

COURSE

“Técnicas Matemáticas para Big Data”

University of Aveiro
2025/2026

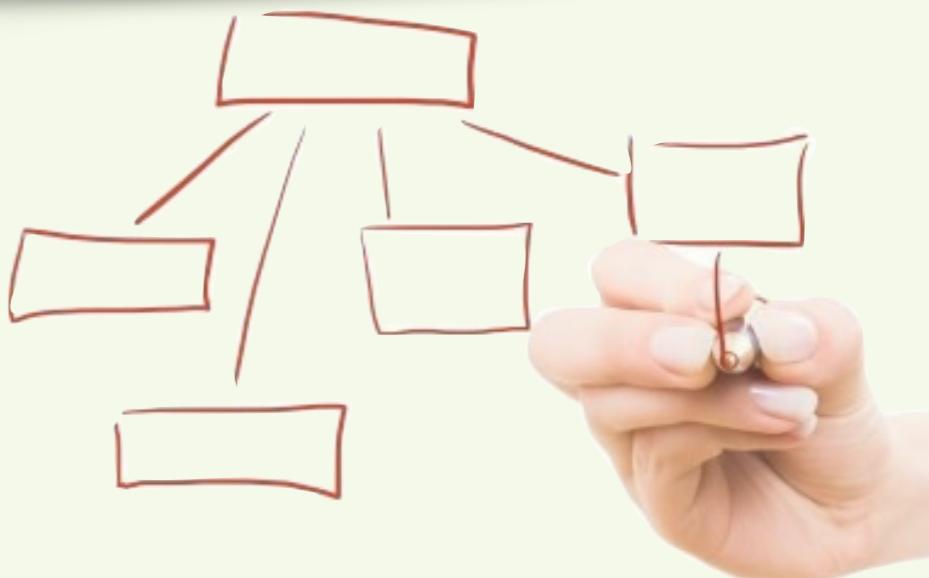


Algorithm 2.1: INSERTION-SORT(A)

```

1 for  $j \leftarrow 2$  to  $A.size$  do
2   key  $\leftarrow A[j]$ 
   // Insert  $A[j]$  into the sorted sequence  $A[1..j - 1]$ 
3    $i \leftarrow j - 1$ 
4   while  $i > 0$  and  $A[i] > key$  do
5      $A[i + 1] \leftarrow A[i]$ 
6      $i \leftarrow i - 1$ 
7    $A[i + 1] \leftarrow key$ 
```

Master in Data Science
Master in Mathematics and Applications



EXTRA INFO - this label will appear in slides with complementary information

Deletion

Insertion

SUMMARY: The course and Introductions to the main requirements and concepts.

[T] About the course — content, practical works, seminars, evaluations.

Introduction to big data: the four V's. Opportunities and challenges. Big data vs machine learning.

Mathematical preliminaries (e.g. Bachmann–Landau notation).

[P] General overview about tools for big data and machine learning: micro-services, IoT (digital twin), ETL, databases (SQL vs noSQL), dashboarding.

Micro-services deploy pipeline for TMBD.

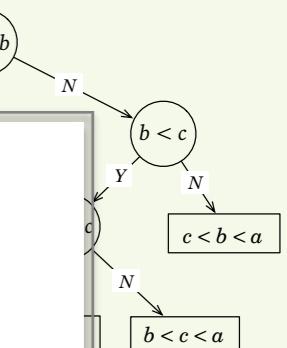
Algorithm 2.1: INSERTION

```

1  for j ← 2 to A.size do
2      key ← A[j]
      // Insert A[j] into A[0..j-1]
3      i ← j - 1
4      while i > 0 and A[i] > key do
5          A[i + 1] ← A[i]
6          i ← i - 1
7      A[i + 1] ← key

```

EXTRA INFO - this label will appear in slides with complementary information



mongoDB

elastic



Big Data vs Machine Learning ?
Any difference ?

Big Data vs Machine Learning ? Any difference ?

Goal: The goal is to turn data into information

Challenges: Capture, curation, time-limitations, storage, search, sharing, transfer, analysis, and visualization of the data.

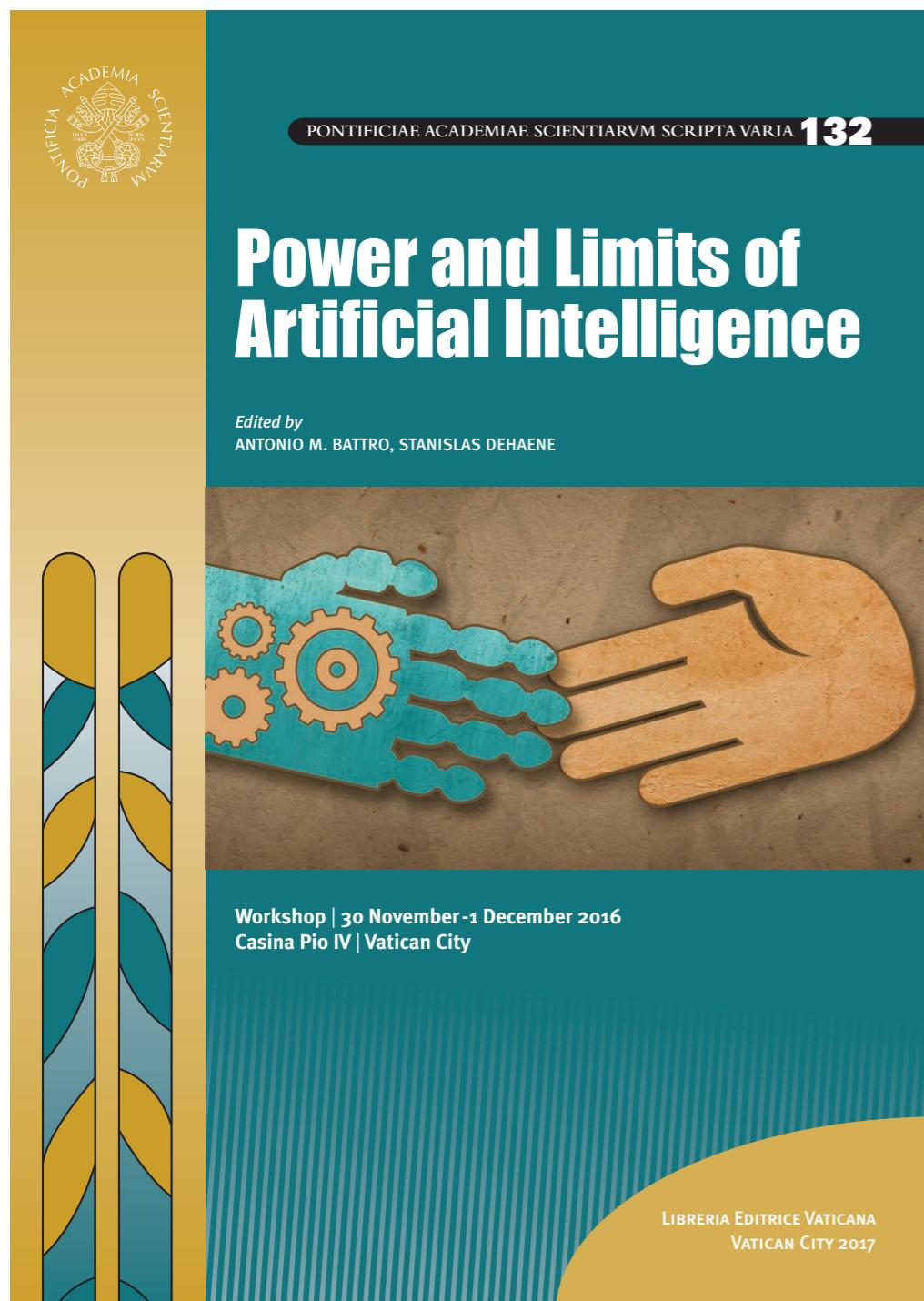
Data can be massive, non-static, multi-modal, incomplete, noisy, non-random, unstructured, dynamic, streaming, ...

Big Data vs Machine Learning ? Any difference ?

Mathematics ?

OR

Computer Science ?



* see dropbox folder

Contents

| | |
|----------------------------|----|
| Preface | 11 |
| Programme | 12 |
| List of Participants | 14 |

► STATE OF THE ART IN ARTIFICIAL INTELLIGENCE, ROBOTICS, BRAIN MODELING, BRAIN-COMPUTER INTERFACES

Artificial Intelligence – Big Achievements and Huge Questions Viewed from Mathematics

| | |
|---------------------|----|
| Cédric Villani..... | 19 |
|---------------------|----|

The Cerebral Cortex: An Evolutionary Breakthrough

| | |
|-------------------|----|
| Wolf Singer | 37 |
|-------------------|----|

Comments: The Ethics of Artificial Intelligence

| | |
|-----------------------|----|
| Stephen Hawking | 50 |
|-----------------------|----|

Optimal Strategies for Decision-Making and Their Neural Basis

| | |
|------------------------|----|
| Alexandre Pouget | 51 |
|------------------------|----|

Motivations and Drives Are Computationally Messy

| | |
|---------------------------------|----|
| Patricia Smith Churchland | 55 |
|---------------------------------|----|

Children and Robots

| | |
|--|----|
| Antonio M. Battro and Magela Fuzatti | 60 |
|--|----|

► PUTATIVE PREROGATIVES OF THE HUMAN BRAIN: EDUCATION, REASONING, CREATIVITY, CONSCIOUSNESS, SENSE OF SELF, ETHICS...COULD THEY BE CAPTURED IN MACHINES?

Ghost In the Machine

| | |
|-------------------|----|
| Olaf Blanke | 69 |
|-------------------|----|

What Is Consciousness, and Could Machines Have It?

| | |
|------------------------|----|
| Stanislas Dehaene..... | 75 |
|------------------------|----|

Big Data everywhere:

Lots of data is being collected and warehoused

- Web data (often user-provided)
- e-commerce, purchases at stores
- Medical data, health care
- Bank/Credit Card transactions
- Social Network
- Traffic, GPS, ...
- Scientific experiments
- ...

- CERN's Large Hydron Collider generates 15 PB a year
- The BRAIN initiatives produce terabytes of data a day
- The Large Synoptic Survey Telescope in Chile will collect 30TB per night. Headed by [Tony Tyson from UC Davis](#)



- YouTube contains 120 million videos and 72 hours of video uploaded every minute.
- Google processes 3.5 billion requests per day
- There is currently an estimate of 3.8 trillion photographs, 10% of them taken in the last year.
- Facebook has about 140 billion images with about 300 million new images a day.
- 2.5PB are flowing through Walmart's databases
- NYSE collects 1 TB each day.

Big Data does not just mean **massive amounts** of data

Big Data also means **complex data**

- Heterogeneous data
- Incomplete data
- Unstructured/semi-structured Data
- Graph Data
- Social Network, Semantic Web
- Streaming Data

Usual tasks for BDt (?!)

- Discovery of useful, possibly unexpected, patterns in data
- Non-trivial extraction of implicit, previously unknown and potentially useful information from data
- Finding outliers (security threat, credit card theft, ...)
- Clustering
- Classification
- Object recognition
- Visualization, dimension reduction
- “Data cleaning”: denoising, smoothing, grouping, ...
- Association Rule Mining (Costumers who buy X often buy Y, Costumer 123 likes product p10)
- Collaborative filtering: users collaborate in filtering information to find information of interest (Amazon, Netflix)

Usual tasks for BDt (?!)

The idea is 100 years old (see Karl Pearson), but its full potential will be unleashed only now.

Example:

In a recent analysis researchers developed a framework for comparing classifiers common in Machine Learning (Boosted decision trees, Random Forests, SVM, KNN, PAM and DLDA) based on a standard series of datasets.

Result: A simple (but mathematically rigorous) method gave better classification results across the data sets than the “glamorous” methods.

The dawning Age of Big Data will make it not just possible but very common (and perhaps necessary?) to validate methods via such meta data analyses.

*Understanding deeply the problem
leads to better results
than
knowing a lot of “fancy” algorithms*

“understand” = have a rigorous mathematical description

WHAT IS BIG DATA ?

WHAT IS MACHINE LEARNING ?

WHAT IS ... ?

Q5.1: What is the meaning of each concept and differences?

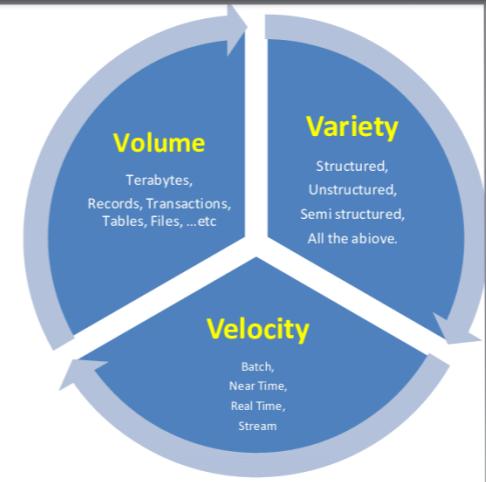
Your ideas: • ...



From Big Data to Big Artificial Intelligence?

Algorithmic Challenges and Opportunities of Big Data

Kristian Kersting¹ · Ulrich Meyer²



What is Big Data?

Some
"definitions"

Definition of big data

The ability to collect huge amounts of data (i.e., database scale) about many different phenomenon (Anonymous)

Large volumes of interrelated data with multiple modalities (Sriraam Natarajan, UT Dallas, USA)

“Big Data” is the body of technology—algorithms, programming systems, and hardware, that stress our abilities to handle the data. A related, and apparently more modern term is “Data Science”, which really means the same thing, but includes the application areas to which “big-data” technology is applied. What it DOESN'T mean is “AI” or “Machine Learning” or “Statistics done right”, as some have claimed (Jeffrey D. Ullman, Stanford, USA)

I define it by size but also by the fact that it is typically collected as a side-effect of some other process (Anonymous)

One or more of: petabytes, high velocity, high complexity—high velocity (more data coming at you faster) predominant (Sofus Macskassy, Branch Metrics, USA)

Large volumes of data from diverse sources that can be used to learn new facts, or predict future events (Anonymous)

When it is faster to transport it on tape than to transfer it via satellite (Anonymous)

Big Data is a collection of data instances which is significantly larger than the data that has existed/been used in a given area before (Anonymous)

Computational approaches to solving computational problems lie on a spectrum of model complexity. Big data is at one end of the spectrum, effectively attempting to solve problems without explicit models. I would therefore define Big Data as the set of methods and efforts that attempt to solve problems without the necessity of explanation or understanding through explicit models (Oliver Brock, TU Berlin, Germany)

Enough data so that only very basic/generic priors (such as in NNs or CNNs) are sufficient ensure generalization to test data (Marc Toussaint, MIT, USA)

Mainly high dimensions. In fact rich SVD spectrum (Nikolaos Vasiloglou, MLtrain, USA)

“Big” is a relative term. A couple of decades ago, datasets with a million elements were considered big. But now, one routinely comes across datasets with trillions of elements (Anonymous)

Opportunities of using Big Data

High robustness in statistical learning (Christian Bauckhage, U. Bonn, Germany)

Almost every field of human endeavor. I see lots of activity and promise in biology and biomedicine. Lots of commercial applications too, as companies like Google or Facebook use their data gathered from billions of people to understand much about human activity, from detecting spam, to knowing what I really mean when I misspell something. I don't see how this question can be answered in a few lines. My favorite example is providing driving directions based on real-time traffic estimates obtained from cell phone location data (Jeffrey D. Ullman, Stanford, USA)

Deeper holistic understanding of complex data (Sofus Macskassy, Branch Metrics, USA)

Quick solutions for problems that we do not understand but for which abundant data exists (Oliver Brock, TU Berlin, Germany)



Understanding and getting the data is the key to solving many problems. It may not sound very complex or challenging on the surface, however is the key to be successful in AI / ML initiatives (Anonymous)

Risks of using Big Data

Privacy concerns when dealing with personal data, particularly for healthcare. Inability to verify the correctness of learned models or discovered patterns when analyzing large datasets, particularly when using “black-box” models. Security issues, dealing with adversaries, information overload, etc. (Anonymous)

Hype and expectation (Sriraam Natarajan, UT Dallas, USA)

too high reliance on what is learned as truth (bad models, bad assumptions, not understanding the data). biased models in ways we do not understand (e.g., gender biases) (Sofus Macskassy, Branch Metrics, USA)

If the learning systems leveraging these Big Data sets are not designed carefully, they may end up codifying our biases and stereotypes (e.g., possible racial bias in automated airport profiling) and this may result in these biases getting even more deeply ingrained (Deepak Ajwani, Nokia Bell Labs, Ireland)

Frankly, I think the biggest risk is that governments, especially in Europe, will worry more about privacy than about the advantages that can come from exploiting “big data”. Privacy is a modern invention. 200 years ago, before the anonymity of cities was available, people didn't imagine that they could keep their lives secret from those around them (Jeffrey D. Ullman, Stanford, USA)

Opportunities of using Big Data

To solve low-level language and image perception tasks. Discover correlations for rare medical conditions (Guy Van den Broeck, UCLA, USA)

There are two very different instances of Big Data: the one makes modelling very easy, because the large data show underlying regularities (all models are wrong and we don't need them any more). This is a huge opportunity for analysis. The other are of the needle in the haystack type and make analysts' work harder (Anonymous)

The opportunity to analyse data streams in real-time moves machine learning into the direction of real-time control. This opens up tremendous applications.(Anonymous)

With the design of more scalable learning algorithms and better parallel hardware, vast amount of collected data can be leveraged for myriad applications, particularly in augmented intelligence (Deepak Ajwani, Nokia Bell Labs, Ireland)

Risks of using Big Data

That science in its current form ceases to exist because people just believe in the power of data but not in the power of understanding (Oliver Brock, TU Berlin, Germany)

How to verify analysis methods and results? How to reproduce results?
(Anonymous)

Algorithms that require large amounts of data will be less useful compared to algorithms that can work on smaller amounts of rare occurrence data (Anonymous)

To get 'stupid' machines that don't reason to yield an answer, that cannot learn from few data (e.g., in personal conversations) (Marc Toussaint, MIT, USA)

- Challenges and opportunities with big data
(Proceedings of the VLDB Endowment 5(12):2032-2033)
- For a Career Path in Big Data
(Edureka 2019 Tech Career Guide)
- Supply Chain Talent of the Future - Findings from the third annual supply chain
(Deloitte)

WHAT IS BIG DATA ?

WHAT IS MACHINE LEARNING ?

WHAT IS ... ?

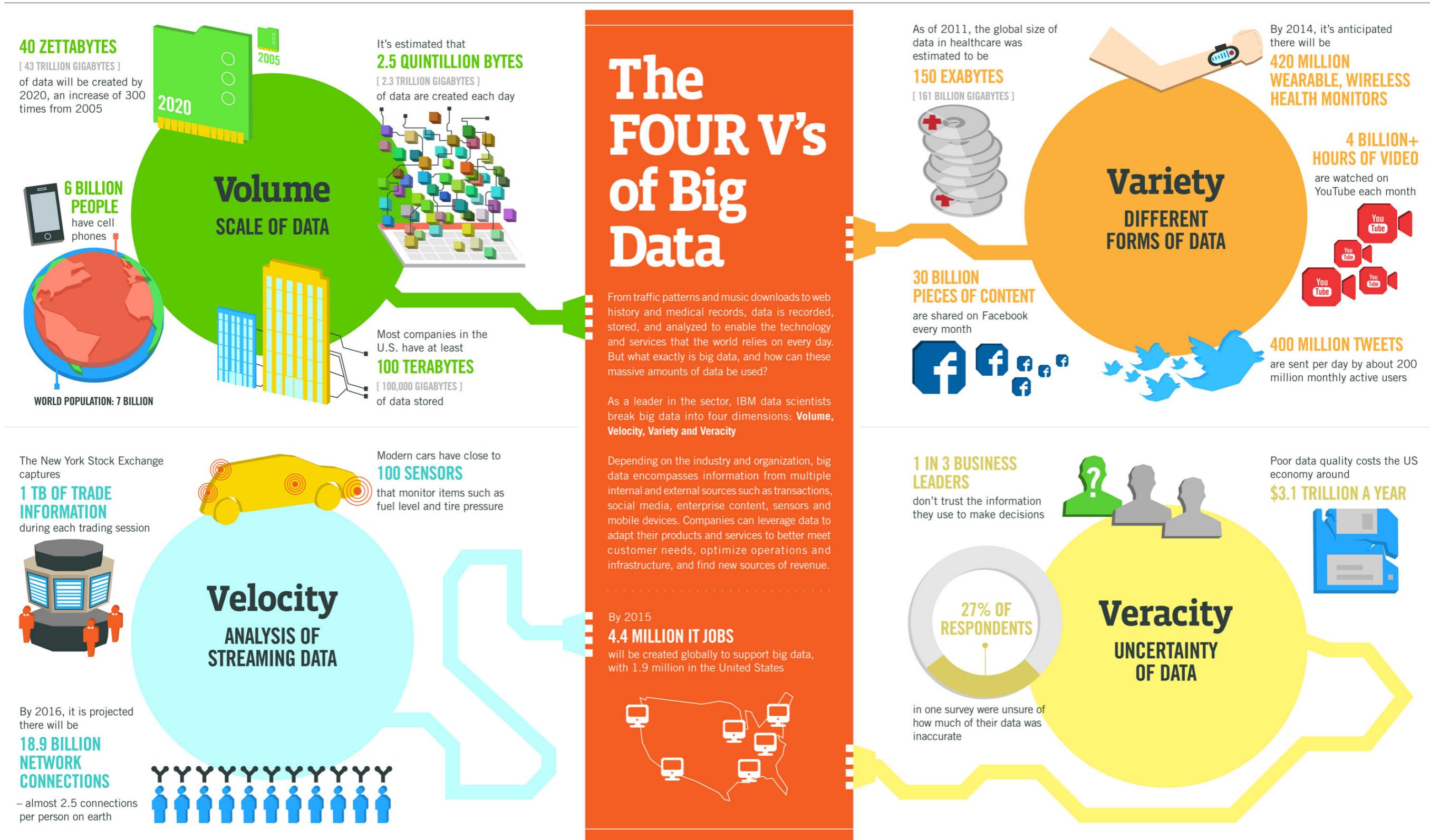
E.Rocha's understanding

BIG DATA: set of technologies and methodologies to (efficiently) deal with the transport, processing and storage of a (nowadays considered) big amount of information.



MACHINE LEARNING: a group of techniques and algorithms to model phenomena based on the information contained in datasets, rooted on mathematical models and optimization techniques, and implemented on well-known computer science frameworks.

What is Big Data?



Now it is common to find an additional V of Big Data: the Value.

IBM

QUESTION

Is AI, ML and BigData “good” or “bad” to humanity?



What is the oldest tool in the human history?

TOOLS FOR CUTTING



Oldest stone tools pre-date earliest humans

By Rebecca Morelle
Science Correspondent, BBC News

They were unearthed from the shores of Lake Turkana in Kenya, and date to 3.3 million years ago.

They are 700,000 years older than any tools found before, even pre-dating the earliest humans in the *Homo* genus.

The find, **reported in Nature**, suggests that more ancient species, such as *Australopithecus afarensis* or *Kenyanthropus platyops*, may have been more sophisticated than was thought.



2017 London Bridge Attack



London Bridge Attack: 7 Killed, 48 Wounded; 3 Suspects Shot Dead by Police



Convicted murderer hailed as 'hero' in London Bridge attack

Man convicted of killing woman, 21, joined others to stop terrorist who went on knife rampage

Dec 02, 2019 06:00 am



LONDON: One of the men who helped to stop the London Bridge terrorist from killing more people was once a killer himself.



A RECENT (NON)STANDARD EXAMPLE



Fresh Cambridge Analytica leak ‘shows global manipulation is out of control’

Company’s work in 68 countries laid bare with release of more than 100,000 documents

Carole Cadwalladr

• @carolecadwalla

Sat 4 Jan 2020 16.55 GMT



42,334

Q5.2: What is the biggest intrusion to your privacy?

Q5.3: Can AI predict if a couple breaks after 4 years?

Q5.4: Can AI model the couple interaction dynamics?

Q5.5: Can AI be used for choosing innovation successfully ?

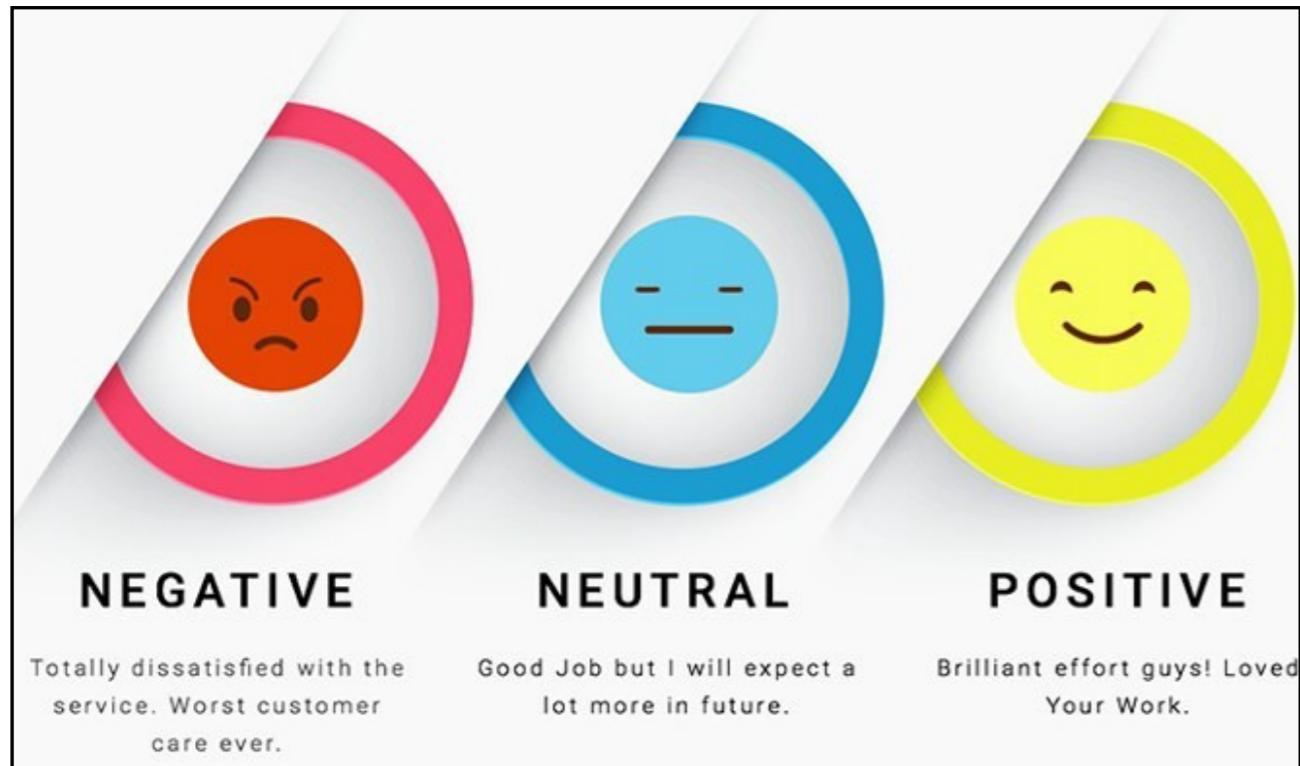
How human sentiments may improve the economy but destroy creativity ? Predicting crowdfunding campaigns...



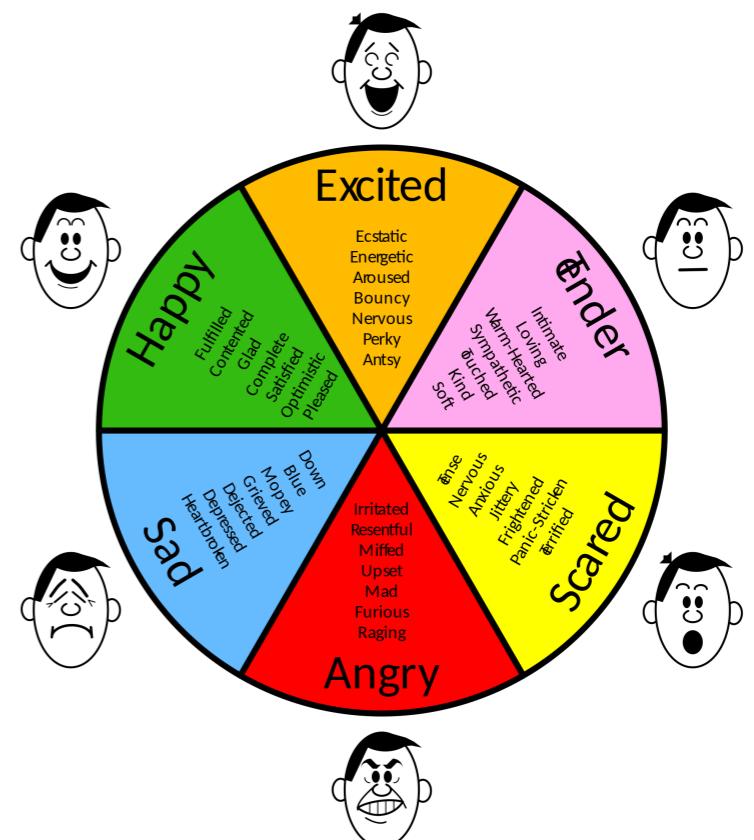
Definição (António Damásio):

- **EMOÇÃO:** conjunto de reacções corporais, automáticas e inconscientes, face a determinados estímulos provenientes do meio onde estamos inseridos;
- **SENTIMENTO:** surge quando tomamos consciência das nossas emoções, isto é, o sentimento dá-se quando as nossas emoções são transferidas para determinadas zonas do nosso cérebro.

Análise de Sentimentos



Análise de Emoções





Análise de Sentimentos (opinion mining)

— Positivo, neutro ou negativo?



- unbelievably disappointing
- Full of zany characters and richly applied satire, and some great plot twists
- this is the greatest screwball comedy ever filmed
- It was pathetic. The worst part about it was the boxing scenes.

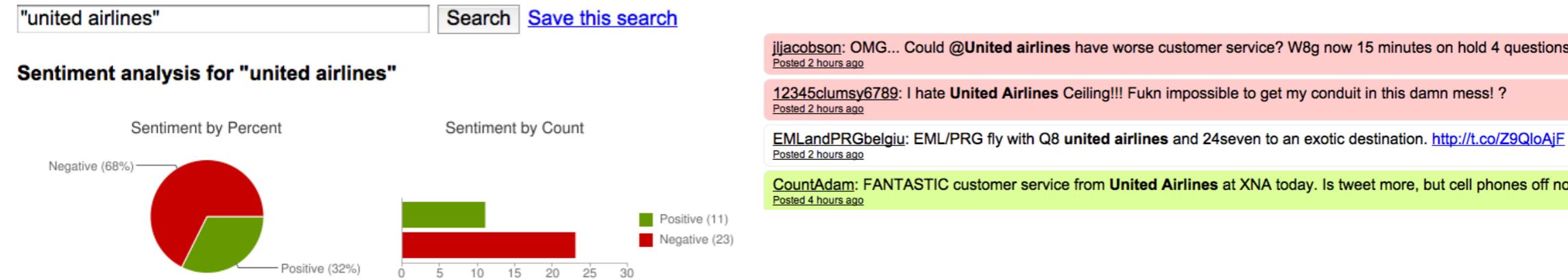
— Várias fontes:

- Google product search
- Bing shopping
- ★ Twitter
- Etc.

Sentiment Analysis (NLP) on Twitter

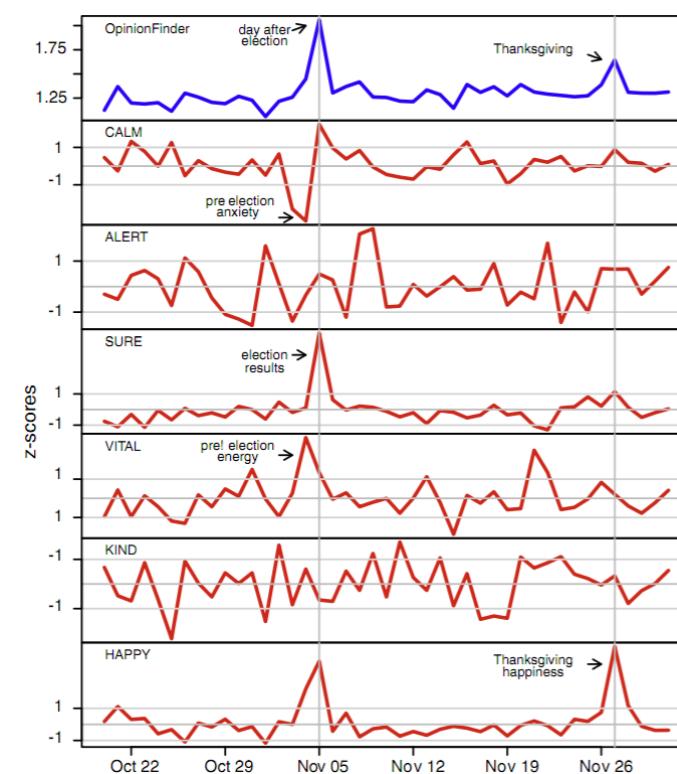
Alec Go, Richa Bhayani, Lei Huang. 2009. Twitter Sentiment Classification using Distant Supervision

Type in a word and we'll highlight the good and the bad

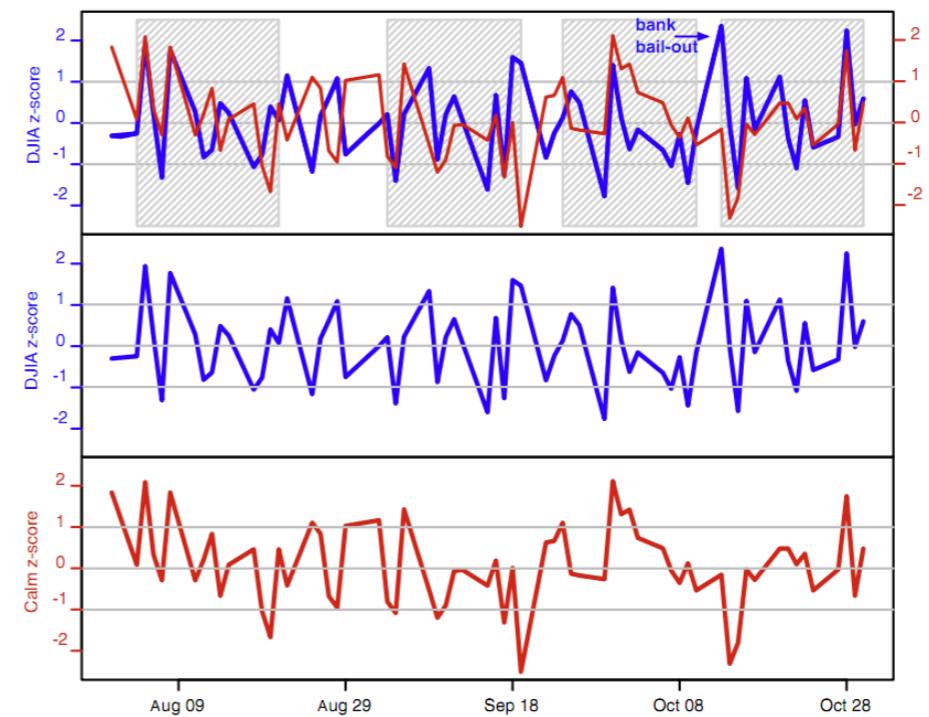


Can we predict Stock Market from Twitter information?

Johan Bollen, Huina Mao, Xiaojun Zeng. 2011. Twitter mood predicts the stock market, Journal of Computational Science 2:1, 1-8. 10.1016/j.jocs.2010.12.007.



FINDING: CALM predicts DJIA 3 days later



An Interpretable Prediction Model based on Campaigns Aspects Emotions

19

work with a GERMAN PRIVATE COMPANY

Commercial Web App

Campaign Goal (\$):

Last Seasons:

EXTRA INFO

20000

3

Show all clusters

Campaign Description:

* For those new to Kickstarter, scroll down to read the FAQ's below. Basically, you pledge a \$ amount to this campaign and only get charged if the project is fully funded. You will receive a Kickstarter survey at the end of the campaign to select reward choices. Watch the Fungisaur's animated origin story below! Hope you enjoyed that as much as we enjoyed making it. This is our first step towards

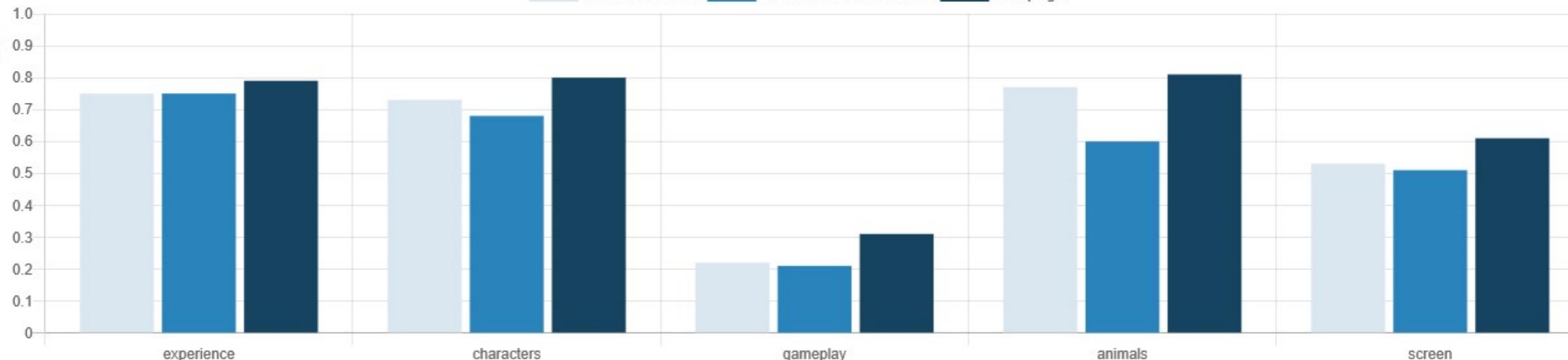
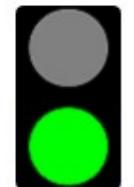
Start

Campaign
Prediction:

Scoring emotional intensity on Clusters.
Benchmarked on campaigns of the last 3 seasons.

Last Seasons All Last Successful Seasons Campaign

Success



| Cluster | experience (0,79) | characters (0,8) | gameplay (0,31) | animals (0,81) | screen (0,61) |
|---|--------------------------|--------------------------|--------------------------|----------------------------|-------------------------------|
| Feature Cluster Need Index Benchmarking | Optimization Opportunity | Optimization Opportunity | Major revisions required | Optimization Opportunity | Minor revisions required |
| Trend Need Index (All) | 0,75 ↓ | 0,73 ↓ | 0,22 ↓ | 0,77 ↓ | 0,53 ↓ |
| Trend Need Index (Successful) | 0,75 ↓ | 0,68 ↓ | 0,21 ↓ | 0,6 ↓ | 0,51 ↓ |
| Lead Mood Season Trend (All) / Lead Mood Season Trend (Successful) / Campaign Lead Mood | Joy / Joy / Joy | Joy / Joy / Joy | Joy / Admiration / Joy | Joy / Disgust / Admiration | Joy / Admiration / Admiration |
| Domain Aspect hits - Needindex | experience: 0,76 | character: 0,73 | adventure: 0,72 | cat: 0,6 | screen: 0,74 |

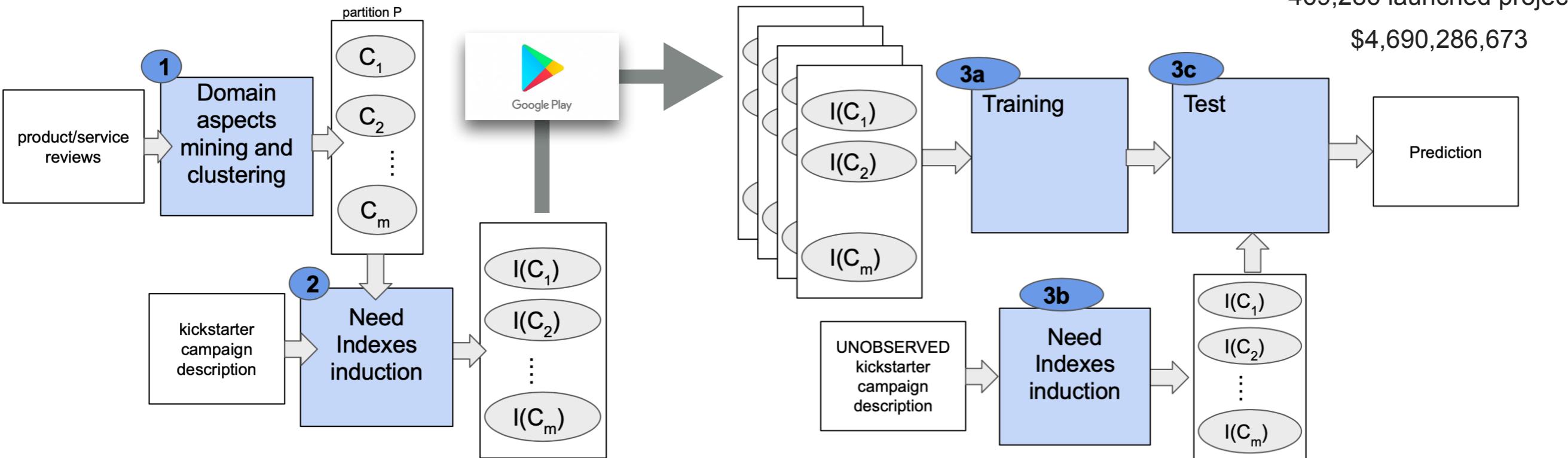
An Interpretable Prediction Model based on Campaigns Aspects Emotions

work with S.Rittinghaus + N.Samsami (GERMANY) and S.Faralli + D.Distante (ITALY)

Kickstarter

469,286 launched projects

\$4,690,286,673



State-of-Art (on 16/01/2020)

To the best of our knowledge, by the art prediction accuracy on the basis of static Kickstarter data is

76,4%

NLP + ML
(python)

NeedIndex

| Unsupervised: $P_MG_2009-2019$ | | | | |
|----------------------------------|------|------|------|------|
| system | P | R | F1 | A |
| MultiLayerPerceptron | 0.75 | 0.87 | 0.81 | 0.87 |
| SupportVectorMachines | 0.75 | 0.87 | 0.81 | 0.87 |
| RandomForestClassifier | 0.78 | 0.85 | 0.80 | 0.85 |
| GradientBoostingClassifier | 0.75 | 0.87 | 0.81 | 0.87 |

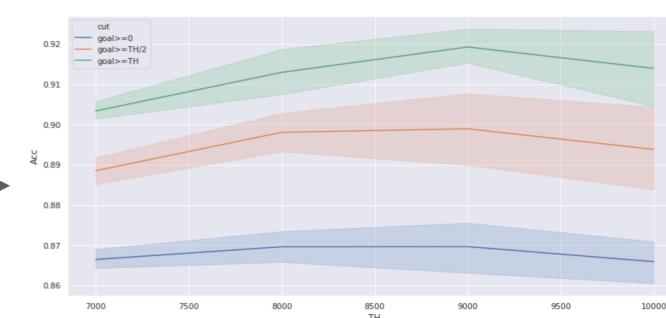
| Supervised: $P^*_MG_2009-2019$ | | | | |
|----------------------------------|------|------|------|------|
| system | P | R | F1 | A |
| MultiLayerPerceptron | 0.81 | 0.90 | 0.85 | 0.90 |
| SupportVectorMachines | 0.81 | 0.90 | 0.85 | 0.90 |
| RandomForestClassifier | 0.81 | 0.89 | 0.85 | 0.89 |
| GradientBoostingClassifier | 0.81 | 0.90 | 0.85 | 0.90 |

| baseline | | | | |
|----------|------|------|------|------|
| system | P | R | F1 | A |
| Random | 0.77 | 0.50 | 0.58 | 0.51 |

recreate all
the pipeline
in **Spark**

Ideas from the
theoretical part
about PDEs

Reproducing Kernel
Banach Spaces
Convexity



Our final accuracy
94%

Best cut: **9000\$**

QUESTION

ML/AI algorithm is (in)parcial? Avoid social discrimination?!





Joy Buolamwini /The Algorithmic Justice League at MIT Media Lab, Joy Buolamwini /The Algorithmic Justice League, 2018/2018,
From the collection of: Barbican Centre

Joy Buolamwini

iconic women, it is time to re-examine how these systems are built and who they truly serve.

machines are neutral, but they aren't. My research uncovered racial bias in AI systems sold by tech giants like IBM, Amazon. Given the task of guessing the gender of a face, all had substantially better on male faces than female faces. The had error rates of no more than 1% for lighter-skinned kinned women, the errors soared to 35%. AI systems from s have failed to correctly classify the faces of Oprah Winfrey, and Serena Williams. When technology denigrates even these

Oprah Winfrey

facial analysis
a complete analysis of facial attributes, including confidence scores.

TYLER

Results

looks like a face 99.9 %

appears to be male 76.5 %

age range 26 - 43 years old

* Request

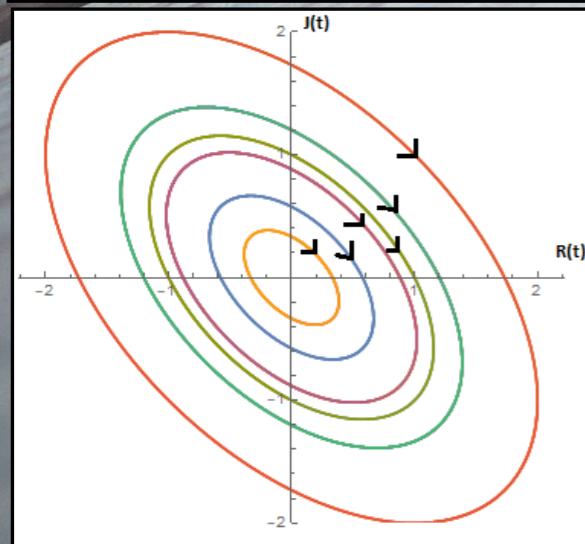
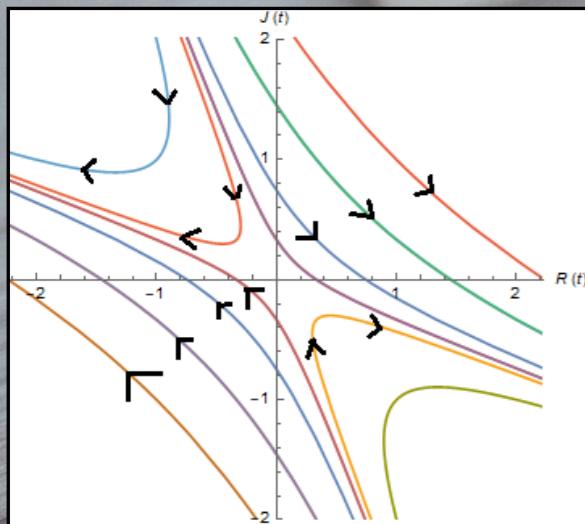
QUESTION

Can we model “love” or “friendship” ?



QUESTION

Can we model “love” or “friendship” ?



There are several mathematical models in the literature, e.g.

$$\dot{R}(t) = \frac{dR(t)}{dt} = aR(t) + bJ(t)$$

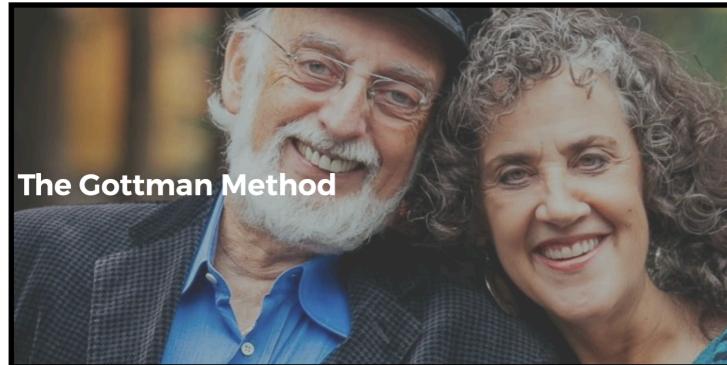
$$\dot{J}(t) = \frac{dJ(t)}{dt} = cR(t) + dJ(t)$$

BG/ML techniques can be used to find parameters a, b, c, d



Can we predict if a couple is going to break after 4 years of relationship?

Trabalho do Gottman Institute (Washington)



Goals and Principles of the Gottman Method

The goals of Gottman Method Couples Therapy are to disarm conflicting verbal communication, increase intimacy, respect, and affection, remove barriers that create a feeling of stagnancy in conflicting situations, and create a heightened sense of empathy and understanding within the context of the relationship.

In our laboratory's multi-method research on marital interaction we have shown, in four longitudinal studies, that we can predict with over 90% accuracy whether a couple will divorce or stay married, and their marital satisfaction if they do stay married (Gottman, 1994; 1999).

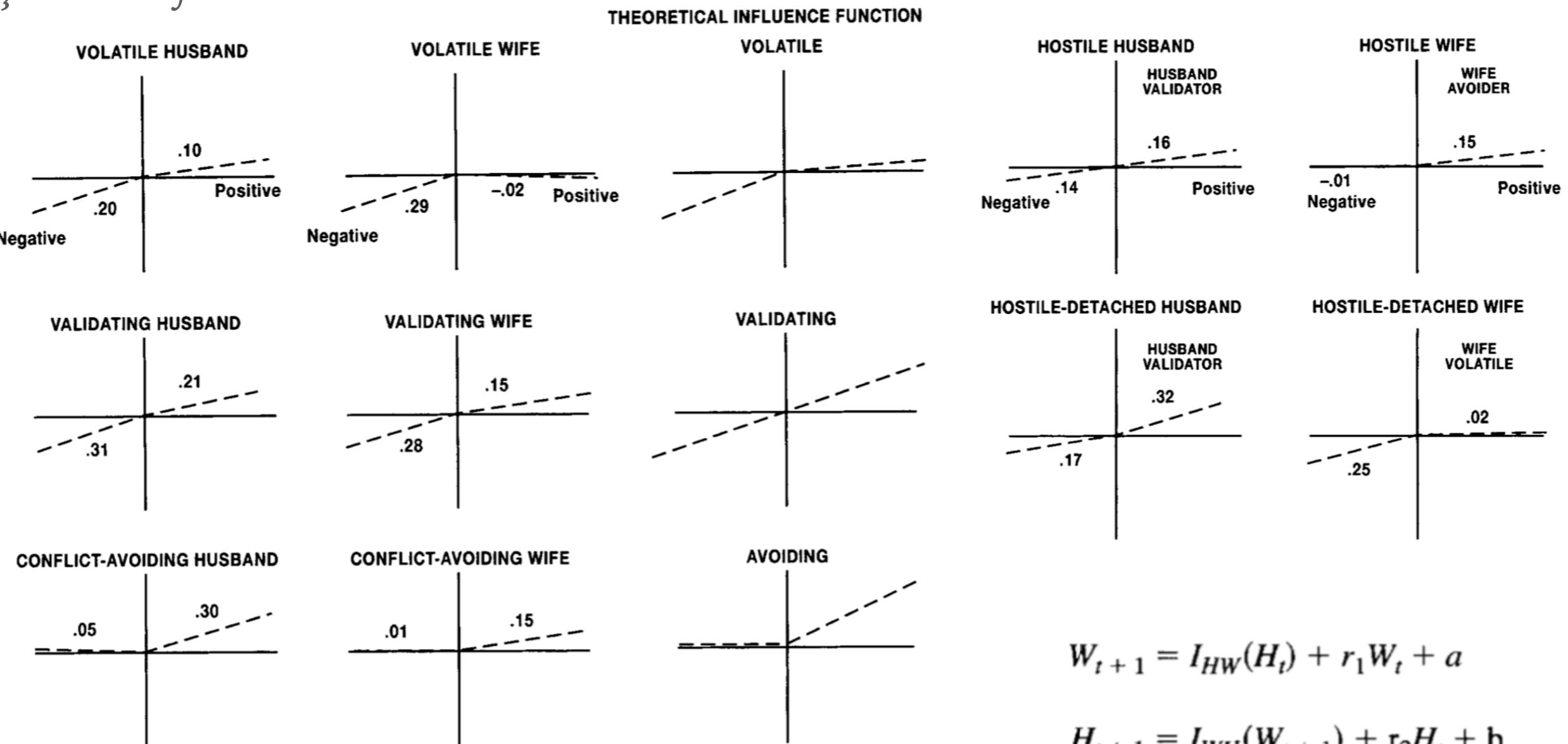
> 90% !!!

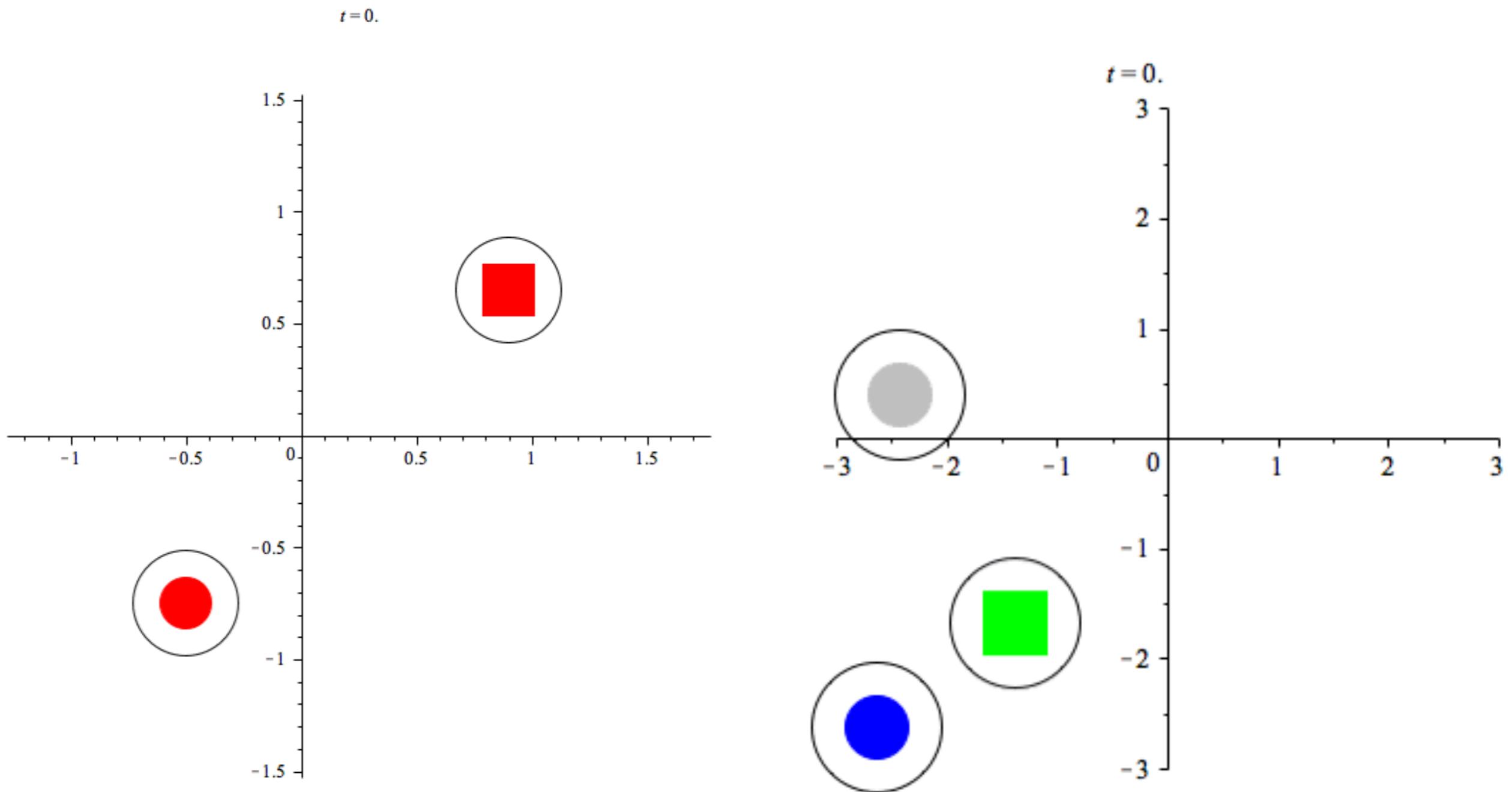
Trabalho do Gottman Institute (Washington)

In our laboratory's multi-method research on marital interaction we have shown, in four longitudinal studies, that we can predict with over 90% accuracy whether a couple will divorce or stay married, and their marital satisfaction if they do stay married (Gottman, 1994; 1999).

> 90% !!!

Funções de influência

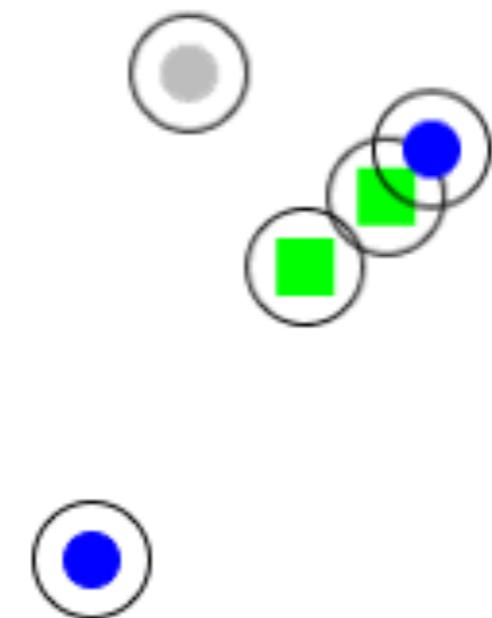


Example 1**LEGENDA:**

verde (volatile), **azul** (validating), **cinzento** (conflict-avoiding), **vermelho** (hostile), **laranja** (hostile-deattached)
bola (woman), quadrado (man)

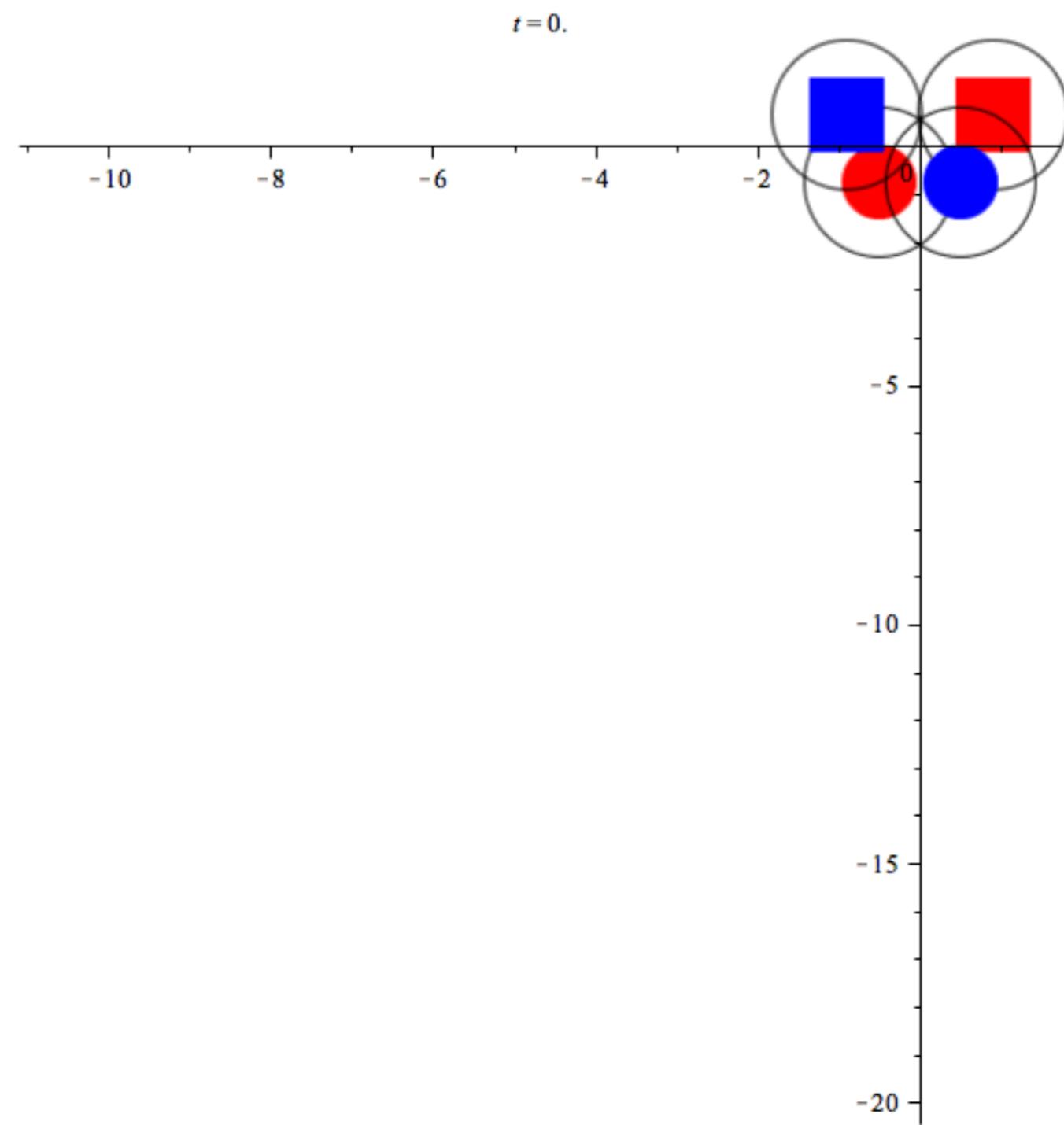
Example 2

$t = 0.$



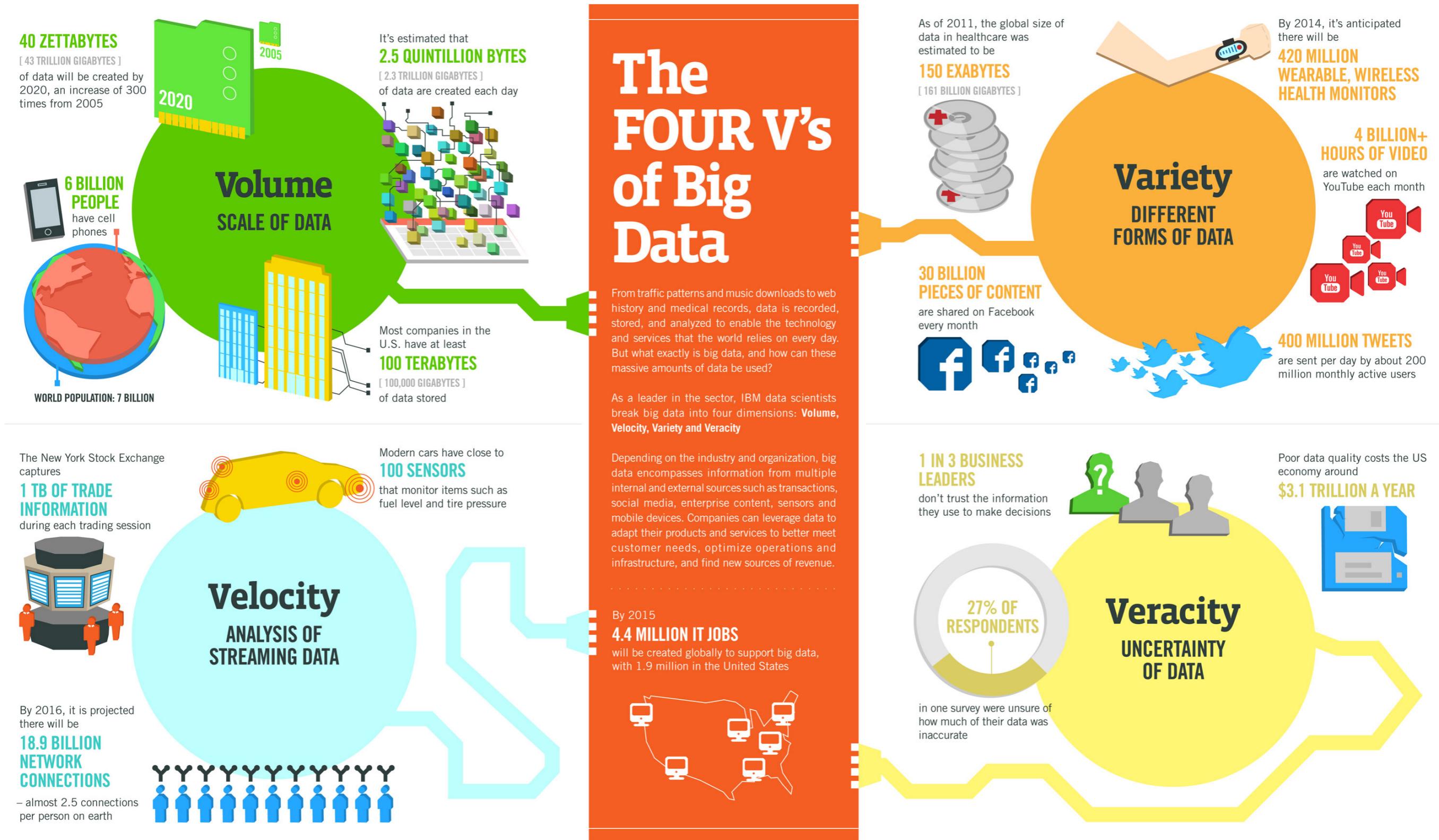
LEGENDA:

verde (volatile), **azul** (validating), **cinzento** (conflict-avoiding), **vermelho** (hostile), **laranja** (hostile-deattached)
bola (woman), quadrado (man)

Example 3**LEGENDA:**

verde (volatile), **azul** (validating), **cinzento** (conflict-avoiding), **vermelho** (hostile), **laranja** (hostile-deattached)
bola (woman), quadrado (man)

What is Big Data?

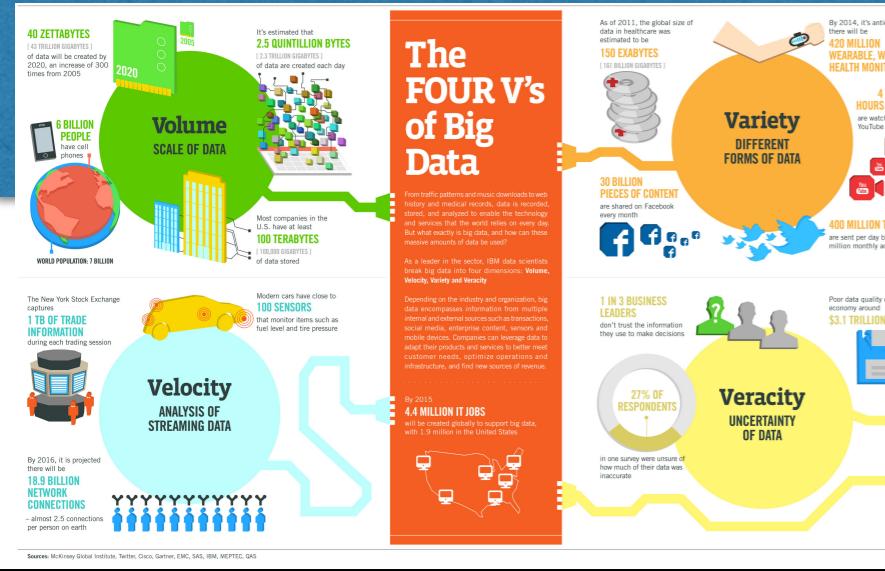


Sources: McKinsey Global Institute, Twitter, Cisco, Gartner, EMC, SAS, IBM, MEPTEC, QAS



What is Big Data?

“Técnicas Matemáticas para Big Data”



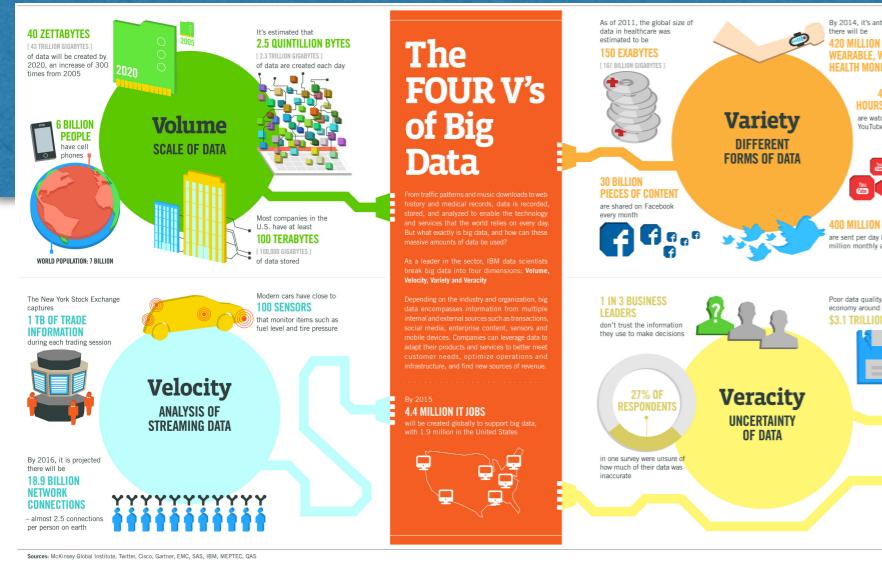
Q: Difficulties?

• Low Veracity

- Veracity, validity, ..., falsification, ethics
- Discrete/Continuous Markov Chains
- Bayesian Networks
- Fuzzy Logics

What is Big Data?

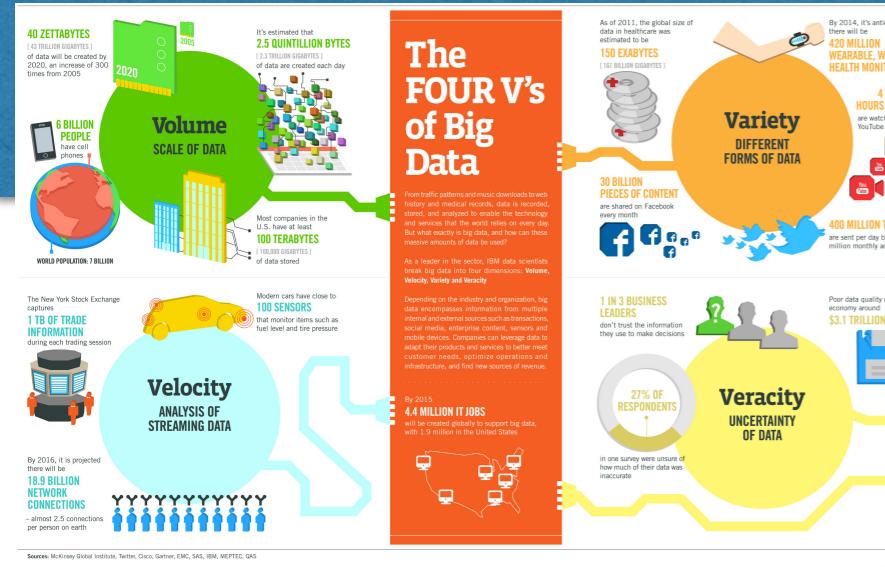
“Técnicas Matemáticas para Big Data”



- Low Veracity
 - Veracity, validity, ..., falsification, ethics
 - Discrete/Continuous Markov Chains
 - Bayesian Networks
 - Fuzzy Logics
- **High Velocity**
 - About data streams (DS)
 - Statistical algorithms for DS

Q: Difficulties?

What is Big Data?



“Técnicas Matemáticas para Big Data”

- Low Veracity
 - Veracity, validity, ..., falsification, ethics
 - Discrete/Continuous Markov Chains
 - Bayesian Networks
 - Fuzzy Logics
- **High Velocity**
 - About data streams (DS)
 - Statistical algorithms for DS

Q: Difficulties?

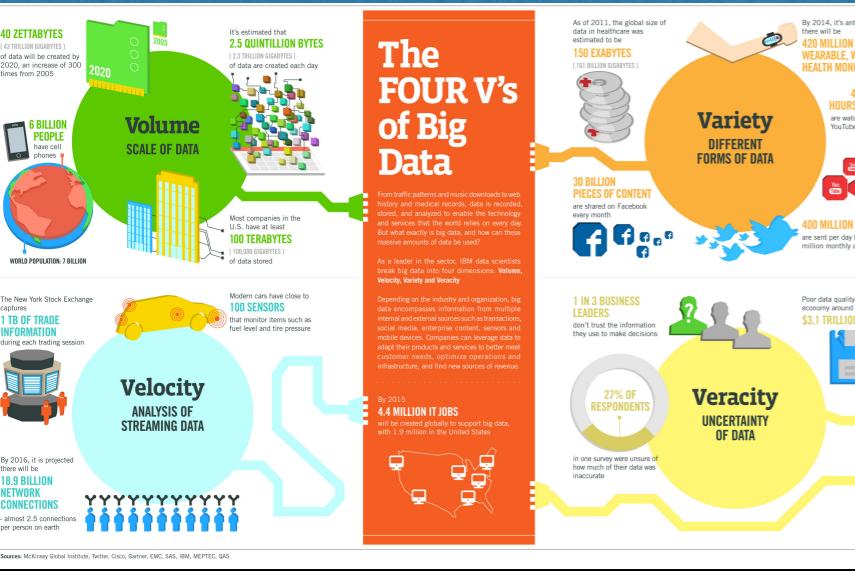
Problem 1: Calculate the mean and standard deviation without keeping the values of the previous temperatures.

$$W = \frac{\sum_{i=1}^n w_i X_i}{\sum_{i=1}^n w_i}$$

$$\sigma = \sqrt{\frac{\sum(x_i - \mu)^2}{N}}$$

What is Big Data?

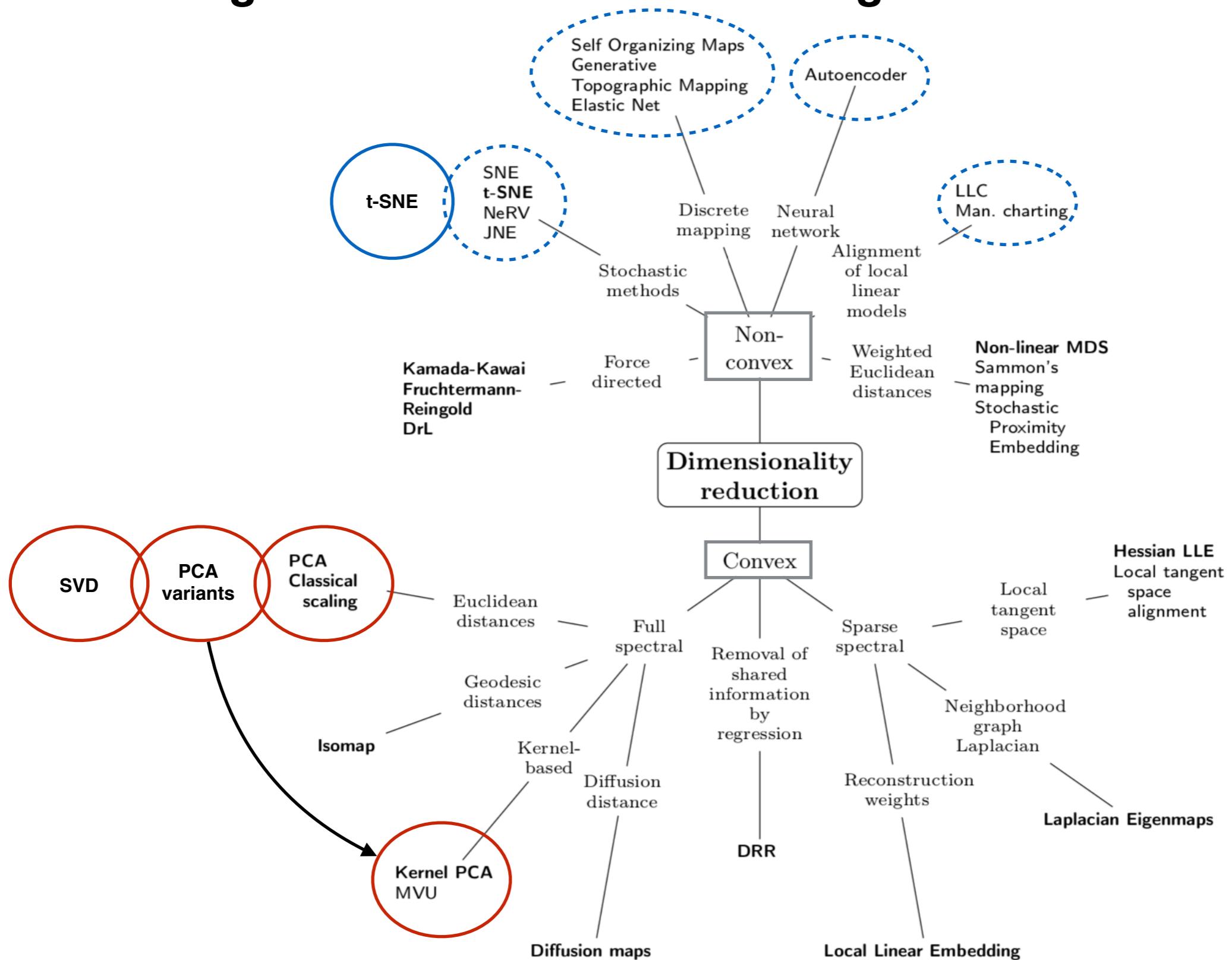
“Técnicas Matemáticas para Big Data”



- Low Veracity
 - Veracity, validity, ..., falsification, ethics
 - Discrete/Continuous Markov Chains
 - Bayesian Networks
 - Fuzzy Logics
- High Velocity
 - About data streams (DS)
 - Statistical algorithms for DS
- **Large Volume**
 - Dimension reduction methods

Q: Difficulties?

Brief Catalog of Dimension Reduction Algorithms



see notes/N01-Higher_Dimensions.pdf

1 Surprises in high dimensions

Our intuition about space is based on two and three dimensions and can often be misleading in high dimensions. It is instructive to analyze the shape and properties of some basic geometric forms, which we understand very well in dimensions two and three, in high dimensions. To that end, we will look at the sphere and the cube as their dimension increases.

1.1 Geometry of the d -dimensional Sphere

Consider the unit sphere in d dimensions. Its volume is given by

$$V(d) = \frac{\pi^{\frac{d}{2}}}{\frac{d}{2} \Gamma\left(\frac{d}{2}\right)}$$

where Γ is the Gamma function. Recall that for positive integers n , $\Gamma(n) = (n - 1)!$. Using Stirling's Formula,

$$\Gamma(n) \sim \sqrt{\frac{2\pi}{n}} \left(\frac{n}{e}\right)^n$$

we can see $\Gamma\left(\frac{d}{2}\right)$ grows much faster than $\pi^{\frac{d}{2}}$, and hence

$$V(d) \rightarrow 0 \quad \text{as} \quad d \rightarrow \infty.$$

Volume of a d-dimension unitary sphere

Consider the unit sphere is d dimensions. Its volume is given by

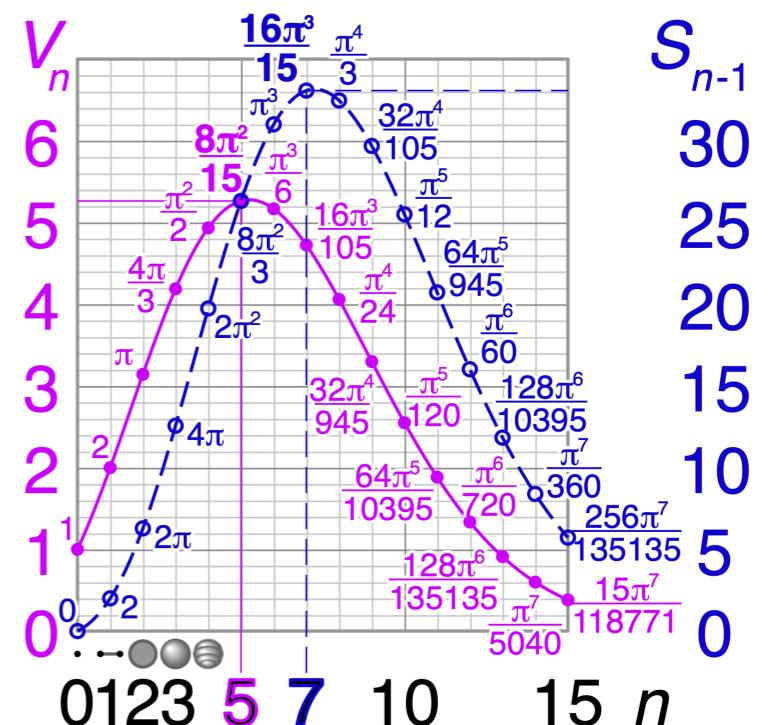
$$V(d) = \frac{\pi^{\frac{d}{2}}}{\frac{d}{2} \Gamma(\frac{d}{2})}$$

where Γ is the Gamma function. Recall that for positive integers n , $\Gamma(n) = (n - 1)!$. Using Stirling's Formula,

$$\Gamma(n) \sim \sqrt{\frac{2\pi}{n}} \left(\frac{n}{e}\right)^n$$

we can see $\Gamma(\frac{d}{2})$ grows much faster than $\pi^{\frac{d}{2}}$, and hence

$$V(d) \rightarrow 0 \quad \text{as} \quad d \rightarrow \infty.$$



Volume of a d-dimension unitary cube

Claim: Most of the volume of the high-dimensional cube is located in its corners.

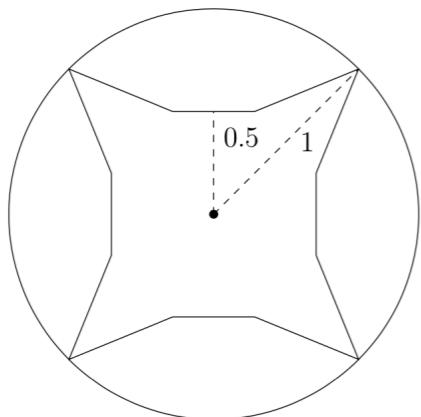


Figure 4: Projections of the 4-dimensional unit sphere and unit cube, centered at the origin (4 of the 16 vertices of the hypercube are shown).

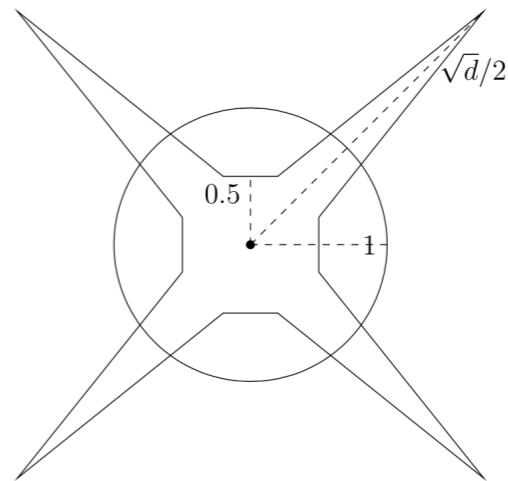
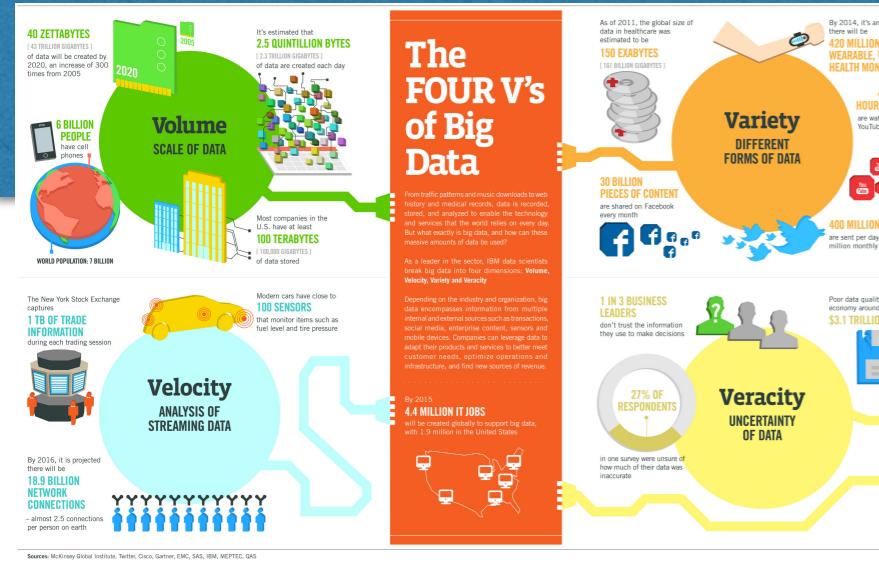


Figure 5: Projections of the d -dimensional unit sphere and unit cube, centered at the origin (4 of the 2^d vertices of the hypercube are shown).

What is Big Data?

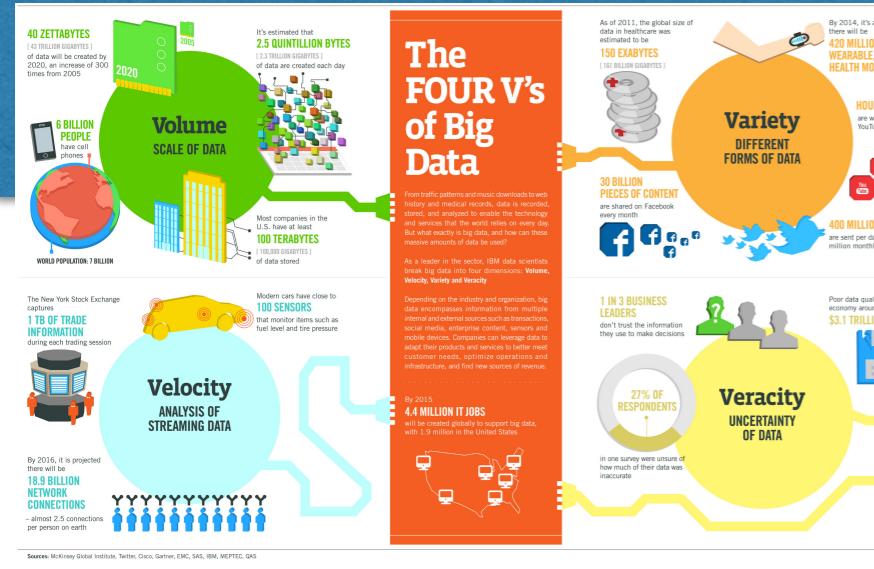


“Técnicas Matemáticas para Big Data”

- Low Veracity
 - Veracity, validity, ..., falsification, ethics
 - Discrete/Continuous Markov Chains
 - Bayesian Networks
 - Fuzzy Logics
- High Velocity
 - About data streams (DS)
 - Probabilistic algorithms for DS
- Large Volume
 - Dimension reduction methods
- **High Variety**
 - PEGs for documents classification
 - Introduction to NLP

Q: Difficulties?

Big Data tools via Docker



HOMEWORK (Basics for Tools Deployment)

YAML (superset of JSON):

<https://www.cloudbees.com/blog/yaml-tutorial-everything-you-need-get-started>

Docker Install (windows):

<https://medium.com/bina-nusantara-it-division/tutorial-for-docker-on-windows-c9af1162f0f2>

Docker Tutorial:

<https://docs.docker.com/get-started/>