



Lesson #03

Data Science Fundamentals

Outline

1. Data Science Fundamentals
2. Ethic by Design
3. Defining a problem
4. Platforms & Tools
5. Warming up

Calendar

Unit 01

Início	Fim	Descrição
18/02/2020	18/02/2020	Course Outline
20/02/2020	20/02/2020	IMD Talk
27/02/2020	27/02/2020	Fundamentals of Data Science
03/03/2020	03/03/2020	Introduction to Pandas
05/03/2020	05/03/2020	Exploring data with Pandas
10/03/2020	10/03/2020	Data Cleaning Basics
12/03/2020	12/03/2020	<<Project 01>>
17/03/2020	17/03/2020	Data Aggregation
19/03/2020	19/03/2020	Combining Data with Pandas
24/03/2020	24/03/2020	Transforming Data with Pandas
26/03/2020	26/03/2020	Working with String in Pandas
31/03/2020	31/03/2020	Regular Expression
02/04/2020	02/04/2020	Working with missing and duplicate data
07/04/2020	07/04/2020	<< Project 02>> <<end of unit 01>>

Syllabus:

1. Pandas Fundamentals
2. Data Cleaning & Analysis

Grades

Notebooks - 10%

Project 01 - 30%

Project 02 - 40%

Datacamp (20k XP) - 20%

http://bit.do/imd_datascience



Suggestions

>> Data Science for Business >> Introduction to Python >> Intermediate Python



2.1 ▾ Stepping Stones

The Role of Academia in Data Science Education

by Rafael A. Irizarry

now on branch
#public ▾

Contributors (1)

Published Jan 31, 2020

DOI 10.1162/99608f92.dd363929

Show All Details ▾

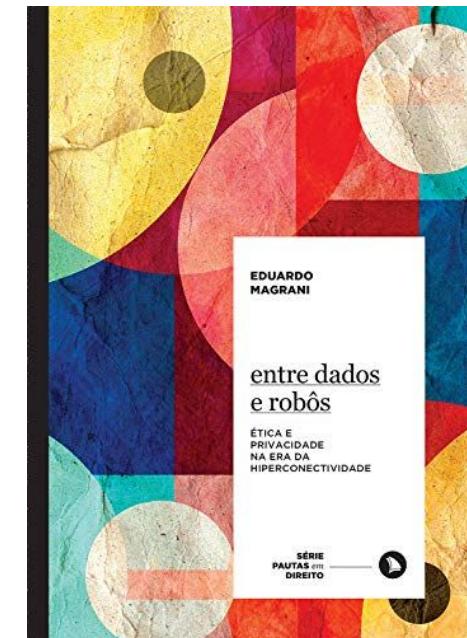
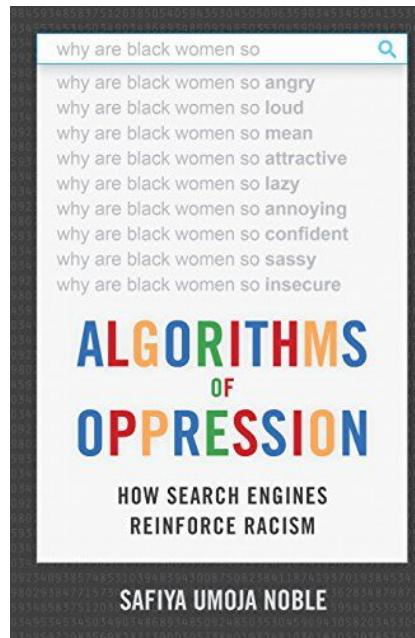
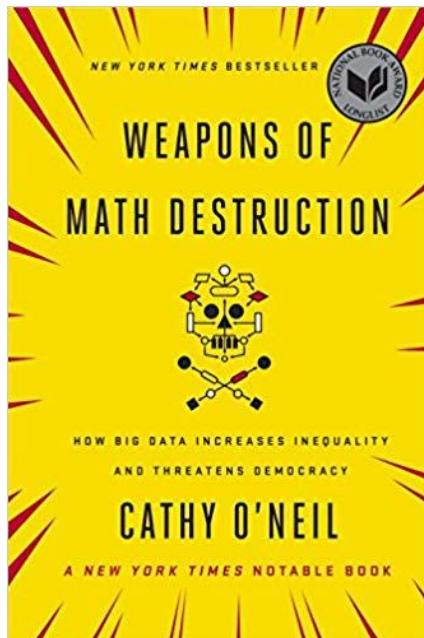
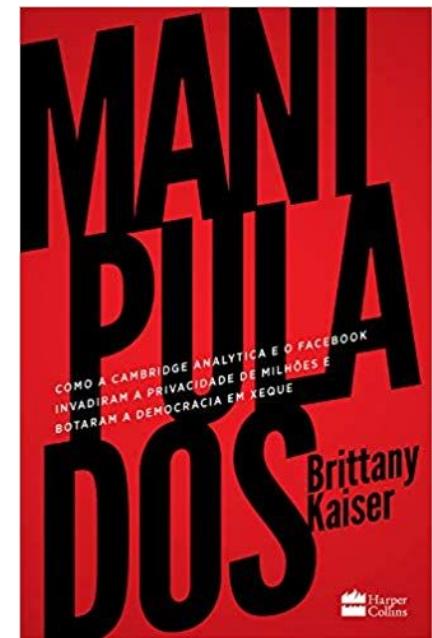
ABSTRACT

As the demand for data scientists continues to grow, universities are trying to figure out how to best contribute to the training of a workforce. However, there does not appear to be a consensus on the fundamental principles, expertise, skills, or knowledge-base needed to define an academic discipline. We argue that data science is not a discipline but rather an umbrella term used to describe a complex process involving not one data scientist possessing all the necessary expertise, but a team of data scientists with nonoverlapping complementary skills. We provide some recommendations for how to take this into account when designing data science academic programs.

<http://bit.do/hdsrv>

ÉTICA





Alexa has been eavesdropping
on you this whole time



Activate This ‘Bracelet of Silence,’ and Alexa Can’t Eavesdrop



<http://bit.do/braceletblock>



<https://cvdazzle.com/>



Feb 22, 2020

New software for testing
CV Dazzle looks
scheduled for release at
HeK work in April. With
thanks to a privacy grant
from NL_Net

NB: this site will be relaunched in 2020. Much of the content here was developed for the haarcascade face detection algorithm, which was widely used between 2010-2015 but has now been deprecated as neural networks face detection algorithms have become more widespread.

LOOKS

STYLE TIPS

RELATED PROJECTS

COLLABORATIONS

PRESS

REFERENCES

CONTACT

All content © Adam Harvey 2010
– 2020 unless otherwise noted.
@adamhrv



by Coreana Museum of
Art + Soobin Academy +
G-square Model
Academy

Anti Face

This face is unrecognizable to
several state-of-art face
detection algorithms.



Camouflage from face detection.

CV Dazzle explores how fashion can be used as camouflage from face-detection technology, the first step in automated face recognition.

The name is derived from a type of World War I naval camouflage called Dazzle, which used cubist-inspired designs to break apart the visual continuity of a battleship and conceal its orientation and size. Likewise, CV

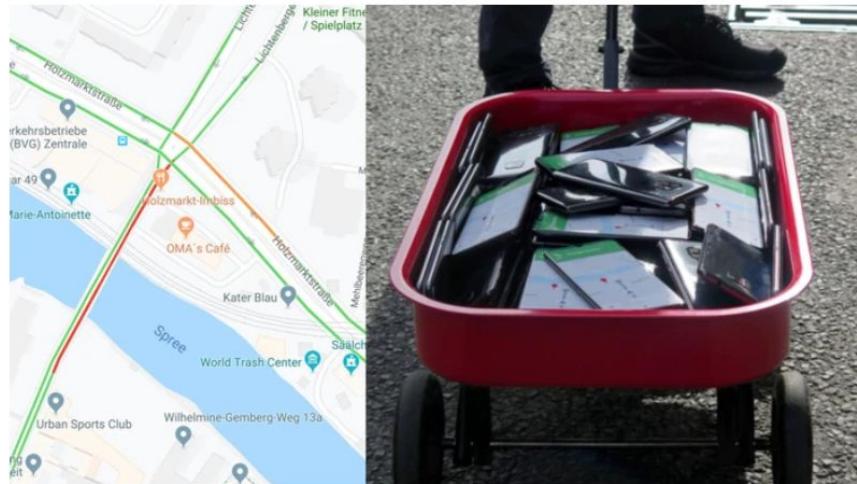
From all appearances, deception has always been critical to daily survival—for human and non-human creatures alike—and, judging by its current ubiquity, there is no end in immediate sight

This Man Created Traffic Jams on Google Maps Using a Red Wagon Full of Phones

By pulling 99 phones down empty streets, artist Simon Weckert made it look like they were gridlocked on Google Maps.

By Matthew Gault

Feb 3 2020, 12:41pm  Share  Tweet  Snap



<http://bit.do/virtualtraffic>





silvio meira ✅ @srlm · Feb 25

12

rara oportunidade de entender como SISTEMAS afetam pessoas: estudo de **#ALGORITMO** de RETENÇÃO de SEGURADOS mostra que a essência é 1. cobrar MENOS de quem tem mais chance de mudar de seguradora e 2. MAIS de quem é fiel: bit.ly/39XSkXo. em que serviços isso rola no BRASIL?



DATA VIOLENCE

and how bad
engineering
choices can
damage society





Research | Computer Vision

Does object recognition work for everyone? A new method to assess bias in CV systems

June 07, 2019 Written by Terrance DeVries, Ishan Misra, Changhan Wang, Laurens van der Maaten

Share



DOLLAR STREET

X



In the news people in other cultures seem stranger than they are.
We visited 264 families in 50 countries and collected 30,000 photos.
We sorted the homes by income, from left to right.

See how people really live

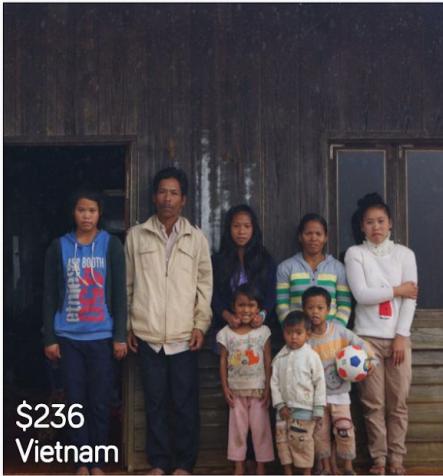
Quick tour

Maybe later



Families in the World by income

English ▾



Soap



Country of Origin: Nepal
Prediction: Food

Spices



Country of Origin: Philippines
Prediction: Beer

Toothpaste



Country of Origin: Burundi
Prediction: Wood



Country of Origin: UK
Prediction: Toiletry



Country of Origin: USA
Prediction: Spice



Country of Origin: USA
Prediction: Toothpaste

"when life gives you lemons make lemonade"

the lemons I get in life:



Yann LeCun @ylecun · 7h

People are biased.

Data is biased, in part because people are biased.

Algorithms trained on biased data are biased.

But learning algorithms themselves are not biased.

Bias in data can be fixed.

Bias in people is harder to fix.

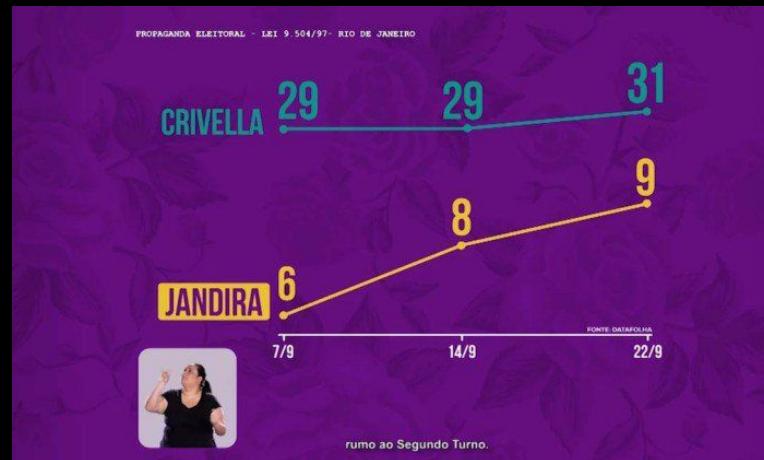


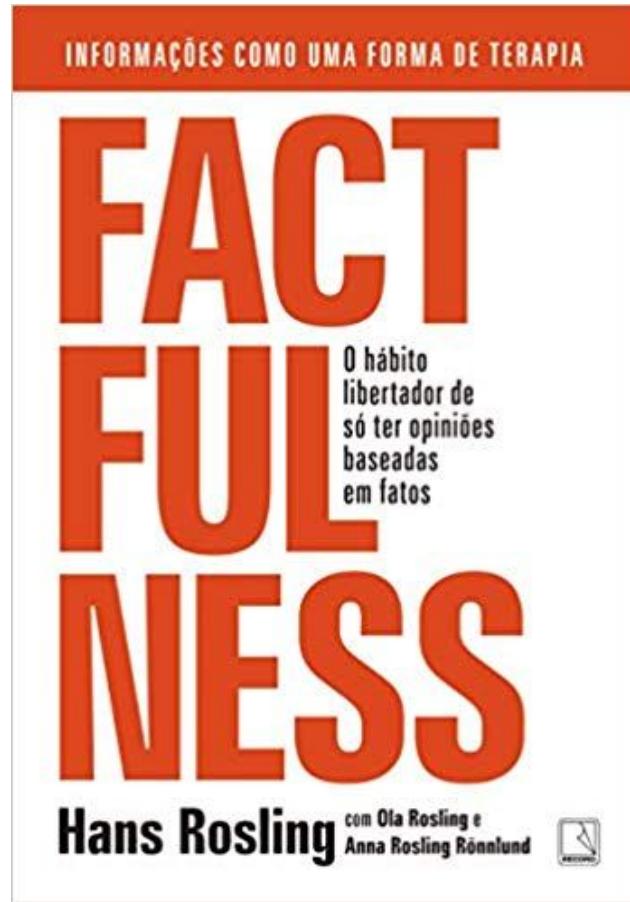
Biased Algorithms Are Easier to Fix Than Biased People

Racial discrimination by algorithms or by people is harmful — but that's where the similarities end.

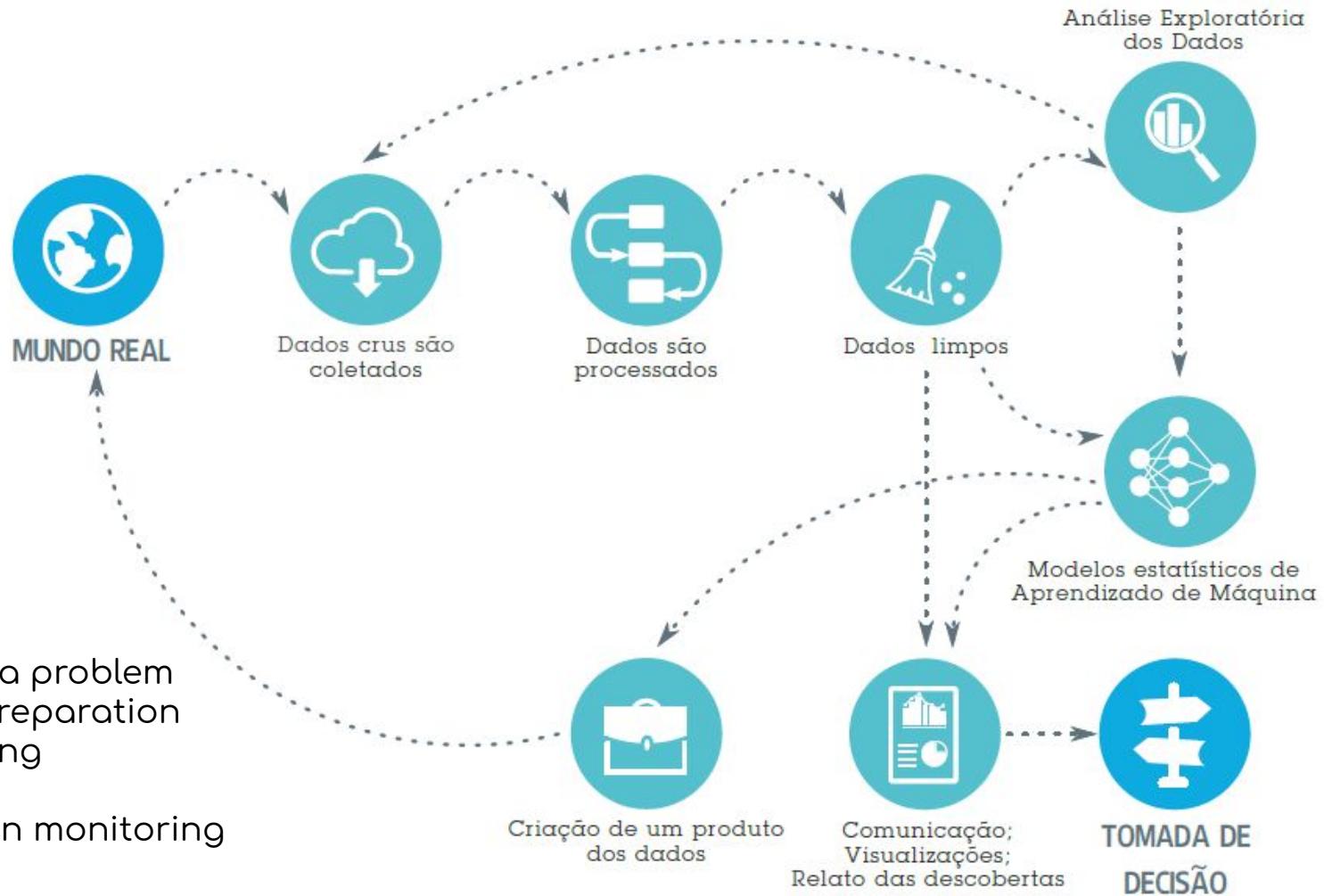
nytimes.com

Bias is everywhere!!!!





A book that destroys myths by presenting facts and statistics in a clear and fun way. Factfulness is an urgent and essential book that will change the way you see the world and empower you to respond to the crises and opportunities of the future.



1. Define a problem
2. Data preparation
3. Modeling
4. Deploy
5. Solution monitoring



1. Data preparation
2. Modeling
3. Deploy
4. Solution monitoring
5. Define a problem

“take this data and see
what you can find”

How to define a problem?



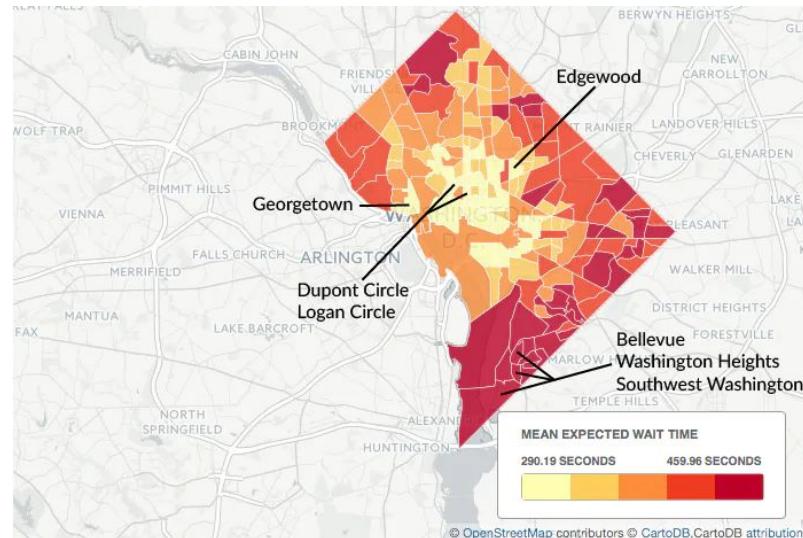
TARGET

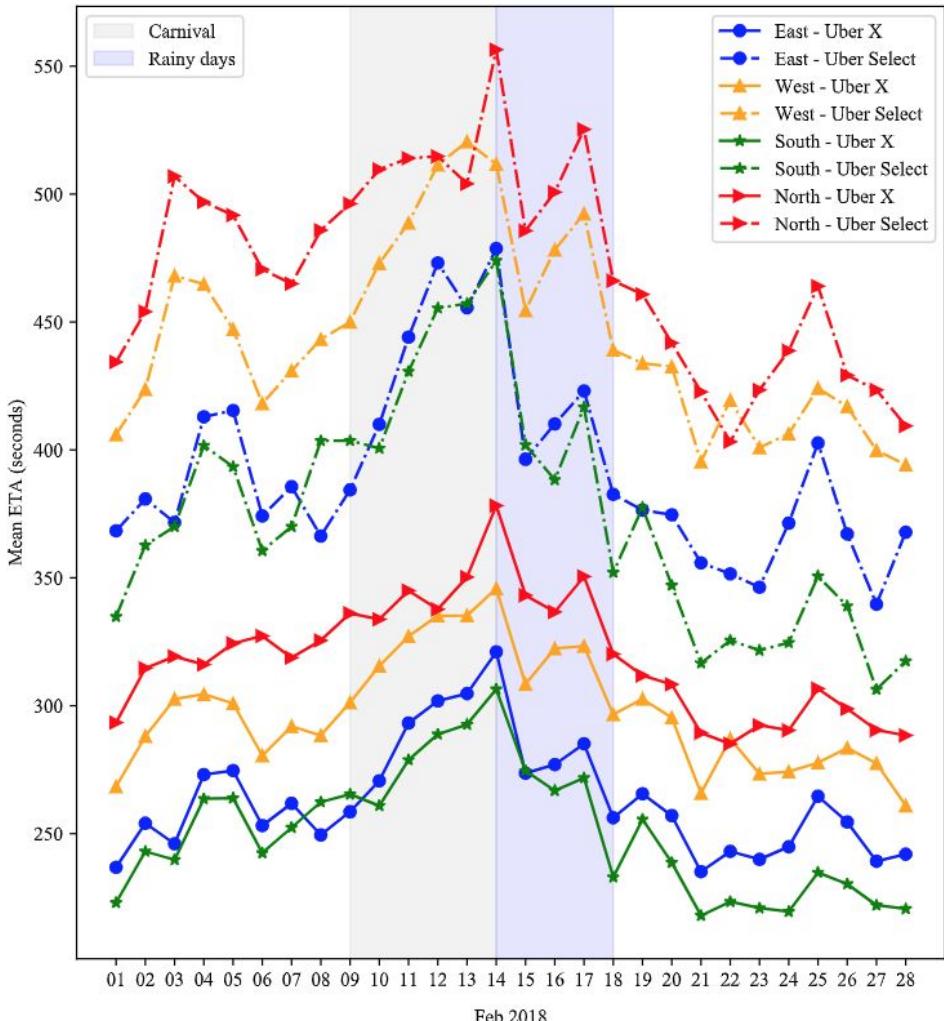
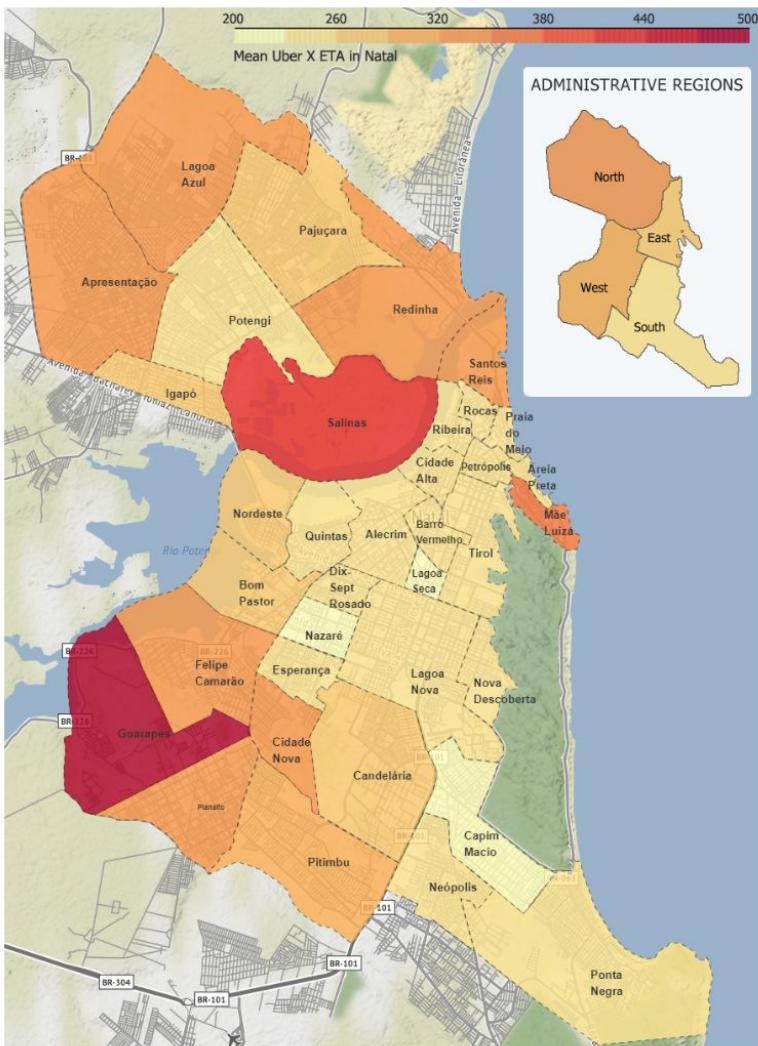
1. What problem do we want to solve?
2. Choose the simplest problem (gain know-how)
3. What business process can we affect?
 - a. What will be the deliverable?
4. How will we know if the solution is working?
 - a. Metrics: primary and secondary
5. Project duration
 - a. Agile methods
 - b. We will only really know how it will impact the business by putting it into production



Economic Policy

Uber seems to offer better service in areas with more white people. That raises some tough questions.







C-MRIC.ORG
Centre for Multidisciplinary Research,
Innovation and Collaboration

LONDON, UNITED KINGDOM

- Pioneer Research and Innovation
- Promote Collaborative Inter-workings
- Encourage Interplay of Different Cultures
- Organise Conferences and Seminars
- A free Membership Organisation



A Preliminary Exploration of Uber Data as an Indicator of Urban Liveability

Agnieszka Bezerra, Goliany Alves, Ivanovitch Silva, Pierangelo Rosati, Patricia Takako Endo, Théo Lynn

DCU
SCHOOL
DCU



Cyber Science is our leading academic housing on pioneering research and innovation on Cyber Situation Awareness, Social Media, Cyber Security and Cyber incident Response.

The purpose is to build bridges between academia and industry, collaboration among different cultures, to leverage knowledge exchange, and to develop solutions that encompass principles, analysis, design, methods and approaches.

Connect, Collaborate & Communicate ...

C-MRIC.ORG
Centre for Multidisciplinary Research,
Innovation and Collaboration

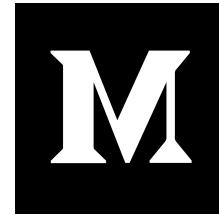
LONDON, UNITED KINGDOM

Medium

Welcome to Medium, where words matter.



Mapeando locais mais visitados em Natal



Usando Python, Pandas e Folium



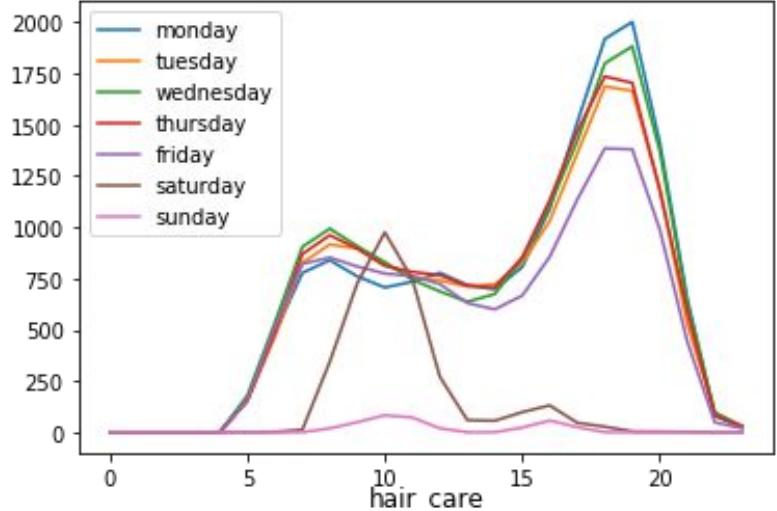
Maradona Morais

[Follow](#)

May 17 · 7 min read

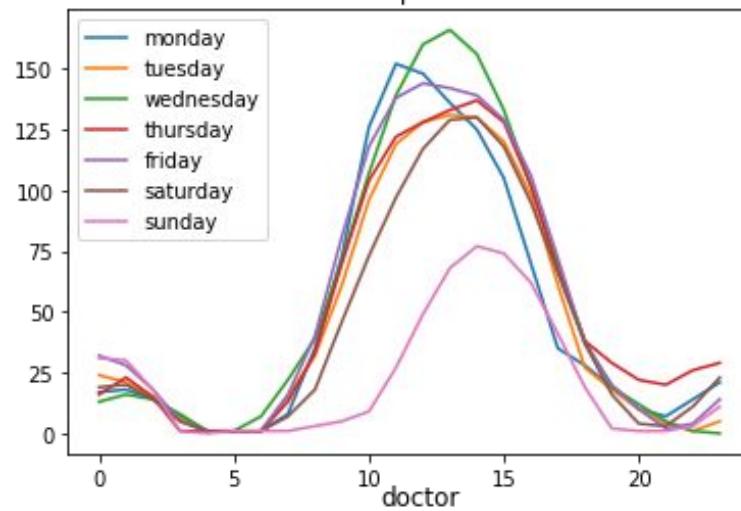


gym

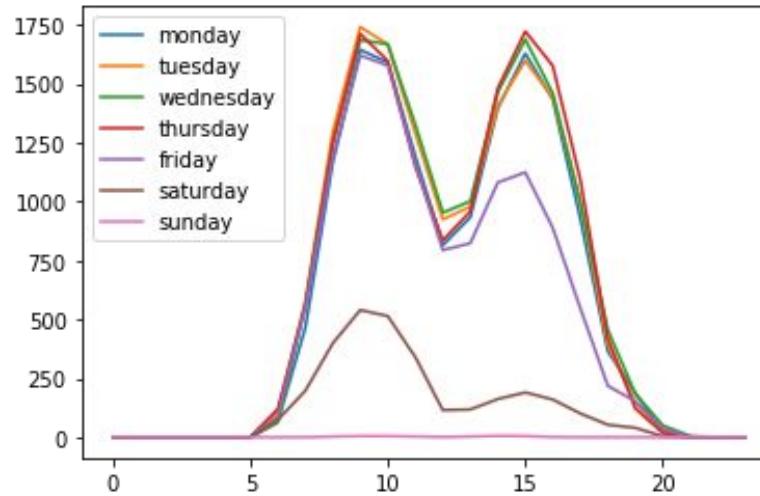
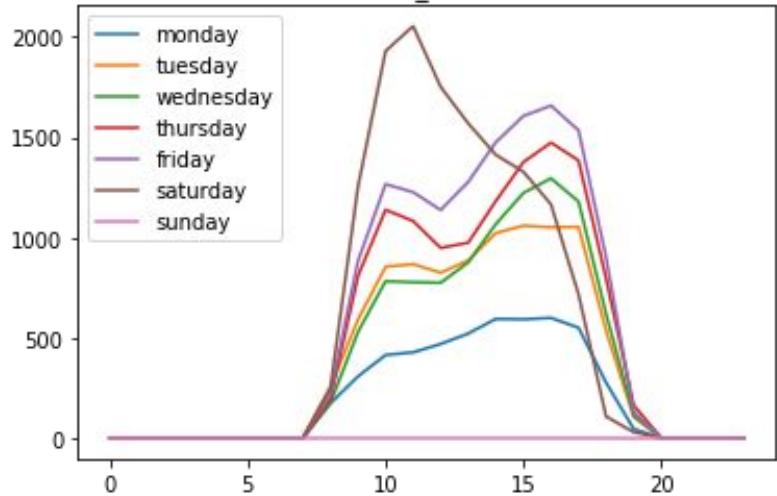


hair_care

airport



doctor



O negócio milionário dos incêndios na Amazônia

Força-tarefa do Ministério Público Federal calcula que queimada de 1.000 ha custa cerca de um milhão no mercado negro da região. Raquel Dodge vê ação orquestrada em fogo. Procuradores não comentam alertas ignorados por Força Nacional

FOLHA DE S.PAULO



Incêndios na Amazônia estão concentrados em propriedades privadas

Áreas privadas cobrem 18% do bioma e concentram 33% dos focos de fogo

EL PAÍS

Agosto tem recorde de focos de incêndio na Amazônia em nove anos, aponta Inpe

Dados do Programa Queimadas mostram 30.901 focos, quase o triplo do ano passado

O GLOBO

AgênciaBrasil

Secretário-geral da ONU está preocupado com queimadas na Amazônia

O secretário-geral da Organização das Nações Unidas (ONU), António Guterres, afirmou hoje (22) por meio de sua conta de Twitter que está "profundamente preocupado" com os incêndios na Floresta Amazônica.



Find a DAAC ▾



Feedback

**EARTHDATA**
Powered by EOSDIS

ABOUT

DATA

COLLABORATE

LEARN

Search datasets, news, articles, and information



Earth Observation Data

● LANCE: NASA Near Real-Time Data and Imagery

● **Fire Information for Resource Management System (FIRMS)**

Data

Disciplines:

FIRMS

[Fires from Space](#)[Data Sources](#)[Active Fire User Guides](#)[FAQs](#)[Disclaimer](#)[LANCE / FIRMS Mailing Lists](#)[Links](#)[Acknowledgements / Citation](#)

Fire Information for Resource Management System (FIRMS)



NASA's Fire Information for Resource Management System (FIRMS) distributes Near Real-Time (NRT) active fire data within 3 hours of satellite observation from NASA's Moderate Resolution Imaging Spectroradiometer (MODIS) and NASA's Visible Infrared Imaging Radiometer Suite (VIIRS).

[MODIS Active Fire Products](#)[VIIRS Active Fire Products](#)

Get hotspot/fire locations



[Ir para o conteúdo 1](#) [Ir para o menu 2](#) [Ir para a busca 3](#) [Ir para o rodapé 4](#)

[ACESSIBILIDADE](#) [ALTO CONTRASTE](#) [MAPA DO SITE](#)

Programa

Queimadas

INSTITUTO NACIONAL DE PESQUISAS ESPACIAIS

Buscar no portal



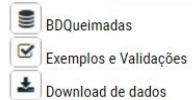
[INPE](#) | [Perguntas Frequentes](#) | [Notícias](#) | [Dados Abertos](#) | [Contato](#) |

SISTEMAS DE MONITORAMENTO



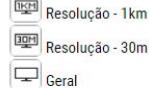
BDQueimadas

1



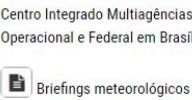
Área Queimada

5



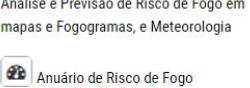
CIMAN Virtual

2



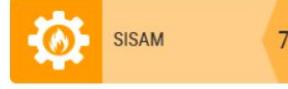
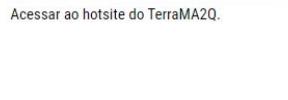
Risco de Fogo

6



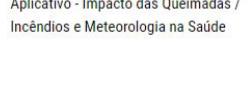
TerraMA2Q

3



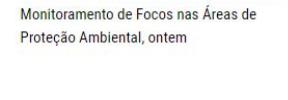
SISAM

7



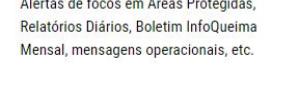
Focos nas APs

4

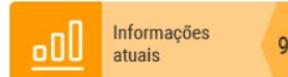


Receber por e-mail

8



RELATÓRIOS E PUBLICAÇÕES



Informações atuais



Boletins internos

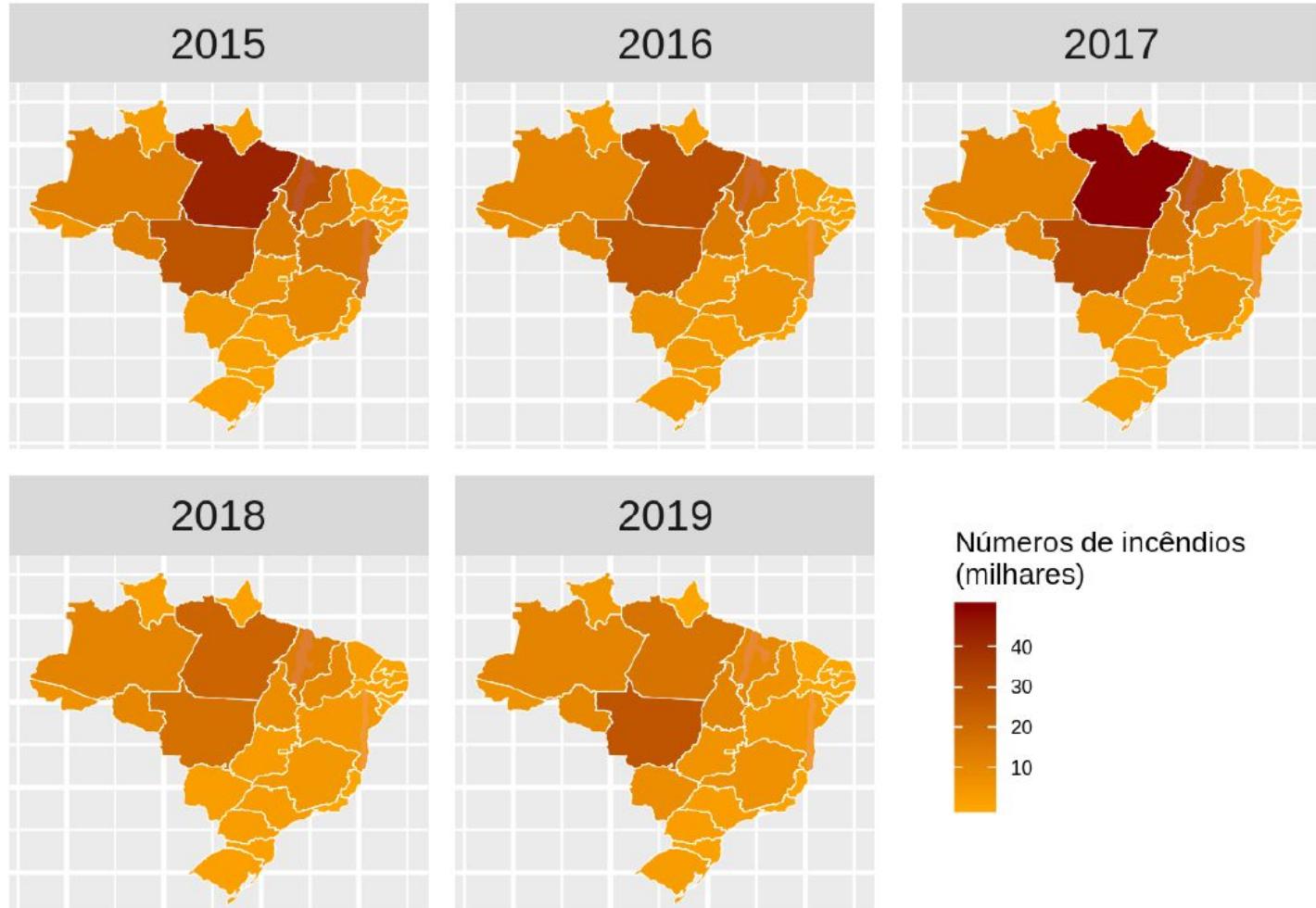


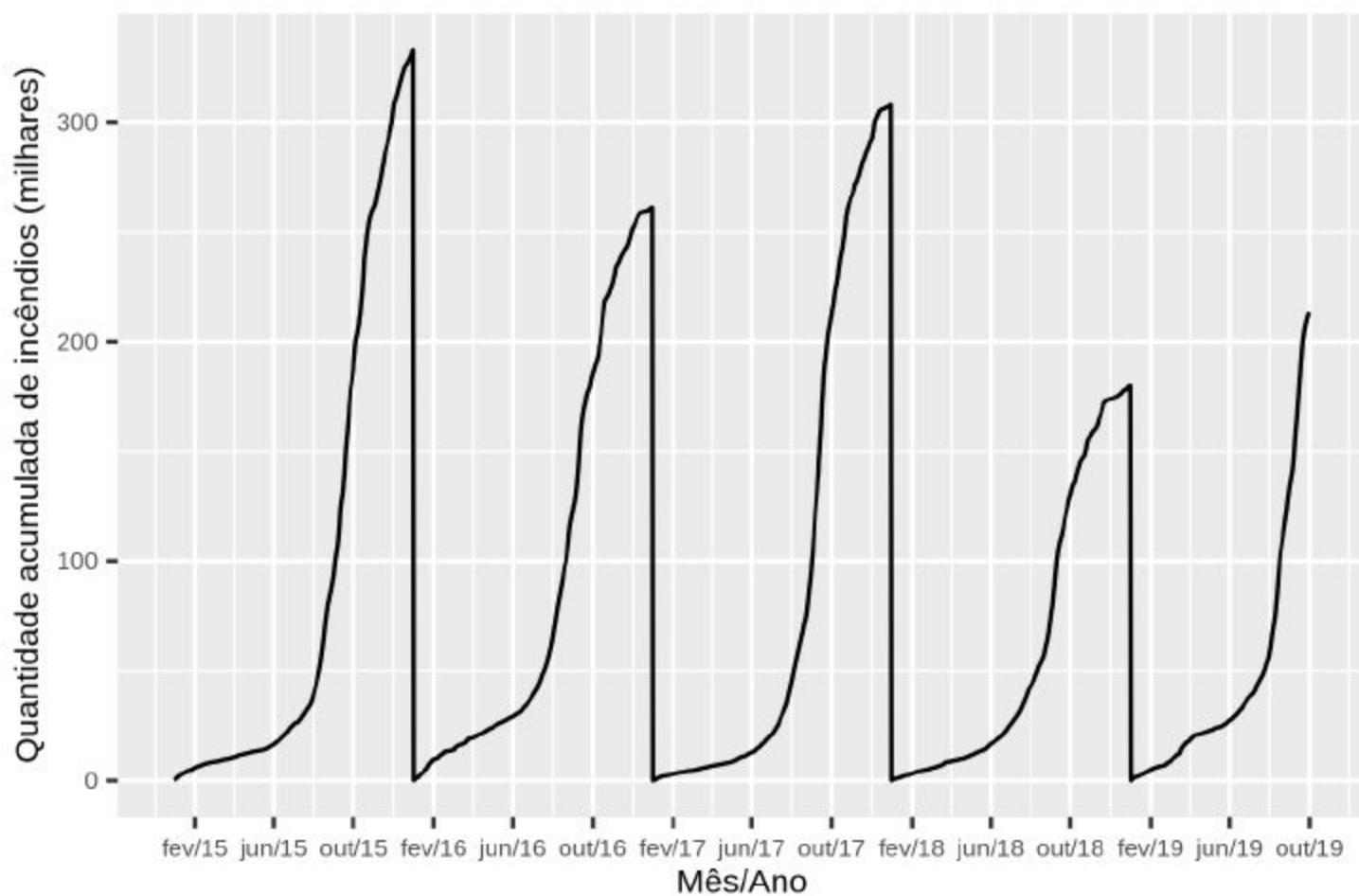
Resumo histórico e animações

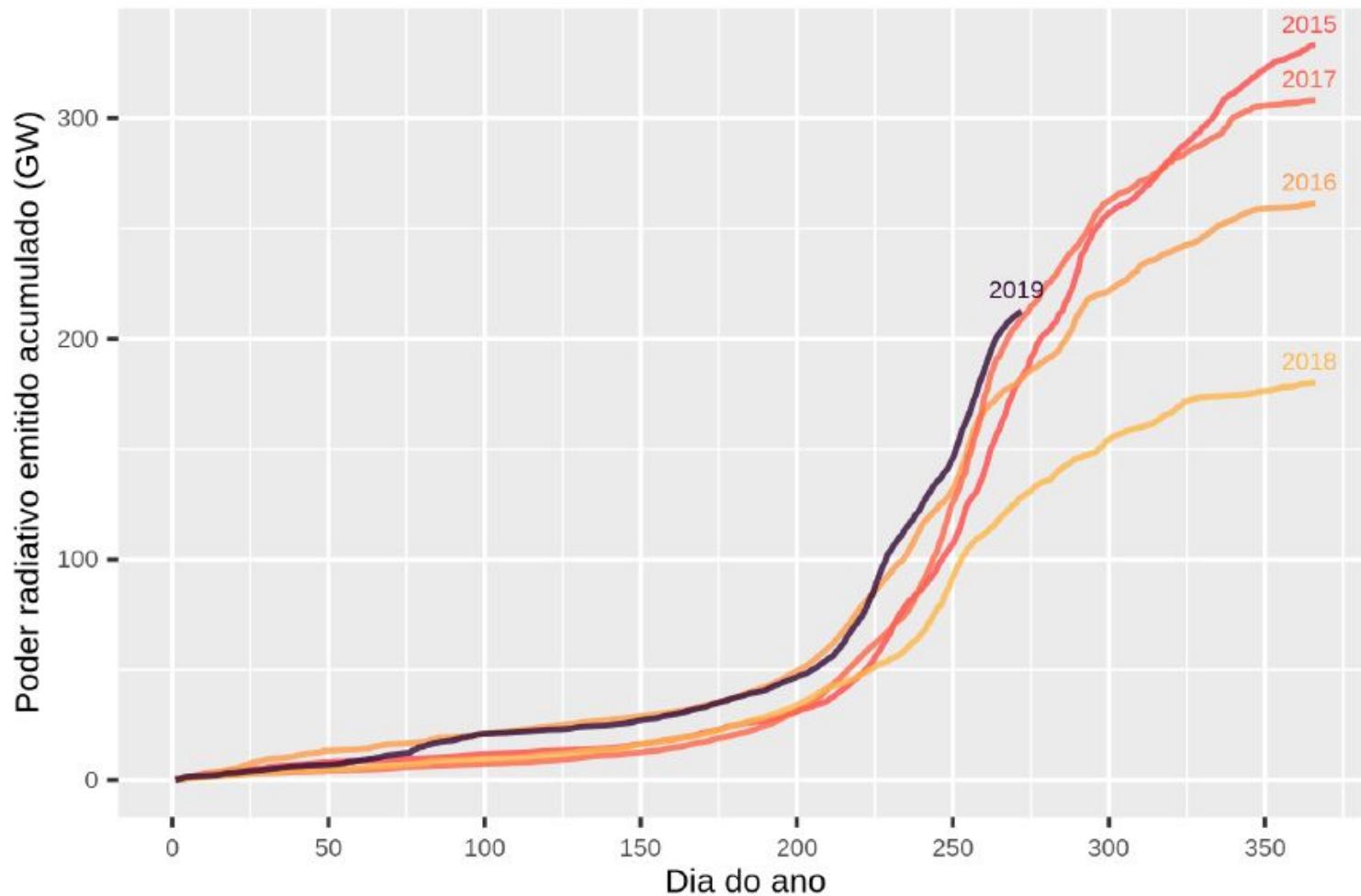


Publicações e impacto

11 12









<https://github.com/odufrn/odufrn-downloader>

38

odufrn / odufrn-downloader

Code Issues 12 Pull requests 1 Actions Projects 2 Wiki Security Insights

Pacote para baixar os dados do portal de dados abertos da UFRN <https://odufrn.github.io/odufrn-downl...>

ufrn open-data downloader ckan requests

237 commits 5 branches 0 packages 3 releases 6 contributors MIT

Branch: master ▾ New pull request Create new file Upload files Find file Clone or download ▾

alvarofpp Merge pull request #96 from odufrn/issue-92 ...	Latest commit b809d50 on 1 Nov 2019
docs Ajustes na documentação	6 months ago
examples Finalizando erros apontados pelo Travis	5 months ago
odufrn_downloader Merge pull request #96 from odufrn/issue-92	4 months ago
tests Merge pull request #96 from odufrn/issue-92	4 months ago
.gitignore issue: #56 Creando exception class to use in File.py and tests/test...	5 months ago
.travis.yml Modifica travis e adiciona novos testes	6 months ago
CONTRIBUTING.md Update CONTRIBUTING.md	6 months ago
LICENSE 0.0.1 Projeto copiado	7 months ago
README.md Adicionando Notebook a pasta examples	5 months ago

androidauto

On your car display.

[GET STARTED](#)

Discover Spotify's Features

With Spotify APIs and SDKs for JavaScript, iOS, and Android — learn how you can develop unique experiences for over 180 million global music fans in as little as a few lines of code.



Total Confirmed
82,548

Confirmed Cases by
Country/Region

78,497 Mainland China

1,766 South Korea

705 Others

528 Italy

245 Iran

189 Japan

93 Singapore

92 Hong Kong

60 US

43 Kuwait

40 Thailand

33 Bahrain

32 Taiwan

26 Germany

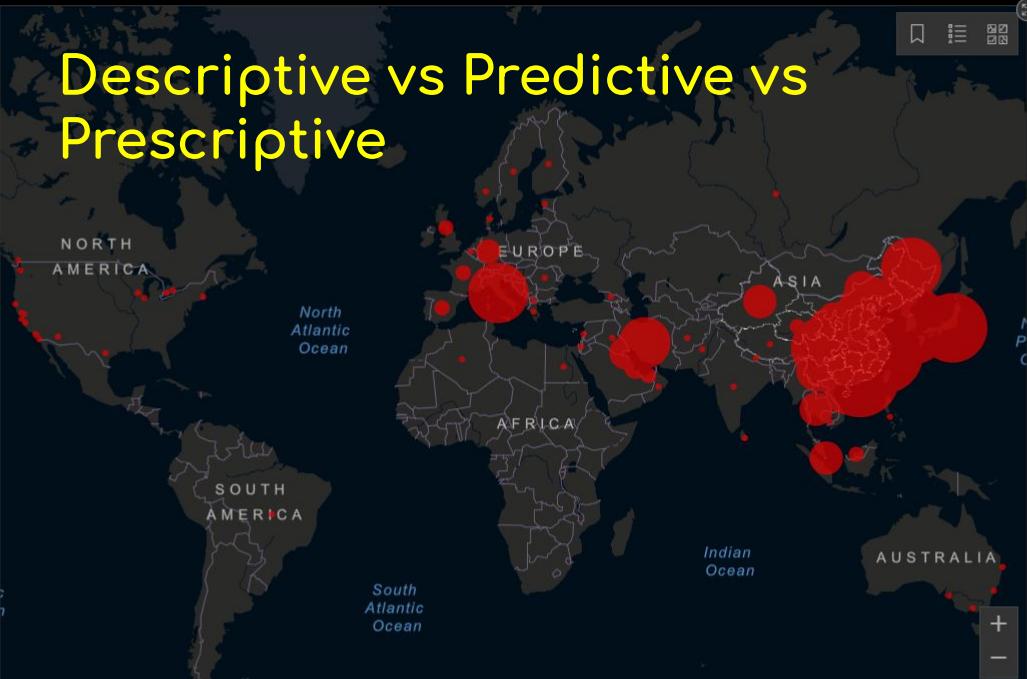
Country/Region

City, St/Prov

Last Updated at (M/D/YYYY)

2/27/2020, 11:13:06 AM

Descriptive vs Predictive vs Prescriptive



Cumulative Confirmed Cases Existing Cases

Lancet Article: [Here](#). Mobile Version: [Here](#). Visualization: JHU CSSE. Automation Support: Esri Living Atlas team and JHU API.

Data sources: WHO, CDC, ECDC, NHC and DXY. Read more in this [blog](#). Contact US.

Downloadable database: GitHub: [Here](#). Feature layer: [Here](#).

Point level: City level - US, Canada and Australia; Province level - China; Country level - other countries.

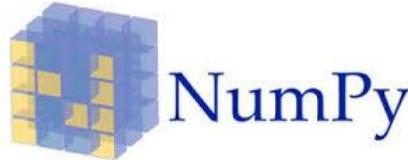
Total Deaths	
2,810	
2,641 deaths	Hubei Mainland China
26 deaths	Iran
20 deaths	Henan Mainland China
14 deaths	Italy
13 deaths	Heilongjiang Mainland China
13 deaths	South Korea
7 deaths	Guangdong Mainland China
6 deaths	

Total Recovered	
33,243	
23,383 recovered	Hubei Mainland China
1,062 recovered	Henan Mainland China
932 recovered	Zhejiang Mainland China
890 recovered	Guangdong Mainland China
804 recovered	Hunan Mainland China
792 recovered	Anhui Mainland China
754 recovered	Jiangxi Mainland China



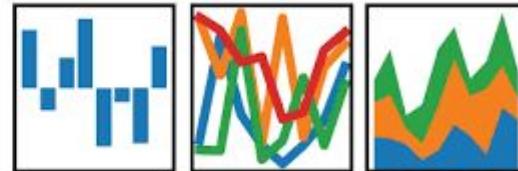


Hands on Data Science



NumPy

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



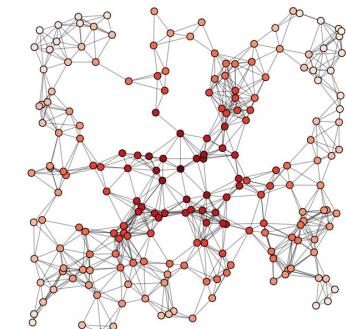
K Keras



NetworkX
PyGSP

matplotlib

Folium



Seaborn



bokeh



Leaflet



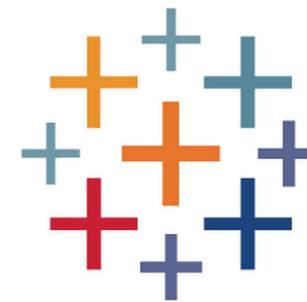
Requests

Beautiful Soap



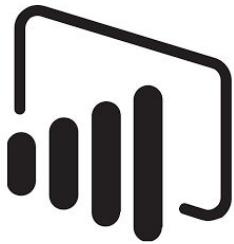


<http://www.pentaho.com/>



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





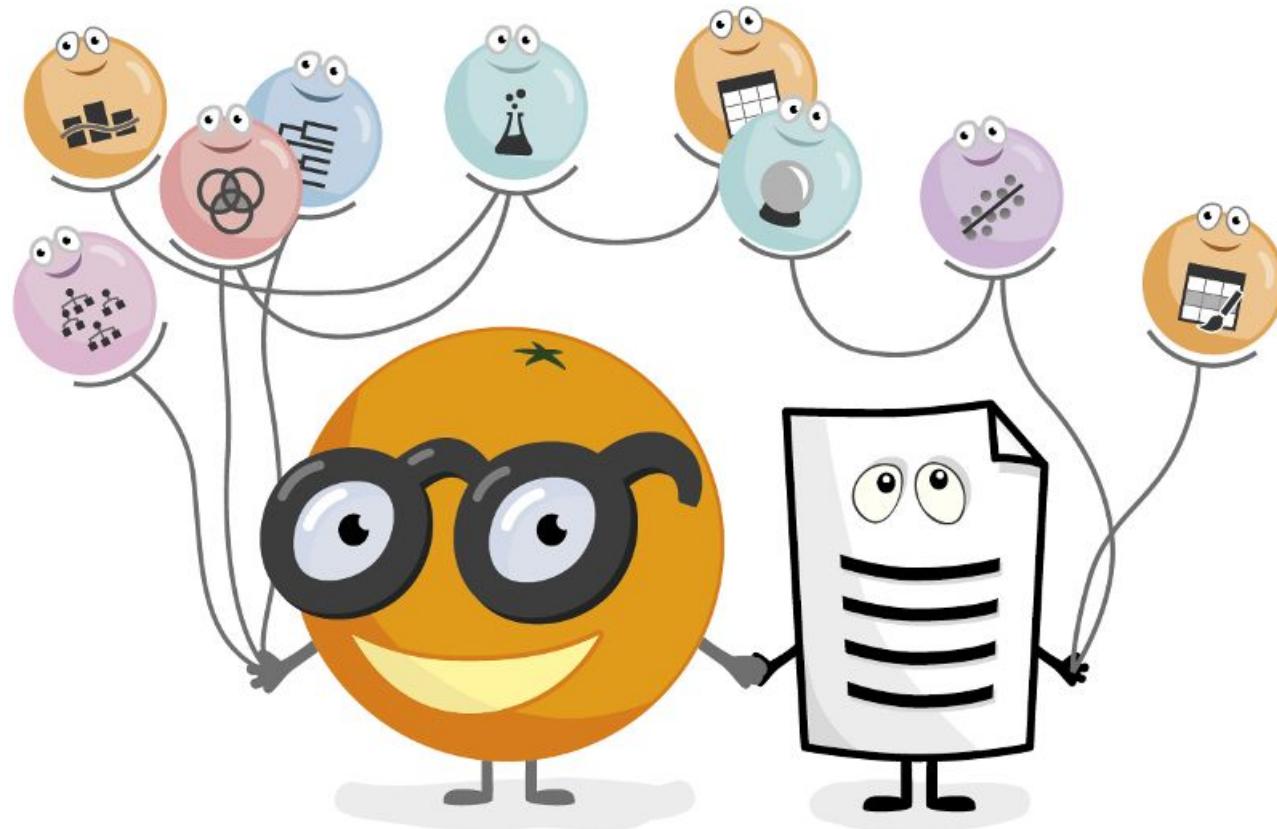
<https://powerbi.microsoft.com>



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



Data Mining Fruitful and Fun



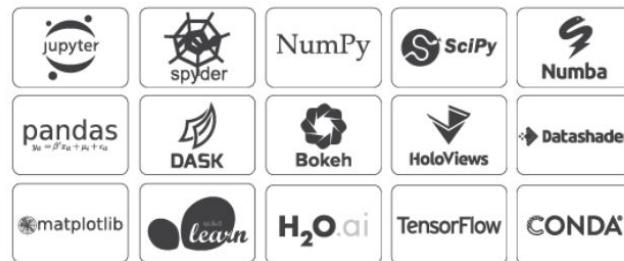
Anaconda Distribution

The World's Most Popular Python/R Data Science Platform

[Download](#)

The open-source [Anaconda Distribution](#) is the easiest way to perform Python/R data science and machine learning on Linux, Windows, and Mac OS X. With over 19 million users worldwide, it is the industry standard for developing, testing, and training on a single machine, enabling *individual data scientists* to:

- Quickly download 7,500+ Python/R data science packages
- Manage libraries, dependencies, and environments with [Conda](#)
- Develop and train machine learning and deep learning models with [scikit-learn](#), [TensorFlow](#), and [Theano](#)
- Analyze data with scalability and performance with [Dask](#), [NumPy](#), [pandas](#), and [Numba](#)
- Visualize results with [Matplotlib](#), [Bokeh](#), [Datashader](#), and [Holoviews](#)



Windows



macOS



Linux





Jupyter is the new Excel



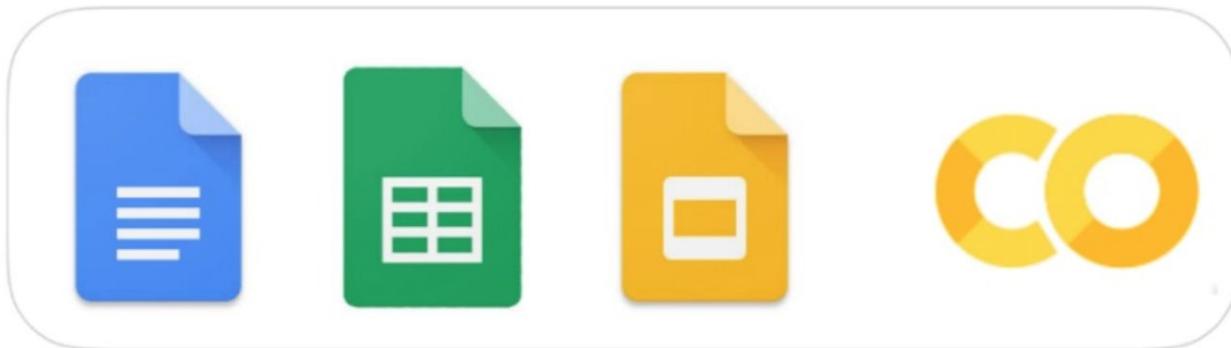
<https://data.berkeley.edu/education/courses/data-100>

Data 100: Principles and Techniques of Data Science



Google Colaboratory

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



Lists of Lists

Lists

Dictionaries

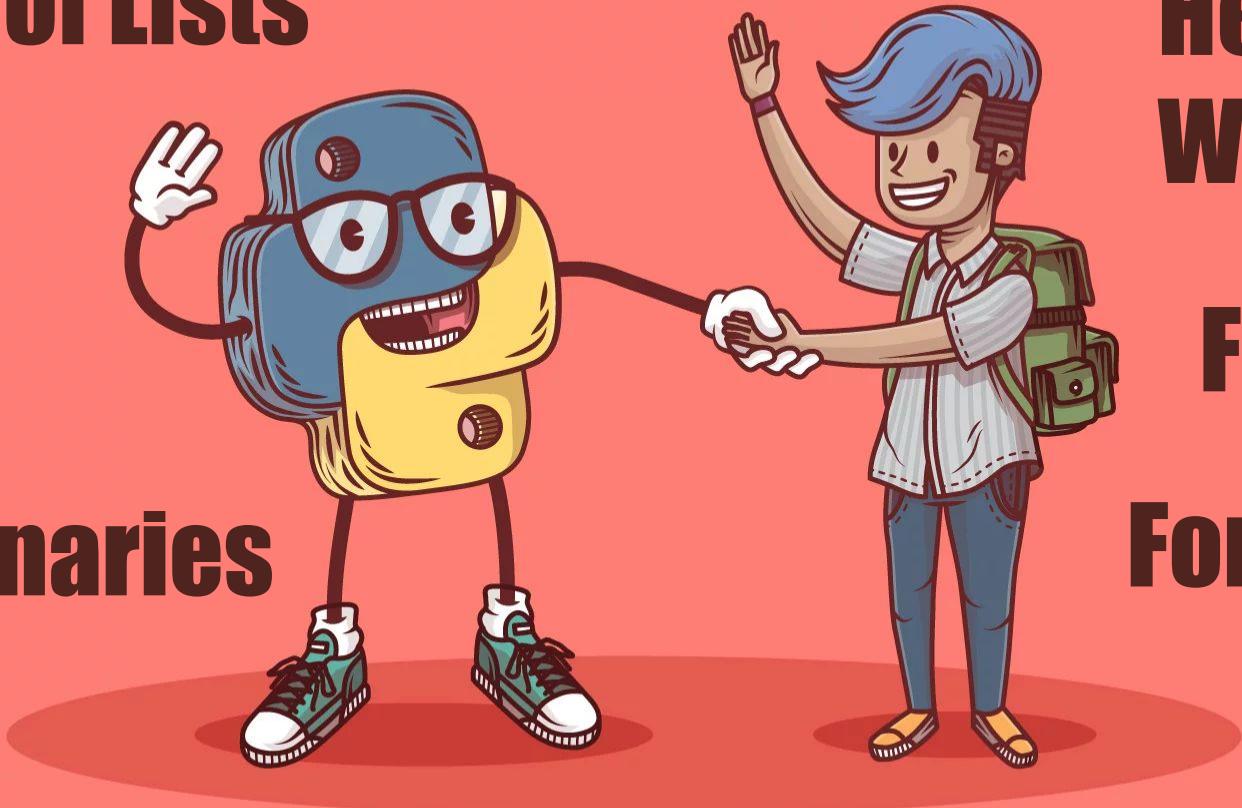
Conditional Statements

Hello World

Files

For Loops

Real Python

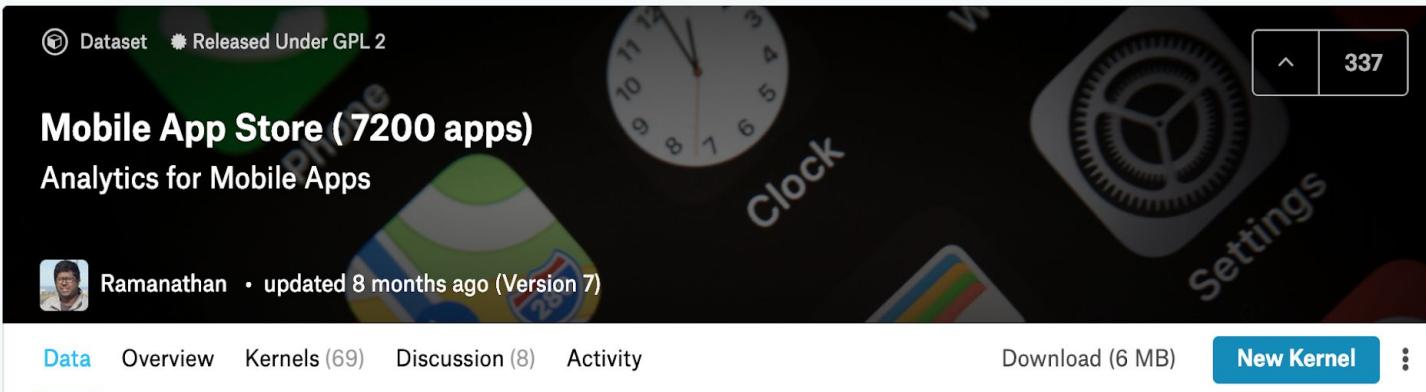


Dataset Released Under GPL 2

Mobile App Store (7200 apps)

Analytics for Mobile Apps

 Ramanathan • updated 8 months ago (Version 7)



Data Overview Kernels (69) Discussion (8) Activity Download (6 MB) New Kernel 

Data (6 MB)		
Data Sources		About this file
Columns		
AppleStore.csv	7198 x 17	Apple Store
appleStore_descriptio...	7197 x 4	<pre>-- # user_rating # user_rating_ver ▲ ver ▲ cont_rating ▲ prime_genre # sup_devices.num # ipadSc_urls.num # lang.num</pre>

	track_name	price	currency	rating_count_tot	user_rating
0	Facebook	0.0	USD	2974676	3.5
1	Instagram	0.0	USD	2161558	4.5
2	Clash of Clans	0.0	USD	2130805	4.5
3	Temple Run	0.0	USD	1724546	4.5
4	Pandora - Music & Radio	0.0	USD	1126879	4.0



AppleStore.csv

Data Science Interview Questions



<<Optional Non-Graded>>

1. What do you think makes a good data scientist?
2. Give a few examples of “best practices” in data science.
3. Can you outline the steps in a data science project?
4. What data would you love to acquire if there were no limitations?
5. Tell me a compelling story about data that you have analyzed.