

Data Science Foundations

Lesson #1 - Outline & Directions

Ivanovitch Silva
July, 2017





Introduction





Ivanovitch Silva (ivan@imd.ufrn.br)

Office: CIVT-A226 2T34, 6T34



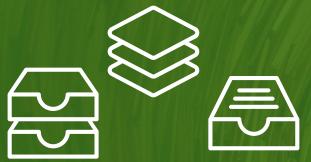
Group Knowledge



Take a Survey



<https://goo.gl/XmRwNj>



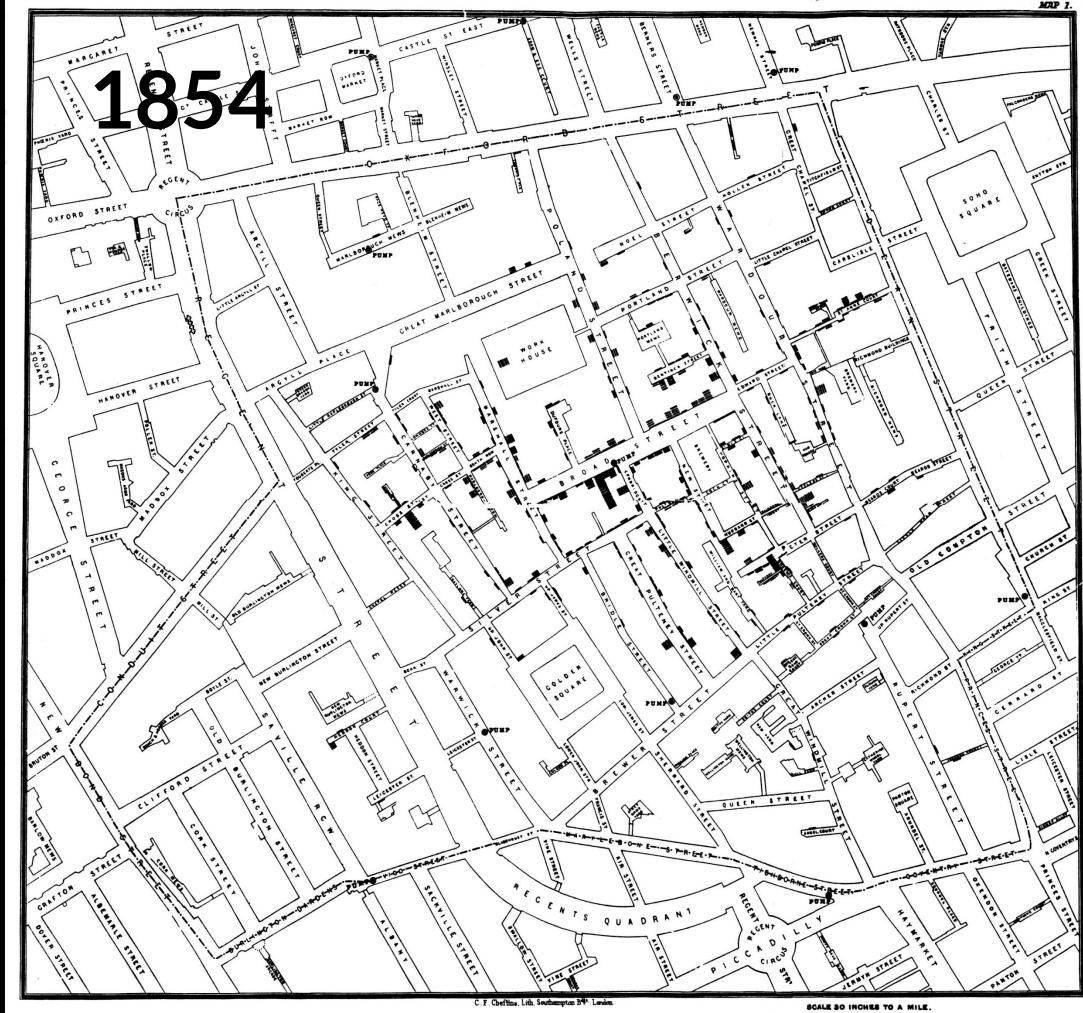
About Data



John Snow,
London, 1854



John Snow



(7-24-5)

Note A.—The Census Year begins June 1, 1880, and ends May 31, 1881.

Note B.—All persons will be included in the Enumeration who were living on the 1st day of June, 1880. No others will be included.

Note C.—Deaths since June 1, 1880, will be OMITTED. Members of Families who have DIED SINCE June 1, 1880, will be INCLUDED.

Note C—Questions Nos. 11, 12, 22 and 23 are not to be asked to persons under 10 years of age.

SCHEDULE I.—Inhabitants in [redacted], in the County of [redacted], State of Wisconsin.

[Signature] *Josephoff* Dated _____

No. of Persons	Name of Person	Relationship to Head of Family	Date of Birth		Color of Hair	Complexion	Markings	Occupation	Value of Real Estate	Value of Personal Estate	Name of School Attended	Name of Religious Denomination	Name of Physician	Name of Hospital		
			Month	Year												
23 (18)	Lark, Stephen	Bachelor	W	As 35	Blonde	1		Bookbinder								Rich (Rich)
	Burke, George	Bachelor	W	As 37	Blonde	1		Bookbinder								Rich (Rich)
		Bachelor	W	As 37	Blonde	1		Clothing Importer								Rich (Rich)
	Blome, Sophia	Wife	W	As 37	Blonde	1		Cook								Rich (Rich)
	Widener, E.	Wife	W	As 37	Blonde	1		Cook								Rich (Rich)
		Wife	W	As 37	Blonde	1		Cook								Rich (Rich)
	Hoppe, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, E.	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
		Bachelor	W	As 18	Blonde	1		Cook								Rich (Rich)
	Widener, Louis	Bachelor														

U.S. census (1880)

- 50 million people
 - Age, sex, occupation, education level
 - It took 8 years to be tabulate



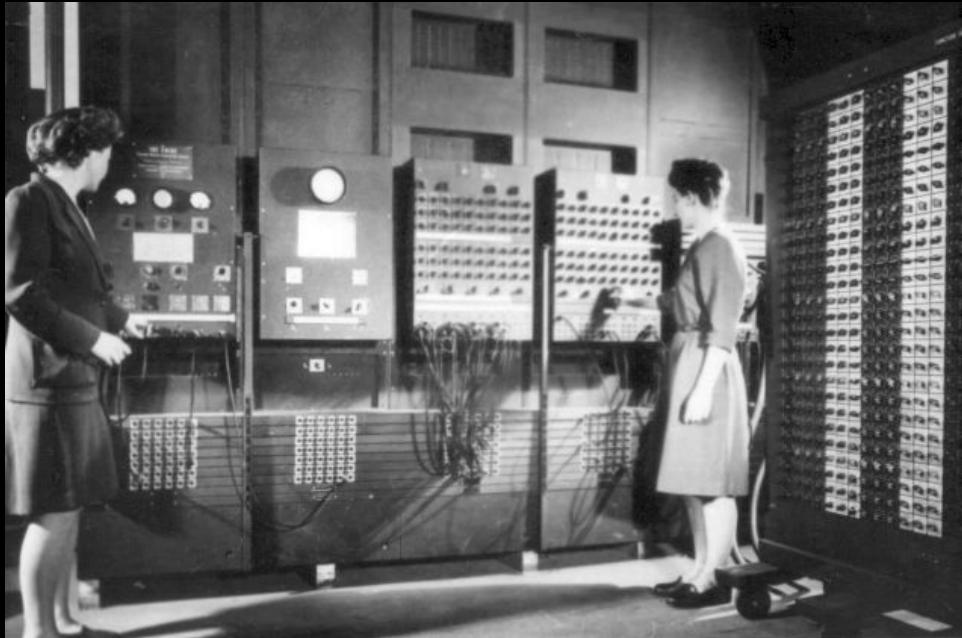
https://en.wikipedia.org/wiki/Tabulating_machine

1	1	3	0	2	4	10	On	S	A	C	E	a	c	e	g	EB	SB	Ch	Sy	U	Sh	Hk	Br	Rm
2	2	4	1	3	E	15	Off	IS	B	D	F	b	d	f	h	SY	X	Fp	Cn	R	X	Al	Cg	Kg
3	0	0	0	0	W	20			0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
A	1	1	1	1		0	25	A	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
B	2	2	2	2		5	30	B	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
C	3	3	3	3		0	3	C	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3	3
D	4	4	4	4		1	4	D	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4	4
E	5	5	5	5		2	C	E	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5	5
F	6	6	6	6		A	D	F	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6	6
G	7	7	7	7		B	E	Q	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7	7
H	8	8	8	8		a	F	H	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8	8
I	9	9	9	9		b	c	I	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9	9

https://en.wikipedia.org/wiki/1890_United_States_Census

U.S. census (1890)

- Tabulating machine.
- 63 million people
- The results were announced after only six weeks of processing.
- Without this invention, experts had estimated, the 1890 census would have taken 13 years to fully tabulate.



The population boom (1932)

- 123 million people.
- Information overload continued with the boom in the US population
- It demanded more thorough and organized record-keeping.



Scientific Knowledge Expands (1962)

- Derek Price publishes "[Science Since Babylon](#)"
- He concludes that the number of new journals has grown exponentially rather than linearly.
- This is now better known as the "[law of exponential increase](#)".



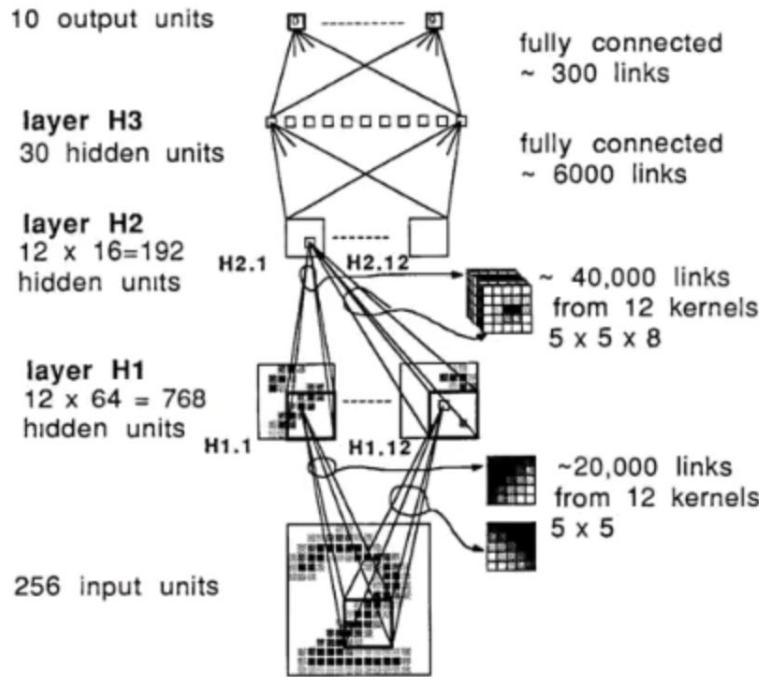
The Rise in Two-Way Communication (1975)

- The Information Flow Census, conducted by the Ministry of Posts and Telecommunications in Japan, started tracking the volume of information circulating in that country
- The study found that information supply greatly exceeded information consumption, and the demand for one-way communication had stagnated.



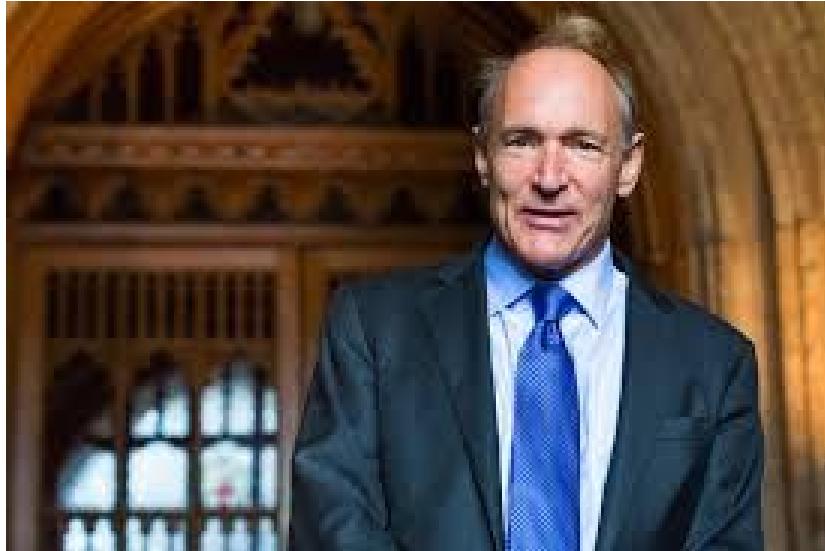
Information Growth and the Broadcasting Industry (1983)

- Companies were beginning to use data to provide answers for better business decisions.
- The massive information growth was credited to the expansion of the broadcasting industry.



Neural Networks debuted (1990)

Matan, Ofer, et al. Handwritten character recognition using neural network architectures. Proceeding of the 4th USPS Advanced Technology Conference.



The bird of the Internet (1991)

- Courtesy of Sr Tim Berners-Lee, data and information can now be posted online for the first time.



The future of data storage (1996)

- Digital storage became more cost-effective for store data than paper.
- Michael Lesk published [How much information is there in the world?](#)
So in only a few years, we will be able [to] save everything



Founded in
1996



The first time (1997)

- The term "Big Data" was used for the first time when researchers M. Cox and D. Ellsworth wrote an article identifying that the rise of data was becoming an issue for current computer systems.



Internet of Things (1999)

- The term “Internet of Things” was coined by British entrepreneur Kevin Ashton, Co-Founder of the Auto-ID Center at MIT, during a presentation linking the idea of Radio Frequency Identification (RFID) in supply chain to the internet world.

An Open Source Solution to the Big Data Explosion (2006)



- Hadoop was created in 2006 out of the necessity for new systems to handle the explosion of data from the web.
- Free to download, use, enhance and improve, Hadoop is a 100% open source way of storing and processing data that enables distributed parallel processing of huge amounts of data across inexpensive, industry-standard servers that both store and process the data, and can scale without limits".



Revolutionary Breakthroughs (2008)

A group of computer science researchers published a paper titled [Big Data Computing: Creating Revolutionary Breakthroughs in Commerce, Science, and Society.](#)

Big-data computing is perhaps the biggest innovation in computing in the last decade

A modest investment by the federal government could greatly accelerate its development and deployment.

Results for #IBMBigData

Save

Top / All

KurtMalueg @KurtMalueg · Feb 17
Looking forward to the tweets from #ibmbigdata!
Expand

Vikas K Manoria @vmanoria · Feb 16
ibm.co/18ksgII - #Bigdata architecture and patterns, Part 3: Understanding the architectural layers of big data solution #IBMBigData
Expand

Vikas K Manoria @vmanoria · Feb 16
ibm.co/16RlcZ8 - #Bigdata architecture and patterns, Part 2: How to know if a big data solution is right for your org. #IBMBigData
Expand

Nancy Kopp-Hensley @nancykoppdw · Feb 14
Followed by Francine Allaire
Big data and analytics transforms Nike Planning #IBMBigData shares/QVb6M via @sharethis
Expand

Marie Ma-Miller @TechMash365 · Feb 14
IBM + Big Data = lovefest
IBM says download 4K movie or 40,000 songs in seconds ibm.co/1gnwXgw #IBMBigData @TechMash365
View summary



#IBMBigData (2011)

IBM introduced a Twitter hashtag, #IBMBigData which expanded on their Big Data themed website that they built in 2008 in an effort to integrate it into their marketing.

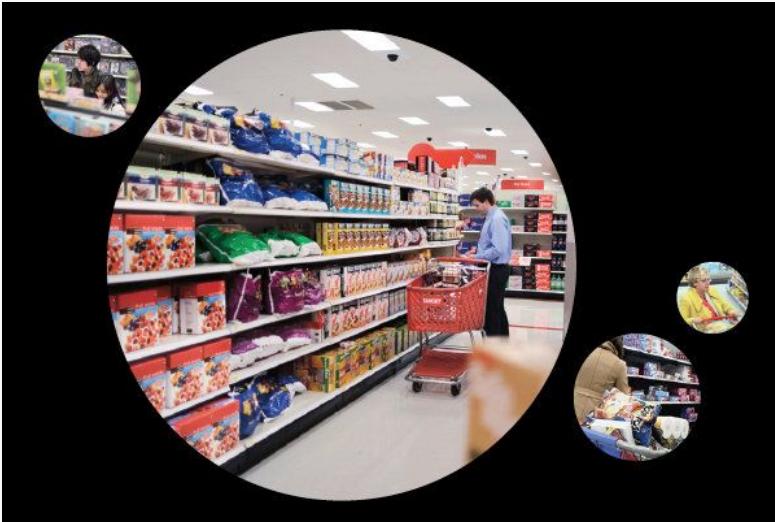




<http://www.zdnet.com/article/google-brain-simulator-teaches-itself-to-recognize-cats/>

Deep Learning experiment at Google (2012)

After seeing 10 million images from YouTube videos within three days, the 16,000-computer network, which had one billion connections, began to recognize cats, even though it had never been taught what a cat looked like.



http://www.nytimes.com/2012/02/19/magazine/shopping-habits.html?pagewanted=1&_r=1&hp



How Companies Learn Your Secrets (2012)

“My daughter got this in the mail!” he said. “She’s still in high school, and you’re sending her coupons for baby clothes and cribs? Are you trying to encourage her to get pregnant?”



#The Year of the Internet of Things (IoT) 2014

The IoT has become a powerful force for business transformation, and its disruptive impact will be felt across all industries and all areas of society.

According to Gartner, there were 3.7 billion connected "things" in use in 2014

Calvin Klein

Knits in myriad varieties, including curly mohair coats and hand-stitched, multipaneled sweaters, were the focus of a disciplined but cozy collection in a soothing palette of earth tones and snow.



Mixed-knit jumpers blocked out in white, black and gray

Turtleneck sweater tops with chunky, sampler-scarf knit panels

Ralph Lauren

A collection distinguished by its sure-handed hybrid of refinement and ease, its muted pastels and buttery fabrics arguing for opulence but in no way overstating the case.

Read more: [Ralph Lauren Plays Polo](#)

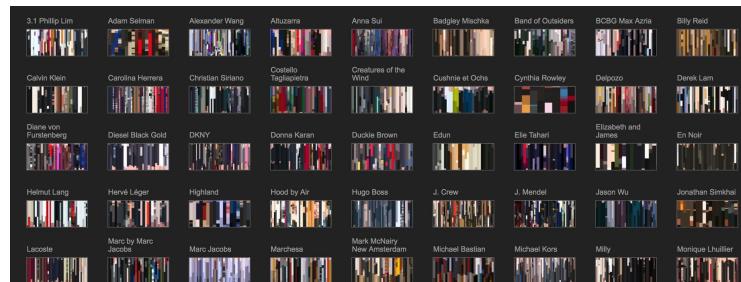


An upscale hippie dress under a distressed bomber

A massive, asymmetrical cape in soft pink



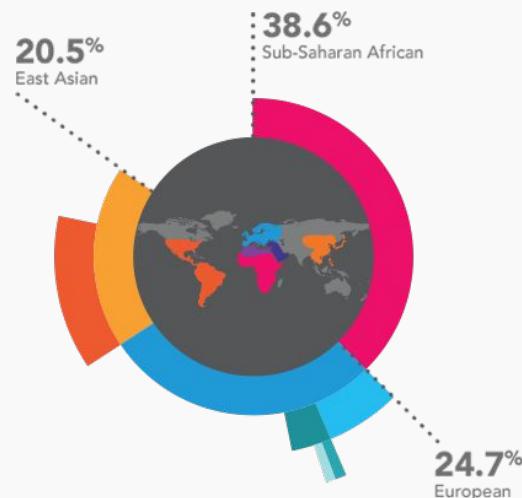
2014





23andMe

<https://www.23andme.com/>

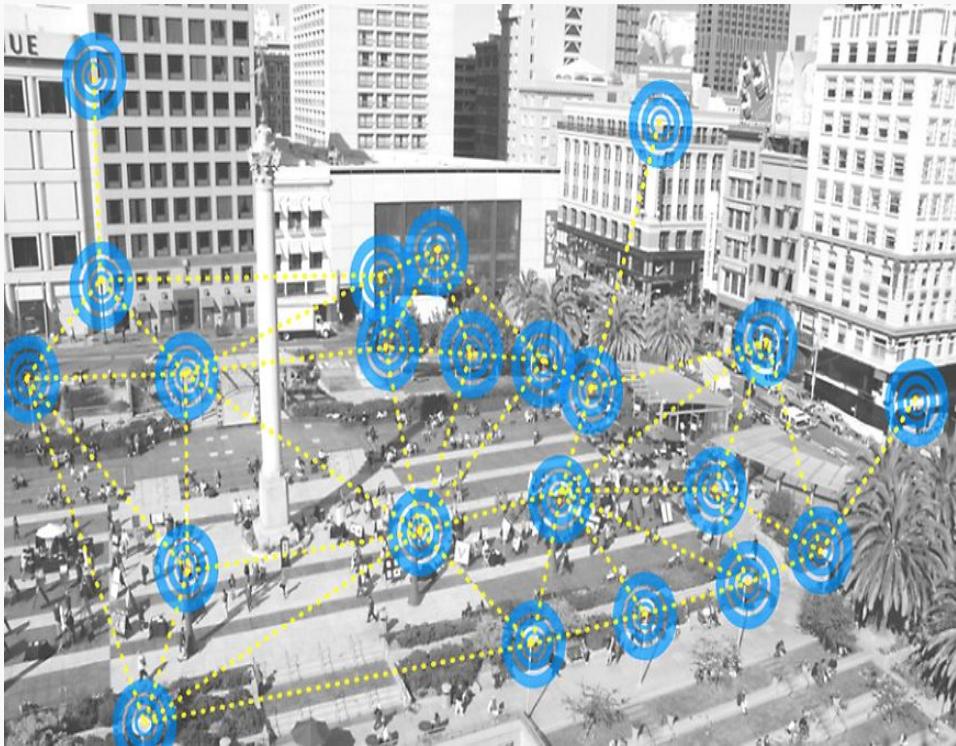


23andMe has been selling a product with both ancestry and health-related components in Canada since October 2014



<http://www.telegraph.co.uk/technology/news/10959864/Germanys-World-Cup-tactics-shaped-by-data.html>

Germany's 12th Man at the World Cup (2014)



<http://www.forbes.com/sites/danielnewman/2016/08/15/big-data-and-the-future-of-smart-cities/#16c07f1d3f2d>

Connect To The Cities Of The Future (2015)

With the growth of our population and the advent of ideas such as big data and the Internet of Things, the natural step cities will take is to become more interconnected.

Not only will this result in brighter streets, but the new lights will also be an interconnected system that will inform the city of each bulb's status.

2016



<http://www.pokemongo.com/>



<https://www.youtube.com/watch?v=NrmMk1Myrxc>



Why Netflix thinks its personalized recommendation engine is worth \$1 billion per year (2016)

"Consumer research suggests that a typical Netflix member loses interest after perhaps 60 to 90 seconds of choosing"

The user either finds something of interest or the risk of the user abandoning our service increases substantially

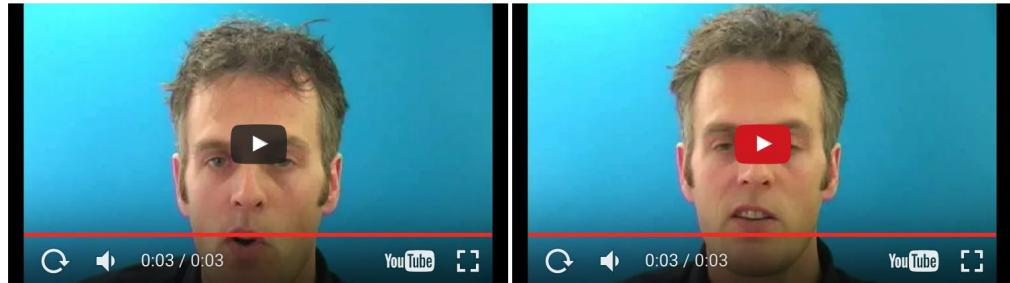


THE YEAR OF INTELLIGENCE



An app for blind people identifies and reads out objects in their surroundings





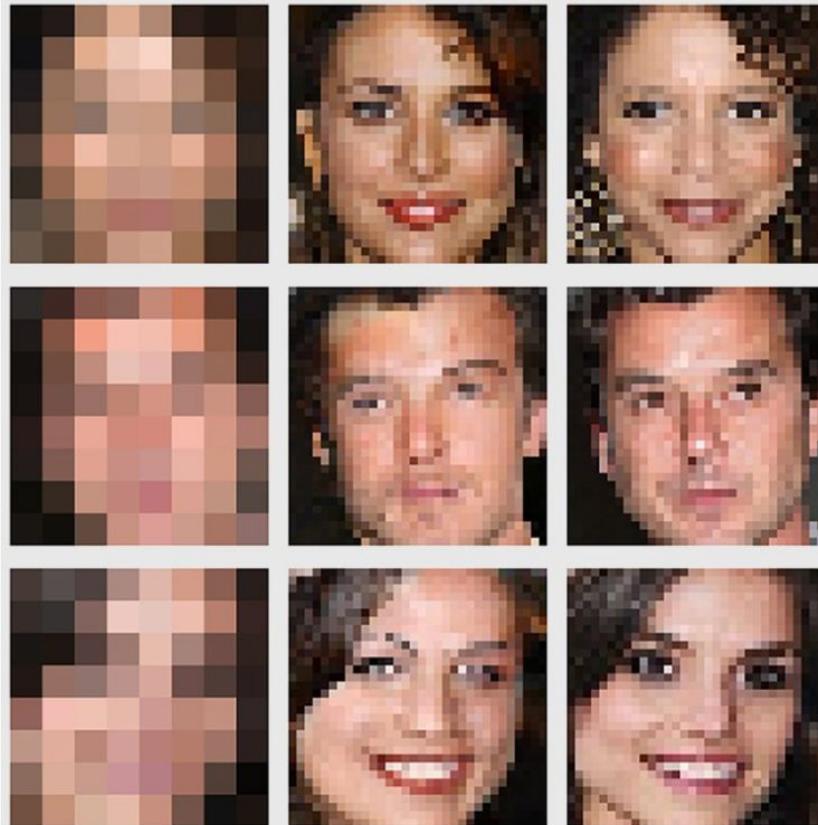
S4



Vid2Speech: Speech Reconstruction from Silent Video

<https://github.com/arielephrat/vid2speech>

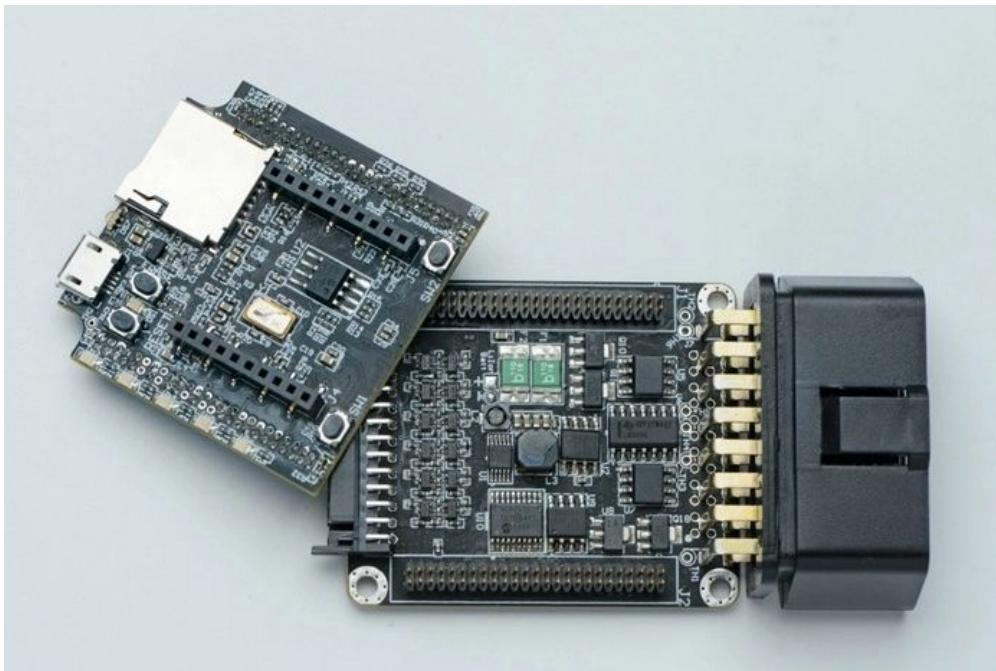
<https://goo.gl/HPYCbK>



Google Uses AI to 'Rebuild' a Portrait from an 8x8-Pixel Image

<http://www.theverge.com/2017/2/7/14532206/google-brain-research-neural-networks-zoom-and-enhance-pixelated-images>

IoT + Car Hacking + Big Data



<http://hackaday.com/2017/02/21/first-look-macchina-m2/>

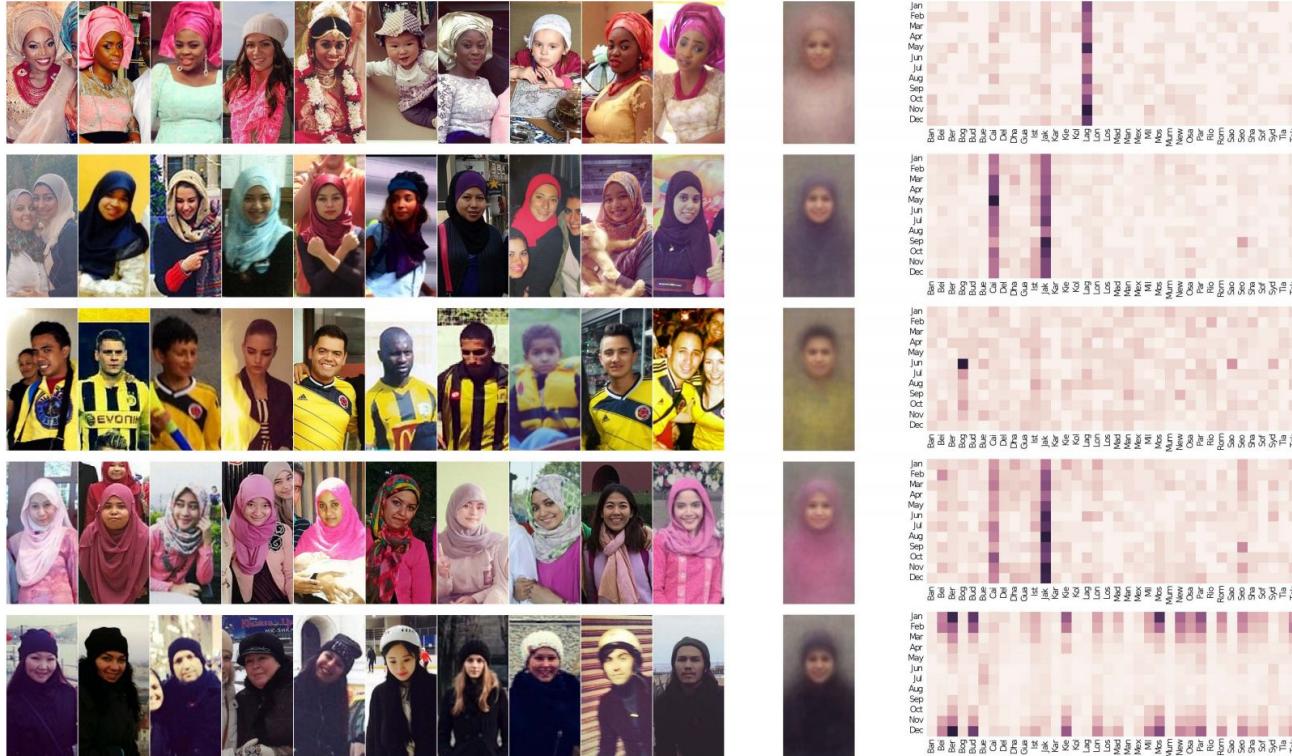
KICKSTARTER

<https://www.kickstarter.com/projects/1029808658/macchina-the-ultimate-tool-for-taking-control-of-y?token=e40d156d>



<https://comma.ai/>

StreetStyle: Exploring world-wide clothing styles from millions of photos



<https://goo.gl/i9rDJX>

Artificial intelligence can now predict suicide with remarkable accuracy



<https://goo.gl/1b1Mr9>



HOSPITAL SÍRIO-LIBANÊS

KUNUMI

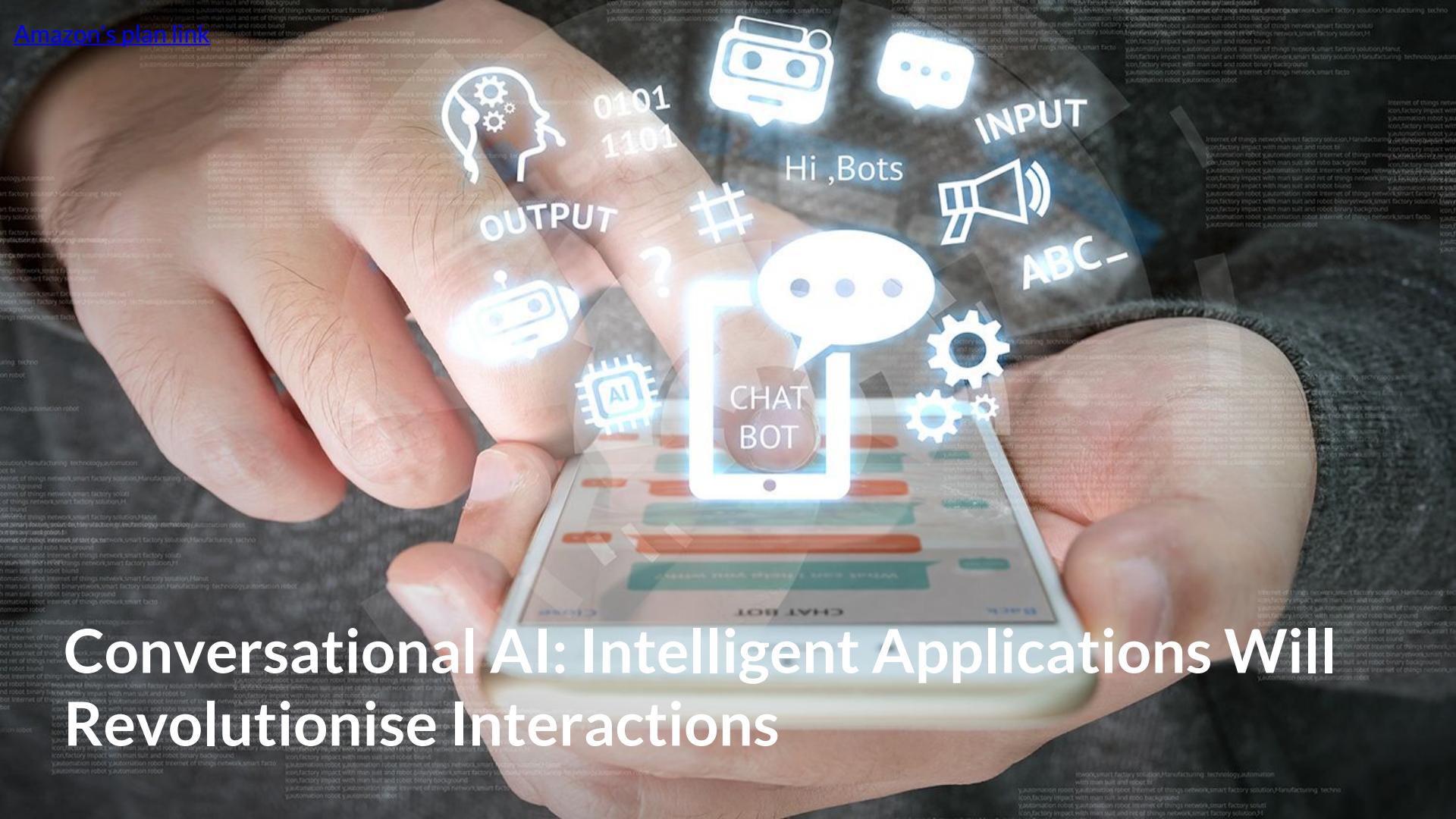
<http://bhetc.org.br/empresas-do-bhetc/kunumi/>



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

[Amazon's plan link](#)

Conversational AI: Intelligent Applications Will Revolutionise Interactions



A Maior Rede de Anúncios Mobile do Brasil

Dia 19 de junho nos encontramos em Cannes

[Assista ao vídeo](#)

BIG DATA & INTELIGÊNCIA ARTIFICIAL

TRANSFORME DADOS EM
INFORMAÇÃO

QUEM SOMOS

NOSSAS SOLUÇÕES

DOMINE O MERCADO COM O PROSPECTA

Tecnologia e conteúdo para o crescimento da sua empresa.
Inteligência para você obter o potencial máximo do mercado.

[DESCUBRA](#)



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



Big Data





68 likes

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

150 likes

- ++ family member

300 likes

- ++ spouse

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

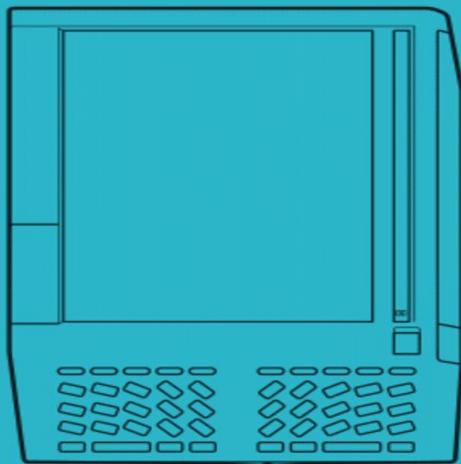
BIG DATA

BIG NUMBERS



MEGABYTES

1 MEGABYTE – APPROXIMATELY 1,000 KILOBYTES
(ACTUALLY 1,024 KB)



256MB
Kindle
first generation



0.004MB
Apple I
RAM in Apple's
first computer, 1976



0.004
Oyster card
London public transport



0.02
Punched
paper tape
largest feasible reel



0.02
Tandy 200 computer
amount of RAM, 1984



0.02
Word document
single page



0.7
Audio cassette
90 min



5
William
Shakespeare
complete works



2.5
War & Peace
Kindle ebook



1.5
3.5" Floppy disk



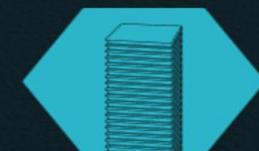
1
ebook
average size



9
Blu-Ray
1 sec at HDTV quality



66
Mosquito genome
DNA for malaria mosquito



20
10,000 pages of text
about 20 novels



0.25GB

Kindle

first generation



0.47GB

**Large Hadron
Collider**

data produced per sec



0.48

YouTube
video uploaded
per sec, 2012



0.51

Raspberry Pi
Model B



0.54

Blu-ray

1 min at HDTV quality



0.76

Single human sperm
all DNA



0.76

Single human egg
all DNA



0.7

CD

80 mins



1

MiniDisc
45 hrs of music



1.5

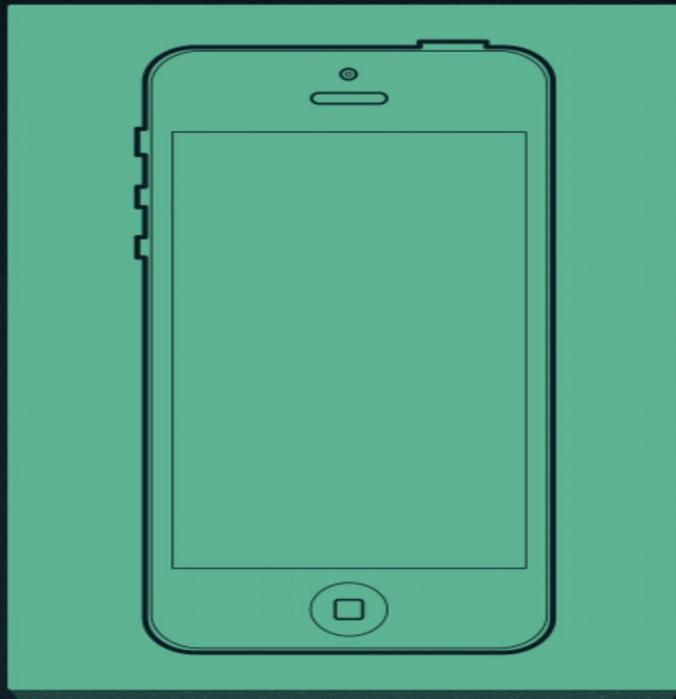
Human body cell
all DNA

0.07
coding
DNA only



2

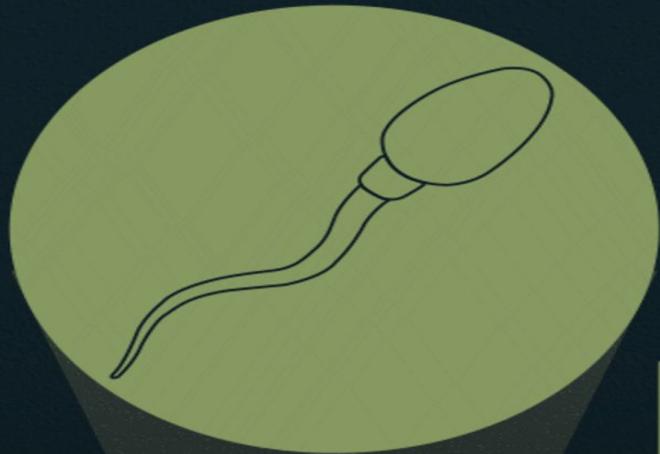
Kindle Paperwhite
total storage



32GB
iPhone 5
total storage

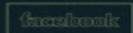
GIGABYTES

1 GIGABYTE - APPROXIMATELY 1,000 MEGABYTES
(ACTUALLY 1,024 MB)



778

The Hobbit
4K digital cinema
high frame rate



3GB
Facebook
photos and videos
stored per sec, 2012



4
Kindle Touch

4k digital cinema



2.3
1 min normal
frame rate
4.6
1 min high
frame rate



26.4
Mad Men
one episode



7
Wikipedia
all current articles
without edit history



32GB
iPhone 5
total storage



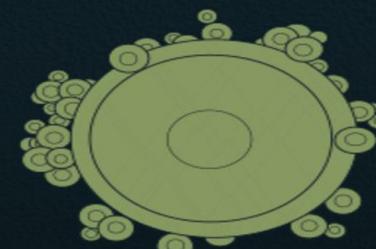
8.5
DVD
dual-layer



20
Wolfgang Amadeus Mozart
complete works
as MP3



128
Blu-ray
max disc capacity



343
Eggs per woman
at birth - about 450

CHANGE
OF SCALE
10:1



1TB
Average
modern
hard disk



1.8TB
Human sperm
DNA created
per man, per sec



1
million novels
500 million pages of text



1.3
Human brain
functional memory
capacity



7
single DNA
sequencing run
data from end-to-end
human DNA sequencing



7.3
Wikipedia
all current articles
with edit history



10
US Library
of Congress
printed collection

Internet traffic



12
for all of 1990

6.3
per sec, 2012

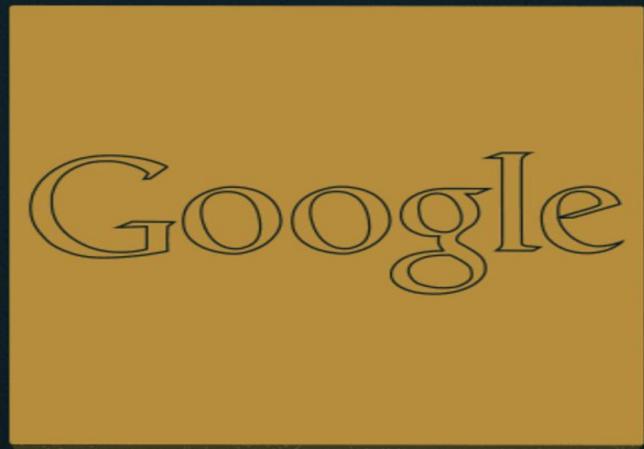
TERABYTES

1 TERABYTE - APPROXIMATELY 1,000,000 MEGABYTES
(ACTUALLY 1,048,576 MB)

facebook

0.18
per min **11**
per hr

274TB
Facebook
photos and videos
stored per day, 2012



20,000TB
Google
data processed
per day 2008



2,220
Synthetic DNA data
storage capacity 1 gram

amazon.com

42TB
Amazon.com
database



47

Human ear
DNA in audio hair cells



304
Human eye
light receptors
per sq mm



651
Eagle's eye
light receptors
per sq mm

CHANGE
OF SCALE
10:1



274TB
Facebook
photos and videos
stored per day, 2012



120
Internet traffic
for all of 1993



900
Mobile internet
traffic per month, 2005



2,000
All US academic
research libraries
printed collection



1,800
Emails sent globally
per day, 2002

PETABYTES

1 PETABYTE – APPROXIMATELY 1,000,000,000 MEGABYTES
(ACTUALLY 1,073,741,824 MB)



3.5PB

Radio programming

globally 2002



16

including repeats



40

Titan supercomputer

storage capacity



70

TV programming

globally 2002



20PB

Google

data processed per day 2008



6

World's largest climate data archive

NOAA National Climatic Data Center



15

Large Hadron Collider

data produced per year



22.8

Internet traffic

for all of 1996



47.4

Telephone calls

globally per day, 2002



45.7

Human nose

smell receptor neurons



152
Human sperm
created per man,
per day



211
Human male ejaculation

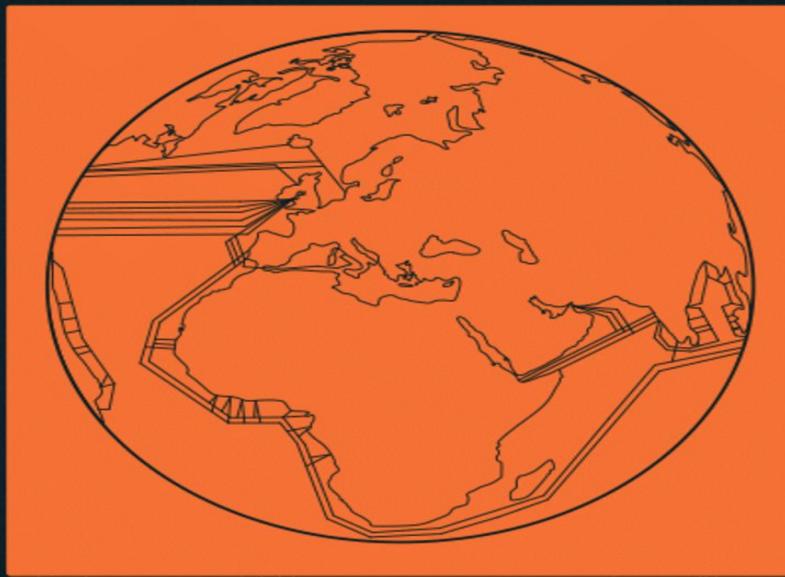


2,400PB
Human skin cells
shed in a month

DIGITAL ANALOGUE ORGANIC

EXABYTES

1 EXABYTE - APPROXIMATELY 1,000,000,000,000 MEGABYTES
(ACTUALLY 1,099,511,627,776 MB)



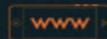
CHANGE
OF SCALE
10:1



2.4EB
human skin cells
shed in a month



0.2
All printed
material on Earth



0.59
Mobile internet
traffic per month, 2011



0.2EB
Human eye
light receptors in retina

Internet
traffic

0.5
per day
2012

1
per month,
2004

4.9
for all of 2002



0.66
Emails sent
globally 2002



1.5
Dog's nose
smell receptor neurons



56
Human sperm
created per man, per year



0.27
TV programming
globally including
repeats, 2002



17.3
Telephone calls
globally for all of 2002

CHANGE
OF SCALE
10:1



0.01 zB
YouTube
hours watched
for all of 2012



0.3 zB
Internet traffic
for all of 2011



0.1
for all of 2008



0.02
for all of 2005



0.05
Neurons in
human brain

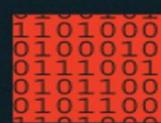
All global data



0.5
2008



0.8
2009



1.2
2010



1.8
2011



2.8
2012

ZETTABYTES

1 ZETTABYTE = APPROXIMATELY 1,000,000,000,000,000 MEGABYTES
(ACTUALLY 1.125,899,906,842,620 MB)



30 zB
All global data
2019 (predicted)

YOTTABYTES

1 YOTTABYTE – APPROXIMATELY 1,000,000,000,000,000,000 MEGABYTES
(ACTUALLY 1,152,921,504,606,850,000 MB)

CHANGE
OF SCALE
1000:1



10,000
All microbes on Earth
unique genetic information

0.03
YB
All global data
2019 (predicted)



0.04
All words
ever spoken
digitized as 16 kHz
16-bit audio



0.09
All cells in
human body
duplicated genetic
information

NOTES

Figures are based on decimal not binary file sizes.
So 1 megabyte = 1,000 kilobytes, not 1,024 kilobytes.

Most organic figures refer to genetic data and are based on multiplying DNA in single cell by number of cells. DNA in two cells is therefore counted twice, though cells within an organism are genetically exact or near-exact copies of one another.

By Information is Beautiful Studio for **FUTURE**

Executive Creative Director David McCandless Creative Director Duncan Swain
Design Matt McLean Research Miriam Quick, Christian Miles

BBC

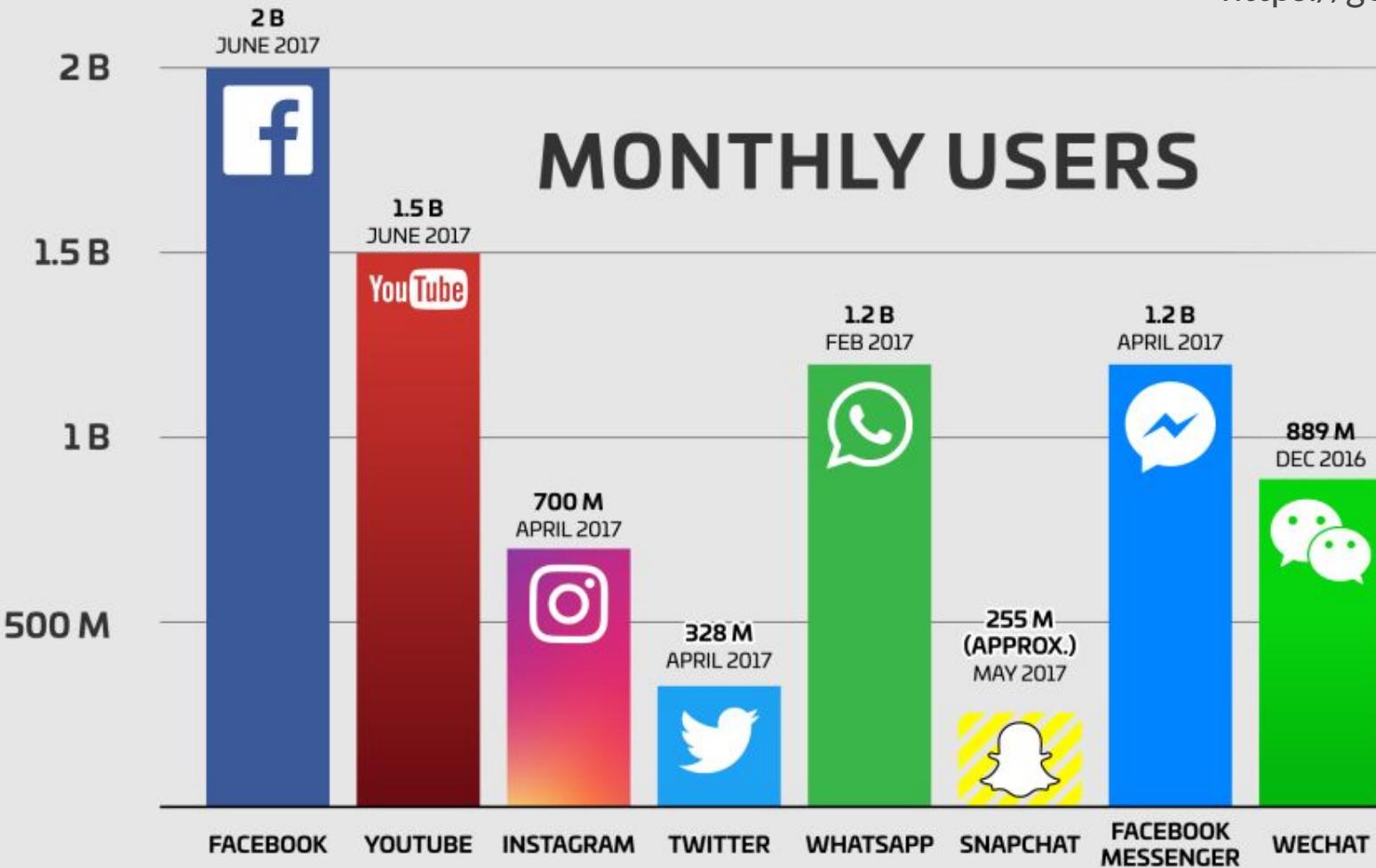
BIG
SCIENCE





BIG

SOCIALMEDIA



BIG

CHALLENGES

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

AUTONOMOUS VEHICLES
4,000 GB
PER DAY... EACH DAY

LIDAR
~10-70 MB
PER SECOND



● Big Data

Search term

+ Compare

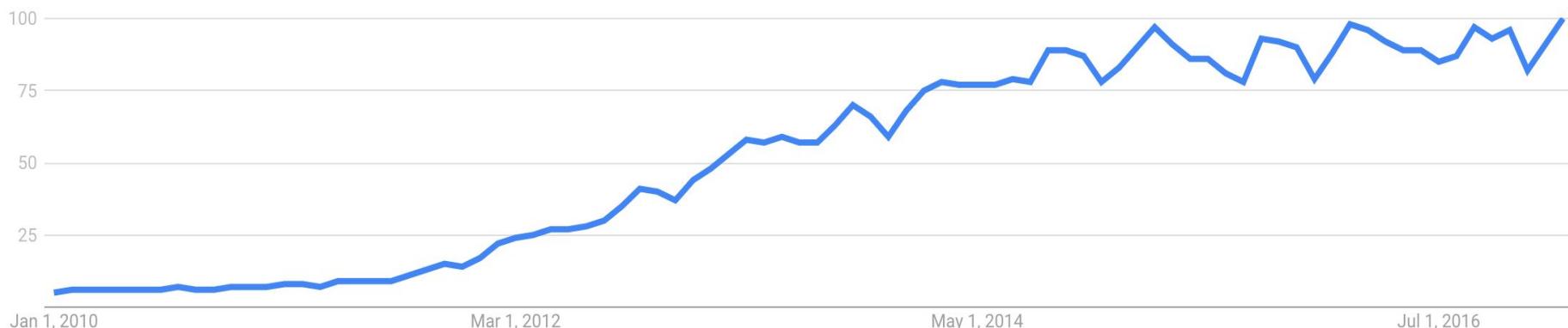
Worldwide ▾

1/1/10 - 2/19/17 ▾

All categories ▾

Web Search ▾

Interest over time ?



BIG DATA LANDSCAPE 2017



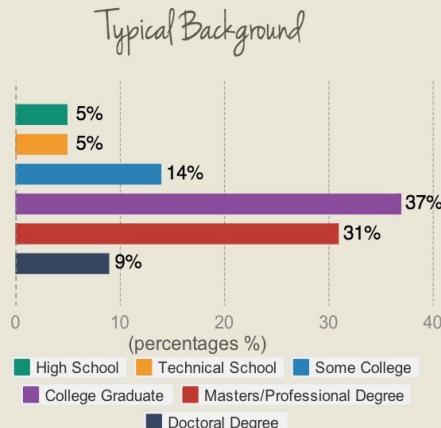
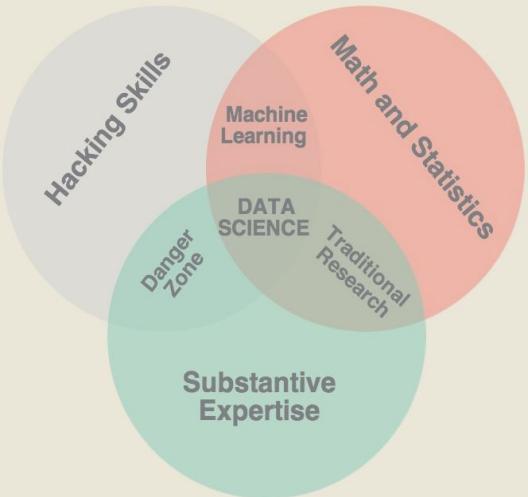
Who studies this stuff?



Data Scientist

in 8 easy steps

What's a data scientist?



A data scientist is someone who is better at statistics than any software engineer and better at software engineering than any statistician.

MODERN DATA SCIENTIST

Data Scientist, the sexiest job of the 21st century, requires a mixture of multidisciplinary skills ranging from an intersection of mathematics, statistics, computer science, communication and business. Finding a data scientist is hard. Finding people who understand who a data scientist is, is equally hard. So here is a little cheat sheet on who the modern data scientist really is.

MATH & STATISTICS

- ★ Machine learning
- ★ Statistical modeling
- ★ Experiment design
- ★ Bayesian inference
- ★ Supervised learning: decision trees, random forests, logistic regression
- ★ Unsupervised learning: clustering, dimensionality reduction
- ★ Optimization: gradient descent and variants



PROGRAMMING & DATABASE

- ★ Computer science fundamentals
- ★ Scripting language e.g. Python
- ★ Statistical computing packages, e.g., R
- ★ Databases: SQL and NoSQL
- ★ Relational algebra
- ★ Parallel databases and parallel query processing
- ★ MapReduce concepts
- ★ Hadoop and Hive/Pig
- ★ Custom reducers
- ★ Experience with xaaS like AWS

DOMAIN KNOWLEDGE & SOFT SKILLS

- ★ Passionate about the business
- ★ Curious about data
- ★ Influence without authority
- ★ Hacker mindset
- ★ Problem solver
- ★ Strategic, proactive, creative, innovative and collaborative

COMMUNICATION & VISUALIZATION

- ★ Able to engage with senior management
- ★ Story telling skills
- ★ Translate data-driven insights into decisions and actions
- ★ Visual art design
- ★ R packages like ggplot or lattice
- ★ Knowledge of any of visualization tools e.g. Flare, D3.js, Tableau

**Harvard
Business
Review**

Data Scientist: The Sexiest Job of the 21st Century

Become a Data Scientist in 8 easy steps

1 Get good at stats, math and machine learning

Math



- > Math Track of Khan Academy
- > Linear Algebra by MIT OpenCourseware



Stats



- > Intro to Statistics by Udacity
- > OpenIntro Statistics



ML



- > Machine Learning by Andrew NG (Stanford Online)
- > Practical Machine Learning by John Hopkins (Coursera)

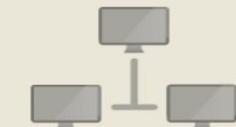
2 Learn to code



Computer Science Fundamentals
> CS50x on edX



Grasp end-to-end development
The things you build will be integrated
into other systems



Choose a first language
> Open Source: R, Python, etc.
> Commercial: SAS, SPSS, etc.



Learn Interactively
> R: DataCamp, tryR
> Python: Codecademy, Google Class



3 Understand databases

As a data scientist student, you will often work with data in text files. However, once you enter the industry, a database is almost always used to store data. It's going to be stored in MySQL, Postgres, MongoDB, Cassandra, etc.



4 Master data munging, visualization and reporting

□ Data cleaning and munging



WHAT

Data munging is the process of converting one "raw" form into another format for more convenient consumption



TOOLS

> Getting and Cleaning data by John Hopkins (Coursera)

DataWrangler alpha

 **data.table**
dplyr

□ Data visualization



WHAT

Data visualization involves the creation and study of the visual representation of data.



TOOLS

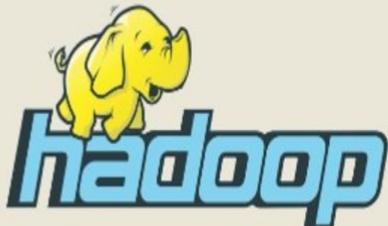
ggvis 

 **vega**

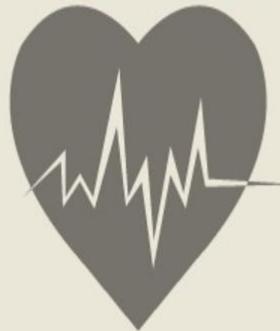
5 Level up with Big Data

When you start operating with data at the scale of the web, the fundamental approach and process of analysis must change. Most data scientists are working on problems that can't be run on single machines. They have large data sets that require distributed processing.

Hadoop is an open-source software framework for storage and large-scale processing of data-sets on clusters of commodity hardware.



MapReduce



MapReduce is this programming paradigm that allows for massive scalability across the servers in a Hadoop cluster.

Apache Spark is Hadoop's speedy Swiss Army knife. It is a fast-running data analysis system that provides real-time data processing functions to Hadoop.



6

Get experience, practice and meet fellow data scientists

Practice makes perfect ...



kaggle

join in
competitions



Meet fellow data
scientists



Have a pet
project



Develop your
intuition

7 Internship, bootcamp or get a job

The best way to find out whether you are a true data scientist or not is to take the bull by the horns and to enter the real-life jungle of data-analysis and science with your freshly acquired skill set.

Internship



BEGINNER

Bootcamp



INTERMEDIATE

Job



ADVANCED

amazon.com



8 Follow and engage with the community

Sites to follow

- > DataTau
- > Kdnuggets
- > fivethirtyeight
- > datascience101
- > r-bloggers

People to follow

- > Hilary Mason
- > David Smith
- > Nate Silver
- > dj patil

Need Data?



<http://mariofilho.com/>

Mistakes to avoid when starting your career in Data Science

While learning Data Science

1. Spending too much time on theory
2. Coding too many algorithms from scratch
3. Jumping into the deep end



<https://goo.gl/CvLVpu>

When applying for a job

1. Having too much technical jargon in a resume
2. Overestimating the value of academic degrees
3. Searching too narrowly
4. Being unprepared to discuss projects

Future or Present?

While studying at Academy, I built ...

Models

BIG DATA



Volume

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



Variety

Data



Data Files
(XML, CSV, Excel, JSON, ...)



Database
(MySQL, Oracle, ...)



API



Sites



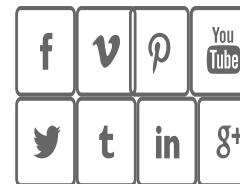
Text and reports



Maps



Image and videos

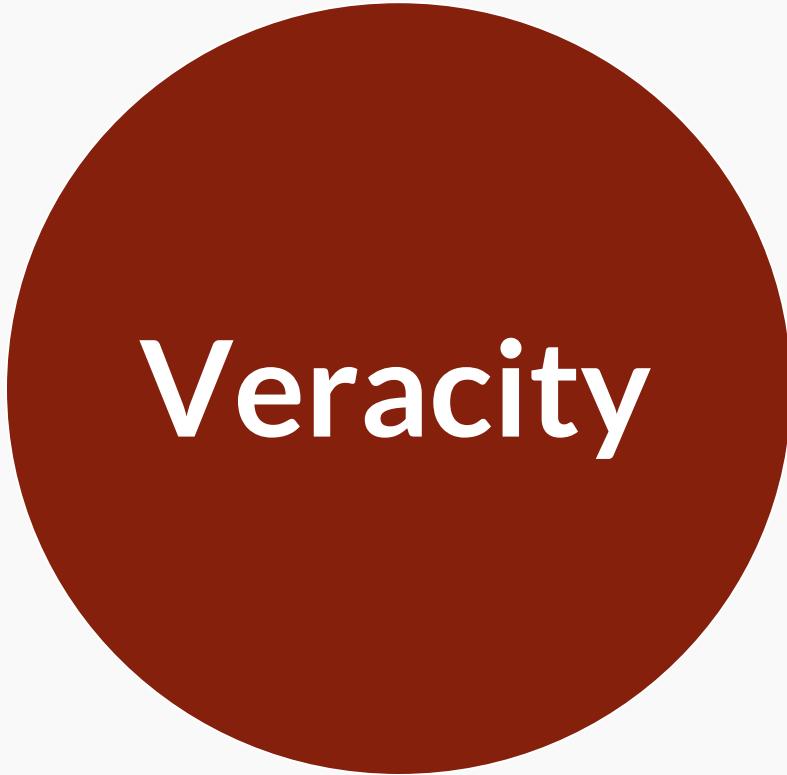


Social Media

Velocity

2017 This Is What Happens In An Internet Minute



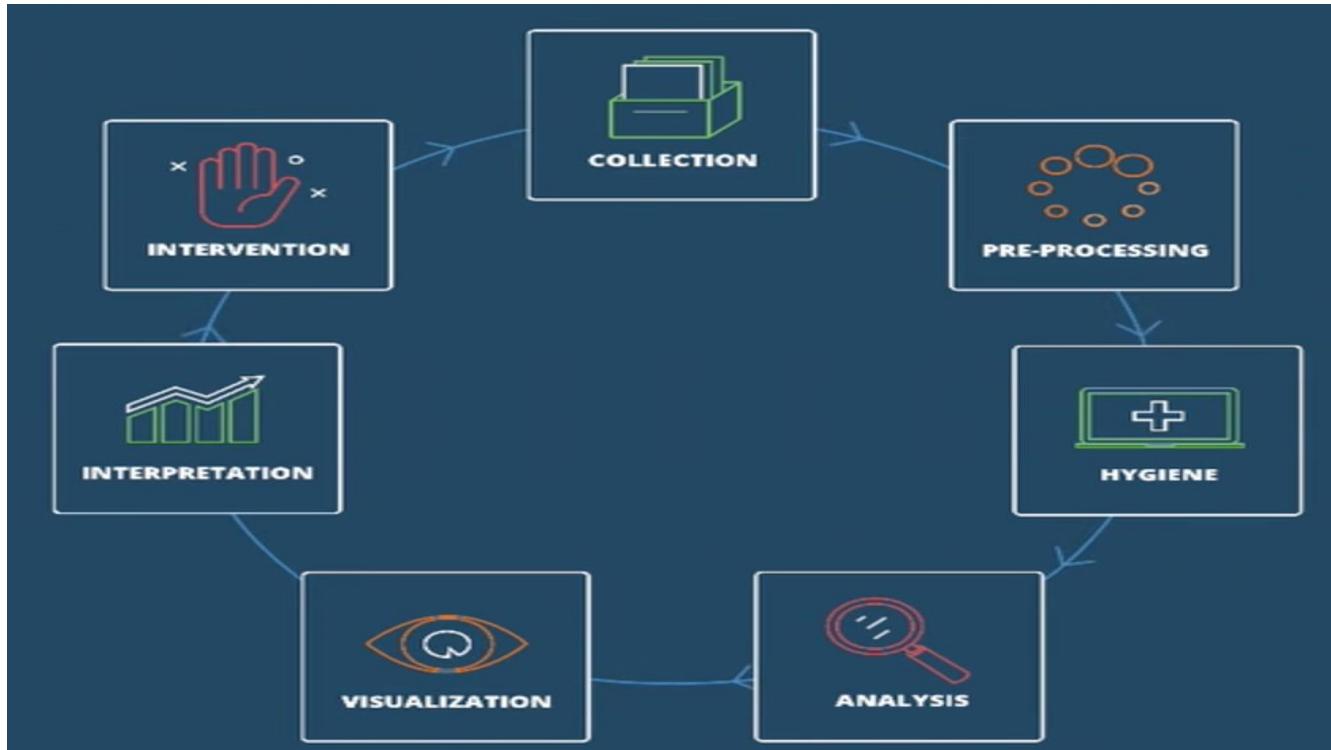


Veracity





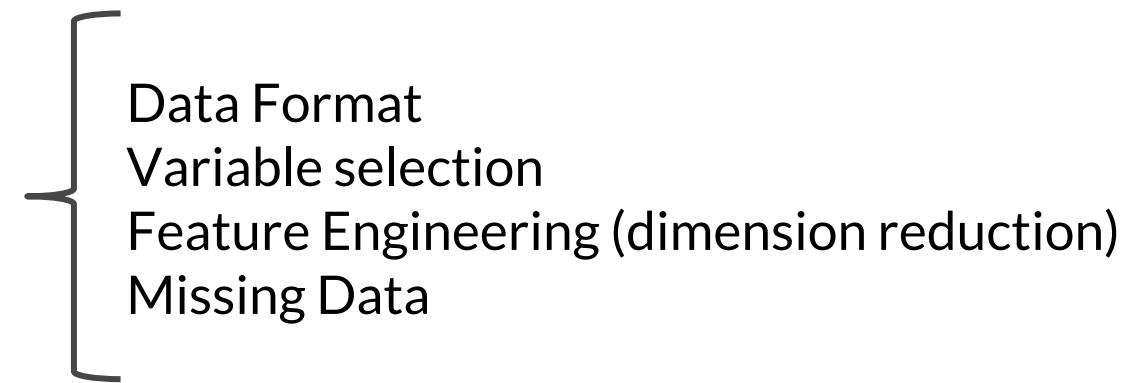
Key steps in a Data Science Analysis



Key steps in a Data Science Analysis

Data Preparation

- Pre-processing
- Hygiene



Key steps in a Data Science Analysis

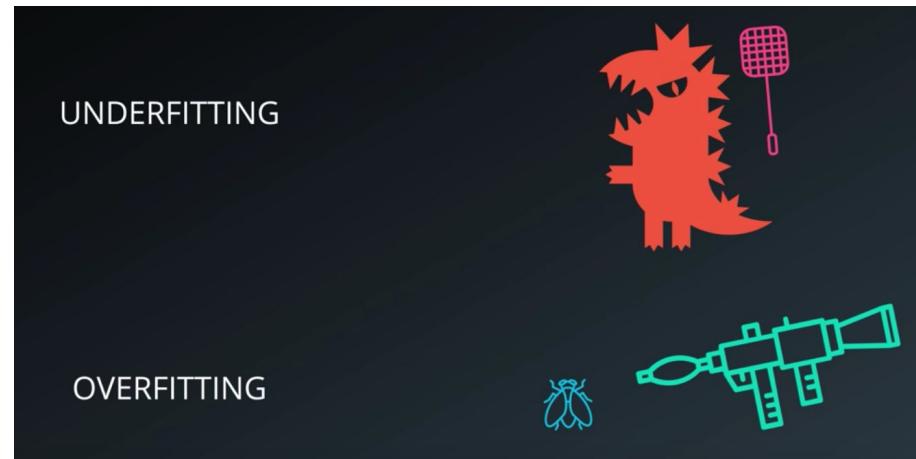
Algorithms Selection

- Analysis

Algorithms	
Unsupervised Learning	k -Means Clustering Principal Component Analysis Association Rules Social Network Analysis
Supervised Learning	Regression Analysis k -Nearest Neighbors Support Vector Machine Decision Tree Random Forests Neural Networks
Reinforcement Learning	Multi-Armed Bandits

Key steps in a Data Science Analysis

Parameter tuning
● Analysis

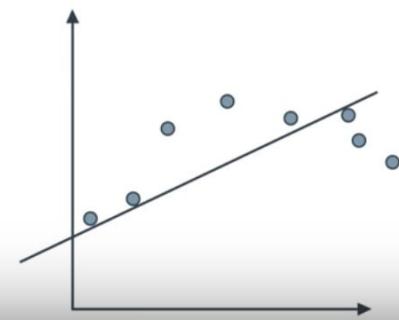
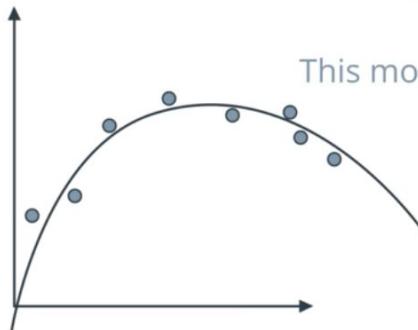


Key steps in a Data Science Analysis

○ UNDERFITTING

Error due to bias

This model will not do well in the training set

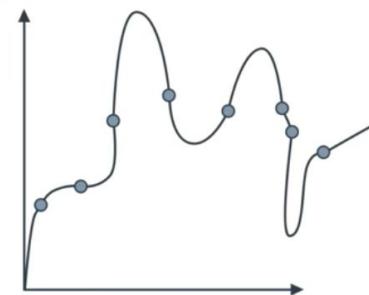
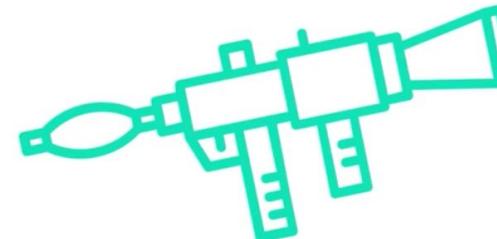
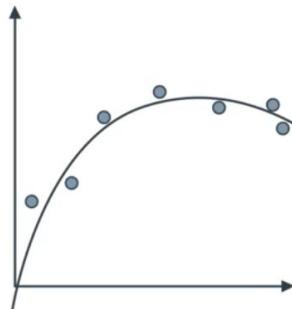


Key steps in a Data Science Analysis

OVERFITTING

Error due to variance

This model performs poorly in the testing set



Key steps in a Data Science Analysis

CONFUSION MATRIX

Evaluating Results

- Visualization
- Interpretation



10, 000 PATIENTS

PATIENTS	DIAGNOSIS	
	Diagnosed Sick	Diagnosed Healthy
Sick	1000 True Positives	200 False Negatives
Healthy	800 False Positives	8000 True Negatives

DECK.GL

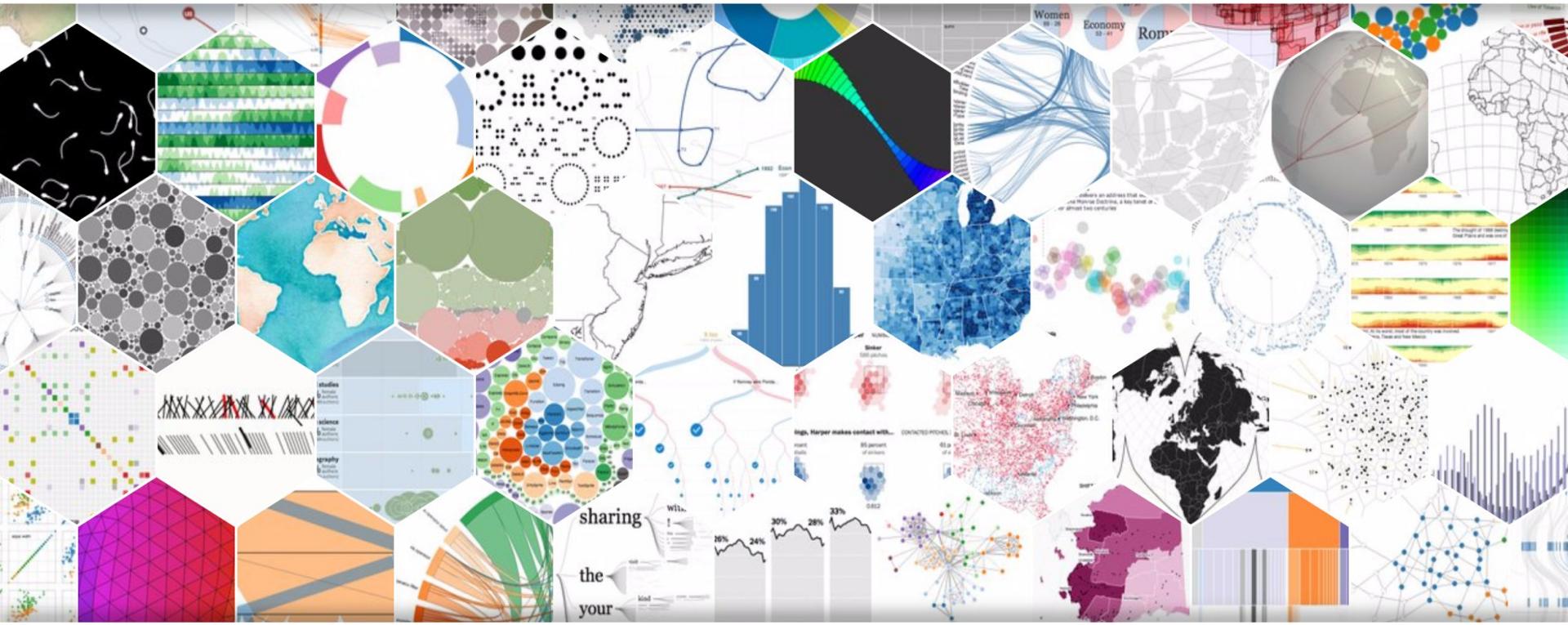
Large-scale WebGL-powered Data Visualization

GET STARTED

30 FPS (4-34)



DB Data-Driven Documents



Search and add articles



fly

map

home

help

about

UI off

full

Article links

Pick an article

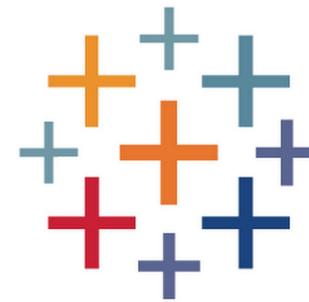
Welcome.

WikiGalaxy is a 3D web experiment that visualizes Wikipedia as a galactic web of information. With it I aim to show the world the beauty and variety of knowledge that is available at our fingertips.

I used 100,000 of 2014's most popular articles, all clustered with hyperlinks. In this world Wikipedia articles are stars, interests are nebulas and you are on a journey through knowledge.

<http://wiki.polyfra.me/>

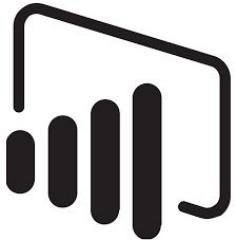
Use the mouse to see a preview of articles in each cluster
Click anywhere on the map to fly there



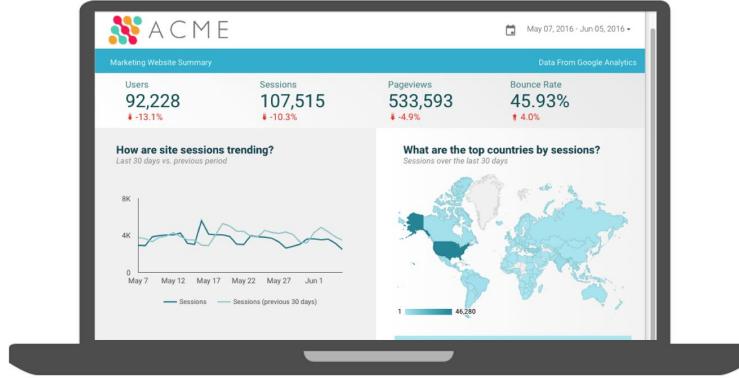
<http://www.pentaho.com/>



<https://www.tableau.com/>



<https://powerbi.microsoft.com>

A laptop screen displaying the Google Data Studio interface. The dashboard is titled "ACME Marketing Website Summary" and includes a "Data From Google Analytics" section with metrics like Users (92,228), Sessions (107,515), Pageviews (533,593), and Bounce Rate (45.93%). It also features two charts: one showing site sessions trending over the last 30 days and another showing the top countries by sessions with a world map.

[https://www.google.com.br/analytics/
data-studio/](https://www.google.com.br/analytics/data-studio/)



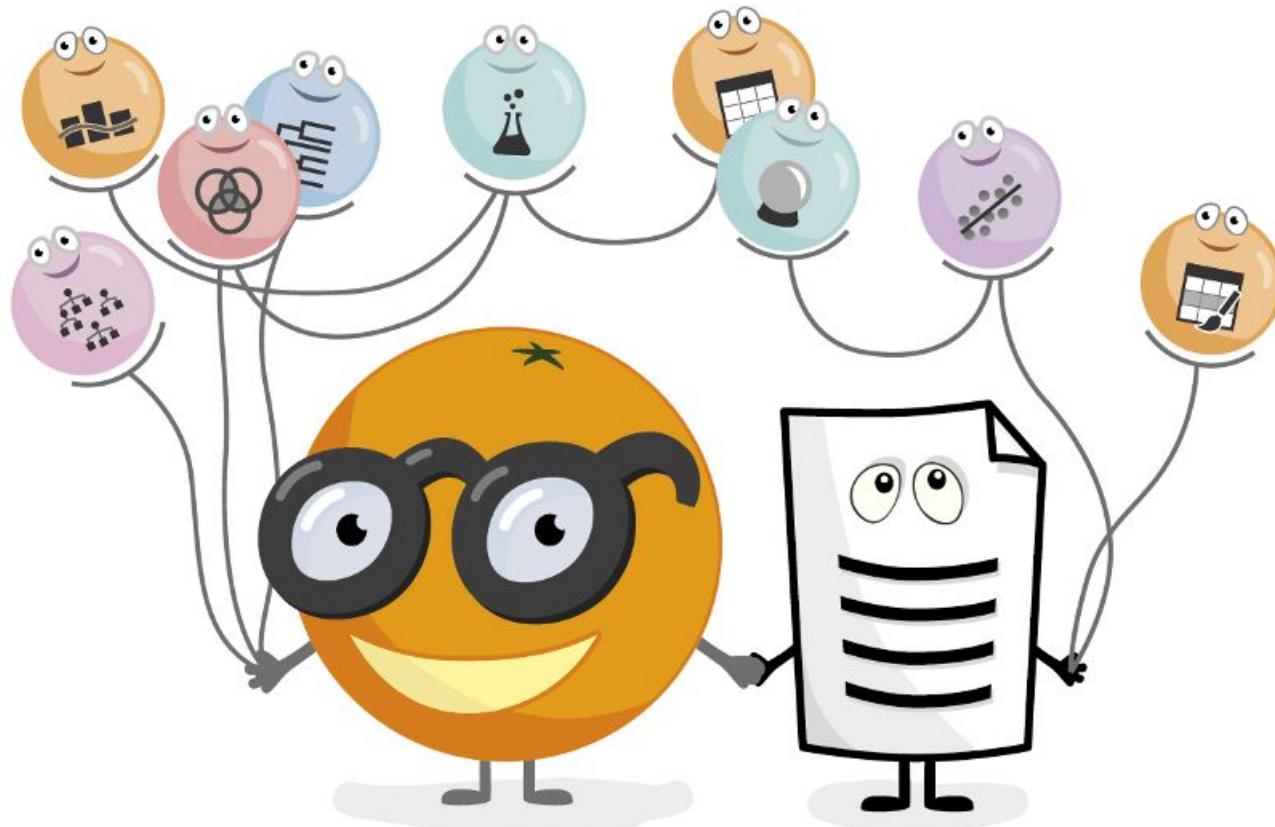
C A R O L



<https://goo.gl/Ndf38Q>



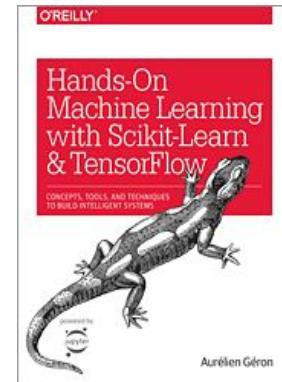
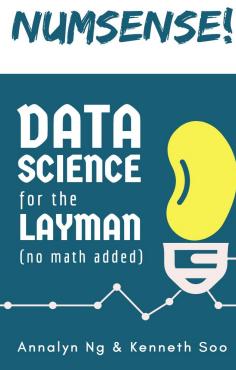
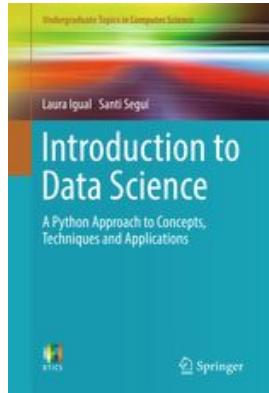
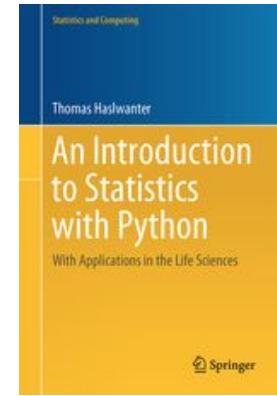
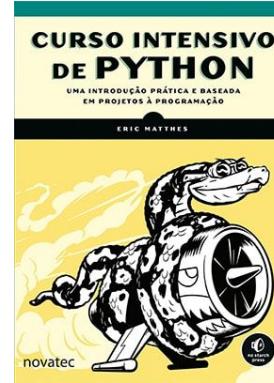
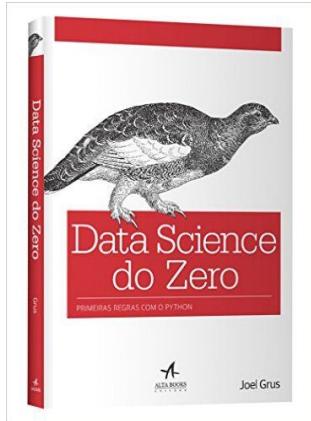
Data Mining Fruitful and Fun



Agenda

- Python Introduction
- Data Analysis and Visualization
- Working with data sources
- Probability & Statistics
- Machine Learning
- Network Analysis
- Working with large datasets
- Advanced topics in data science

References



References



Data Analyst

facebook mongoDB

This block contains a thumbnail image for the Data Analyst course, which shows a purple background with abstract data points and connections. Below the thumbnail, the course name 'Data Analyst' is displayed in a white box, along with social media icons for Facebook and MongoDB.



Machine Learning Engineer



kaggle



References

The screenshot shows the DataCamp website homepage. At the top, there is a navigation bar with links for Home, Courses, Tracks (beta), Pricing, Business, Community, Sign in, and a prominent 'Create Free Account' button. Below the navigation bar, a large banner features the text 'THE EASIEST WAY TO Learn Data Science Online'. It includes two main calls-to-action: 'Start Learning R' and 'Start Learning Python'. To the right of the banner is a 'Create Your Free Account' form. This form has fields for email (with placeholder 'ivan@imd.ufrn.br'), password (with placeholder '.....'), and social media links for LinkedIn ('in'), Facebook ('f'), and Google+ ('G+'). A large 'Get Started' button is at the bottom of the form. The background of the page features a grid of course thumbnails.

References

The screenshot shows the Dataquest website's homepage. At the top, there is a dark header bar with the Dataquest logo, a 'Dashboard' link, a 'GET HELP' button, a notifications icon (0), and a user profile for 'Ivanovitch'. Below the header, the main content area has a light gray background. It features a large title 'Become a Data Scientist' in bold, dark gray font. Underneath the title is a subtext: 'Our hands-on method teaches you all the skills you need to become a data scientist or data analyst.' Below that is another subtext: 'Learn by writing code, working with data, and building projects in your browser.' To the right of the text, there is a stylized illustration of a rocket ship launching from clouds, with a teal trajectory line and small green stars.

Become a Data Scientist

Our hands-on method teaches you all the skills you need to become a data scientist or data analyst.

Learn by writing code, working with data, and building projects in your browser.

References



The image shows the homepage of the Data Science Academy website. At the top left is the logo, which consists of a network of colored dots (blue, green, yellow) connected by lines. To the right of the logo is the text "Data Science Academy". Along the top edge are several navigation links: "Início", "Contato", "Blog", "Pagamento", "Todos os Cursos" (which is highlighted in blue), "Sobre Nós", "Metodologia", "Trabalhe Conosco", "FAQ", "Clientes", and "App's". Below the header is a large banner image of a woman sitting at a desk, looking down at her laptop screen. Overlaid on this image is the text "Data Science Academy" in large white letters. In the bottom right corner of the banner is a small Brazilian flag icon. At the bottom center of the page is a blue button with the text "Inscreva-se Agora". In the bottom right corner of the entire page are two small circular arrows, one pointing left and one pointing right.

Data Science Academy

Tecnologia e formação profissional para ampliar sua empregabilidade de forma ilimitada e online!

Inscreva-se Agora

References



MASSACHUSETTS INSTITUTE OF TECHNOLOGY

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.

The MIT Experimental Learning logo is visible in the top right corner of the slide.

References

Stanford | ONLINE



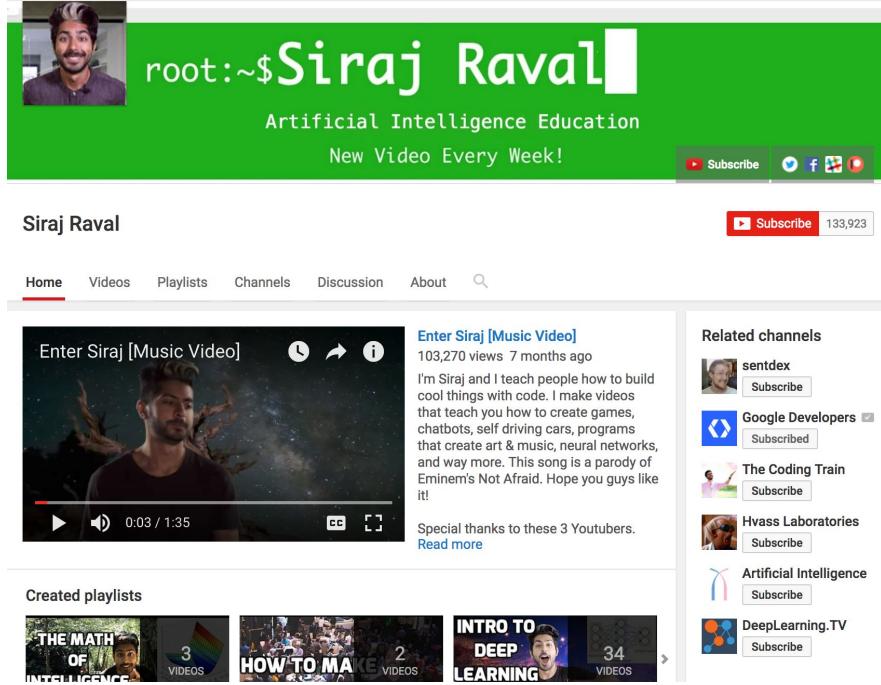
BEHIND AND BEYOND BIG DATA

STARTS ONLINE: 06/12/17
AT STANFORD: 07/25/17 - 07/28/17

[APPLY NOW](#)



References



Siraj Raval

[Home](#) [Videos](#) [Playlists](#) [Channels](#) [Discussion](#) [About](#) [Search](#)

Enter Siraj [Music Video]
103,270 views 7 months ago

I'm Siraj and I teach people how to build cool things with code. I make videos that teach you how to create games, chatbots, self driving cars, programs that create art & music, neural networks, and way more. This song is a parody of Eminem's Not Afraid. Hope you guys like it!

Special thanks to these 3 You tubers.
[Read more](#)

Created playlists

- THE MATH OF INTELLIGENCE** 3 VIDEOS
- HOW TO MAKE** 2 VIDEOS
- INTRO TO DEEP LEARNING** 34 VIDEOS

Related channels

- sentdex** [Subscribe](#)
- Google Developers** [Subscribed](#)
- The Coding Train** [Subscribe](#)
- Hvass Laboratories** [Subscribe](#)
- Artificial Intelligence** [Subscribe](#)
- DeepLearning.TV** [Subscribe](#)



<http://www.andrewng.org/>



<http://mariofilho.com/>



References

<https://elitedatascience.com/>

<https://algobeans.com/>

<https://www.datacamp.com/home>

<https://www.dataquest.io/dashboard>

<https://www.datascienceacademy.com.br/>

<http://www.bigdatabusiness.com.br/>

<https://datafloq.com/>

<https://github.com/data-8/>

<http://www.informationisbeautiful.net/>

<http://www.datapedia.info/public/>

<http://ckan.imd.ufrn.br/>

<http://dados.ufrn.br/>

<https://fivethirtyeight.com/>

<https://news.ycombinator.com/>

<http://machinelearningmastery.com/blog/>

References

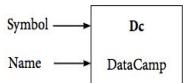
- How do I learn python in deep?
 - <https://www.hackerrank.com/>
 - <https://www.codewars.com/>
 - <https://br.codecombat.com/>
 - <https://www.hackerearth.com>
 - <https://www.drivendata.org/>
 - <https://www.kaggle.com/>
 - <https://goo.gl/WhLvs9>



The Periodic Table of Data Science

An overview of key companies, resources and tools in data science (as of 4/12/2017)

<https://goo.gl/eecQ2V>



Dc	Ga	Sd
DataCamp	General Assembly	Strata Data

Courses	Data
Boot camps	Projects & Challenges, Competitions
Conferences	Programming Languages & Distributions

Search & Data Management	Collaboration
Machine Learning & Stats	Community & Q&A
Data Visualization & Reporting	

News, Newsletters & Blogs
Podcasts

Kdn	Ibd
KDnuggets	insideBIGDATA
Rb	Pp
R-Bloggers	PlanetPython
Hn	Dt
HackerNews	DataTau
Dsc	Dsr
Data Science Central	Data Science Roundup
Dsw	Or
Data Science Weekly	O'Reilly
Dr	Pw
Data Elixir	Python Weekly
Rw	Pd
R Weekly	Partially Derivative
Bds	Tm
Becoming a Data Scientist	Talking Machines
Ds	Dsk
Data Stories	Data Skeptic
Ld	Ns
Linear Digressions	Not So Standard Deviations

Ex	Di	Tc
Edx	Data Incubator	Tableau Conference

Py	Js	Vb	Pgs	Sli	Ah	W	Bml	Kn	Sm	Pb	Obi	Shn	Ddl	De
Python	JavaScript	Visual Basic	PostgreSQL	SQLite	Apache Hadoop	Weka	BigML	KNIME	Spark MLlib	Power BI	Oracle BI	Shiny	Domino Data Lab	Data Science Experience
R	Cp	Sc	Ar	Bq	Hw	O	Dar	Lib	Ho	Bo	Alt	Mpl	Nt	Rs
R	C++	Scala	Amazon Redshift	Google BigQuery	Hortonworks	Oracle	DataRobot	LIBSVM	H2O	BusinessObjects	Alteryx	Matplotlib	Nteract	Rstudio
S	Pl	Ca	Hb	Td	Cl	Mss	Microsoft SQL server	Rm	Mat	Th	Sp	Sav	Ply	Ro
SQL	Perl	Cassandra	HBase	Teradata	Cloudera	Microsoft SQL server	RapidMiner	Mathematica	Theano	Spotfire	SAS Visual Analytics	Plotly	Rodeo	Beaker Notebook
B	Mr	P	Mdb	To	Aem	Spl	Cho	Mah	Aml	Ql	Po	Me	Spy	Ze
Bash	Microsoft R Open	Pig	Mongo DB	Toad	Amazon Elastic Mapreduce	Splunk	Chorus	Mahout	Azure Machine Learning	Qlikview	PowerPivot	Microsoft Excel	Spyder	Apache Zeppelin
Mtl	Cy	Im	K	Ms	Mar	Sr	Tf	St	D	Co	Gch	Pe	Dst	Ju
Matlab	Canopy	Impala	Kafka	MySQL	MapR	Solr	Tensorflow	Stata	D3	Cognos	Google Charts	Pentaho	Data Science Studio	Jupyter
J	An	Sp	Hi	Idb	Lu	El	Sk	Da	My	Aa	T	B	Db	Gh
Java	Anaconda	Spark	Hive	IBM DB2	Lucene	ElasticSearch	Scikit-Learn	Dato/GraphLab	MicroStrategy	Adobe Analytics	Tableau	Bokeh	Databricks notebook	Github

Tt	Dsj	Icd
TeamTreeHouse	Data Science Dojo	IEEE International Conference on Data Mining

Data.world	Quandl	FiveThirtyEight	Socrata	Google Public	Data.gov	Kaggle
Statista	Uci UCI Machine Learning Repository	Wb World Bank	At Academic Torrents	Bf Buzzfeed	Dk DataKind	Dd DrivenData

Reddit	Stack Overflow	Cross Validated	Quora	Analytics Vidhya	Data Science Stack Exchange
Mu	Rdm				
Meetup	RDataMining				



Lesson #1