
Uma experiência em Ciência de Dados aplicada a Projetos Sociais

— Laura de Oliveira F. Moraes —
PESC/COPPE/UFRJ

Data Science For Social Good

Summer Fellowship



Data Science For Social Good Europe

Summer Fellowship 2018

CASCAIS

NOVA
NOVA SCHOOL OF
BUSINESS & ECONOMICS



“The goal for the program was to find people who are interested in data and analytics and want to use those skills to help society”

-- Rayid Ghani, Director of Data Science for Social Good Fellowship

https://www.uchicago.edu/features/tackling_citys_challenges_with_data/

Data Science For Social Good

Summer Fellowship



Data Science For Social Good Europe

Summer Fellowship 2018

CASCAIS

N^{OVA} NOVA SCHOOL OF
BUSINESS & ECONOMICS



“The goal for the program was to find people who are interested in data and analytics and want to use those skills to help society”

-- Rayid Ghani, Director of Data Science for Social Good Fellowship

https://www.uchicago.edu/features/tackling_city_challenges_with_data/

Educação

Saúde

Meio-ambiente

Justiça Criminal

Transporte

Segurança Pública

Serviços Sociais

Desenvolvimento Econ.

Rayid Ghani



- Cientista de Dados Chefe da campanha do Obama em 2012
- 10 anos como Pesquisador na Accenture

Rayid Ghani



~ **40** estudantes/ano
~ **12** semanas
~ **9** mentores
~ **12** projetos/ano



→ Cientista de Dados Chefe da campanha do Obama em 2012

→ 10 anos como Pesquisador na Accenture





Data Science For Social Good
Summer Fellowship

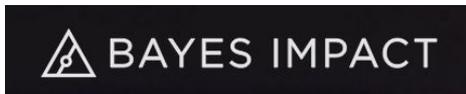


UNIVERSITY of WASHINGTON



IBM Social Good Fellowship

SoGood 2016



**Volunteerism for the Data Generation:
Using Data Science Superpowers for Social Good**



Partnership for Social Good



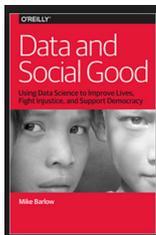
Statistics for Social Good

We're a group of Stanford students, researchers, and faculty exploring the potential to promote social good through effective data analysis.

#Data4Good



DO GOOD
DATA 2015



Data Science For Social Good Europe
Summer Fellowship 2018



Problema importante

Problemas sociais crônicos com ações de impacto.

Dados

Dados estão disponíveis e acessíveis durante o projeto.

+

Parceiro comprometido

Possui a experiência diária e é especialista nos dados.

Problema desafiador e solucionável



Pranjali Bajaj

Quantitative Methods in the Social Sciences
Columbia University



Andrew Bell

Mathematical Sciences
Clemson University



Ruqian Chen

Mathematics
University of Washington



Bruno Del Papa

Physics/Neuroscience
Frankfurt Institute for Advanced Studies and Max Planck Institute for Brain Research, Goethe University Frankfurt



Anne Driscoll

Statistical Science
Duke University



Jordan Kupersmith

Masters of Information and Data Science
University of California at Berkeley



Indu Manickam

Electrical & Computer Engineering
Rice University



Kaushik Mohan

Applied Statistics for Social Science Research
New York University



Laura Moraes

Computer and Systems Engineering Program
Universidade Federal do Rio de Janeiro



Harsh Nisar

Information Communication Technology ICT
Dhirubhai Ambani Institute Information Communication Technology (DA-IICT)



Alexander Rich

Psychology, New York University



Mélisande Teng

Applied Mathematics,
CentraleSupélec / ENS Paris-Saclay



Can Udomcharoenchaikit

Computer Engineering
Chulalongkorn University



Orsolya Vásárhelyi

Center for Network Science
Central European University

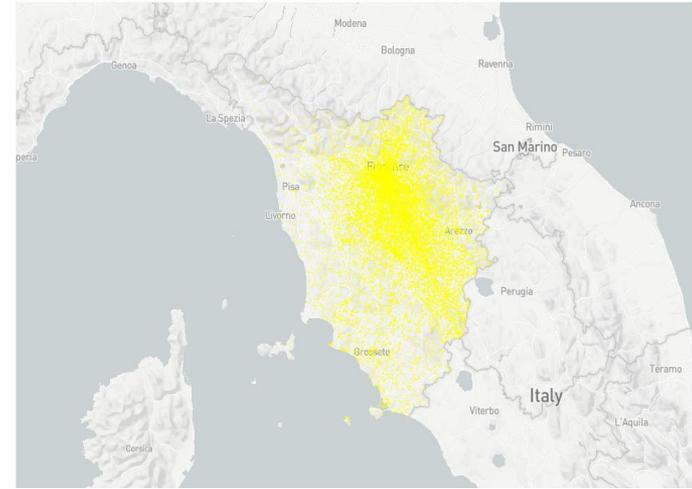


Yanbing Wang

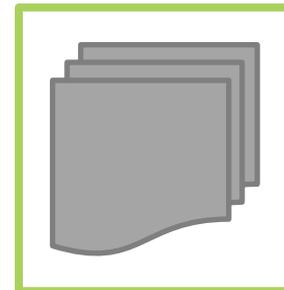
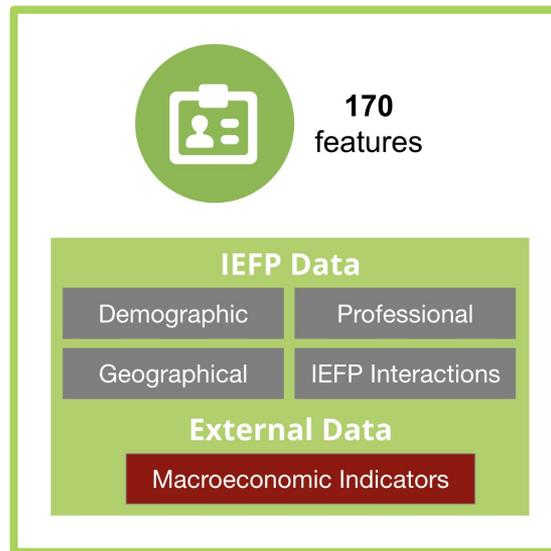
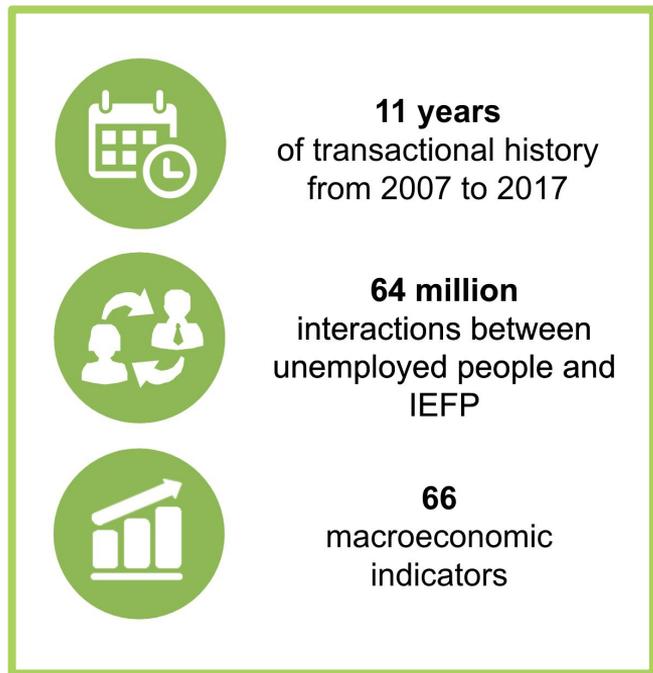
Agricultural Economics
Purdue University

15
estudantes

4
projetos



Dados



Data Story





Fizemos um *brainstorm* dos atributos...

FEAT. ENG.

100% i

Demographic

- OK - Genders
- OK - Age: ① continuous, ② categorical
- OK - Civil Status
- OK - Education
 - Professional
 - Educational
- OK - Dependents
- OK - Nationality
- * - Disability
 - boolean
 - types of disability
 - ? severity
- ✓ - Unemployment Benefit
- OK - RSI: Universal Benefits

Transactional

- Reinscription
- Reason of enrollment
- Whether person was kicked out
 - Reason for getting kicked out
- ✓ - How many times a person entered the system
 - Motivation proxies
 - Attendance rate of X
 - Kicked out from intervention
 - Medical Wish also
- * - How many times person was out of system
 - # of interventions (past year / 3 yr / 15 yr) ok
 - it's a description

Professional

- OK - Tempo Practica
- OK - All experience
- UNCP: Last Job
 - Tabo: % NULLS?
- ✓ - CAE: Too many industries
 - Need to collapse / cut off
- OK - Employment Status in that month
 - Months unemployed
 - was LTV in the past all time
 - D/ CPP → Δ CPA → Professions they aspire
 - CNP
 - ANTERIOR
- OK - Past-time / Full time
- * - has-prof-cent

Geographic

- Where they live (3 levels)
- Where they are looking for jobs
- Make X3 Geographic features for 3 levels
- Inland / Coastal
- Urban / Rural
- Movement from Rural → Urban
- How many FregZ

Economic

- ✓ - recession Y/N
- Municipality
 - purchasing power
 - salary
 - NUTS2
 - labour force particip. rate
 - ...
- Unemployment rate (low granular/freq. possible)
- ruling party (L/R)
- PRANTAL?!
- VAB

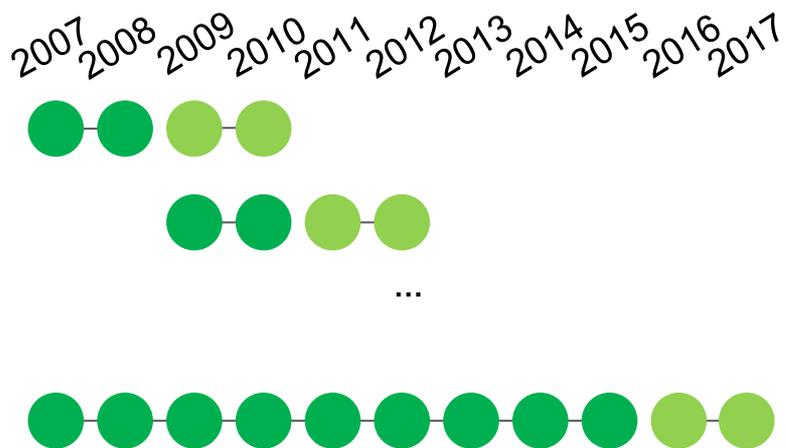
functions:

- min \sum
- max \sum
- aggre Δ
- avg

TODO:

- import utilities de inscrição
- % NULLS
- scrape economic data? for someone else
- Database size? (load X+Y)

Validação cruzada temporal





Classificação / Data de referência do modelo: 30/04/2017 / 2753993

Risco:

0

Perfil 3

Região:

Área

Metropolitana

de

Lisboa

↑ Fatores de aumento do risco

- ^ Idade: 55.0
- ^ Foi DLD no passado: Sim
- ^ Apresentações nos últimos 5 anos: 4.0
- ^ Intervenções nos últimos 5 anos: 22.0
- ^ Proporção de profissionais socialmente mais valorizados: 60.0
- ^ Indicador de confiança dos consumidores: -7.0

↓ Fatores de diminuição de risco

- ∨ Desempregado: Não
- ∨ Número de registos com o IEFP: 8.0
- ∨ Intervenções nos últimos 3 anos: 19.0
- ∨ Total de meses desempregado: 8.0

 Estatísticas1 apresentações
entre 30/04/2016 e
31/12/20179 intervenções entre
30/04/2016 e
31/12/2017 Histórico de Risco



Análise por Região

Região

Norte

Capacidade
para
Perfil 1

5000

Capacidade
para
Perfil 2

20000

Ver estatísticas

Valores limites mínimos de risco

Perfil 1: 49

Perfil 2: 33

Usar esses valores



“So much of the computing agenda today is motivated by the problems that the big Internet companies face,” he observed. “But if all these kids are doing that, they are not working on the things that really matter.” -- Rayid

Obrigada!

Additional Material

DSSG 2018: General Overview of the Summer

Data Science For Social Good Europe

Summer Fellowship 2018

CASCAIS

N.O.V.A.

UNIVERSITY OF
LIVORNO

UNIVERSITY OF
CHICAGO

Week 1	Week 2	Week 3	Week 4	Week 5	Week 6	Week 7	Week 8	Week 9	Week 10	Week 11	Week 12
May 28 th - June 1 st	June 4 th - 8 th	June 11 th - 15 th	June 18 th - 22 nd	June 25 th - 29 th	July 2 nd - 6 th	July 9 th - 13 th	July 16 th - 20 th	July 23 rd - 27 th	July 30 th - August 3 rd	August 6 th - 10 th	August 13 th - 17 th
Orientation	Project Scoping & Data Discovery	Project Scoping & Data Discovery	Data Discovery	Data Science Pipeline Development	Data Science Pipeline Development	Data Science Pipeline Development	Data Science Pipeline Development	Data Science Pipeline Development	Data Science Pipeline Development	Handover to Partner & Transition Planning	Presentations & Transition
Make New Friends <small>(both human and technical)</small>	Data Audit & Exploration	Finalize Project Scope & Data Stories	Pipeline & Technical Plan	Iteration 1 Build End-to-End Code Pipeline Focus on end-to-end structure	Iteration 1 Build End-to-End Code Pipeline Focus on end-to-end structure	Iteration 2 End-to-End Code Pipeline Focus on feature development	Iteration 2 End-to-End Code Pipeline Focus on feature development	Iteration 3 End-to-End Code Pipeline Focus on evaluation, results and initial front-end demo	Iteration 3 End-to-End Code Pipeline Focus on evaluation, results and initial front-end demo	Final Deliverable Acceptance & Handover to Partner	Final Event Presentations & Handover to Partner
	Descriptive Statistics	Project Scope and Data Stories	Outline of Pipeline	Initial Pipeline First Version	Early Results	Prioritized Feature List	Interpretable Model (by Partner)	Results: Across models features & metrics	Actionable Deliverable (e.g. UI Demo)	Final Report	Poster
	Questions on Data for the Partner	Project Charter approved by partner with Project Plan	Mockups of actionable deliverable			Mid-Summer Partner Visit	Reproducible Data Science Pipeline			Signed off final deliverables	Presentation
	Data Dictionary	Presentation for Partner with initial insights and Data Stories	Technical Plan							Code developed during the project (e.g. clean Github repo)	