



# CIÊNCIA DE DADOS E SUA APLICAÇÃO EM SAÚDE

**Prof. Marcos Barreto**

AtylmoLab / LaSiD  
Departamento de Ciência da Computação  
Instituto de Matemática e Estatística  
Universidade Federal da Bahia (UFBA)

Salvador, 22 de março de 2018

# Popularização da Ciência de Dados



=



# Popularização da Ciência de Dados





# Ciência de dados

[ocultar]

Origem: Wikipédia, a enciclopédia livre.

**Ciência de dados** (em **inglês**: *data science*) é uma área interdisciplinar voltada para o estudo e a análise de dados, estruturados ou não, que visa a extração de conhecimento ou *insights* para possíveis tomadas de decisão, de maneira similar à **mineração de dados**. Ciência de dados alia **big data** e **machine learning**, além de técnicas de outras áreas interdisciplinares como estatística, economia, engenharia e outros subcampos da **computação** como: **banco de dados** e **análise de agrupamentos** (*cluster analysis*). A ciência de dados é um campo que já existe há 30 anos, porém ganhou mais destaque nos últimos anos devido a alguns fatores como: o surgimento e popularização do Big Data e o desenvolvimento de áreas como o machine learning. A ciência de dados pode, por exemplo, transformar essa grande quantidade de dados brutos em *insights* de negócios, e com isso, auxiliar empresas em tomadas de decisões para atingir melhores resultados.<sup>[1]</sup>

LIFE  
=  
WHAT WE THINK WE KNOW  
+  
RANDOM UNKNOWN

REAL DATA  
=  
A MODEL EXPLAINING DATA  
+  
RANDOM ERROR

Is Data Science a science?

<https://towardsdatascience.com/data-science-and-ai-for-business-data-analysts-64f28a5d7ff2>

# História (base Estatística)

Data analysis, and the parts of statistics which adhere to it, must...take on the characteristics of science rather than those of mathematics... data analysis is intrinsically an empirical science...

"[Data is] a representation of facts or ideas in a formalized manner capable of being communicated or manipulated by some process."

"...more emphasis needed to be placed on using data to suggest hypotheses to test and that Exploratory Data Analysis and Confirmatory Data Analysis should proceed side by side."



1962 (John Turkey)  
*The future of data analysis*

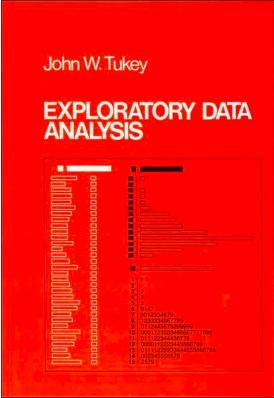
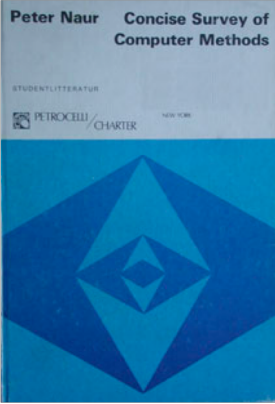
1968 (IFIP)  
Datalogy

1974 (Peter Naur)  
*Concise Survey of Computer Methods*

1977 (John Turkey)  
*Exploratory data analysis*

1977 (ISI)  
*International Association for Statistical Computing*

Datalogy, the science of data and of data processes and its place in education...



"...link traditional statistical methodology, modern computer technology, and the knowledge of domain experts in order to convert data into information and knowledge."

# História (base Computação + BI)

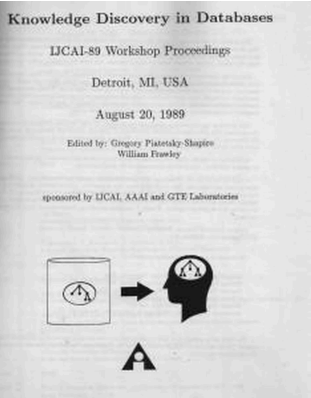
“Companies are collecting mountains of information about you, crunching it to predict how likely you are to buy a product, and using that knowledge to craft a marketing message precisely calibrated to get you to do so...”

The classification societies have variously used the terms data analysis, data mining, and data science in their publications.

“Historically, the notion of finding useful patterns in data has been given a variety of names, including data mining, knowledge extraction, information discovery, information harvesting, data archeology, and data pattern processing...”



1989 (Piatetsky-Shapiro)  
1st KDD workshop



1994 (Business Week)

Business  
**Database Marketing**  
Jonathan Berry  
5 de setembro de 1994 01:00 BRT

1996 (IFCS)  
International Federation of Classification Societies

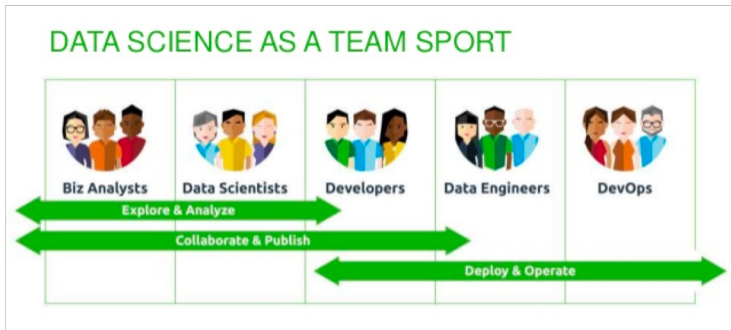
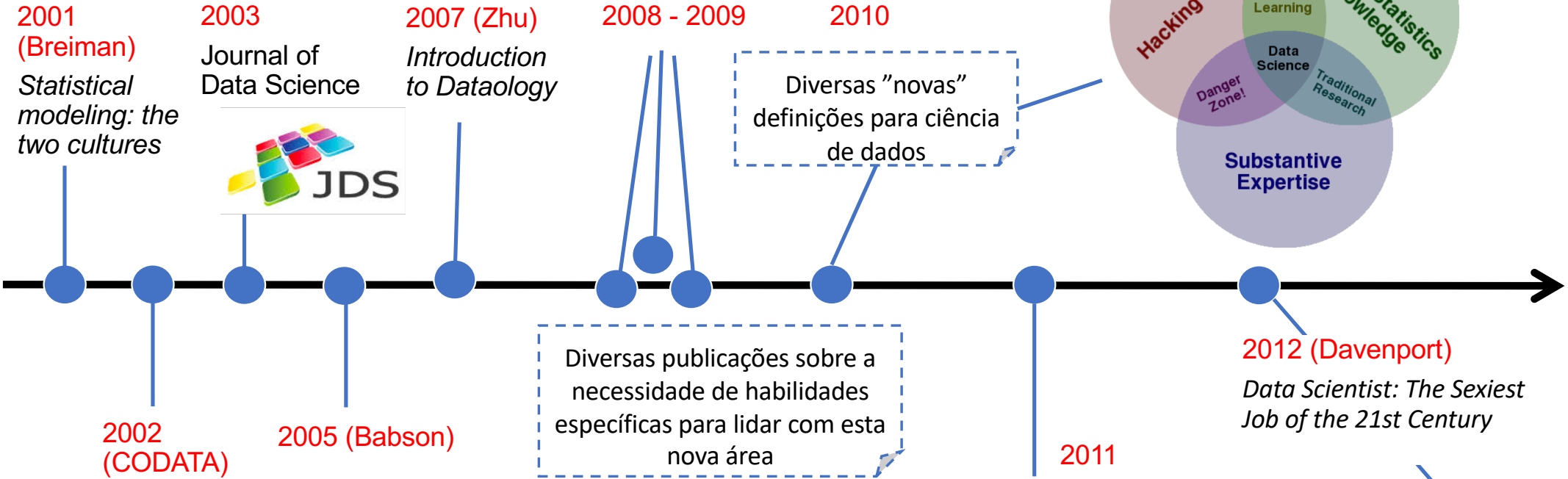


1996 (Fayyad - Shapiro)  
From data mining to KDD paper

2001 (William Cleveland)  
Data Science Action Plan

It is a plan “to enlarge the major areas of technical work of the field of statistics.”

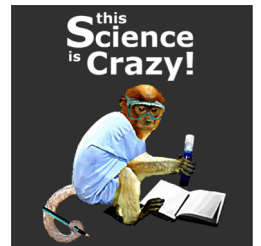
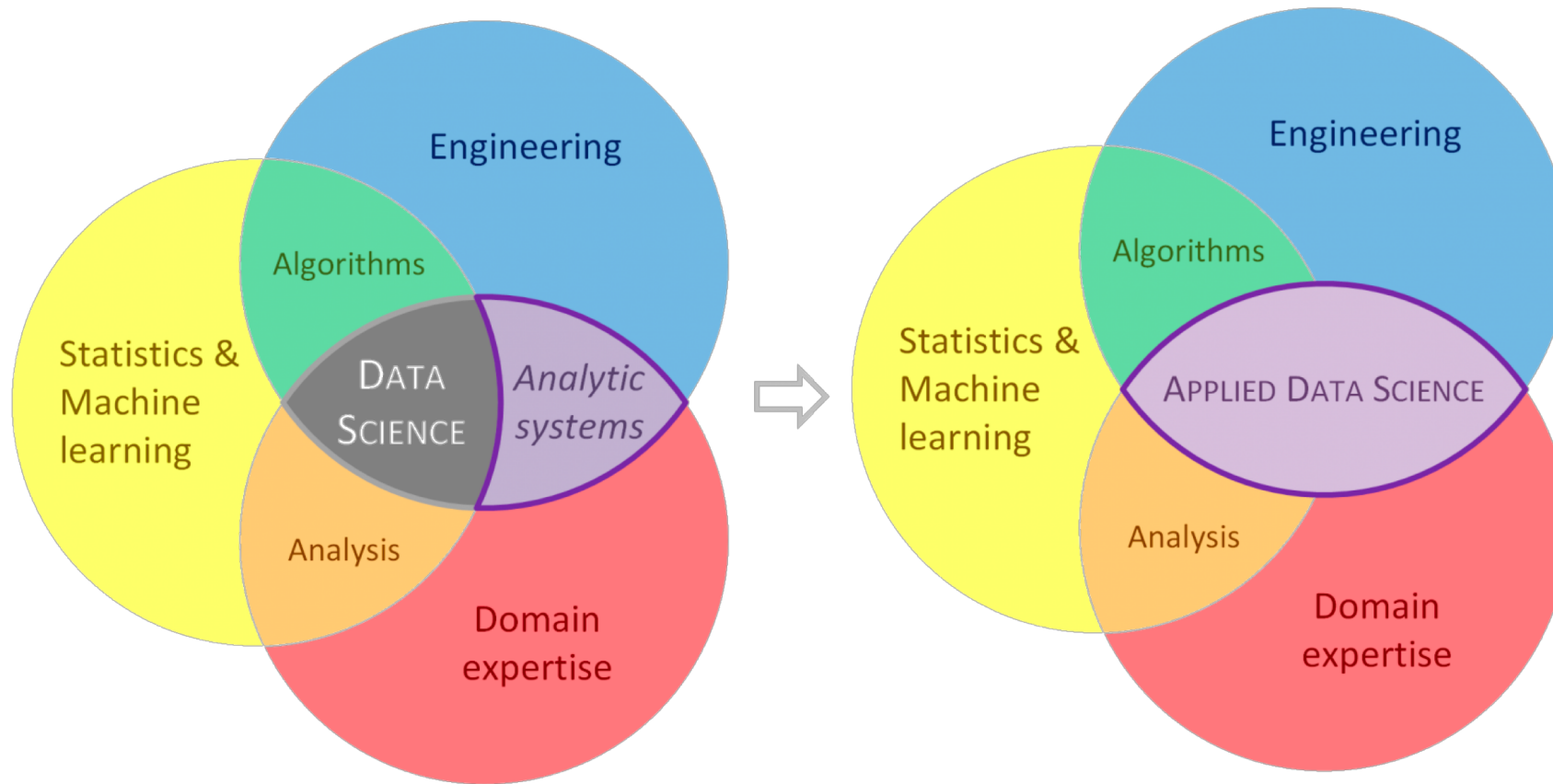
# História (base Computação + BI)



<https://www.forbes.com/sites/gilpress/2013/05/28/a-very-short-history-of-data-science/#78699be255cf>



# Definindo Ciência de Dados

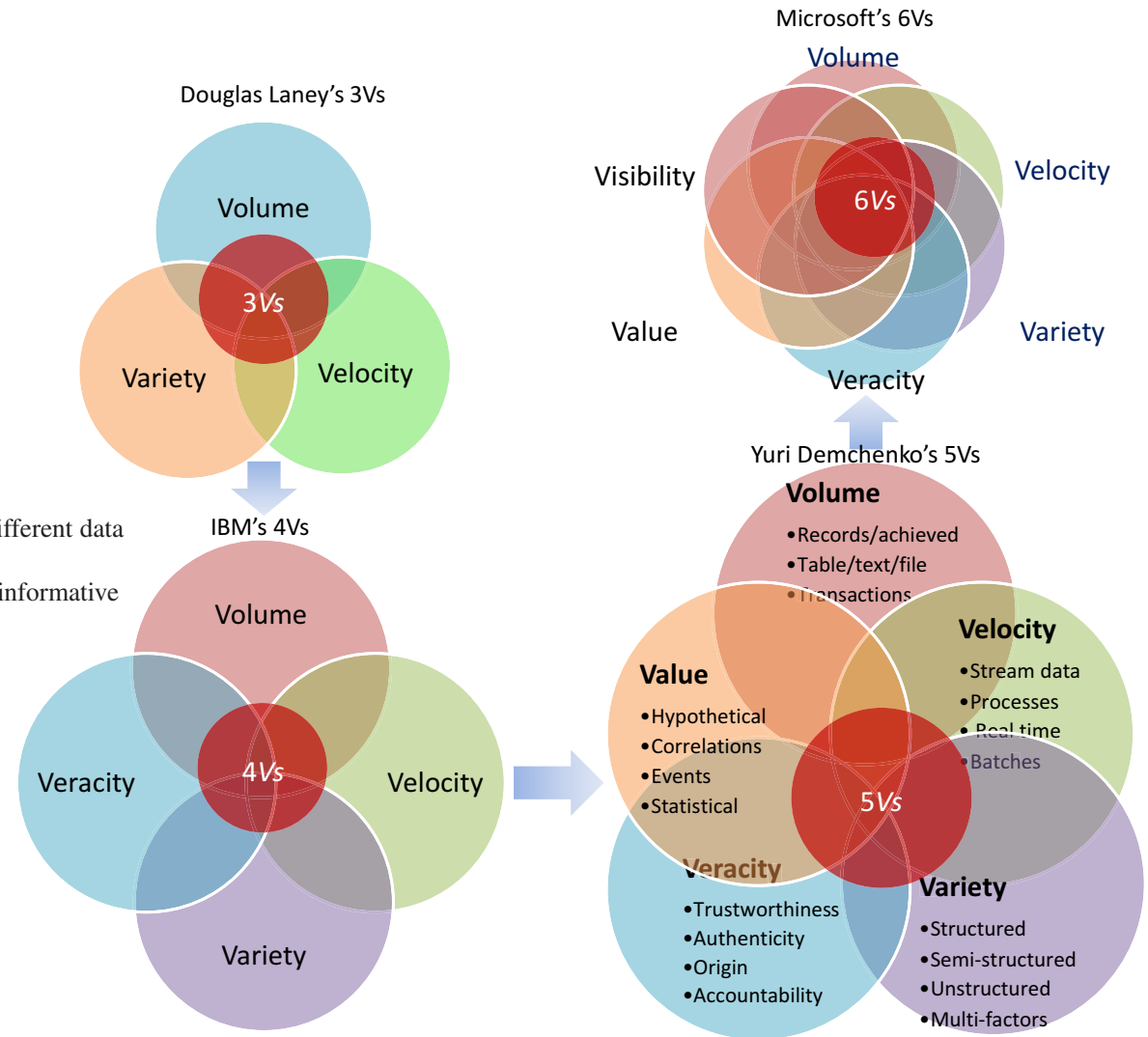
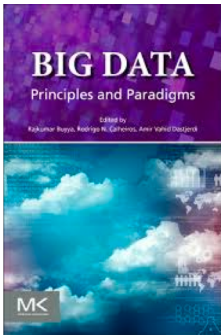


# BIG DATA



# Big Data

1. Volume stands for scale of data
2. Velocity denotes the analysis of streaming data
3. Variety indicates different forms of data
4. Veracity focuses on trustworthiness of data sources
5. Variability refers to the complexity of data set. In comparison with “Variety” (or different data format), it means the number of variables in data sets
6. Visibility emphasizes that you need to have a full picture of data in order to make informative decision

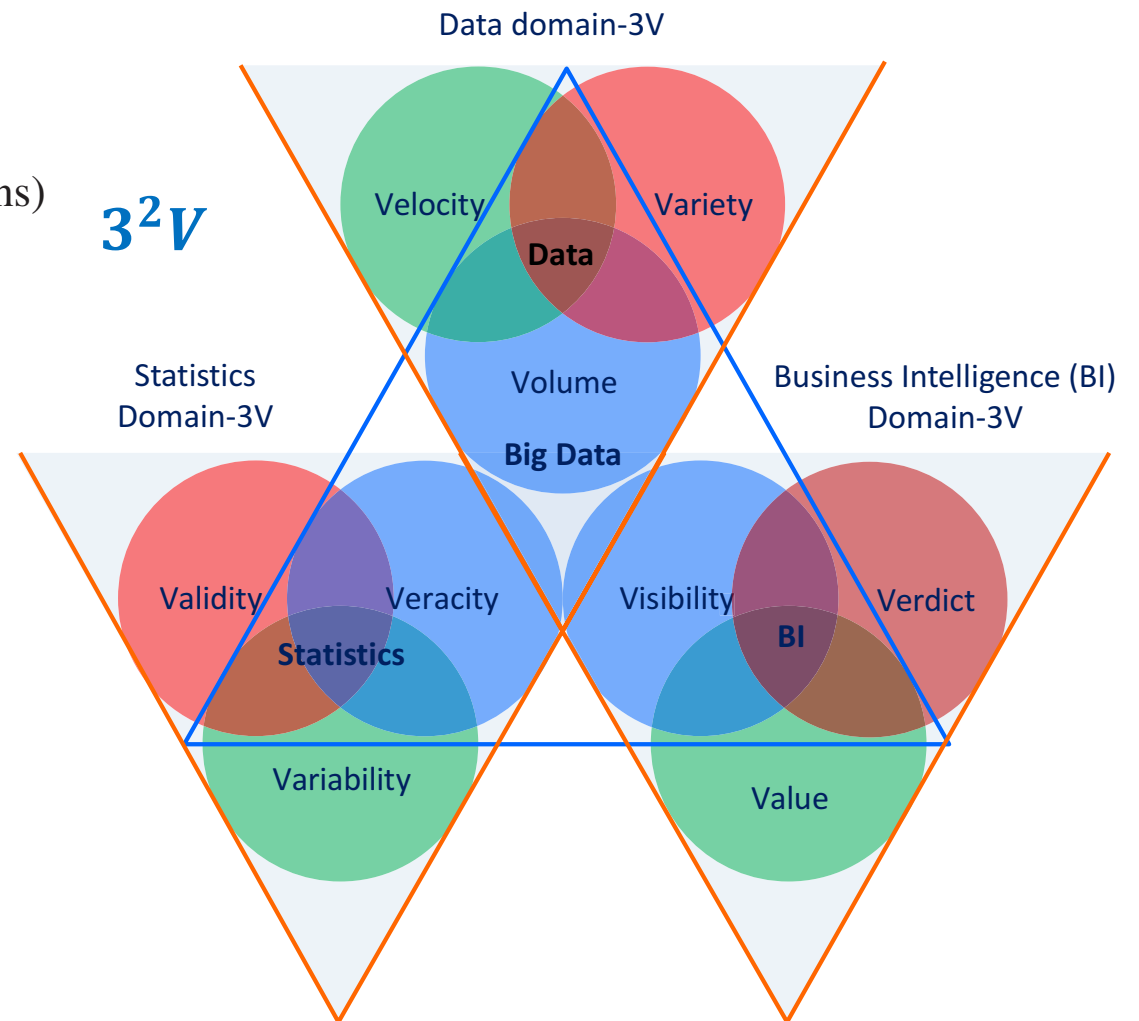


**FIG. 3**

From 3Vs, 4Vs, 5Vs, and 6Vs big data definition.

# Big Data

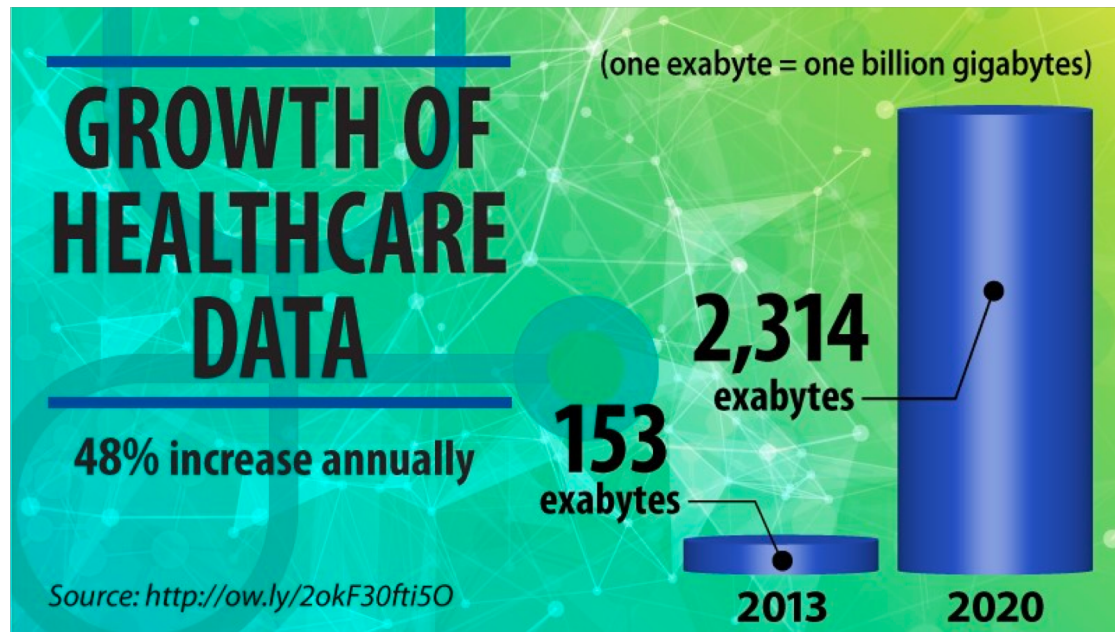
- Data domain (searching for patterns)
- Business intelligence domain (making predictions)
- Statistical domain (making assumptions)



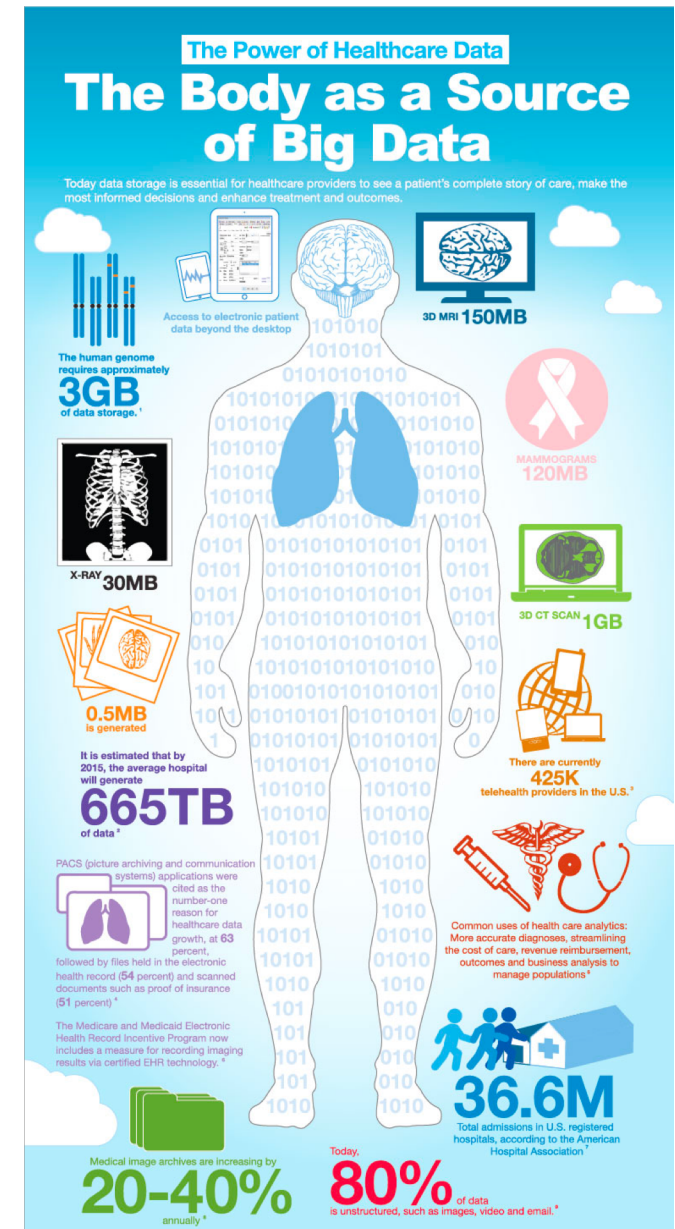
**FIG. 5**

3<sup>2</sup>Vs Venn diagrams in hierarchical model.

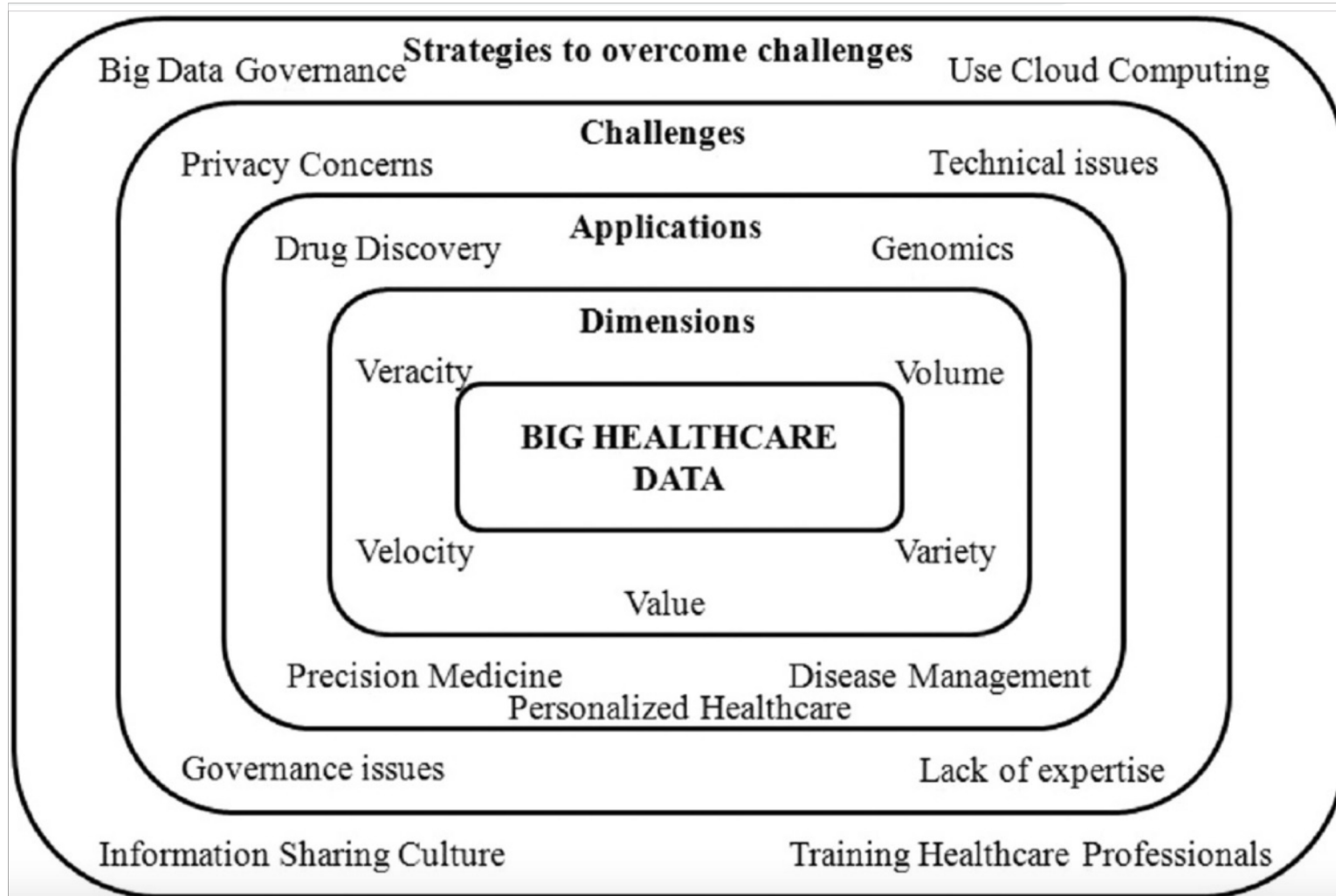
# Big Data em Saúde



<http://sites.ieee.org/futuredirections/2018/05/18/the-future-of-health-care-is-tied-to-big-data/>

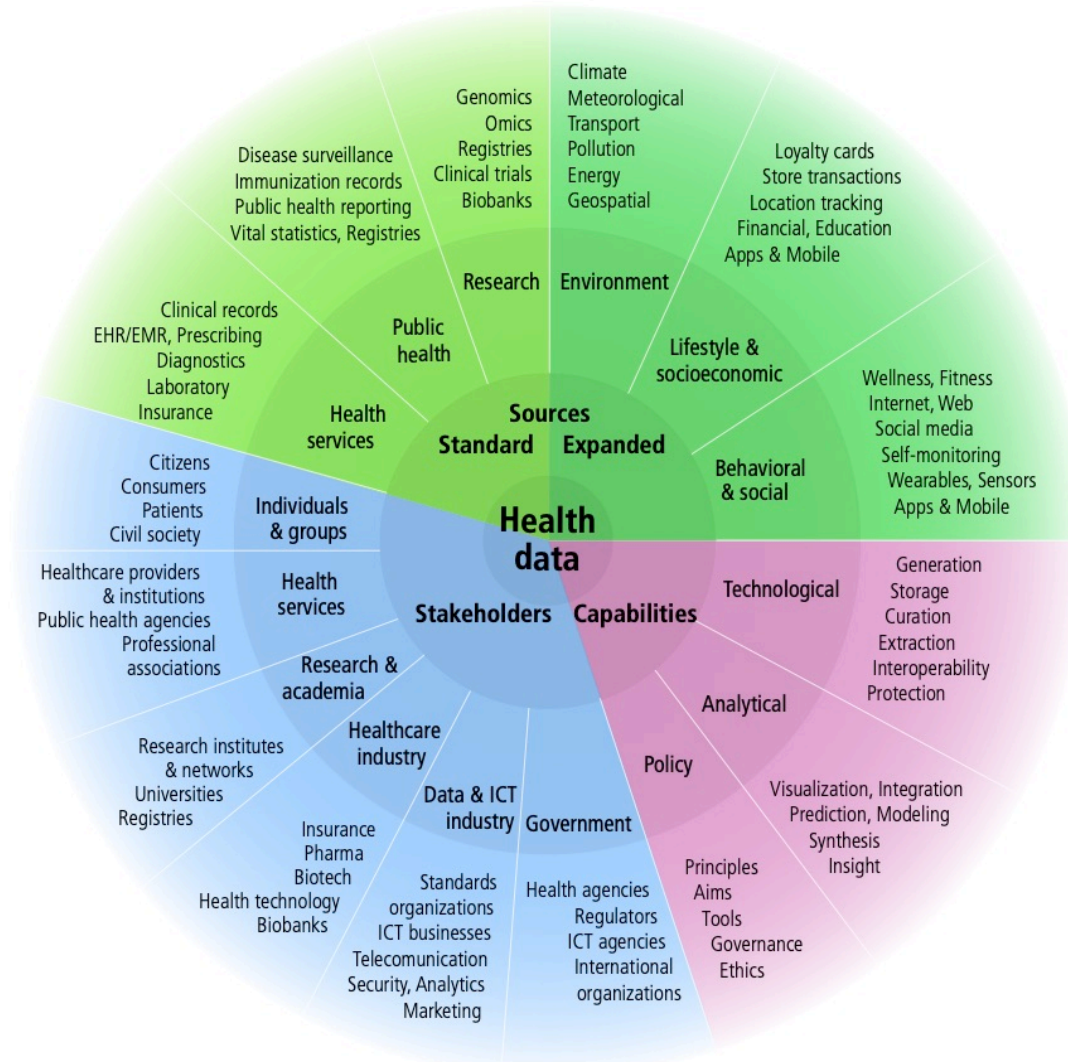


# Big Data em Saúde



# DADOS DE SAÚDE

## Evolving health data ecosystem



# Dados de saúde

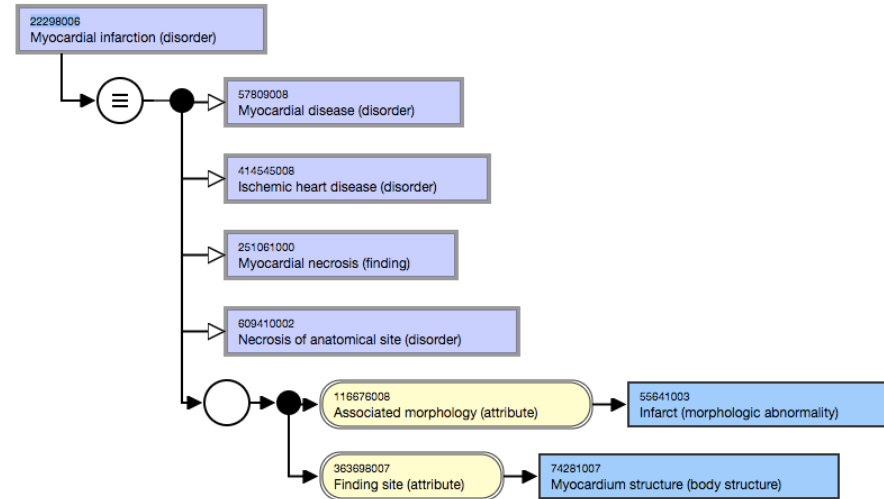
1738.00 663v100,  
J45, J44, 66YL.11,  
G20.00, 662O.00,  
1738.00 1682.00,  
I50, 06,  
116676008, I21.00

## Structured

Orientation plus  
Stabilized plus

Make rs174546(C;C)
Make rs174546(C;T)
Make rs174546(T;T)

Reference GRCh38 38.1/141  
Chromosome 11  
Position 61802358  
Gene FADS1  
is a snp  
is mentioned by  
dbSNP rs174546



- Heart failure**
- Excl.:** complicating:
- abortion or ectopic or molar pregnancy ([O00-O07](#), [O08.8](#))
  - obstetric surgery and procedures ([O75.4](#))
- due to hypertension ([I11.0](#))
- with renal disease ([I13.-](#))
- following cardiac surgery or due to presence of cardiac prosthesis ([I97.1](#))
- neonatal cardiac failure ([P29.0](#))
- 3.0 Congestive heart failure**  
Congestive heart disease  
Right ventricular failure (secondary to left heart failure)
- 3.1 Left ventricular failure**  
Cardiac asthma  
Left heart failure  
Oedema of lung | with mention of heart disease NOS or heart failure  
Pulmonary oedema
- 3.9 Heart failure, unspecified**  
Cardiac, heart or myocardial failure NOS



# Dados de saúde

1738.00 663v100,  
J45, J44, 66YL.11,  
G20.00, 662O.00,  
1738.00 1682.00,  
I50, 06,  
116676008, I21.00

Structured

220, 110, 0.002, 1,  
200, 3, 2, 2.1, 2.01,  
20, 1, 99092, 1.2,  
99, 123, 6, 23.2,  
878, 9901, 11,  
203.1

Semi-structured

Local Code

OBX||NM|123^RBC^MyHosp|26453-1^Erythrocytes [# /volume] in Blood^LN|4.82|10\*6/uL

LOINC Code

# Dados de saúde

1738.00 663v100,  
J45, J44, 66YL.11,  
G20.00, 662O.00,  
1738.00 1682.00,  
I50, 06,  
116676008, I21.00

Structured

220, 110, 0.002, 1,  
200, 3, 2, 2.1, 2.01,  
20, 1, 99092, 1.2,  
99, 123, 6, 23.2,  
878, 9901, 11,  
203.1

Semi-structured

~~~~~  
~~~~~  
~~~~~  
~~~~~  
~~~~~  
~~~~~

Unstructured

76 yo man with h/o HTN, DM, and sleep apnea who presented to the ED complaining of chest pain. He states that the pain began the day before and consisted of a sharp pain that lasted around 30 seconds, followed by a dull pain that would last around 2 minutes. The pain was located over his left chest area somewhat near his shoulder. The onset of pain came while the patient was walking in his home. He did not sit and rest during the pain, but continued to do household chores. Later on in the afternoon he went to the gym where he walked 1 mile on the treadmill, rode the bike for 5 minutes, and swam in the pool. After returning from the gym he did some work out in the yard, cutting back some vines. He did not have any reoccurrences of chest pain while at the gym or later in the evening. The following morning (of his presentation to the ED) he noticed the pain as he was getting out of bed. Once again it was a dull pain, preceded by a short interval of a sharp pain. The patient did experience some tingling in his right arm after the pain ceased. He continued to have several episodes of the pain throughout the morning, so his daughter-in-law decided to take him to the ED around 12:30pm. The painful episodes did not increase in intensity or severity during this time.

# Dados de saúde

1738.00 663v100,  
J45, J44, 66YL.11,  
G20.00, 662O.00,  
1738.00 1682.00,  
I50, 06,  
116676008, I21.00

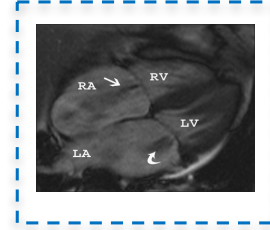
Structured

220, 110, 0.002, 1,  
200, 3, 2, 2.1, 2.01,  
20, 1, 99092, 1.2,  
99, 123, 6, 23.2,  
878, 9901, 11,  
203.1

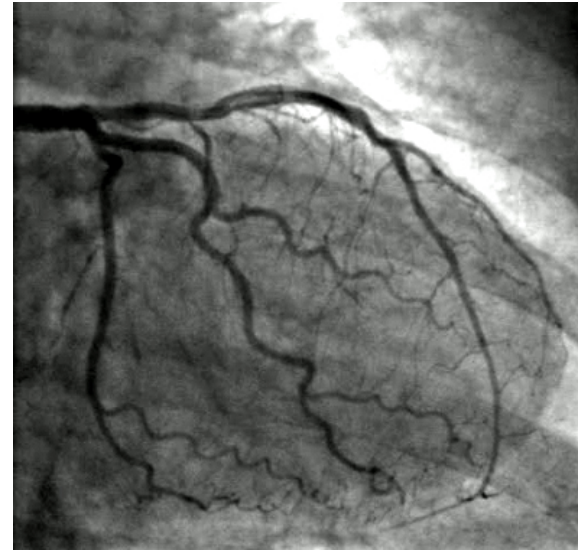
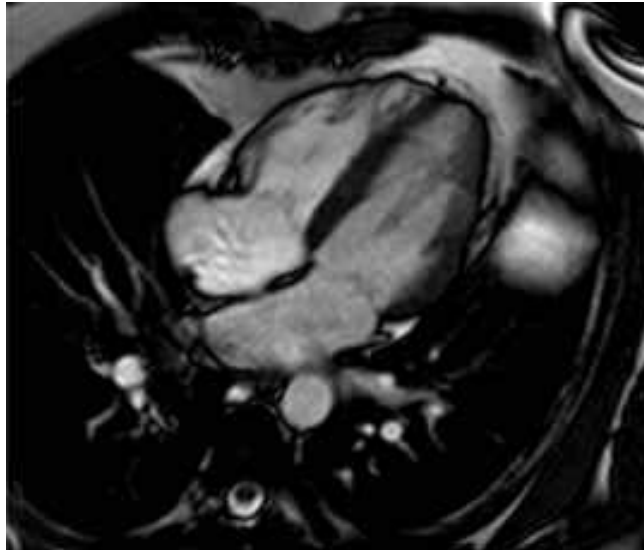
Semi-structured

~~~~~  
~~~~~  
~~~~~  
~~~~~  
~~~~~

Unstructured



Binary



# Dados de saúde

```
1738.00 663v100,  
J45, J44, 66YL.11,  
G20.00, 662O.00,  
1738.00 1682.00,  
I50, 06,  
116676008, I21.00
```

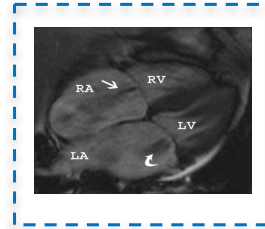
Structured

```
220, 110, 0.002, 1,  
200, 3, 2, 2.1, 2.01,  
20, 1, 99092, 1.2,  
99, 123, 6, 23.2,  
878, 9901, 11,  
203.1
```

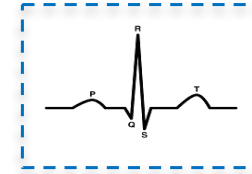
Semi-structured

```
~~~~~  
~~~~~  
~~~~~  
~~~~~  
~~~~~  
~~~~~
```

Unstructured



Binary



Streaming

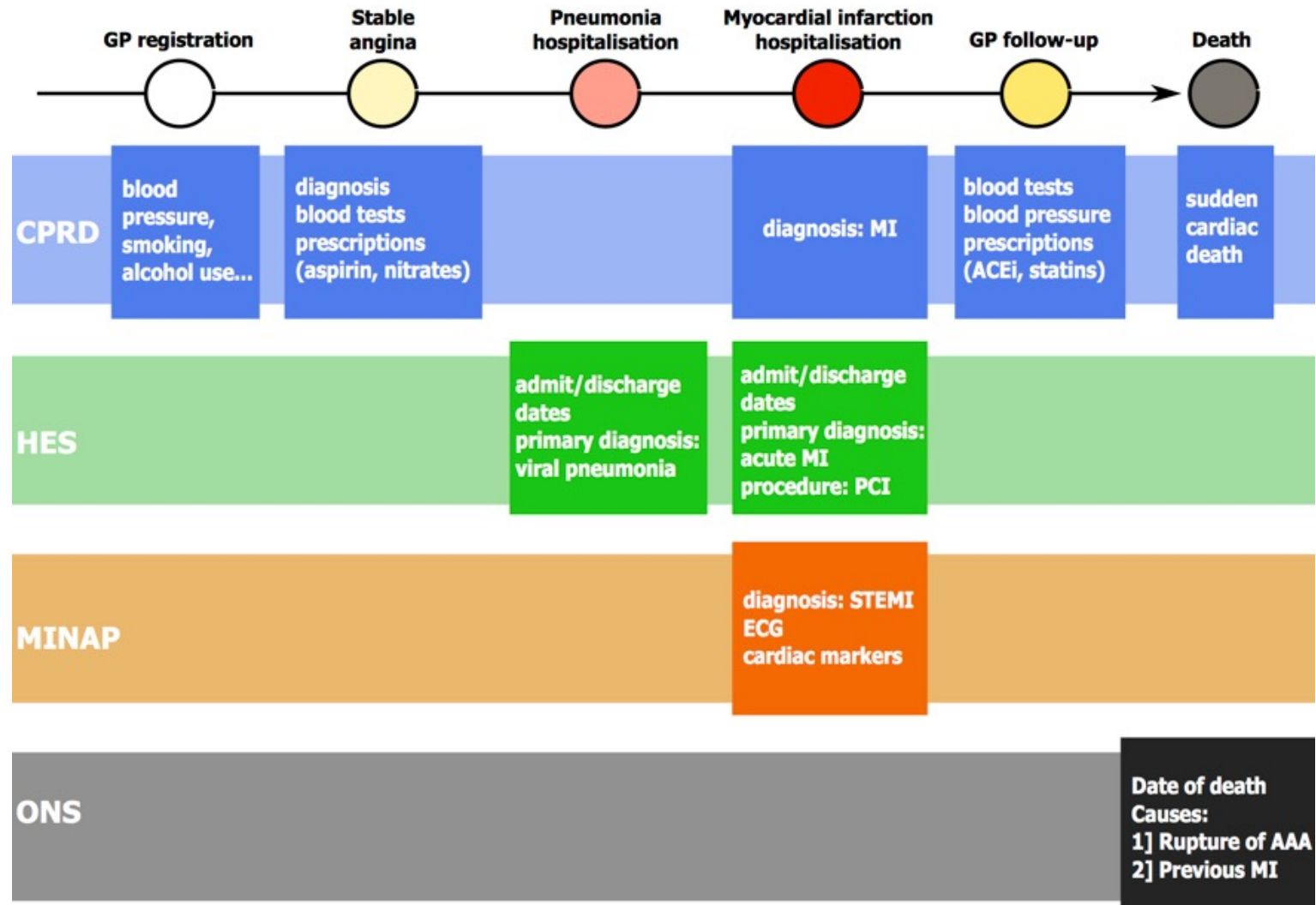


## APLICAÇÕES

- ✓ Estudos epidemiológicos de base populacional (coorte)
  - ✓ Análises clínicas
  - ✓ Descoberta de novos fármacos
  - ✓ Cura de doenças / resistência antimicrobiana
  - ✓ Medicina personalizada
  - ✓ Prevenção de visitas desnecessárias ao médico
- 
- ✓ Análise preditiva em saúde
  - ✓ Registros eletrônicos (*electronic health records – EHR*)
  - ✓ Monitoramento em tempo real
  - ✓ Sistemas de suporte à decisão
  - ✓ Medicina assistida por computação / *mobile health* (m-Health)

# Aplicações

# Population-based clinical epidemiology



# Aplicações

## Type-2 Diabetes and 12 CVDs

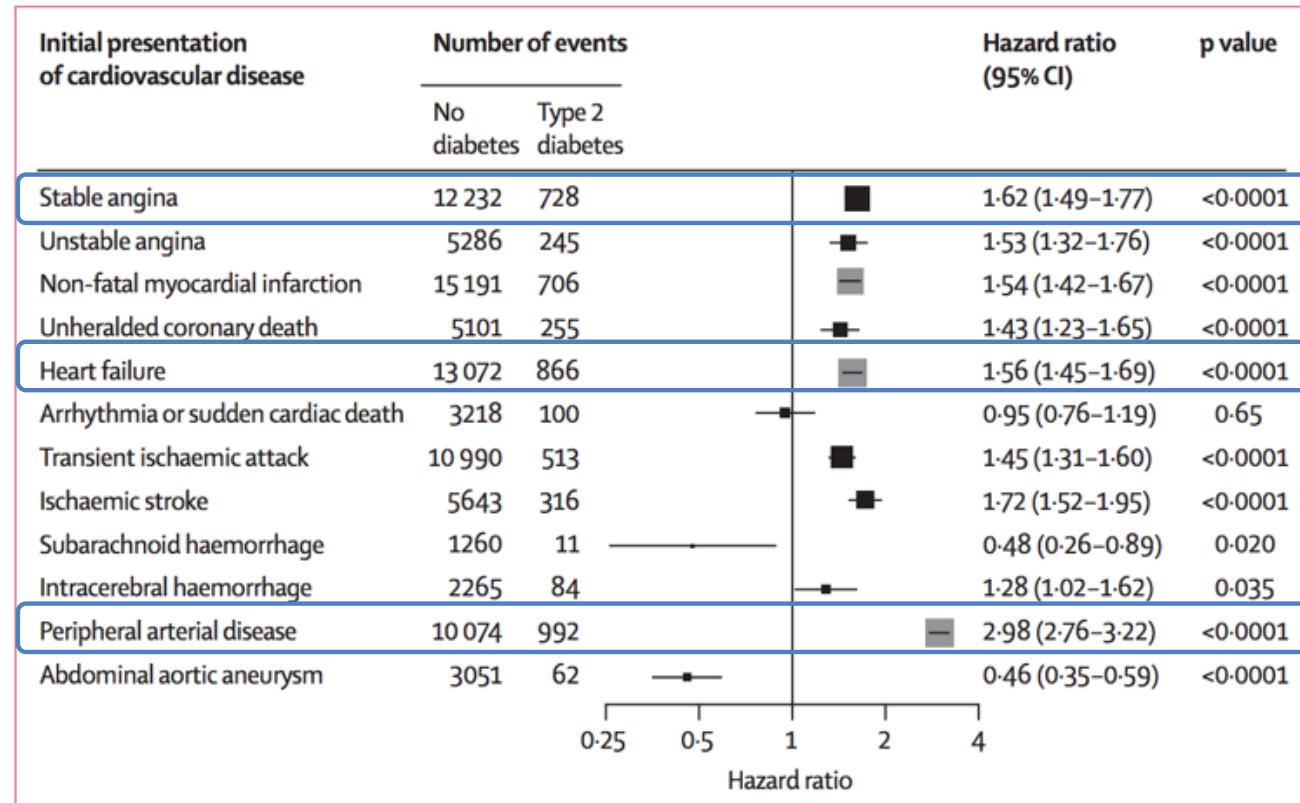


Figure 3: Association of type 2 diabetes with 12 cardiovascular diseases in patients aged  $\geq 30$  years

# Aplicações

# Clinical trials pipelines

**Problem:** A lot of medical care is educated guesses

**Opportunity:** Decisions based on what happened to people like you.

**My Patient**

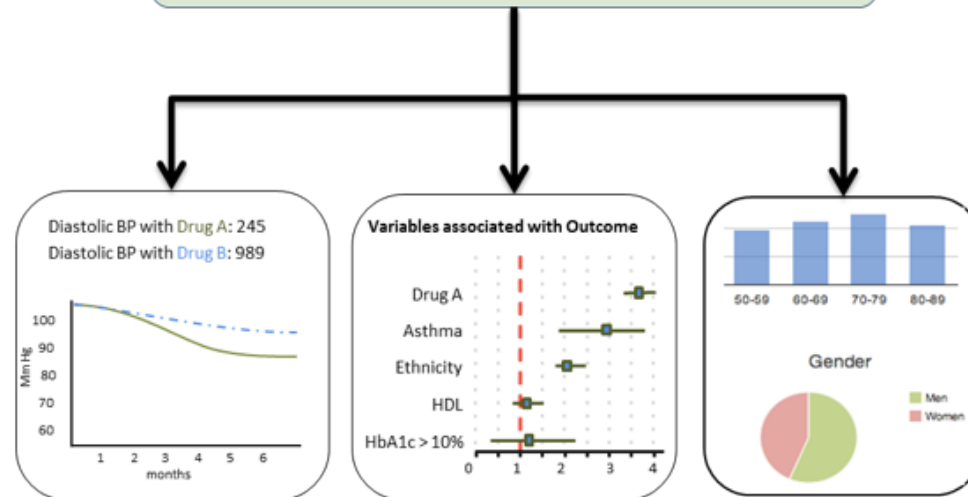
A 55 year old female of Vietnamese heritage with known asthma presents to her physician with new onset moderate hypertension

**Intervention**

antihypertensives

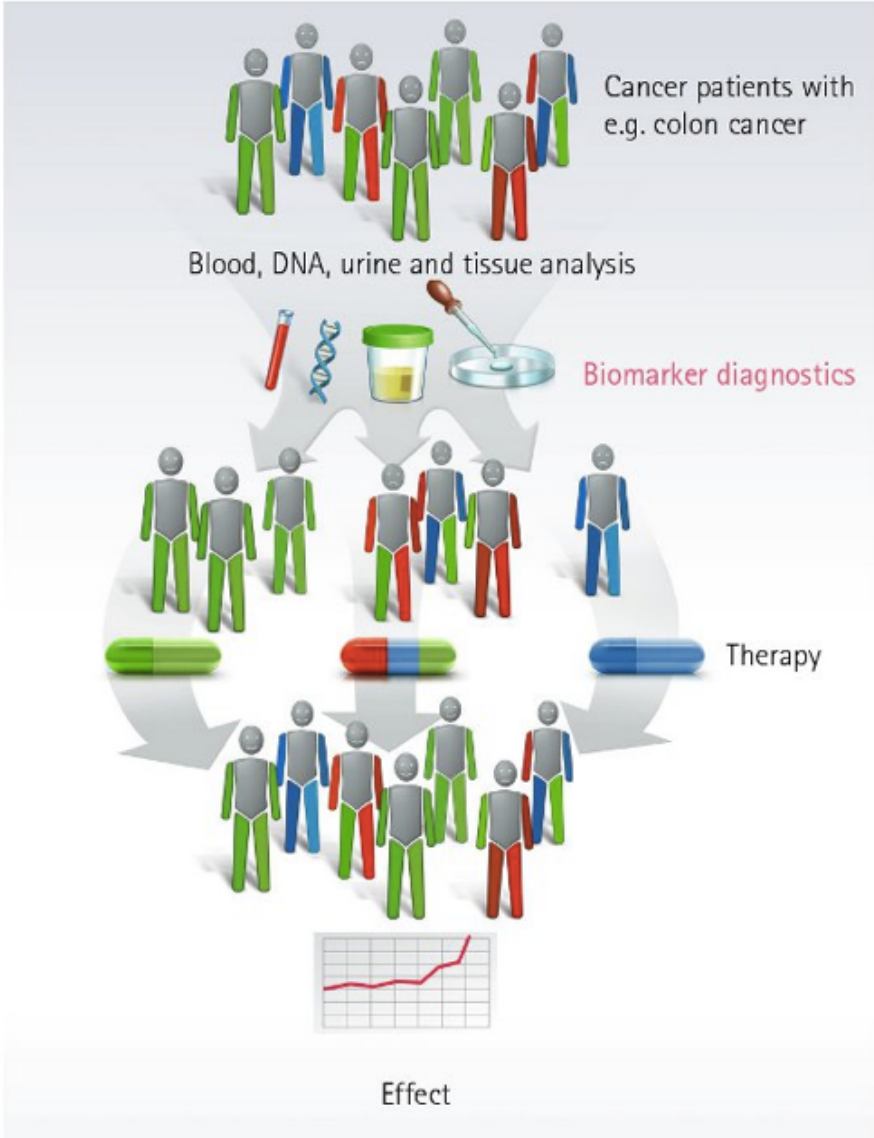
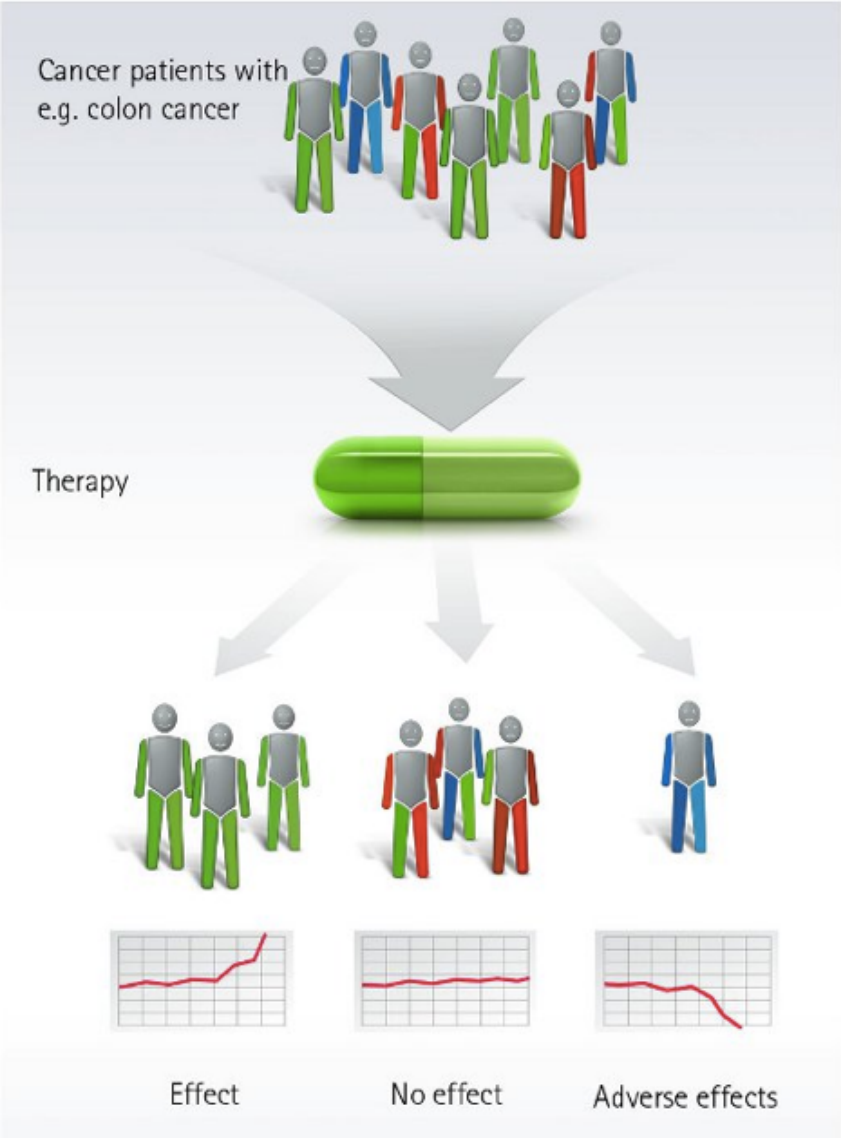
**Outcome**

Diastolic pressure < 90 mm Hg





# Aplicações

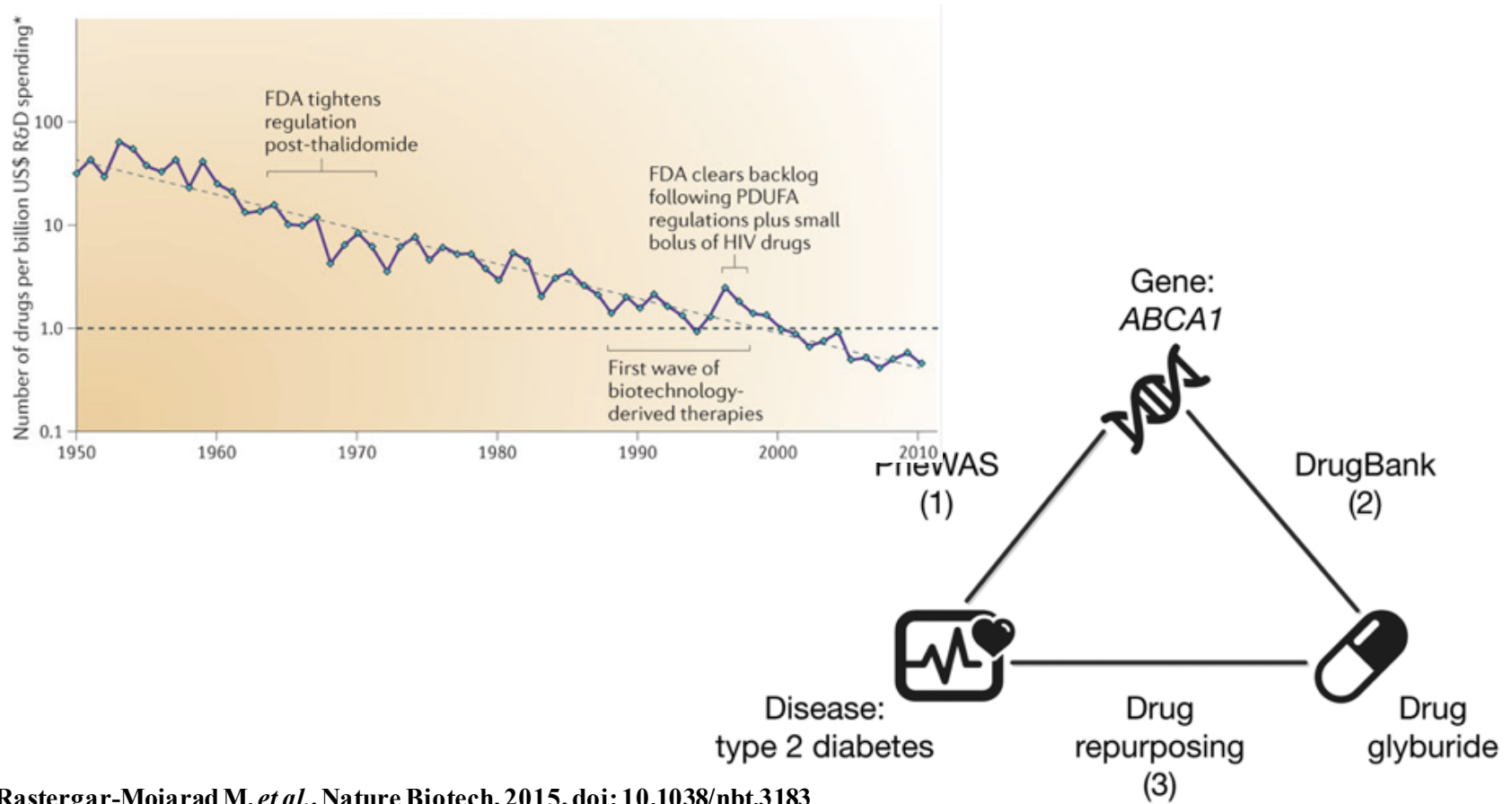


Precision medicine

# Aplicações

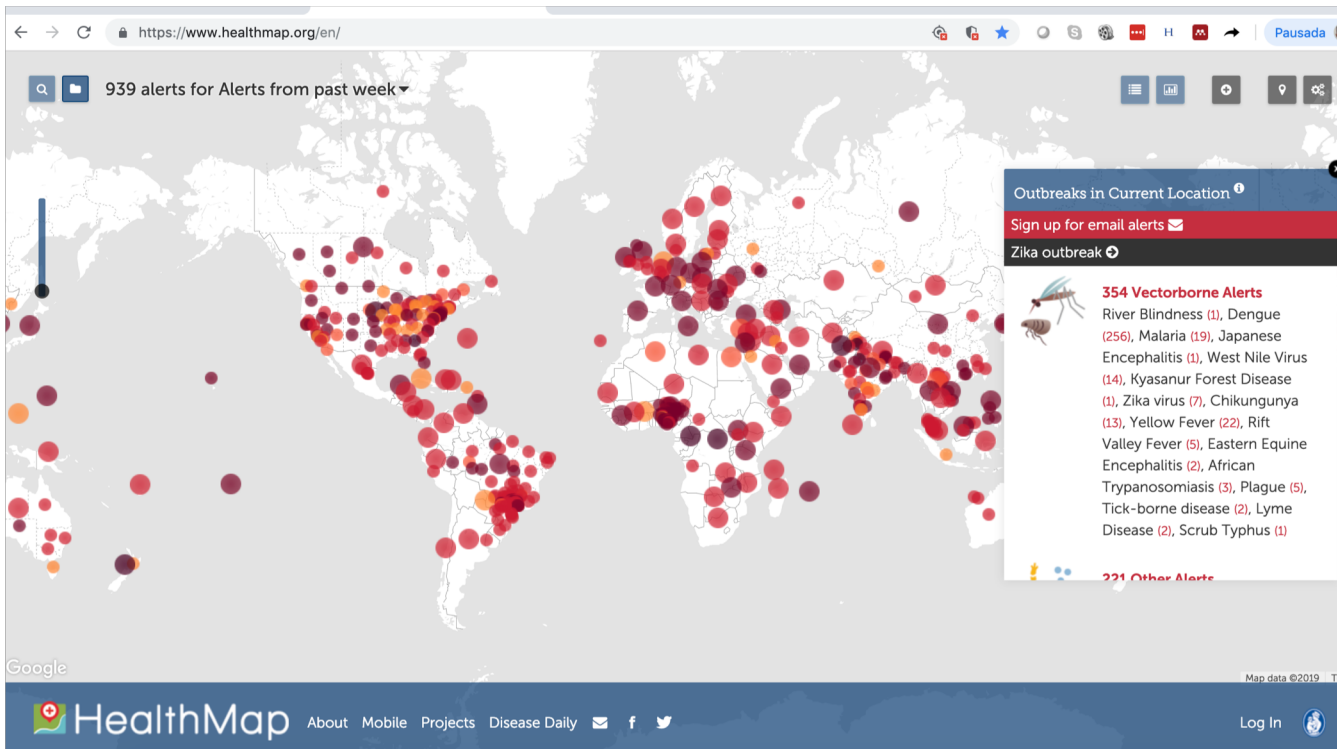
# Drug discovery and repositioning

**Challenge:** costs (5-11bn USD), time (17 years)



Rastegar-Mojarad M. *et al.*, *Nature Biotech*, 2015, doi: 10.1038/nbt.3183  
Scannell J.W. *et al.*, *Nat Review Drug Discovery*, 2012, doi: 10.1038/nrd3681

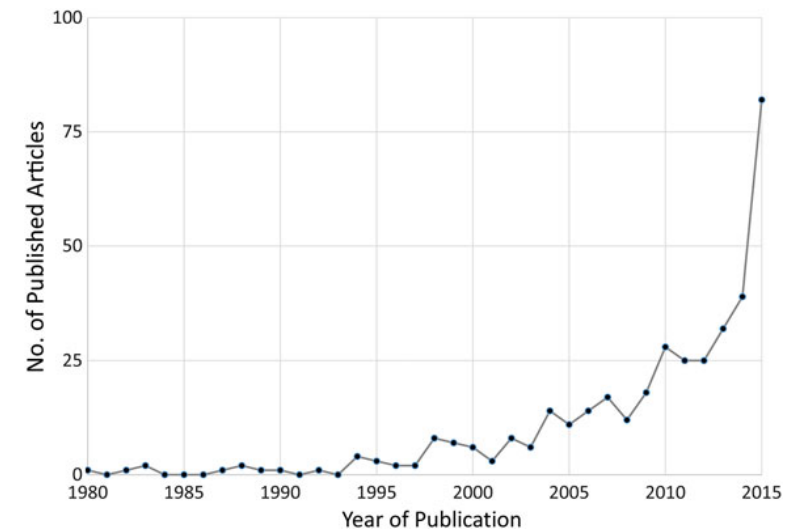
# Aplicações



## Big Data for Infectious Disease Surveillance and Modeling

Shweta Bansal,<sup>1,2</sup> Gerardo Chowell,<sup>1,3</sup> Lone Simonsen,<sup>1,4</sup> Alessandro Vespignani,<sup>5</sup> and Cécile Viboud<sup>1</sup>

<sup>1</sup>Fogarty International Center, National Institutes of Health, Bethesda, Maryland; <sup>2</sup>Department of Biology, Georgetown University, Washington D.C.; <sup>3</sup>School of Public Health, Georgia State University, Atlanta; <sup>4</sup>Department of Public Health, University of Copenhagen, Denmark; and <sup>5</sup>Network Science Institute, Northeastern University, Boston, Massachusetts

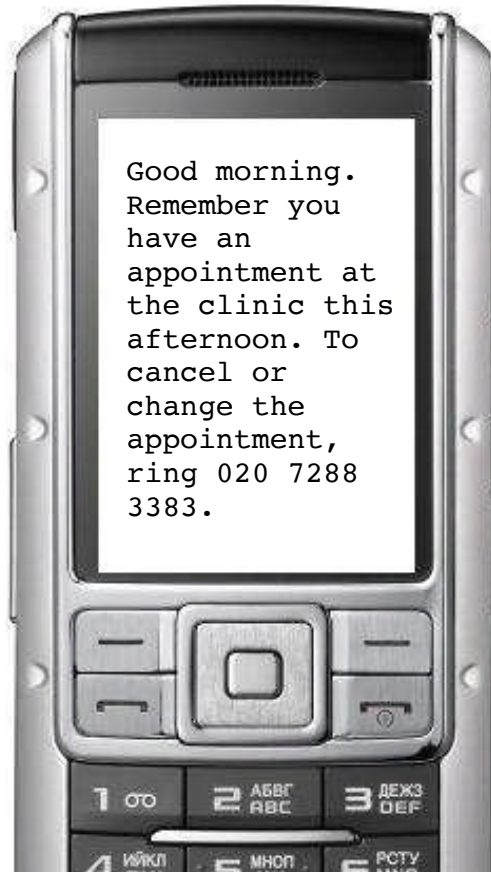


**Figure 1.** Exponential increase since the early 2000s in publications at the intersection of big data and infectious diseases. Annual trends in the number of publications were identified through a Scopus search for articles published between 1980 and 2015, using the following keywords: (big data AND infectious diseases) OR (big data AND epidemics) OR (digital epidemiology AND infectious diseases).

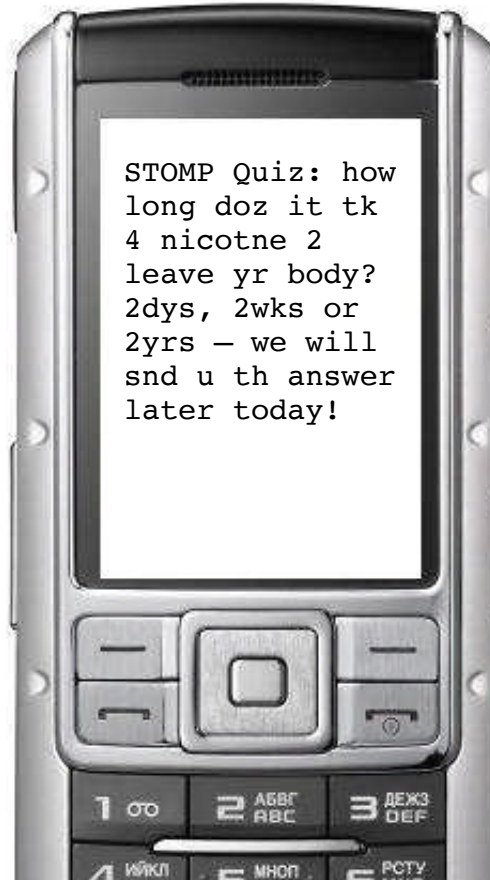
# Aplicações

## m-Health

Dr Henry Potts: h.potts@ucl.ac.uk



Appointment reminder



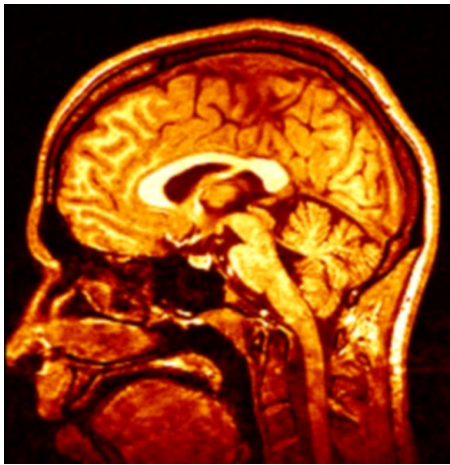
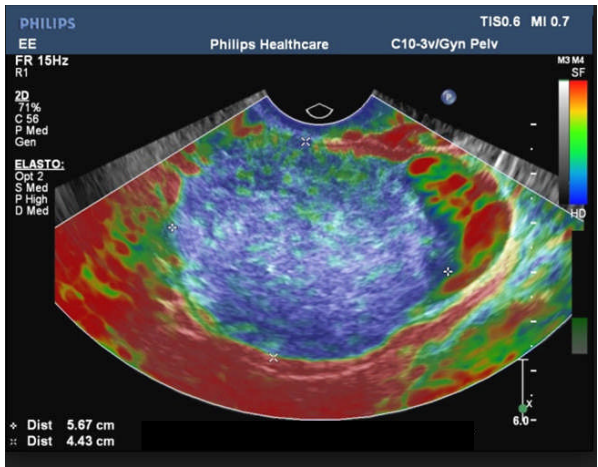
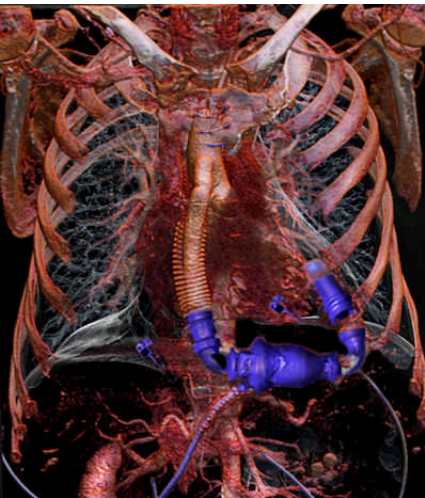
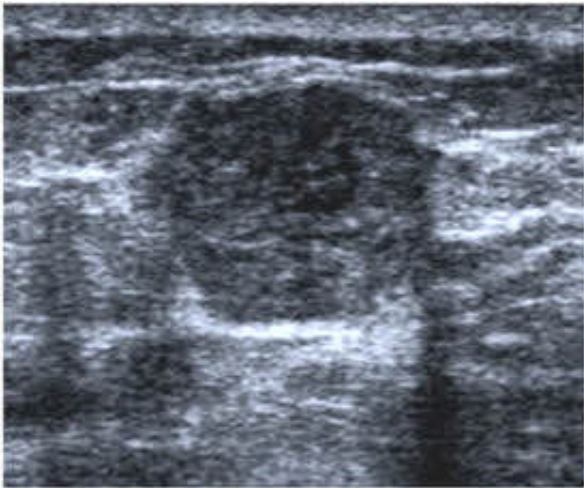
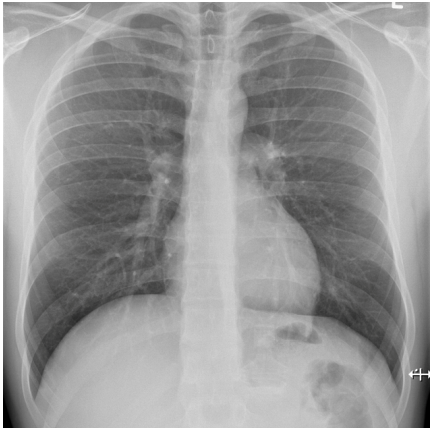
Behaviour change



Personal health record

# Aplicações

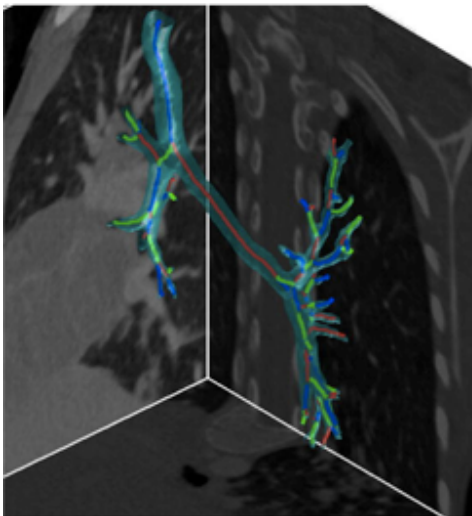
# Medical image processing



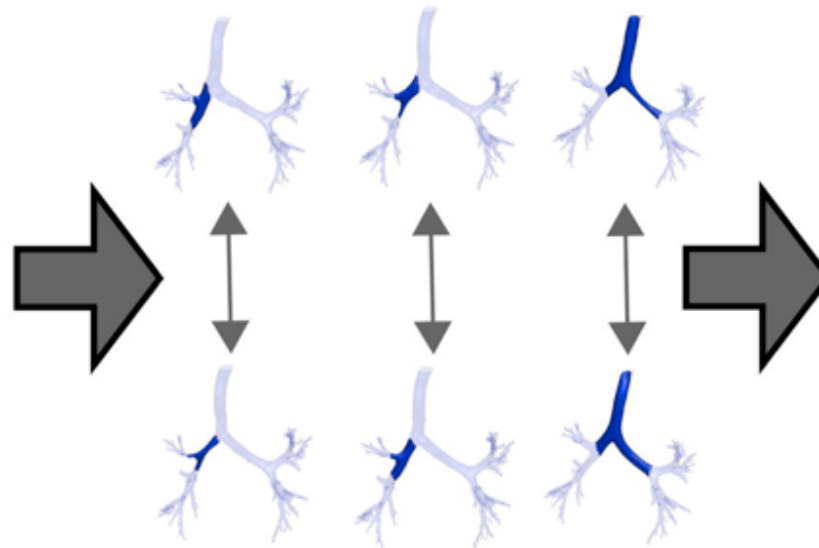
# Aplicações

## Medical image processing

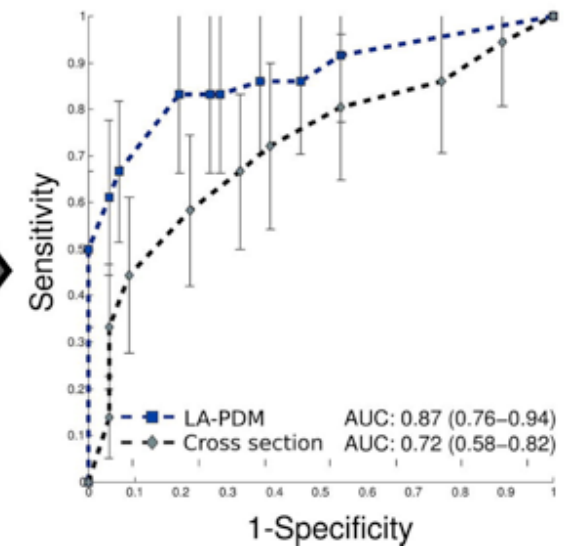
Airway segmentation from CT



Population analysis of airway regions



Classification of pathological variation

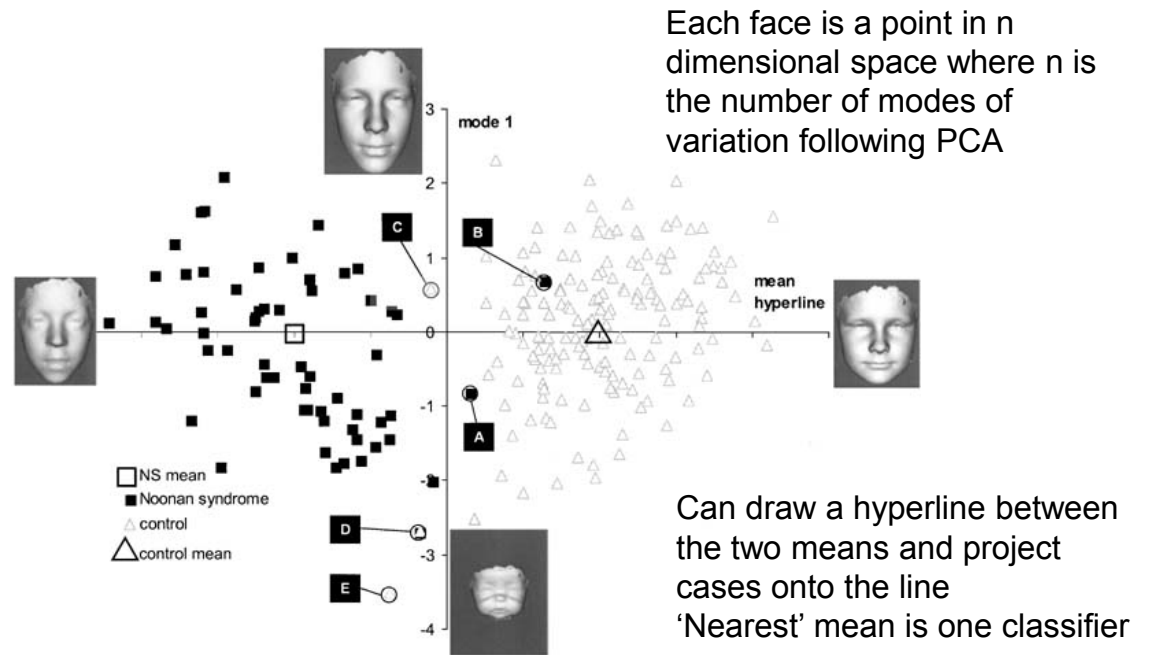


# Aplicações

Each row shows a series of views of the mean face for a different syndrome



## Hammond et al. Using a database of 3D facial scans to identify genetic disorders



# DESAFIOS

## CASO CLÍNICO

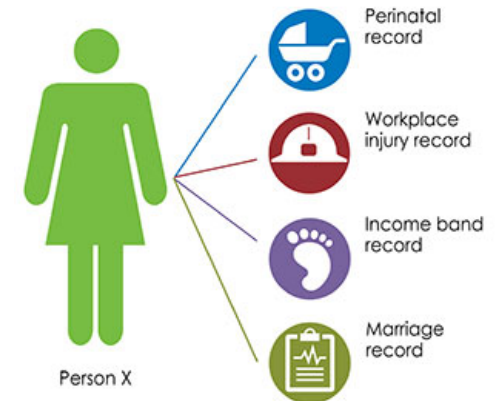
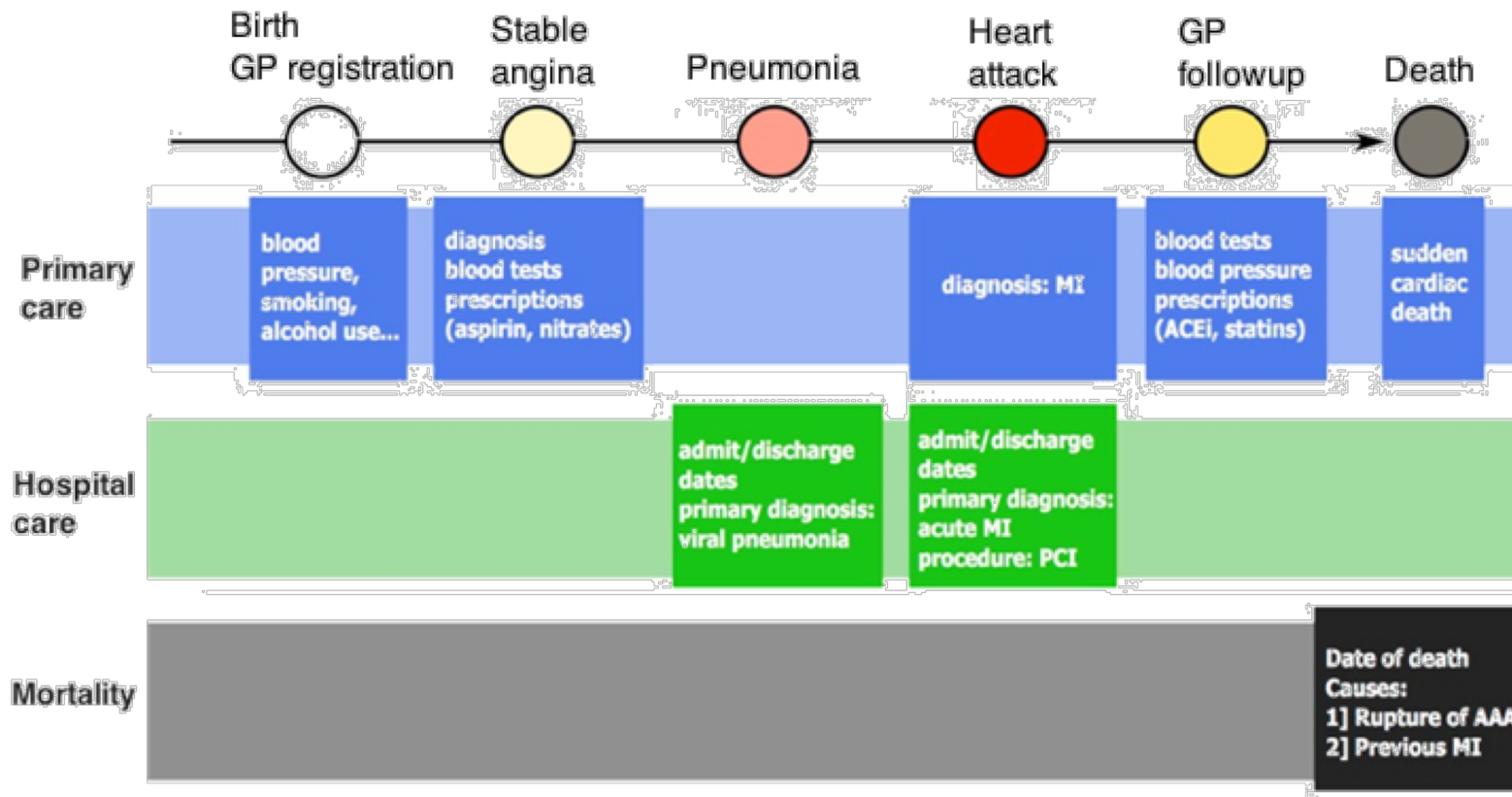
► Caminho do Paciente





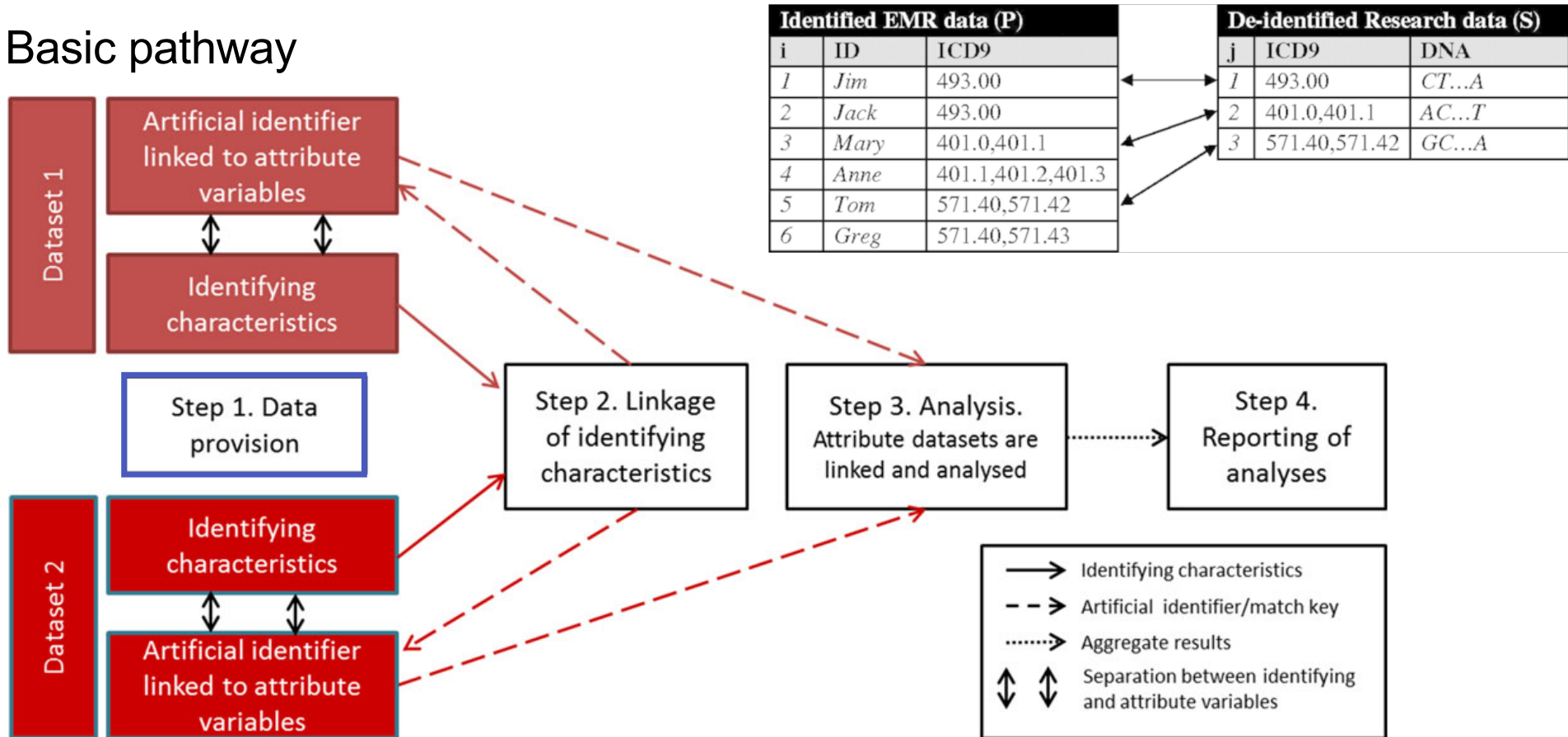
# Vinculação de dados (*data linkage*)

- To bring together electronic records containing information from different sources about an entity (individual, organisation, location etc).



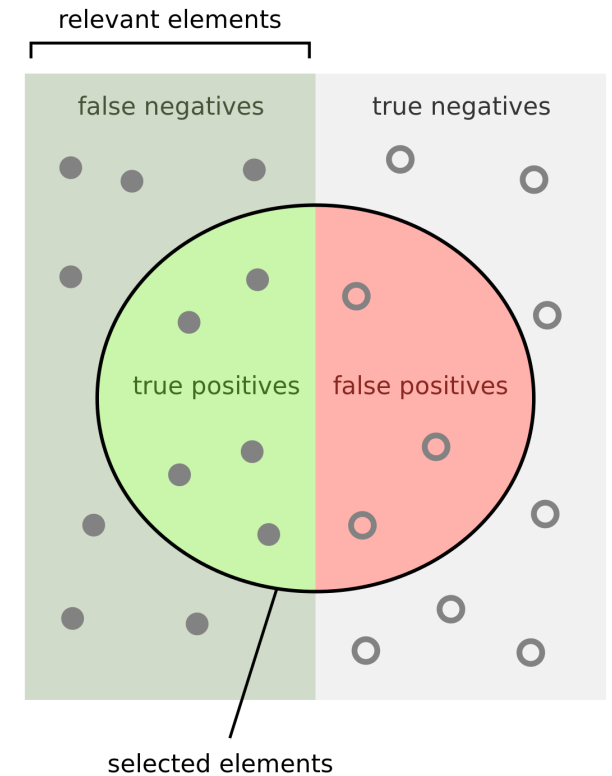
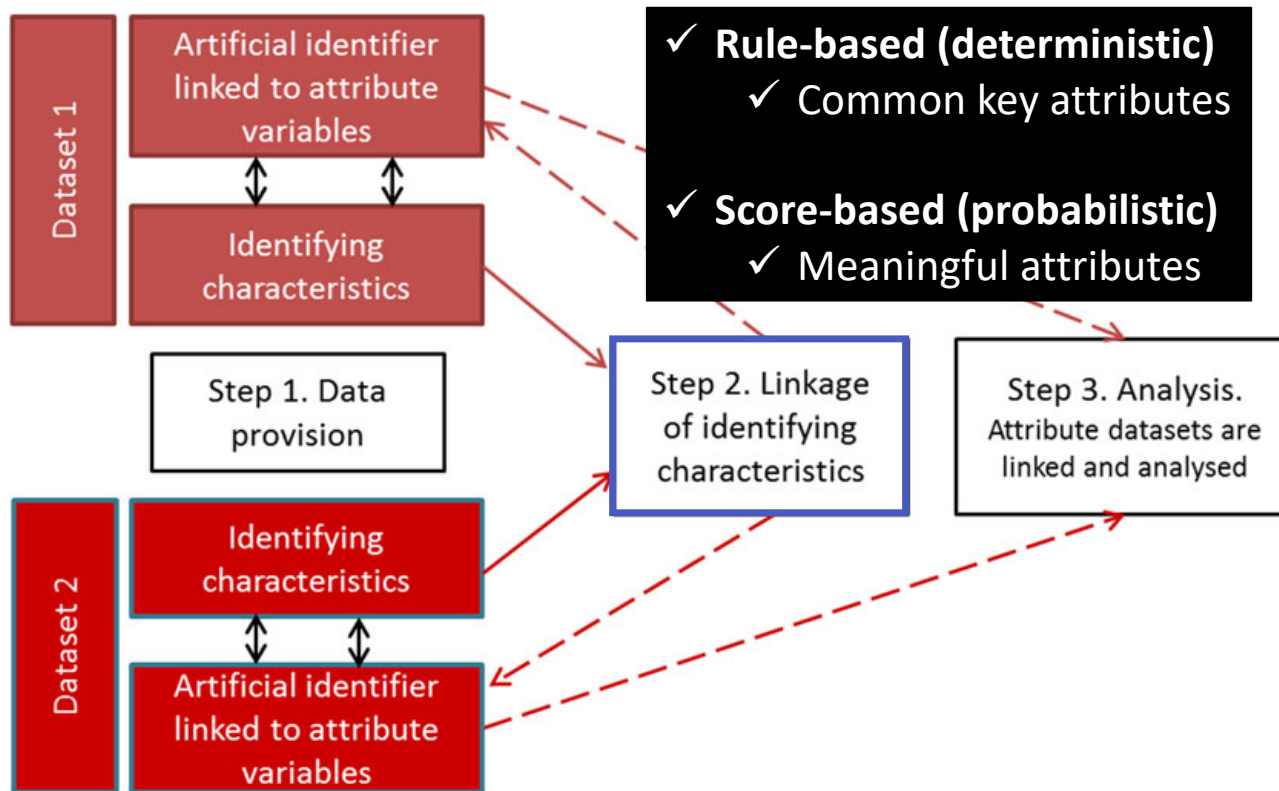
# Vinculação de dados (*data linkage*)

- Basic pathway



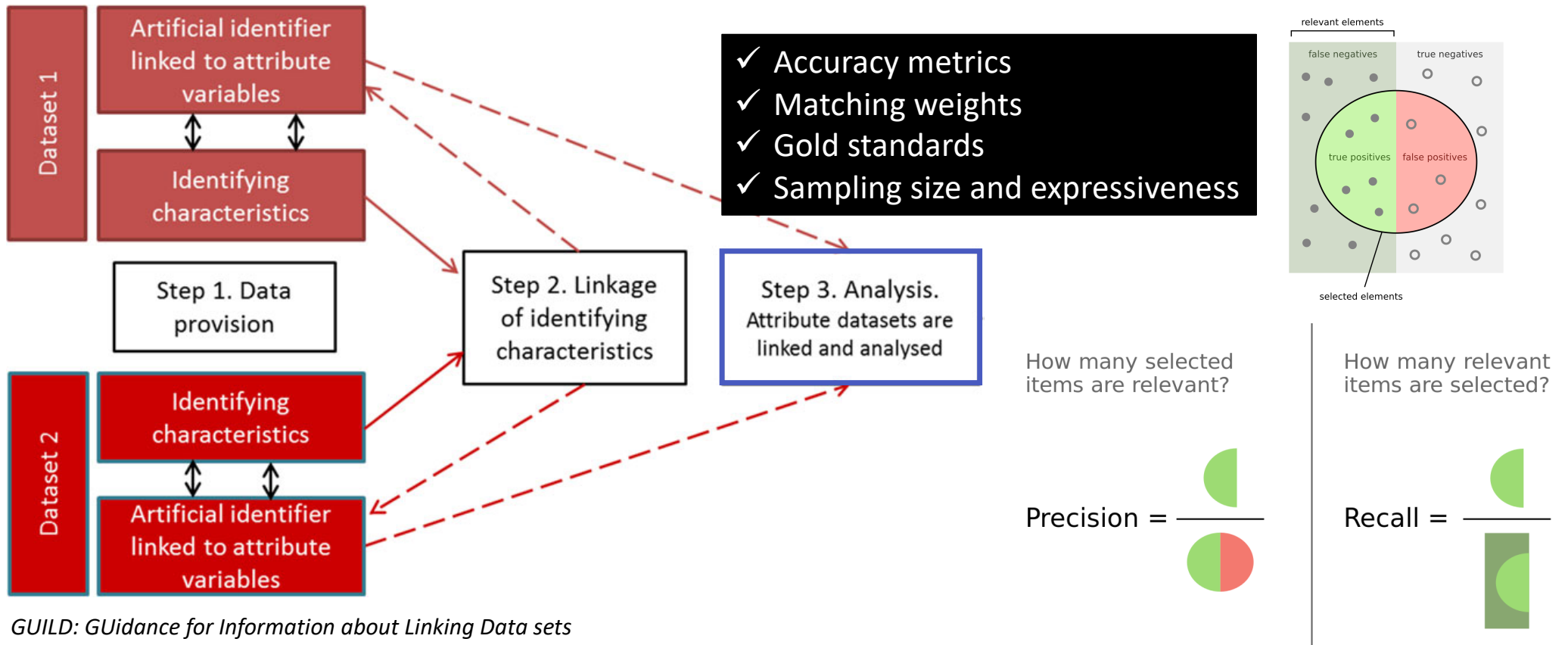
# Vinculação de dados (*data linkage*)

- Pairwise comparison

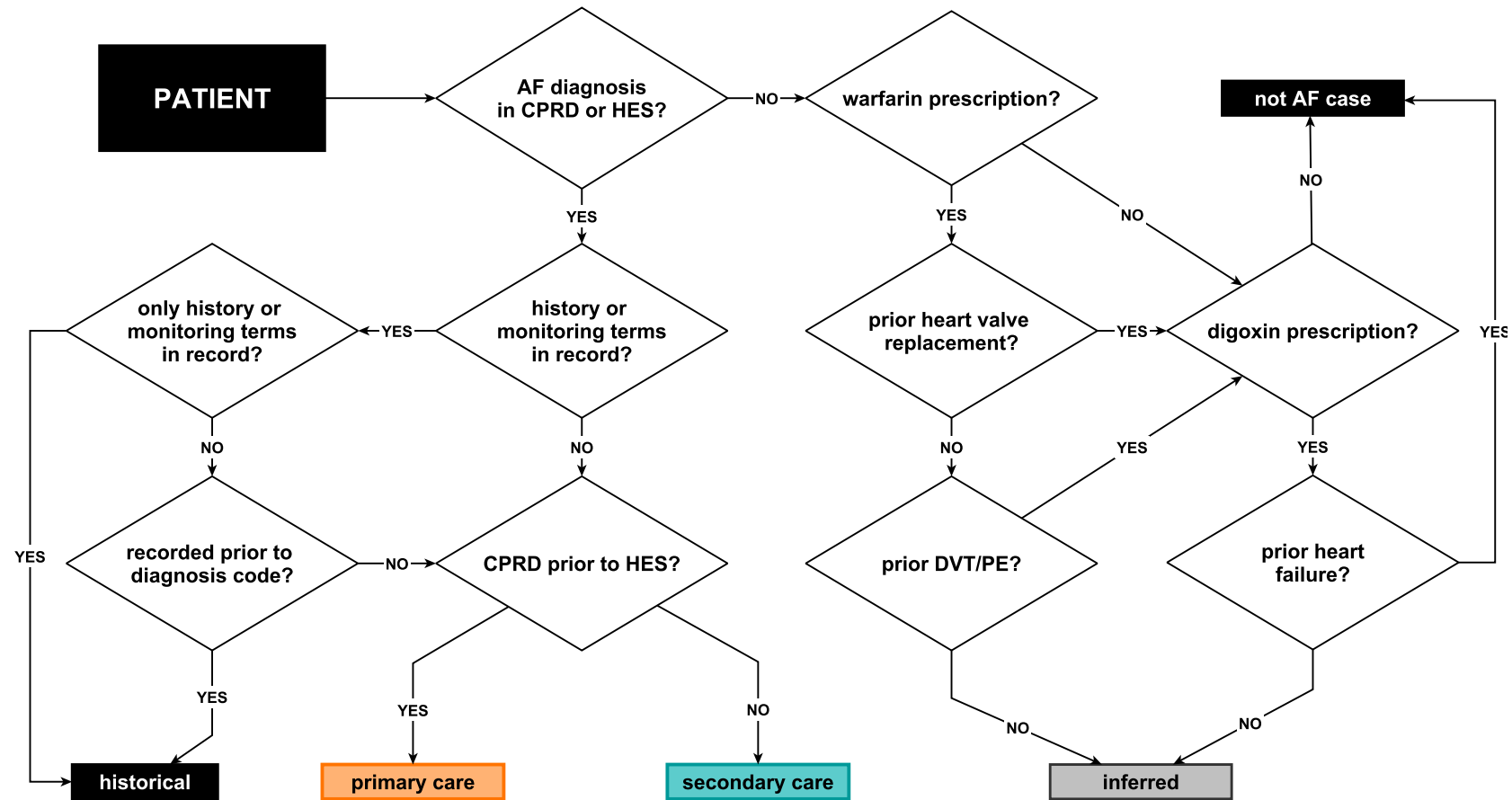


# Vinculação de dados (*data linkage*)

- Accuracy ascertainment



# Extratratificação de pacientes (*phenotyping*)



# Análise de dados temporais



Review

## A Systematic Review on Healthcare Analytics: Application and Theoretical Perspective of Data Mining

Md Saiful Islam<sup>1</sup>, Md Mahmudul Hasan<sup>1</sup>, Xiaoyi Wang<sup>1</sup>, Hayley D. Germack<sup>1,2,3</sup> and Md Noor-E-Alam<sup>1,\*</sup>

### Big Data Analytics in Healthcare – Pattern Mining of Temporal Clinical Events

Svetla Boytcheva<sup>1</sup>, Galia Angelova<sup>1</sup>, Dimitar Tcharaktchiev<sup>2</sup>, Zhivko Angelov<sup>3</sup>

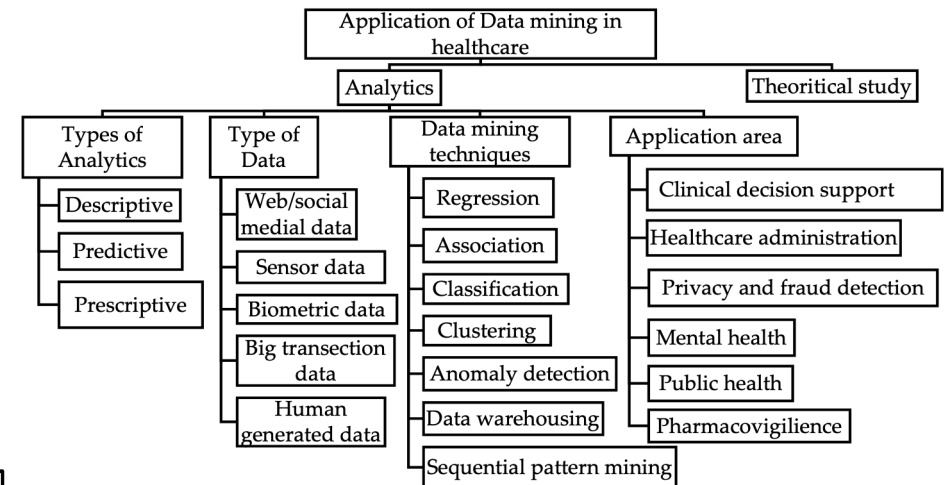
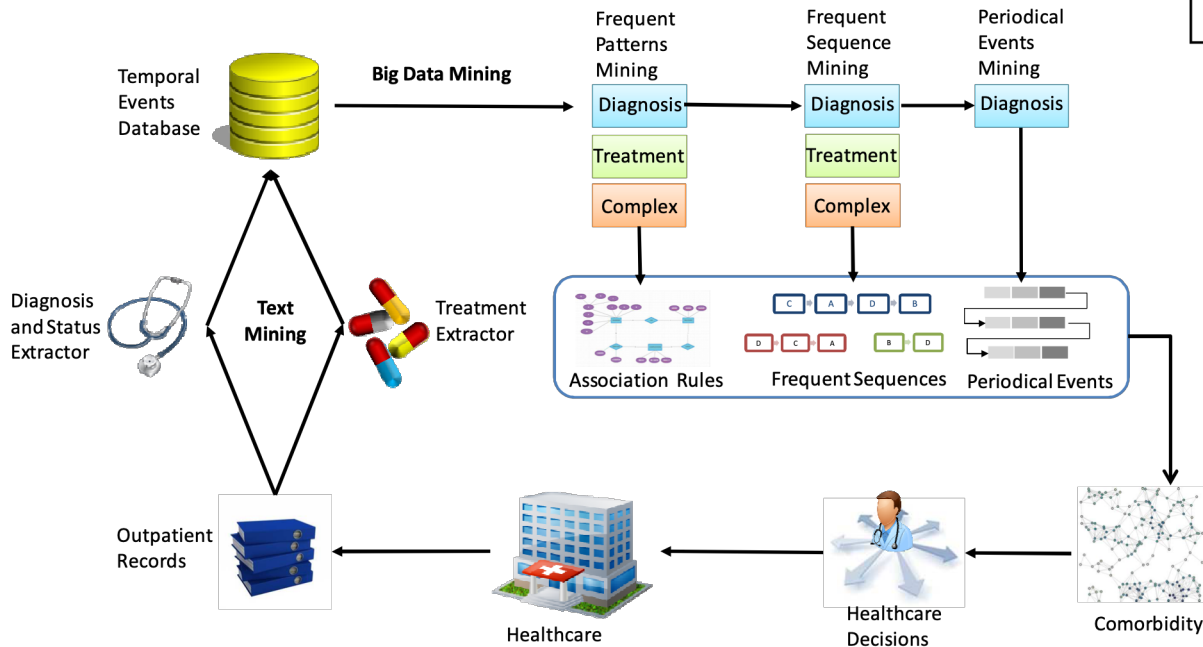
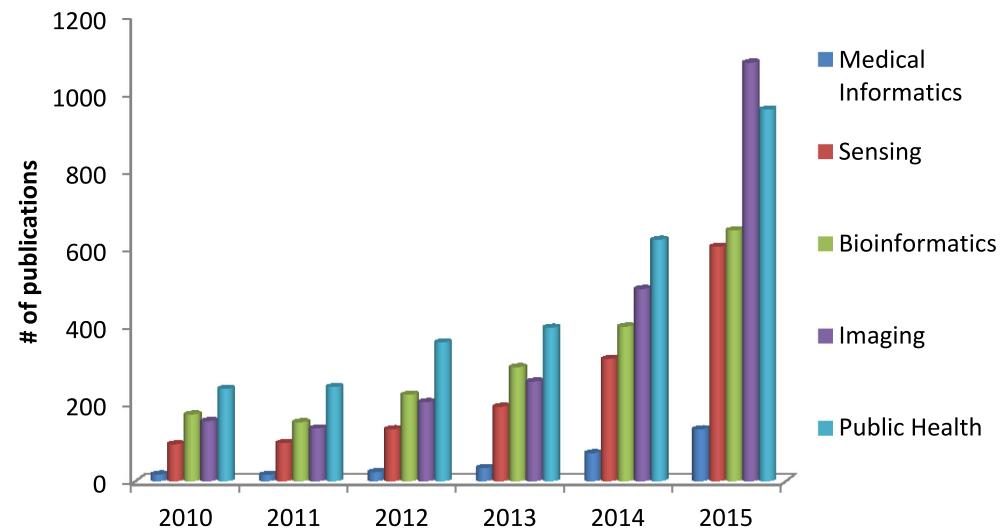


Figure 3. Classification scheme of the literature.

# Análise preditiva

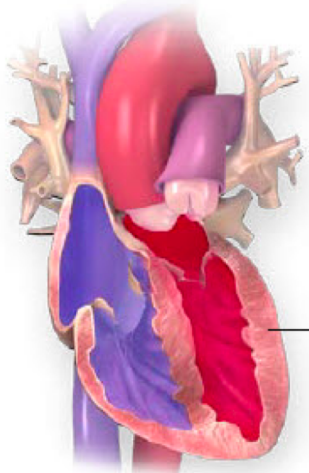
## Deep Learning for Health Informatics

Daniele Ravi, Charence Wong, Fani Deligianni, Melissa Berthelot, Javier Andreu-Perez, Benny Lo,  
and Guang-Zhong Yang, *Fellow, IEEE*



**Fig. 1.** Distribution of published papers that use deep learning in subareas of health informatics. Publication statistics are obtained from Google Scholar; the search phrase is defined as the subfield name with the exact phrase *deep learning* and at least one of *medical* or *health* appearing, e.g., “public health” “deep learning” medical OR health.

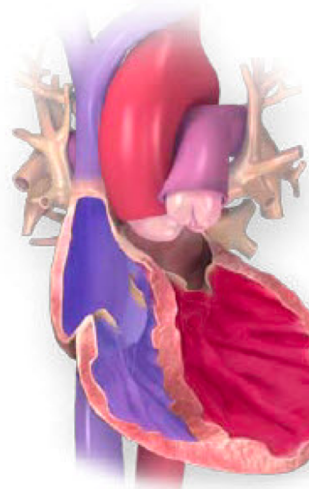
# Análise preditiva



## The Normal Heart

has strong muscular walls which contract to pump blood out to all parts of the body.

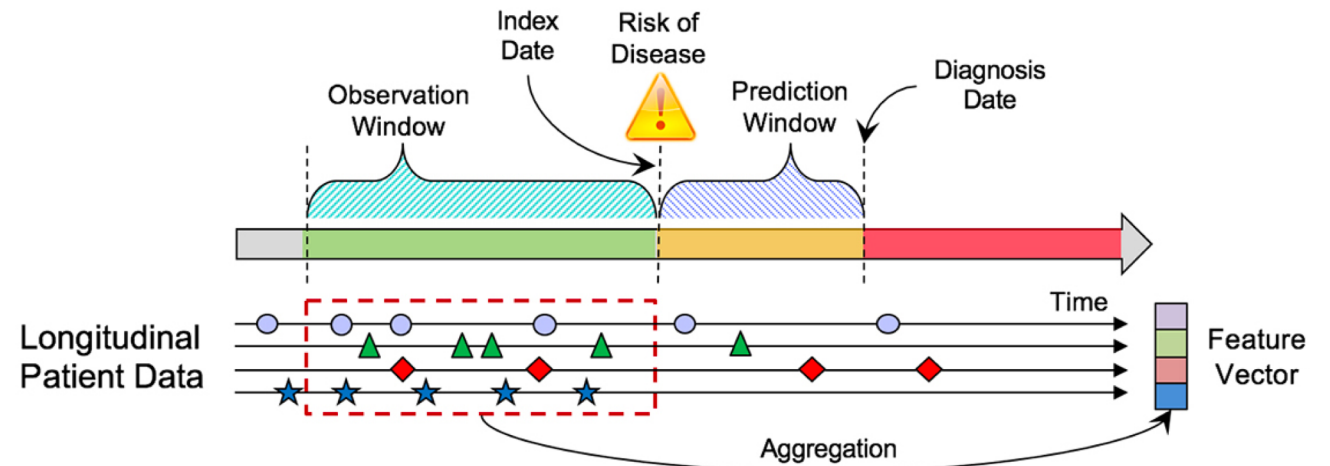
Heart muscle pumps blood out of the left ventricle.



## Heart Failure

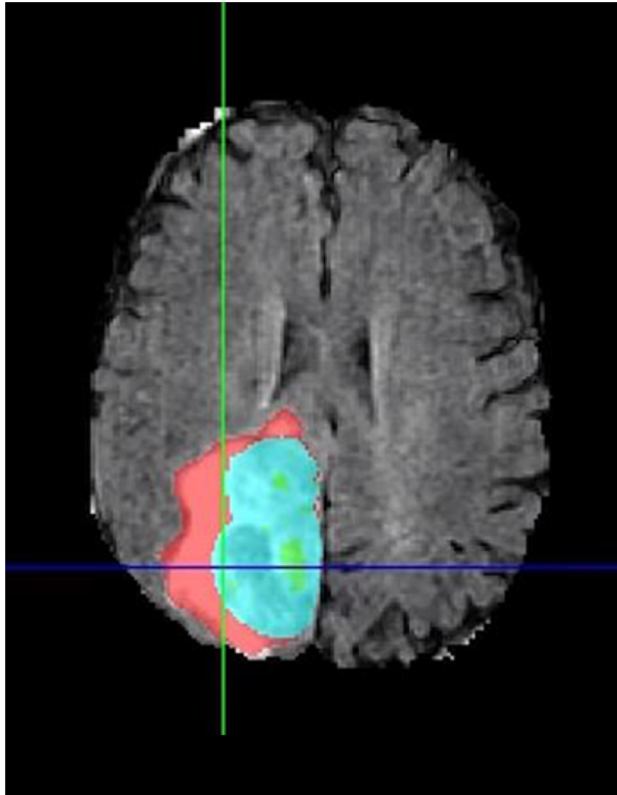
is a condition that causes the muscle in the heart wall to slowly weaken and enlarge, preventing the heart from pumping enough blood.

Weakened muscle prevents left ventricle from pumping enough blood.



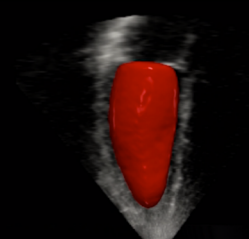


# Processamento de imagens

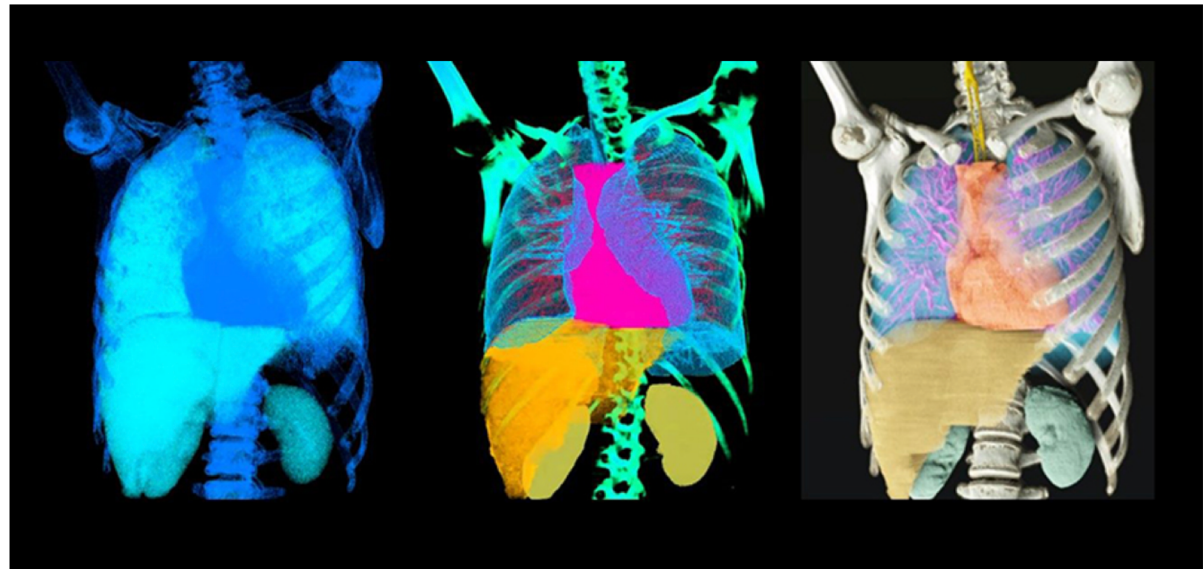


**NVIDIA CLARA PLATFORM**  
Intelligent Compute Platform for Medical Imaging

[DOWNLOAD NOW](#)



The NVIDIA CLARA PLATFORM logo features a stylized red heart shape with a white outline, set against a dark background with a subtle grid pattern.



# Mineração visual de dados

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



## Brazil Health Data - Online Brazil Epidemiology and Real World Data

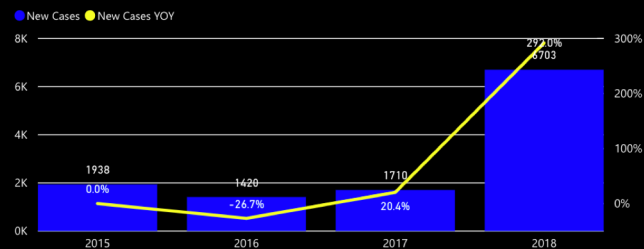
Home Cardiovascular Diseases Infectious Diseases Endocrinology Psychiatry Nervous System Diseases Oncology Dermatology Rheumatology

### Malária Cases

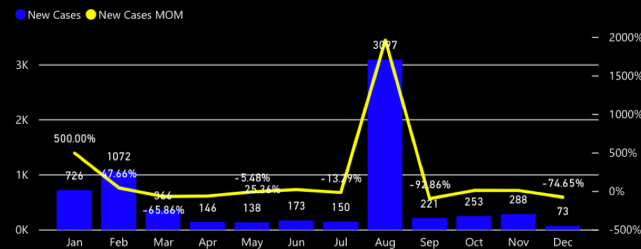
CID10 Codes - B50, B51, B52, B53, B54

2015 2016 2017 2018

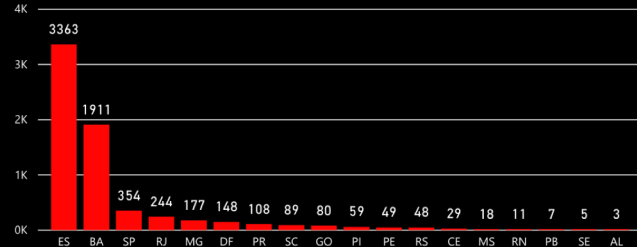
New Cases and New Cases YOY by Year



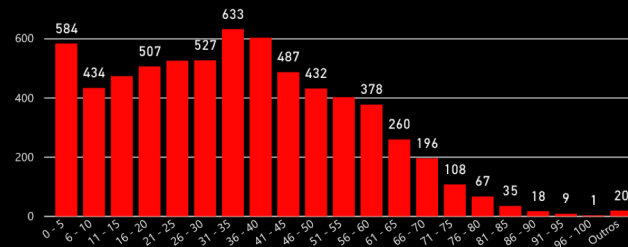
New Cases and New Cases MOM by Month



New Cases by State



New Cases by Age Group



Source: DataSUS

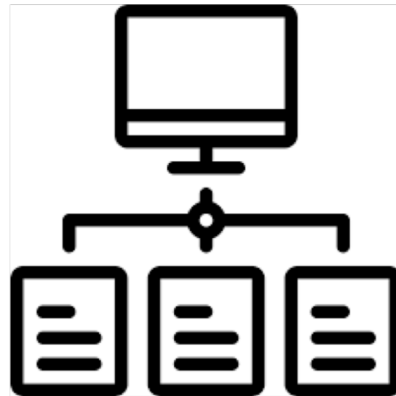


[www.atyimolab.ufba.br](http://www.atyimolab.ufba.br)

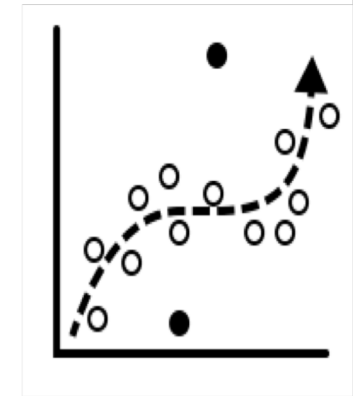
## What we do



cloud  
robotics

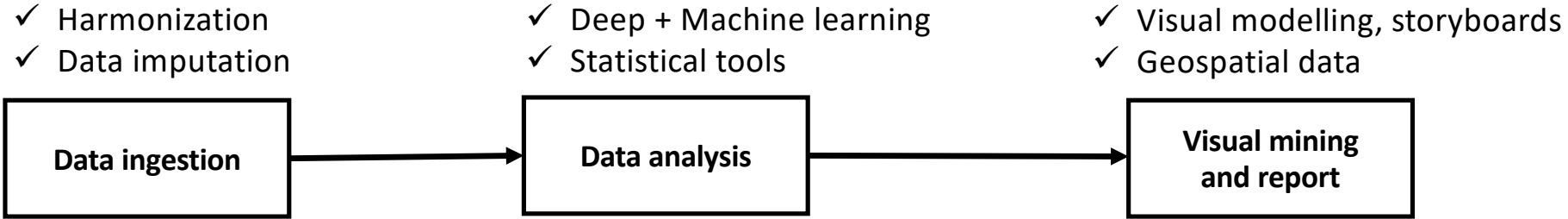
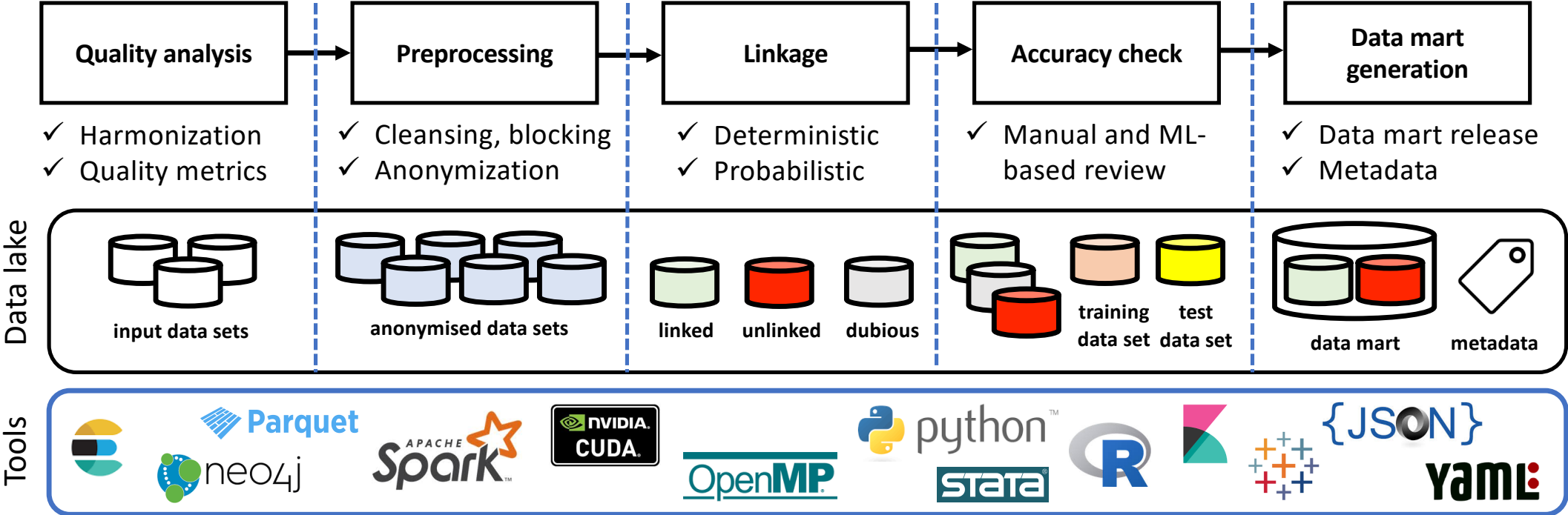


hybrid  
parallel  
computing



big data  
linkage &  
analytics

# AtyImo – Data linkage platform



Data analytics pipeline

# Brazilian governmental databases

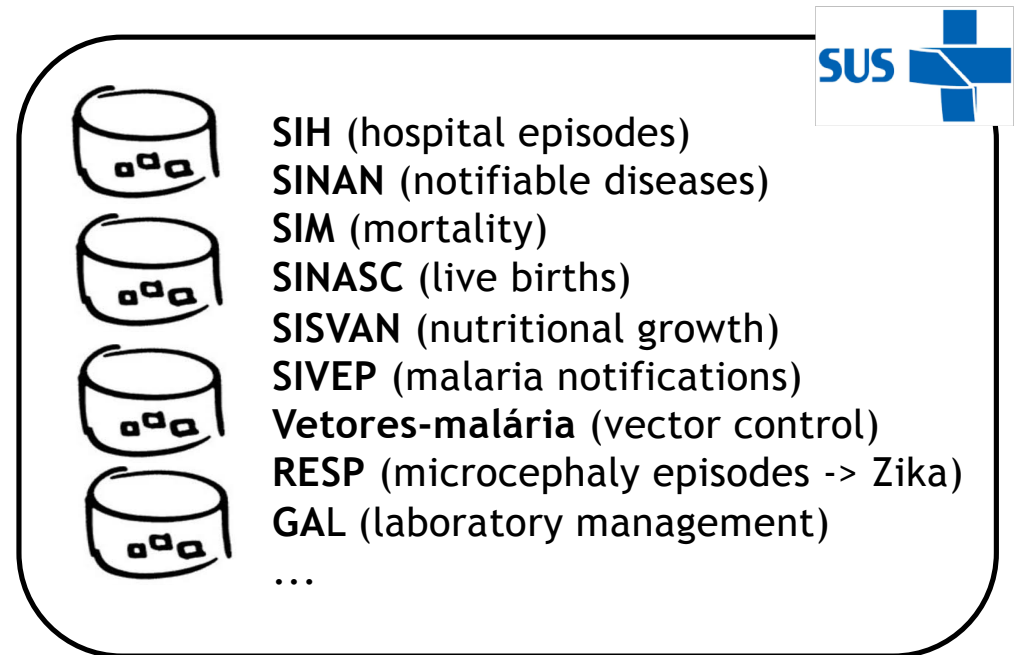
## Social programmes

- ✓ Targeted to poor and extremely poor families.
- ✓ Cadastro Único: central registrar for all programmes.



## Public health system (SUS)

- ✓ Big and complex public health system.
  - from primary care to specialised transplantations.
- ✓ Used by approximately 77% of the Brazilian population (164 million people).



# Existing research platforms we contribute to

## ✓ The 100 Million Cohort



- ✓ Baseline: CadastroÚnico, 2001-2015, **114 million individuals** x 367 attributes.
- ✓ Cohort: baseline + Bolsa Família (cash transfers) + Housing (MCMV), 2001 – 2015, **400 million records**, 3,000 attributes.
- ✓ Used by +20 projects assessing the effects of social programmes on health outcomes.

## ✓ Zika surveillance (+ microcephaly)



- ✓ Birth cohort, 2001 – 2030,  $\cong$  80 million records.
- ✓ Morbidity, mortality, socioeconomic and service data.
- ✓ Focus on the triple epidemic (Zika, Dengue and Chikungunya) and health/educational outcomes related to microcephaly.



## ✓ Malaria linkage & analytics



- ✓ Malaria episodes (>5 million records) + mortality + socioeconomic + climate data, 2000 – 2018.
- ✓ Focus on i) data aggregation and ii) epidemic forecasting.



BILL & MELINDA GATES foundation

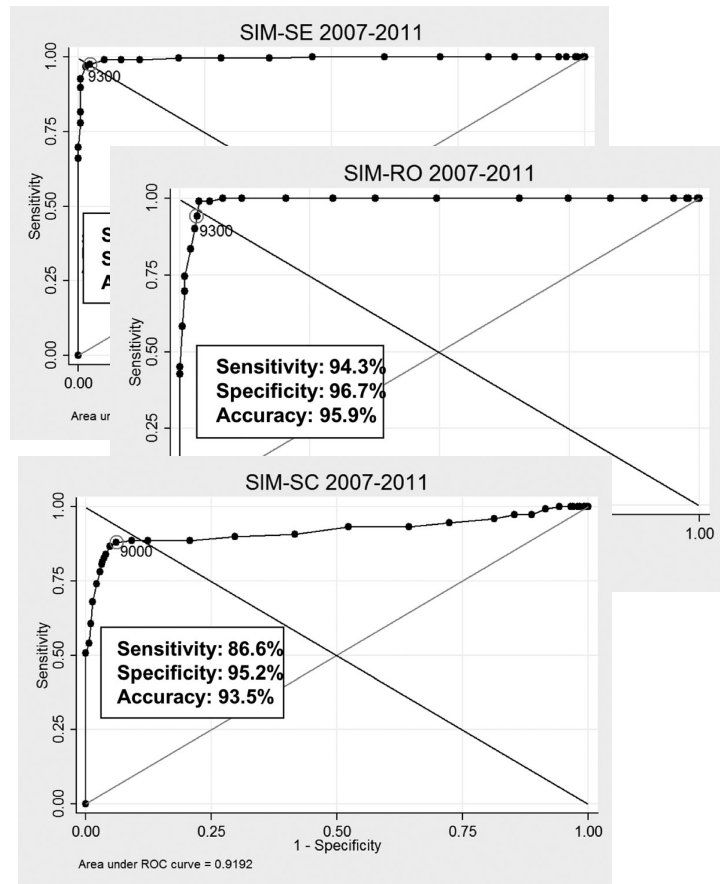


# Example results

346 IEEE JOURNAL OF BIOMEDICAL AND HEALTH INFORMATICS, VOL. 22, NO. 2, MARCH 2018 

## On the Accuracy and Scalability of Probabilistic Data Linkage Over the Brazilian 114 Million Cohort

Robespierre Pita<sup>1</sup>, Clícia Pinto, Samila Sena, Rosemeire Fiaccone, Leila Amorim, Sandra Reis, Maurício L. Barreto<sup>2</sup>, Spiros Denaxas, and Marcos Ennes Barreto



## Exploring hybrid parallel systems for probabilistic record linkage

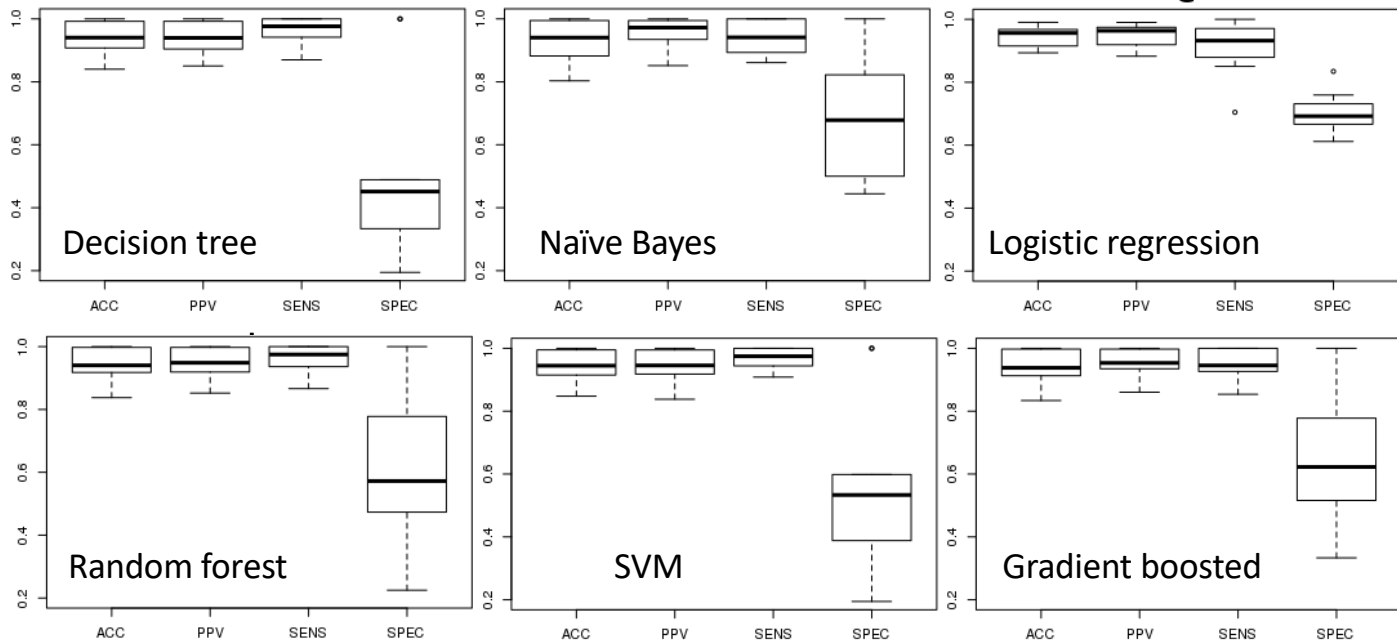
Murilo Boratto<sup>1</sup> · Pedro Alonso<sup>2</sup> · Clícia Pinto<sup>3</sup> · Pedro Melo<sup>3</sup> · Marcos Barreto<sup>3</sup> · Spiros Denaxas<sup>4</sup>

J Supercomput  
<https://doi.org/10.1007/s11227-018-2328-3>

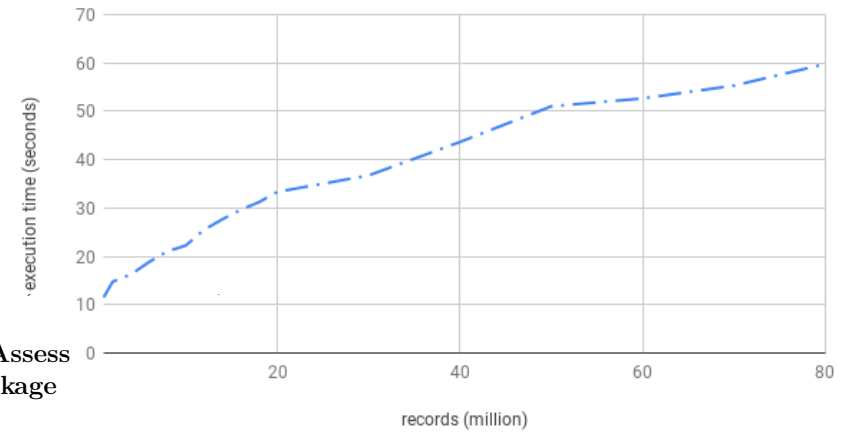
DOI: 10.1007/978-3-319-64283-3\_16

## A Machine Learning Trainable Model to Assess the Accuracy of Probabilistic Record Linkage

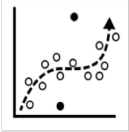
Robespierre Pita<sup>1</sup>, Everton Mendonça<sup>1</sup>, Sandra Reis<sup>2</sup>, Marcos Barreto<sup>1,3</sup>, and Spiros Denaxas<sup>3</sup>



## Hybrid Execution Time







# data linkage & analytics

Latin America Data Science Workshop  
AUGUST 27TH - VLDD 2018 WORKSHOP - RIO DE JANEIRO, BRAZIL

## Applying term frequency-based indexing to improve scalability and accuracy of probabilistic data linkage

Robespierre Pita<sup>1,2</sup>, Luan Menezes<sup>1,2</sup>, Marcos E. Barreto<sup>1,2</sup>

Figura 2. Term frequency-based approach used in Atylmo.

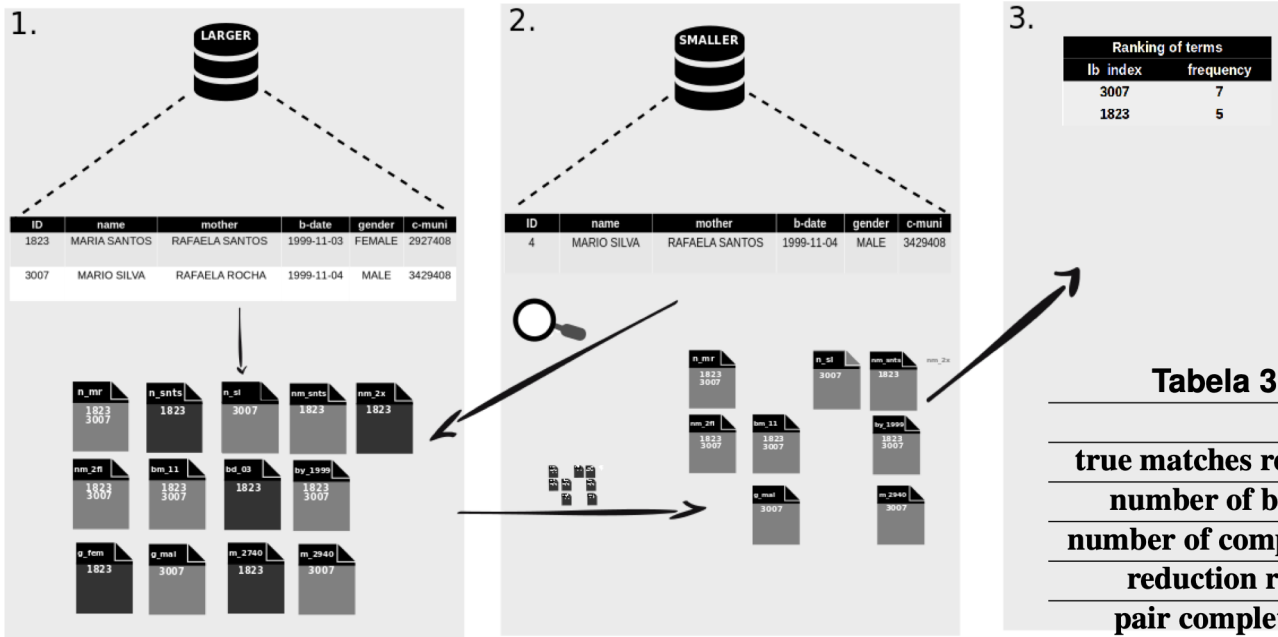


Tabela 1. Gold standard data set used for validation.

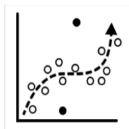
|                         |                    |
|-------------------------|--------------------|
| SIM                     | 6,458 records      |
| SINASC                  | 13,046 records     |
| Total of comparisons    | 84,251,068 records |
| Expected true positives | 3,030 records      |

Tabela 2. Size of generated blocks for each indexing technique

| method      | predicate 1 |           | predicate 2 |           | term frequency |           |
|-------------|-------------|-----------|-------------|-----------|----------------|-----------|
|             | <i>sb</i>   | <i>lb</i> | <i>sb</i>   | <i>lb</i> | <i>sb</i>      | <i>lb</i> |
| <b>min</b>  | 1           | 1         | 1           | 1         | 1              | 100       |
| <b>med</b>  | 24          | 51        | 2           | 2         | 1              | 100       |
| <b>mean</b> | 43          | 88.38     | 8.289       | 11.57     | 1              | 100       |
| <b>max</b>  | 1855        | 41528     | 87          | 611       | 1              | 100       |

Tabela 3. Results of each indexing technique used.

|                               | predicate 1 | predicate 2 | term frequency |
|-------------------------------|-------------|-------------|----------------|
| <b>true matches retrieved</b> | 2,382       | 3,018       | 3,020          |
| <b>number of blocks</b>       | 5,806       | 6,432       | 6,458          |
| <b>number of comparisons</b>  | 44,406,049  | 29,111,755  | 645,800        |
| <b>reduction ratio</b>        | 0.472       | 0,654       | 0,992          |
| <b>pair completeness</b>      | 0.786       | 0.996       | 0.996          |





data linkage & analytics


## ✓ Integrating socioeconomic and healthcare data to combat malaria


✓ Phase I: November 2016 - October 2018

✓ Focus on i) data aggregation and ii) epidemic forecasting.

|                                                                                   |                       |
|-----------------------------------------------------------------------------------|-----------------------|
|  | ✓ Coverage: 2003-2018 |
|                                                                                   | ✓ Records: 5,340,564  |
| <b>SIVEP</b>                                                                      | ✓ Attributes: 52      |

|                                                                                     |                       |
|-------------------------------------------------------------------------------------|-----------------------|
|  | ✓ Coverage: 2003-2018 |
|                                                                                     | ✓ Records: 1,004      |
| <b>SIM</b>                                                                          | ✓ Attributes: 37      |

|                                                                                    |                       |
|------------------------------------------------------------------------------------|-----------------------|
|  | ✓ Coverage: 2003-2018 |
|                                                                                    | ✓ Records: 46,170     |
| <b>SINAN</b>                                                                       | ✓ Attributes: 20      |

|                                                                                      |                       |
|--------------------------------------------------------------------------------------|-----------------------|
|  | ✓ Coverage: 2003-2018 |
|                                                                                      | ✓ Records: 5,570      |
| <b>Climate</b>                                                                       | ✓ Attributes: 5       |





- Informações Gerais
- » Base de dados
- » Dicionário de dados
- Mineração de dados <
- Mineração Visual de dados <
- Estadística <
- Análise Univariada
- Séries Temporais
- Graficos de Controle
- Análise Bivariada
- Operacional <
- Analytics <

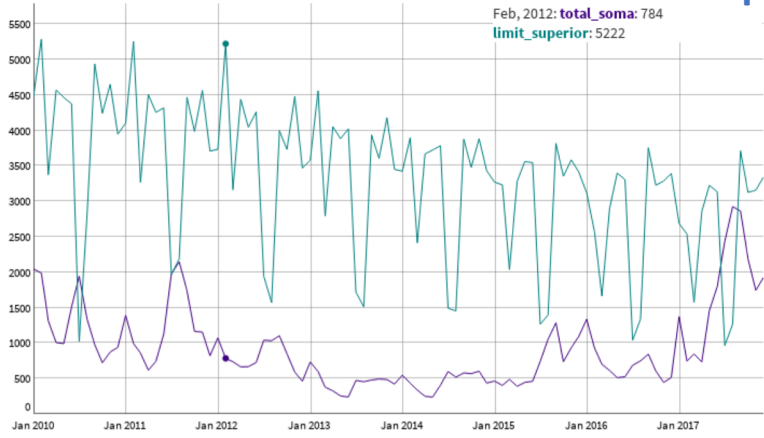
Input

Select the Variable:

3° Quartil

Manaus-AM

Grafico de controle



[http://200.128.60.86:3838/shiny\\_integracao/](http://200.128.60.86:3838/shiny_integracao/)

Descriptive analytics

- Estadística <
- Análise Univariada
- Séries Temporais
- Graficos de Controle
- Análise Bivariada
- Operacional <
- Analytics <

Climate Variable:

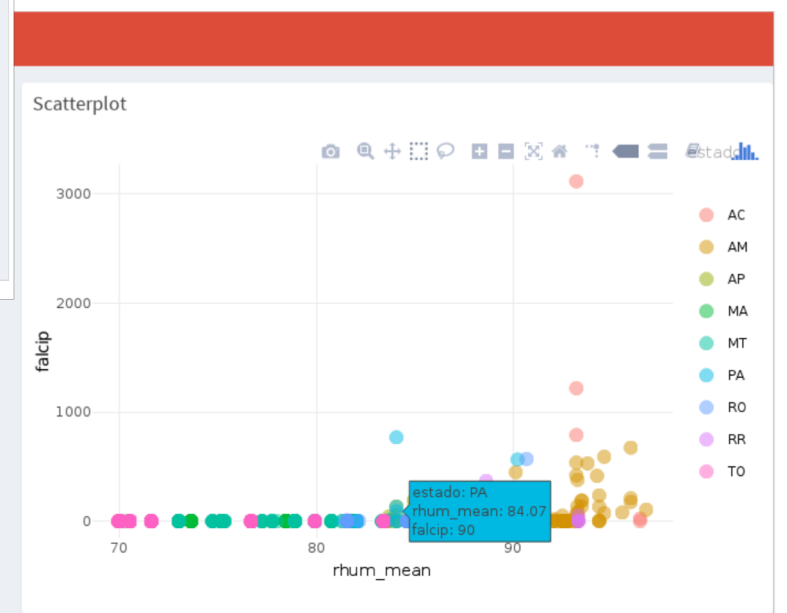
Umidade do Ar

Coloring by:

Estado

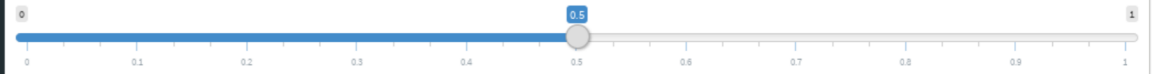
Year:

2015



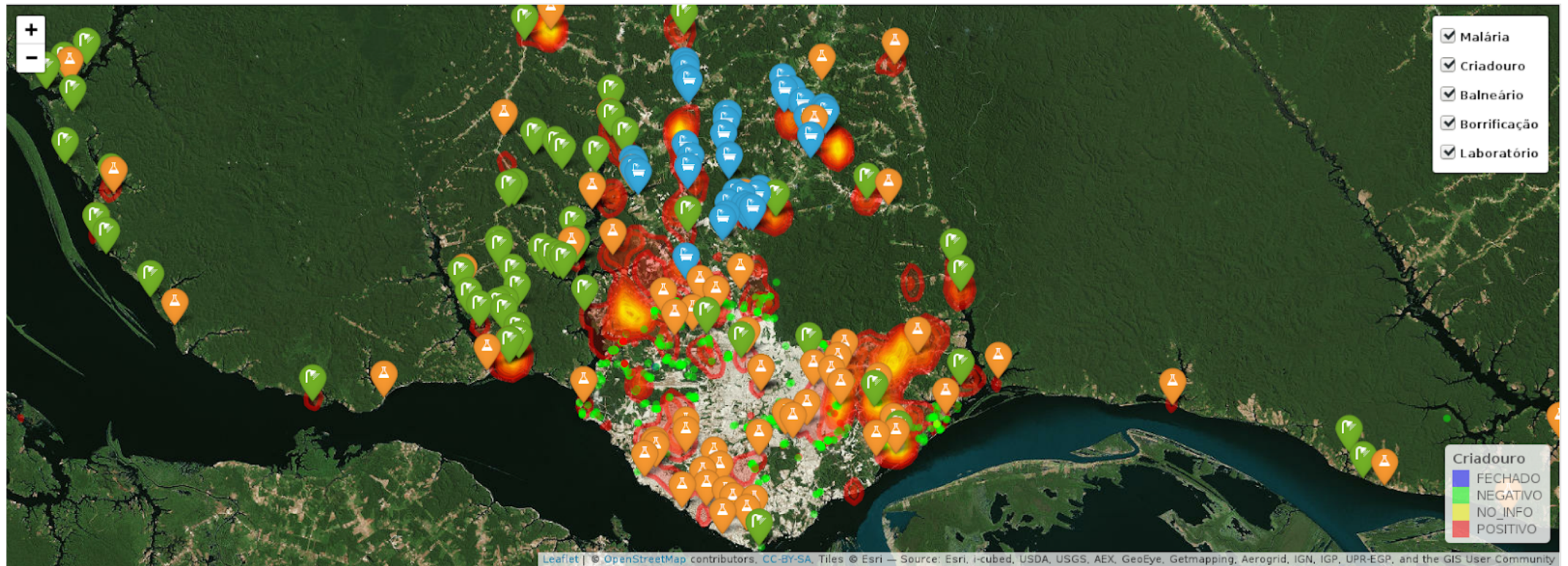
More Inputs

Transparency:




Descriptive analytics

## Example: multilayer visual mining



# Existing results



**MEDTROP**  
54º CONGRESSO DA SOCIEDADE BRASILEIRA DE MEDICINA TROPICAL  
02 a 05 Setembro 2018 - Centro de Convenções de Pernambuco | (Recife PE)

**INTEGRAÇÃO E MINERAÇÃO VISUAL DE DADOS PARA ESTUDO DA MALÁRIA NO BRASIL**

**ALBERTO SIRONI<sup>1</sup>, JURACY BERTOLDO JUNIOR<sup>1</sup>, MARCOS E. BARRETO<sup>1</sup>, VANDERSON SAMPAIO<sup>2</sup>, ANDRÉ SIQUEIRA<sup>3</sup>**

<sup>1</sup> DEPARTAMENTO DE CIÊNCIA DA COMPUTAÇÃO, UNIVERSIDADE FEDERAL DA BAHIA (SALVADOR, BA), <sup>2</sup> FUNDAÇÃO DE VIGILÂNCIA EM SAÚDE DO AMAZONAS (MANAUS, AM), <sup>3</sup> INSTITUTO NACIONAL DE INFECTOLOGIA EVANDRO CHAGAS (RIO DE JANEIRO, RJ)



## Título

INTERACTIVE DATA VISUALIZATION OF MALARIA USING R SHINY

## Resumo

Data visualization consists in representing data in some systematic form including attributes and variables for the unit of information. A simple and quick information can highlight possible errors with data just as it helps uncover interesting trends. There are traditional and new approaches for visualization methods. Data coming from different sources (SIVEP, IBGE, NOAA, MDS) were considered and are potentially useful for the effective understanding of malaria in the Amazon region. Using data of malaria, we illustrated the process of exploratory analysis and traditional visualization tools that can make the evaluation of high dimensional data available and feasible. Exploratory analysis tools (univariate and bivariate) were used to enhance the data visualization techniques. We believe that data integration and the exploration of visualization tools, such as those available using R Shiny, can assist decision making, provide significant contribution for understanding several processes and make massive amounts of available data in several systems become useful.

## Autores

André Alves Ferreira Mendes, Rosemeire Leovigildo Fiaccone, Leila Denise Alves Ferreira Amorim, Marcos Ennes Barreto, Juracy Bertoldo Santos Junior, Alberto Sironi, Marcos Aurelio Eustorgio Filho

International Journal of Population Data Science (2018) 3:3:453

## International Journal of Population Data Science



Journal Website: [www.ijpds.org](http://www.ijpds.org)

## Linking surveillance and climate data to combat malaria

Sironi, A<sup>1</sup>, Barreto, M<sup>1</sup>, Bertoldo, J<sup>1</sup>, Conceição, D<sup>1</sup>, and Sampaio, V<sup>2</sup>



## Frontiers of Engineering for Development symposium:

*Engineers as healthcare practitioners*

**Ho Chi Minh City, Vietnam  
30 October to 2 November 2018**



**CLOSER 2019**  
9<sup>th</sup> INTERNATIONAL CONFERENCE ON CLOUD COMPUTING AND SERVICES SCIENCE  
HERAKLION, CRETE - GREECE 2 - 4 MAY, 2019

# Current research / projects

## • The 100 million Brazilian cohort

Funding:

BILL & MELINDA  
GATES foundation



- Principal investigators: Maurício Barreto (FIOCRUZ BA), Gerson Penna (FIOCRUZ DF), Laura Rodrigues (London School of Hygiene and Tropical Medicine), Liam Smeeth (London School of Hygiene and Tropical Medicine).
- Period: 2015-2019
- Scope: i) integration of socioeconomic data from CadastroÚnico and Bolsa Família (conditional cash transfer programme) databases to build a huge population-based cohort covering the period 2007-2015. Current cohort size is **114 million records**; ii) design a probabilistic data linkage tool (**AtyImo**) to link this cohort with Public Health databases and generate domain-specific data from several epidemiological studies on HIV, tuberculosis, leprosy etc; iii) propose and validate statistical approaches to probabilistic linkage of huge datasets; iv) promote technology transfer and capacity building on big data integration.  
More information [here](#).



## • Long-term surveillance platform for Zika virus and microcephaly

Funding:



- Principal investigators: Maurício Barreto (FIOCRUZ/BA), Maria Glória Teixeira (UFBA), Cláudio Henriques Maierovisch (Ministry of Health).
- Period: 2016-2020
- Scope: i) Design a cohort based on live births (from SINASC database) from 2001 to 2030; ii) assessment of health and educational outcomes related to Zika virus and microcephaly.  
More information [here](#).

## • The 100 million Brazilian linked data and datacentre.

Funding:



- Principal investigators: Mauricio Barreto (FIOCRUZ/BA), Laura Rodrigues (London School of Hygiene and Tropical Medicine).
- Period: 2017-2022
- Scope: i) link electronic health records from Brazilian governmental databases; ii) build the **CIDACS** datacentre and its public interface.  
More information [here](#).

# Current research / projects

- **Design of a scientific repository (data lake) for big data applications**

Funding:



- Principal investigator: Marcos Barreto (UFBA).
- Period: 2016-2019
- Scope: Design and deployment of a data repository (data lake) for big data applications. The first prototype comprises malaria surveillance data to support predictive analytics.

- **Treating heterogeneity and uncertainty in data integration: case study on Brazilian databases.**

Funding:



- Principal investigators: Marcos Barreto (UFBA), Spiros Denaxas (Farr Institute of Health Informatics Research).
- Period: 2016-2018
- Scope: i) design and validation of a data integration model and related computing tools addressing heterogeneity, uncertainty and scalability targeted to big data integration; ii) support for some Brazil-UK ongoing projects: the 100 million cohort, the surveillance platform for Zika and microcephaly, and predictive analytics methods applied to Malaria data (Post-doctoral proposal).  
More information [here](#).

# Current research / projects

- **Integrating socioeconomic and health data to combat malaria.**

- Principal investigator: Marcos Barreto (UFBA), Spiros Denaxas (Farr Institute of Health Informatics Research).
  - Period: 2016-2018
  - Scope: i) develop a platform to integrate surveillance data from Malaria with other sources (socioeconomic and public health data); ii) design and validate predictive analytics methods to help on Malaria elimination.
- More information [here](#).

Funding:

BILL & MELINDA  
GATES foundation



- **Early childhood development friendly index: assessing the enabling environment for Nurturing Care.**

- Principal investigator: Muriel Gubert (UnB).
- Team: Marcos Barreto (UFBA), Gabriela Buccini (Yale School of Public Health), Rafael Perez-Escamilla (Yale School Of Public Health), Sonia Isoyama Venancio (Health Institute of São Paulo)
- Period: 2018-2020
- Scope: This project aims to develop an ECD (Early Childhood Development) friendly index (ECD-FI), based on a core set of evidence-based Nurturing Care indicators, to assess the enabling environment and promote ECD at the municipality level by monitoring and identifying opportunities to scale up ECD programs locally. More information [here](#).

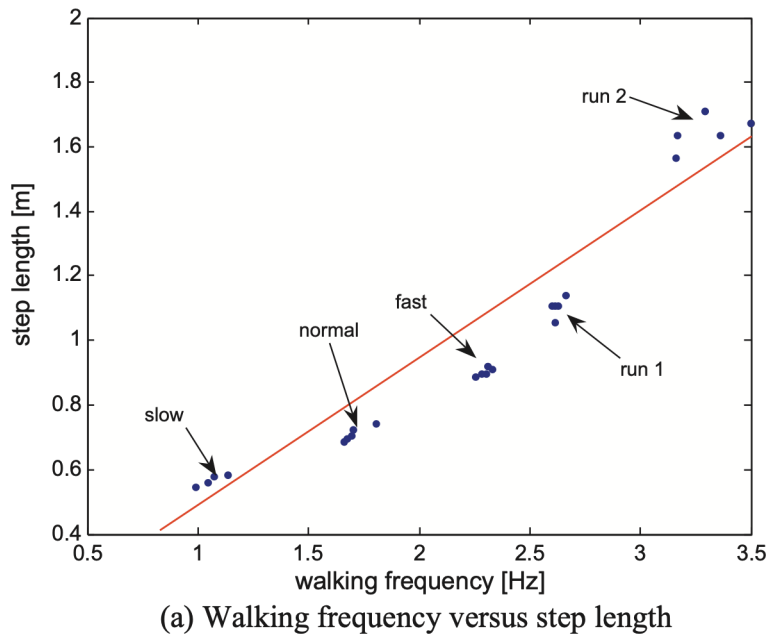
Funding:

BILL & MELINDA  
GATES foundation





# Current research / projects



- **Standardisation of wearable-based algorithms for healthcare applications in developing countries.**

- Principal investigator: Alan Godfrey (Northumbria University, Newcastle-upon-Tyne, UK).
- Team: Rodrigo Vitorio (UNESP), Marcos Barreto (UFBA), Azad Hussain (University of Birmingham), Clara Aranda-Jan (University College London)
- Period: 2018-2019
- Scope: This proposal aims to develop a novel standardised framework to better inform algorithms for a more harmonised gait assessment in Parkinson's disease (PD), particularly for developing countries where guidance is lacking. This project will lead to the design of an online simulation to test algorithms. Additionally, it will outline an educational process for all clinicians to better understand the functionality of wearables/algorithms and resulting outcomes. This will better guide PD assessment for sustainable health, promoting and encouraging low-cost wearables as routine diagnostics in developing countries. This framework will also be adapted to the needs of those in developed regions.

Funding:



Fig. 1. Schematic of the human gait cycle and the spatiotemporal parameters validated in this study. Specifically, we validated the IMU system's ability to measure the stance percent, swing percent, stride duration (gait cycle time), stride length, and step duration in addition to the speed and cadence of the cycle.

# Current research

- **Stratification of patients suffering from myalgic encephalomyelitis/chronic fatigue syndrome.**

- Principal investigator: Marcos Barreto (UFBA).
- Team: Nuno Sepulveda (London School of Hygiene & Tropical Medicine), Robespierre Pita (UFBA)
- Period: 2019-2020
- Scope: This study aims at to stratify ME/CFS patients into different clusters (or symptom subtypes). The respective objectives are the following: i) to distinguish ME/CFS patients from those suffering from multiple sclerosis (MS); ii) to identify sets of clinical symptoms that could characterize different clusters of ME/CFS patients; iii) to identify the best (or exclusive) predictive symptoms for CFS and compare the results with those obtained from different statistical/computational methods; (iv) to compare the stability of patients stratification using baseline and follow-up data.

Support:





**OBRIGADO!**

Contato:  
[marcosb@ufba.br](mailto:marcosb@ufba.br)  
[www.atyimolab.ufba.br](http://www.atyimolab.ufba.br)