# Data and Algorithmic **Bias** in the **Web**

**Ri**c**ardo Baeza-Yates**
**C**alifornia, **C**atalonia, **C**hile

*Web Science 2016, Hannover, Germany, May 2016*

---

## All Data has Bias

- Gender
- Racial
- Sexual
- Religious
- Social
- Linguistic
- Geographic
- Political
- Educational
- Economic
- Technological

- from Noise or Spam
- Validity (e.g. temporal)
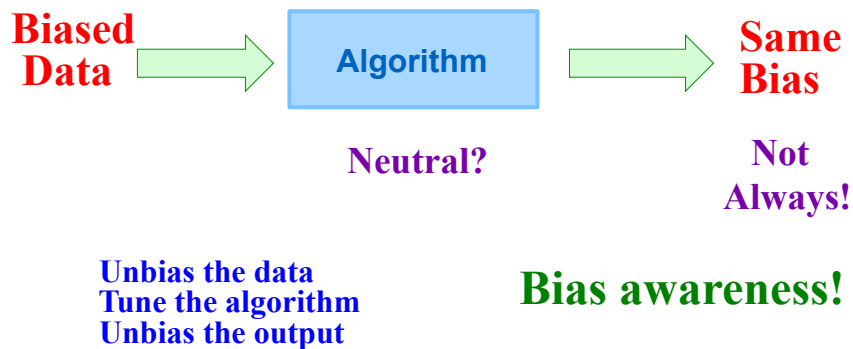- Completeness
- Gathering process
- ….

**However many people extrapolate results to the whole population (e.g., social media analysis)**

**In addition there is bias when measuring bias as well as bias towards measuring it!**

# Yes, We Live in a (Very) Biased World!



# A Non-Technical Question

**Biased Data** → **Algorithm** → **Same Bias**

**Neutral?**     **Not Always!**

**Unbias the data**
**Tune the algorithm**
**Unbias the output**

**Bias awareness!**
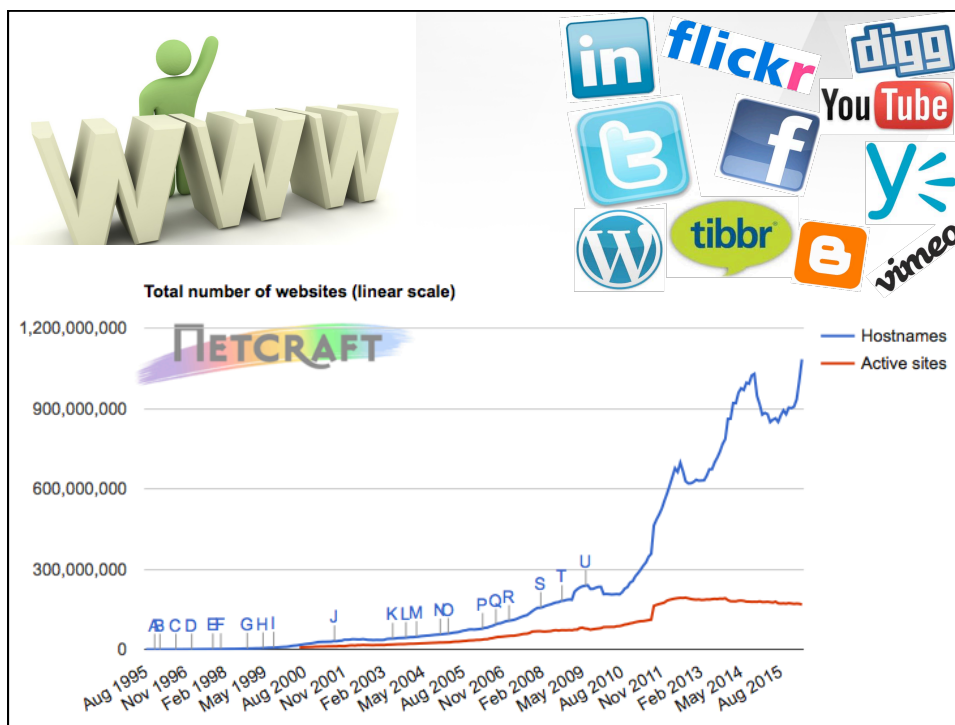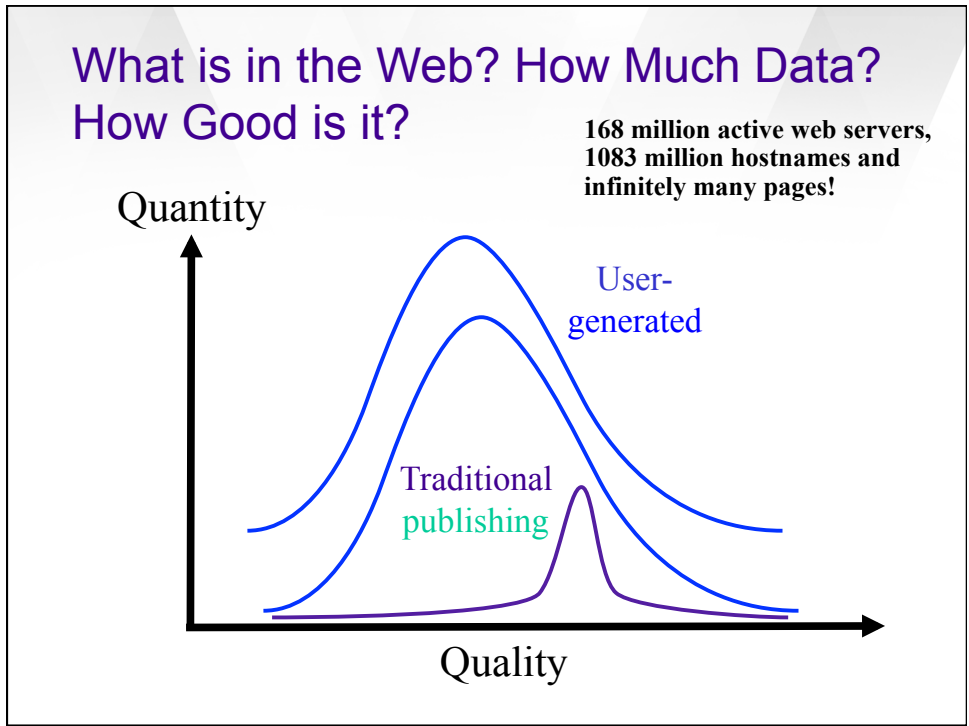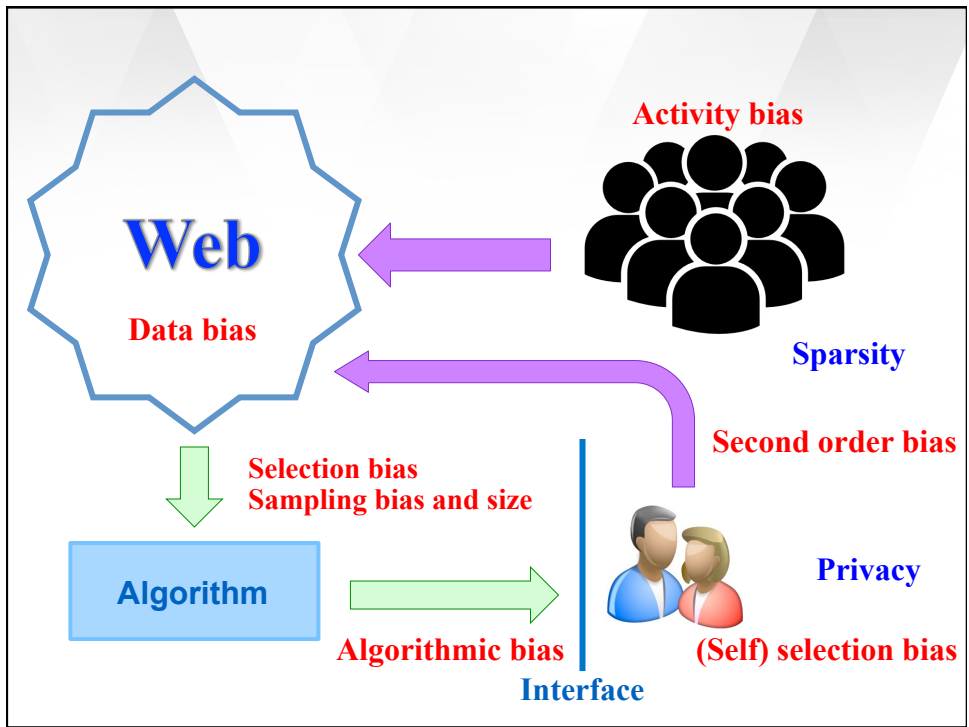
# Big Data and Bias

- The quality of any algorithm is bounded by the quality of the data that uses
- Data bias awareness
- Algorithmic fairness
- Key issues for machine learning
  - Uniformity of data properties
    - In the Web, distributions resemble a power law
  - Uniformity of error
  - Data sample methodology
    - E.g., sample size to see infrequent events or sampling bias issues

12

Activity bias

**Web**

Data bias

Sparsity

Second order bias

Selection bias
Sampling bias and size

Algorithm

Privacy

Algorithmic bias

(Self) selection bias

Interface

---

# What is in the Web? How Much Data? How Good is it?

**168 million active web servers, 1083 million hostnames and infinitely many pages!**



Quantity

User-generated

Traditional publishing

Quality

# What else is in the Web?

---

# Noise and Spam

- Noise may come from many places:
  - Instruments that measure (e.g., IoT)
  - How we interpret the data

- Spam is everywhere

- Fight both with the wisdom of the crowds

# Web Spam

- Deceiving text, links, clicks…
        due to an economic incentive
- Depending on the goal and the data,
    spam is easier to generate
- Depending on the type & target data,
    spam is easier to fight
- Disincentives for spammers?
  - Social
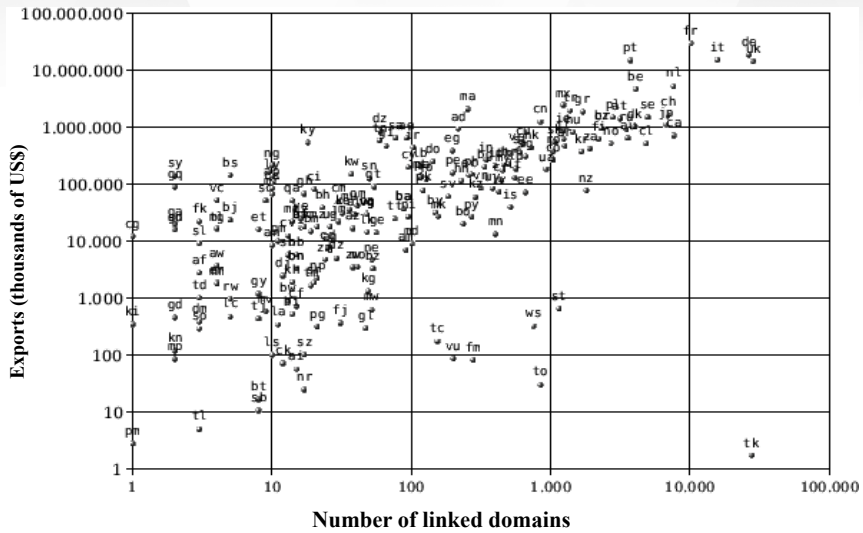  - Economic

Web Spam is NOT Mail Spam

# Data Bias and Redundancy

- There is any dependency in the data?
- There is any duplication?
  - Lexical duplication in the Web is around 25%
  - Semantic duplication is larger (more later)
- Any other biases? Many!
  - Web structure (economic, cultural)
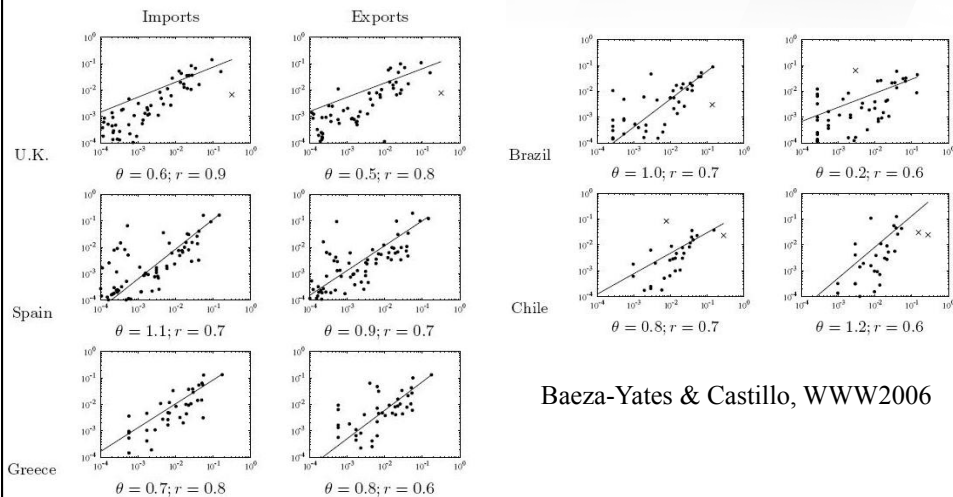  - Web content (linguistic, geography, gender)

# Economic Bias in Links



[Baeza-Yates, Castillo & López. Characteristics of the Web of Spain.
The Information Professional (Spanish), 2006, vol. 15, n. 1, pp. 6-17]
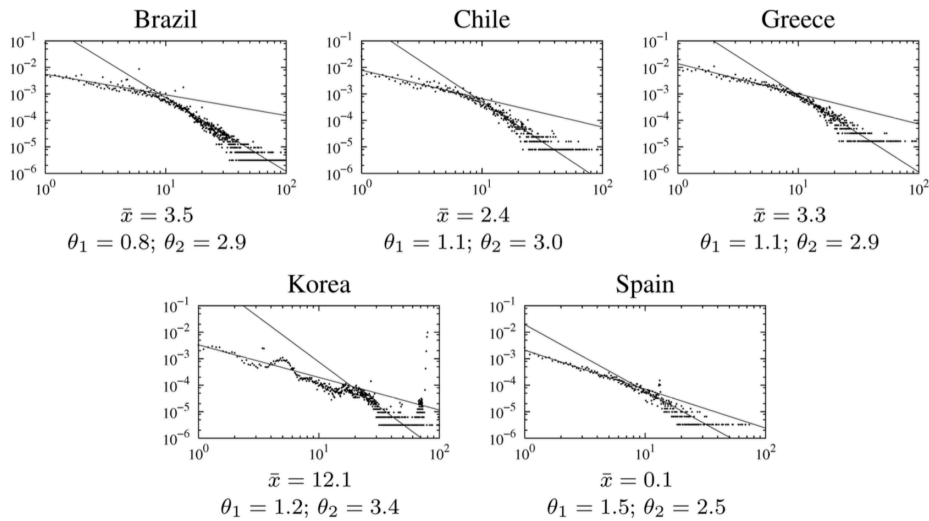
36

---

# Exports/Imports vs. Domain Links



Baeza-Yates & Castillo, WWW2006

37

# Website Structure

Shame

Minimal effort

### Brazil
$\bar{x} = 3.5$
$\theta_1 = 0.8; \theta_2 = 2.9$

### Chile
$\bar{x} = 2.4$
$\theta_1 = 1.1; \theta_2 = 3.0$

### Greece
$\bar{x} = 3.3$
$\theta_1 = 1.1; \theta_2 = 2.9$

### Korea
$\bar{x} = 12.1$
$\theta_1 = 1.2; \theta_2 = 3.4$

### Spain
$\bar{x} = 0.1$
$\theta_1 = 1.5; \theta_2 = 2.5$

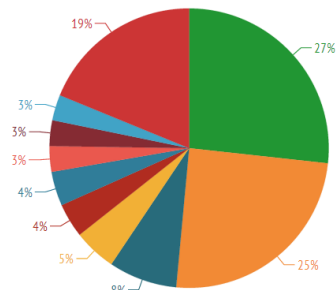[Baeza-Yates, Castillo, Efthimiadis, TOIT 2007]

---

# Linguistic Bias

Top 25 World Languages
- Chinese, Mandarin
- Spanish
- English
- Hindi

**Top Ten Languages in the Internet in millions of users - November 2015**
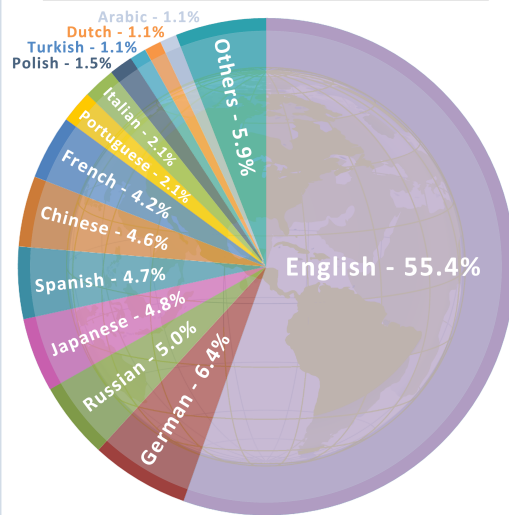
## Languages on the Web

English, Chinese Mandarin, Spanish, Japanese, Portuguese
German, Arabic, French, Russian, Other

### The Languages of Web Content
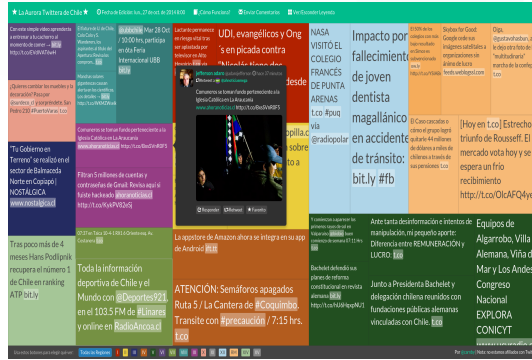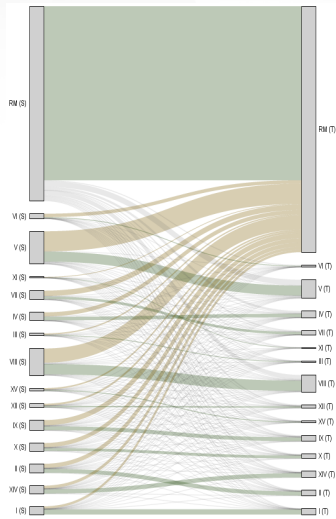The percentage of the top 1 million websites available in various languages

- Arabic - 1.1%
- Dutch - 1.1%
- Turkish - 1.1%
- Polish - 1.5%
- Italian - 2.1%
- Portuguese - 2.1%
- French - 4.2%
- Chinese - 4.6%
- Spanish - 4.7%
- Japanese - 4.8%
- Russian - 5.0%
- German - 6.4%
- Others - 5.9%
- English - 55.4%

Source: Web Technology Surveys
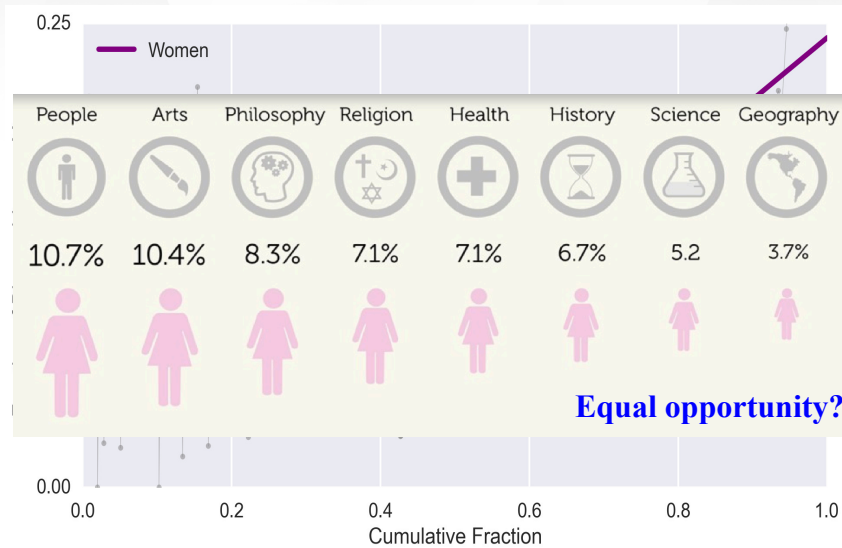
language connect

Insight   Digital   Audio   Branding   Text

Web: languageconnect.net  /  E-mail: info@languageconnect.net  /  Follow us: @lconnect

# Geographical Bias



[E. Graells-Garrido and M. Lalmas, "Balancing diversity to counter-measure geographical centralization in microblogging platforms", ACM Hypertext'14]

---

# Gender Bias

**Systemic bias?**



| People | Arts | Philosophy | Religion | Health | History | Science | Geography |
|--------|------|-----------|----------|--------|---------|---------|-----------|
| 10.7% | 10.4% | 8.3% | 7.1% | 7.1% | 6.7% | 5.2 | 3.7% |

0.25

Women

**Equal opportunity?**

0.00

0.0          0.2          0.4          0.6          0.8          1.0

Cumulative Fraction

[E. Graells-Garrido et al,. "First Women, Second Sex: Gender Bias in Wikipedia", ACM Hypertext'15]
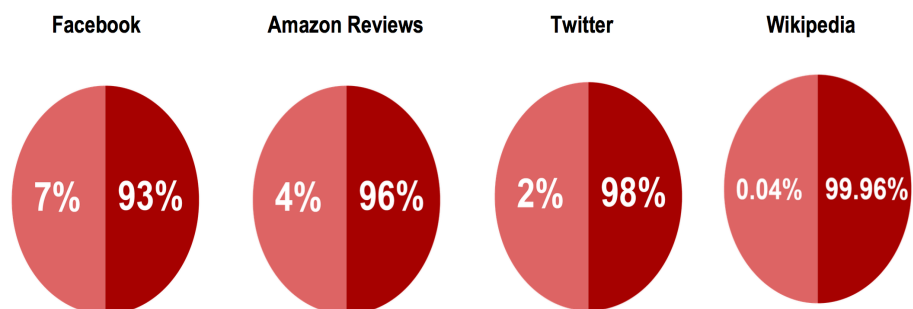
## Activity Bias: Wisdom of a Few?

- The Web already is influenced by small groups

  - "0.05% of the user population, attract almost
    50% of all attention within Twitter" (50K users)
    [Wu, Hofman, Mason & Watts, WWW 2011]

- We explored this issue further with four different datasets:
  1. a large one from Twitter (2011),
  2. a small one from Facebook (2009),
  3. Amazon reviews (2013), and
  4. Wikipedia editors (2015).

- Digital desert: the content that is never seen

[Baeza-Yates & Saez-Trumper, ACM Hypertext 2015]

---

## Examples

# How many users produce most of the content?

| Facebook | Amazon Reviews | Twitter | Wikipedia |
|----------|----------------|---------|-----------|



7% 93%    4% 96%    2% 98%    0.04% 99.96%

[Baeza-Yates & Saez-Trumper, ACM Hypertext 2015]

**theguardian**

ort   football   opinion   culture   business   lifestyle   fashion   environment   tech   travel                    ≡ all sections

Amazon sues 1,000 'fake reviewers'

**October 2015**

Online retailer files lawsuit in US against people whose names it says it does not know, claiming they offer reviews for sale

# Amazon Continues Their Crusade Against Fake Reviews

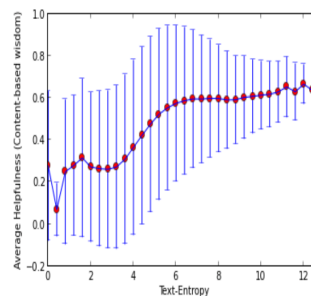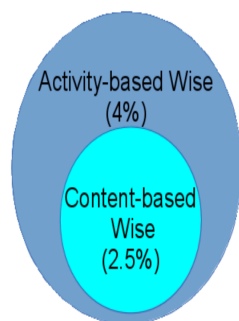By Tyler Lee on 04/26/2016 05:07 PDT

---

# Quality of Content?

- Adding content implies adding wisdom?
- We use Amazon's reviews helpfulness
- We computed the text entropy
- Content-based-wise users
- How many of those users are being paid?

# Digital Desert

- 1.1% of the Twitter content is never seen.*
- 31% of articles added/edited in May 2014 in wikipedia, were not visited in June.



# Bias in the Interface



**Presentation bias**

**Position bias**
**Ranking bias**

**Irrelevant alternatives bias**

**Social bias**

**Interaction bias**

# Presentation Bias

- Interaction data will be biased to what is shown
- In recommender systems, items recommended will get more clicks than items not recommended
- In search systems top ranked results will get more clicks than other results
  - Ranking bias
  - Interaction bias

[Dupret & Piwowarski, SIGIR 2008]
[Chapelle & Zhang, WWW 2009]



# Social Bias



[WHY AMAZON'S RATINGS MIGHT MISLEAD YOU; The Story of Herding Effects
Ting Wang and Dashun Wang, Big Data, 2014]

# Independence of Irrelevant Alternatives



# Irrelevant Alternatives (IIA)

- IIA does not always hold
- If that is the case, nested logit should be used instead of multinomial logit
- Optimal quadratic algorithm to recover trees for a nested decision process [Benson et al, WWW 2016]
- Statistical tests to check if a nested model works (95%)

| Dataset | Test | | | |
|---------|------|------|------|------|
| | SB | MSB | AMSB | CSB |
| RESTAURANTS | 0.087 | 0.066 | 0.076 | 0.041 |
| JAPANESECUISINE | 0.325 | 0.238 | 0.316 | 0.093 |
| LASTFMARTISTS | 0.106 | 0.102 | 0.129 | 0.049 |
| LASTFMGENRE | 0.300 | 0.143 | 0.284 | 0.094 |

- How much can be explained by bias?

# Extreme Algorithmic Bias



# Second Order Bias in Web Content



**Ranking bias**
**Redundancy grows (35%)**

Person

Query

New

Web content is redundant

Clicks in results are biased to
the ranking and the interaction

Search results

**[Baeza-Yates, Pereira & Ziviani,
Geneological Trees in the Web, WWW 2008]**

# The Long Tail: Sparsity

- Why there is a long tail?
- Sampling in the tail
- When the crowd dominates
- Empowering the tail

[Anatomy of the long tail: Ordinary People with Extraordinary Tastes, Goel, Broder, Gabrilovich, Pang; WSDM 2010]

**Most measures in the Web follow a power law**

---

# Sample Size?

- If we want to estimate the frequency of queries that appear with probability at least *p* with a certain relative error $\epsilon$ we can use the standard binomial error formula √(1-p)/*np* which works well for *p* near ½

- Better is the Agresti-Coull technique (also called *take 2*) which gives:

$$n \geq Z_{1-\alpha/2}^2 \left( \frac{p'(1-p')}{\epsilon^2} - 1 \right)$$
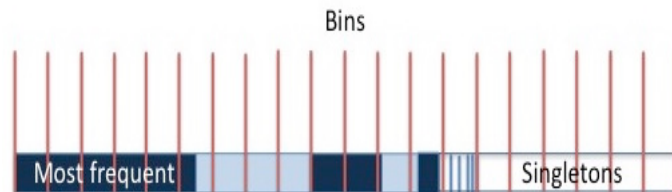
- where *Z* is the inverse of the standard normal distribution, $1-\alpha$ is the confidence interval and $p' = p + Z^2/2$

- If *p* = 0.1, $1-\alpha$ is 90% and $\epsilon$ is 10%, we get *n* = 2342. The standard formula gives *n* = 900.

[Baeza-Yates, SIGIR 2015, Industry track]

# Incremental Sampling

- Main goal: make good samples consistent across time

- Simple idea based in stratified sampling: bins + random start point



Bins

Most frequent          Singletons

- Bin size can be found by binary search starting with a good approximation if a query frequency model is used ($b < V/n$)

- This perfectly mimics the head of the distribution, but not the tail

- Change the bins in the tail to get the right distribution

# Fixing the Tail

- To mimic the tail we change the binning size when we reach a query frequency of $b/2$

- If we want a singleton ratio of $\beta = S/V$ we recalculate the binning size as

$$b' = (1 - \beta)(Q - Q')/(\beta V')$$

- where Q' and V' are the partial vocabulary size and volume before changing the bin size.

# Stratified Sampling Example

---

# When the Crowd Dominates

Kills the long tail

Personalization "facets":
- Language (not always)
- Location
- Semantic facets per user
  - Query intent prediction in search

# Empowering the Tail

The Filter "Bubble", Eli Pariser

- Avoid the Poor get Poorer Syndrome
- Avoid the Echo Chamber
- How to expose opposite views?

Solutions:
  - Diversity
  - Novelty
  - Serendipity

**Cold start problem solution:**
**Explore & Exploit**

# Aggregating in theTail

- Exploit the context (and deep learning!)

  91% accuracy to predict the next app you will use
  [Baeza-Yates et al, WSDM 2015]

- Personalization vs. Contextualization
  Recall that user interaction is another long tail

Interests



People

Photoset          User Photo Streams          Photo-POI Mapping          Timed Paths

[De Choudhury et al, ACM HT 2010]

# Crowdsourcing Data: Good Paths



SHORTEST          HAPPY

BEAUTY          QUIET

[Quercia et al, ACM HT 2014]          83

# Regions from Pictures



[Thomee et al, Demo at CHI 2014]

---

# Privacy 101:
# AOL Query Logs Release Incident

A Face Is Exposed for AOL Searcher No. 4417749,
By MICHAEL BARBARO and TOM ZELLER Jr,
The New York Times, Aug 9  2006



- No. 4417749 conducted hundreds of searches over a three-month period on topics ranging from "numb fingers" to "60 single men".
- Other queries: "landscapers in Lilburn, Ga," several people with the last name Arnold and "homes sold in shadow lake subdivision gwinnett county georgia."
- Data trail led to Thelma Arnold, a 62-year-old widow who lives in Lilburn, Ga., frequently researches her friends' medical ailments and loves her three dogs.

## Netflix settles privacy lawsuit, ditches $1 million contest

Netflix's next recommendation engine may not come from its community after all …

by **Jacqui Cheng** - Mar 12, 2010 10:04pm CET

Netflix has canceled its $1 million contest aimed at finding a better recommendation engine in the wake of a privacy lawsuit settlement. The company informed its users today via the company blog, noting that it had "reached an understanding" with the Federal Trade Commission, leading it to ditch the Netflix Prize contest.

Netflix first announced the contest—actually the sequel to its original contest—in August of 2009. The goal was to crowdsource its active user base to write a more intelligent recommendation engine based on users' past rentals. This is something Netflix already does, of course, but there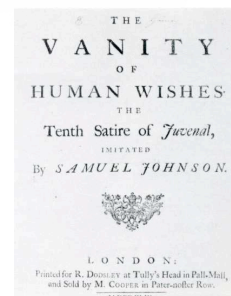's always room for improvement; the company wanted to find the Next Big Thing™ by offering $1 million to the person with the best algorithm.

Part of the contest involved Netflix disclosing what it considered to be anonymized user data to those trying to come up with solutions. This, however, led to a lawsuit by a closeted lesbian mother who argued that Netflix had not sufficiently anonymized the information and that she (among others) could be easily outed due to her own rental history. Indeed, within weeks of the data being released, researchers had found a way to use an external data source to decode an individual's viewing history with surprising accuracy, but Netflix did not immediately withdraw the contest.

87

---

## Risks of Privacy in Query Logs

- Profile [Jones, Kumar, Pang, Tompkins, CIKM 2007]
  - Gender:          84%
  - Age (±10):        79%
  - Location (ZIP3): 35%
- Vanity Queries [Jones et al, CIKM 2008]
  - Partial name:  8.9%
  - Complete:       1.2%
- More information:
  - A Survey of query log privacy-enhancing techniques from a policy perspective [Cooper, ACM TWEB 2008]
- A good anonymization technique is still an open problem

THE
VANITY
OF
HUMAN WISHES.
THE
Tenth Satire of *Juvenal*,
IMITATED
By *SAMUEL JOHNSON*.

LONDON:
Printed for R. Dodsley at Tully's Head in Pall-Mall,
and Sold by M. Cooper in Pater-noster Row.
MDCC.XLIX.

# Privacy Awareness

- How our privacy changes when we change our social network?
- Information gain to predict a private attribute based on public data
- Each user may have a promiscuity score

**Confirm** **Delete Request**

- Example: new friendship request

  Promiscuity( me ) > Promiscuity( new) ● (green)

  Promiscuity( me ) ≥ Promiscuity( new ) + max-gain-I-allow ● (yellow)

  Promiscuity( me ) < Promiscuity( new ) + max-gain-I-allow ● (red)

  Related work by [Estivill-Castro & Nettleton; Singh, ASONAM 2015]

  Cancel **Wagner Meira Jr.** Post

  **Ricardo Baeza-Yates**
  Wagner's friends of friends

  Is this a comment that they should see?

---

# The Web Works Thanks to Bias!

- Web traffic
  › Local caching
  › Proxy/Akamai caching
- Search engines
  › Answer caching
  › Essential web pages
    • 25% queries can be answered with less than 1% of the URLs
      [BY,Boldi, Chierichetti, WWW 2015]
- E-Commerce
  › Most revenue comes from few items

**Activity bias**

**(Self) selection bias**

# Web Data

- A mirror of ourselves, the good, the bad and the ugly

- The web amplifies everything, good or bad, but always leaves traces

- We have to be aware of the biases and contrarrest them

- We have to be aware of our privacy

    **Big Data** of **People** is huge…..
        …..  but is tiny compared to the future
        **Big Data** of the **Internet of Things (IoT)**


# It's Hard to Get Data to Tell the Truth

- The blindness of the averages
    - Look at distributions

- Absolute vs. relative
    - Income per capita vs. Inequality

- Local vs. global optimization
    - Teams competing without knowing, uncorrelated criteria

- You can always see/torture data as you wish
    › 61 analysts, 29 teams: 20 yes and 9 no  (Univ. of Virginia, COS)

**Same Data, Different Conclusions**

Twenty-nine research teams were given the same set of soccer data and asked to determine if referees are more likely to give red cards to dark-skinned players. Each team used a different statistical method, and each found a different relationship between skin color and red cards.

Referees are **three times as likely** to give red cards to dark-skinned players

**Statistically significant** results showing referees are more likely to give red cards to dark-skinned players

Twice as likely

95% CONFIDENCE INTERVAL

ONE RESEARCH TEAM

**Equally likely**

**Non-significant** results

FIVETHIRTYEIGHT

SOURCE: BRIAN NOSEK ET AL.

---



**ASIST 2012 Book of the Year Award**

Modern Information Retrieval
the concepts and technology behind search
Second edition

Ricardo Baeza-Yates
Berthier Ribeiro-Neto

**Questions?**

**Biased Questions?**

Contact: **rbaeza@acm.org**

**www.baeza.cl**

**@polarbearby**