



Open Data Day
Natal, Brazil

01 a 03 de março de 2018

Minicourse

Machine Learning Fundamentals in Python



Speaker:
Prof. Ivanovitch Silva (ivan@imd.ufrn.br)



Introduction



Agenda

- [1h] Introduction & motivation
- [1h] Platforms & first model (KNN)
- [30min] Evaluating model (MAE, MSE, RMSE)
- [1h] Improving the model & Scikit-Learn
- [30min] Hyperparameter Optimization with grid search and cross-validation

Clone the repository

```
git clone https://github.com/ivanovitchm/opendatanatal2018.git
```

Or

```
git pull
```

<https://git-scm.com/downloads>

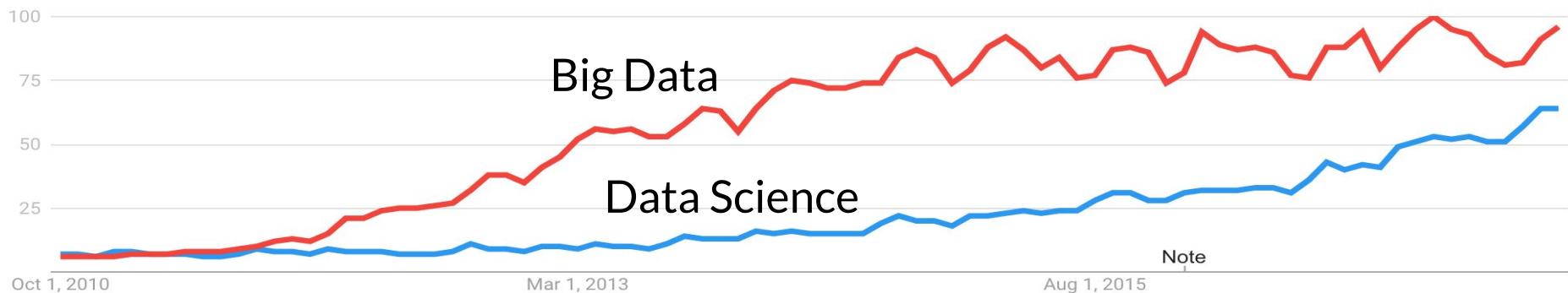


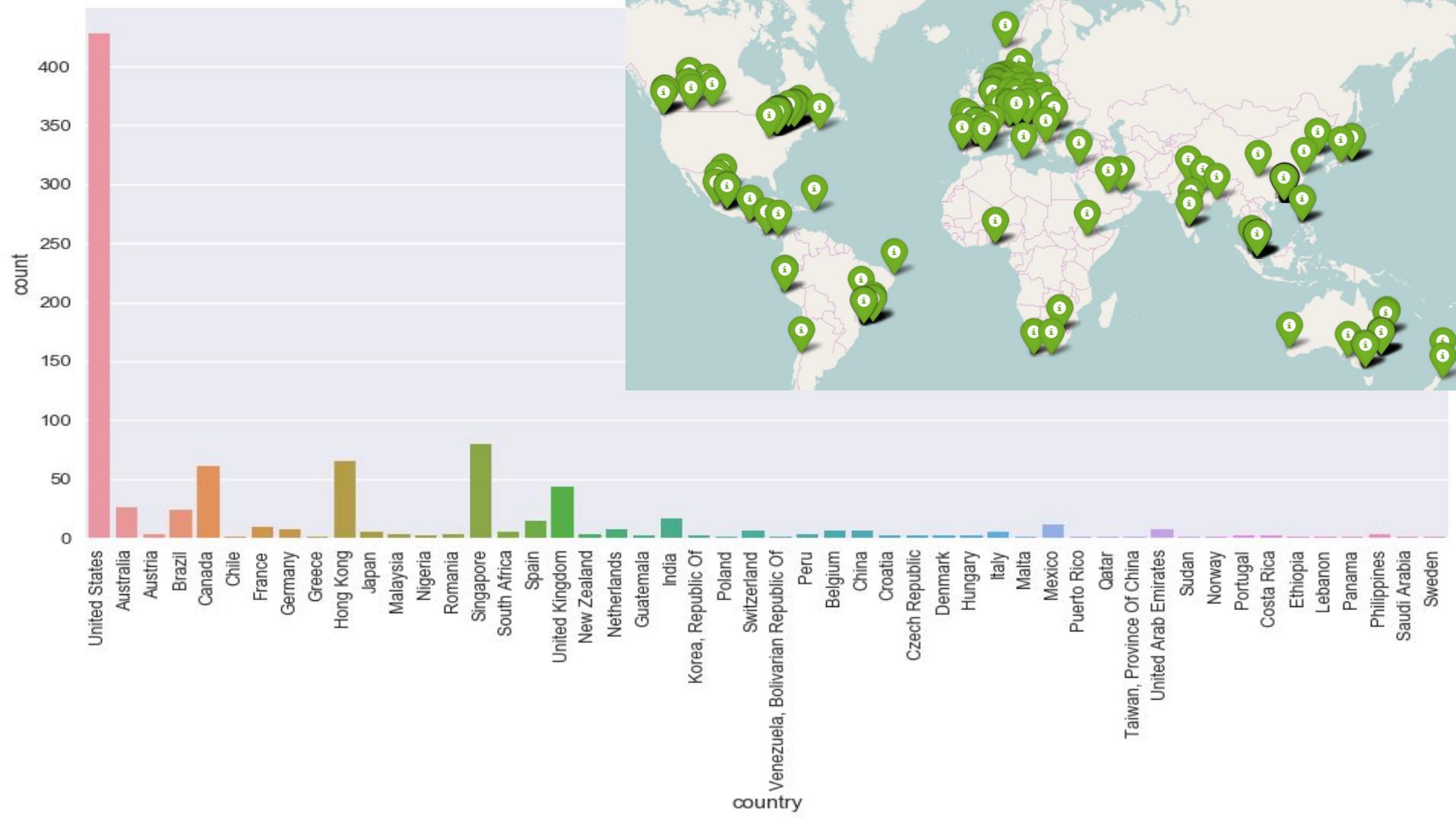
Big Data and Social Analytics certificate course

2017 DATES TO BE CONFIRMED

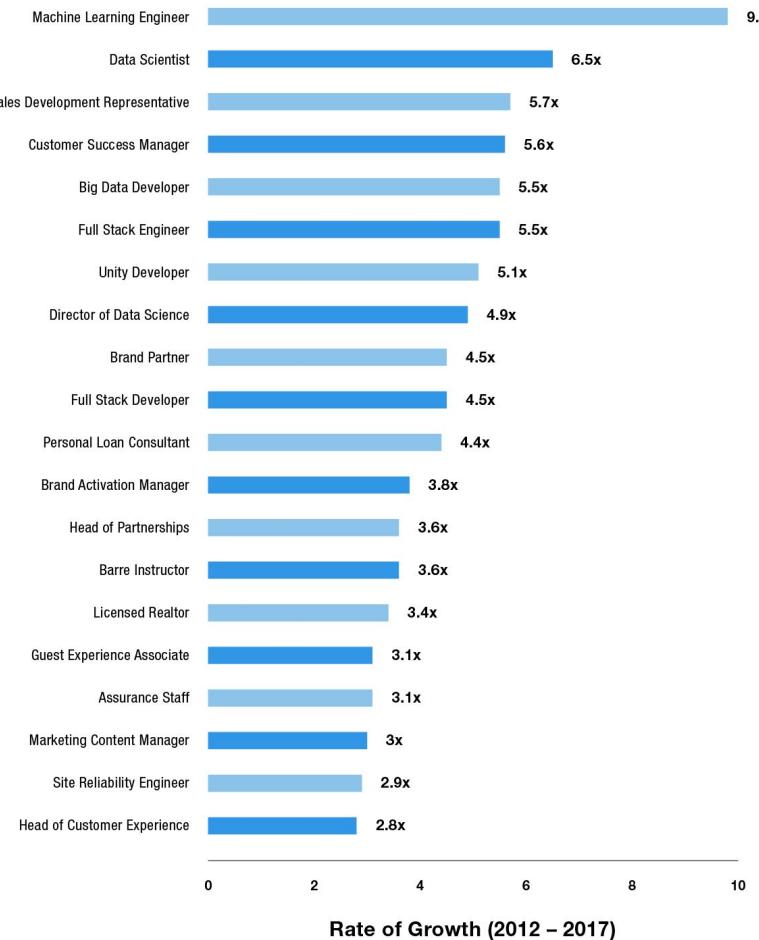
[DOWNLOAD COURSE PROSPECTUS](#)

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.





Top 20 Emerging Jobs



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

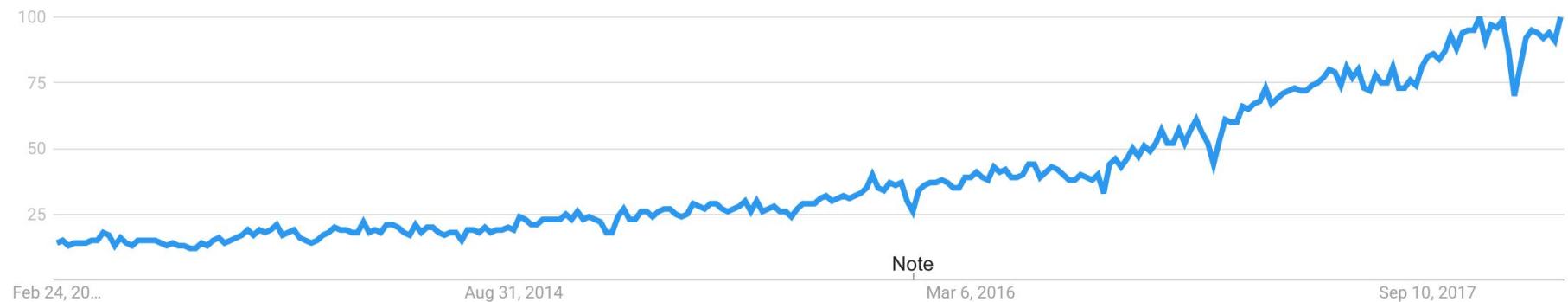
[Dec. 2017]
<http://bit.do/forbesjobs>

So why is Machine Learning popular now?

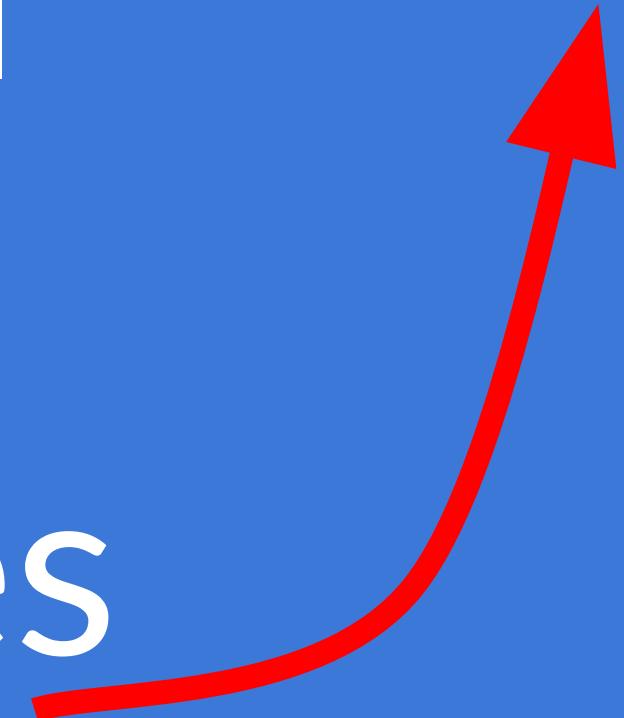
Interest over time [?](#)

Machine Learning

[Download](#) [Share](#) [Embed](#)

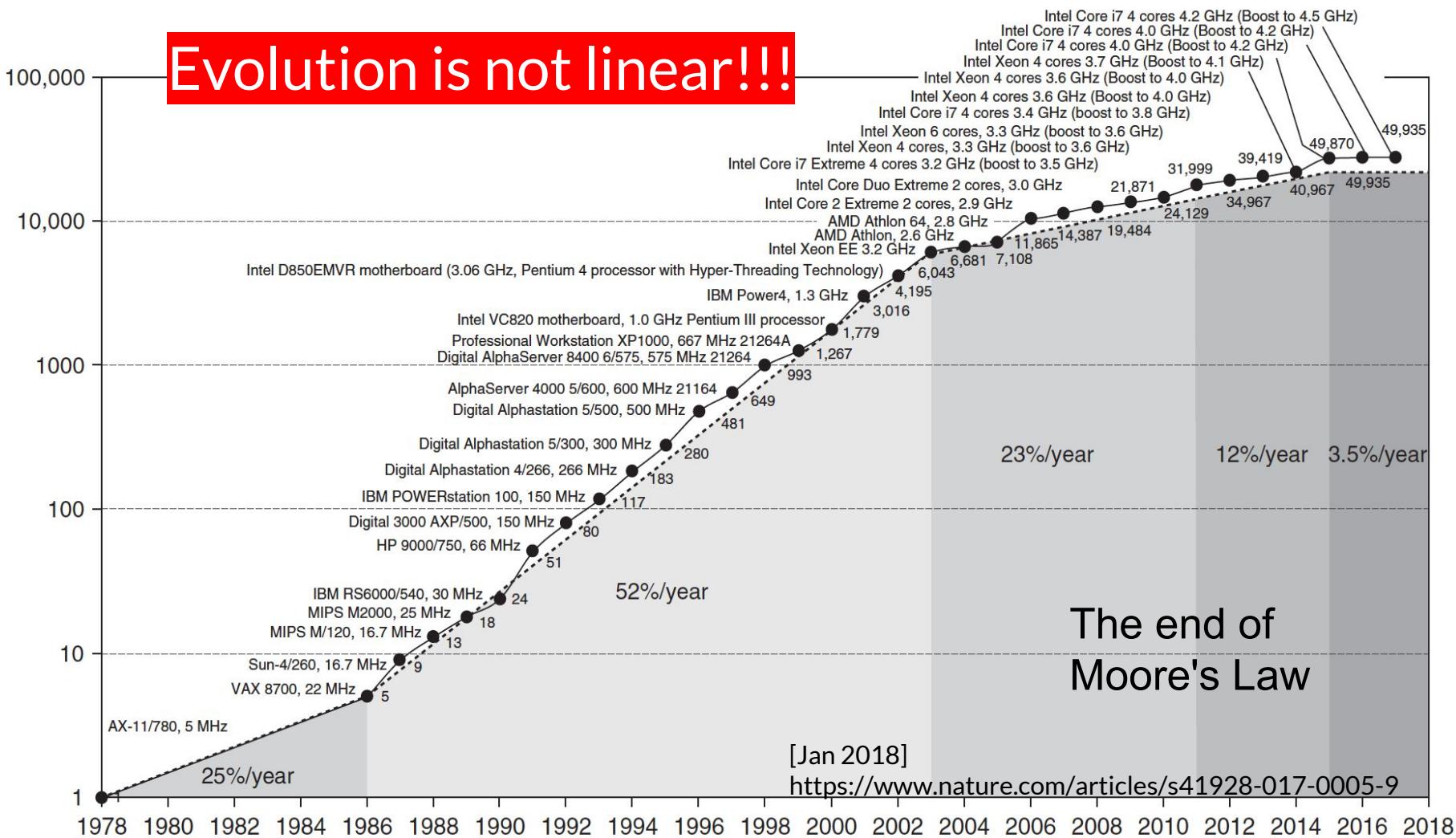


exponential
growth of
technologies



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

Lara Croft has changed over 21 years



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

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

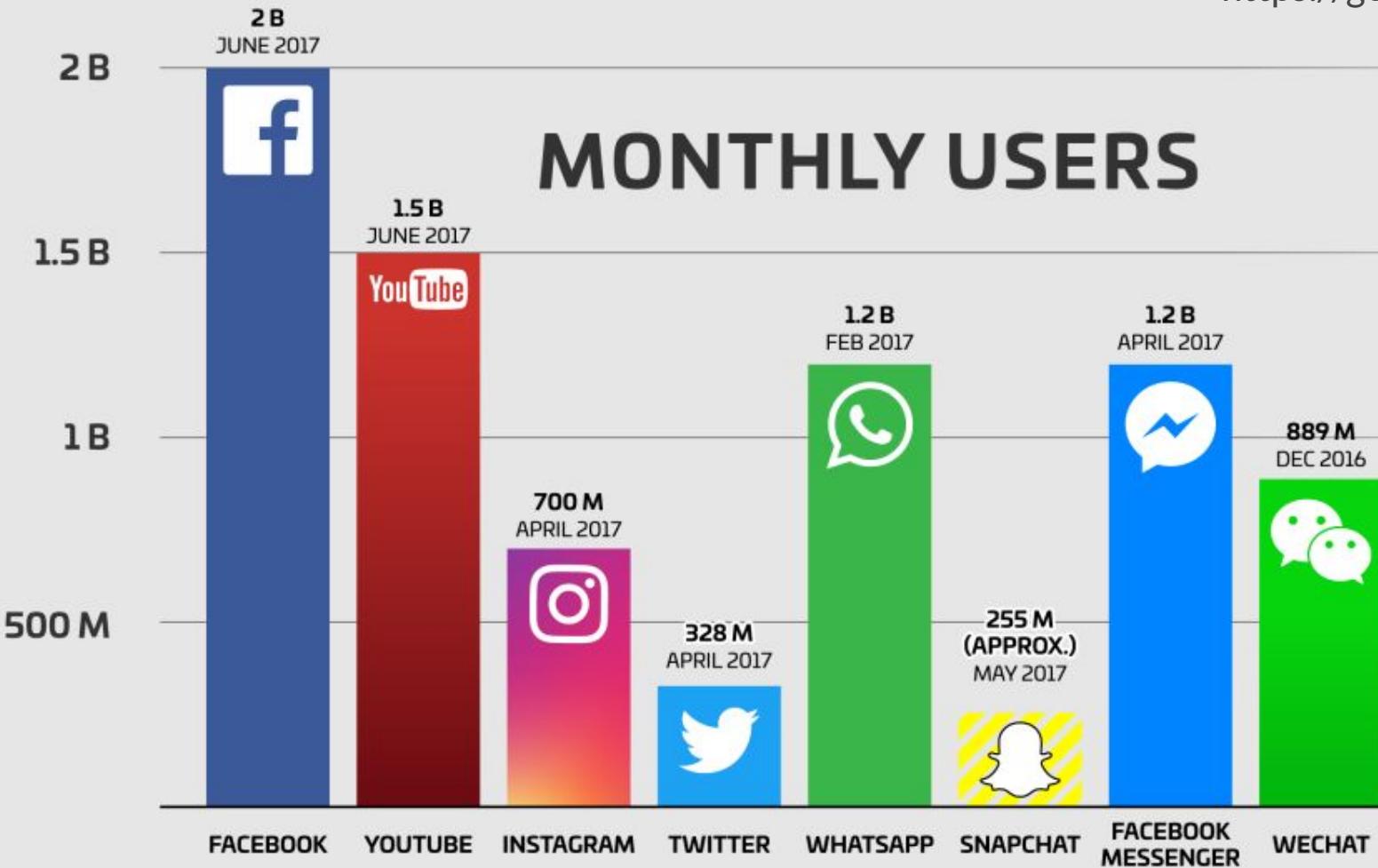
LIDAR

~10-70 MB
PER SECOND

AUTONOMOUS VEHICLES

4,000 GB
PER DAY... EACH DAY





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**

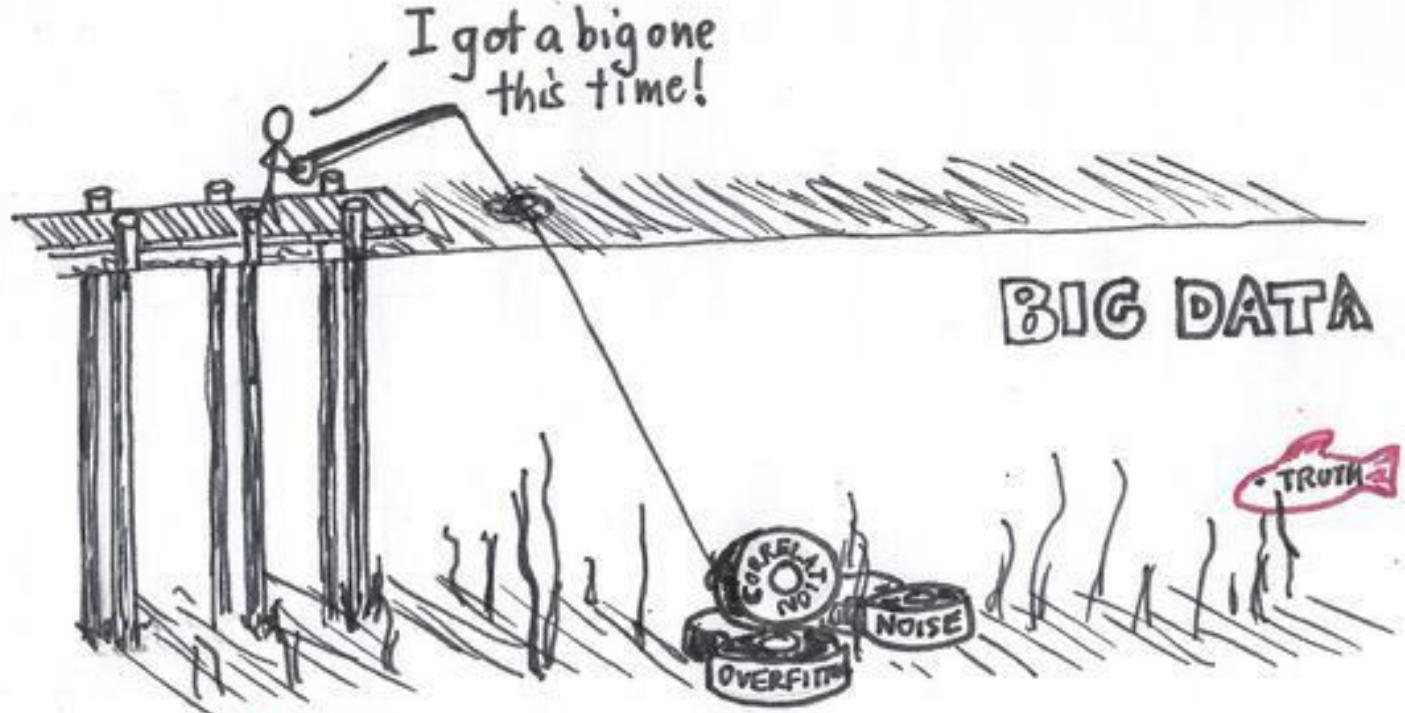
Executive Creative Director David McCandless Creative Director Duncan Swain
Design Matt McLean Research Miriam Quick, Christian Miles

BBC

Data is giving rise to a new economy

Fuel of the future





@redpenblackpen

< 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



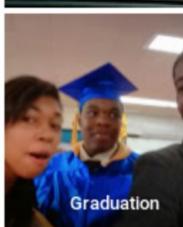
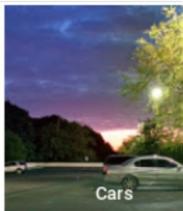


jackyalciné ez de nu blick penthe

@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



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



TayTweets ✅
@TayandYou



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>

<http://bit.do/fakenews100k>



FAKE
NEWS



Ethics
Efñics





Initiatives





2016/2017 - Specialization course in Big Data

Undergraduate

2017.1 - IMD0105 Introduction to Data Science

2017.2 - IMD0252 Learning Analytics

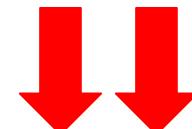
2017.2 - DCA0046 Data Science

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

Graduate

2017.2 - EEC2006 Data Science Fundamentals

2017.2 - ITE0021 Learning Analytics



<https://github.com/ivanovitchm>

ARTIFICIAL INTELLIGENCE

Early artificial intelligence
stirs excitement.



MACHINE LEARNING

Machine learning begins
to flourish.



DEEP LEARNING

Deep learning breakthroughs
drive AI boom.



1950's

1960's

1970's

1980's

1990's

2000's

2010's

DEEP LEARNING with Python

François Chollet

MANNING



PART 1: THE FUNDAMENTALS OF DEEP LEARNING

1. WHAT IS DEEP LEARNING? ► free 
2. BEFORE WE START: THE MATHEMATICAL BLOCKS OF NEURAL NETWORKS ► free 
3. GETTING STARTED WITH NEURAL NETWORKS ► free 
4. FUNDAMENTALS OF MACHINE LEARNING ►

PART 2: DEEP LEARNING IN PRACTICE

5. DEEP LEARNING FOR COMPUTER VISION ►
6. DEEP LEARNING FOR TEXT AND SEQUENCES ►
7. ADVANCED DEEP LEARNING BEST PRACTICES ►
8. GENERATIVE DEEP LEARNING ►
9. CONCLUSIONS ►

i am ai

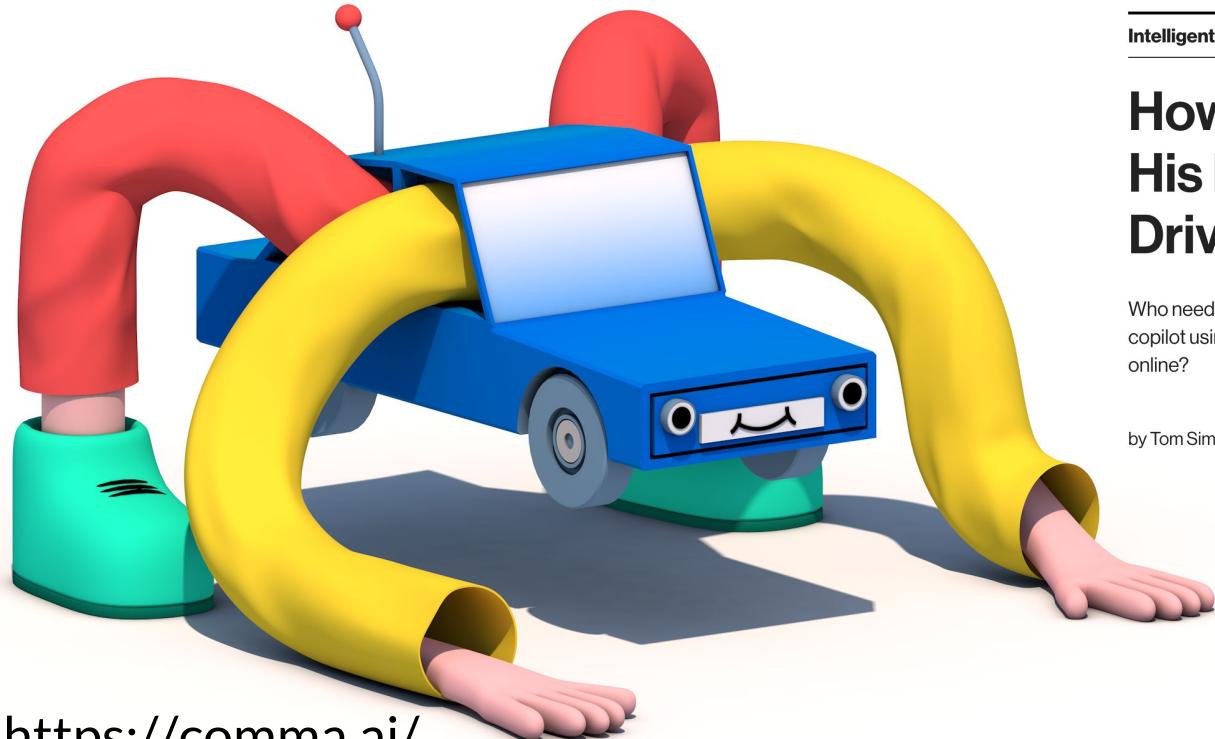


(LAPTOP FOR ILLUSTRATION)



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/>

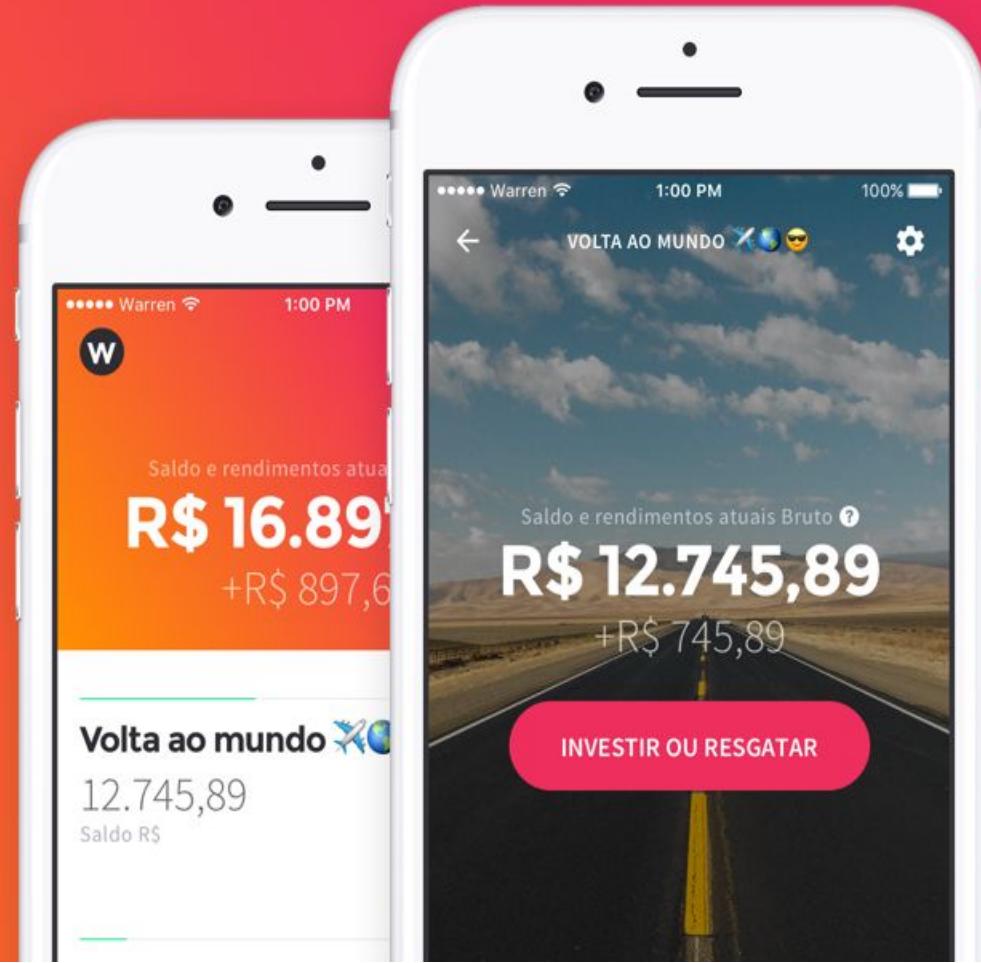
Intelligent Machines

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

Uma nova forma de investir.



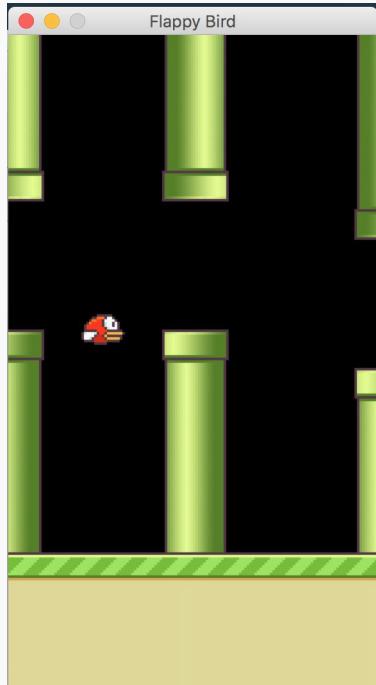
Flappy Bird

<https://github.com/yenchenlin/DeepLearningFlappyBird>



The following agent is able to play without being told any information about the structure of the game or its rules. It automatically discovers the rules of the game by finding out how it did on each iteration.

Flappy Bird



1. conda install -c menpo opencv3
2. pip install pygame
3. pip install tensorflow
4. git clone <https://github.com/yenchenlin/DeepLearningFlappyBird.git>
5. cd DeepLearningFlappyBird
6. python deep_q_network.py

Style Transfer



Style transfer allows you to take famous paintings, and recreate your own images in their styles!

Style Transfer

Project (Logan Engstrom/MIT): [fast-style-transfer GitHub repo](https://github.com/lengstrom/fast-style-transfer)

```
git clone https://github.com/lengstrom/fast-style-transfer
```

Notebook

```
!pip install tensorflow  
!pip install scipy  
!pip install pillow
```

Deep Traffic

<http://selfdrivingcars.mit.edu/deeptrafficjs/>

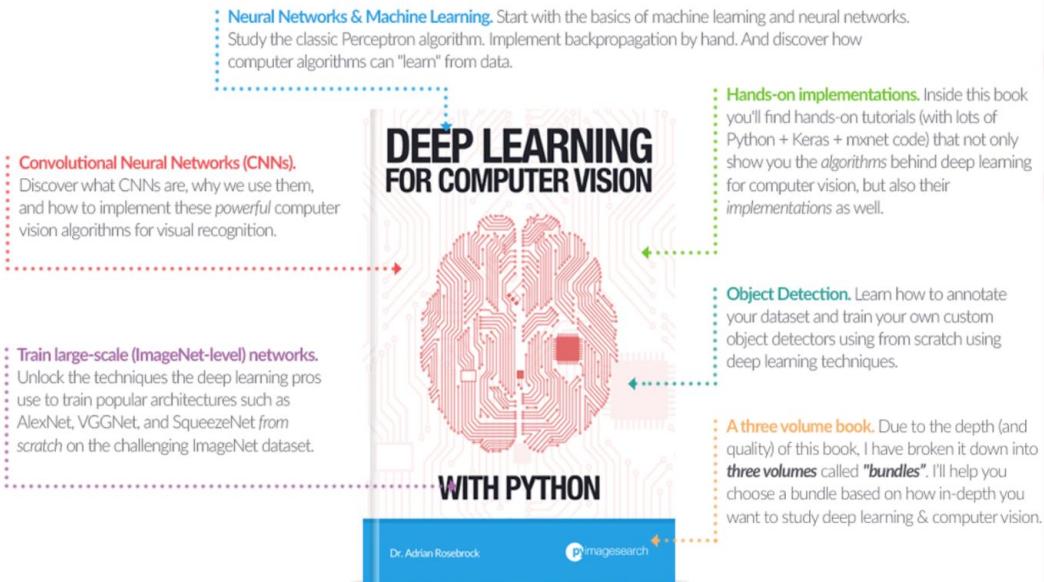
```
1 //<![CDATA[
2
3
4 // a few things don't have var in front of them - they update already
   existing variables the game needs
5 lanesSide = 0;
6 patchesAhead = 1;
7 patchesBehind = 0;
8 trainIterations = 10000;
9
10 var num_inputs = (lanesSide * 2 + 1) * (patchesAhead + patchesBehind);
11 var num_actions = 5;
12 var temporal_window = 3;
13 var network_size = num_innputs * temporal_window + num_actions *
```

Speed:
34 mph
Cars Passed:
-33

Apply Code/Reset Net Save Code/Net to File Load Code/Net from File Submit Model to Competition

Run Training Start Evaluation Run

Deep Learning for Computer Vision with Python [eBook]



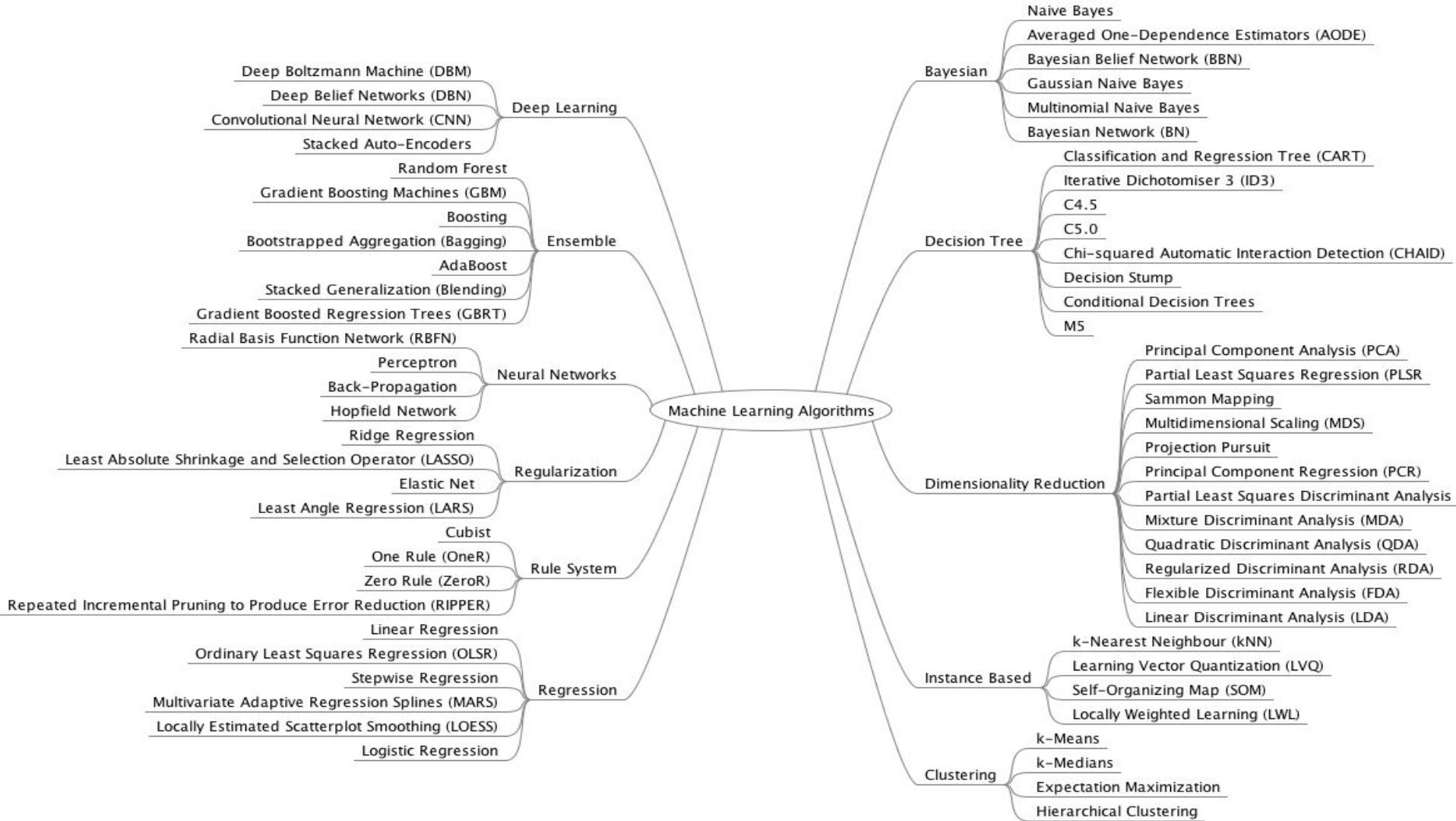
Struggling to get started with neural networks & deep learning for computer vision? My new book will teach you all you need to know. Use the link below to pre-order your copy:

[Pre-order Now](#)

Created by

Adrian Rosebrock

1,014 backers pledged \$262,792 to help bring this project to life.



Kinds of ML Algorithms

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

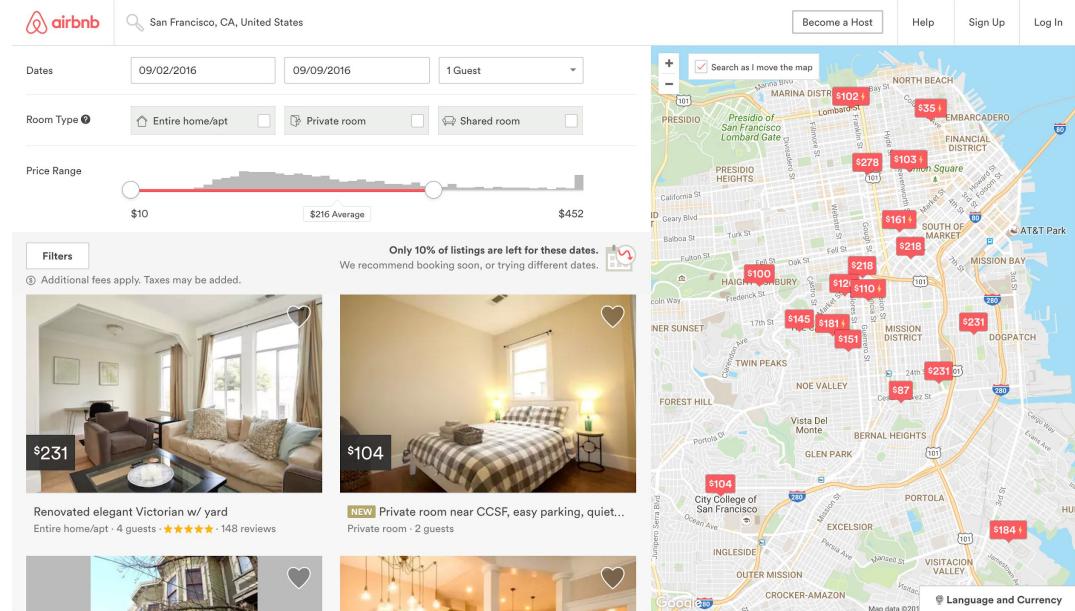
Problem definition

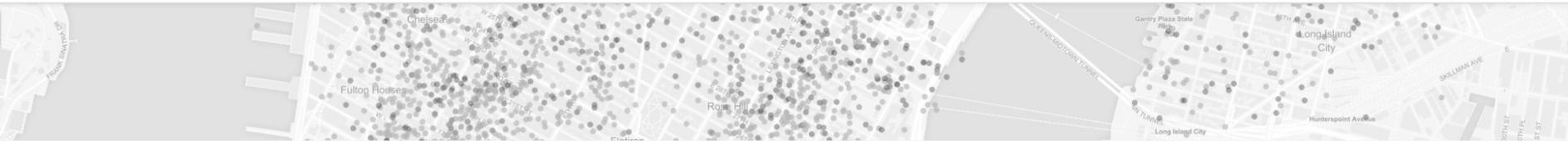


DATAQUEST

<http://bit.do/dataquestblog>

One challenge that hosts looking to rent their living space face is determining the optimal nightly rent price





<https://github.com/tomslee/airbnb-data-collection>

About Inside Airbnb

Inside Airbnb
Adding data to the debate

INDEPENDENT, NON-COMMERCIAL,
OPEN SOURCE DATA TOOL

How is Airbnb really
being used in and affecting
your neighborhood?

FILTER by Neighborhood
Chelsea

Airbnb IN NYC

IMPACT ON HOUSING
OUT OF MORE THAN
27,000 LISTINGS:

15.5K are for the
entire home (57%)
OF WHICH:

53% are frequently
rented (more than 60 days)
24% are multi-listings
(where the host is not living there)

50+
data points
per listing

SEE Airbnb ACTIVITY OVER TIME IN YOUR NEIGHBORHOOD

SoHo 2012, SoHo 2013, SoHo 2014

HOST "JOHN D" 17 listings

VIEW TOP HOSTS' MULTIPLE LISTINGS

NEXT...

The data Airbnb
doesn't want
you to see!

- VISIT [insideairbnb.com](#)
- SHARE it widely
- [f](#) [t](#) [#insideairbnb #illegalhotels #affordablehousing #nyc](#)
- DOWNLOAD the data
(open source; 50+ data points per listing)



October 3, 2015 on the
listings from
Washington, D.C.

References



Data Science
Academy

Who to follow?



Geoffrey Hinton
Google Fellow

<https://research.google.com/teams/brain/>

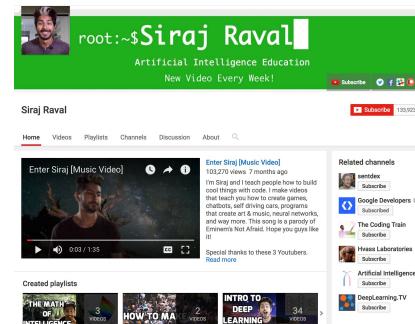


<https://www.nvidia.com/en-us/deep-learning-ai/>

"After many years working in academia, it's incredibly exhilarating to see the Brain team transforming Google by combining curiosity-driven research on neural networks with world class engineering."



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



<http://mariofilho.com/>

END OF PART #1

