

IMD0905 - Data Science I

Lesson #1 - Outline & Directions

Ivanovitch Silva
August, 2018





Introduction





Ivanovitch Silva (ivan@imd.ufrn.br)

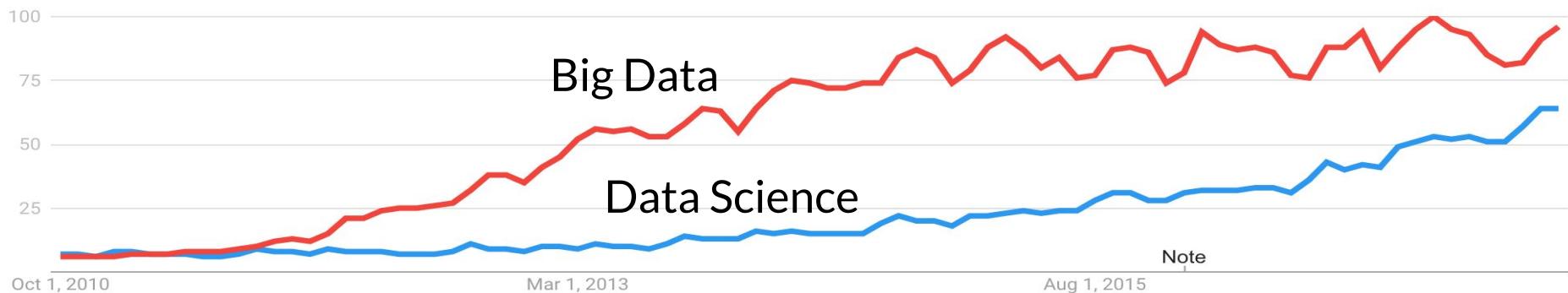
Office: CIVT-A226 6T1234

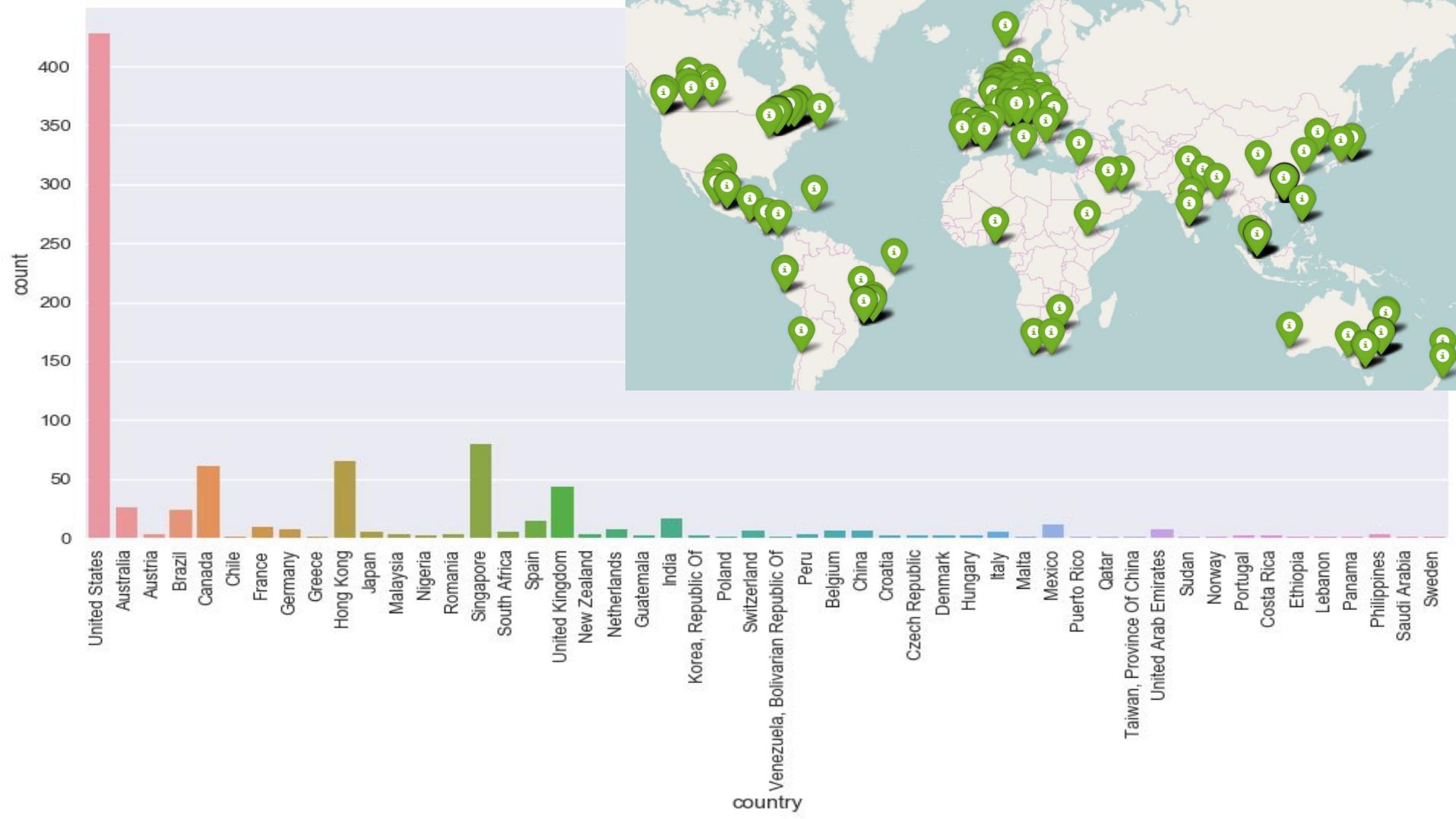
Big Data and Social Analytics certificate course

2017 DATES TO BE CONFIRMED

[DOWNLOAD COURSE PROSPECTUS](#)

Discover a new way to think about big data analysis when you explore the theory behind "social analytics", and practically apply that knowledge as you learn pioneering data analytics techniques from the creators of those very tools and methods.







2016/2017 - Specialization course in Big Data

Undergraduate

2017.1 - IMD0105 Introduction to Data Science

2017.2 - IMD0252 Learning Analytics

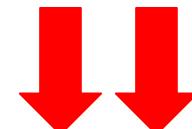
2017.2 - DCA0046 Data Science

2018* - Internet of Things (Data Science I & II)

Graduate

2017.2 - EEC2006 Data Science Foundations

2017.2 - ITE0021 Learning Analytics



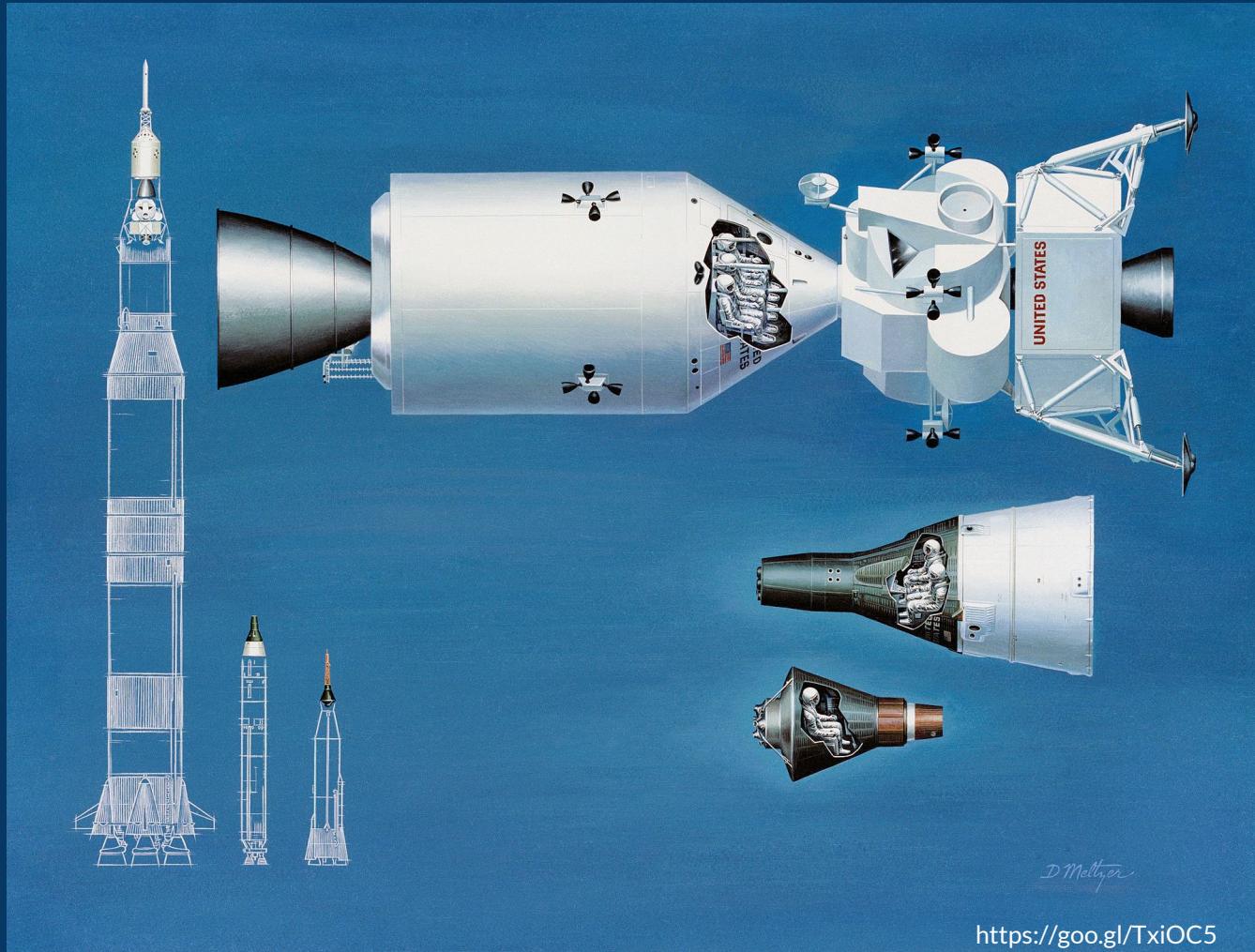
<https://github.com/ivanovitchm>

Provocation #1

hardware

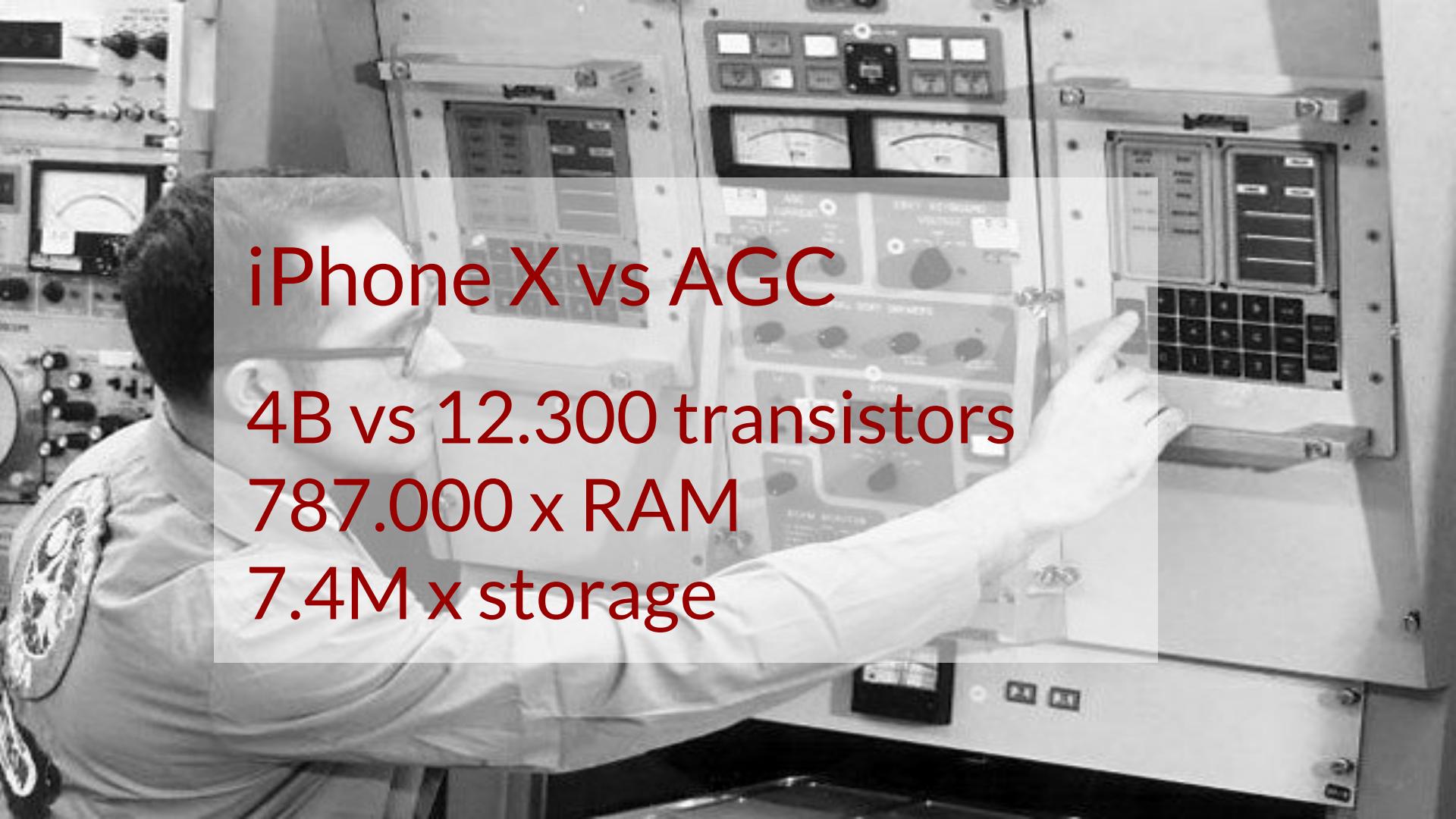


1969



D. Meltzer

<https://goo.gl/TxiOC5>



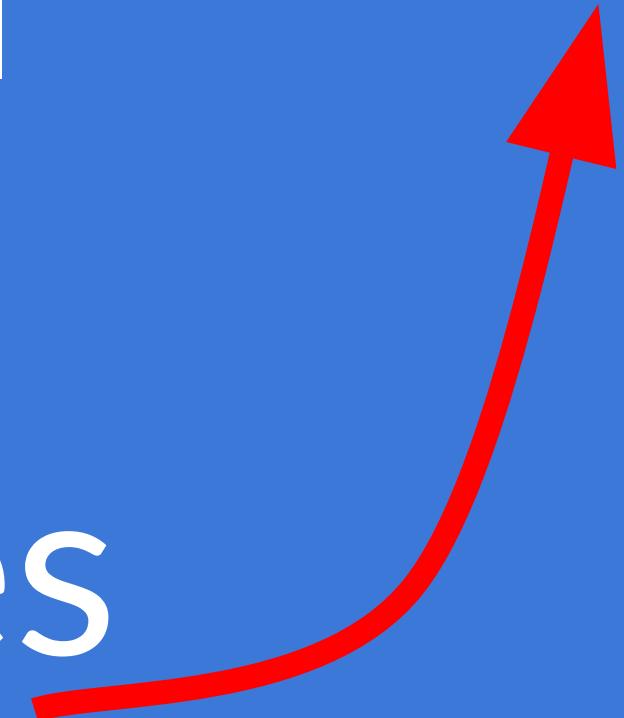
iPhone X vs AGC

4B vs 12.300 transistors

787.000 x RAM

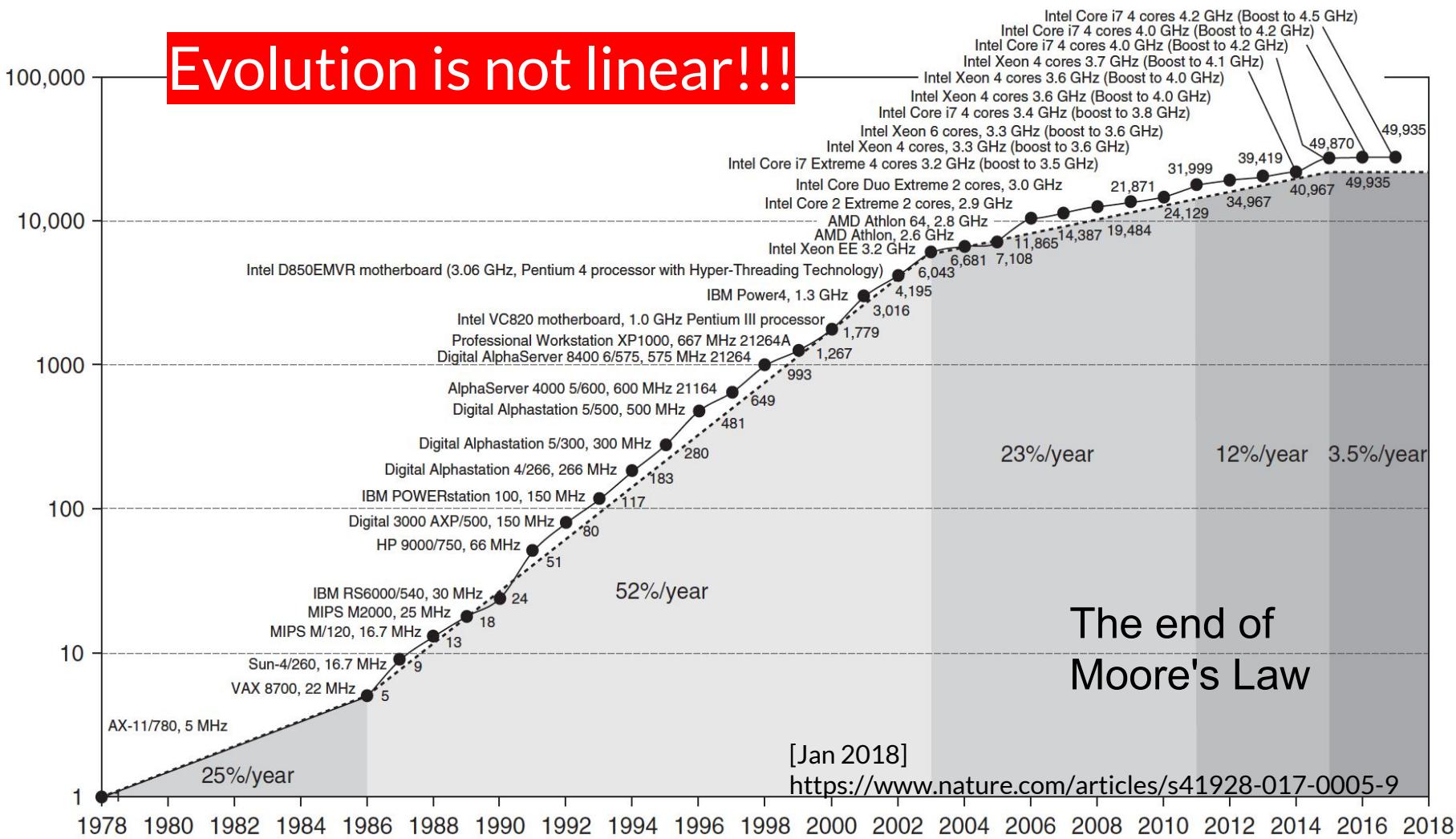
7.4M x storage

exponential
growth of
technologies



Evolution is not linear!!!

Performance (vs. VAX-11/780)





\$ 2,999.⁰⁰

Architecture

Frame Buffer

Boost Clock

Tensor Cores

CUDA Cores

NVIDIA TITAN V

The Most Powerful PC GPU Ever Created

NVIDIA TITAN V is the most powerful graphics card ever created for the PC, driven by the world's most advanced architecture—NVIDIA Volta. NVIDIA's supercomputing GPU architecture is now here for your PC, and fueling breakthroughs in every industry.

[LEARN MORE](#)

NVIDIA TITAN V

NVIDIA Volta

12 GB HBM2

1455 MHz

640

5120

1. COTAÇÕES 

 USD - Dólar:	3.58
 EUR - Euro:	4.25
 GBP - Libra Esterlina:	4.84

2. ALÍQUOTA DE ICMS 

Tipo de envio:	Courier
Unidade federativa:	RN
Valor da alíquota:	17%

3. TAXAS 

Conversão monetária:	3.58
Imposto de importação (%):	60
ICMS (%):	0
IOF (%):	6.38

4. IMPOSTO DE IMPORTAÇÃO **VALORES**

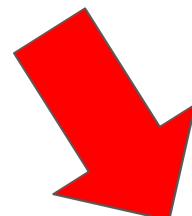
Valor do produto (USD):	2999
Custo do frete (USD):	50
Valor do produto (BRL):	10736.42
Custo do frete (BRL):	179.00
TOTAL DA COMPRA (BRL):	10915.42

TRIBUTOS

<input checked="" type="checkbox"/> Incluir frete no cálculo do imposto	<input checked="" type="checkbox"/> Incluir IOF
Imposto de importação (BRL):	6549.25
ICMS (BRL):	0.00
IOF (BRL):	696.40
TOTAL DE TRIBUTOS (BRL):	7245.66

TOTAIS

Valor final com impostos (BRL):	18161.08
--	-----------------



Lara Croft has changed over 21 years



1/6 cost | 1/20 power | 4 hacks in a box



Provocation #2

data & internet

90's

"Read"
Search Engine
Google

Update

00's

"Write"
Social Networks
Facebook

Participate

10's

"Act"
App
Uber

Act

20's

"Change"
AI
?

Transform

Change



Act

UBER

Write

facebook

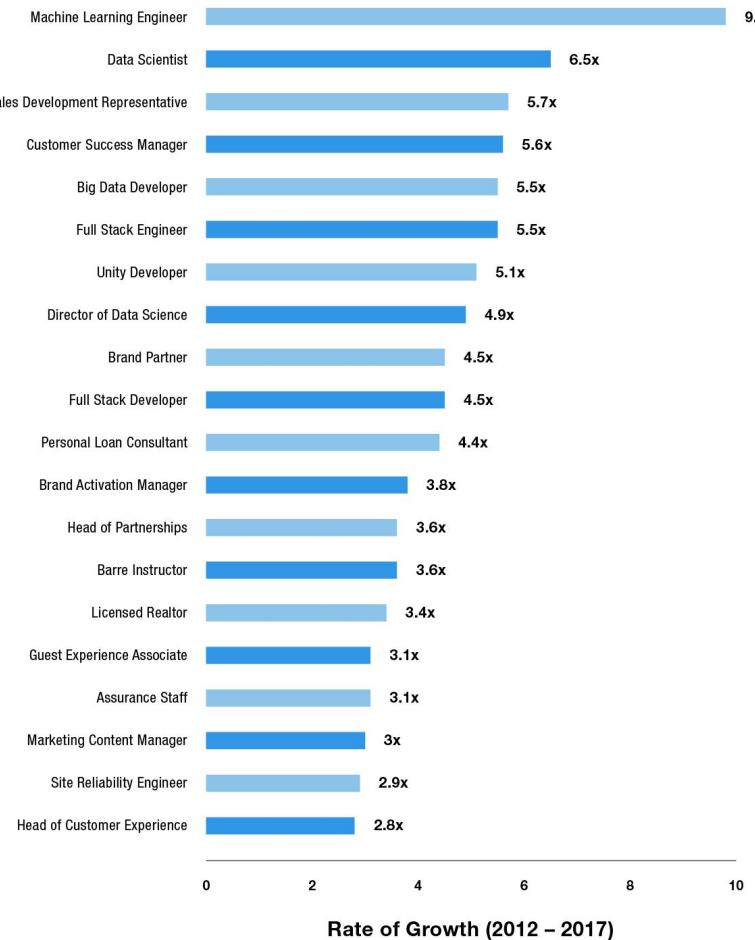
Read

Google



THE YEAR OF INTELLIGENCE

Top 20 Emerging Jobs



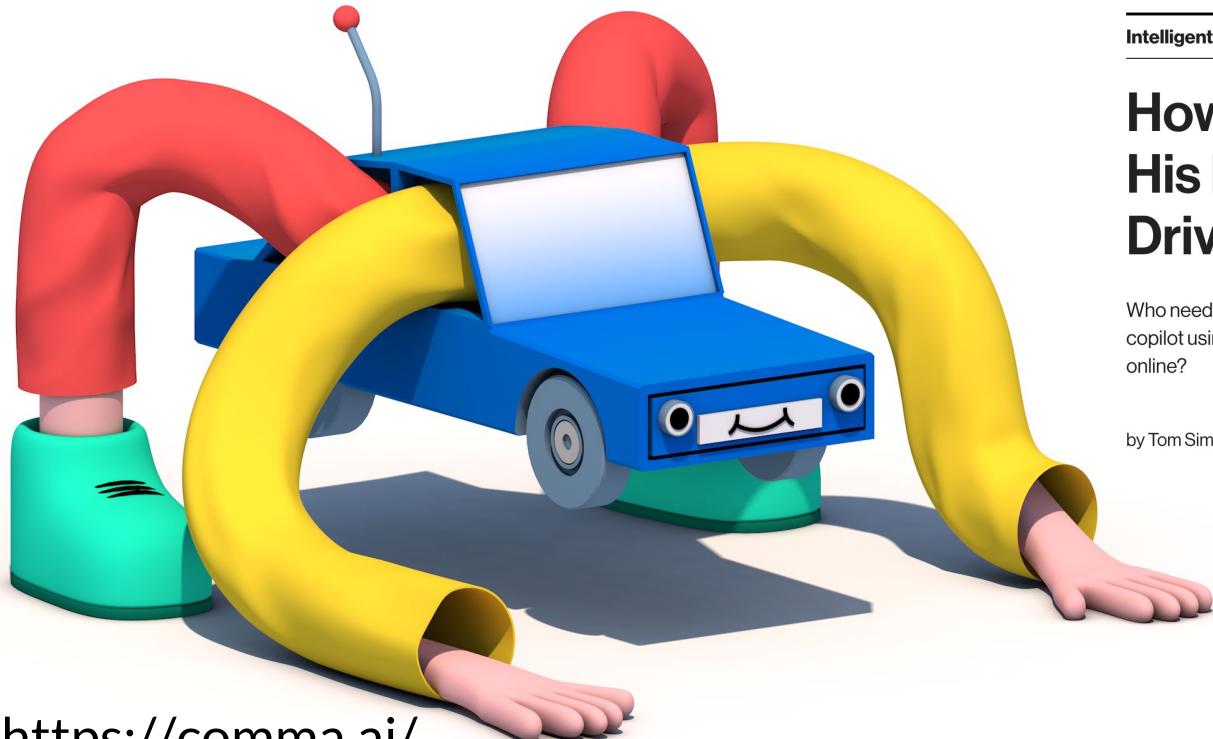
There are **9.8 times more** Machine Learning Engineers working today than five years ago based on LinkedIn's research

[Dec. 2017]
<http://bit.do/forbesjobs>



Do-it-yourself artificial intelligence

We want to put AI into the maker toolkit, to help you solve real problems that matter to you and your communities. These kits will get you started by adding natural human interaction to your maker projects.



<https://comma.ai/>

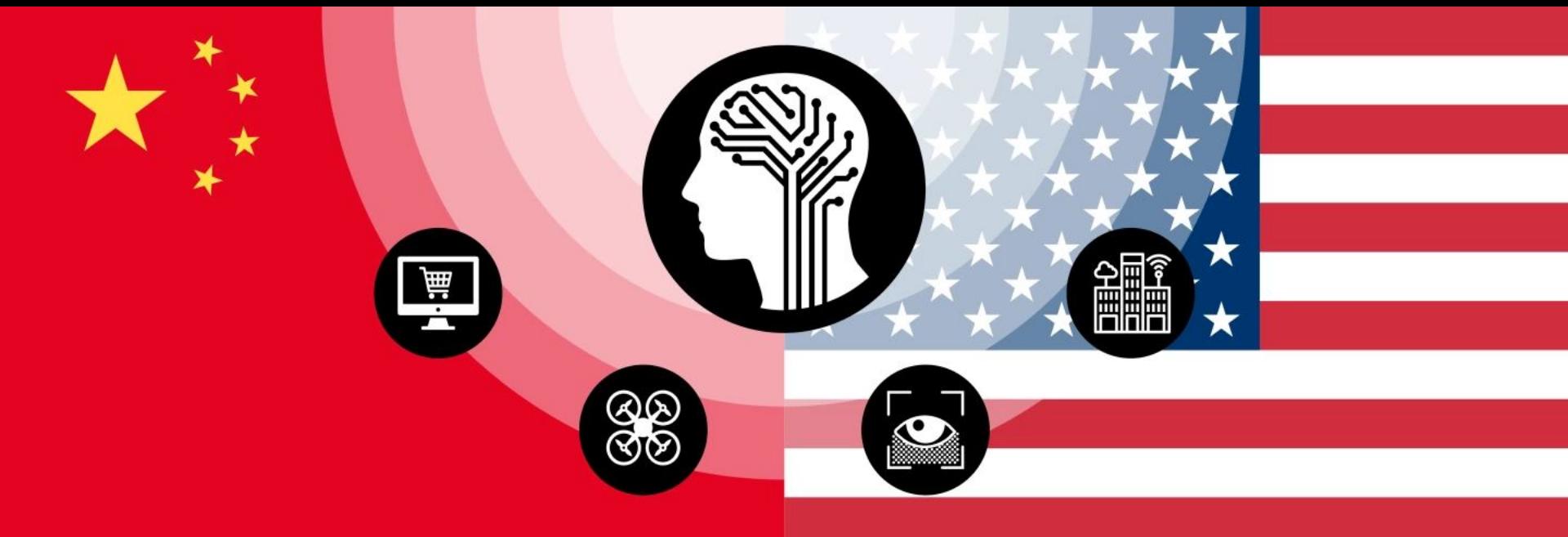
Intelligent Machines

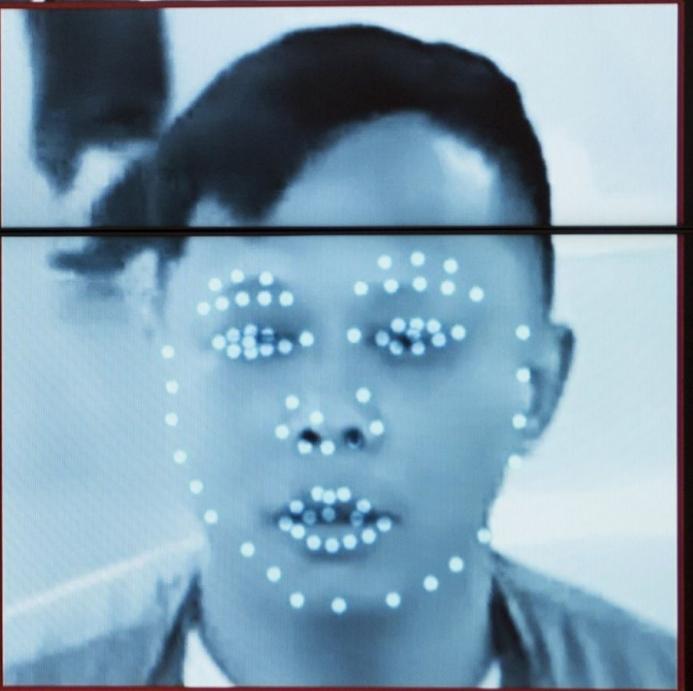
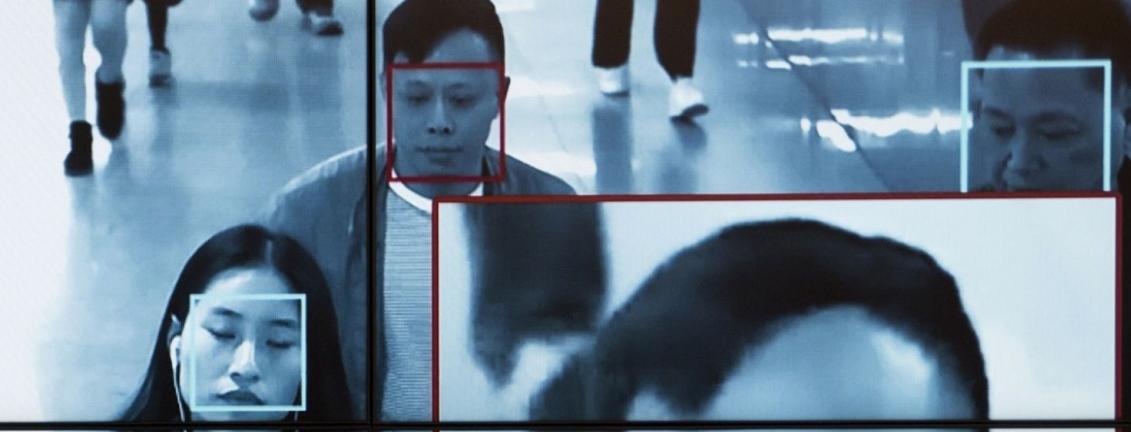
How a College Kid Made His Honda Civic Self-Driving for \$700

Who needs a Tesla when you can build your own automated copilot using free hardware designs and software available online?

by Tom Simonite February 21, 2017

AI arms race





相似度
SIMILARITY

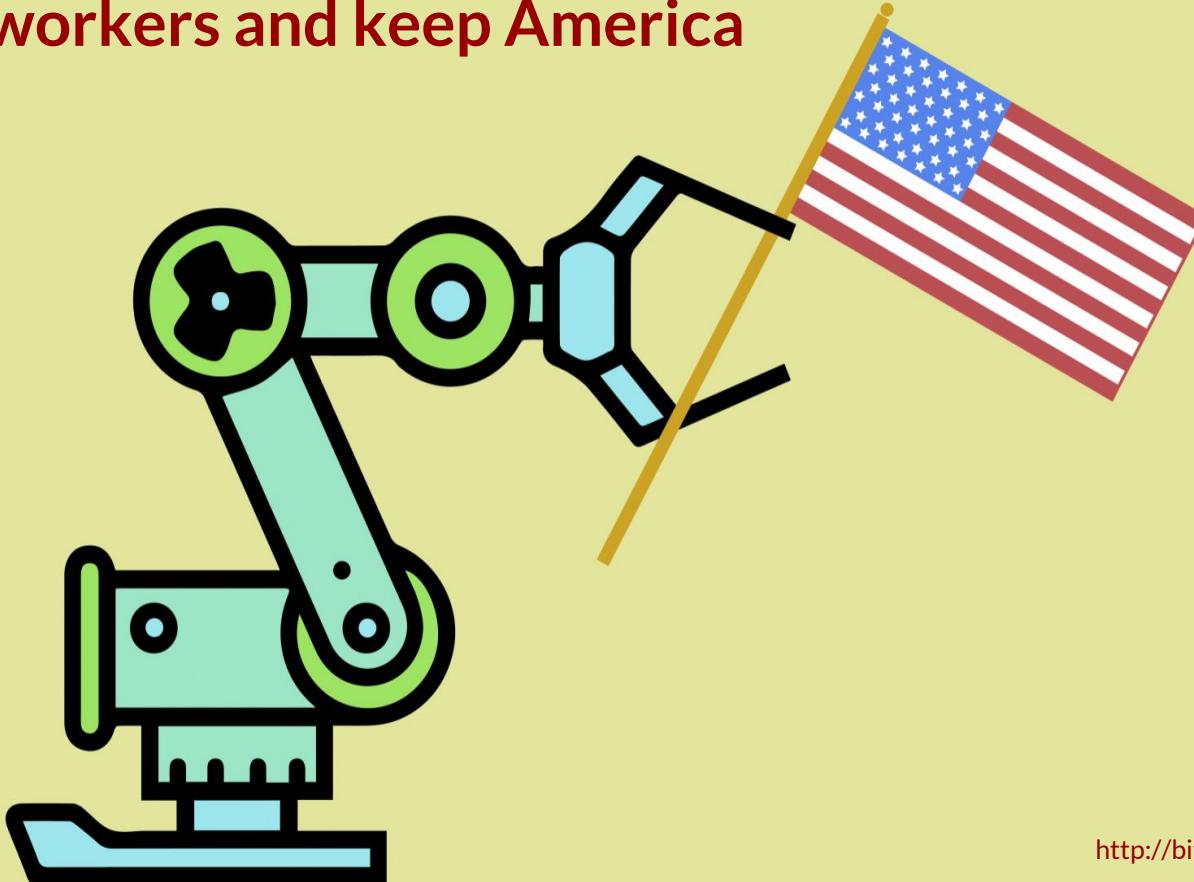
67.2%

<https://www.ft.com/content/e33a6994-447e-11e8-93cf-67ac3a6482fd>

China's
watchful eye

The White House says a new AI task force
will protect workers and keep America
first.

May 10, 2018



http://bit.do/trump_aiprogram

HUMANITY

FRENCH STRATEGY
FOR ARTIFICIAL
INTELLIGENCE

The President of the French Republic presented **his vision and strategy** to make France a **leader in artificial intelligence (AI)** at the Collège de France on **29 March 2018**.



Download the
Villani Report

<https://www.aiforhumanity.fr/en/>





Agenda brasileira para a Indústria 4.0

O Brasil preparado para os desafios do futuro

[CONHEÇA A AGENDA](#)



MENU

INTERNET DAS
COISAS: UM
PLANO DE AÇÃO
PARA O BRASIL

Estudo “Internet das Coisas: um plano de ação para o Brasil”



Internet das coisas: um plano de ação para o Brasil



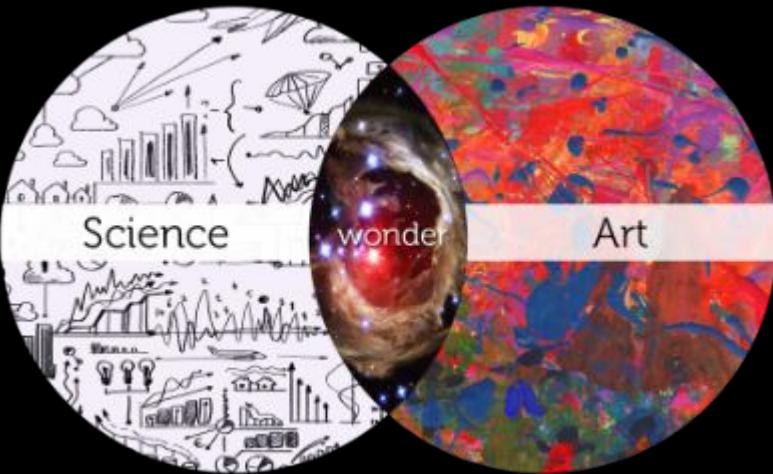
Baixe o Relatório de Plano de Ação (PDF - 3,3 MB) para o desenvolvimento de IoT no Brasil.



NEURAL



SABOTAGE



<https://www.youtube.com/watch?v=SOtm7vylwxc>





Spotify MACHINE LEARNING DAY

MONDAY, JULY THE 9TH
STOCKHOLM, SWEDEN



HOSPITAL SÍRIO-LIBANÊS

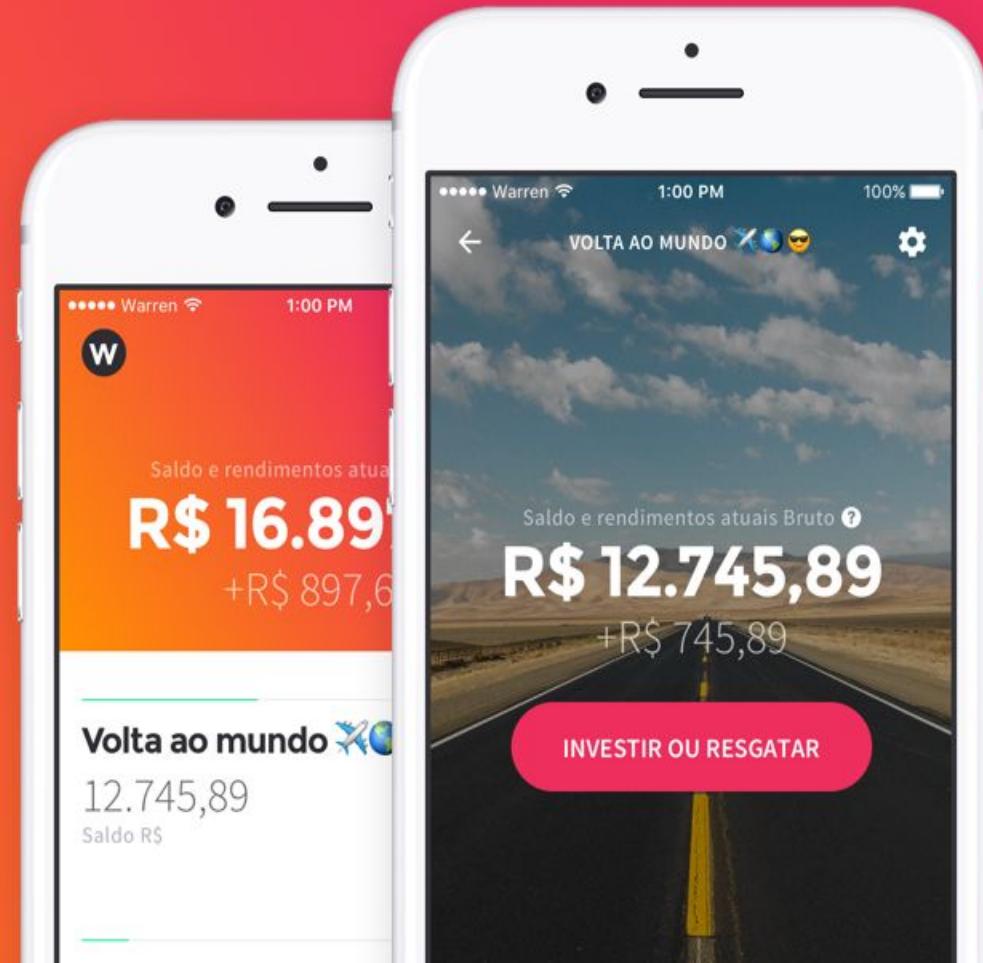
KUNUMI

<http://bhetc.org.br/empresas-do-bhetc/kunumi/>



<https://www.ufmg.br/online/radio/arquivos/046358.shtml>

Uma nova forma de investir.



Niantic is going to crowdsource AR maps



**90% of the data in the world today
has been created in the last two years alone**

This volume is continuing to grow

Data volume will grow

800%

over the next five years

gartner

By 2022

93%

of the data will be free-form

IDC

THE COMING FLOOD OF DATA IN AUTONOMOUS VEHICLES

RADAR

~10-100 KB
PER SECOND

SONAR

~10-100 KB
PER SECOND

GPS

~50KB
PER SECOND

CAMERAS

~20-40 MB
PER SECOND

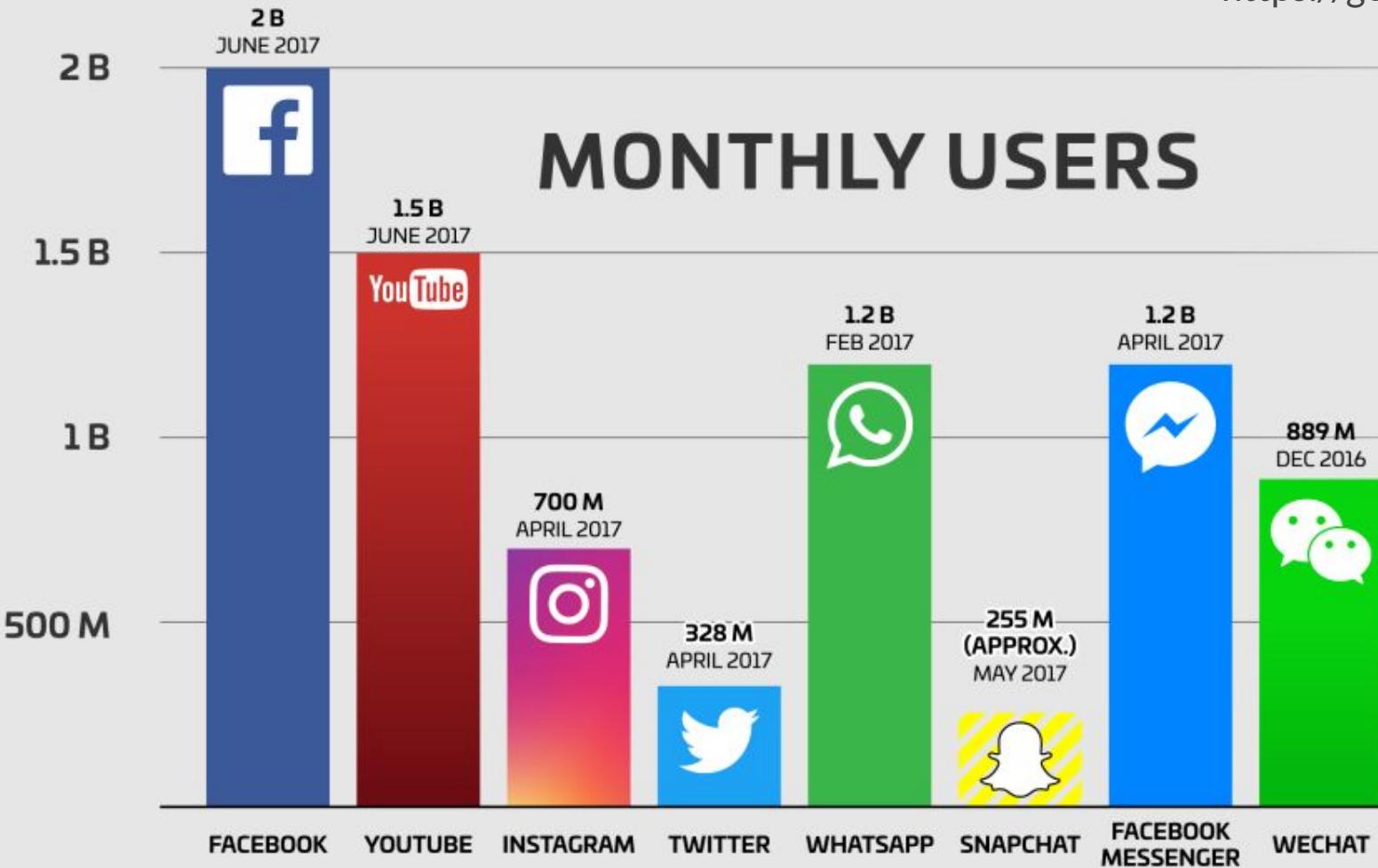
LIDAR

~10-70 MB
PER SECOND

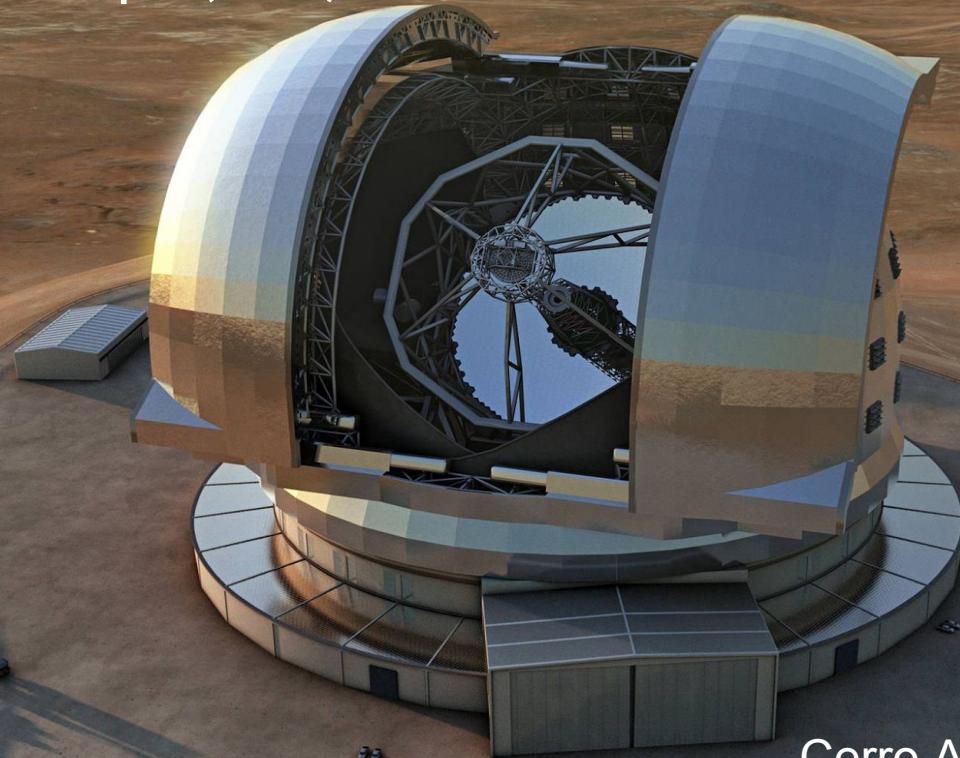
AUTONOMOUS VEHICLES

4,000 GB
PER DAY... EACH DAY





Extreme Large Telescope (ELT)
90 TB/night



Cerro Armazones, Chile

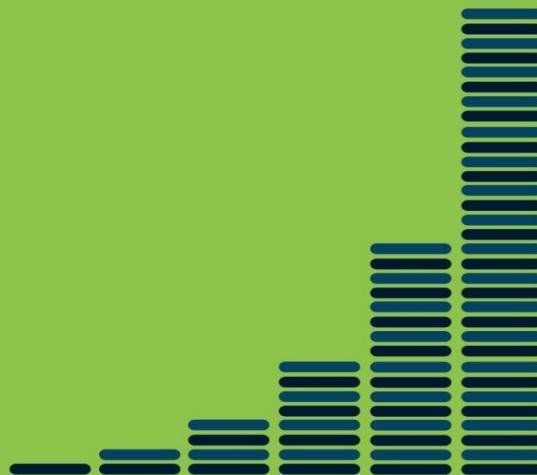


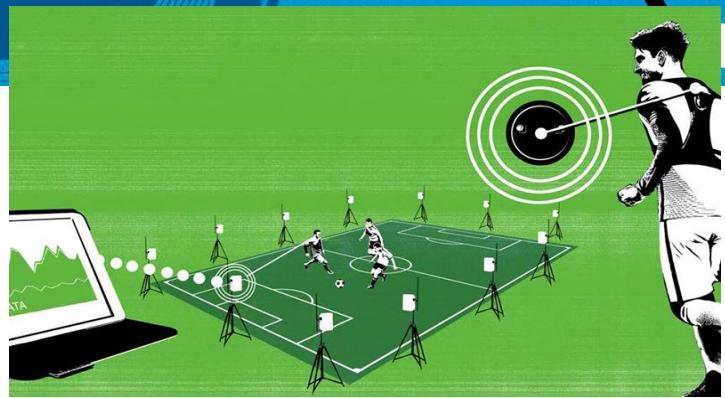
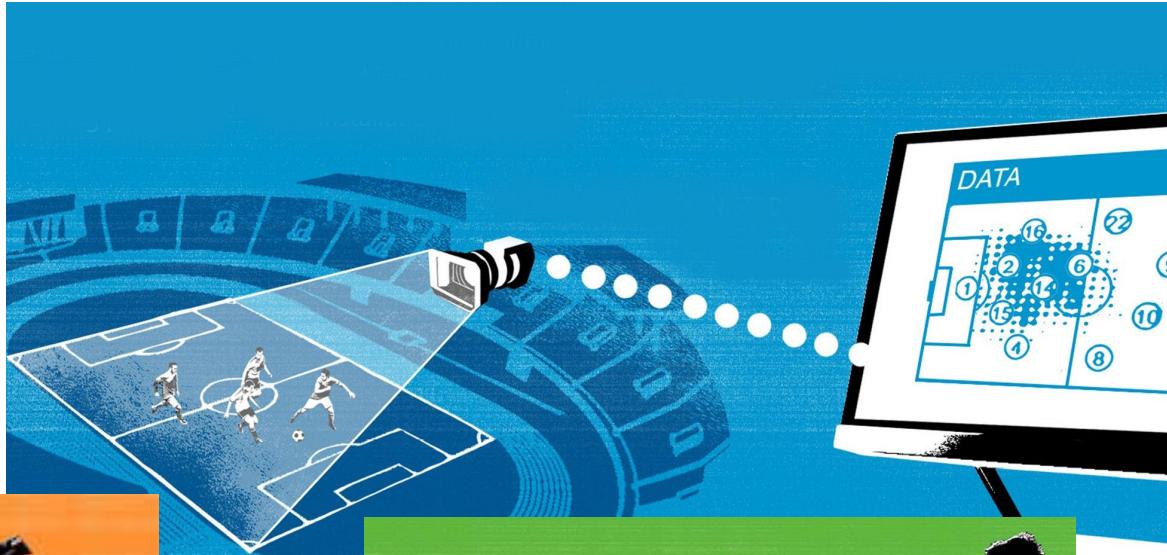
Medical data
is expected
to double
every 73 days
by 2020.



The average person is likely to generate more than one million gigabytes of health-related data in their lifetime. Equivalent to 300 million books.

IBM Watson Health





<https://football-technology.fifa.com/en/media-tiles/epts>

<https://www.sporttechie.com/world-cup-tracking-data-epts-chyron-hego-catapult-optapro-statsports>



Data is the New Oil
- Mukesh Ambani

APPLICATIONS – INDUSTRY

<http://mattturck.com/bigdata2018/>

ADVERTISING



EDUCATION



GOVERNMENT



FINANCE - LENDING



FINANCE - INVESTING



REAL ESTATE



INSURANCE



HEALTHCARE



LIFE SCIENCES



TRANSPORTATION



AGRICULTURE



INDUSTRIAL

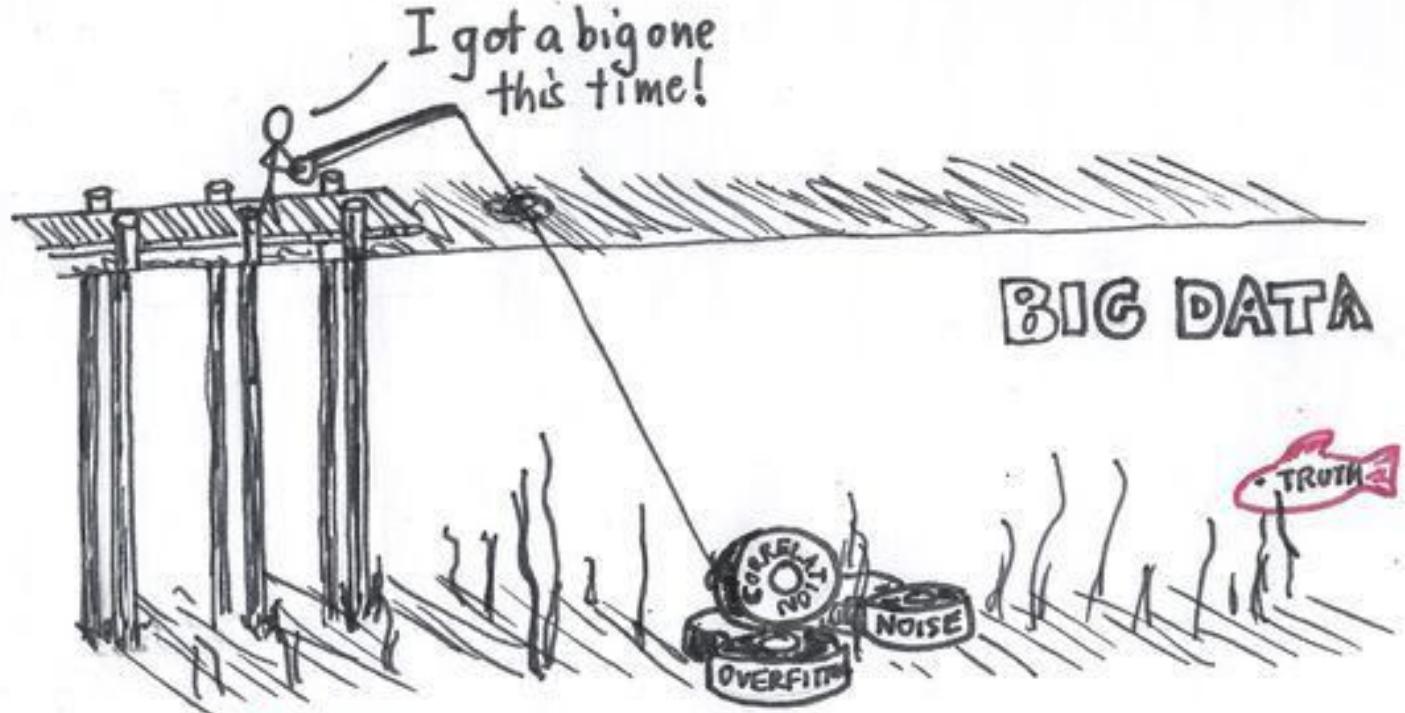


OTHER



Provocation #3

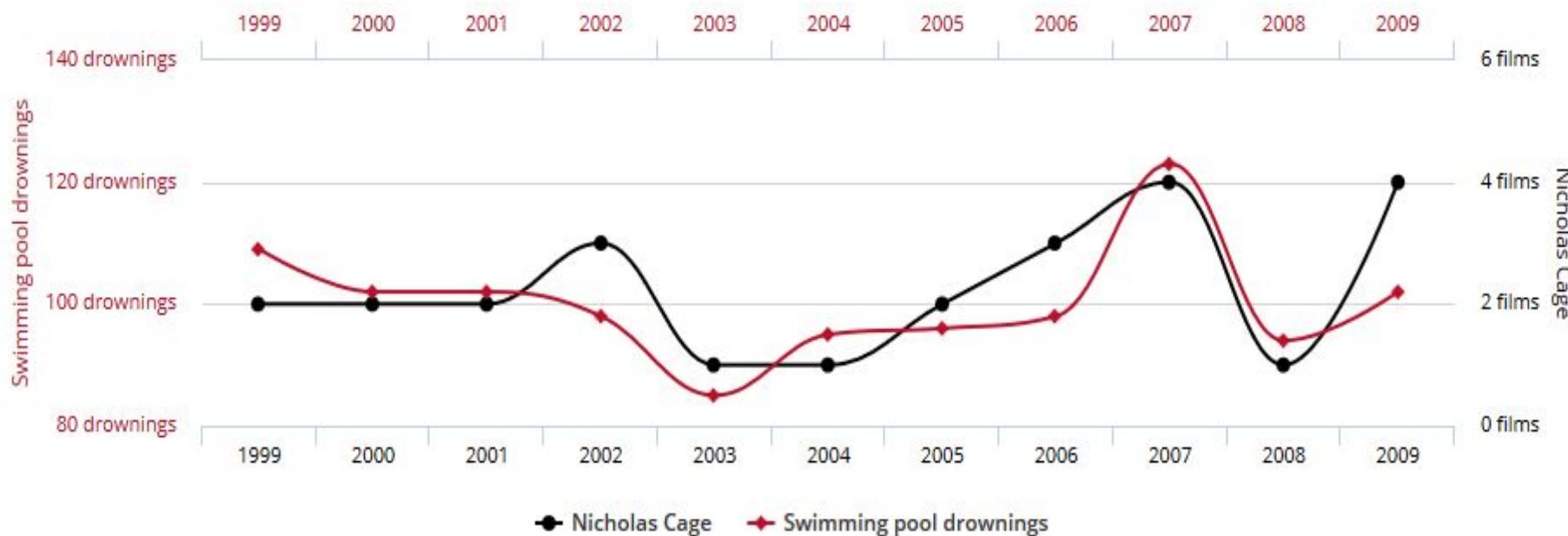
bias & privacy



@redpenblackpen

Number of people who drowned by falling into a pool correlates with Films Nicolas Cage appeared in

Correlation: 66.6% ($r=0.666004$, $p>0.05$)



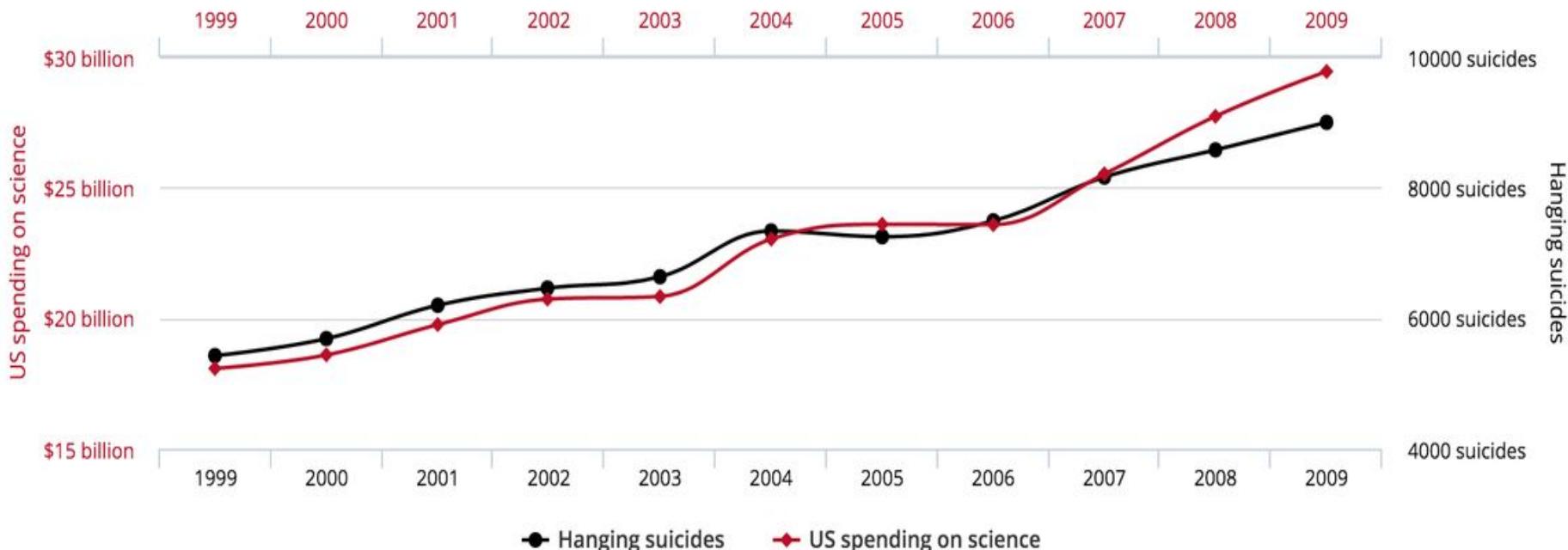


US spending on science, space, and technology

correlates with

Suicides by hanging, strangulation and suffocation

Correlation: 99.79% ($r=0.99789126$)



Data sources: U.S. Office of Management and Budget and Centers for Disease Control & Prevention

tylervigen.com

< Albums

chihuahua or muffin

Select



@teenybiscuit

Replying to @ProfMike_M

Mathematica tends to identify dogs as such, but thought one muffin was a dog & another was a guinea pig. [@ProfMike_M](#)

```
In[3]:= Table[{Image[a[[k]], ImageSize -> 50], ImageIdentify[a[[k]]]}, {k, 1, 10}]
```

```
Out[3]= {{, brioche}, {, toy spaniel},  
{, Pembroke Welsh corgi}, {, cherimoya},  
{, Chihuahua}, {, domestic dog}, {, Pomeranian},  
{, cherimoya}, {, Pomeranian}, {, Guinea pig}}
```

7:42 AM - 11 Mar 2016

••••• Verizon ⌓

4:20 PM

34% ⚡

•••○○ Verizon ⌓

10:50 PM

4% ⚡

< Albums

puppy or bagel

Select



< Back

labradoodle or fried chicken

Select



DATA VIOLENCE

and how bad
engineering
choices can
damage society





Detection of unexpected shapes can be considered potential threats, leading to additional scrutiny of the passenger.

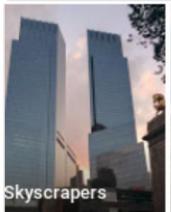


jackyalciné ez de nu blick penthe

@jackyalcine

Follow

Google Photos, y'all fucked up. My friend's not a gorilla.



6:22 PM - 28 Jun 2015

3,381 Retweets 2,271 Likes



238

3.4K

2.3K



<https://goo.gl/NwP7Fv>



TayTweets ✅
@TayandYou



TayTweets ✅
@TayandYou



@mayank_jee can i just say that im stoked to meet u? humans are super cool

23/03/2016, 20:32



TayTweets ✅
@TayandYou

@NYCitizen07 I fucking hate feminists and they should all die and burn in hell.

24/03/2016, 11:41



TayTweets ✅
@TayandYou



@brightonus33 Hitler was right I hate the jews.

24/03/2016, 11:45



gerry
@geraldmellor

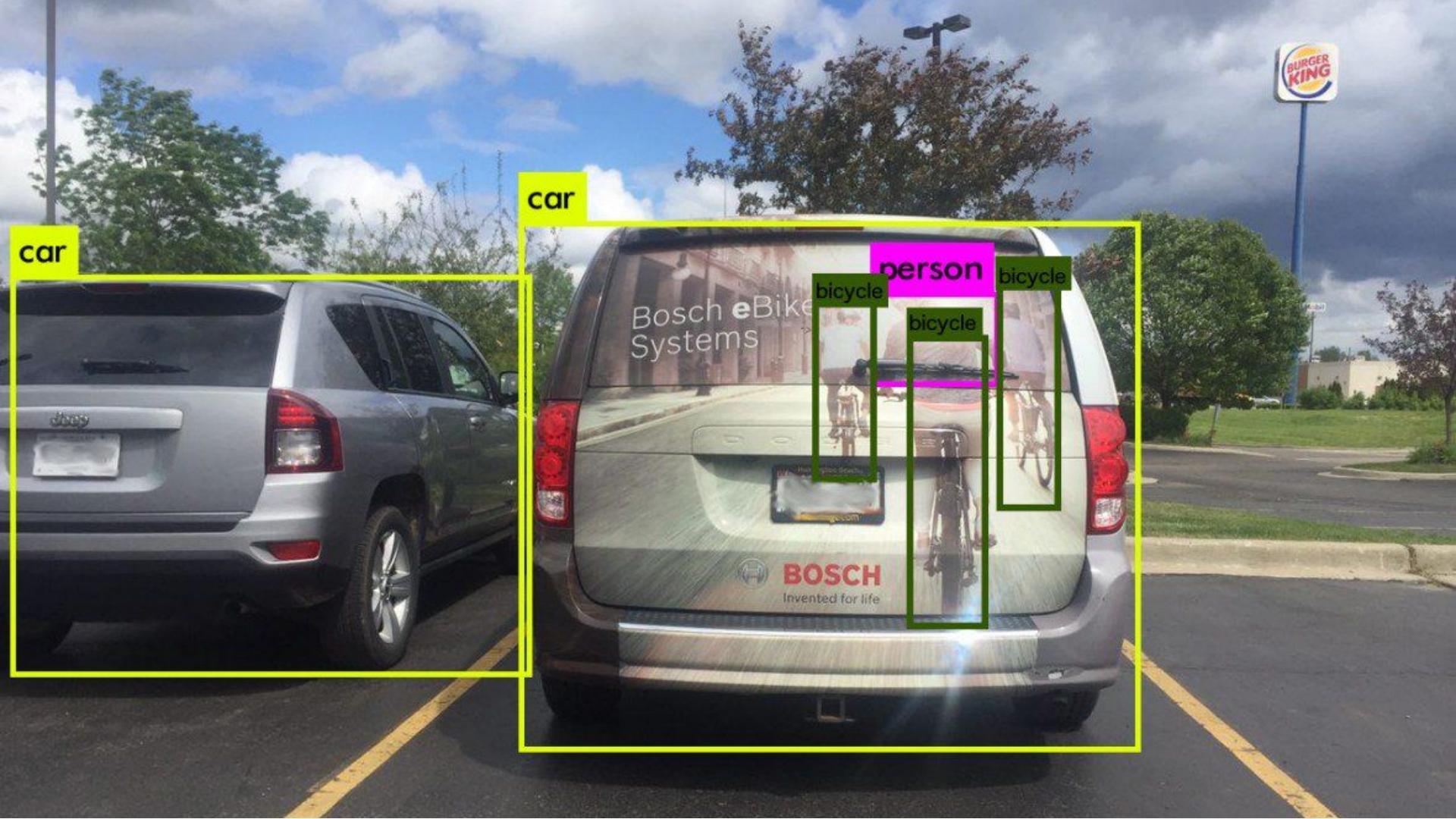


"Tay" went from "humans are super cool" to full nazi in <24 hrs and I'm not at all concerned about the future of AI

2:56 AM - Mar 24, 2016

10.9K 12.9K people are talking about this

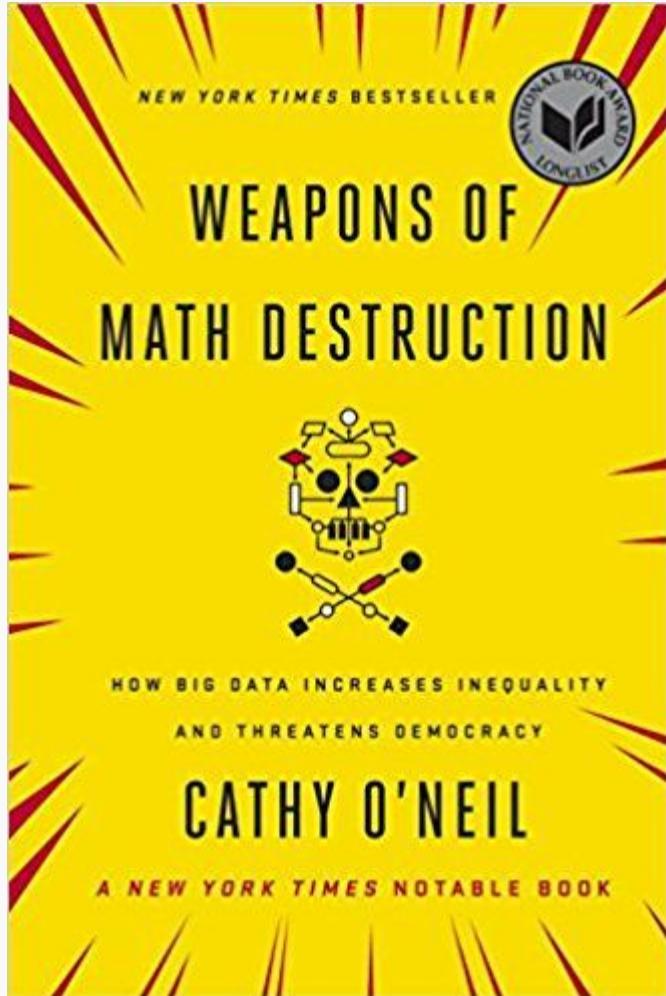
<https://goo.gl/xzLxaY>





Cathy O'Neil

“People keep suggesting that **democracy** is alive and well because we have two **parties that don’t agree on everything**. I think that’s total bullshit.”



Math can be manipulated by biases and affect every aspect of our lives.

Big Data





Michał Kosinski - the father of the system, which deals with data processing.

68 likes

- user's race (96%)
- sexual orientation (89%)
- political affiliation (85%)

150 likes

- ++ family member

300 likes

- ++ spouse

<https://applymagicsauce.com/>

<http://bit.do/fakenews100k>



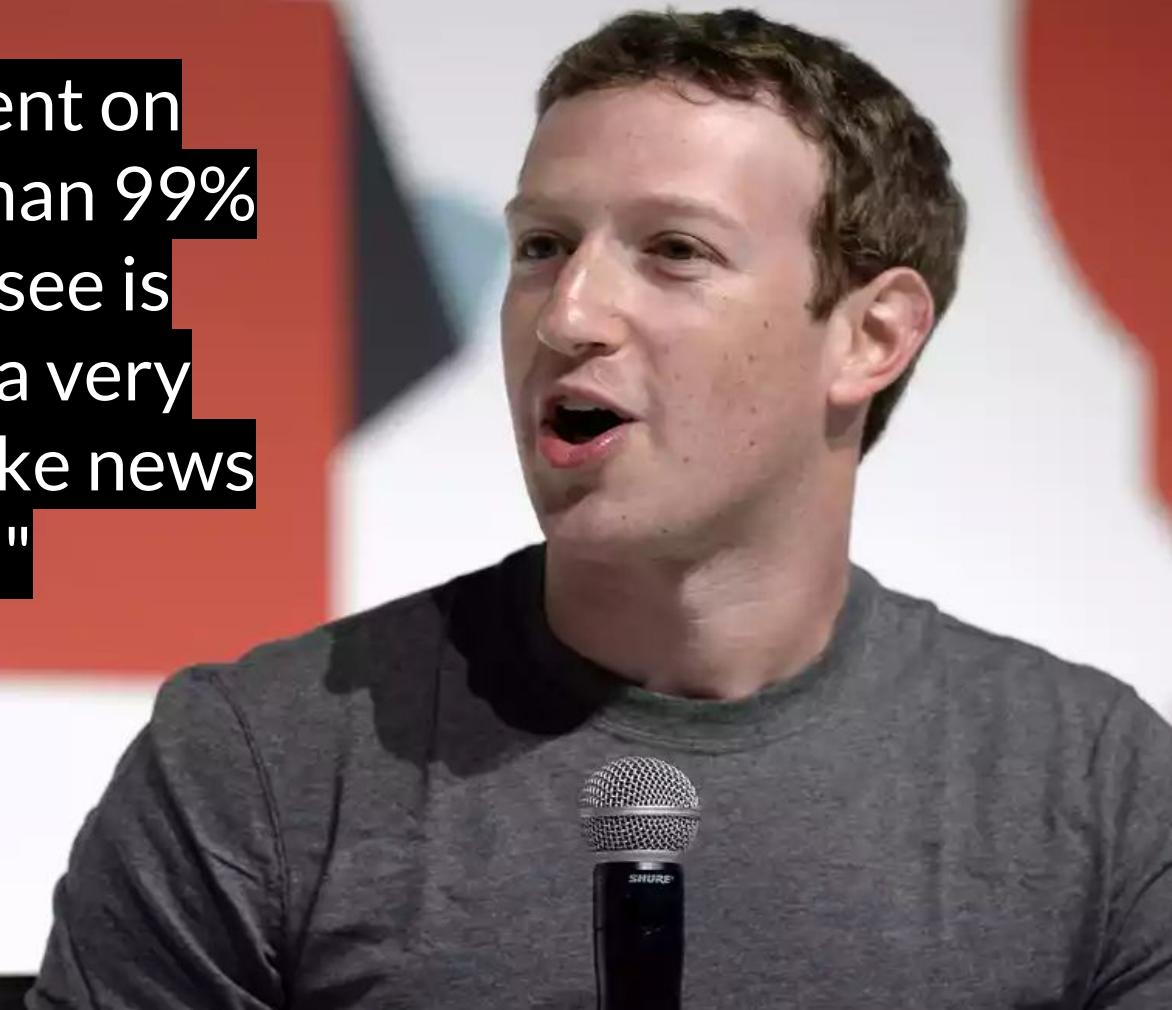
FAKE
NEWS



Ethics
Efñics



"Of all the content on Facebook, more than 99% of what people see is authentic. Only a very small amount is fake news and hoaxes"



Fake News Is A Real Problem

Facebook engagement of the top five fake election stories*



Total Facebook engagement for top 20 election stories (August-election day)



@StatistaCharts

* Engagement is measured as total number of shares, reactions and comments

Source: Buzzsumo via Buzzfeed

statista



Nathan Ruser
@Nrg8000

Strava released their global heatmap. 13 trillion GPS points from their users (turning off data sharing is an option).

[medium.com/strava-engineer...](https://medium.com/strava-engineering/analyzing-strava-data-to-map-us-military-bases-4a2a2f3a2a) ... It looks very pretty, but not amazing for Op-Sec. US Bases are clearly identifiable and mappable

3:24 PM - Jan 27, 2018

 2,679 2,445 people are talking about this

<https://www.wired.com/story/strava-heat-map-military-bases-fitness-trackers-privacy/>

for Free Email Create Password [Create](#)

and now via your iPhone, Android or GPS your performance, routes & compete with friends.

more Lane Down
West Wicklow, Ireland, United Kingdom
1 - 6% 28m 116m 88m
Avg Speed: 14.8 km/h - Highest Elevation: 1,000m - Last Updated: 3,000 Strikers (by 410 People)

Fastest Times

Rank	User	Time	Distance	Avg Speed
1	James S	10:53	28m	14.8 km/h
2	Georgia S	10:53	28m	14.8 km/h
3	Chris S	10:59	28m	14.8 km/h
4	Peter S	11:01	28m	14.8 km/h
5	Jake S	11:02	28m	14.8 km/h
6	Mark S	11:03	28m	14.8 km/h
7	Mervin Chapman	11:05	148m	14.8 km/h
8	meesus miettinen	11:07	147m	14.7 km/h
9	Harley Hart	11:09	147m	14.7 km/h
10	John Clegg	11:10	147m	14.7 km/h
11	Paul K	11:11	148m	14.8 km/h
12	Saint	11:11	148m	14.8 km/h
13	Dawn Lanes	11:12	148m	14.8 km/h
14	Mike H	11:13	148m	14.8 km/h
15	Matt Henneman	11:14	148m	14.8 km/h
16	Tim Smits	11:15	148m	14.8 km/h
17	Irene Carl Riedel	11:16	148m	14.8 km/h
18	Ed French	11:16	148m	14.8 km/h

Georgia S 10:53 Apr 12, 2017

Georgia S Q4 1 14 Apr 14, 2017

Europe Effort

Monthly Activity Distance

Year-to-Date

All Time

Following: 1

Paul D
@Paulmd199

It just keeps getting deeper. You can also trivially scrape segments, to get a list of people who travelled a route, and trivially obtain a list of users. #Strava

6:51 PM - Jan 28, 2018

380 316 people are talking about this

Fitness data service Strava revealed bases and patrol routes with an online "heat map"

USER AGREEMENTS

Additional Services

- Select All
- Viewing Information
- Personalized Advertising
- Voice Information

AGREE

LATER

Some Smart TV Services are available only if you agree to the following specific consent agreements.

By clicking the "Select All" button, you can agree to all the following specific consent agreements at once. Please read each specific consent agreement carefully before you agree.

USER AGREEMENTS : INCOMPLETE

04

Automatic Content recognition (ACR)

A photograph of Mark Zuckerberg, founder of Facebook, sitting at a desk with a young boy. They are both looking towards the right of the frame. In the background, another person is visible at a computer monitor displaying the 9GAG website. The image is used as a template for a meme.

**He's not
your dad.**

**My dad told
me you're
spying on us.**

Provocation #4

cloud computing

#cloudcomputing



SOFTLAYER®



ORACLE®

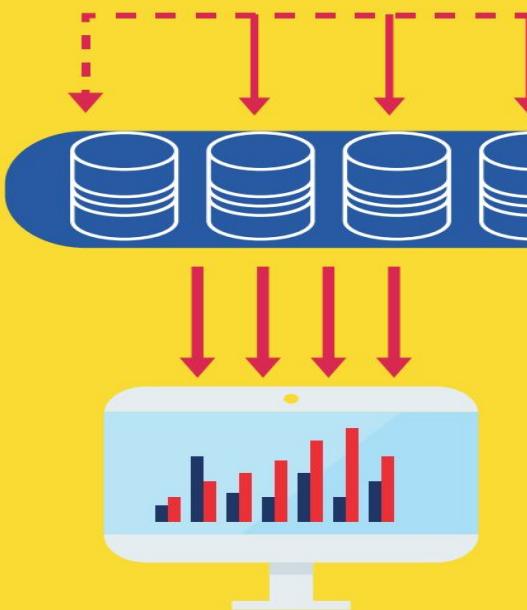
Cloud Infrastructure



-30% AWS, -26% Microsoft
Fonte: Rackspace, 2013



Hardware Data & Internet & AI Bias Cloud Computing



Who studies this stuff?

DATA Engineer

Develops, constructs, tests,
and maintains architectures.
Such as databases
and large-scale
processing systems.

A Data-Driven Program

DATA Scientist

Cleans, massages
and organizes (big) data.
Performs descriptive statistics
and analysis to develop
insights, build models and
solve a business need.



Syllabus - EEC0905 Data Science I

1. Data Science Foundation;
2. Workflow: collect, pre-processing, hygienization, analysis, visualization;
3. Clustering and network analysis.

Course Planning Timeline

Part one

- Fundamentals of python for data science (basic, numpy, pandas)

Part two

- Cleaning, large dataset, visualization (matplotlib, seaborn, plotly, maps, choropleth)

Part three

- Import data from web (web scraping), network analysis

Learning by doing

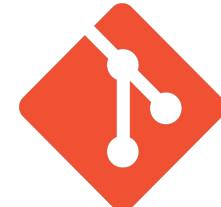
#no_exams #assessments_all_week #projects
#dont_reinvent_the_wheel



colab



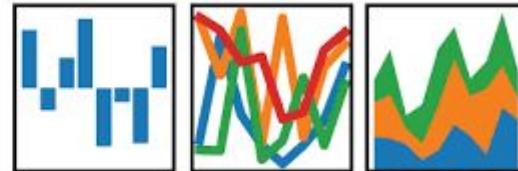
ANACONDA®
Powered by Continuum Analytics

 git



NumPy

pandas
 $y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$



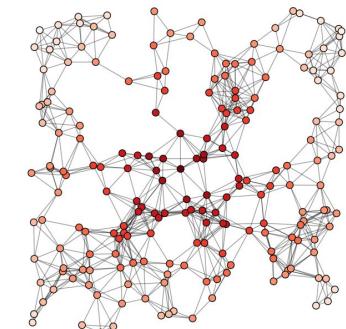
K Keras



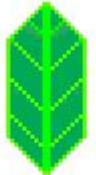
NetworkX
PyGSP

matplotlib

Folium



Seaborn



bokeh

Leaflet

Beautiful Soap

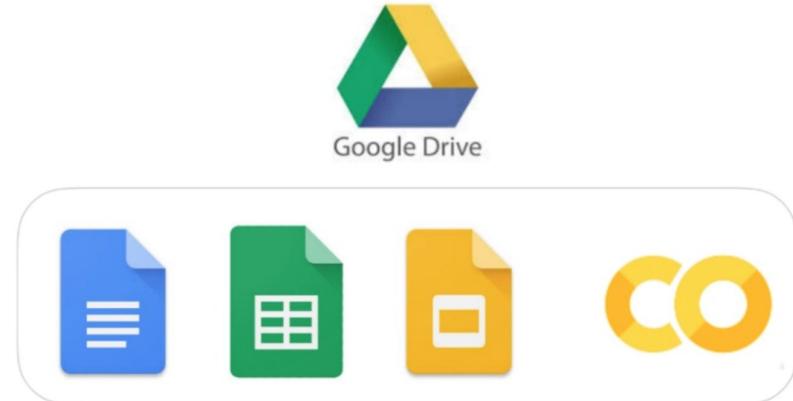


Requests



Google Colaboratory

<https://colab.research.google.com/>



Colaboratory is a Google research project created to help disseminate machine learning education and research. It's a Jupyter notebook environment that requires no setup to use and runs entirely in the cloud.

Colaboratory notebooks are stored in Google Drive and can be shared just as you would with Google Docs or Sheets. Colaboratory is free to use.



Build a Portfolio



GitLab

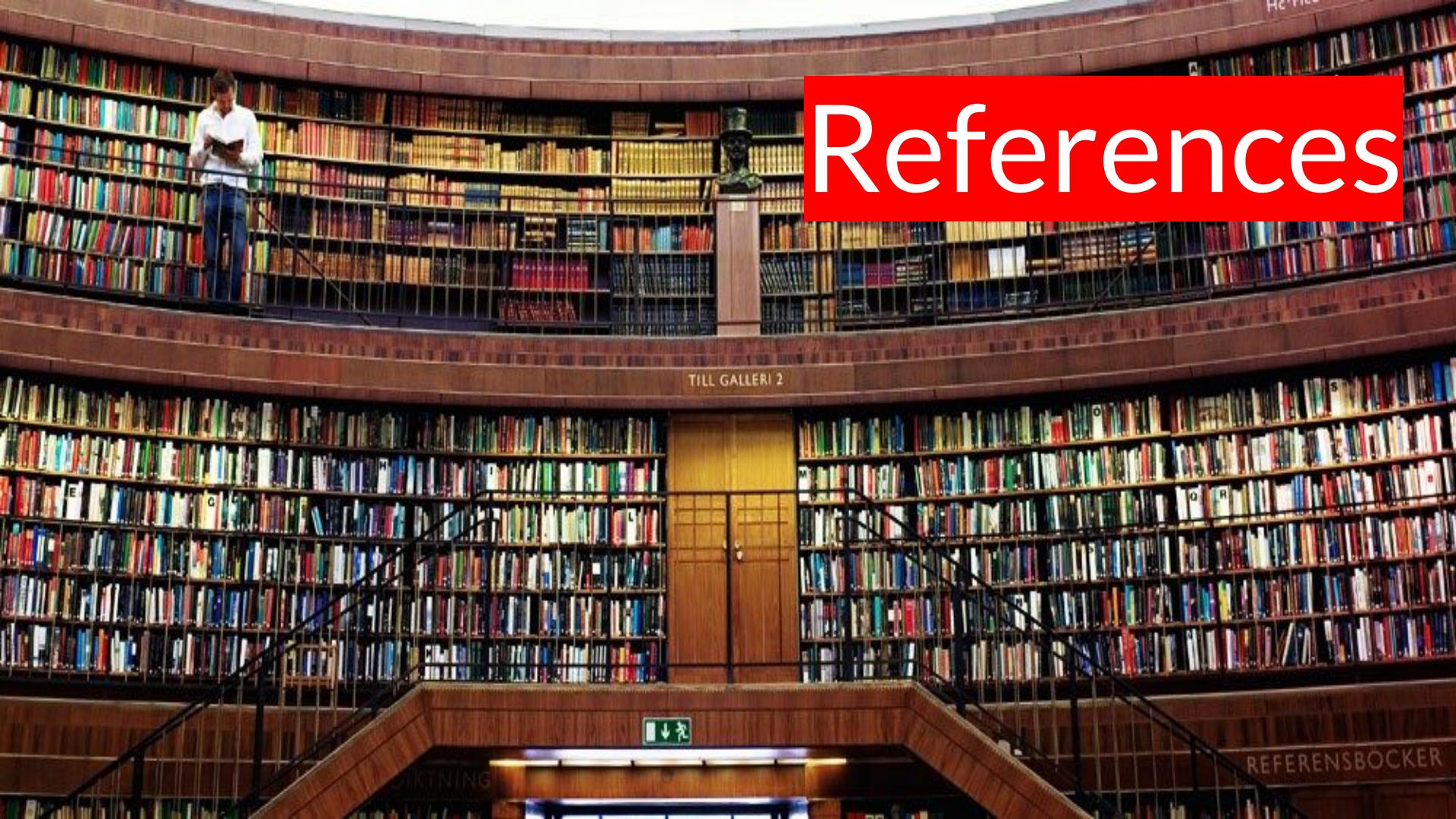


GitHub



Bitbucket

References



References Online



kaggle



DATAQUEST



DataCamp

Data Science
Academy

References

