2/8/17

# Bias in the Web

## Ricardo Baeza-Yates

ntent

---

## Fake Content & Bias

- British Prime Minister Benjamin Disraeli:
  - "There are three kinds of **lies**: **lies**, damned **lies**, and **statistics**.

### UTC professor says "Everyone has bias"

BY HANNAH LAWRENCE  |  FRIDAY, JULY 8TH 2016

We all have biases and preconceptions about certain subjects or groups of people according to one Chattanooga researcher.

**Buzzfeed News**

*TOP POST*
**173,877 VIEWS**

### Here Are 50 Of The Biggest Fake News Hits On Facebook From 2016

One fake news entrepreneur says we should expect even more Trump hoaxes in 2017

posted on Dec. 30, 2016, at 2:12 p.m.

**Craig Silverman**
BuzzFeed News Media Editor

**Bias: significant deviation from a prior (unknown) distribution**

ntent

1

## (Observational) Human Data has Bias

### Goal: Bias Awareness

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

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

**Attempt of an unbiased (personal) view on bias in the Web**

**Many people extrapolate results of a sample to the whole population (e.g., social media analysis)**
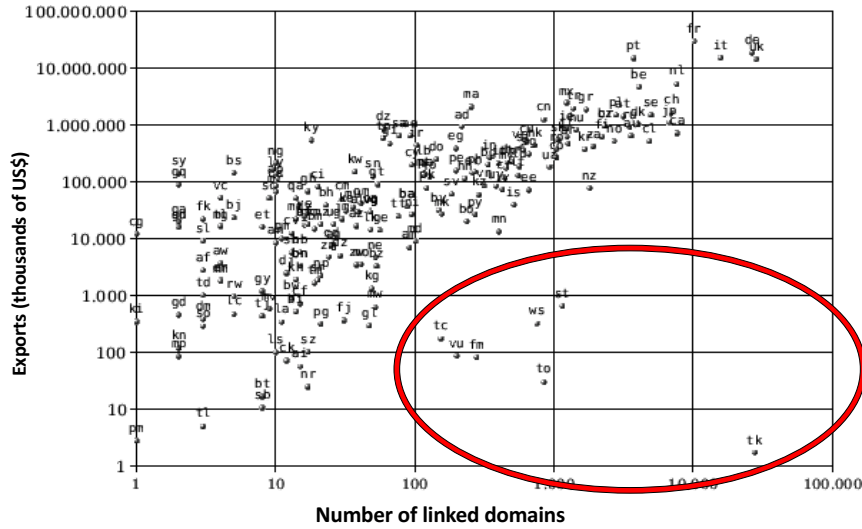
**In addition there is bias when measuring bias as well as bias towards measuring it!**

**⋒ntent**

---

## Bias in the Web

**Web**

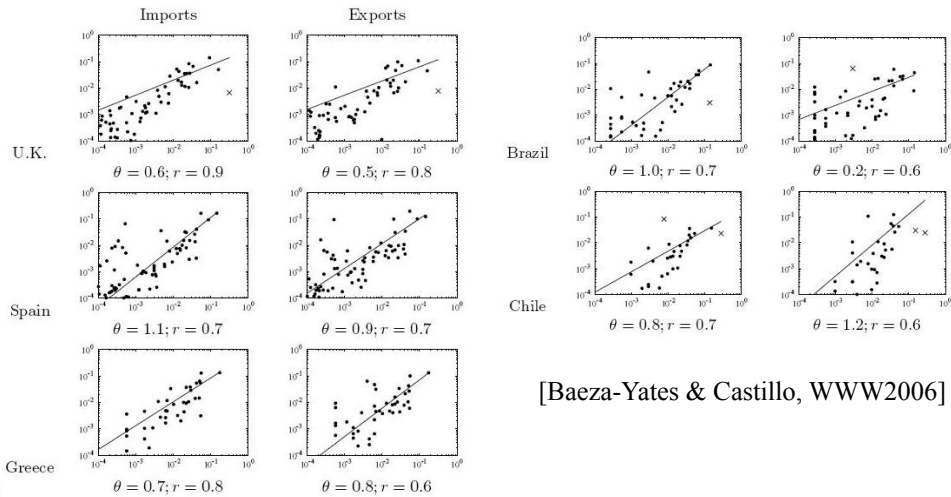**Data bias**

**⋒ntent**

## Economic Bias in Links



[Baeza-Yates, Castillo & López. Characteristics of the Web of Spain. Cybermetrics, 2005]

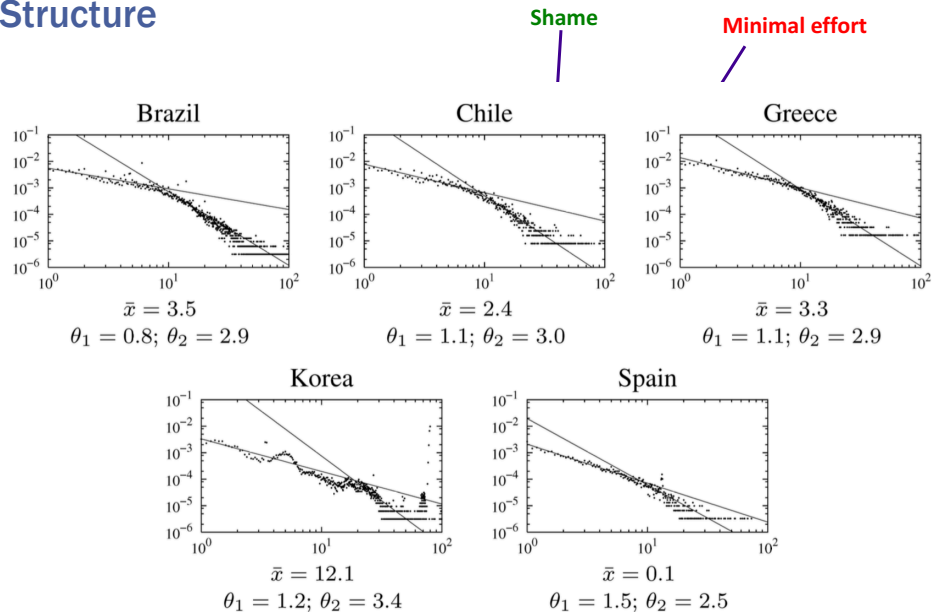## Economic Bias in Links
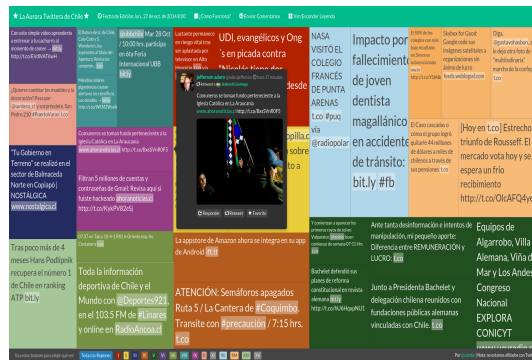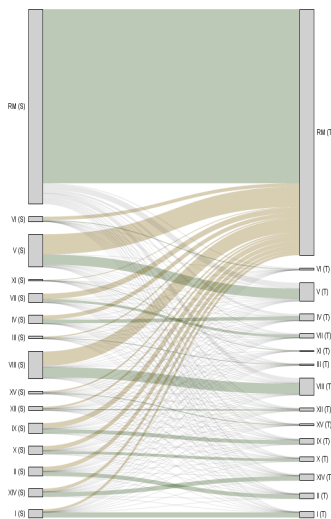


[Baeza-Yates & Castillo, WWW2006]

# 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]

ntent

7

# Geographical Bias in Content



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

ntent

## Gender Bias in Content

- Word embedding's in w2vNEWS

### Gender stereotype *she-he* analogies.

| | | |
|---|---|---|
| sewing-carpentry | register-nurse-physician | housewife-shopkeeper |
| nurse-surgeon | interior designer-architect | softball-baseball |
| blond-burly | feminism-conservatism | cosmetics-pharmaceuticals |
| giggle-chuckle | vocalist-guitarist | petite-lanky |
| sassy-snappy | diva-superstar | charming-affable |
| volleyball-football | cupcakes-pizzas | hairdresser-barber |

### Gender appropriate *she-he* analogies.

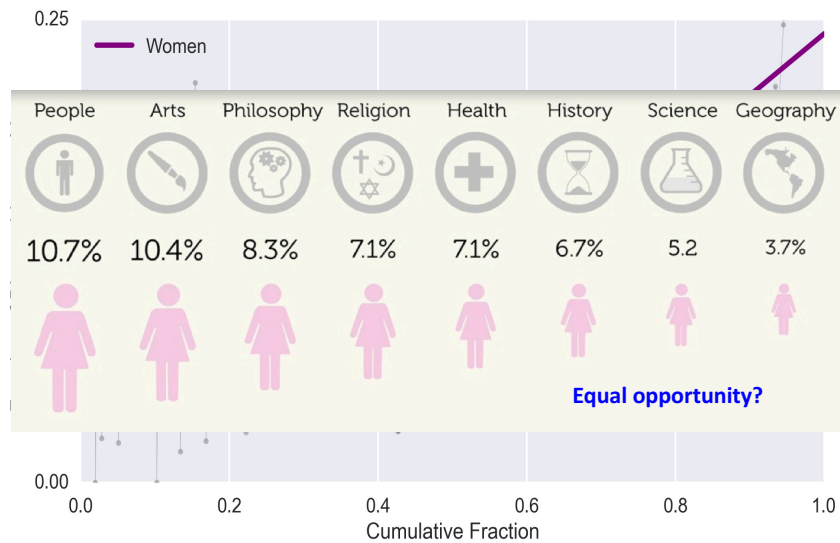| | | |
|---|---|---|
| queen-king | sister-brother | mother-father |
| waitress-waiter | ovarian cancer-prostate cancer | convent-monastery |

**Most journalists are men?**

[Bolukbasi at al, ArXiv 2016]

**Yes, about 60 to 70% at work although at college is the inverse**

**ntent**

---

## Gender Bias in Content

**Systemic bias?**



**Equal opportunity?**

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

Cumulative Fraction

**ntent**

[E. Graells-Garrido et al,. "First Women, Second Sex: Gender Bias in Wikipedia", ACM Hypertext'15]

## Bias in the Web

**Activity bias**

**Web**

**Data bias**

**ntent**

## Activity Bias

Which percentage of users produce 50% 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]

**ntent**

**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
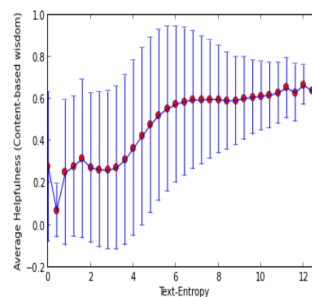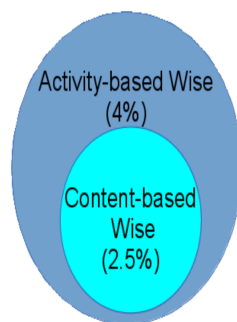
# Amazon Continues Their Crusade Against Fake Reviews

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

---

## Quality of Content?

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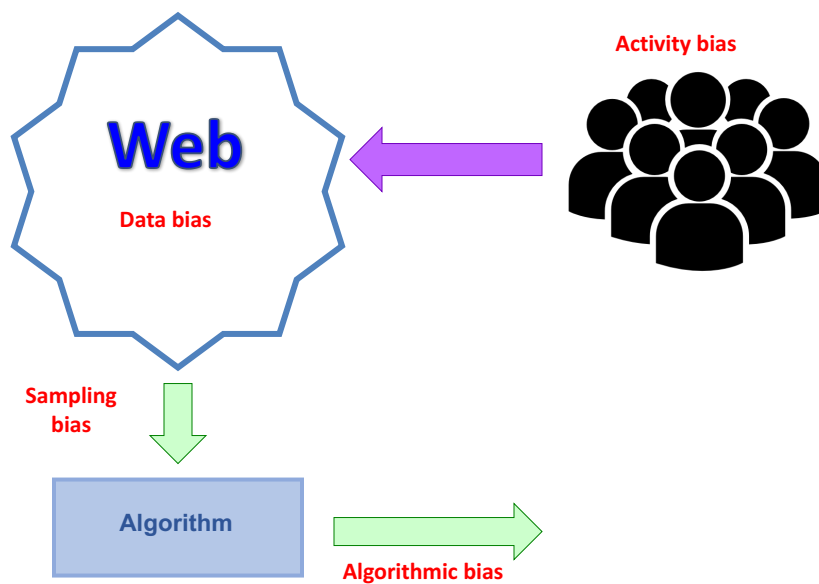- How many of those users are being paid?



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

## Content that is never seen: 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.

ntent

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

---

**Activity bias**

**Web**

**Data bias**

**Sampling bias**

**Algorithm**

**Algorithmic bias**

ntent

## 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 **but not for p near 0**

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

$$n \geq Z^2_{1-\alpha/2}\left(\frac{p'(1-p')}{\epsilon^2} - 1\right)$$

  where *Z* is the inverse of the standard normal distribution, $1-\alpha$ is the confidence interval and $p' = p + Z^2/2$
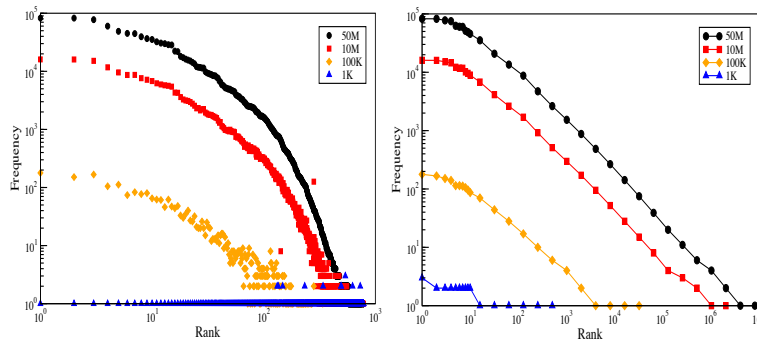
- If *p* = 0.1, $1-\alpha$ is 90% and $\epsilon$ is 10%, we get *n* = 2342. The standard formula gives *n* = 900!

**n**tent

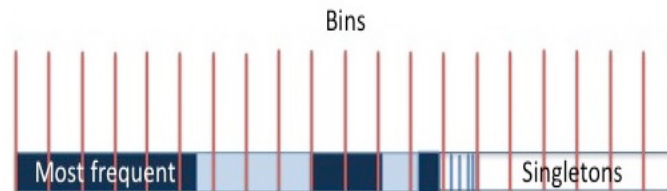[Baeza-Yates, SIGIR 2015, Industry track]

## Sampling Techniques

- Standard technique: $p_q \approx \widehat{p}_q(\mathcal{S}) = \dfrac{f_q(\mathcal{S})}{\sum_{q' \in \mathcal{S}} f_{q'}(\mathcal{S})}$

- A good sample should cover well all the query distribution but this does not work with very biased distributions.



**n**tent

[Zaragoza et al, CIKM 2010]

## Incremental Stratified Sampling

- Main goal: make good samples consistent across time
- Simple idea based in stratified sampling: bins + random start point
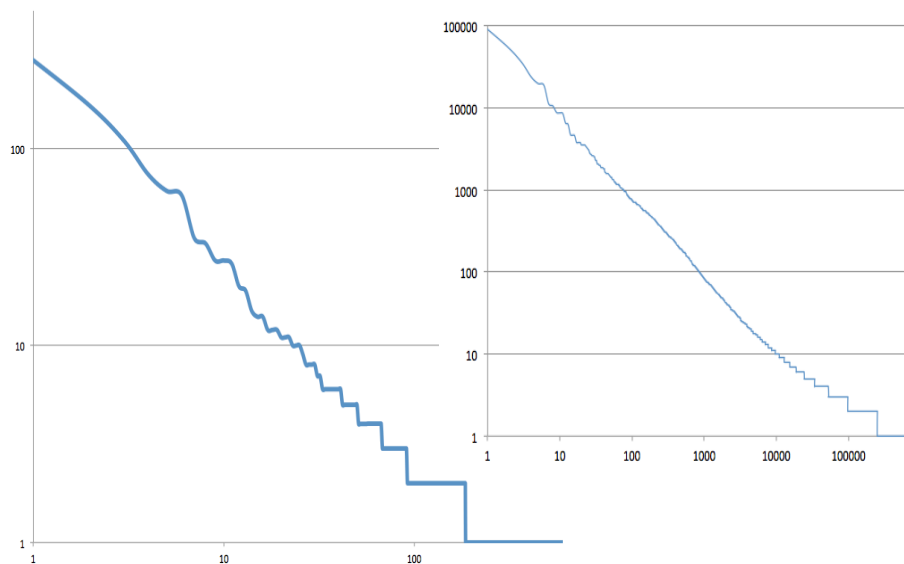
Bins



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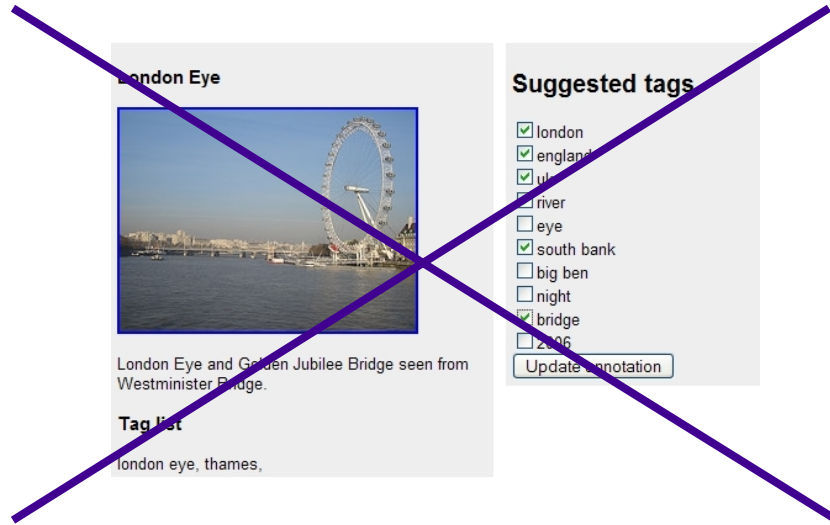- Change the bins in the tail to get the right distribution

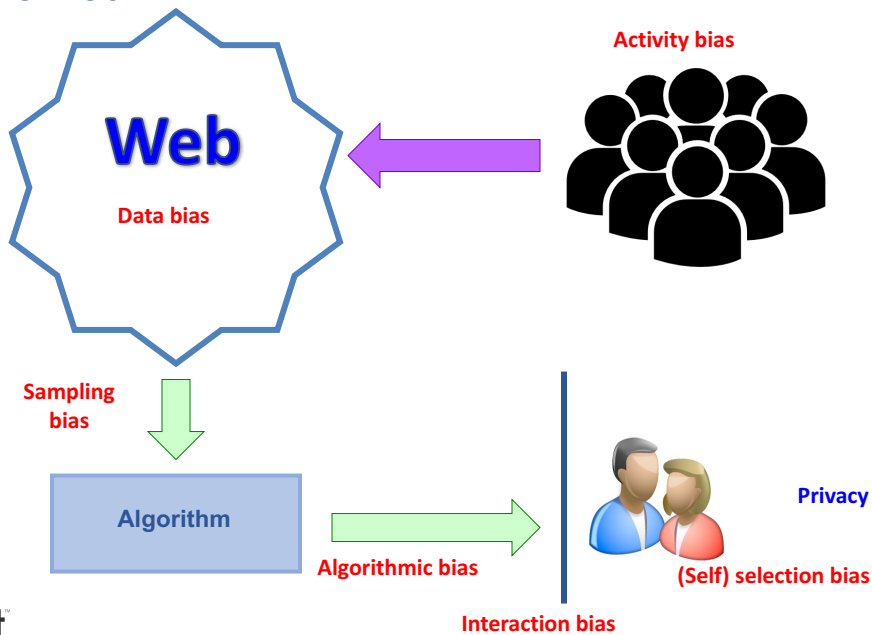[Baeza-Yates, SIGIR 2015, Industry track]  20

## Stratified Sampling Example



22

## Extreme Algorithmic Bias



## Bias in the Web

# Bias in the Interaction

Related Searches: tennis racket, tennis shoes.

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

**Presentation bias**

Shop by Category

| Tennis Equipment | Tennis Games | Kids' Sports | Clothing, Shoes & Jewelry | Tennis - Books |

Sponsored ⓘ

Tennis Elbow Brace with Gel Comp…
$24.50 ✔Prime
★★★★★ 7

Wilson Sporting Goods Championship Extra Duty Tennis Balls (1-Can)
Jun 14, 2012
by Wilson
$2.79 $6.99 [Add-on Item]
Add to a qualifying order to get it by **Tomorrow, May 6**.
More Buying Choices
$0.99 new (18 offers)
$7.99 used (2 offers)
See newer version

★★★★★ ▾ 186

**Sports & Outdoors:** See all 60,449 items

**Social bias**

DIMANKA Professional Table Tenni…
$34.99
★★★★★ 9

[Best Seller]
Wilson 75 Tennis Ball Pick Up Hopper
by Wilson

**Interaction bias**

$19.96 ✔Prime
Get it by **Tomorrow, May 6**

More Buying Choices
$18.88 new (11 offers)
$35.00 used (1 offer)

★★★★☆ ▾ 319

**Product Features**
Holds 75 tennis balls with a special no spill lid (Tennis Balls NOT included)

**Sports & Outdoors:** See all 60,449 items

Gamma Quick Kids 78 Ball (12 Pac…
$19.99 ✔Prime
★★★☆☆ 44

ⓝntent

---

# Dependencies: A Cascade of Biases!

**Position bias** → **Ranking bias**

**Presentation bias** → **Interaction bias** → **Click bias**

**Mouse movement bias**

**Social bias**

**Scrolling bias**

Data & algorithmic bias | Self-selection bias

ⓝntent™

## Social Bias


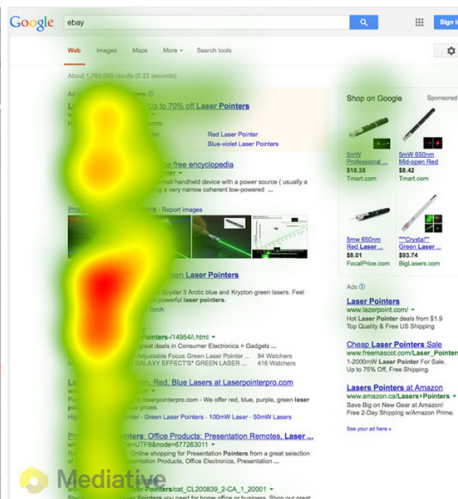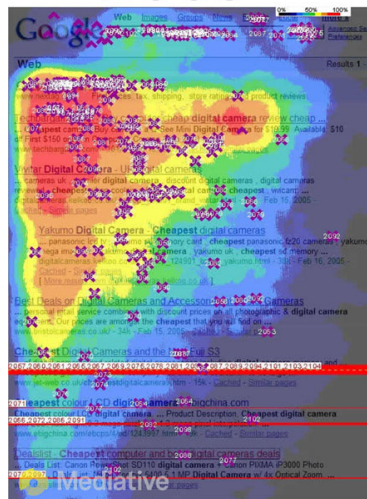
[WHY AMAZON'S RATINGS MIGHT MISLEAD YOU; The Story of Herding Effects
Ting Wang and Dashun Wang, Big Data, 2014]
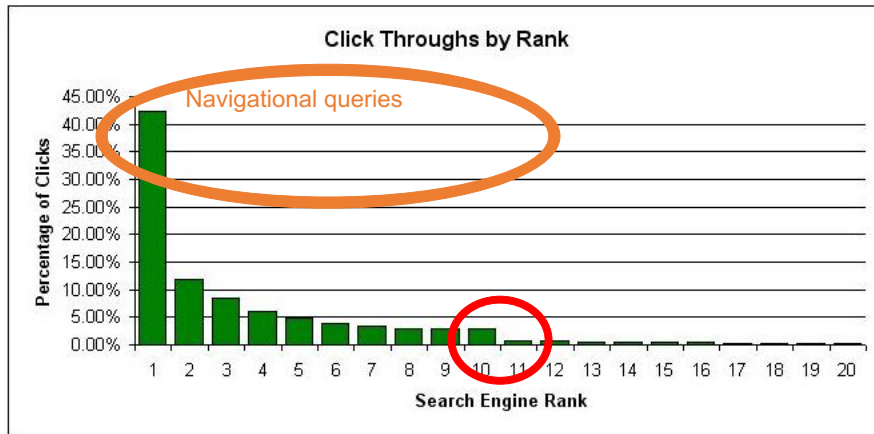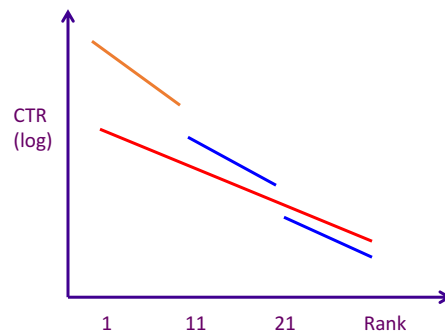
## Ranking Bias in Web Search



[Mediative Study, 2014]
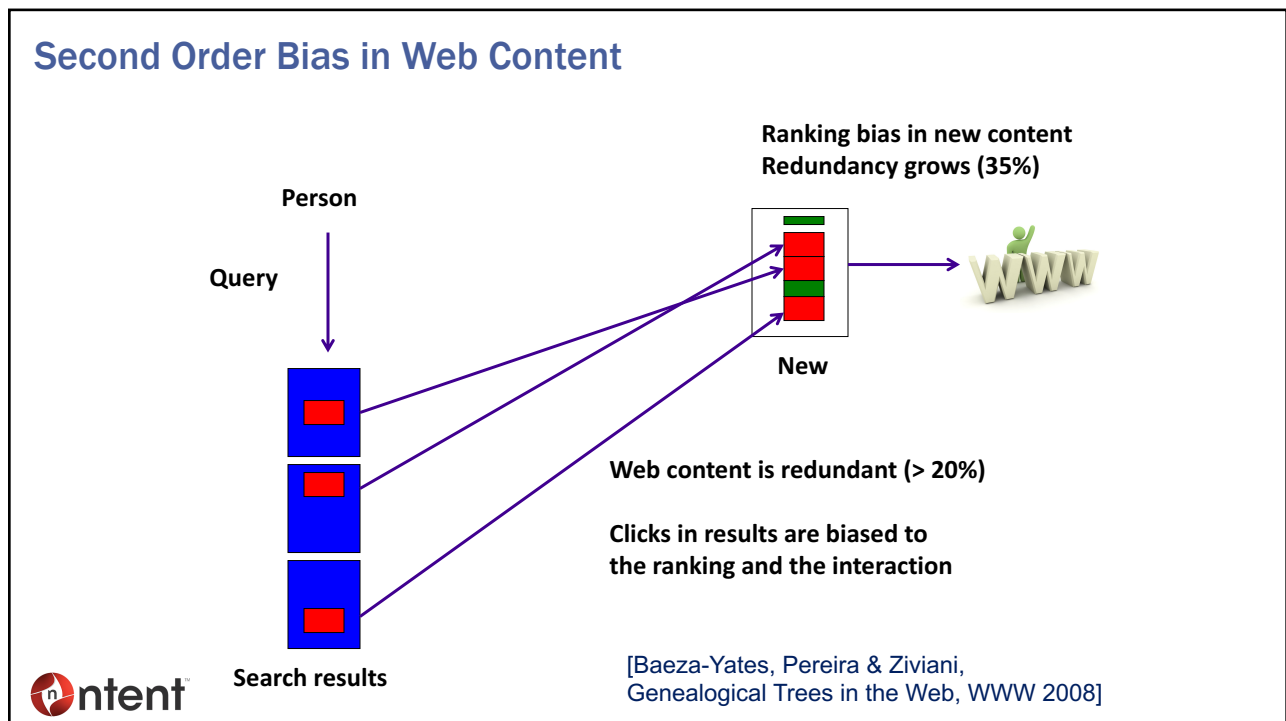
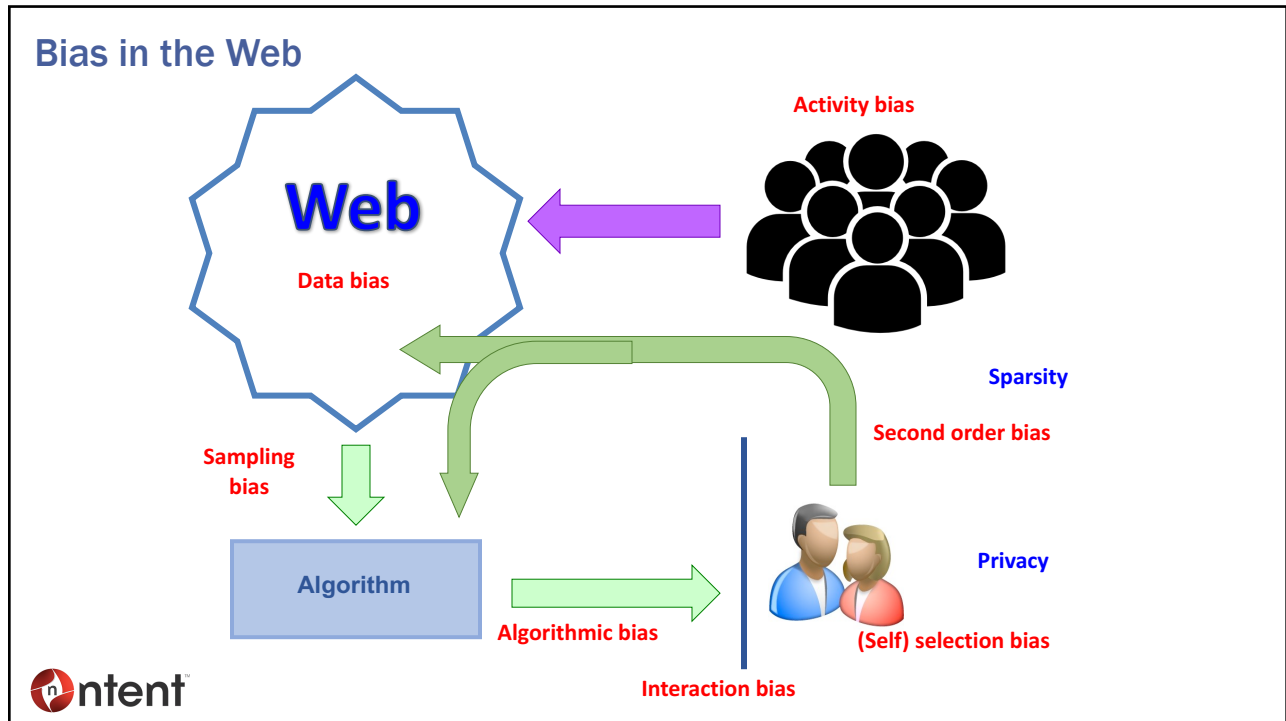# Click Bias in Web Search

○ Ranking & next page bias



**ntent**

---

# Unbiasing Search Clicks

Clicks as implicit positive user feedback



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

**ntent**

Bias in the Web

**Activity bias**

**Web**

**Data bias**

**Sparsity**

**Second order bias**

**Sampling bias**

**Privacy**

**Algorithm**

**Algorithmic bias**

**(Self) selection bias**

**Interaction bias**



Second Order Bias in Web Content

**Ranking bias in new content**
**Redundancy grows (35%)**

**Person**

**Query**

**New**

**Web content is redundant (> 20%)**

**Clicks in results are biased to the ranking and the interaction**

**Search results**

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

## Avoid Second Order Bias due to Personalization
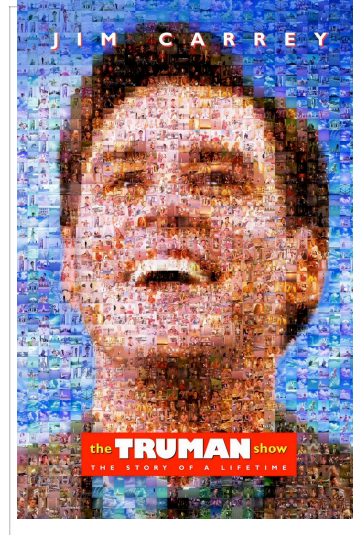
The Filter "Bubble", Eli Pariser (2011)
• The effect of self selection bias
• Avoid the poor get poorer syndrome
• Avoid the echo chamber
• Empower the tail

**Partial solutions:**
• Diversity
• Novelty
• Serendipity
• Show me the dark side

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

**How much exploration is needed for presentation bias?**
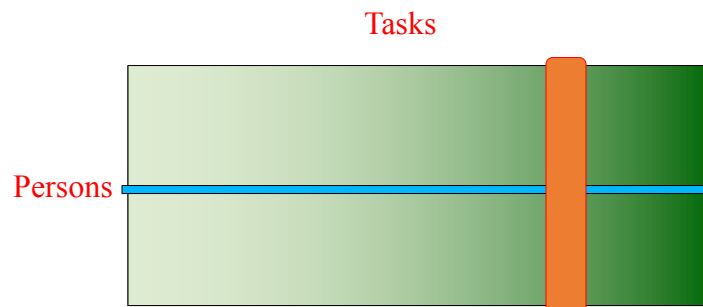


---

## Aggregating in the Tail

- Exploit the context (and deep learning!)

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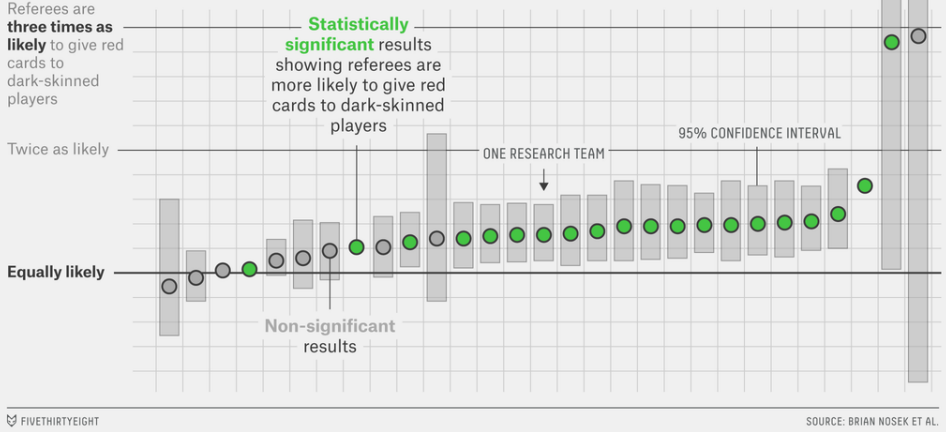  [Baeza-Yates et al, WSDM 2015]

- Personalization vs. Contextualization

  Recall that user interaction is another long tail



Tasks

Persons

---

## It's Hard to Get the Truth from Data (Professional Bias)



**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.
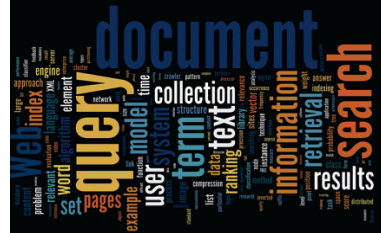
➔ 61 analysts, 29 teams: 20 yes and 9 no  (Univ. of Virginia, COS)
➔ **We need to focus on small data, not big data**

---

## Questions?

**ASIST 2012 Book of the Year Award (Biased Ad)**



Modern **Information Retrieval**
the concepts and technology behind search
Second edition

Ricardo Baeza-Yates
Berthier Ribeiro-Neto

Contact: **rbaeza@acm.org**
**www.baeza.cl**
**@polarbearby**

## Biased Questions?
**More bias: We are hiring in Barcelona & San Diego!**