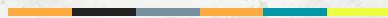


IMD0905 - Data Science I

Lesson #1 - Outline & Directions

Ivanovitch Silva
February, 2019





Introduction





Ivanovitch Silva (ivan@imd.ufrn.br)
3T1234

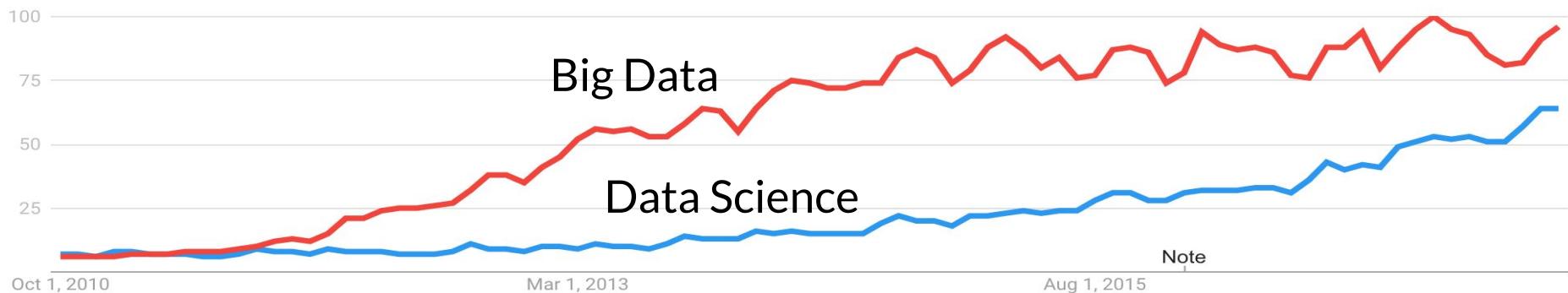


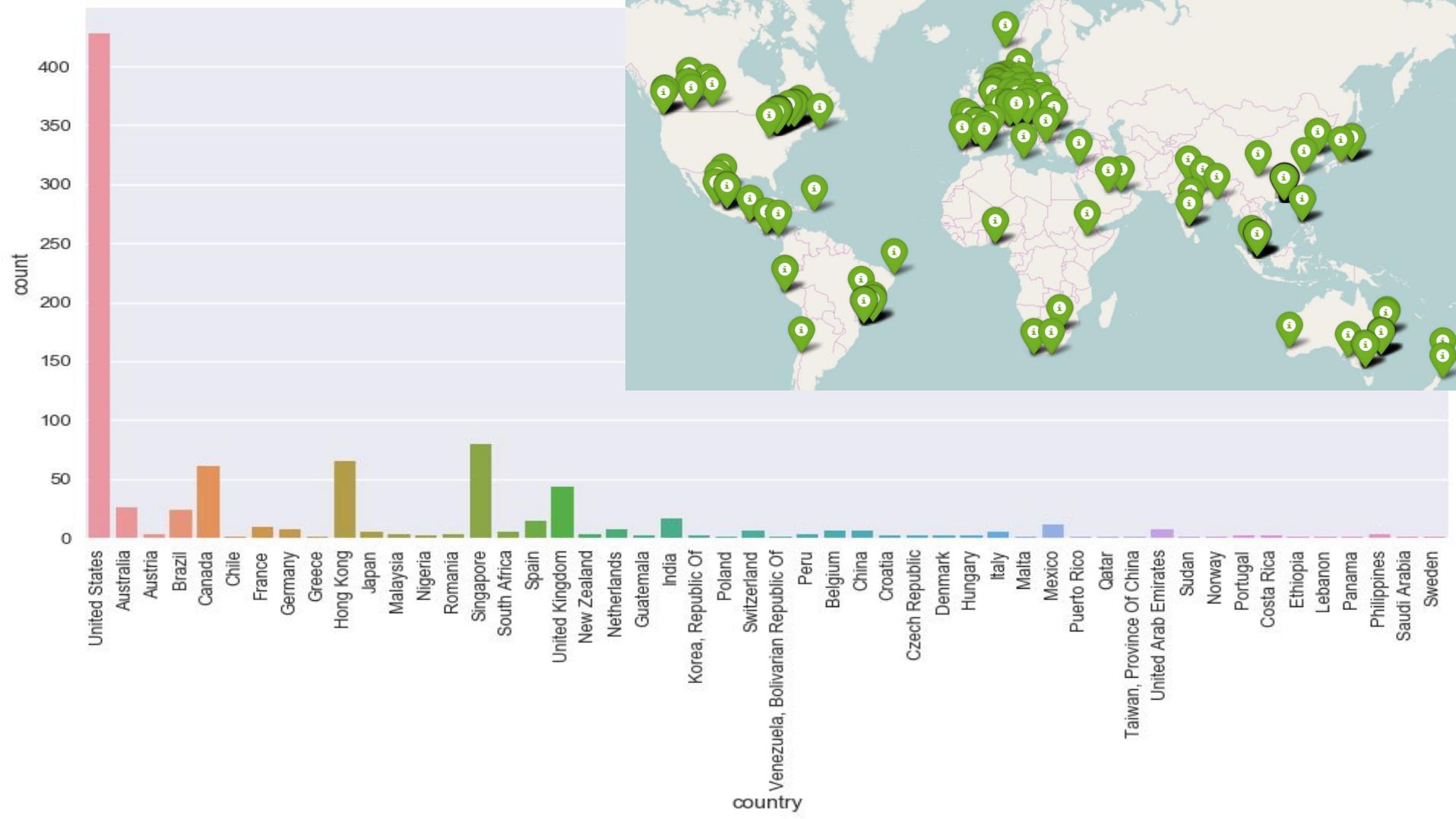
Big Data and Social Analytics certificate course

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.





**Insanity is doing the same
thing, over and over again, but
expecting different results.**

Albert Einstein



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

2018.2 - IMD0905 Data Science I

Graduate

2017.2 - EEC2006 Data Science Foundations

2017.2 - ITE0021 Learning Analytics

2018.2 - EEC1509 Machine Learning

2019 - PES



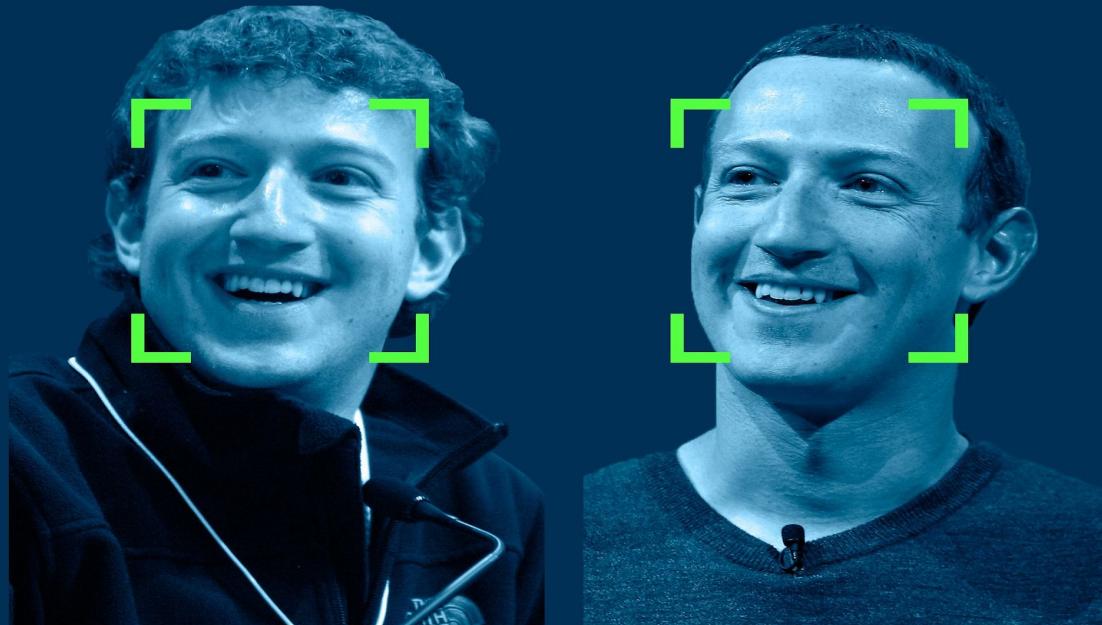
COURSE Expectations

infringing



The background image shows a close-up of a wall that has been painted yellow. The paint is heavily peeling and cracking, especially towards the bottom right, revealing a dark, textured surface underneath. A rough brick wall is visible at the bottom left.

#10yearchallenge



#10YearsChallenge



2009



2019





\$0.00
2009



\$3,363.36
2019



2009



2019



*BIG
changes*

125 kbps

R\$ 59,90
2009



30,15



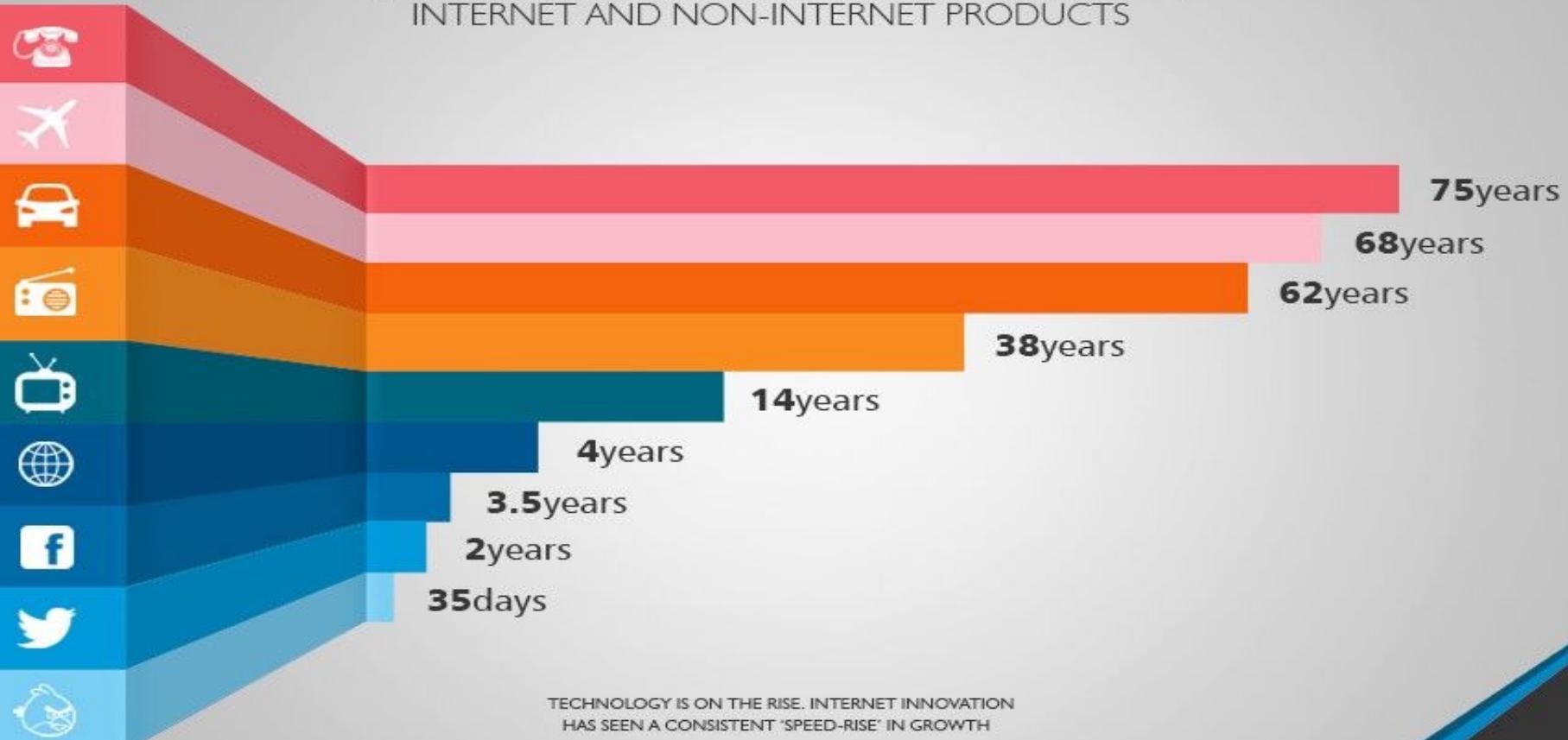
30Mbps
R\$ 89,90
2019

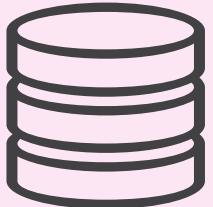
*What do they
have in common?*

exponential
growth of
technologies

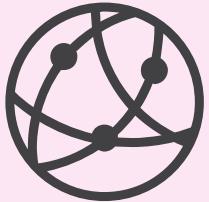


REACHING 50 MILLION USERS: THE JOURNEY OF INTERNET AND NON-INTERNET PRODUCTS

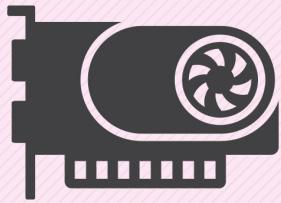




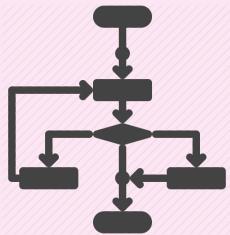
DATA



INTERNET



HARDWARE



ALGORITHM



data comes from everywhere



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



Database
(MySQL, Oracle, ...)



API



Sites



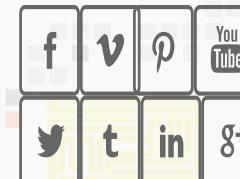
Text and reports



Maps



Image and videos



Social Media



Facebook

Monthly Active Users: **2.2 Billion** Daily Active Users: **1.4 Billion** Founded: **2004**

Photos uploaded daily: **300 Million** Video views daily: **8 Billion** Rank: **#1**



WhatsApp

Monthly Active Users: **700 Million** Daily Active Users: **320 Million** Founded: **2009**

New users daily: **1 Million** Messages sent daily: **43 Billion** Rank: **#4**



YouTube

Monthly Active Users: **1.5 Billion** Daily Active Users: **30 Million** Founded: **2005**

Video views daily: **5 Billion** Average visit length: **40 min.** Rank: **#2**



Google+

Monthly Active Users: **395 Million** Total Registered Users: **2 Billion** Founded: **2011**

U.S. based users: **55%** Ages 15-34 users: **28%** Rank: **#5**



Instagram

Monthly Active Users: **800 Million** Daily Active Users: **500 Million** Founded: **2010**

Photos uploaded daily: **95 Million** Stories daily: **250 Million** Rank: **#3**

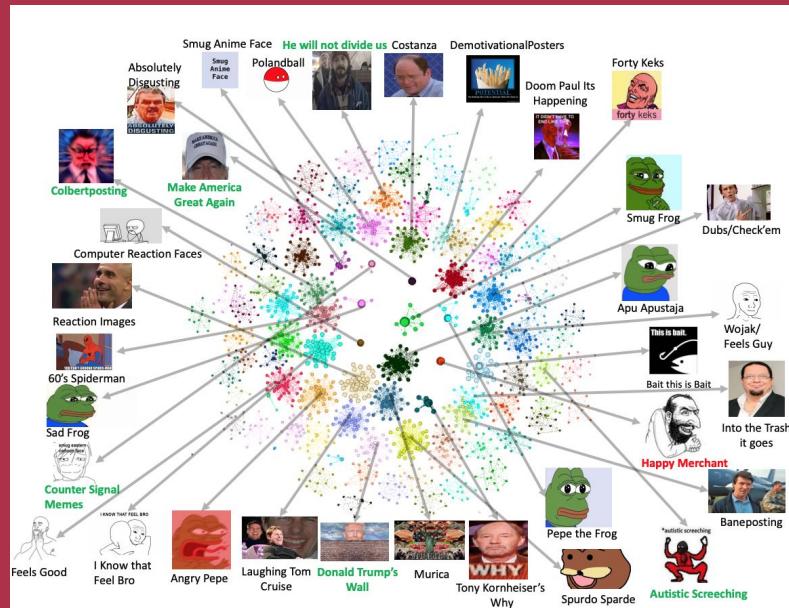
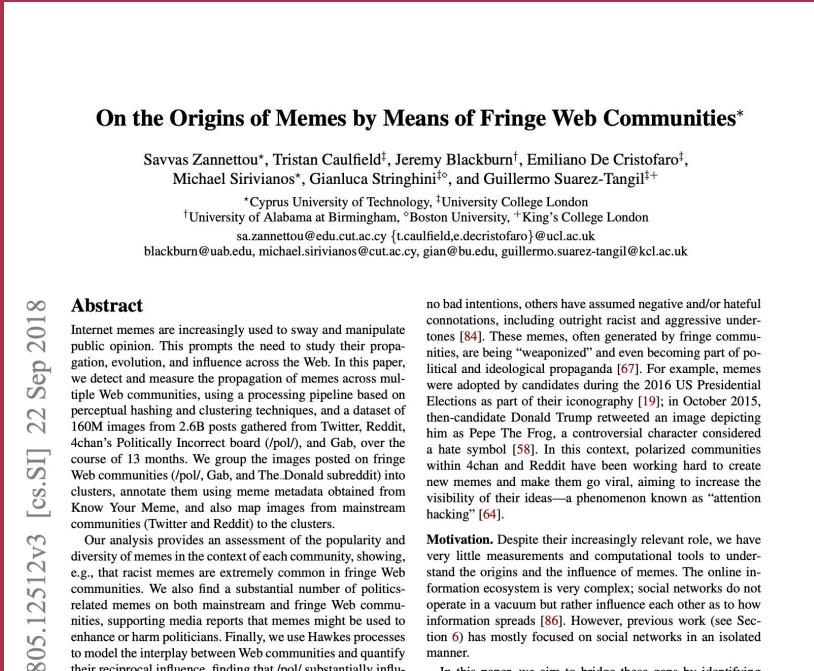


Twitter

Monthly Active Users: **330 Million** Daily Active Users: **100 Million** Founded: **2006**

Tweets published daily: **140 Million** New accounts daily: **460,000** Rank: **#6**

This is where internet
memes come from



Abstract

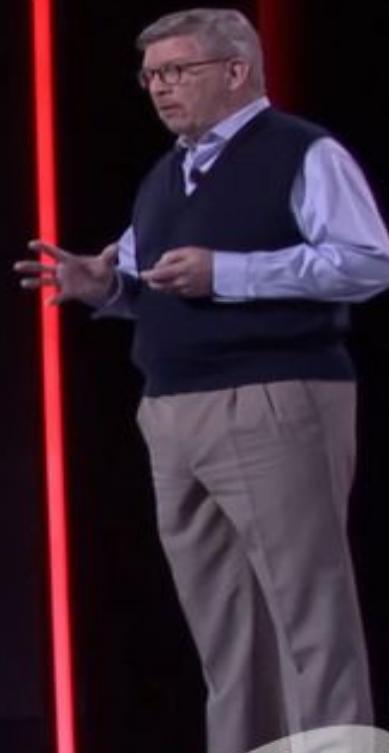
Internet memes are increasingly used to sway and manipulate public opinion. This prompts the need to study their propagation, evolution, and influence across the Web. In this paper, we detect and measure the propagation of memes across multiple Web communities, using a processing pipeline based on perceptual hashing and clustering techniques, and a dataset of 160M images from 2.6B posts gathered from Twitter, Reddit, 4chan's Politically Incorrect board (*/pol/*), and Gab, over the course of 13 months. We group the images posted on fringe Web communities (*/pol/*, Gab, and The Donald subreddit) into clusters, annotate them using meme metadata obtained from Know Your Meme, and also map images from mainstream communities (Twitter and Reddit) to the clusters.

Our analysis provides an assessment of the popularity and diversity of memes in the context of each community, showing, e.g., that racist memes are extremely common in fringe Web communities. We also find a substantial number of politics-related memes on both mainstream and fringe Web communities, supporting media reports that memes might be used to enhance or harm politicians. Finally, we use Hawkes processes to model the interplay between Web communities and quantify their reciprocal influence, finding that *local* substantially infl

no bad intentions, others have assumed negative and/or hateful connotations, including outright racist and aggressive undertones [84]. These memes, often generated by fringe communities, are being “weaponized” and even becoming part of political and ideological propaganda [67]. For example, memes were adopted by candidates during the 2016 US Presidential Elections as part of their iconography [19]; in October 2015, then-candidate Donald Trump retweeted an image depicting him as Pepe The Frog, a controversial character considered a hate symbol [58]. In this context, polarized communities within 4chan and Reddit have been working hard to create new memes and make them go viral, aiming to increase the visibility of their ideas—a phenomenon known as “attention hacking” [64].

Motivation. Despite their increasingly relevant role, we have very little measurements and computational tools to understand the origins and the influence of memes. The online information ecosystem is very complex; social networks do not operate in a vacuum but rather influence each other as to how information spreads [86]. However, previous work (see Section 6) has mostly focused on social networks in an isolated manner.





motorsport.com

How it works



1

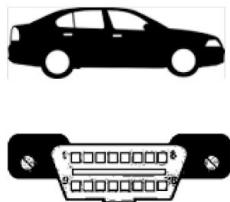


2



3

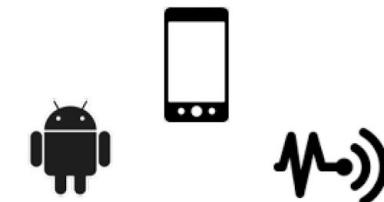
a) Vehicular connection module



{OBD API}

Bluetooth

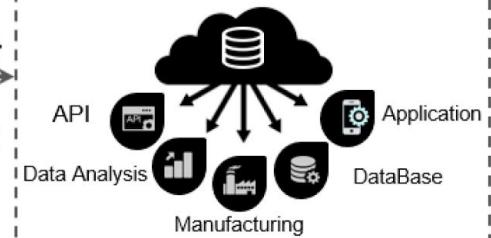
b) Data capture module



{REST API}

(GPRS, 3G or 4G) /Wi-Fi

c) Data storage module



Extremoz



189

256

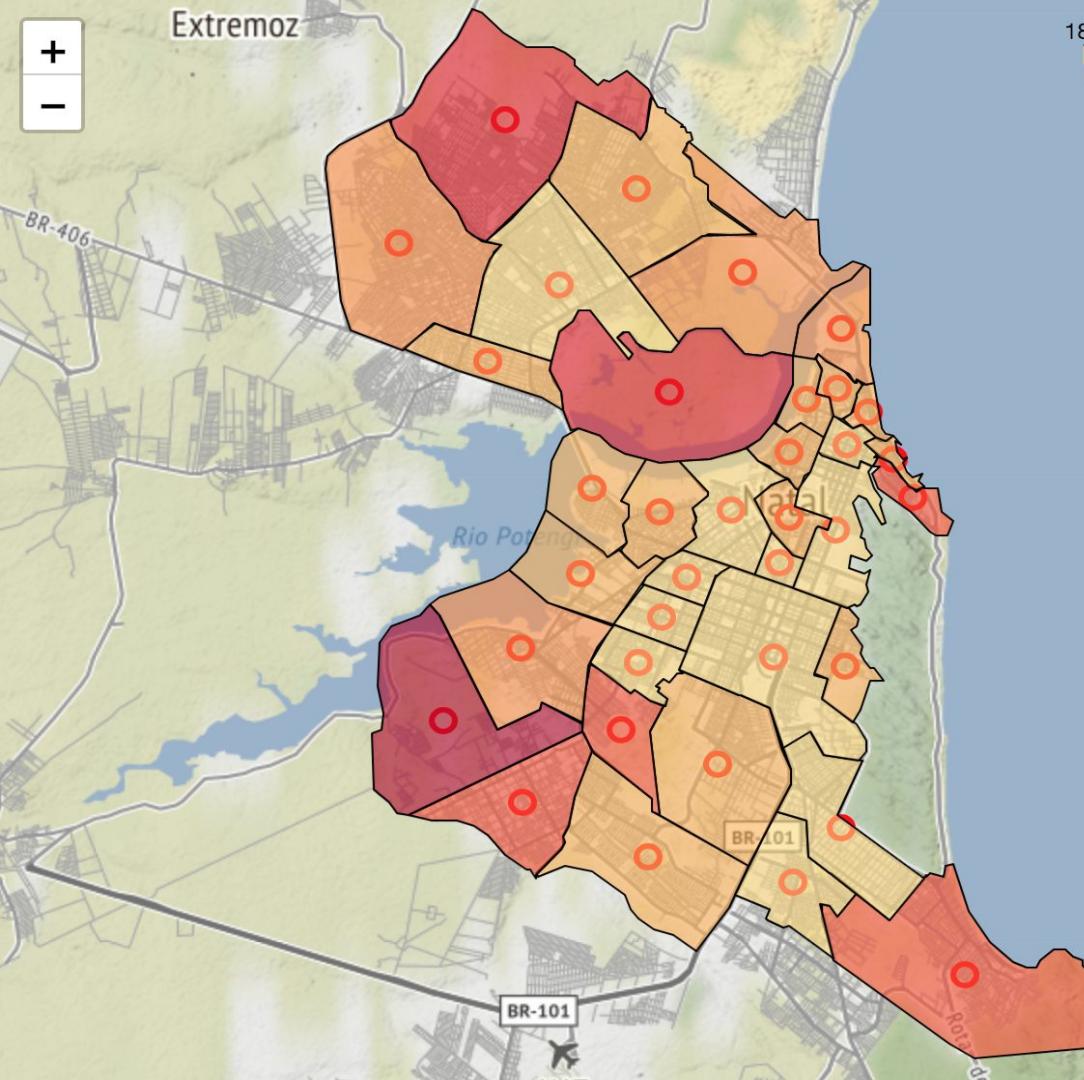
323

390

457

523

Uber ETA (avg.) during 23 days (18 oct to 31 oct)



BR-101

Leaflet

Measuring Depression Symptom Severity from Spoken Language and 3D Facial Expressions

Albert Haque¹ Michelle Guo¹ Adam S Miner^{2,3} Li Fei-Fei¹

¹Department of Computer Science, Stanford University

²Department of Psychiatry and Behavioral Sciences, Stanford University

³Department of Health Research and Policy, Stanford University

Abstract

With more than 300 million people depressed worldwide, depression is a global problem. Due to access barriers such as social stigma, cost, and treatment availability, 60% of mentally-ill adults do not receive any mental health services. Effective and efficient diagnosis relies on detecting clinical symptoms of depression. Automatic detection of depressive symptoms would potentially improve diagnostic accuracy and availability, leading to faster intervention. In this work, we present a machine learning method for measuring the severity of depressive symptoms. Our multi-modal method uses 3D facial expressions and spoken language, commonly available from modern cell phones. It demonstrates an average error of 3.67 points (15.3% relative) on the clinically-validated Patient Health Questionnaire (PHQ) scale. For detecting major depressive disorder, our model demonstrates 83.3% sensitivity and 82.6% specificity. Overall, this paper shows how speech recognition, computer vision, and natural language processing can be combined to assist mental health patients and practitioners. This technology could be deployed to cell phones worldwide and facilitate low-cost universal access to mental health care.

Intelligent Machines

Your smartphone's AI algorithms could tell if you are depressed

Smartphones that are used to track our faces and voices could also help lower the barrier to mental-health diagnosis and treatment.

by Will Knight December 3, 2018

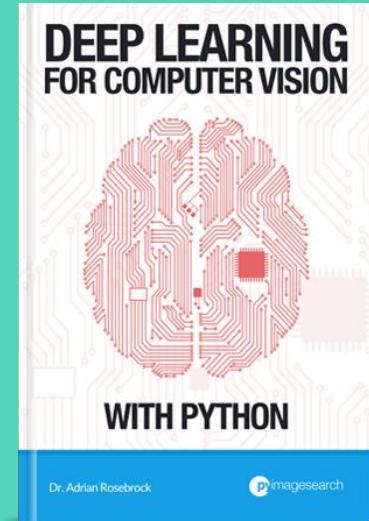


<https://www.ibm.com/blogs/research/2018/09/agropad/>

Ping a Good Doctor



ROBÔ PARA IDENTIFICAÇÃO DE FALHAS EM VEÍCULOS



McCormick & IBM Announce Collaboration

Pioneering the Use of Artificial Intelligence in Flavor and Food Product Development



DATA VIOLENCE

and how bad
engineering
choices can
damage society



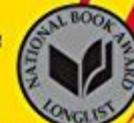
COMO MENTIR com ESTATÍSTICA

Darrell Huff

ILUSTRADO POR
Irving Geis
"MAIS RELEVANTE QUE NUNCA."
BILL GATES



NEW YORK TIMES BESTSELLER



WEAPONS OF MATH DESTRUCTION



HOW BIG DATA INCREASES INEQUALITY
AND THREATENS DEMOCRACY

CATHY O'NEIL

A NEW YORK TIMES NOTABLE BOOK





$$\begin{aligned} & A^* \\ & \hat{y}_i = \beta_0 + \sum_{j=1}^p x_{ij} \beta_j \\ & \text{Loss Function: } \frac{1}{2} \left(y_i - \hat{y}_i \right)^2 \\ & \text{Optimization: } \min_{\beta_0, \beta_1, \dots, \beta_p} \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 \end{aligned}$$

$$A = U\Sigma V^T$$

$$\begin{aligned} & A = U\Sigma V^T \\ & \text{SVD Components: } U, \Sigma, V \\ & \text{Matrix Factorization: } \hat{A}_{ij} = \sum_{k=1}^p U_{ik} \Sigma_{kk}^{-1} V_{kj} \\ & \text{Loss Function: } \sum_{i,j} \left(A_{ij} - \hat{A}_{ij} \right)^2 \end{aligned}$$

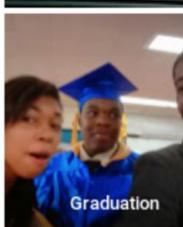


jackyalciné ez de nu blick penthe

@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>



宁波交警行人闯红灯事无安

<https://gizmodo.uol.com.br/reconhecimento-facial-falha-china/>



**Even Amazon says it wants face
recognition to be regulated**

< 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% ⚡

◀ Albums

puppy or bagel

Select



•••○○ Verizon ⌓

10:50 PM

4% ⚡

◀ Back

labradoodle or fried chicken

Select



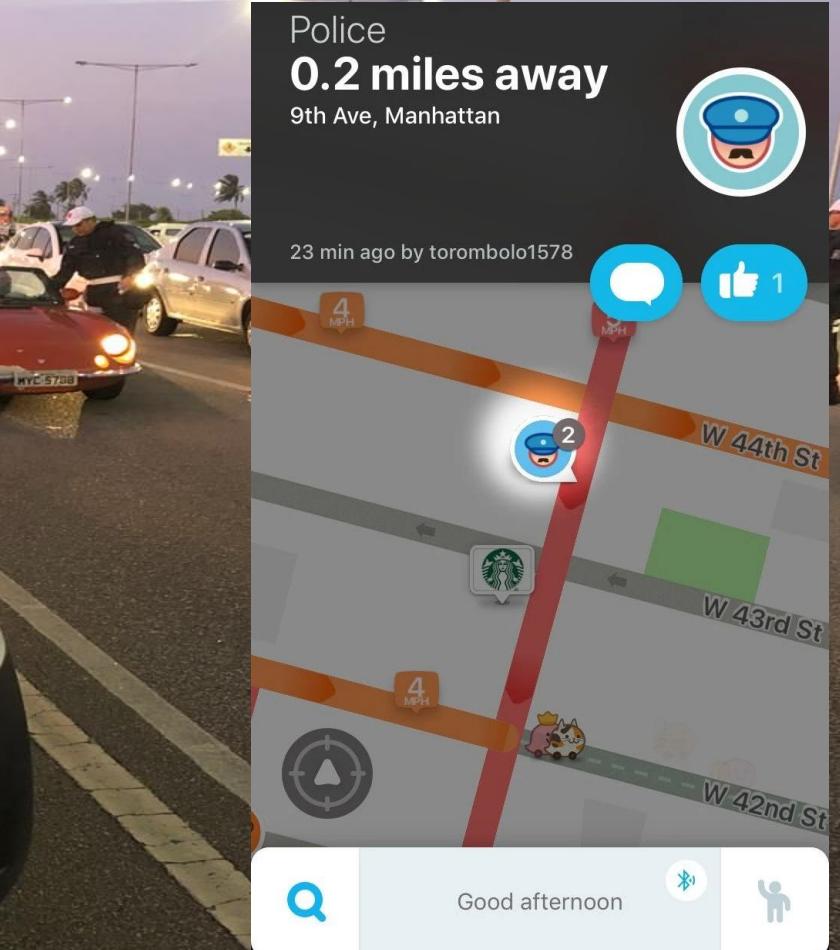


Nice science,
but don't
forget about
the ethics

**Many popular
iPhone apps
secretly record
your screen
without asking**



Google and Waze Must Stop Sharing Drunk-Driving Checkpoints



<https://www.nytimes.com/2019/02/06/nyregion/waze-nypd-location.html>

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.



MODERN DATA SCIENTIST

Data Scientist, the sexiest job of 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

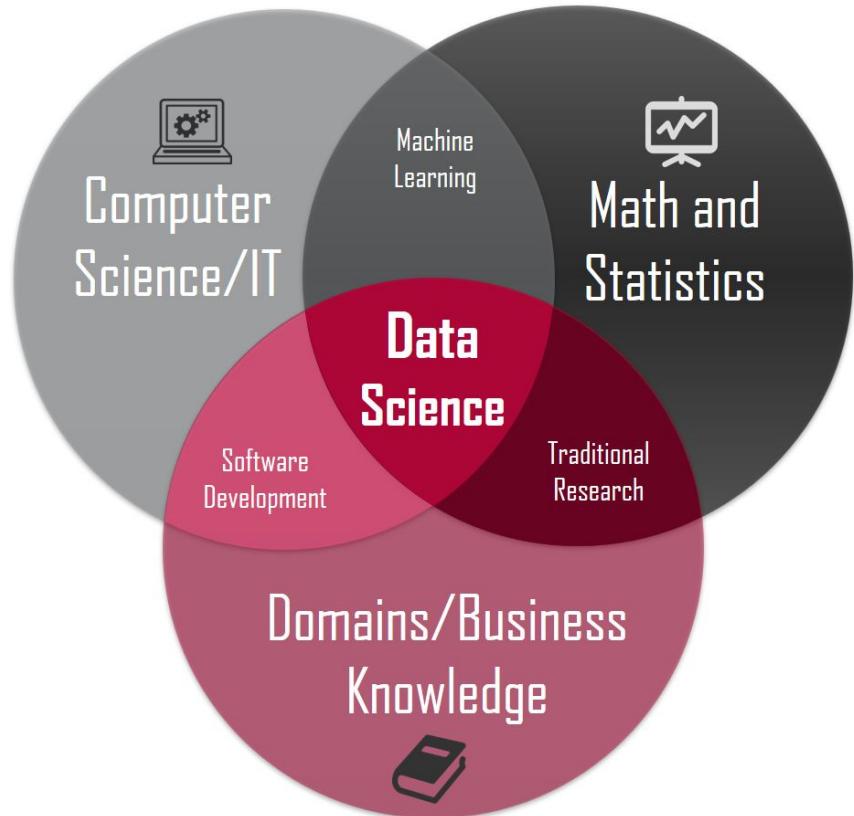
- ★ Computer science fundamentals
- ★ Scripting language e.g. Python
- ★ Statistical computing package e.g. R
- ★ Databases SQL and NoSQL
- ★ Relational algebra
- ★ Parallel databases and parallel query processing
- ★ MapReduce concepts
- ★ Hadoop and Hive/Pig
- ★ Custom reducers
- ★ Experience with 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 visualization tools e.g. Flare, D3.js, Tableau

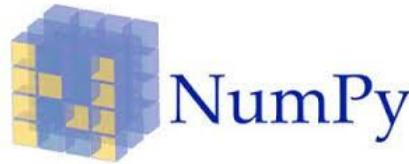


jobs

Learning by doing

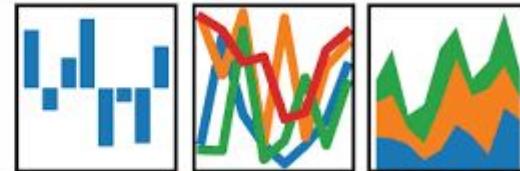
#no_exams #assessments_all_time #projects
#dont_reinvent_the_wheel





NumPy

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



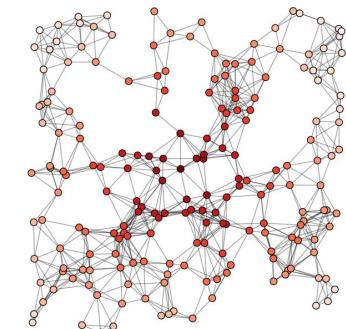
K Keras



NetworkX
PyGSP

matplotlib

Folium



Seaborn



bokeh



Leaflet



Requests

Beautiful Soap





IMD0905 Data Science I Syllabus - 2019.1

Data Science Foundation Data Pipeline

- Store
- Collect
- Cleaning
- Analyse
- Visualize

Clustering and network analysis

HARD WORKING part #01

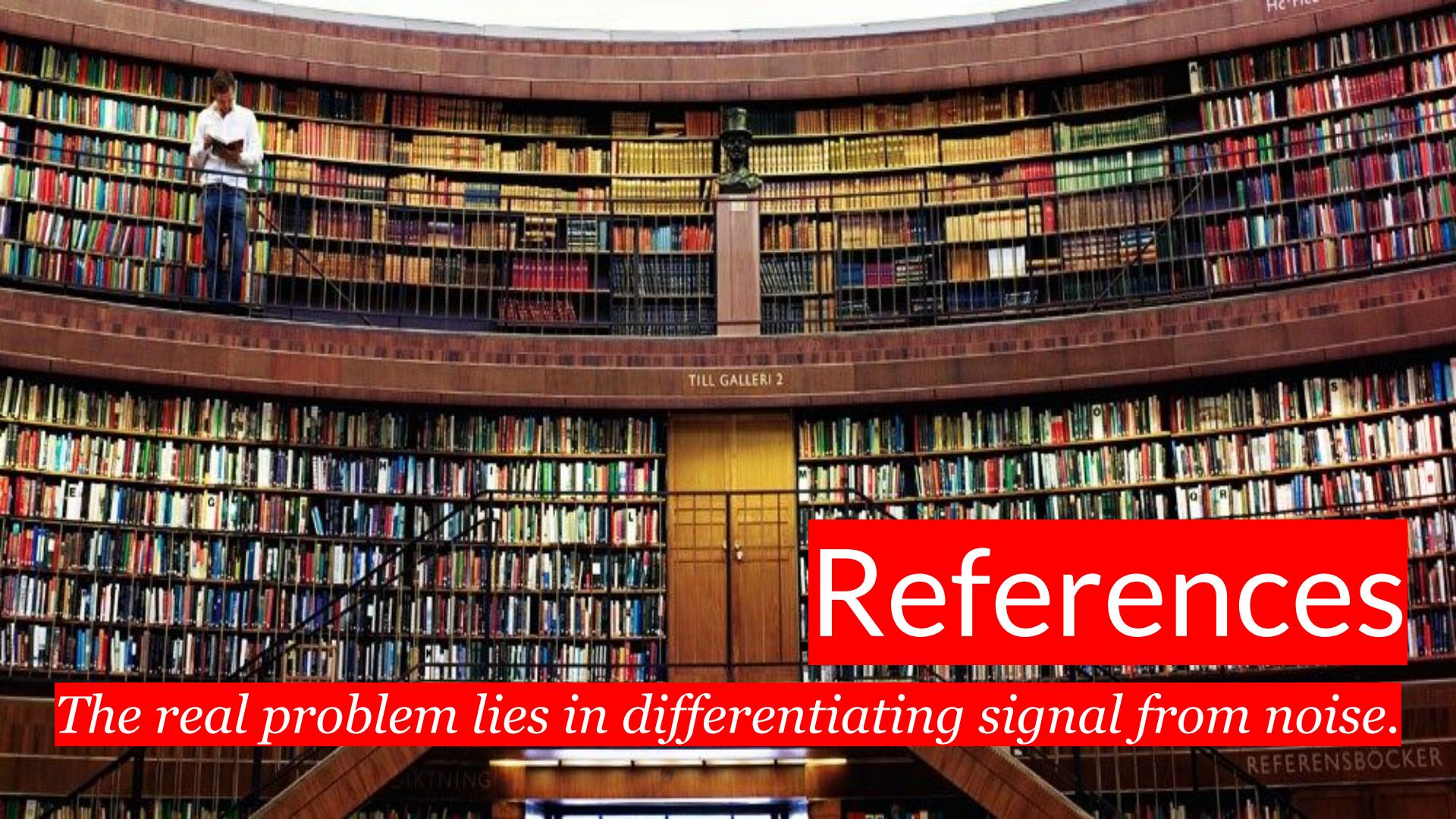
- Python Crash Course
- Pandas
- 02 Assessments (individual, group[2])
- Datacamp Courses
 - 16 courses
- End: 25/03

KEEP ON RUNNING part #02

- Visualization using maps
- Data cleaning
- 01 Big Assessment (group[3-4])
- Datacamp Courses
 - 1 course
 - 10k points (whatever)
- End: 29/04

FINAL SPRINT part #03

- Optimizing Dataframe
- Python & Database
- Web scraping Introduction
- Network Analysis: mining social web
- 01 Assessment about Optimization & Database (group[2])
- Final Project - Web scraping and network analysis (group)
- Datacamp Courses
 - 5 courses
- End: 26/06



References

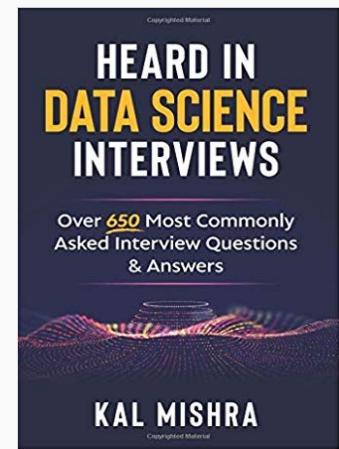
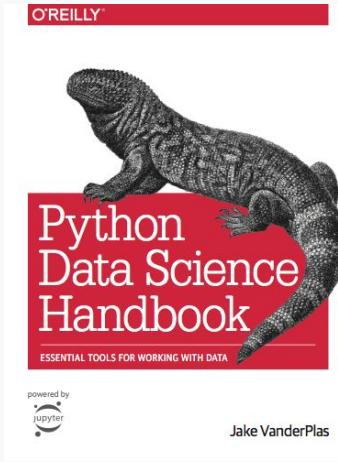
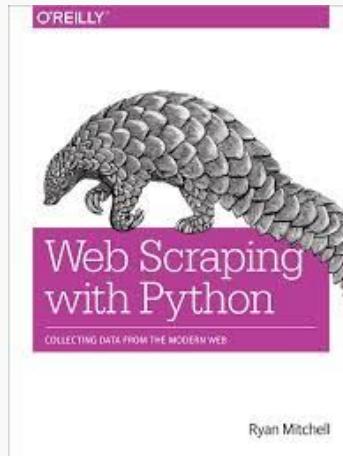
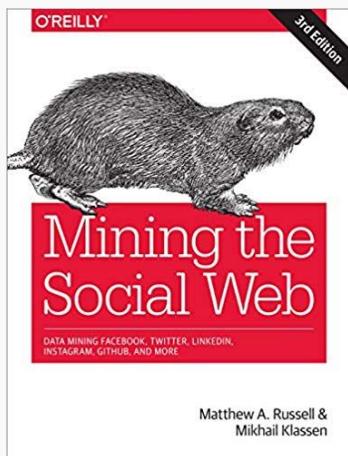
The real problem lies in differentiating signal from noise.

References Online

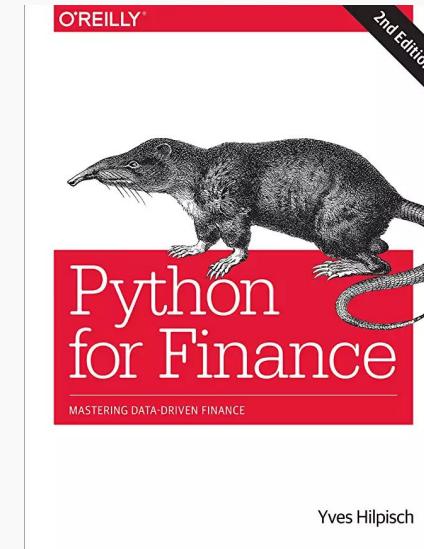
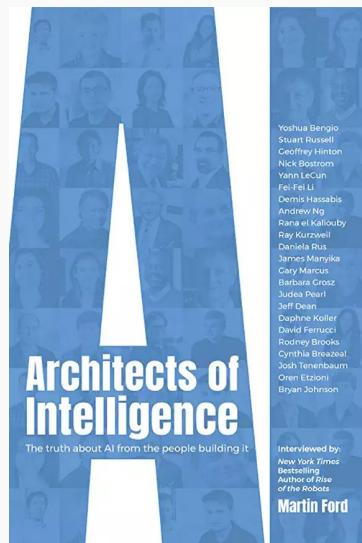


Data Science
Academy

References



References (out of scope)





Faça parte da maior comunidade de Data Science do Brasil!

Cresça junto com profissionais de data science, machine learning, big data e inteligência artificial

Inscreve-se na newsletter quinzenal

Digite seu melhor email

Inscreve-se



Participar também da comunidade no Slack





http://bit.do/ds_podcast_intro

hipsters
ponto tech



#134

Primeiros Passos em Data Science: Do Excel e BI ao Python

<https://dataelixir.com>

Subscribe

Join now to get the latest help articles from our team.

Enter your e-mail

Subscribe

