# Intelligent Decision Making in the Era of Semantic Web and Big Data

## António Grilo

Faculdade de Ciências e Tecnologia da Universidade Nova de Lisboa

& UNIDEMI

acbg@fct.unl.pt

# Agenda

- Living in the Era of Big Numbers

- The Concept: Web Competitive Intelligence

# Living in the Era of Big Numbers
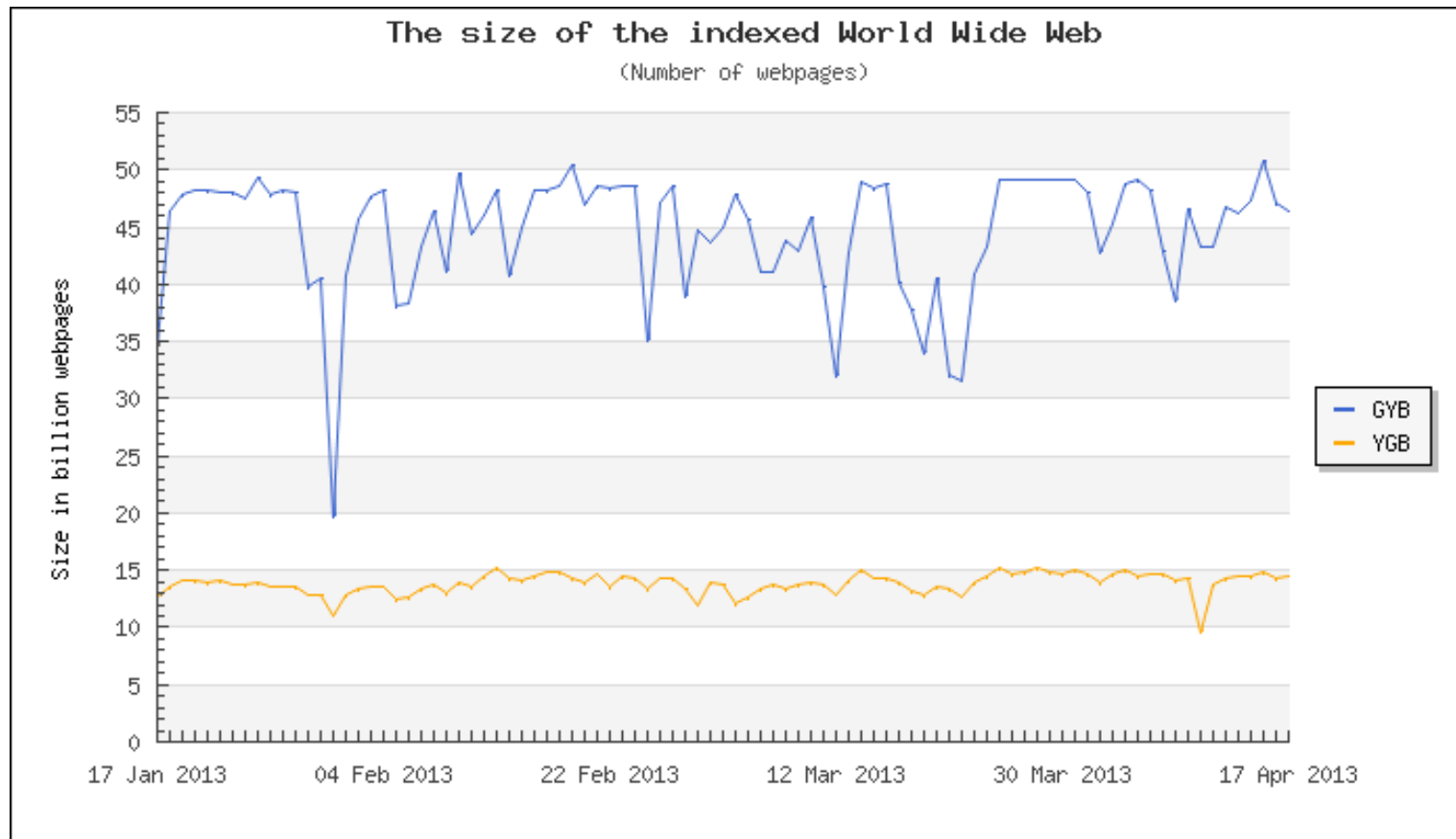
# The Internet has 50 Billion Webpages



The size of the indexed World Wide Web
(Number of webpages)

**Source:** http://www.worldwidewebsize.com/

# B2C E-Commerce of 1 Trillion US Dollars

**Top 5 Countries, Ranked by B2C Ecommerce Sales, 2011-2013**
*billions*

**1. US\***
- $301.69
- $343.43
- $384.80

**2. China\*\***
- $56.69
- $110.04
- $181.62

**3. UK**
- $109.03
- $124.76
- $141.53

**4. Japan**
- $112.78
- $127.82
- $140.35

**5. Germany**
- $38.08
- $47.00
- $53.00

■ 2011     ■ 2012     ■ 2013

**Source:** www.emarketer.com - Ecommerce Sales Topped $1 Trillion for First Time in 2012

# Web Analytics are Common

# Sophisticated Web Analytics

# The Emergence of Social Networks

*"**Social media is a group of platforms and tools** that users employ to share information, photos, videos, and other contents."*

*(Turban and Lai, 2011)*

*"Social technologies are products and services that enable **social interactions in the digital realm and provide distributed rights to communicate and add, modify, or consume content**. They include social media, Web 2.0, and enterprise collaboration technologies."*

*(McKinsey Quarterly, November 2012)*

# Facebook Users in 2013

1.11 Billion people using the site each month (+23 % 2012)

665 million active users each day on average

751 million from a mobile device each month, (+54 % 2012).

(…)

1 million users by the end of 2004.

# Facebook Users in 2013

There are more than 1.5 Billion social network users…

…which make up to 80% of total internet users.

70% of companies use social networks

90% recognize the benefit

Workers spend 28 hours every week writing emails, searching information, and collaborating internally.

# Social Network Platforms



| Platform | Percentage |
|----------|-----------|
| Facebook | 92% |
| Twitter | 82% |
| LinkedIn | 73% |
| Blogs | 61% |
| YouTube | 57% |
| Google+ | 40% |

**Source:** Social marketing industry report: How marketers are using social to grow their business, Stelzner MA, 2012

# Social Media Analytics

# Social Media Analytics - Targeting



28/06/13        KES - IDT/IIMSS 2013

# Social Media Analytics – Sentiment Analysis

# Social Media Analytics – Sentiment Analysis

Dionnova: @DJFreshSA I-phone; wen u type pepsi it corrects to capital letter P, but type coke or **coca cola** it d
Posted: 21 seconds ago

MaxaEnPointe: That unofficial **Coca Cola** Zim ad is soooo messed up it's funny. I wasn't offended though! http://
Posted: 7 minutes ago

rimonator: @Econsultancy BRILLIANT how **coca cola** use their social media
Posted: 8 minutes ago

YetNaive: I love **coca cola** so much. It's like the only dark soda I actually enjoy.
Posted: 10 minutes ago

MarthaaMor_77: **Coca cola** cherry is a heavenly creation
Posted: 10 minutes ago

aiR_La: RT @khuul_khidd: ?@dRealest_felix: **Coca cola** RT @Questionnier: A drink you're addicted to? #QnA
Posted: 14 minutes ago

Nomusa_DM: **Coca cola** RT"@Questionnier: A drink you're addicted to? #QnA"
Posted: 16 minutes ago

coca_cola_fan: **#CocaCola** reveals plans to sell its drinks in "greener" plastic bottles in China. http://t.co/SSQh4
Posted: 19 minutes ago

# Social Media Analytics – Flu & Twitter

# Social Media Analytics – Who's Important

# Semantic Web

The Semantic Web is the extension of the World Wide Web that enables people to share content beyond the boundaries of applications and website. **Semantic Web is a web that is able to describe things in a way that computers applications can understand**.

The Semantic Web describes the **relationships between things** (like A is a part of B and Y is a member of  Z) and the **properties of things** (like size, weight, age, and price)

"*If HTML and the Web made all the online documents look like one huge* **book***, RDF, schema, and inference languages will make all the data in the world look like one huge* **database**"

*Tim Berners-Lee, Weaving the Web, 1999*

# Semantic Web for Dummies…

Source: www.semanticwebexplained.co.uk

# Semantic Web:  No Big Numbers Yet…

**50 Billion mobile wireless devices** connected to the Internet across the globe

**Total number of devices connected to the Internet** in some way could reach 500 Billion.



Source: OECD (2012), "Machine-to-Machine Communications: Connecting Billions of Devices", *OECD Digital Economy Papers*, No. 192,

# "Internet of Things"

**Internet of Things** is mainly associated with applications that involve Radio Frequency Identification (RFID). These make use of so called tags, tiny chips with antennae that start to transmit data when they come in contact with an electromagnetic field.

**Machine to Machine** communication (M2M) describes devices that are connected to the Internet, using a variety of fixed and wireless networks and communicate with each other and the wider world.

**Embedded Wireless** has been coined, for a variety of applications where wireless cellular communication is used to connect any device that is not a phone

**Smart** is used in conjunction with various words such **as Living, Cities, Metering, Grids, Water Levy and Lighting** to describe a variety of applications that make use of inexpensive communication to improve the delivery of services.

# Internet of Things Applications by Mobility and Dispersion

|  | **Fixed** | **Mobile** |
|---|---|---|
| **Dispersed** | **Smart Grid, Meter, City**<br><br>**Remote monitoring** | **Car automation**<br><br>**eHealth**<br><br>**Logistics**<br><br>**Portable consumer electronics** |
| **Concentrated** | **Smart Home**<br><br>**Factory automation**<br><br>**eHealth** | **On-site logistics** |

Source: OECD (2012), "Machine-to-Machine Communications: Connecting Billions of Devices", *OECD Digital Economy Papers*, No. 192,

# Big Data

Big data usually includes data sets with sizes beyond the ability of commonly used software tools to capture, curate, manage, and process the data within a tolerable elapsed time.

Source: Snijders, C., Matzat, U., & Reips, U.-D. (2012). 'Big Data': Big gaps of knowledge in the field of Internet science. *International Journal of Internet Science*

Big data are high **volume**, high **velocity**, and/or high **variety** information assets that require new forms of processing to enable enhanced decision making, insight discovery and process optimization.

Source: Douglas, Laney. "The Importance of 'Big Data': A Definition". Gartner, 2012

Big Data mostly uses **inductive statistics with data with low information density whose huge volume** allow to infer laws and thus giving (with the limits of inference reasoning) to Big Data some **predictive capabilities**.

# Where does Big Data come from ?

**12+ TBs**
of tweet data
every day

**25+ TBs** of
log data every
day

**30 billion**
RFID tags today
(1.3B in 2005)

**76 million** smart
meters in 2009…
200M by 2014

**4.6 billion**
camera
phones
world
wide

**100s of
millions of
GPS**
enabled
devices
sold
annually

**2+ billion**
people
on the
Web by
end
2011

Source: IBM Research, 2013

# Big Data is Still Very Complex do Implement

KES - IDT/IIMSS 2013

Source: IBM Research, 2013

# Data Analysis or Decision Models?

# The Concept: Web Competitive Intelligence

# Web Competitive Intelligence

The basic objective is the creation of a simple methodology to develop, implement and manage Competitive Intelligence tools based on information collected from the Web.

Web-based competitive intelligence information is based on **automatic gathering, filtering, search and transformation of information** in the Web using a combination **of crawlers, wrappers and ontologies.**

Data from the Web

# Proposed Framework



Crawler Wrapper

Ontology Storage

Data Analysis & Decision Models

# Proposed Framework



DATA

FILTER

ORGANIZE

SEARCH

INTELLIGENCE

**Multiple-Configuration:**

**New Paradigm: Information flows to the user.**

User

**Notifications**
email
web plataform
instante messaging
intranets

II SÉRIE

**DIÁRIO DA REPÚBLICA**

Quarta-Feira, 28 de Abril de 2010

[www.dre.pt](www.dre.pt)

# PDFs contain information to populate the ontology.

MODELO DE ANÚNCIO DO CONCURSO PÚBLICO

1 - IDENTIFICAÇÃO E CONTACTOS DA ENTIDADE ADJUDICANTE
NIF e designação da entidade adjudicante:
505456010 - Município da Amadora
Serviço/Órgão/Pessoa de contacto: Departamento de Obras Municipais
Endereço: Travessa Vasco da Gama, nº 7
Código postal: 2701 833
Localidade: Amadora
Telefone: 00351 214369000
Fax: 00351 214927837
Endereço Eletrónico: obras.municipais@cm-amadora.pt

2 - OBJETO DO CONTRATO
Designação do contrato: Empreitada nº 9/12 - "Escola EB1/JI Moinhos da Funcheira (ex-Mina 9) - Execução de Obras de Beneficiação"
Descrição sucinta do objeto do contrato: A empreitada consiste na conservação e impermeabilização da Escola EB1/JI Moinhos da Funcheira
Tipo de Contrato: Empreitada de Obras Públicas
Valor do preço base do procedimento 331182.35 EUR
Classificação CPV (Vocabulário Comum para os Contratos Públicos)
Objeto principal
Vocabulário principal: 45214200

6 - LOCAL DA EXECUÇÃO DO CONTRATO
Escola EB1/JI Moinhos da Funcheira, freguesia de S. Brás
País:  PORTUGAL
Distrito: Lisboa
Concelho: Amadora
Código NUTS: PT171

# Method

# Public Tender Ontology

KES - IDT/IIMSS 2013

# Data Reasoning and Querying

Past data and autonomous collection of current instances results in all data related to public tender markets.

Ontologies query and reasoning engines results in a "smart" data base able to answer questions in line with the conceptual knowledge.

**REASONING**: given axioms and restrictions, the engine is able to compute conclusions. Ex.: given the notion the NIF is unique per entity, all entity with NIF X are the same. No need to inject more data then the NIF to identify the public tender publisher.

**QUERY**: given the axioms and restrictions, the engine can compute query without the prior concept of data tables structures and relations. Ex.: Entities that have published PT from cpv 45xxxxxx, in Lisbon, in the last month. All PTs in Porto, above 1.000.000 Euros celling value, with CPV 33xxxxxx.

# Caracterization of the Data

- Developed in Visual Basic.

- Ontology building supported in Protégé.

- Data collected from 2010, 2011 and until 2012.

- 20.000 (aprox.) gathered documents.

- Continuous and autonomous functioning.

# Example of Data Analysis on Public Tenders



Public Tenders in Portuguese E-Procurement Platforms

# Example of Data Analysis on Public Tenders

Select client. Only clients with more than 20 PT available.

### INSTITUTO NACIONAL DE SAUDE DR. RICARDO JORGE I.P.

| | Total | 2010 | 2011 | 2012 |
|---|---|---|---|---|
| PT (#) | 87 | 10 | 50 | 27 |
| PT (%) | 0,56% | 0,15% | 0,75% | 1,19% |
| ∑ BV | 8.768.560,67 € | 649.737,47 € | 5.376.845,98 € | 2.741.977,22 € |
| ∑ BV (%) | 0,08% | 0,01% | 0,11% | 0,22% |
| Average BV | 100.788,05 € | 64.973,75 € | 107.536,92 € | 101.554,71 € |

| Preferred CPV | | | |
|---|---|---|---|
| | PT (#) | Main Cat. | Description |
| 1º | 55 | 33 | Equipamento médico, medicamentos e produtos para cuidados pessoais |
| 2º | 17 | 24 | Produtos químicos |
| 3º | 7 | 31 | Maquinaria, aparelhagem, equipamento e consumíveis eléctricos; iluminação |

| Preferred NUT codes | | |
|---|---|---|
| | PT (#) | Main Cat. |
| 1º | 87 | PT171 |
| 2º | - | - |
| 3º | - | - |

# Example of Data Analysis on Public Tenders



Location of Tenders

# Case Study – Eficiency of Facebook Posts in Hotels Using DEA

Facebook public data was collected. Characterizes the publishing behavior of Facebook pages marked as "Hotel".

- Sample includes 50 most popular "Hotel" pages in Facebook (by Fans number)

- Fans range between 540.000 and 25.000

- All post history, from each page was collected.

- Each post is defined by its type and total number of shares, likes, and comments.

- 78,000 were collected.

- In 2012 all page were active. Hence 38,000 post were used to comparison analysis.

# Page fans (descending)



# Posts per page: same order (random posting strategies)

# Posts Activity by Week day and Type

Posts per weekday: Publishing activity drops on weekends.



| Monday | Tuesday | Wednesday | Thursday | Friday | Saturday | Sunday |
|--------|---------|-----------|----------|--------|----------|--------|
| 6044 | 5656 | 5727 | 6211 | 6187 | 4363 | 3887 |

Posts per type: Photo leads



| photo | video | status | link | offer | question | swf |
|-------|-------|--------|------|-------|----------|-----|
| 24344 | 1326 | 2327 | 9635 | 55 | 100 | 288 |

# Models that Classify the Page Efficiency Considering that High Interaction Levels Represent Good Output of Marketing Strategy.

**Inputs:**

- Number of photo posts done in 2012
- Number of status posts done in 2012
- Number of link posts done in 2012
- Number of video posts done in 2012

**DEA models 1.X:**

**Outputs:**

- Total number of shares of all posts
- Total number of likes of all posts
- Total number of comments of all posts

**DEA models 2.X:**

**Outputs:**

- Ratio of total number of shares of all posts over page fans.
- Ratio of total number of likes of all posts over page fans.
- Ratio of total number of comments of all posts over page fans.

# DEA Models

Four BCC DEA models where design for each output type 1.X and 2.x
They matchup different BCC parameterizations.

DEA INPUTS (Ex: 15 from the 50 pages)

| # | Hotel | INPUTS | | | | Outputs 1.X | | | Outputs 2.X | | |
|---|-------|--------|--------|-------|--------|--------|--------|--------|--------|--------|--------|
| | | Photos | Status | Links | Videos | Shares | Likes | Comm. | Shares | Likes | Comm. |
| 1 | AO Hostels | 115 | 82 | 25 | 8 | 483 | 8496 | 4363 | 1.95% | 34.37% | 17.65% |
| 2 | ARIA Resort & Casino | 572 | 89 | 301 | 13 | 35818 | 681219 | 24712 | 6.61% | 125.80% | 4.56% |
| 3 | Atlantis The Palm Dubai | 448 | 12 | 62 | 35 | 45173 | 336411 | 16855 | 19.23% | 143.22% | 7.18% |
| 4 | Bellagio Las Vegas | 337 | 23 | 53 | 5 | 26998 | 271123 | 13449 | 7.71% | 77.45% | 3.84% |
| 5 | Big Cedar Lodge Official Page | 300 | 38 | 151 | 8 | 2472 | 32638 | 3481 | 9.90% | 130.67% | 13.94% |
| 6 | Caesars Palace | 455 | 53 | 171 | 40 | 19796 | 228561 | 13265 | 6.25% | 72.11% | 4.19% |
| 7 | Casa Andina Hotels | 360 | 14 | 25 | 19 | 9308 | 58605 | 4471 | 20.17% | 126.97% | 9.69% |
| 8 | Courtyard by Marriott Aguadilla | 73 | 16 | 4 | 1 | 1703 | 22029 | 643 | 3.36% | 43.49% | 1.27% |
| 9 | Cove Haven Entertainment Reso | 412 | 267 | 143 | 20 | 1425 | 39564 | 7317 | 5.54% | 153.89% | 28.46% |
| 10 | Danubius Hotels Group | 196 | 2 | 144 | 10 | 3972 | 38097 | 2190 | 9.58% | 91.92% | 5.28% |
| 11 | El Conquistador Resort | 647 | 117 | 94 | 31 | 3686 | 72452 | 3820 | 10.38% | 204.02% | 10.76% |
| 12 | French Lick Resort | 132 | 70 | 128 | 22 | 2183 | 20206 | 2592 | 7.43% | 68.76% | 8.82% |
| 13 | Grand Sierra Resort and Casino | 393 | 28 | 128 | 40 | 3958 | 48728 | 6030 | 9.57% | 117.82% | 14.58% |
| 14 | Great Wolf Lodge | 388 | 51 | 217 | 61 | 7440 | 103926 | 9963 | 1.87% | 26.12% | 2.50% |
| 15 | Hard Rock Hotel and Casino Las | 1052 | 148 | 335 | 23 | 14101 | 158669 | 9982 | 9.57% | 107.63% | 6.77% |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |

# Example of DEA Models Results: Model 1.1

| # | Hotels | Efficiency | SLACKS | | | | | | |
|---|--------|------------|--------|--------|-------|--------|--------|--------|-------|
|   |        |            | Photos | Status | Links | Videos | Shares | Likes | Comm. |
| 1 | AO Hostels | 1.000 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | ARIA Resort & Casino | 1.000 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | Atlantis The Palm Dubai | 1.000 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | Bellagio Las Vegas | 1.000 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5 | Big Cedar Lodge Official Page | 0.342 | 8 | 0 | 24 | 2 | 829 | 30950 | 561 |
| 6 | Caesars Palace | 0.741 | 0 | 16 | 74 | 25 | 7202 | 42562 | 184 |
| 7 | Casa Andina Hotels | 1.000 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 8 | Courtyard by Marriott Aguadilla Hotel & Casino | 1.000 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 9 | Cove Haven Entertainment Resorts | 0.818 | 0 | 195 | 64 | 11 | 25573 | 231559 | 6132 |
| 10 | Danubius Hotels Group | 1.000 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 11 | El Conquistador Resort | 0.564 | 28 | 43 | 0 | 12 | 23312 | 198671 | 9629 |
| 12 | French Lick Resort | 0.720 | 0 | 37 | 64 | 15 | 1118 | 43382 | 1450 |
| 13 | Grand Sierra Resort and Casino | 0.858 | 0 | 1 | 57 | 29 | 23040 | 222395 | 7419 |
| 14 | Great Wolf Lodge | 0.869 | 0 | 21 | 135 | 48 | 19558 | 167197 | 3486 |
| 15 | Hard Rock Hotel and Casino Las Vegas | 0.320 | 0 | 24 | 54 | 2 | 12897 | 112454 | 3467 |
| … | … | … | … | … | … | … | … | … | … |

Slacks can be interpret as: the over investment in publishing (inputs) or the lack of results (output) in order to achieve equivalent efficiency compared to best practices.

# DEA Models Result: All models. Combined Analysis Sustain Consistent Conclusions.

- 13 and 18 are consistent efficient Facebook pages.
- Ratio interactions demonstrates that large pages are in fact inefficient.
- Model 1.1 and 2,1 are very "forgiving".

| # | Fans | Hotels | Absolute Interations | | | | Relative Interaction | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | 1.1 | 1.2 | 1.3 | 1.4 | 2.1 | 2.2 | 2.3 | 2.4 |
| 1 | 541492 | ARIA Resort & Casino | 1.000 | 1.000 | 1.000 | 1.000 | 0.444 | 0.255 | 0.255 | 0.211 |
| 2 | 453806 | Holiday Inn | 1.000 | 0.380 | 0.297 | 0.297 | 0.333 | 0.172 | 0.054 | 0.054 |
| 3 | 397837 | Great Wolf Lodge | 0.869 | 0.516 | 0.516 | 0.442 | 0.255 | 0.104 | 0.100 | 0.100 |
| 4 | 350056 | Bellagio Las Vegas | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 | 0.411 | 0.397 | 0.397 |
| 5 | 316960 | Caesars Palace | 0.741 | 0.633 | 0.633 | 0.593 | 0.558 | 0.160 | 0.160 | 0.160 |
| 6 | 238724 | Mandalay Bay Resort and Casino | 1.000 | 0.752 | 0.741 | 0.741 | 1.000 | 0.315 | 0.315 | 0.305 |
| 7 | 234883 | Atlantis The Palm Dubai | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 | 0.784 | 0.784 | 0.724 |
| 8 | 212354 | Ushua\u00efa Ibiza Beach Hotel (C | 1.000 | 1.000 | 1.000 | 0.868 | 0.886 | 0.614 | 0.614 | 0.338 |
| 9 | 197343 | Sakura Hotel & Hostel in Tokyo Jap | 1.000 | 0.705 | 0.705 | 0.651 | 0.513 | 0.292 | 0.292 | 0.245 |
| 10 | 179008 | Resorts World Genting | 0.349 | 0.155 | 0.150 | 0.150 | 0.118 | 0.052 | 0.052 | 0.052 |
| 11 | 178922 | Planet Hollywood Resort & Casino | 0.789 | 0.259 | 0.241 | 0.241 | 0.286 | 0.116 | 0.104 | 0.104 |
| 12 | 176212 | Mazagan Beach Resort | 1.000 | 0.186 | 0.173 | 0.173 | 0.139 | 0.112 | 0.094 | 0.094 |
| **13** | **171890** | **Vital Hotel Westfalen Therme Spa** | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** |
| 14 | 169475 | Pearl Continental Karachi | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 | 0.634 | 0.634 | 0.551 |
| 15 | 156678 | Palms Casino Resort | 0.955 | 0.237 | 0.235 | 0.235 | 0.269 | 0.126 | 0.119 | 0.119 |
| 16 | 151413 | The Cosmopolitan of Las Vegas | 1.000 | 1.000 | 1.000 | 1.000 | 1.000 | 0.527 | 0.527 | 0.488 |
| 17 | 147419 | Hard Rock Hotel and Casino Las Ve | 0.320 | 0.200 | 0.200 | 0.199 | 0.241 | 0.117 | 0.117 | 0.115 |
| **18** | **119188** | **Horta da Moura** | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** | **1.000** |

Objective:
- Gather Facebook data
- Model Post Life Cycle (LC)
- Design intelligent algorithms to detect behavior outliers.

Development:
- Large data set collection.
- Model design
- Operationalize algorithm
    - Design ontology system
    - Populate
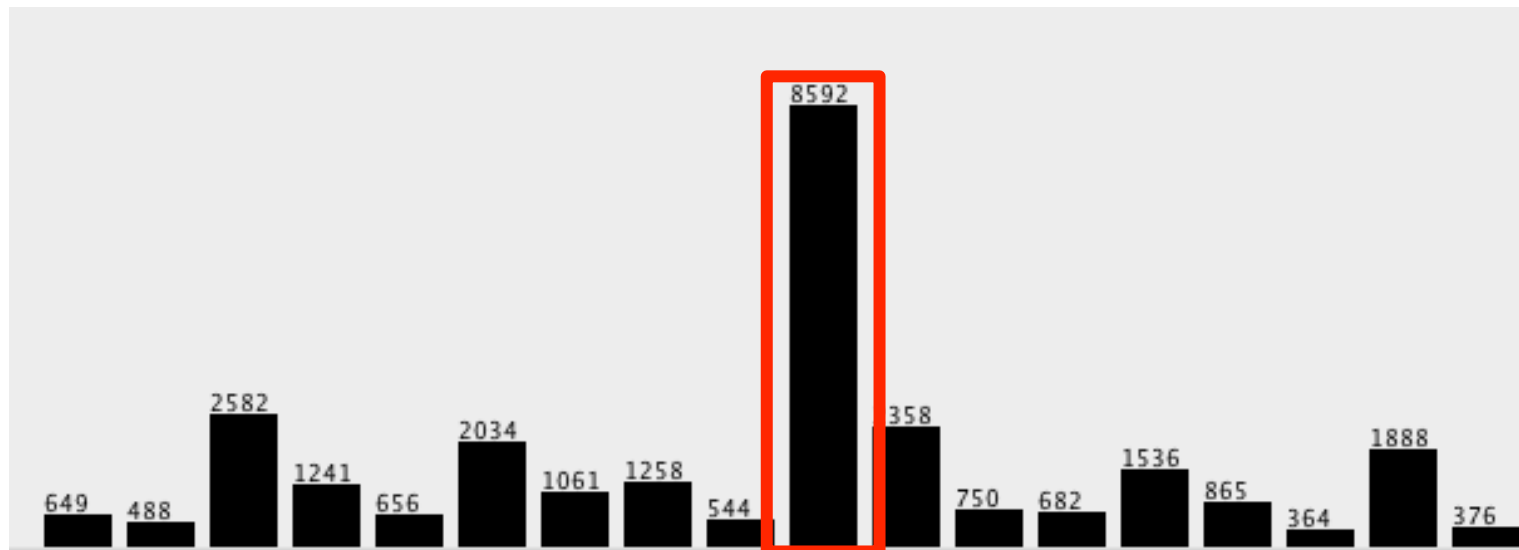    - Let program infer when posts have uncommon behavior

# Facebook public data was collected. Its characterizes the publishing behavior of Facebook pages from 18 categories.

---

- Sample includes 560 most popular pages.
- Fans range between 189.660 and 450 thousand fans.
- Posts publish between 21st December 2012 and 21st January 2013.
- Collection gather post history, i.e., time series data of each one.
- 680.000 lines data set.
- 25.450 valid posts (deleted posts by page administrators were excluded).
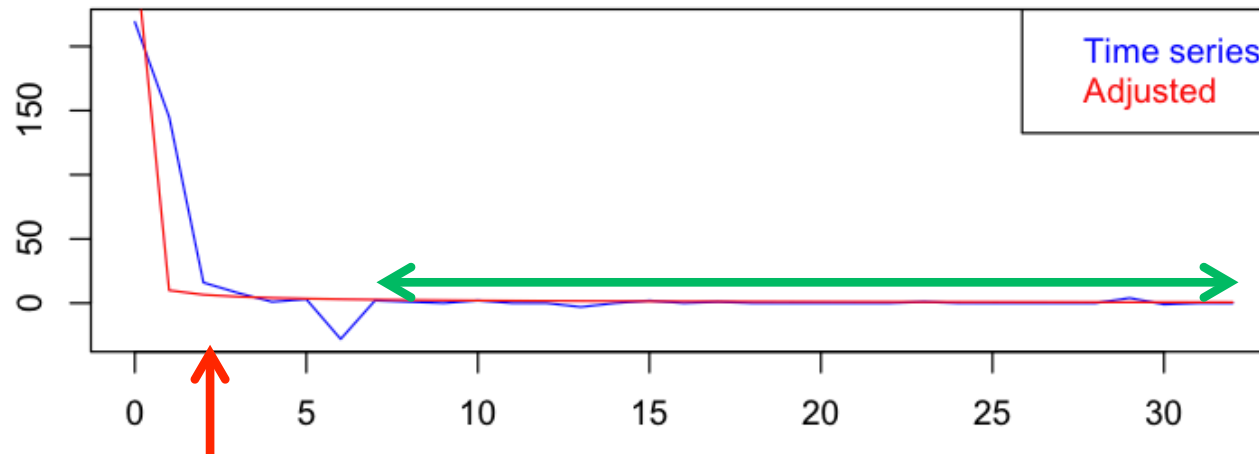
Posts per category.



Media pages are the most active.

# Typical post life cycle (95%).
## Negative exponential fitting with logarithmic time transformation.
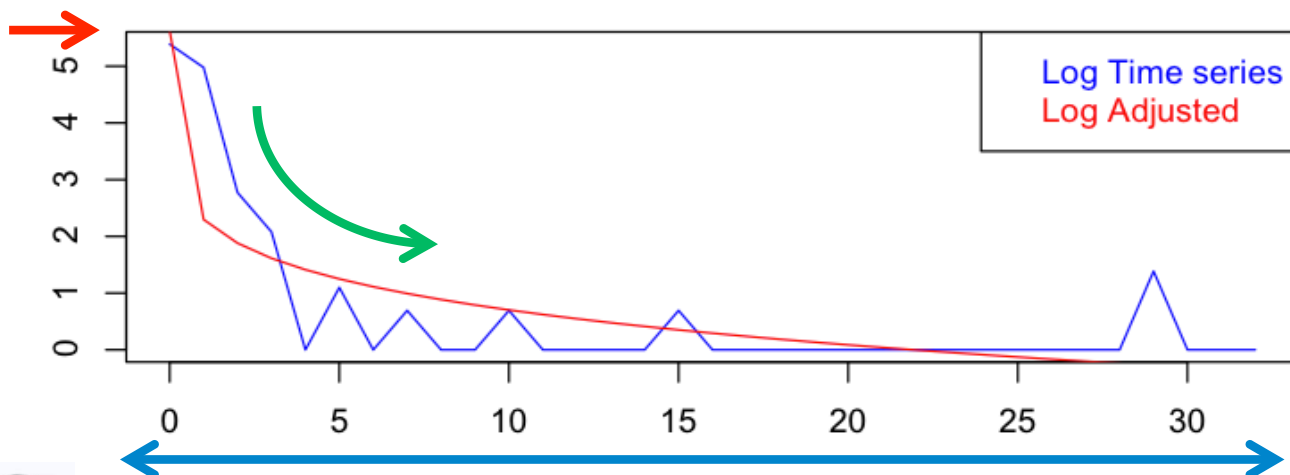
**Comments Real Data**

Real data

One to two days fierce activity

Flat (no activity) after initial days.

LC Model

**Comments Log Data**

$\alpha$ : how high it starts.
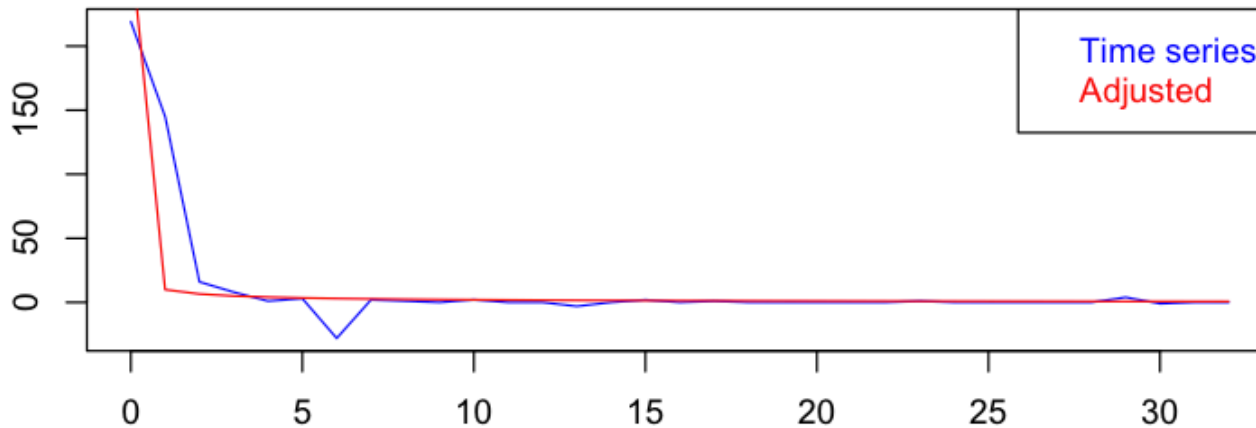
$\beta$: how fast it decays in time.

R: Time transformation for fitting purposes.

Adjst: how well the model predicts LC.
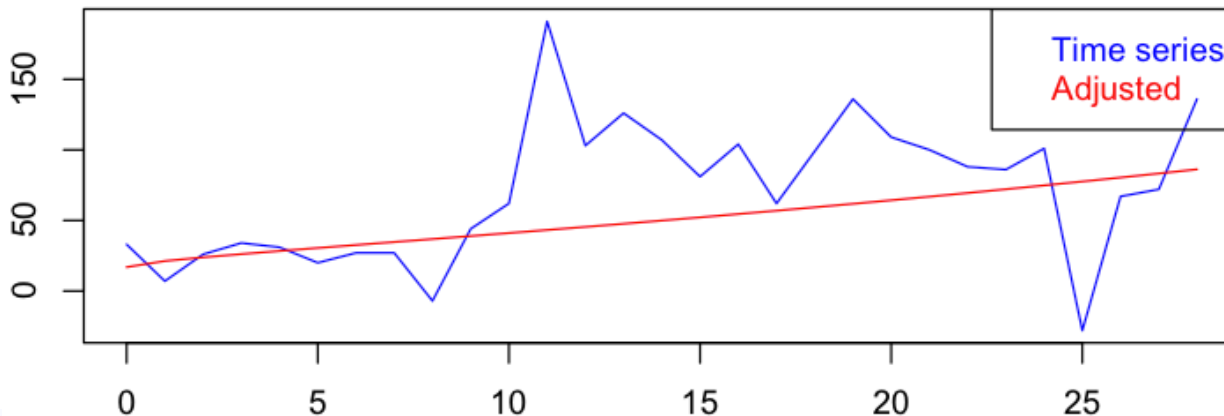
# Typical vs. outlier post life cycle.

## TYPICAL LC

**Comments Real Data**



- α: can be high or low.
- β: highly negative constant.
- R: can be high or low
- Adjst: adjust well (> 0.75)

## OUTLIER LC

**Comments Real Data**



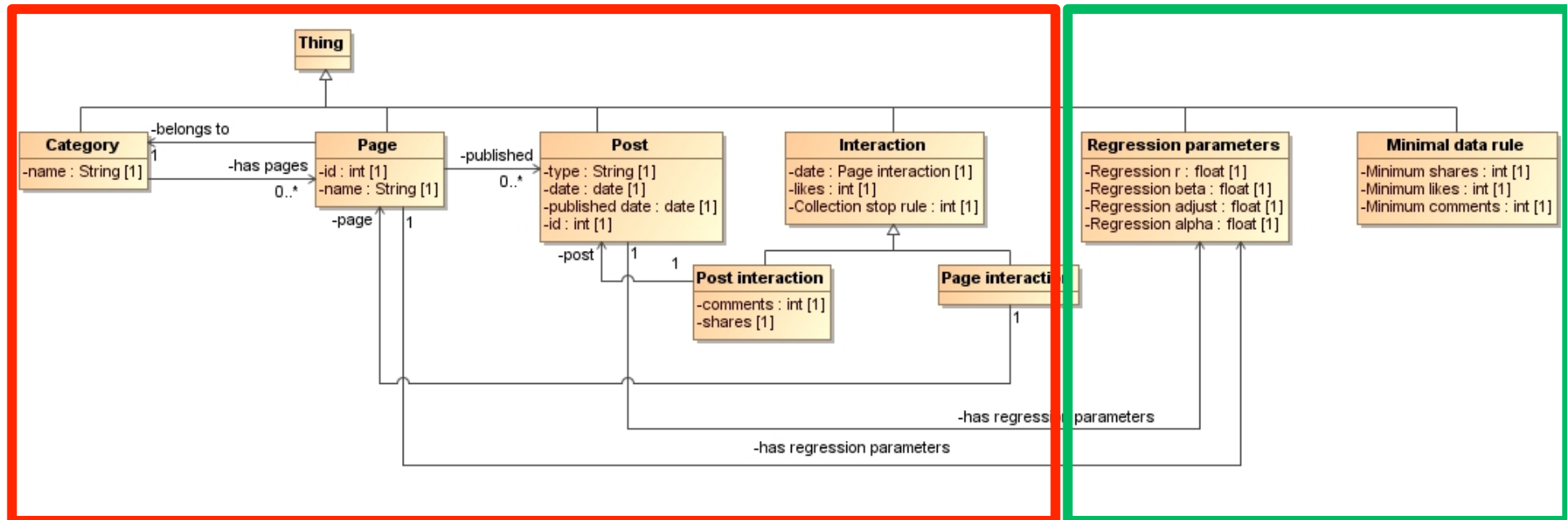- α: can be high or low
- β: negative near zero or positive.
- R: can be high or low.
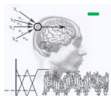- Adjst: does not adjust well (< 0.75)

# Facebook Post Life Cycle Intelligent Algorithm – The Ontology



It represents knowledge about Facebook pages:

- Pages have fans, belong to categories and publish posts.
- Posts are categorized by types, have like, shares, and comments responses.
- Page fans and posts interactions are recorded each day.
- Given a minimum post life time (ex.: 3 days) the regression model calculates regression models indicators.
- Regression indicators are then related to a post like, share or comment time series.

# Facebook Post Life Cycle Intelligent Algorithm – The Ontology

A collection algorithm is responsible for instantiating the ontology.

Given the regression indicators described before is possible to detect posts that had unusual behaviors. In other words by representing a knowledge concept of a unusual post, is possible to ask the ontology to detect all behavior outliers. The rules by which a post is considered and outlier can be parameterized and improved over time. Examples of outlier post knowledge concept:

1. Post interaction that: has **beta > 0** *AND* **adjustment < 0.7**
2. Post interaction that: has **adjustment < 0.5**
3. Post interaction that: has **beta > 0** *OR* **adjustment < 0.6**
4. Post interaction that: [**has beta > 0** *AND* **adjustment < 0.8**] *OR* [**has adjustment < 0.5**]
5. Etc…

Posts are classified has having a typical or outlier behavior. Outliers should be analyzed and requires special attention. They may represent controversial content, interactive features, have viral nature, generated contestation, etc…

# Conclusions

- We are living in the Era of Big numbers

- We must learn how to converge Data, Analysis and Decision Making