

# Introduction and Overview

APAM E4990  
Modeling Social Data

Jake Hofman

Columbia University

January 25, 2019

# Course overview

Modeling social data requires an understanding of:

- ① How to obtain data produced by (online) human interactions
- ② What questions we typically ask about human-generated data
- ③ How to make these questions precise and quantitative
- ④ How to interpret and communicate results

# Questions

Many long-standing questions in the social sciences are notoriously difficult to answer, e.g.:

- “Who says what to whom in what channel with what effect”? (Laswell, 1948)
- How do ideas and technology spread through cultures? (Rogers, 1962)
- How do new forms of communication affect society? (Singer, 1970)
- ...

# Questions

Typically difficult to observe the relevant information via conventional methods

## EMOTIONS MAPPED BY NEW GEOGRAPHY

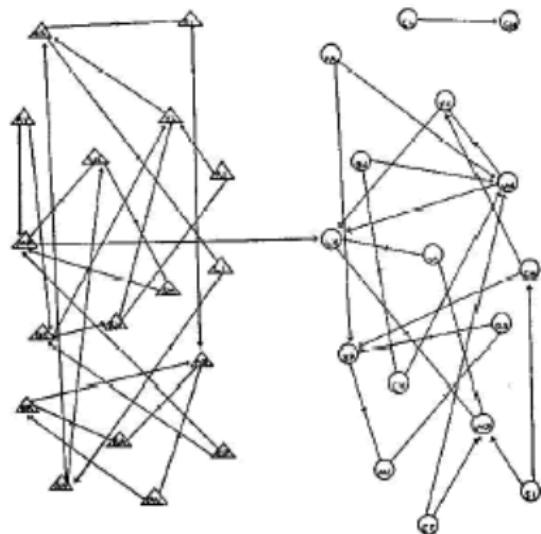
Charts Seek to Portray the Psychological Currents of Human Relationships.

### FIRST STUDIES EXHIBITED

Colored Lines Show Likes and Dislikes of Individuals and of Groups.

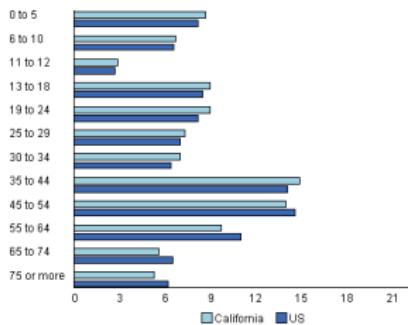
### MANY MISFITS REVEALED

Moreno, 1933

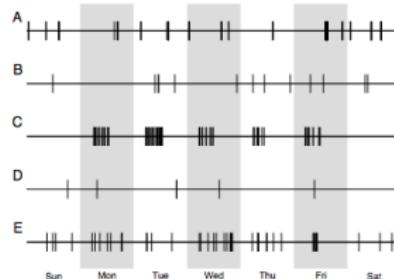


# Large-scale data

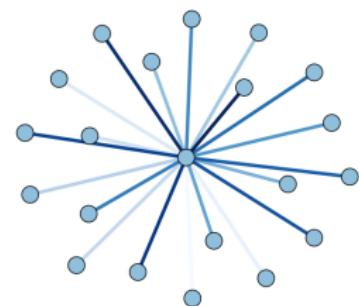
Recently available electronic data provide an unprecedented opportunity to address these questions at scale



Demographic



Behavioral



Network

# Computational social science

An emerging discipline at the intersection of the social sciences,  
statistics, and computer science

# Computational social science

An emerging discipline at the intersection of the **social sciences**,  
statistics, and computer science

(motivating questions)

# Computational social science

An emerging discipline at the intersection of the social sciences,  
statistics, and computer science

(fitting large, potentially sparse models)

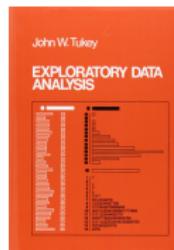
# Computational social science

An emerging discipline at the intersection of the social sciences,  
statistics, and computer science

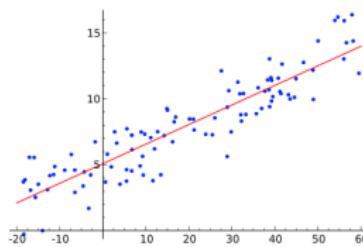
(parallel processing for filtering and aggregating data)

# Topics

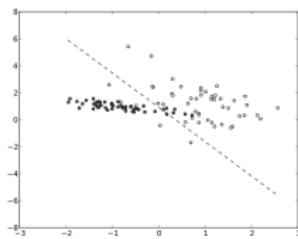
## Exploratory Data Analysis



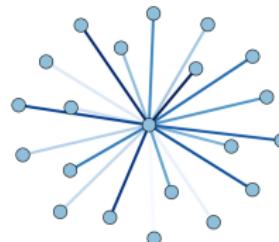
## Regression



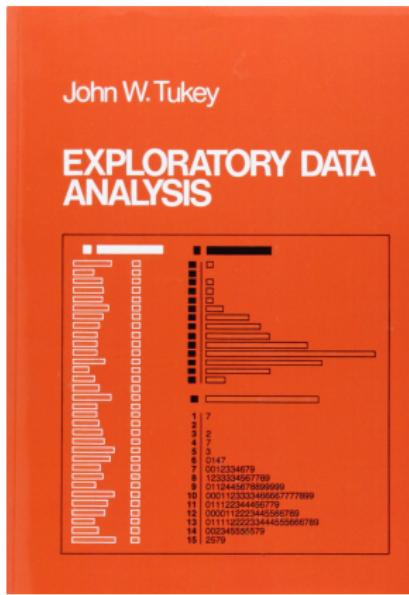
## Classification



## Networks

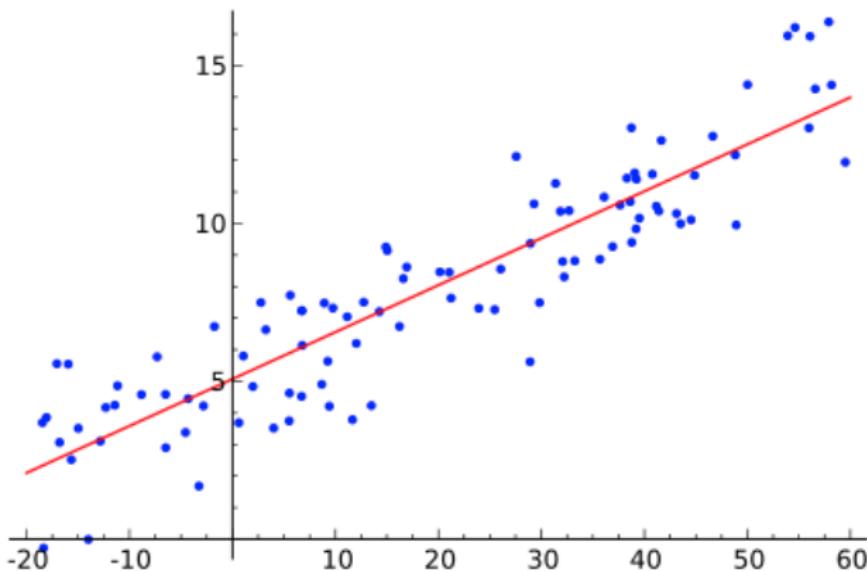


# Exploratory Data Analysis



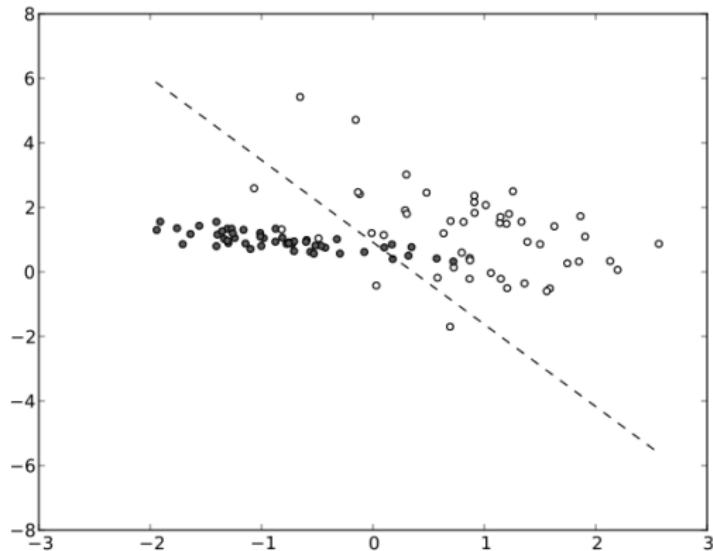
(a.k.a. counting and plotting things)

# Regression



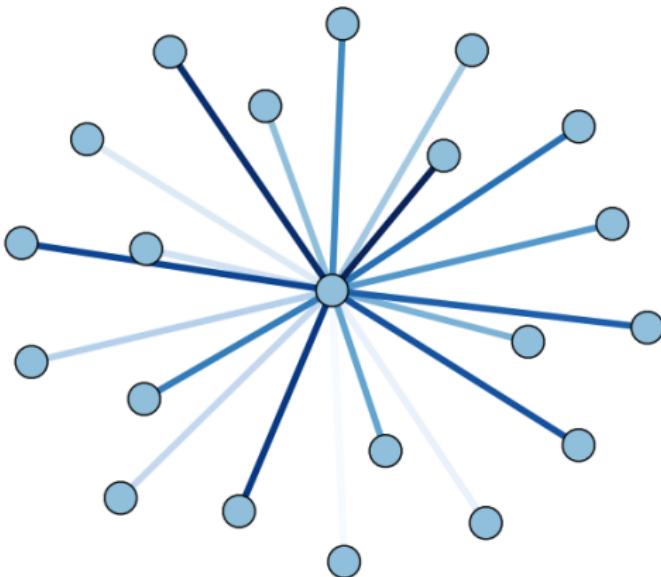
(a.k.a. modeling continuous things)

# Classification



(a.k.a. modeling discrete things)

# Networks



(a.k.a. counting complicated things)

# Prediction and explanation

Important to view prediction and explanation as compliments,  
not substitutes

Computer science

$$\hat{y}$$

Predict

vs  
and

Social science

$$\hat{\beta}$$

Explain

Otherwise it can be difficult to make long-term progress in advancing social science

# The clean real story

*"We have a habit in writing articles published in scientific journals to make the work as finished as possible, to cover all the tracks, to not worry about the blind alleys or to describe how you had the wrong idea first, and so on. So there isn't any place to publish, in a dignified manner, what you actually did in order to get to do the work ..."*

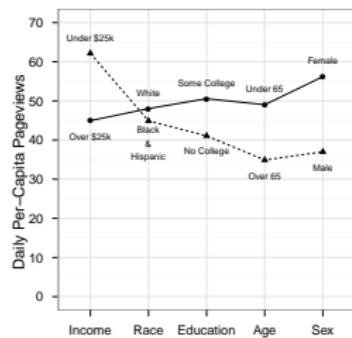
-Richard Feynman  
*Nobel Lecture<sup>1</sup>, 1965*

---

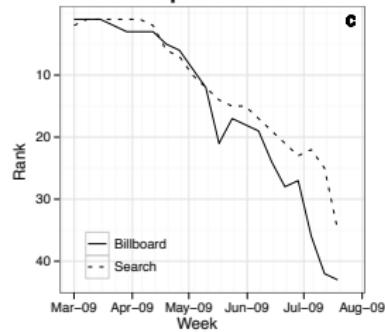
<sup>1</sup><http://bit.ly/feynmannobel>

# Case studies

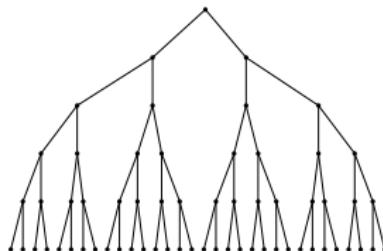
## Web demographics



## Search predictions

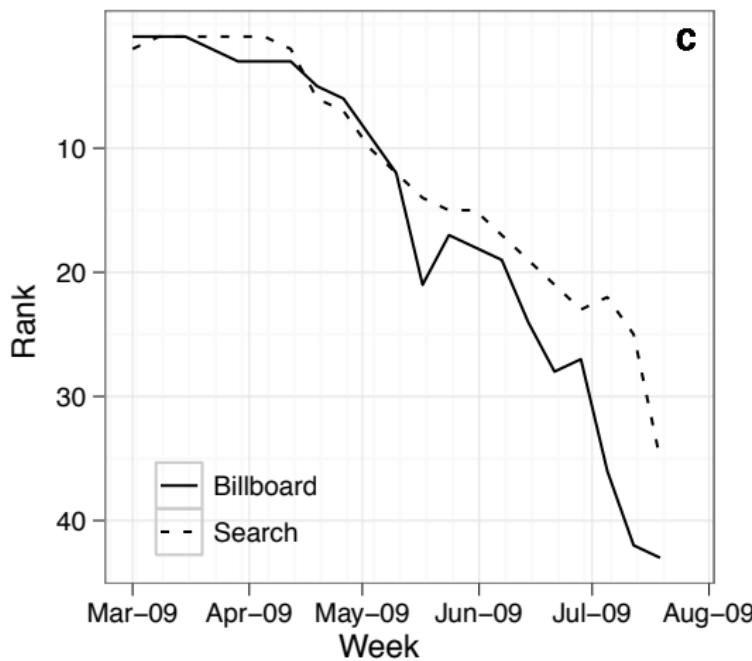


## Viral hits



# Predicting consumer activity with Web search

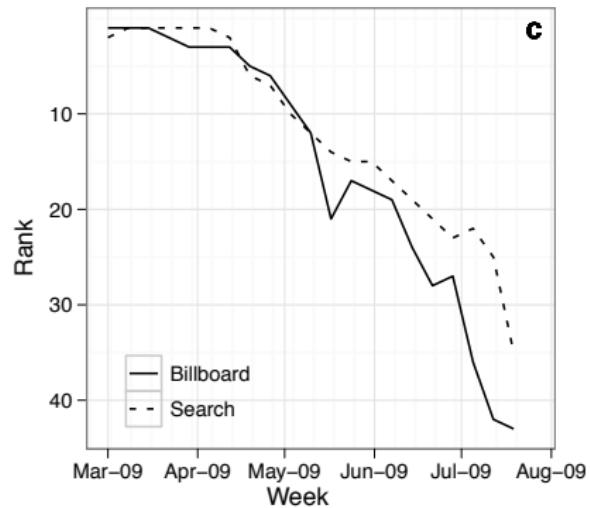
with Sharad Goel, Sébastien Lahaie, David Pennock, Duncan Watts



# Search predictions

## Motivation

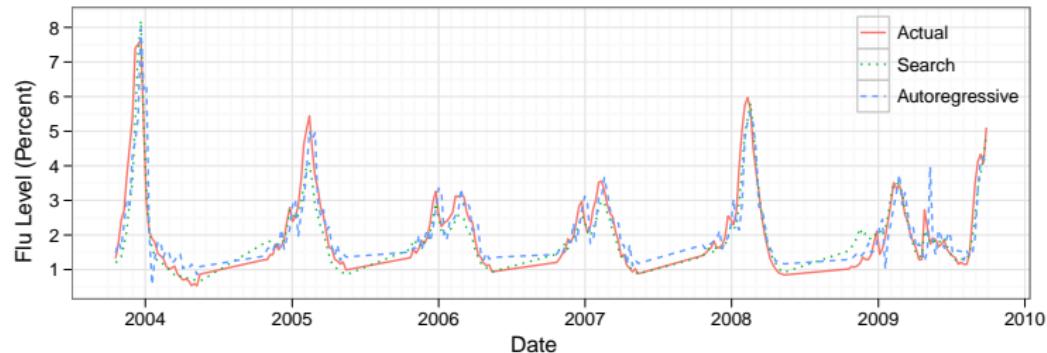
Does collective search activity provide useful predictive signal about real-world outcomes?



# Search predictions

## Motivation

Past work mainly focuses on predicting the present<sup>2</sup> and ignores baseline models trained on publicly available data

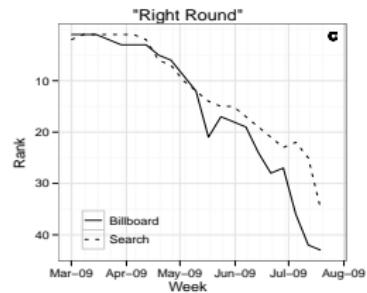
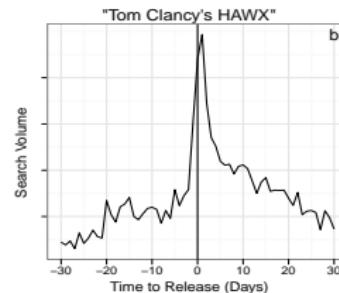
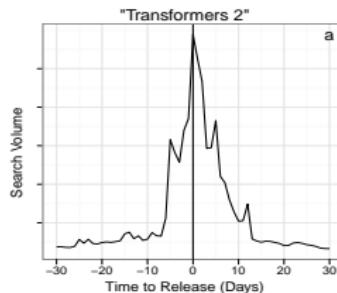


<sup>2</sup>Varian, 2009

# Search predictions

## Motivation

We predict future sales for movies, video games, and music



# Search predictions

## Search models

For movies and video games, predict opening weekend box office and first month sales, respectively:

$$\log(\text{revenue}) = \beta_0 + \beta_1 \log(\text{search}) + \epsilon$$

For music, predict following week's Billboard Hot 100 rank:

$$\text{billboard}_{t+1} = \beta_0 + \beta_1 \text{search}_t + \beta_2 \text{search}_{t-1} + \epsilon$$

## Search predictions

## Search volume

Web Images Video Local Shopping News More ▾

no country  Options ▾

 QuickApps

 SafeSearch - On

509,000,000 results for no country:

 Show All

 Wikipedia

 IMDb

 MySpace

 NY Daily News

 GameSpot

 Sponsored Results

Also try: [no country for old men](#), [no country for old men ending](#), [more...](#)

**No Country for Old Men (film) - Wikipedia, the ...**  
[Plot](#) | [Cast and characters](#) | [Themes and style](#) | [Production](#)  
No Country for Old Men is a 2007 American crime thriller directed by Joel Coen and Ethan Coen, and starring Tommy Lee Jones, Javier Bardem, and Josh Brolin. The film was adapted from...  
[en.wikipedia.org/wiki/No\\_Country\\_for\\_Old\\_Men\\_\(film\)](http://en.wikipedia.org/wiki/No_Country_for_Old_Men_(film)) - Cached



**No Country for Old Men (2007) - IMDb**  
Violence and mayhem ensue after a hunter stumbles upon some dead bodies, a stash of heroin and more than \$2 million in cash near the Rio Grande. With Tommy Lee Jones ...  
[www.imdb.com/title/tt0477348](http://www.imdb.com/title/tt0477348) Cached

**no country | Free Music, Tour Dates, Photos, Videos**  
no country's official profile including the latest music, albums, songs, music videos and more updates.  
[www.myspace.com/nocountrytheband](http://www.myspace.com/nocountrytheband) - Cached

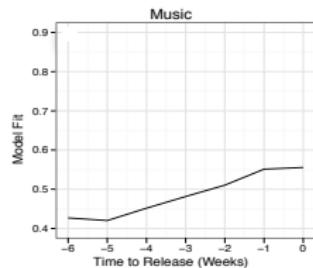
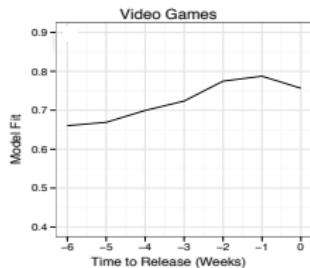
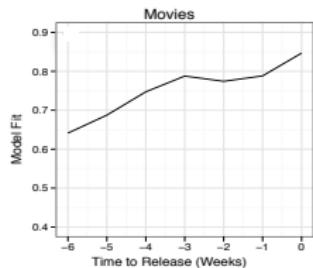
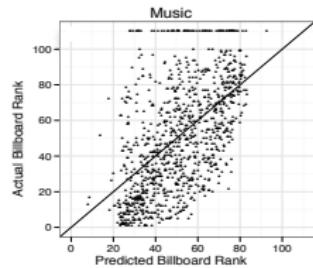
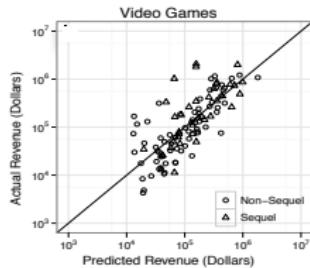
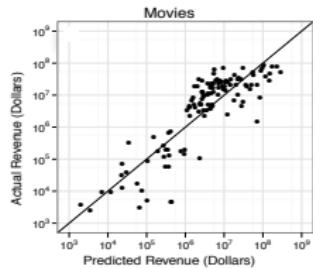
**No Country - Video Results**


# Search predictions

## Search models

Search activity is **predictive** for movies, video games, and music weeks to months in advance



# Search predictions

## Baseline models

For movies, use budget, number of opening screens and Hollywood Stock Exchange:

$$\log(\text{revenue}) = \beta_0 + \beta_1 \log(\text{budget}) + \beta_2 \log(\text{screens}) + \beta_3 \log(\text{hsx}) + \epsilon$$

# Search predictions

## Baseline models

For video games, use critic ratings and predecessor sales (sequels only):

$$\log(\text{revenue}) = \beta_0 + \beta_1 \text{rating} + \beta_2 \log(\text{predecessor}) + \epsilon$$

# Search predictions

## Baseline models

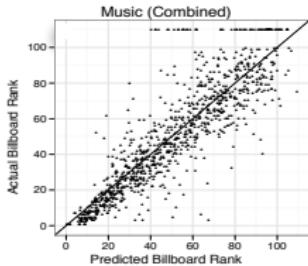
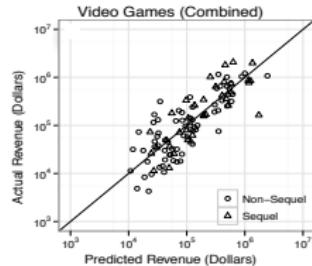
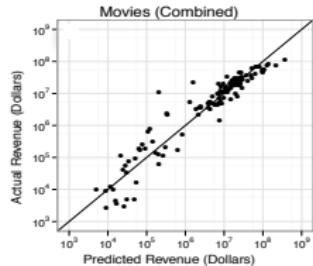
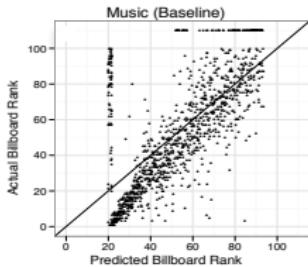
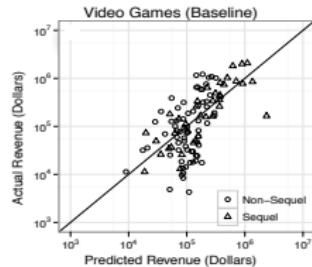
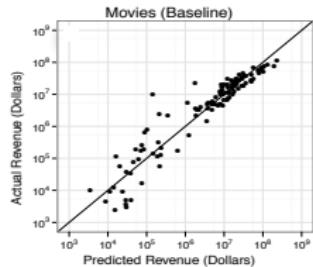
For **music**, use an autoregressive model with the previously available rank:

$$\text{billboard}_{t+1} = \beta_0 + \beta_1 \text{billboard}_{t-1} + \epsilon$$

# Search predictions

Baseline + combined models

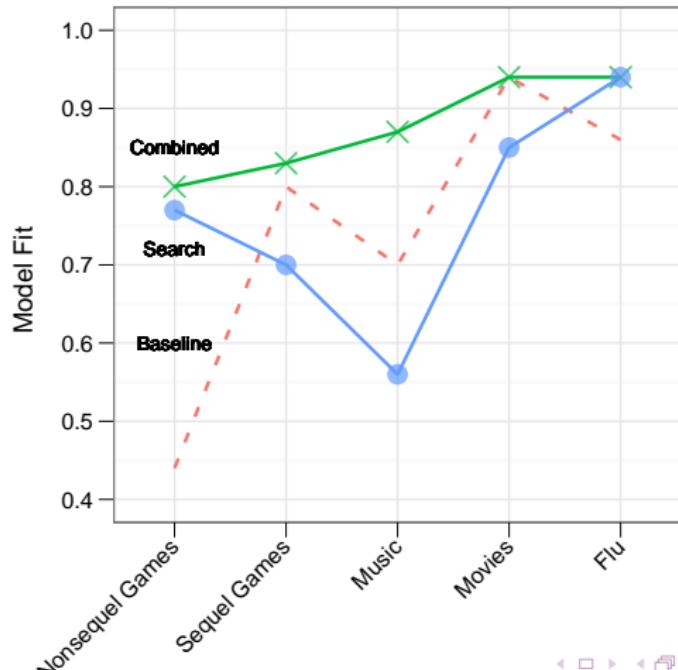
Baseline models are often surprisingly good



# Search predictions

## Model comparison

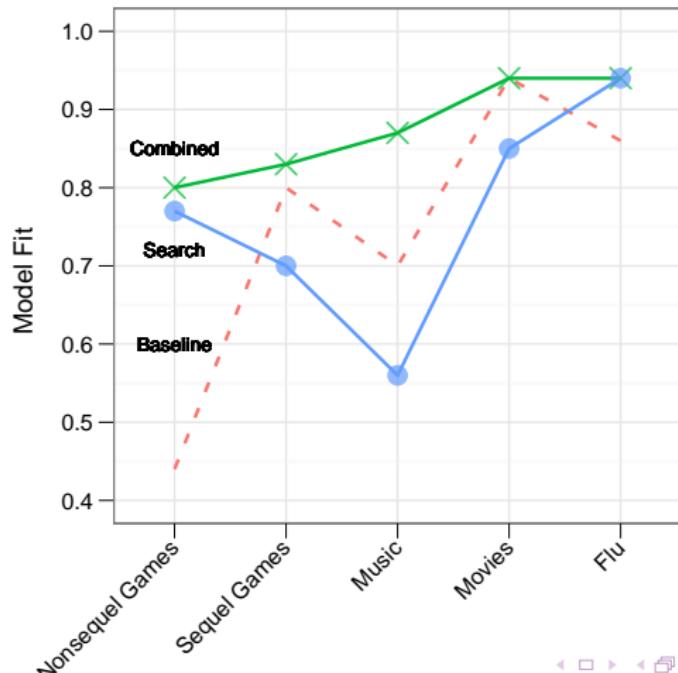
For movies, search is outperformed by the baseline and of little marginal value



# Search predictions

## Model comparison

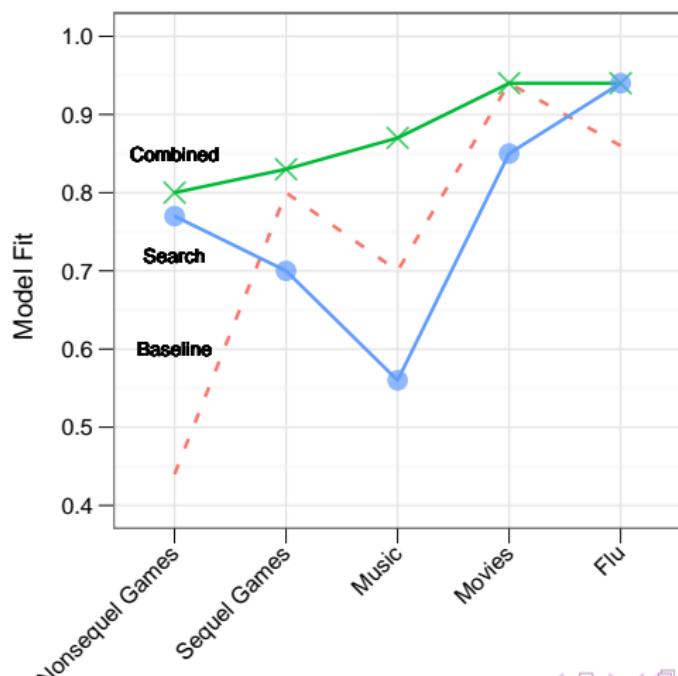
For video games, search helps substantially for non-sequels, less so for sequels



# Search predictions

## Model comparison

For music, the addition of search yields a substantially better combined model



# Search predictions

## Summary

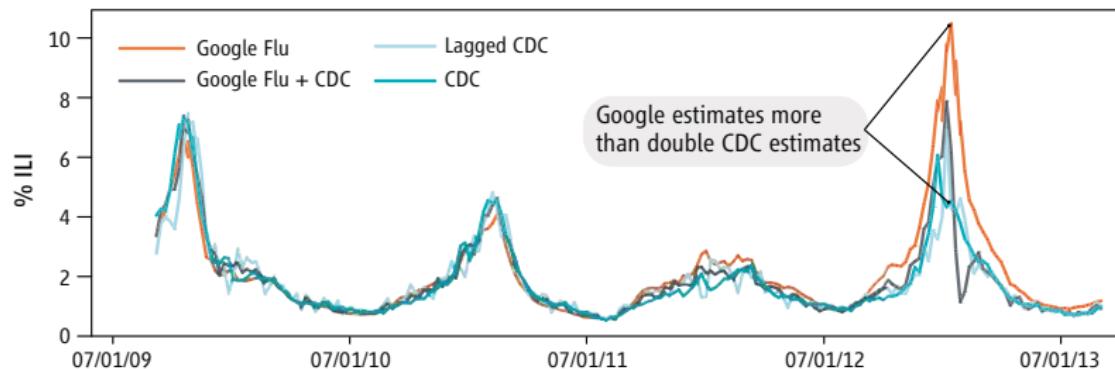
- Relative performance and **value** of search varies across domains
- Search provides a **fast**, **convenient**, and **flexible signal** across domains
- “Predicting consumer activity with Web search”  
Goel, Hofman, Lahaie, Pennock & Watts, PNAS 2010

## BIG DATA

# The Parable of Google Flu: Traps in Big Data Analysis

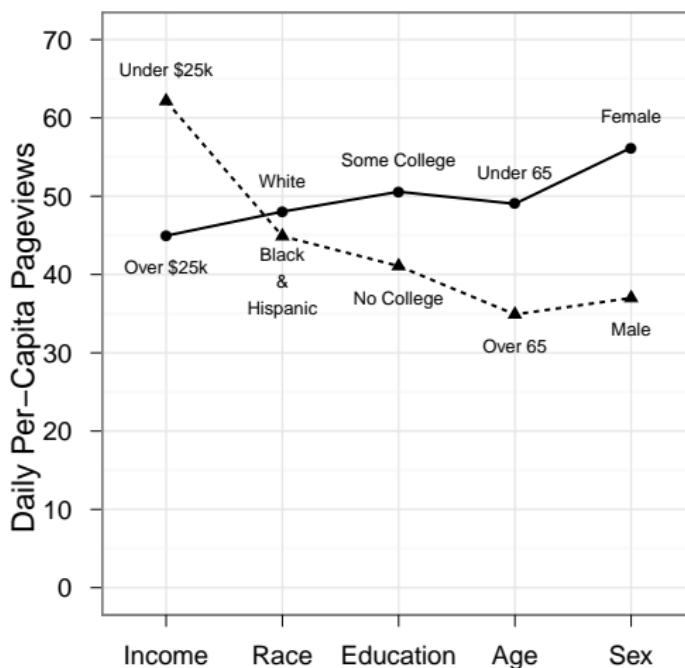
David Lazer,<sup>1,2\*</sup> Ryan Kennedy,<sup>1,3,4</sup> Gary King,<sup>3</sup> Alessandro Vespignani<sup>3,5,6</sup>

Large errors in flu prediction were largely avoidable, which offers lessons for the use of big data.



# Demographic diversity on the Web

with Irmak Sirer and Sharad Goel (ICWSM 2012)



# Motivation

Science 17 April 1998:  
Vol. 280 no. 5362 pp. 390-391  
DOI: 10.1126/science.280.5362.390

< Prev | Table of Contents | Next >

POLICY

INFORMATION ACCESS

## Bridging the Racial Divide on the Internet

Donna L. Hoffman and Thomas P. Novak

 Author Affiliations

The Internet is expected to do no less than transform society (1); its use has been increasing exponentially since 1994 (2). But are all members of our society equally likely to have access to the Internet and thus participate in the rewards of this transformation? Here we present findings both obvious and surprising from a recent survey of Internet access and discuss their implications for social science research and public policy.

Previous work is largely **survey-based** and focuses on group-level differences in online **access**

# Motivation

*"As of January 1997, we estimate that 5.2 million African Americans and 40.8 million whites have ever used the Web, and that 1.4 million African Americans and 20.3 million whites used the Web in the past week."*

-Hoffman & Novak (1998)

# Motivation

Focus on activity instead of access



How diverse is the Web?

To what extent do online experiences vary across demographic groups?

## nielsen MegaPanel

- Representative sample of **265,000 individuals** in the US, paid via the Nielsen MegaPanel<sup>3</sup>
- Log of **anonymized, complete browsing activity** from June 2009 through May 2010 (URLs viewed, timestamps, etc.)
- Detailed individual and household **demographic information** (age, education, income, race, sex, etc.)

---

<sup>3</sup>Special thanks to Mainak Mazumdar

# Data

```
# ls -alh nielsen_megapanel.tar
-rw-r--r-- 100G Jul 17 13:00 nielsen_megapanel.tar
```

# Data

```
# ls -alh nielsen_megapanel.tar  
-rw-r--r-- 100G Jul 17 13:00 nielsen_megapanel.tar
```

- Normalize pageviews to at most three domain levels, sans www  
e.g. www.yahoo.com → yahoo.com,  
us.mg2.mail.yahoo.com/neo/launch → mail.yahoo.com

# Data

```
# ls -alh nielsen_megapanel.tar  
-rw-r--r-- 100G Jul 17 13:00 nielsen_megapanel.tar
```

- Normalize pageviews to at most three domain levels, sans www  
e.g. www.yahoo.com → yahoo.com,  
us.mg2.mail.yahoo.com/neo/launch → mail.yahoo.com
- Restrict to top 100k (out of 9M+ total) most popular sites  
(by unique visitors)

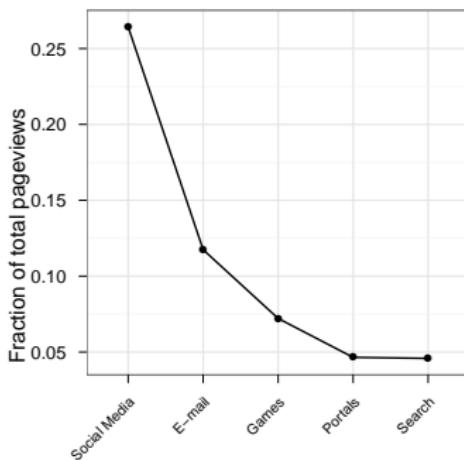
# Data

```
# ls -alh nielsen_megapanel.tar  
-rw-r--r-- 100G Jul 17 13:00 nielsen_megapanel.tar
```

- Normalize pageviews to at most three domain levels, sans www  
e.g. www.yahoo.com → yahoo.com,  
us.mg2.mail.yahoo.com/neo/launch → mail.yahoo.com
- Restrict to top 100k (out of 9M+ total) most popular sites  
(by unique visitors)
- Aggregate activity at the site, group, and user levels

# Aggregate usage patterns

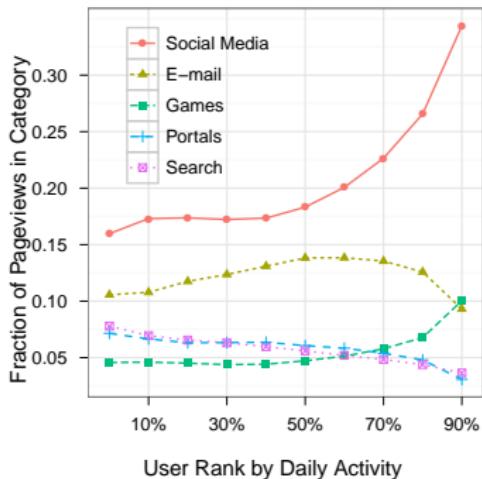
How do users distribute their time across different categories?



All groups spend the **majority** of their **time** in the top five most **popular** categories

# Aggregate usage patterns

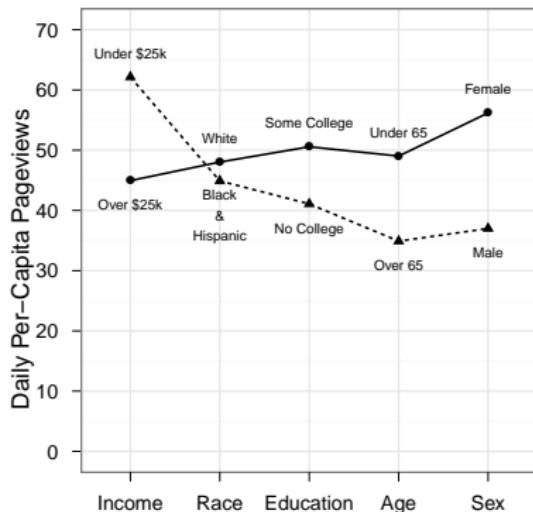
How do users distribute their time across different categories?



Highly active users devote nearly twice as much of their time to social media relative to typical individuals

# Group-level activity

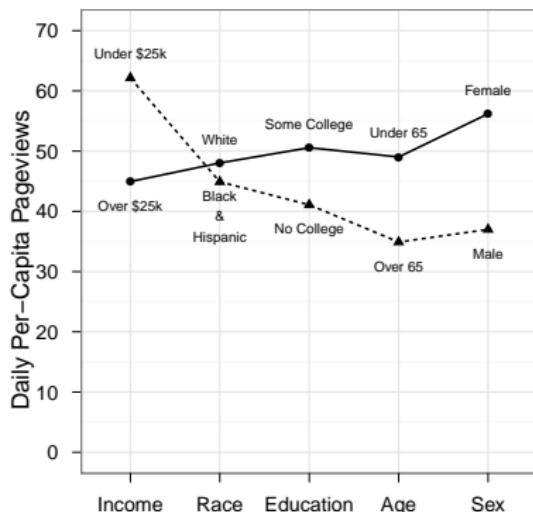
How does browsing activity vary at the group level?



Large differences exist even at the aggregate level  
(e.g. women on average generate 40% more pageviews than men)

# Group-level activity

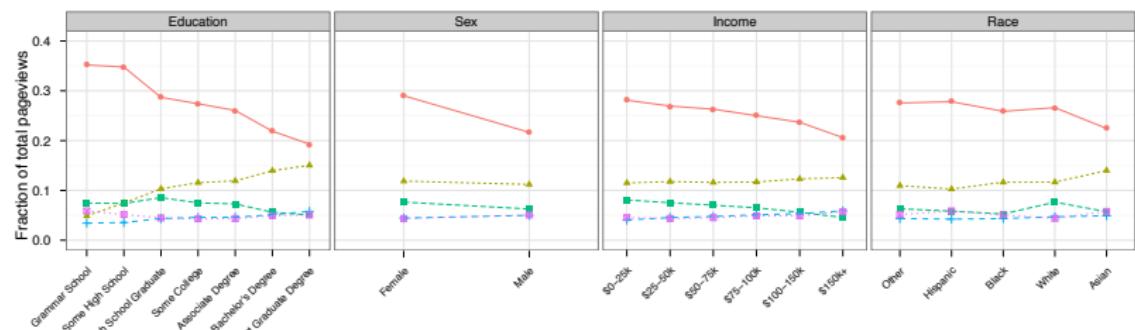
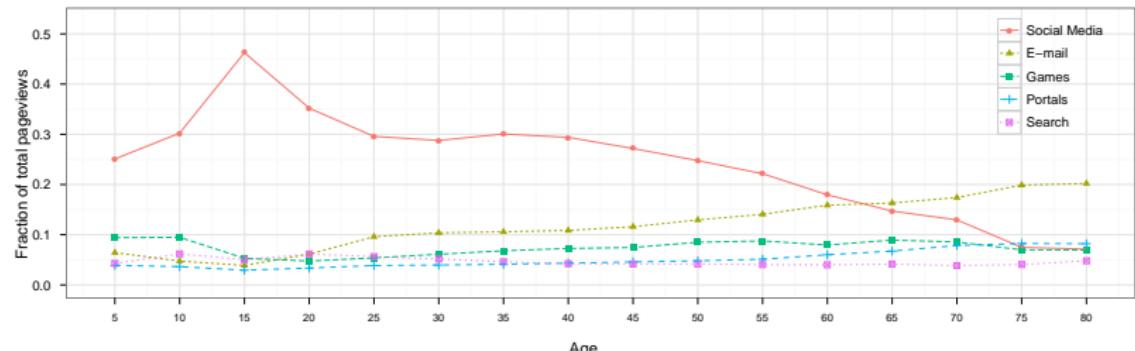
How does browsing activity vary at the group level?



Younger and more educated individuals are both more likely to access the Web and more active once they do

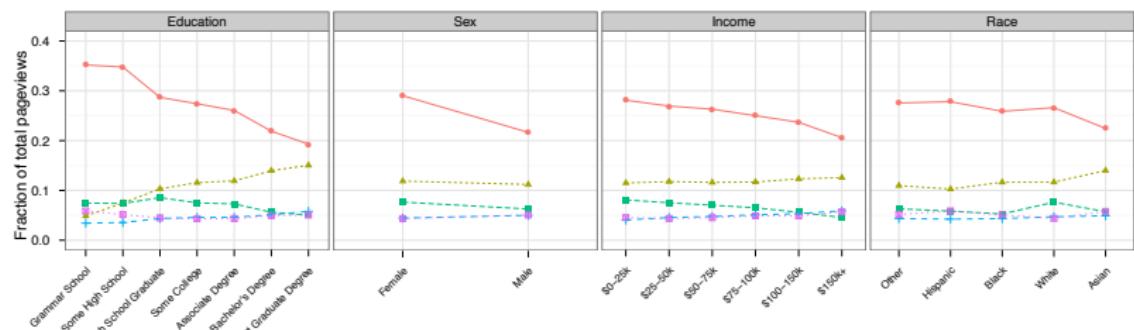
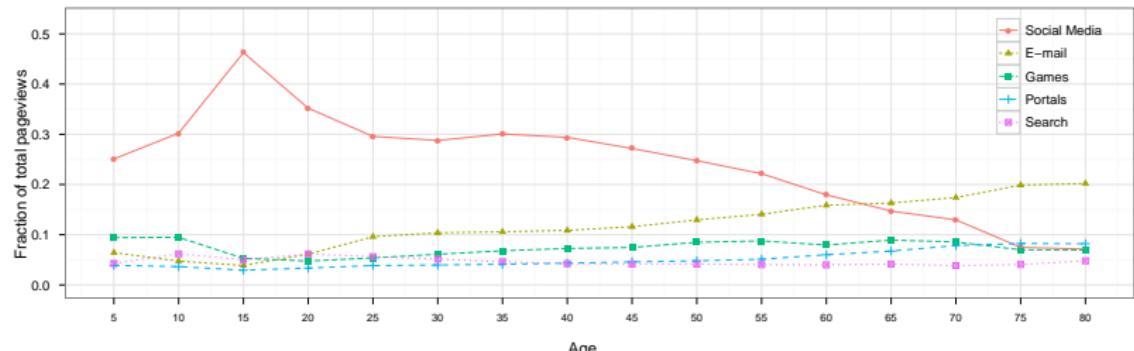
# Group-level activity

All demographic groups spend the majority of their time in the same categories



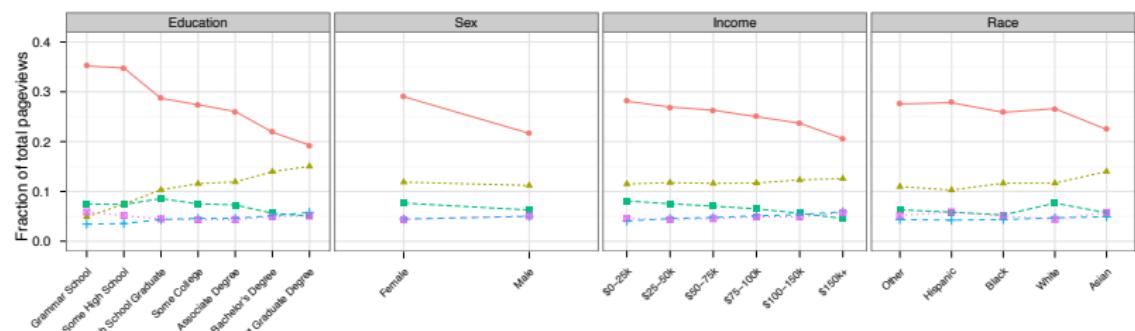
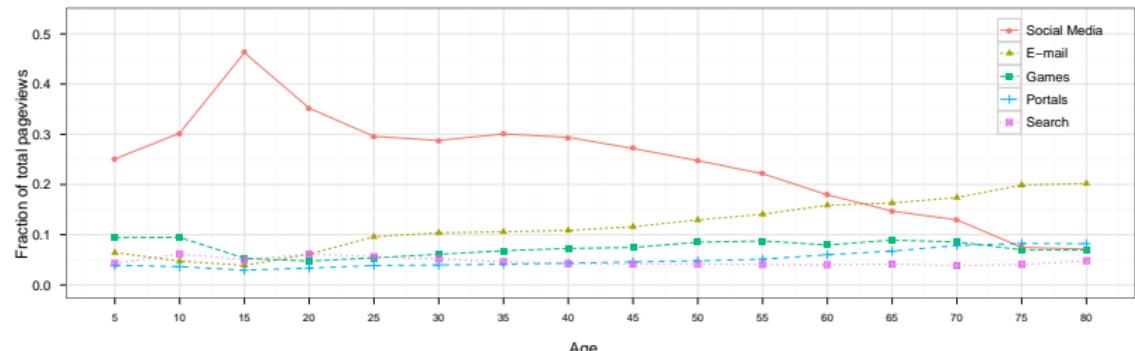
# Group-level activity

Older, more educated, male, wealthier, and Asian Internet users spend a smaller fraction of their time on social media



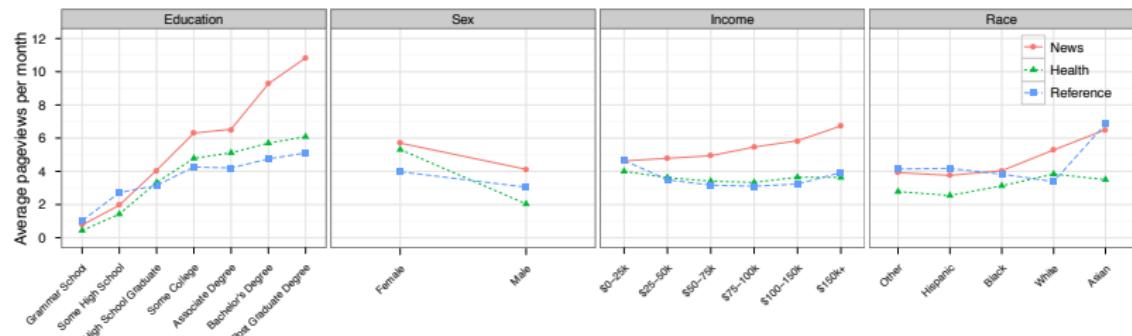
# Group-level activity

Lower social media use by these groups is often accompanied by higher e-mail volume



# Revisiting the digital divide

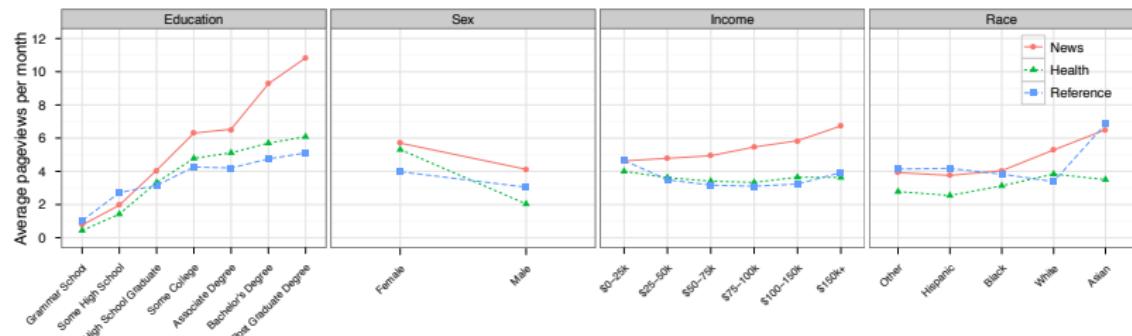
How does usage of news, health, and reference vary with demographics?



Post-graduates spend **three times** as much time on **health** sites than adults with only some high school education

# Revisiting the digital divide

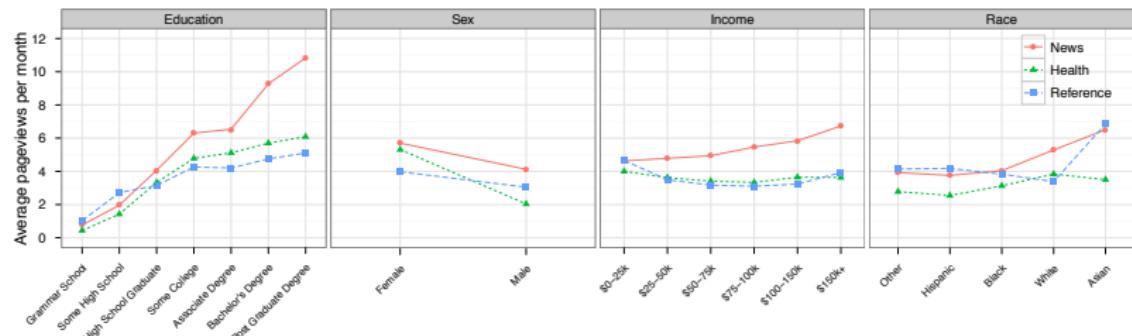
How does usage of news, health, and reference vary with demographics?



Asians spend more than 50% more time browsing online news than do other race groups

# Revisiting the digital divide

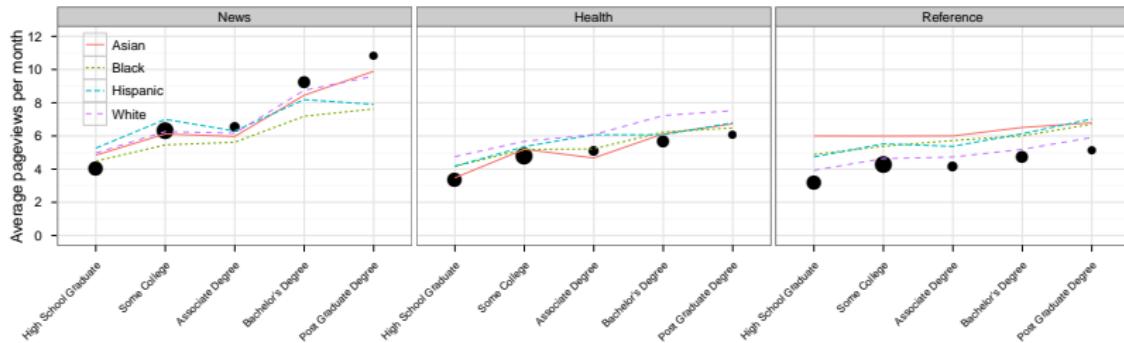
How does usage of news, health, and reference vary with demographics?



Even when less educated and less wealthy groups gain access to the Web, they utilize these resources relatively infrequently

# Revisiting the digital divide

How does usage of news, health, and reference vary with demographics?

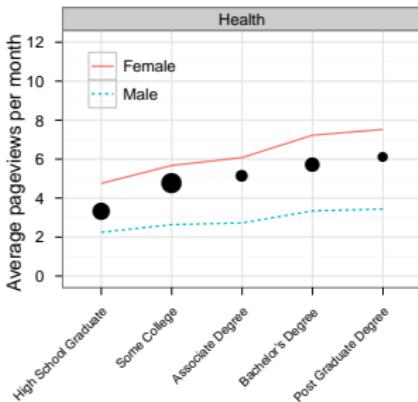


Controlling for other variables, effects of race and gender largely disappear, while education continues to have large effect

$$p_i = \sum_j \alpha_j x_{ij} + \sum_j \sum_k \beta_{jk} x_{ij} x_{ik} + \sum_j \gamma_j x_{ij}^2 + \epsilon_i$$

# Revisiting the digital divide

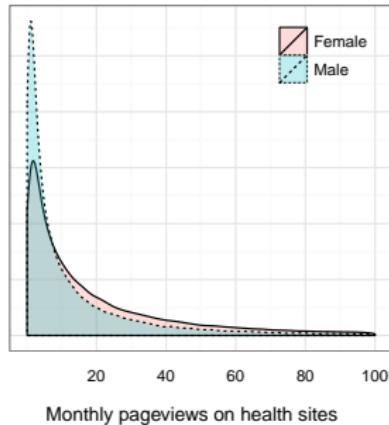
How does usage of news, health, and reference vary with demographics?



However, women spend considerably more time on health sites compared to men

# Revisiting the digital divide

How does usage of news, health, and reference vary with demographics?



However, women spend considerably more time on health sites compared to men, although means can be misleading

# Individual-level prediction

How well can one predict an individual's demographics from their browsing activity?

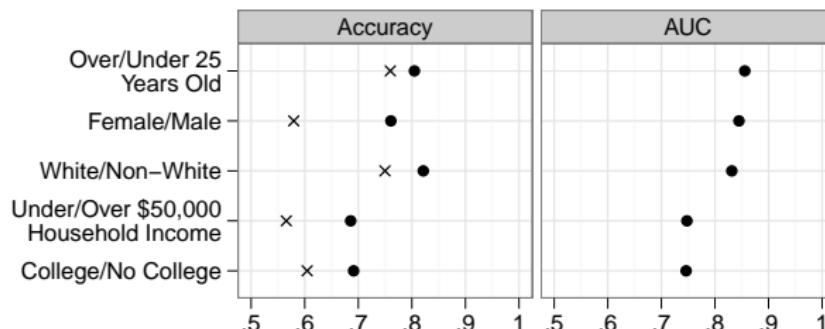
- Represent each user by the set of sites visited
- Fit linear models<sup>4</sup> to predict majority/minority for each attribute on 80% of users
- Tune model parameters using a 10% validation set
- Evaluate final performance on held-out 10% test set

---

<sup>4</sup><http://bit.ly/svmpref>

# Individual-level prediction

Reasonable (~70-85%) accuracy and AUC across all attributes



# Individual-level prediction

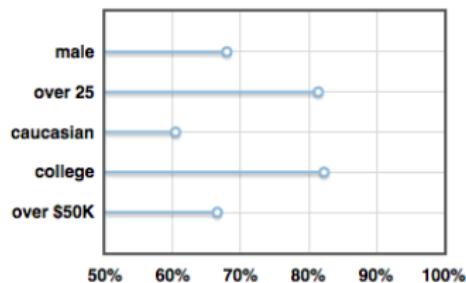
## Highly-weighted sites under the fitted models

	Large positive weight	Large negative weight
Female	winsters.com lancome-usa.com	sports.yahoo.com espn.go.com
White	marlboro.com cmt.com	mediatakeout.com bet.com
College Educated	news.yahoo.com linkedin.com	youtube.com myspace.com
Over 25 Years Old	evite.com classmates.com	addictinggames.com youtube.com
Household Income Under \$50,000	eharmony.com tracfone.com	rownine.com matrixdirect.com

# Individual-level prediction

## Proof of concept browser demo

From the 28 sites we found in your browser history, it appears that you're a **caucasian male** who is **over 25** years old with a **college** education earning **over \$50K** per year.



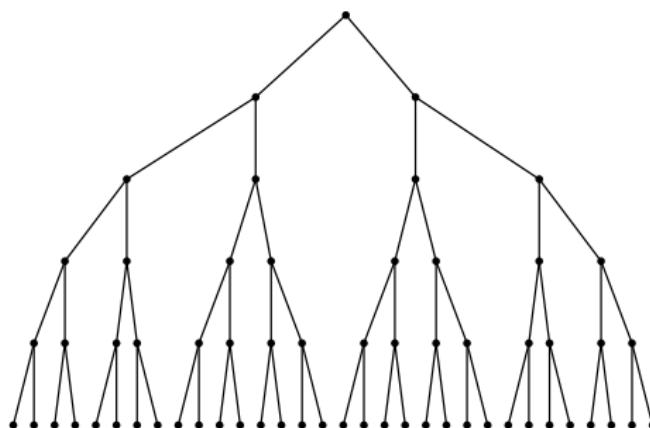
<http://bit.ly/surfpreds>  
(deprecated)

# Summary

- Highly active users spend disproportionately more of their time on social media and less on e-mail relative to the overall population
- Access to research, news, and healthcare is strongly related to education, not as closely to ethnicity
- User demographics can be inferred from browsing activity with reasonable accuracy
- “Who Does What on the Web”, Goel, Hofman & Sirer, ICWSM 2012

# The structural virality of online diffusion

with Ashton Anderson, Sharad Goel, Duncan Watts (Management Science 2015)



# “Going Viral”?

## viral

[Contents](#) [show]

### English

#### Etymology

From the stem of [virus](#) with suffix [-al](#).

#### Pronunciation

- IPA: /'vɪrəl/
- Rhymes: -aɪrəl

#### Adjective

**viral** (*not comparable*)

1. (*virology*) Of or relating to a biological [virus](#).  
*viral DNA*
2. (*virology*) Caused by a virus.  
*viral infection*
3. (*computing*) Of the nature of an informatic virus; able to spread copies of itself to other computers.
4. (*advertising and marketing*) Spread by word of mouth, with minimal intervention in order to create [buzz](#) and interest.

#### Derived terms

- [go viral](#)
- [viral marketing](#)

# "Going Viral"?

A MORE ET STUDIO ELVCID ANDAE  
ueritatis hac fabri ipsius diffundebant Vittenbergae. Prædicti te  
R. P. Martinus Luther, Artii & S. Theologie Magistri, eius  
deinde ibidem lectorum Ordinario. Quare petit ut qui non pos-  
sunt uerbis praefentes nobiscum disceptare, agant id literis ab-  
fentes. In nomine domini nostri Iesu Christi. Amen.

i) Ominus & Magister nosfer Iesu Christo, di-  
cendo penitentiam agite &c. omnem uitam fi-  
delium penitentiam esse uoluit.  
Quod uerbū penitentia de penitentia sacra-  
mentali, i.e. confessione & latitudinis quae  
sacerdotum ministerio celebratur) non po-  
telli intelligi.

ii) Non tamen sola inuidit interior, immo interior nulla est, nisi  
foris operetur tardis carnis mortificationes.

iii) Manciū ipso pena donec maneat odiūm fui, i.e. penitentia uera  
intus) collice usq; ad introitum regni celorum.

v) Papa non uult nec porf; illas penas remittere: præter eas,  
quas arbitrio uel suo uel canonum imposuit.

vj) Papa nō potest remittere ullum culpā, nisi declarato & appro-  
bando remissam a deo. Aut certe remittēdo casus referuntur  
sibi, quibus contemp̄tis culpa profusa remanet.

vij) Nulli profut remittit deus culpam, quia simul cum subiicit  
humilitatem in omnibus sacerdoti suo uicario.

vij) Canones penitentiales solum inuentibus sunt impositi: nihilq;  
mortuoris, secundū eodem debet imponi.

ix) Inde bene nobis facit spirituflamus in Papa: excipido in fu-  
re decreta tempore articulum mortis & necessitatis,

x) Indo te & male faciunt sacerdotes ij, qui mortuoris penitentias  
canonicas in purgatorium referunt.

xj) Zizania illa de mutanda pena Canonica in penā purgato-  
riū, uidenter certe dormientibus Episcopis feminata.  
Olla penae canonicæ nō post, fed ante absolutionem im-  
ponabantur, tanq; tentamenta uerae contritionis.

xi) Mortuū, per mortem omnes soluant, & legibus canonū mor-  
tuū fam funt, habentes tunc carū relaxationem.

xii) Imperfecta sanitas seu charitas morituri, necessario fecunt fert  
magnum timorem, tuncq; maiorem, quanto minor fuerit ipsa.

xv) Eficū timor & horror, sicut est, le folio (ut alia tacetam) facere po-  
nam purgatorij, cum sit proximus delerationis horrois.

xvi) Videntur, infernus, purgatoriorum, celum differre, sicut & despe-  
ratio, prope desperatio, securitas differunt.

xvii) Neccliarum uidetur animabus in purgatorio sicut minui hor-  
rorum, ita angeri charactarum.

xviii) Nec probatum uidetur ullis, aut rationibus, aut scripturis, q; sint  
extra statum mortis seu augendae charactarū.

xix) Nec hoc probatum esse uidetur, q; sint de sua beatitudine certæ  
& securitas, saltem oēs, sicut nos certissimi sumus.

xx) Igī Papa per remissiōnē plenariā omnū penarū, non simpli-  
citer omniū intelligit, sed a seipso et modo imponit.

xxi) Errantiaq; indulgentiarū predicatores ij, qui dicunt per Pa-  
pa indulgentias, hominē ab omni pena solui & saluari.

xxii) Quia nullam remittit animalibus in purgatorio, quā in hac ui-  
ta debuit senti secundum Canonem soluare.

xxiii) Si remissio illa omnī omnina penari pōe alicui dati; certū  
est eam nō nisi perfectissimis i.e. pateficiis dati.

xxiv) Palli ob id necesse est, maiores parē populi; per indifferentē  
illam & magnificam penae solute promissione.

xxv) Quale potestē habet Papa in purgatoriū gnatice talē haber-  
et libet Episcopus & cura in sua dioceſi, & parochia sp̄litter.

i) Optime facit Papa, q; nō potestate clausi (quā nullam habet)  
sed per modum suffragij, dat animalibus remissionem.

j) Hominē predican, qui huius, ut factus nūmus in cīdam tē-  
nient, euolare dicunt animām.

tj) Centū et nūmo in cīdam tēniente, augeri questum & auaci-  
onē polle, suffragij aut ecclēsiae est in arbitrio dei solis.

ii) Quis scit q; omnes animas in purgatorio uelint redire, sicut de  
sancto Seuerino & pachali factum narratur?

v) Nullus securus est de ueritate sue contritionis; multo minus

# “Going Viral”?

*“Therefore we ... wish to proceed with great care as is proper, and to cut off the advance of this **plague** and **cancerous disease** so it will not **spread** any further ...”<sup>5</sup>*

-Pope Leo X  
*Exsurge Domine* (1520)

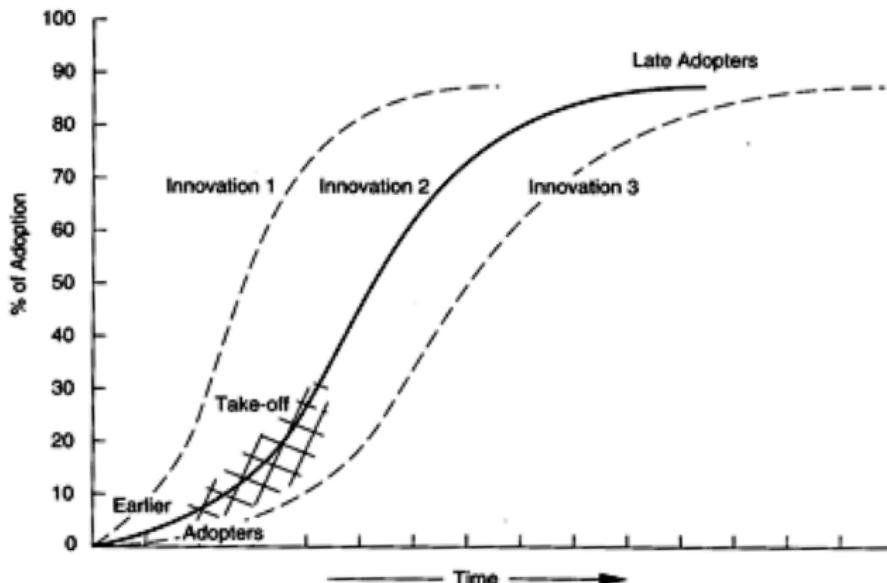
---

<sup>5</sup><http://www.economist.com/node/21541719>

# “Going Viral”?

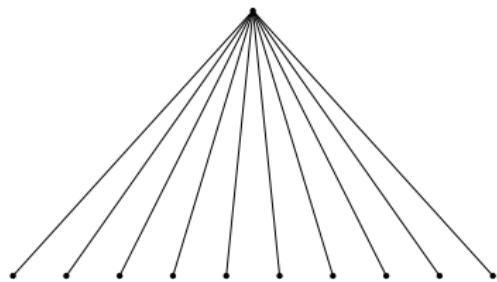
**FIGURE 6.5** Shapes of curves of diffusion for innovations that spread over various periods of time

source: Everett M. Rogers, *Diffusion of Innovations*, 3rd ed. (New York: Free Press, 1963), p. 11.



Rogers (1962), Bass (1969)

# “Going viral”?



CNNMoney.com @CNNMoney

3 May

Dow crosses 15,000 for the first time, fueled by strong jobs report.

cnnmoney.com/bkgnews

 Retweeted by CNN Breaking News

[Collapse](#) [Reply](#) [Retweet](#) [Favorite](#) [More](#)

273

49

RETWEETS FAVORITES



10:23 AM - 3 May 13 · Details

# “Going viral”?

**CWB Brasil queremos novamente o show da Banda Restart em Curitiba - Paraná**

Created 12 months ago by @PeLuMoraComigo

## Description

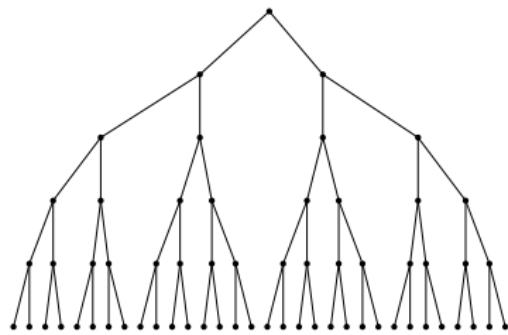
Pedimos atenciosamente a CWB Brasil novamente o show da Banda Restart em Curitiba. Desde o dia 29 de

2,352  
signatures so far ...

**Sign ►**

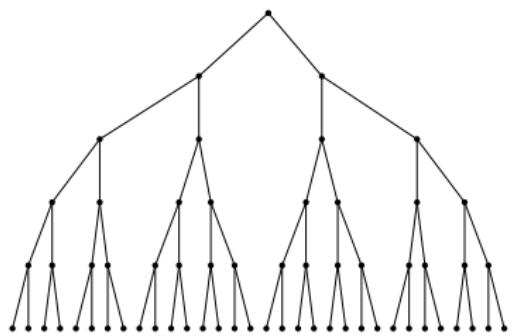
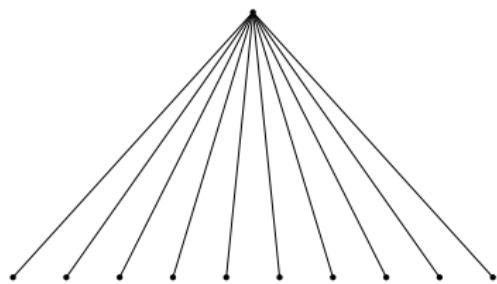
Tweet my signature  
#Twittion CWB Brasil queremos novamente o show da Banda Restart em Curitiba - Paraná http://twittion.com/lopypv

Follow @Twittion



# “Going viral”?

How do popular things become popular?



# Data

- Examined one year of tweets from July 2011 to July 2012

# Data

- Examined one year of tweets from July 2011 to July 2012
- Restricted to 1.4 billion tweets containing links to top news, videos, images, and petitions sites

# Data

- Examined one year of tweets from July 2011 to July 2012
- Restricted to 1.4 billion tweets containing links to top news, videos, images, and petitions sites
- Aggregated tweets by URL, resulting in 1 billion distinct “events”

# Data

- Examined one year of tweets from July 2011 to July 2012
- Restricted to 1.4 billion tweets containing links to top news, videos, images, and petitions sites
- Aggregated tweets by URL, resulting in 1 billion distinct “events”
- Crawled friend list of each adopter

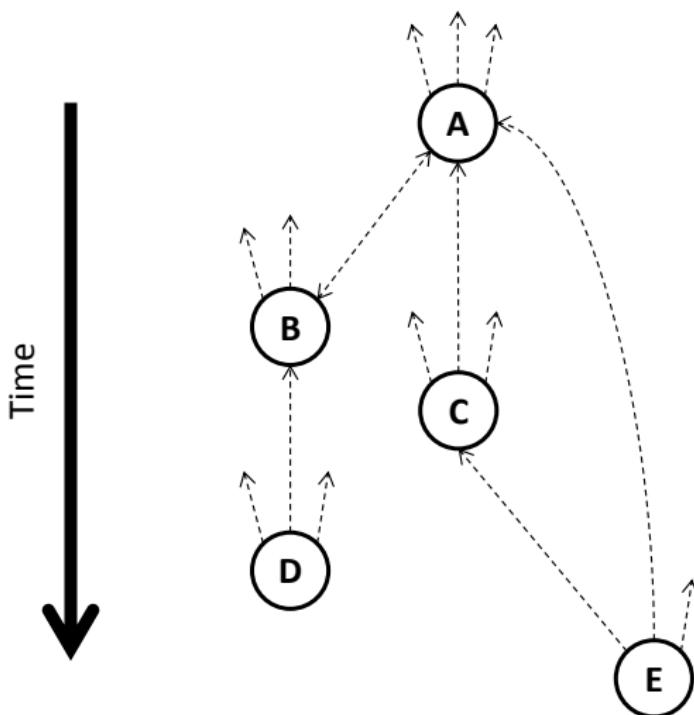
# Data

- Examined one year of tweets from July 2011 to July 2012
- Restricted to 1.4 billion tweets containing links to top news, videos, images, and petitions sites
- Aggregated tweets by URL, resulting in 1 billion distinct “events”
- Crawled friend list of each adopter
- Inferred “who got what from whom” to construct diffusion trees

# Data

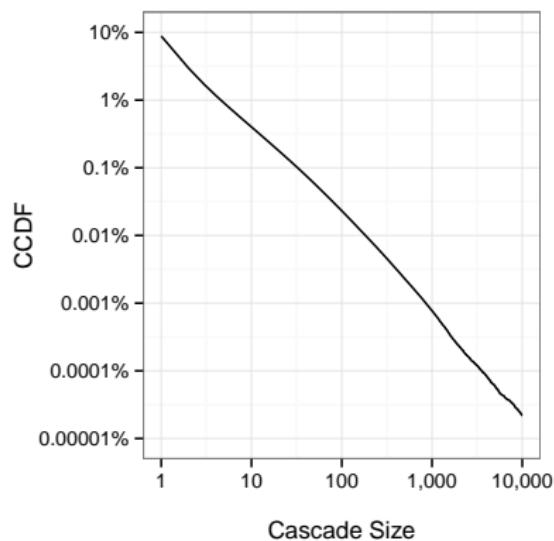
- Examined one year of tweets from July 2011 to July 2012
- Restricted to 1.4 billion tweets containing links to top news, videos, images, and petitions sites
- Aggregated tweets by URL, resulting in 1 billion distinct “events”
- Crawled friend list of each adopter
- Inferred “who got what from whom” to construct diffusion trees
- Characterized size and structure of trees

# The Structural Virality of Online Diffusion



# Information diffusion

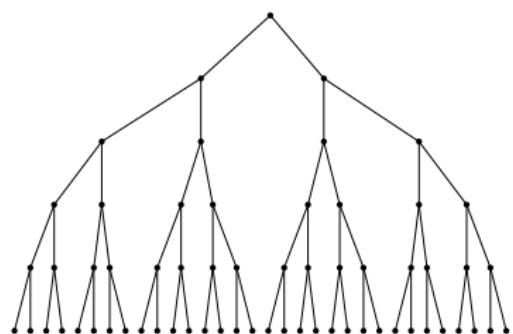
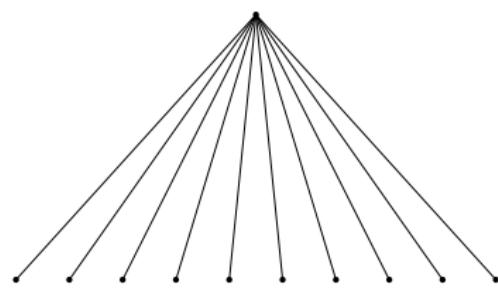
## Cascade size distribution



Focus on the rare hits that get at least 100 adoptions

# Quantifying structure

Measure the **average distance** between **all pairs of nodes**<sup>6</sup>



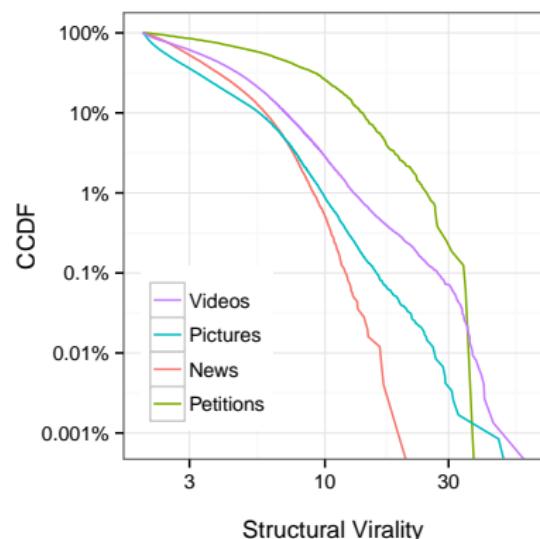
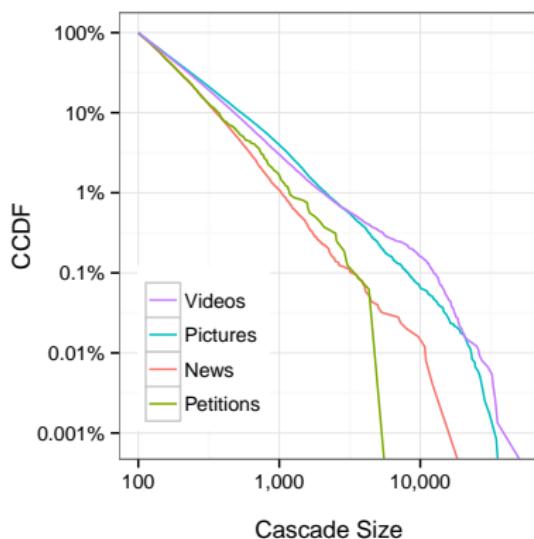
---

<sup>6</sup>Weiner (1947); correlated with other possible metrics

# Information diffusion

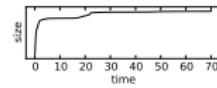
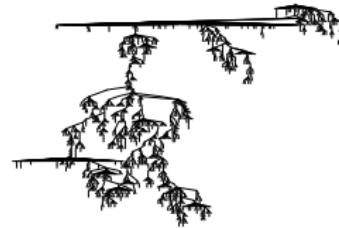
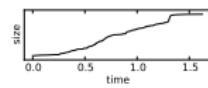
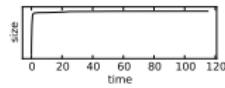
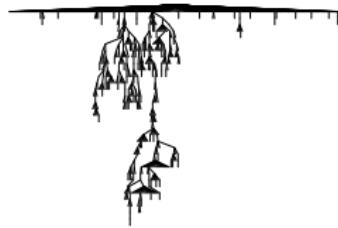
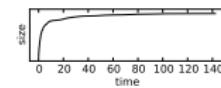
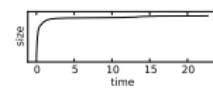
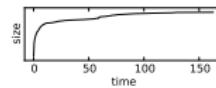
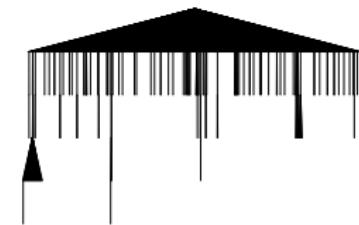
Size and virality by category

Remarkable structural diversity across categories



# Information diffusion

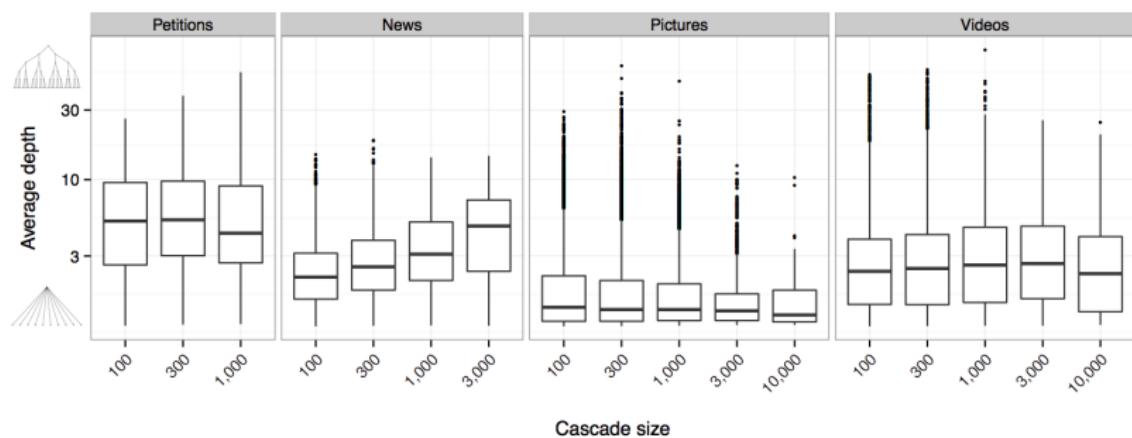
## Structural diversity



# Information diffusion

## Structural diversity

Size is relatively poor predictive of structure



Popular ≠ Viral

# Information diffusion

## Summary

- Most cascades fail, resulting in fewer than two adoptions, on average
- Of the hits that do succeed, we observe a wide range of diverse diffusion structures
- It's difficult to say how something spread given only its popularity
- "The structural virality of online diffusion", Anderson, Goel, Hofman & Watts (Management Science 2015)

# 1. Ask good questions

There's nothing interesting in the data without them

## 2. Think before you code

5 minutes at the whiteboard is worth an hour at the keyboard

### 3. Keep the answers simple

Exploratory data analysis and linear models go a long way

## 4. Replication is key

Otherwise it's easy to get fooled by randomness and difficult to assess progress