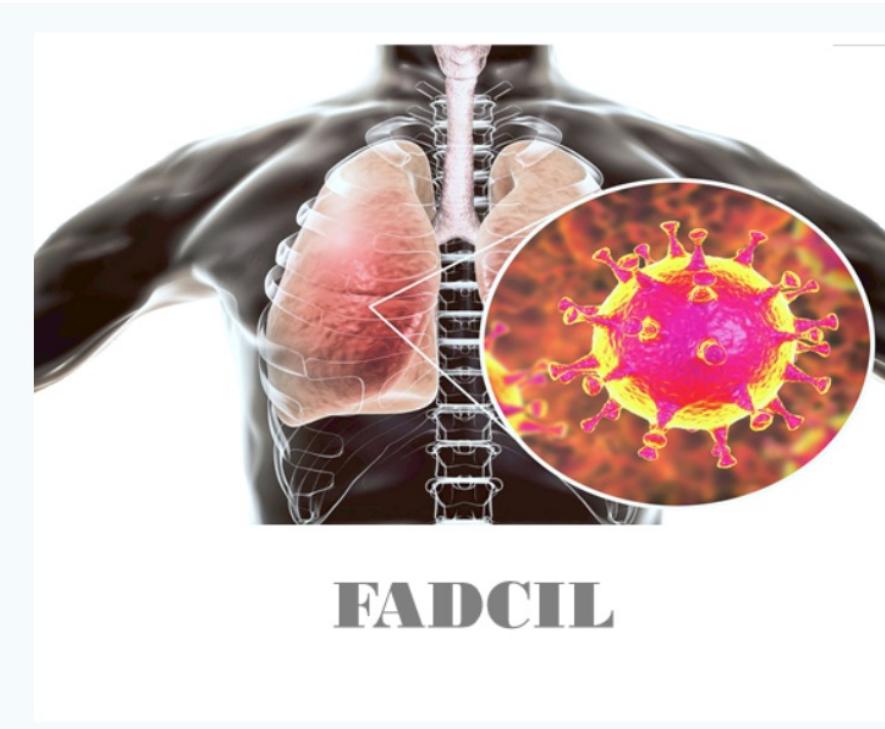
Dataiku

[Dataiku](#) is a New York-headquartered start-up that was founded by Clément Stenac, Florian Douetteau, Marc Batty and Thomas Cabrol in 2013. It has built a centralised data platform that supports businesses working with analytics at scale and enterprise AI.



**TOP
LATIN AMERICA**

STARTUPS - 2020 | Awarded by
STARTUPCITY



FADCIL

*Data
Science*

Data Science

Birdie is an **innovative startup** that helps companies understand millions of consumers' opinions



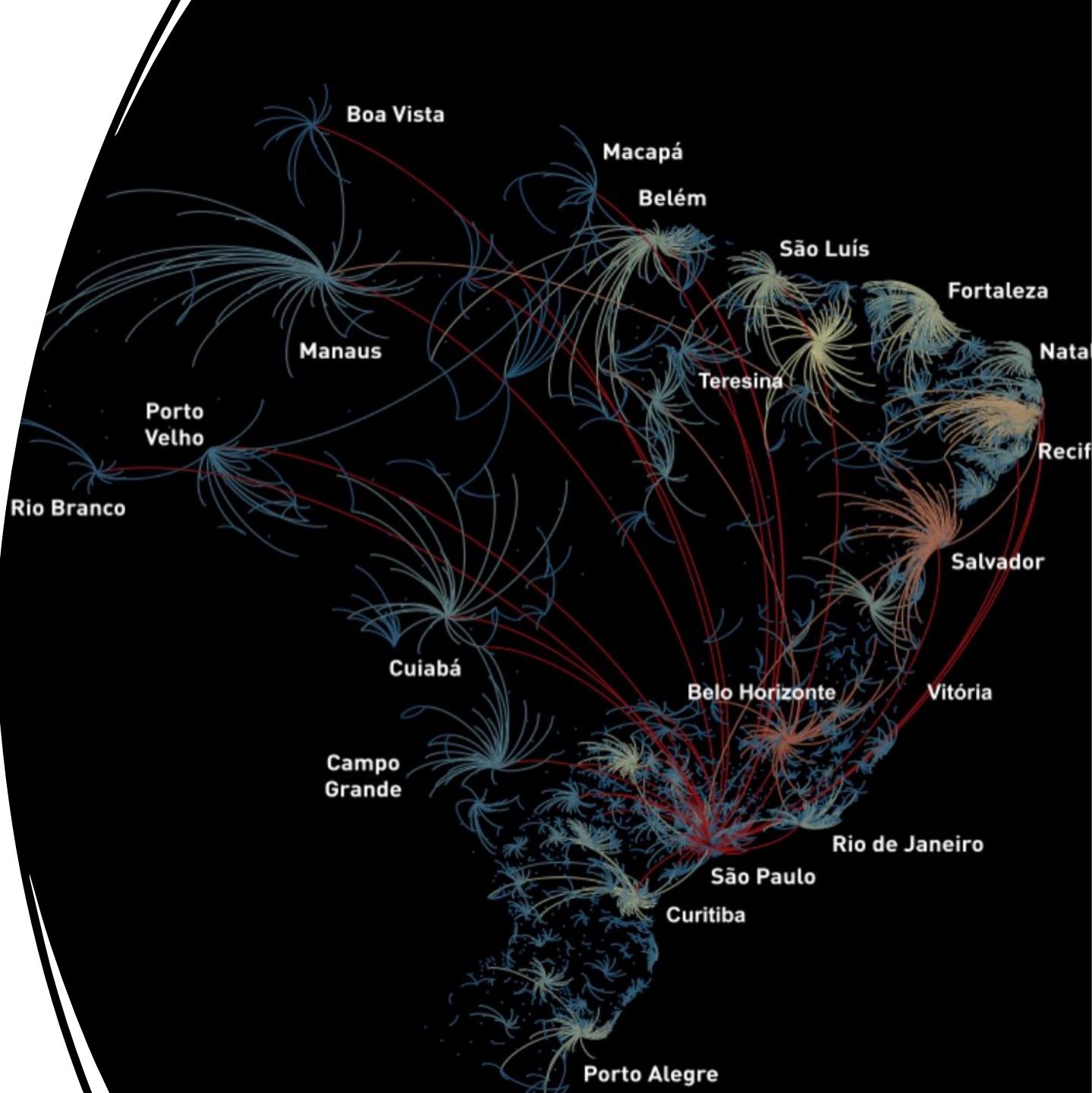
crunchbase news

Data Science COVID-19

scientific reports

 OPEN The impact of super-spreader cities, highways, and intensive care availability in the early stages of the COVID-19 epidemic in Brazil

Miguel A. L. Nicolelis^{1,2,3,4,5,6}, Rafael L. G. Raimundo⁷, Pedro S. Peixoto⁸ & Cecilia S. Andreazzi⁹



BrainCare

Invenção brasileira pode trazer novo sinal vital para triagem do paciente

Criado por Sérgio Mascarenhas, sensor da Brain4care traz método não invasivo para monitorar problemas neurológicos e agora parte para conquistar uma fatia estimada de US\$ 1 bilhão no mercado dos EUA

Por **Laura Pancini**

Publicado em: 22/11/2021 às 11h18

🕒 Tempo de leitura: **3 min**



Sensor da Brain4care é totalmente não invasivo e envia relatório em tempo real para aplicativo (Brain4care/Divulgação)

THE DATA SCIENCE HIERARCHY OF NEEDS

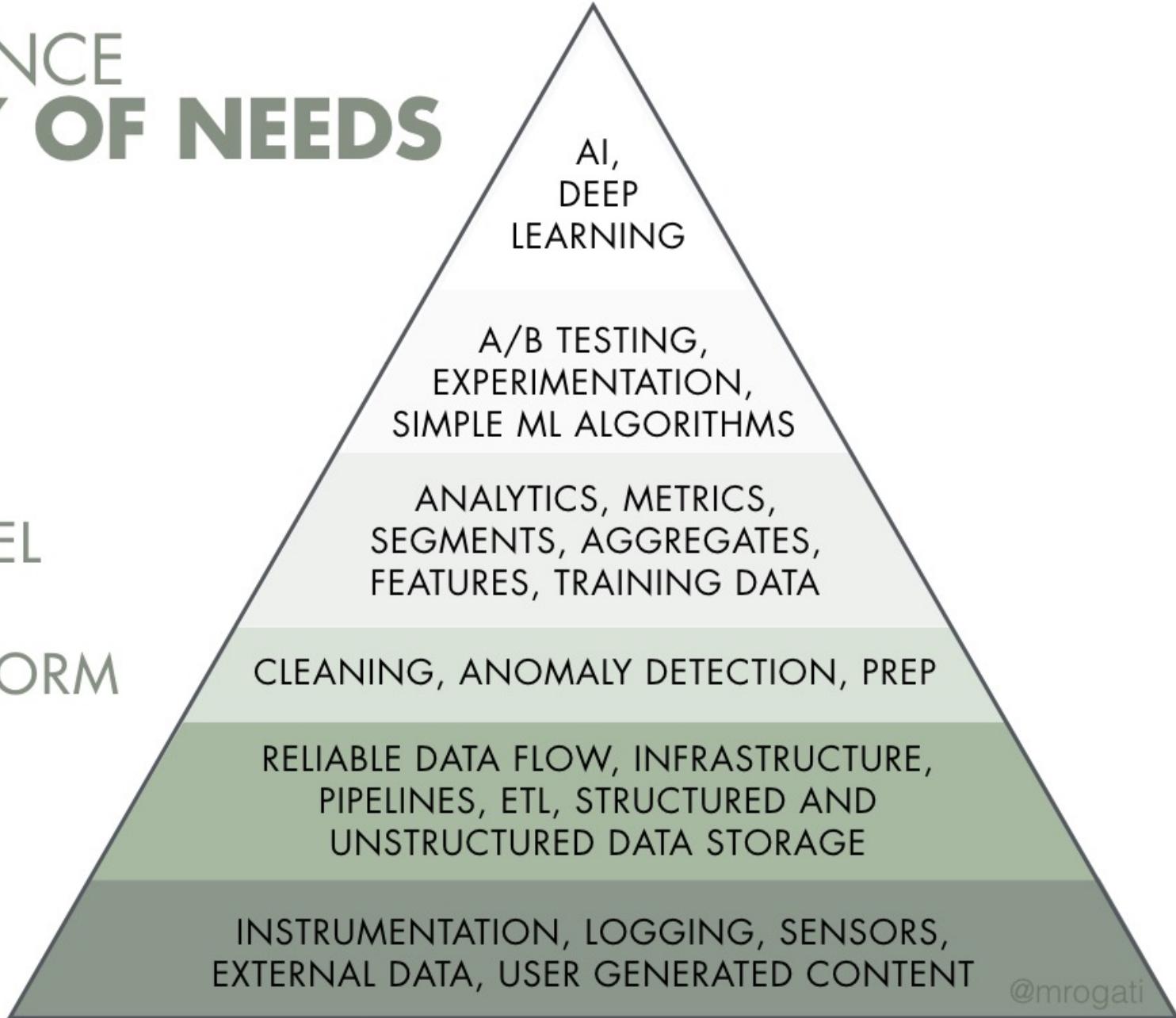
LEARN/OPTIMIZE

AGGREGATE/LABEL

EXPLORE/TRANSFORM

MOVE/STORE

COLLECT



Jobs!



Hot Companies Hiring Data Scientists

November 11, 2019 | Posted by Jillian Kramer



Inspiration

- Intro to Data Science by Steve Brunton (2019) / University of Washington
 - eigensteve.com
- Ciência de Dados: Introdução e Motivação por Francisco Rodrigues (2020) / Universidade de São Paulo
 - <https://www.youtube.com/watch?v=Im2lagDGDAU&t=636s>
- Notas de aula do Prof. Andre Carvalho (2016) / Universidade de São Paulo