

data science @ The New York Times

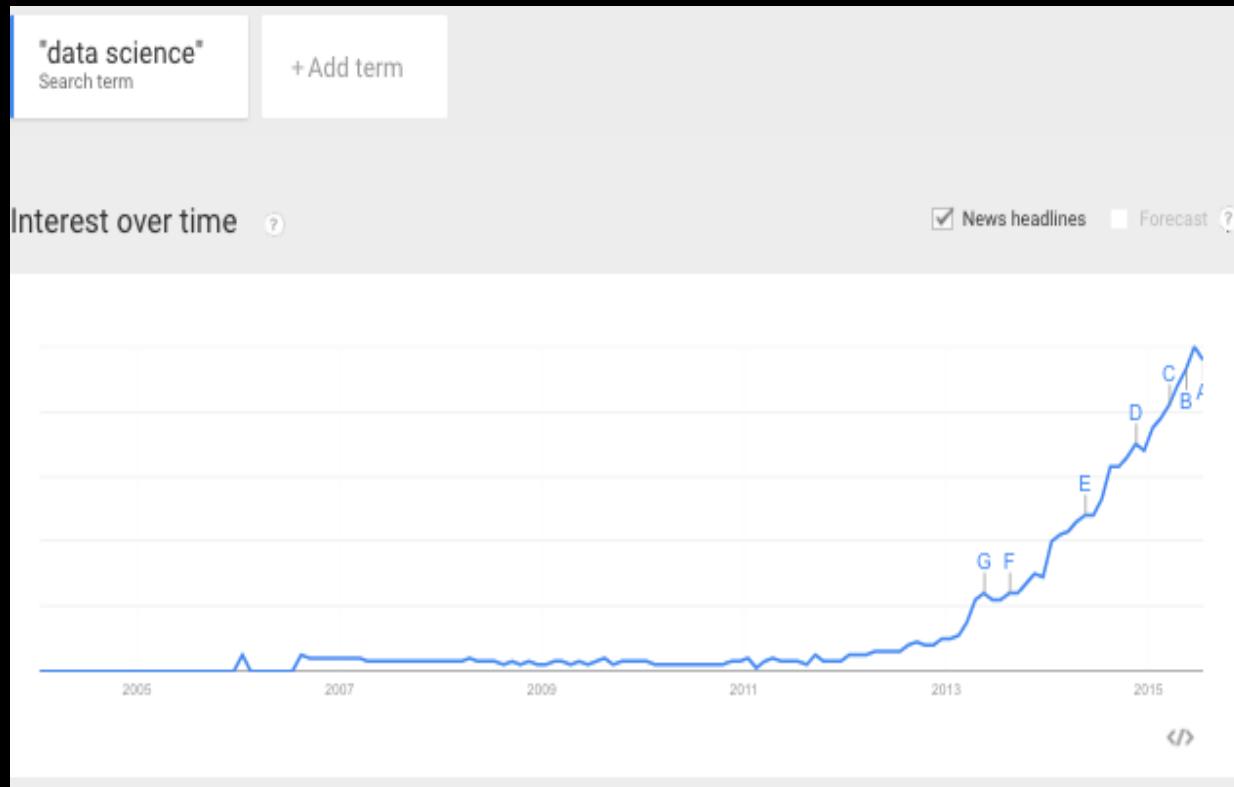


chris.wiggins@columbia.edu  
chris.wiggins@nytimes.com  
@chrishwiggins

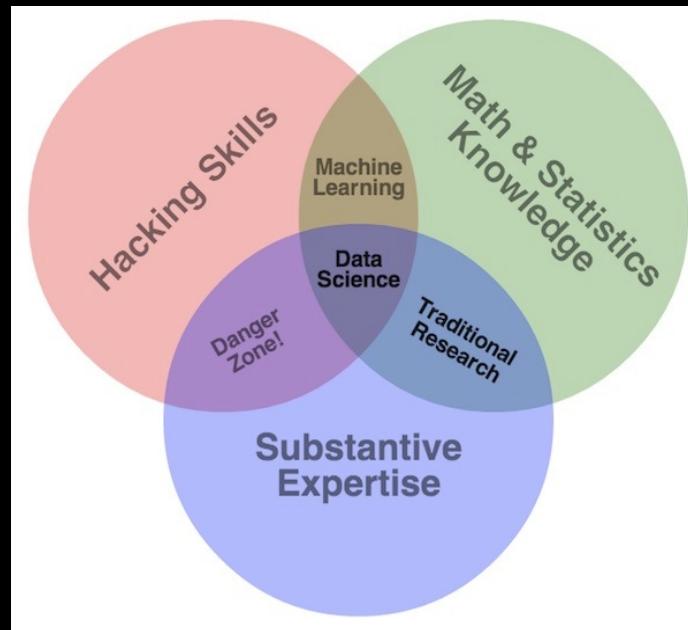
references: <http://bit.ly/stanf16>

data science

# data science: searches

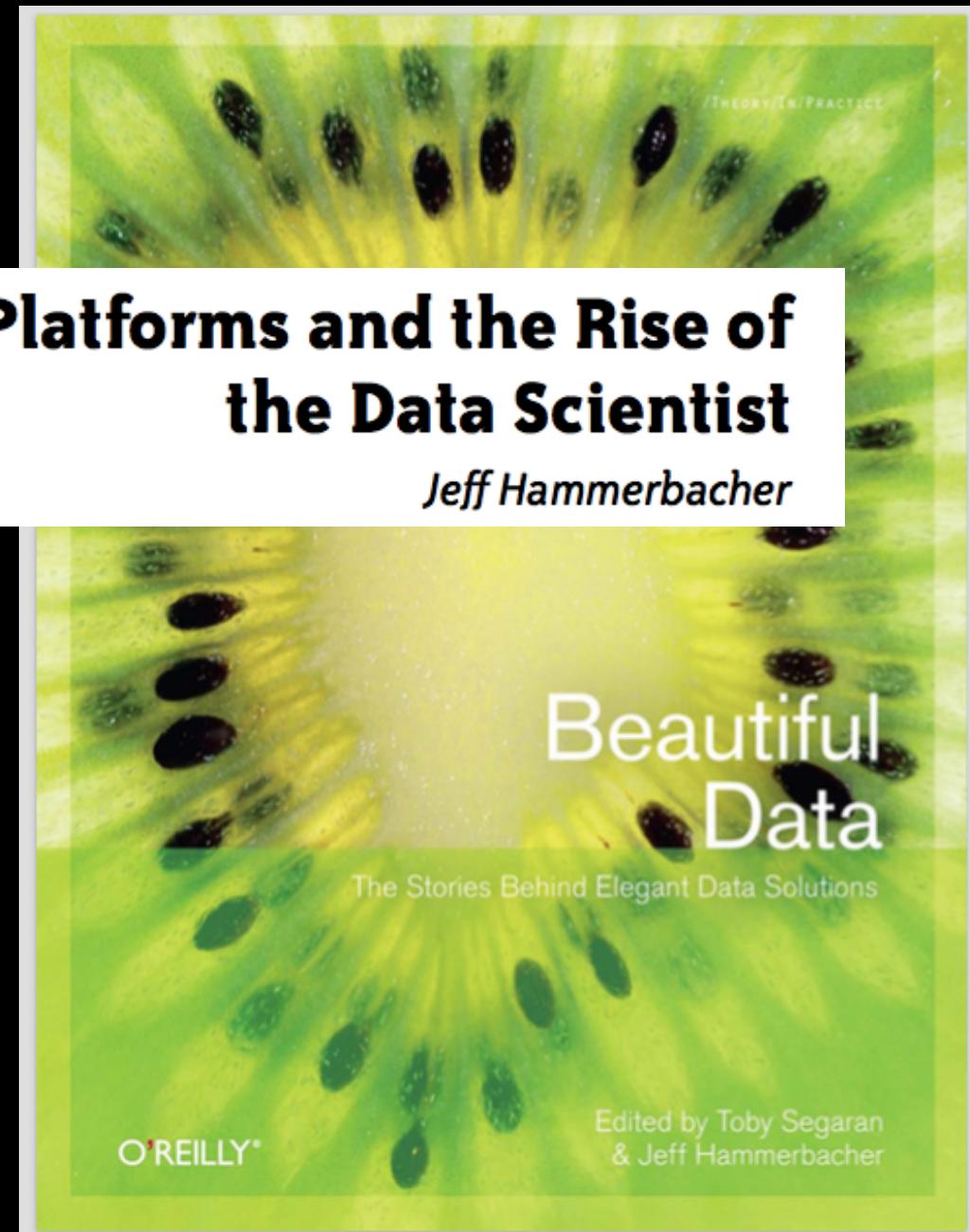


## data science: mindset & toolset



drew conway, 2010

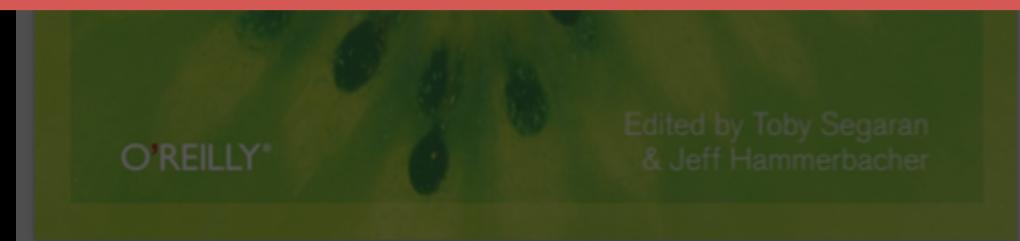
modern history:  
2009



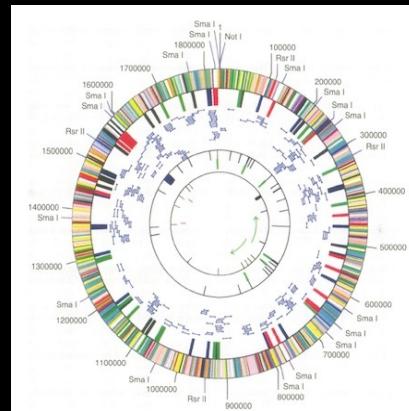
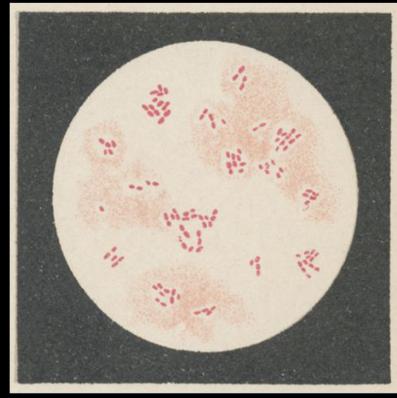
## Information Platforms and the Rise of the Data Scientist

At Facebook, we felt that traditional titles such as Business Analyst, Statistician, Engineer, and Research Scientist didn't quite capture what we were after for our team. The workload for the role was diverse: on any given day, a team member could author a multistage processing pipeline in Python, design a hypothesis test, perform a regression analysis over data samples with R, design and implement an algorithm for some data-intensive product or service in Hadoop, or communicate the results of our analyses to other members of the organization in a clear and concise fashion. To capture the skill set required to perform this multitude of tasks, we created the role of "Data Scientist."

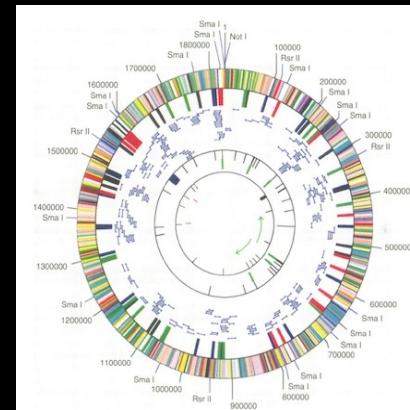
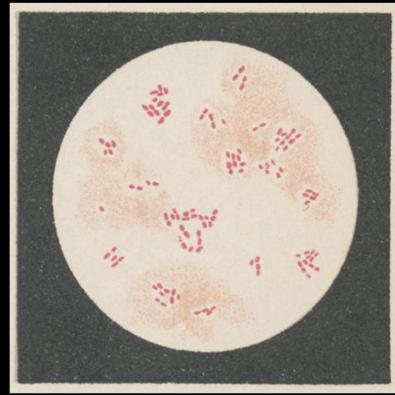
2009



## biology: 1892 vs. 1995

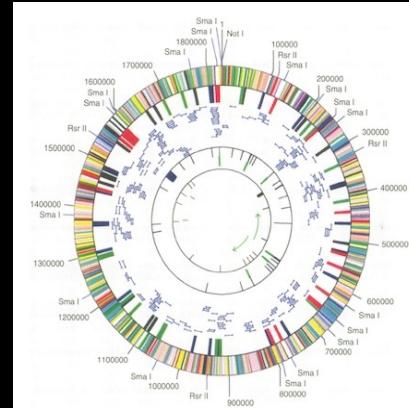


biology: 1892 vs. 1995



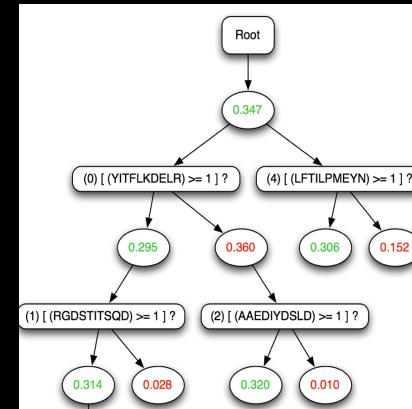
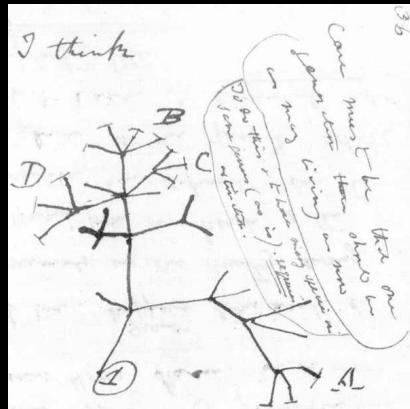
biology changed for good.

biology: 1892 vs. 1995



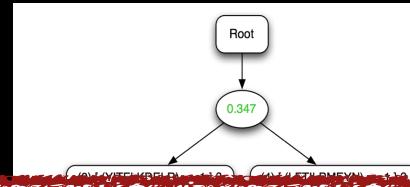
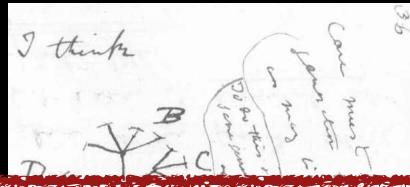
new toolset, new mindset

## genetics: 1837 vs. 2012



ML toolset; data science mindset

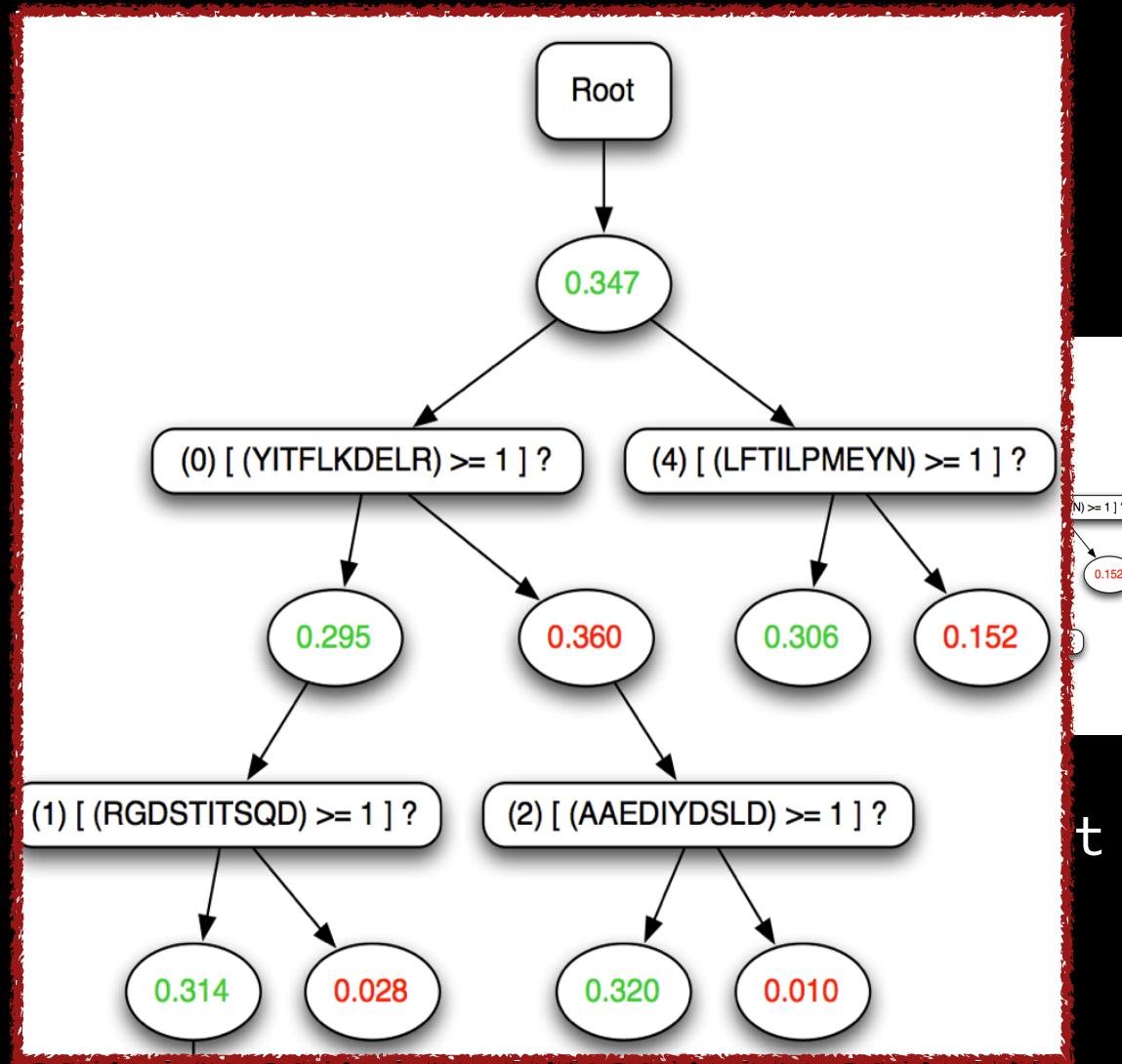
genetics: 1837 vs. 2012



*Statistical Science*  
2001, Vol. 16, No. 3, 199–231

# Statistical Modeling: The Two Cultures

Leo Breiman



*data science: mindset & toolset*



1851

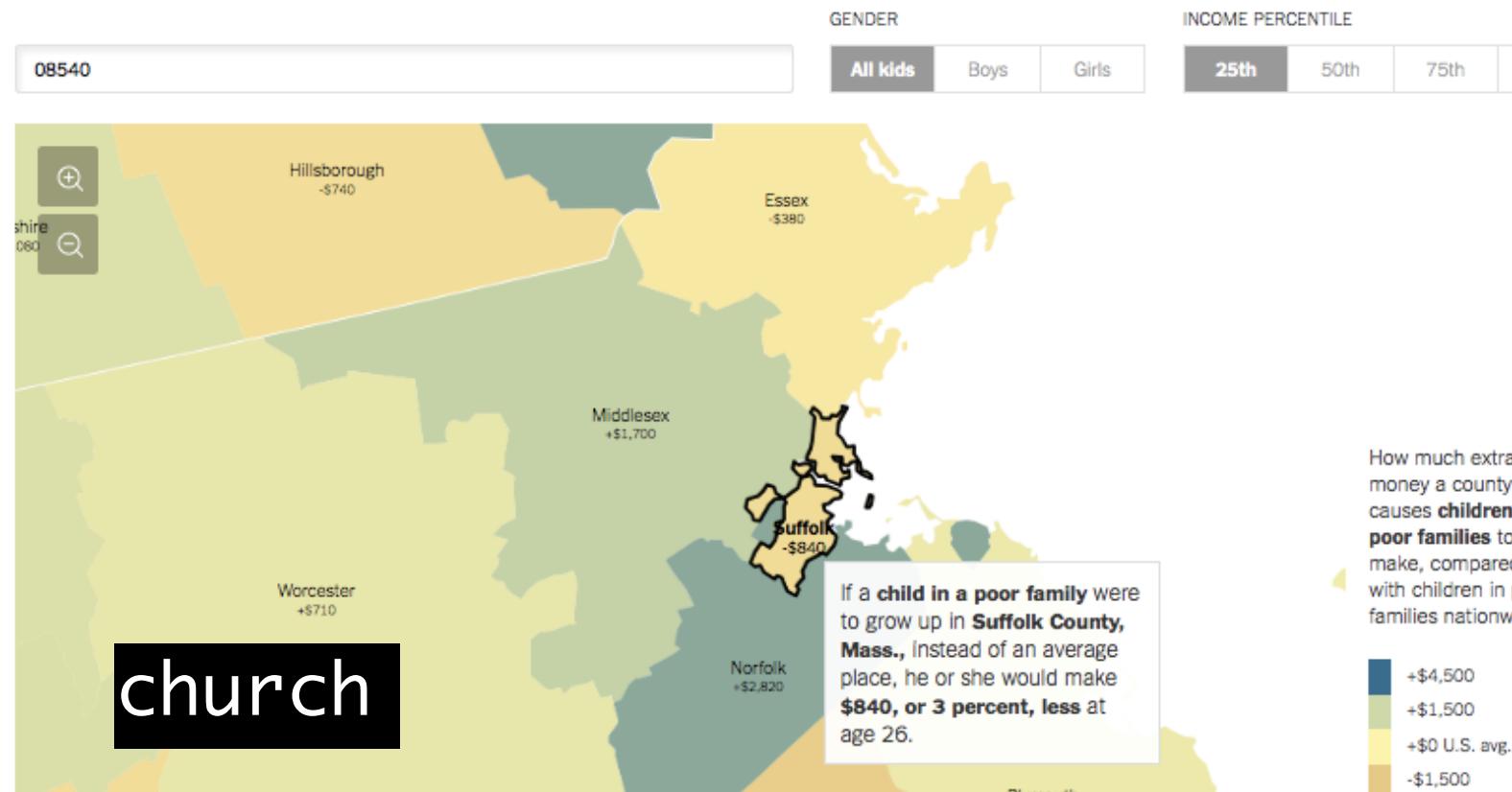
news: 20th century

church

state

## The Best and Worst Places to Grow Up: How Your Area Compares

Children who grow up in some places go on to earn much more than they would if they grew up elsewhere. MAY 4, 2015 | RELATED ARTICLE



TheUpshot/leo-senate-model

<https://github.com/TheUpshot/leo-senate-model>

**GitHub** This repository Search or type a command

PUBLIC TheUpshot / leo-senate-model

church

Code and data for The Upshot's Senate model. <http://www.nytimes.com/newsgraphics/2014/senate-model/>

12 commits 1 branch 0 releases 3 contributors

branch: master / leo-senate-model / +

changing default parameters

joshkatz authored 2 hours ago latest commit 30e1af96c9

File	Description	Time
data-publisher	Include directories required for the script to generate output	11 hours ago
fundamentals	Rename file (.r -> .R) for case-sensitive filesystems (e.g. Linux extN).	11 hours ago
model	Remove dependence on the authors' directory structure	11 hours ago
output	Include directories required for the script to generate output	11 hours ago
.gitignore	Leo lives	15 hours ago
LICENSE	Like grownups	4 hours ago
README.md	added sample data output to README.md	8 hours ago
master-public.R	changing default parameters	2 hours ago

[www.nytimes.com/interactive/2014/12/29/us/year-in-interactive-storytelling.html](http://www.nytimes.com/interactive/2014/12/29/us/year-in-interactive-storytelling.html)

U.S.

# 2014 The Year in Interactive Storytelling, Graphics and Multimedia

Multimedia Stories | Data Visualization

[www.nytimes.com/newsgraphics/2013/12/30/year-in-interactive-storytelling/](http://www.nytimes.com/newsgraphics/2013/12/30/year-in-interactive-storytelling/)

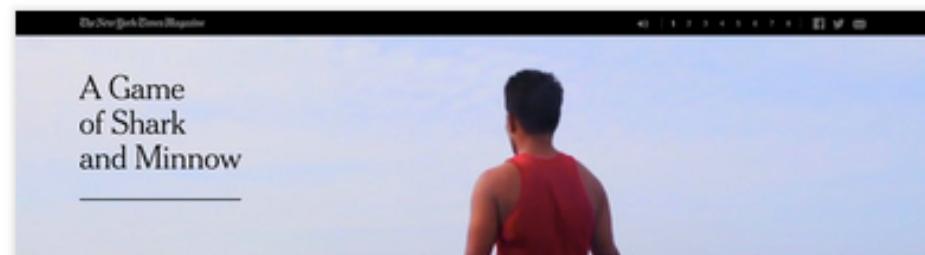
The New York Times

# 2013: The Year in Interactive Storytelling

Multimedia Stories | Data Visualization | Explanatory Graphics | Breaking News | Visual and Interactive Features

## Multimedia Stories

From a ship in the South China Sea to the cost of health care in the United States, the range of subjects here is broad, but the common thread is the form of storytelling — an integration of text, video, photography and graphics.



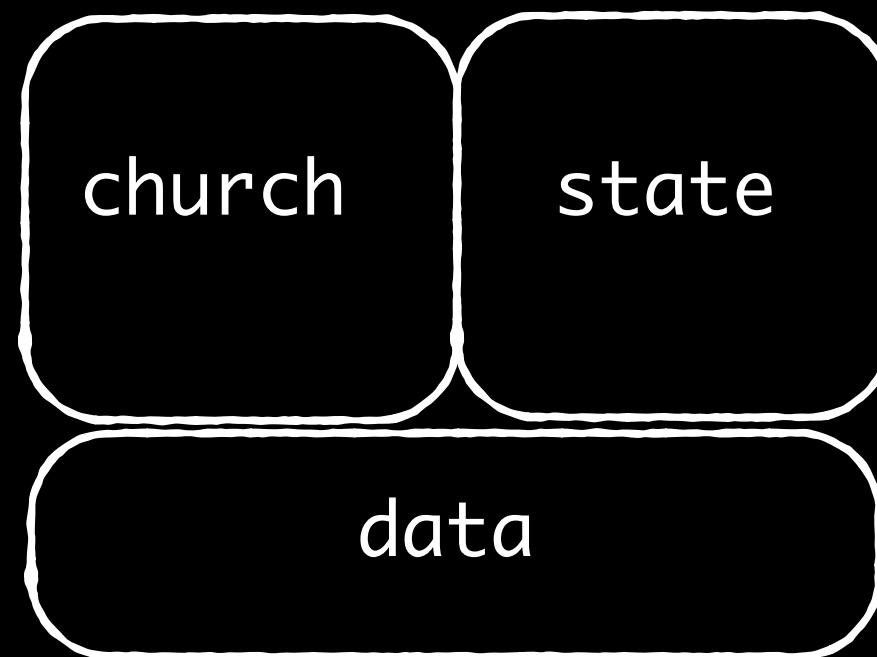
church

news: 20th century

church

state

news: 21st century



# newspapering: 1851 vs. 1996



1851

## The New York Times Introduces a Web Site

By PETER H. LEWIS

Published: January 22, 1996

The New York Times begins publishing daily on the World Wide Web today, offering readers around the world immediate access to most of the daily newspaper's contents.

The New York Times on the Web, as the electronic publication is known, contains most of the news and feature articles from the current day's printed newspaper, classified advertising, reporting that does not appear in the newspaper, and interactive features including the newspaper's crossword puzzle.

1996

**1,615,934** site-wide views over the last hour

**1,257,958** average Sunday New York Times print circulation

**554** stories written over the last 24 hours

**206** countries with visitors in the past 25 minutes

**243,192** words written in the last 24 hours

**65** New York Times newspaper print sites globally

**733** page views from India in the last 10 minutes



The lobby of The New York Times Building/Nic Lehoux

2015

**01/09/2014**

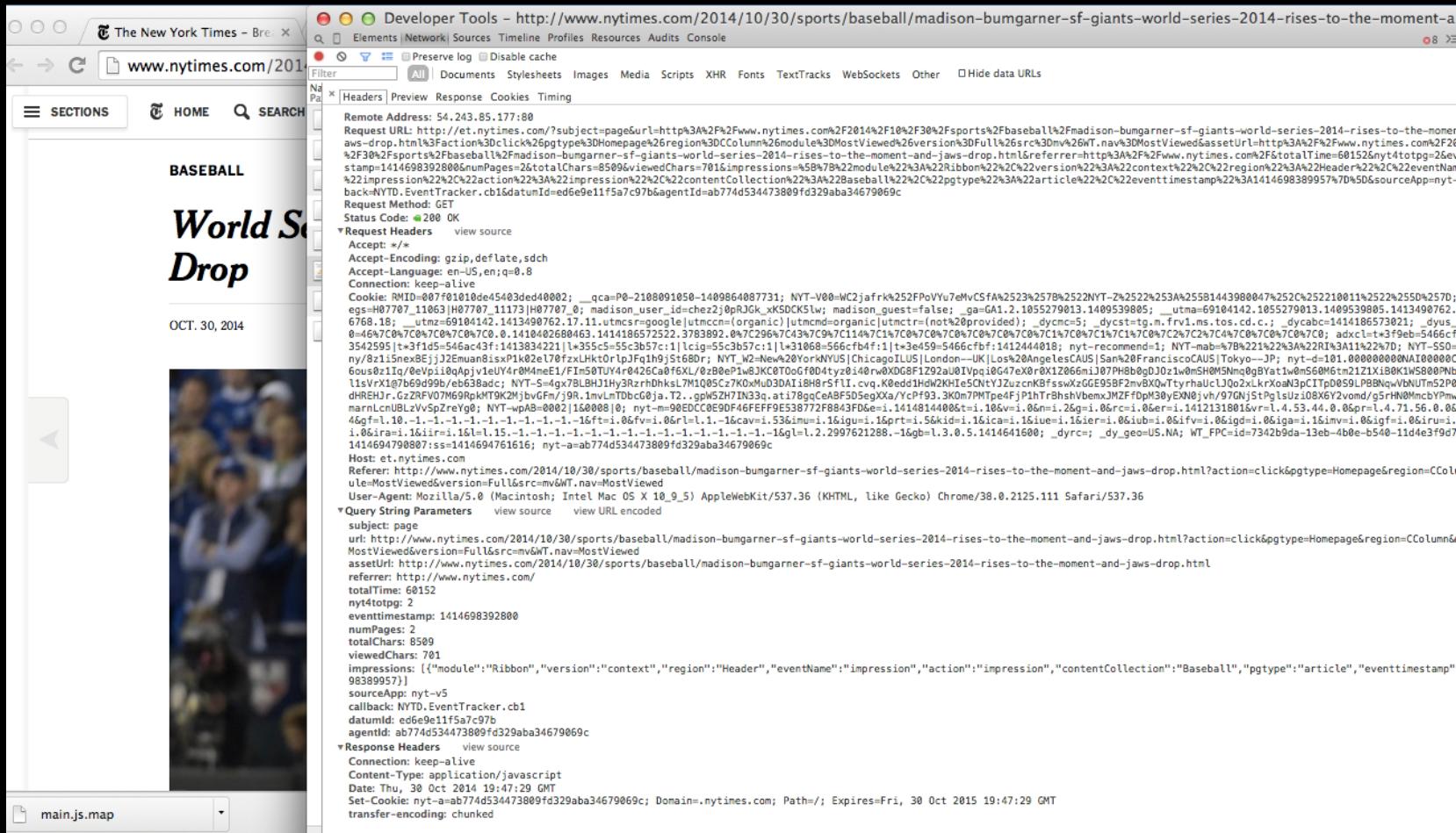
The New York Times Company to Webcast Fourth-Quarter and Full-Year 2013 Earnings Conference Call »

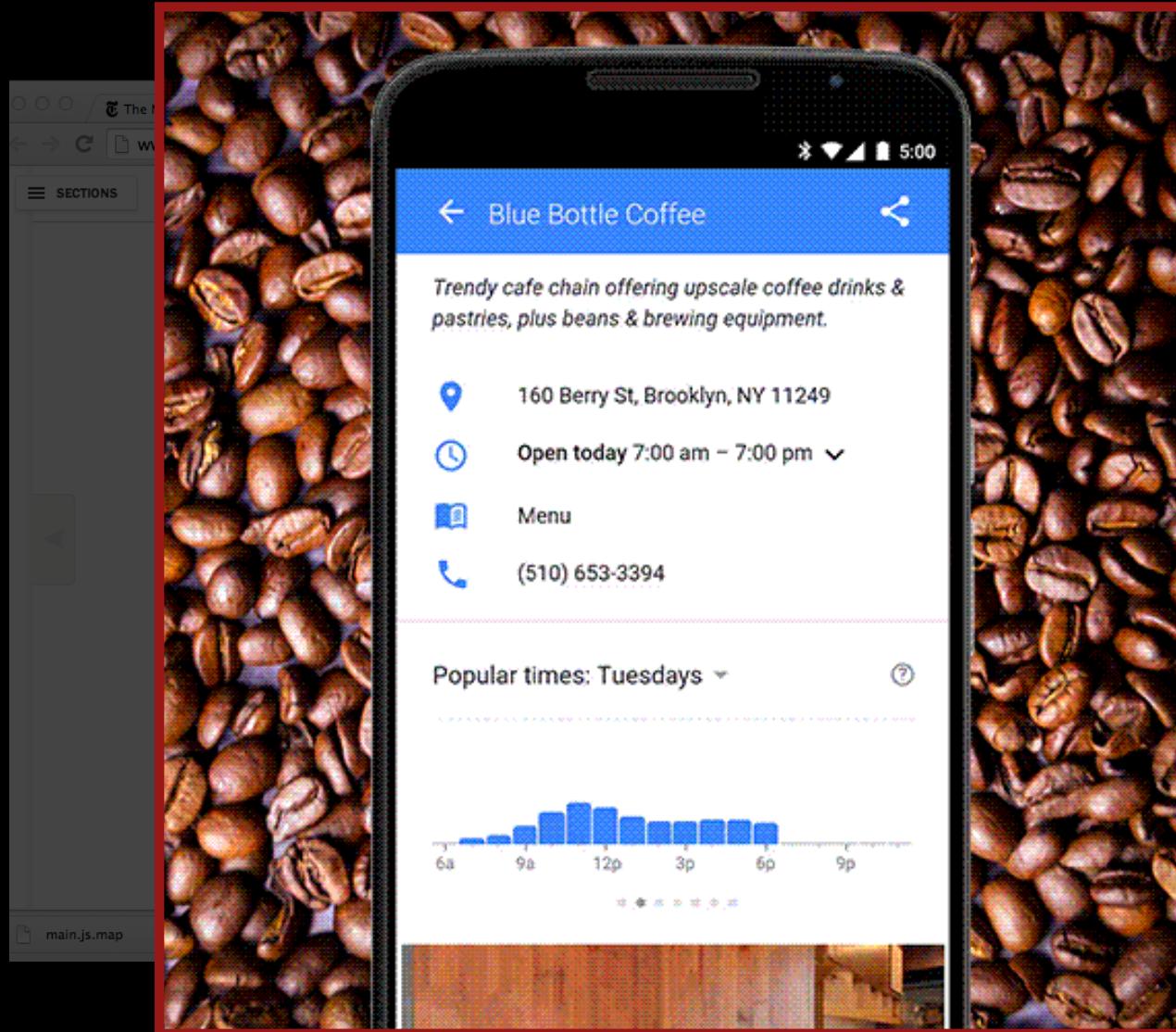
**01/02/2014**

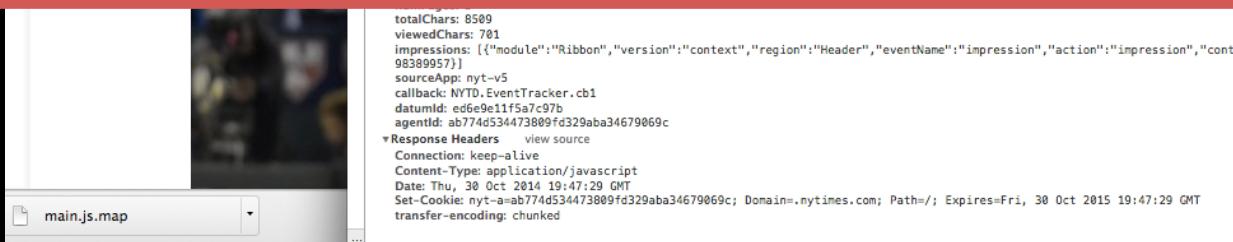
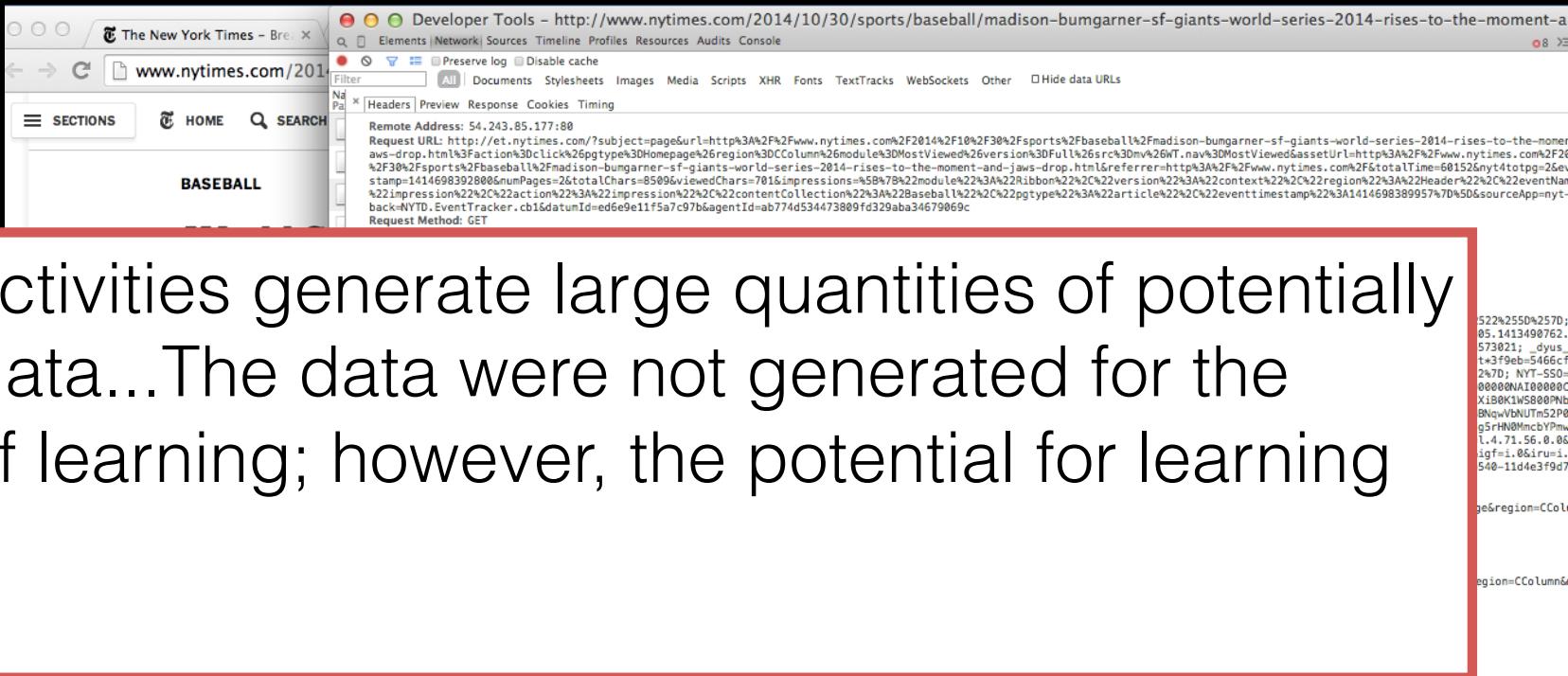
The New York Times to Introduce Redesign of NYTimes.com to All Users Jan. 8 »

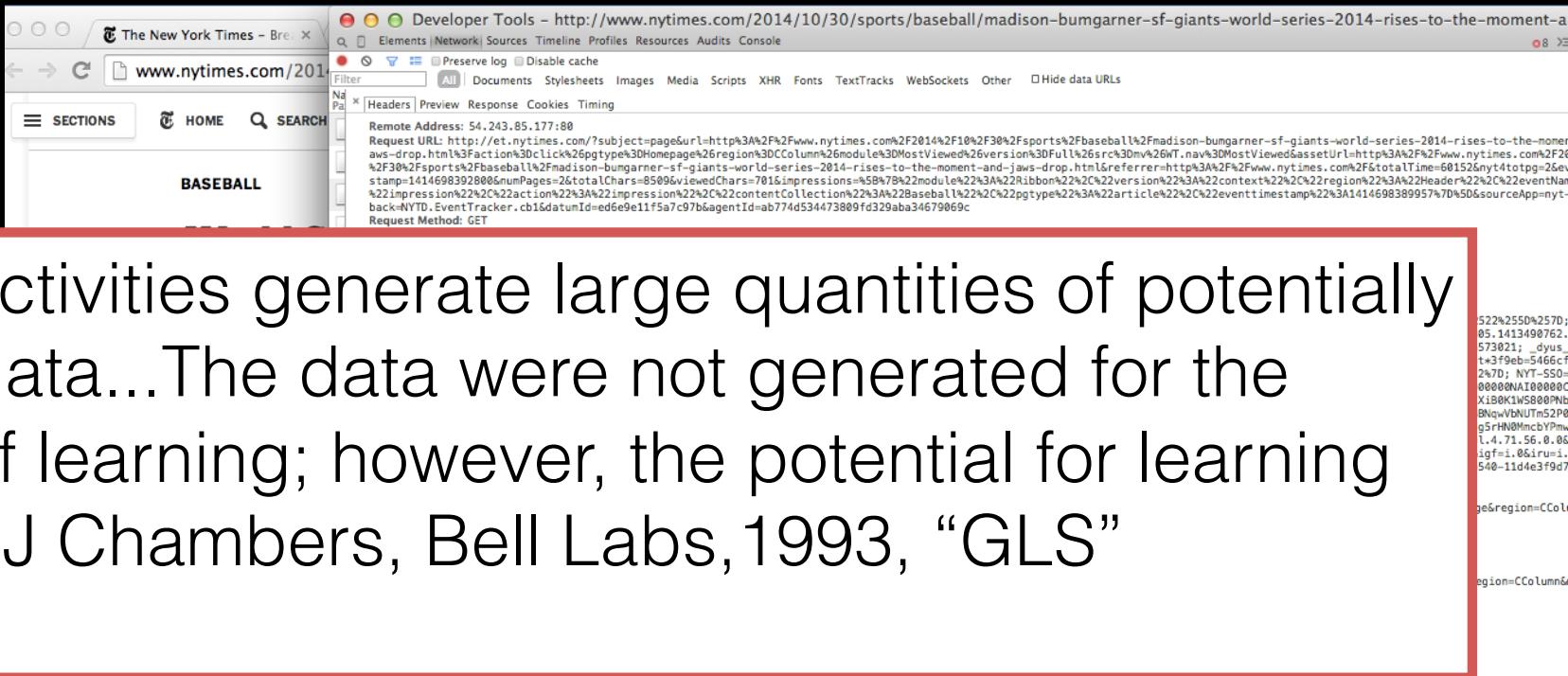
**12/12/2013**

The New York Times Company Declares Regular Quarterly Dividend »









data science: the web

data science: the web

is your “online presence”

data science: the web

is a microscope

data science: the web  
is an experimental tool

data science: the web

is an optimization tool

# newspapering: 1851 vs. 1996 vs. 2008



1851

## The New York Times Introduces a Web Site

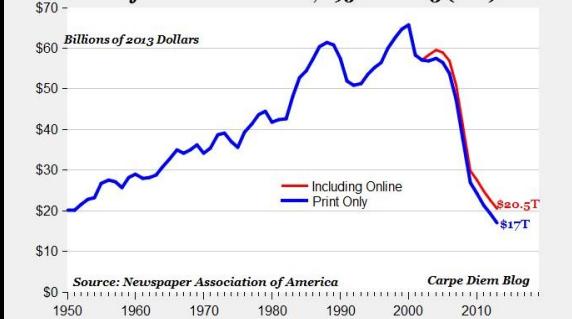
By PETER H. LEWIS  
Published: January 22, 1996

The New York Times begins publishing daily on the World Wide Web today, offering readers around the world immediate access to most of the daily newspaper's contents.

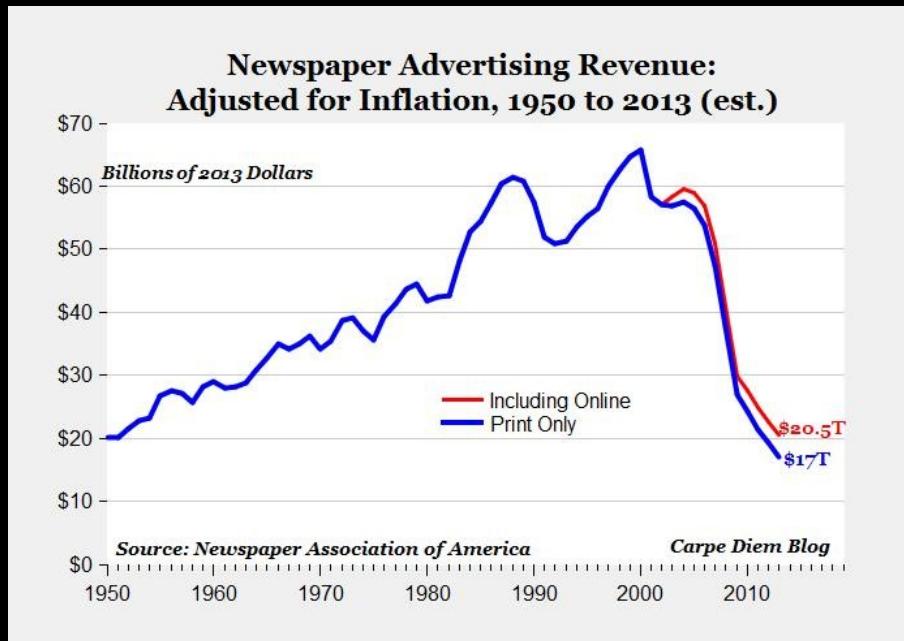
The New York Times on the Web, as the electronic publication is known, contains most of the news and feature articles from the current day's printed newspaper, classified advertising, reporting that does not appear in the newspaper, and interactive features including the newspaper's crossword puzzle.

1996

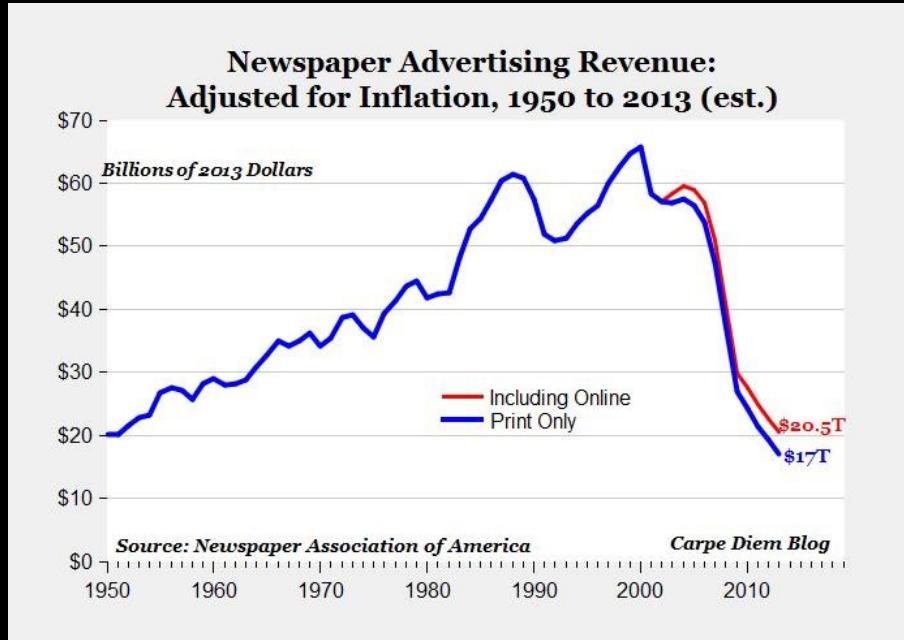
Newspaper Advertising Revenue:  
Adjusted for Inflation, 1950 to 2013 (est.)



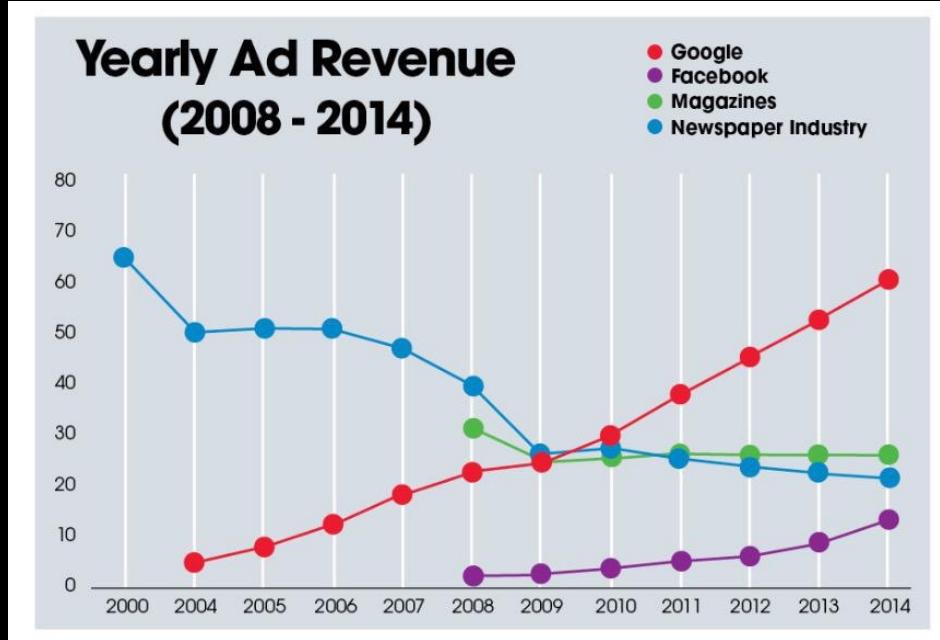
2008



“a startup is a temporary organization in search of a repeatable and scalable business model” –Steve Blank



every publisher is now a startup



every publisher is now a startup

The screenshot shows a web browser window with the following details:

- Title Bar:** Media Websites Battle Falter
- Address Bar:** www.nytimes.com/2016/04/18/business/media-websites-battle-falteringad-
- Header:** BUSINESS DAY | Media Websites Battle Faltering Ad Revenue and Traffic
- Text Content:** Advertisers adjusted spending accordingly. In the first quarter of 2016, 85 cents of every new dollar spent in online advertising will go to Google or Facebook, said Brian Nowak, a Morgan Stanley analyst.

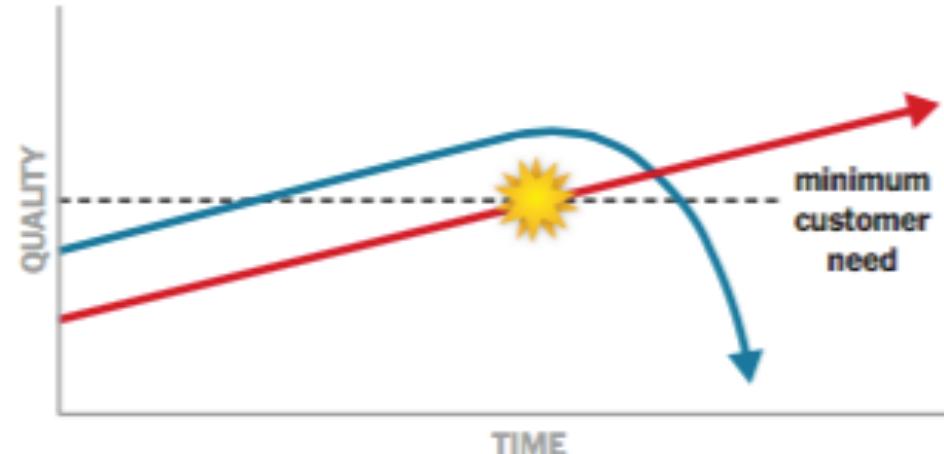
every publisher is now a startup

The New York Times

Innovation

**3.** Over time, **disruptors** improve their product, usually by adapting a new technology. The **flash-point** comes when their products become “good enough” for most customers.

They are now poised to grow by taking market share from **incumbents**.

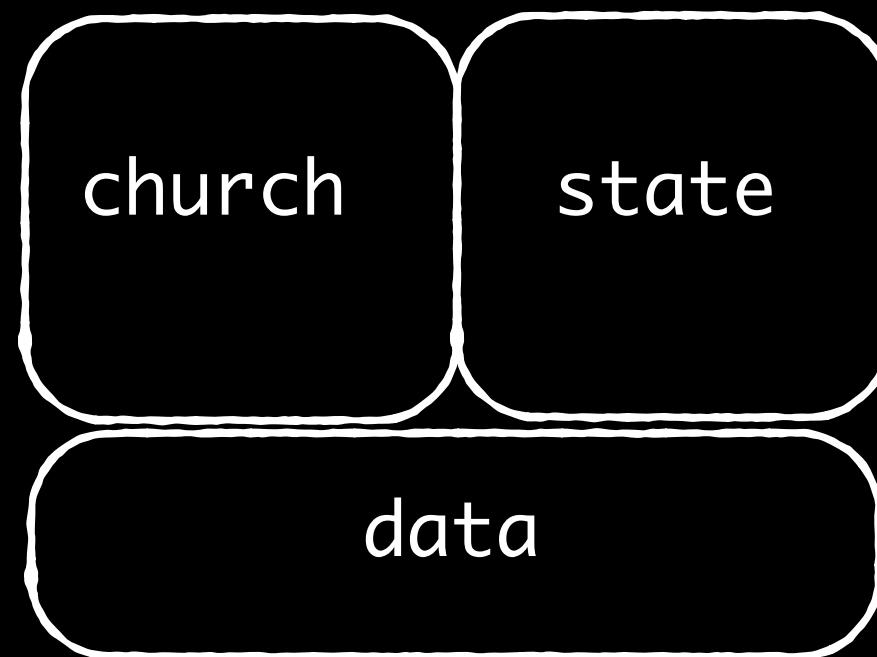


#### HALLMARKS OF DISRUPTIVE INNOVATORS

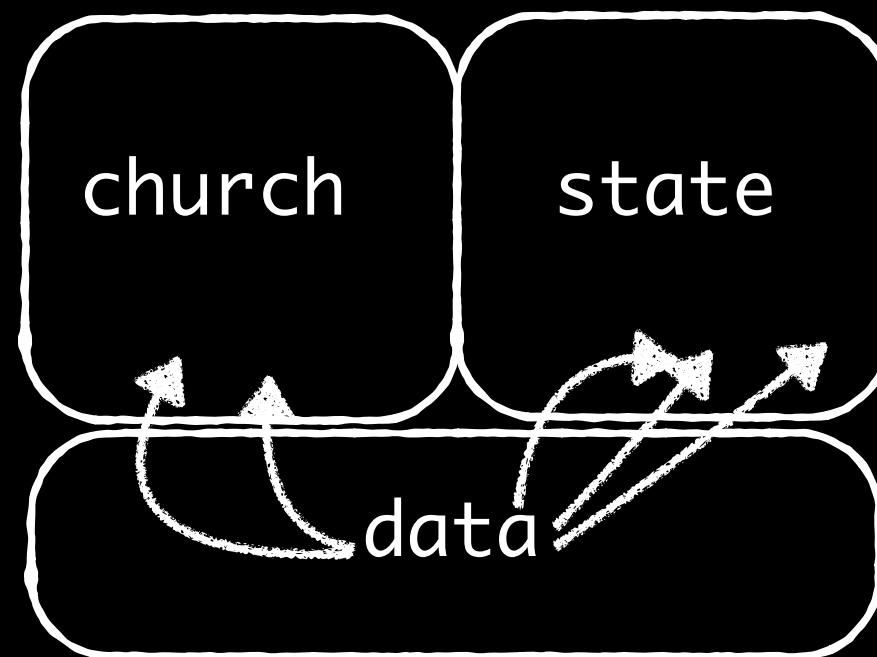
- Introduced by an “outsider”
- Less expensive than existing products
- Targeting underserved or new markets
- Initially inferior to existing products
- Advanced by an enabling technology



news: 21st century



# news: 21st century



learnings

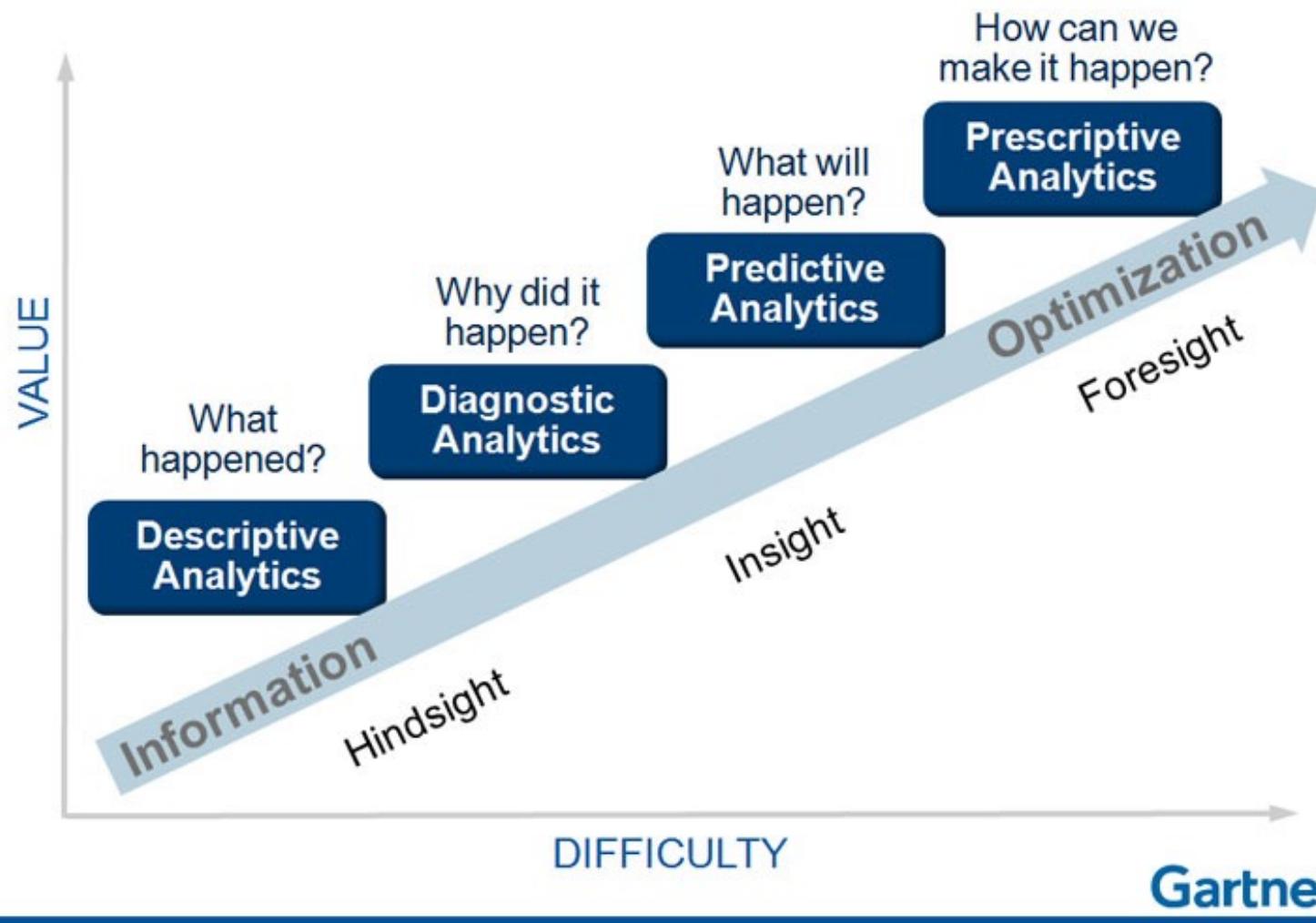
# learnings

- predictive modeling
- descriptive modeling
- prescriptive modeling

(actually ML, shhh...)

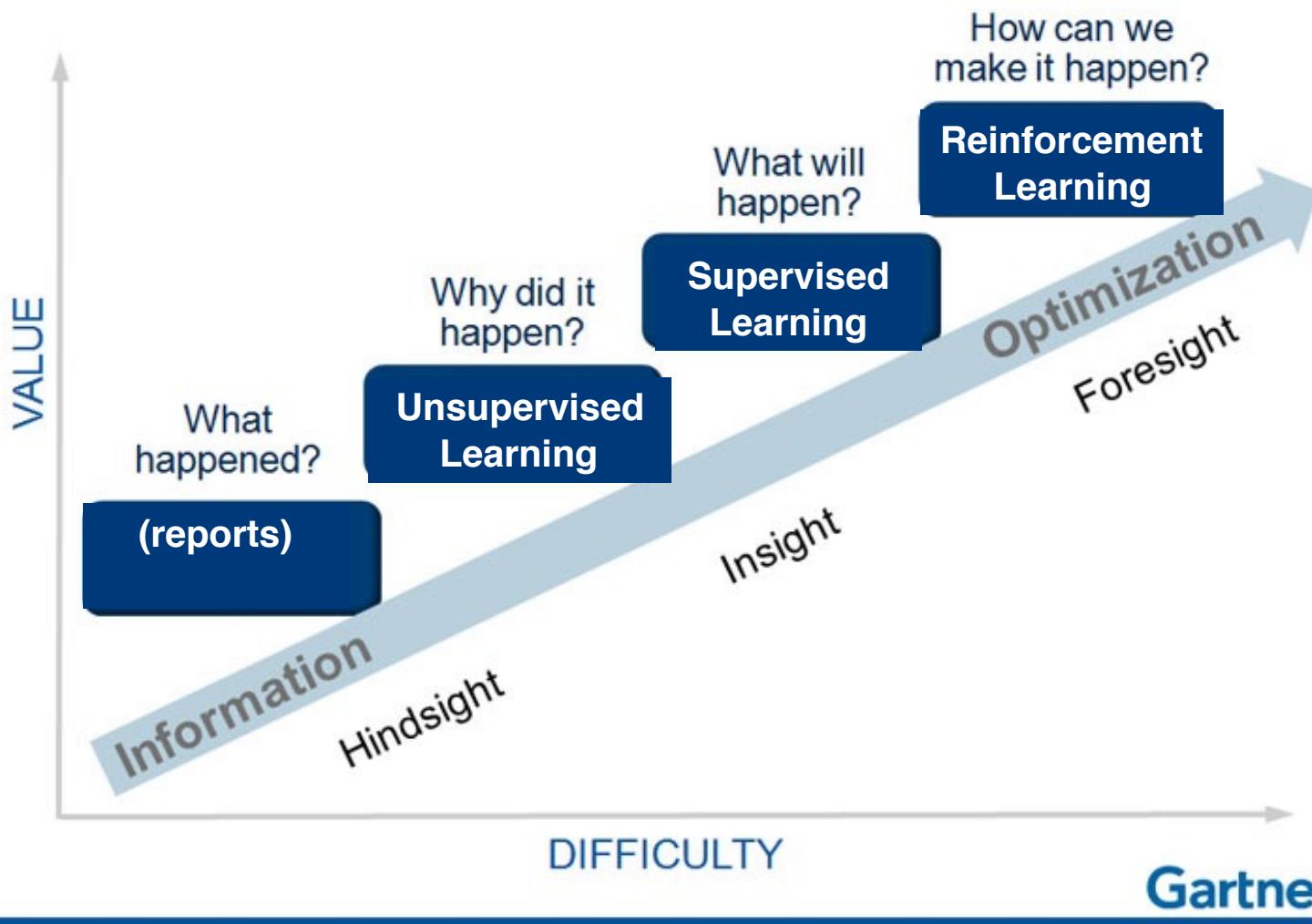
- (supervised learning)
- (unsupervised learning)
- (reinforcement learning)

# Analytic Value Escalator



h/t michael littman

# Analytic Value Escalator

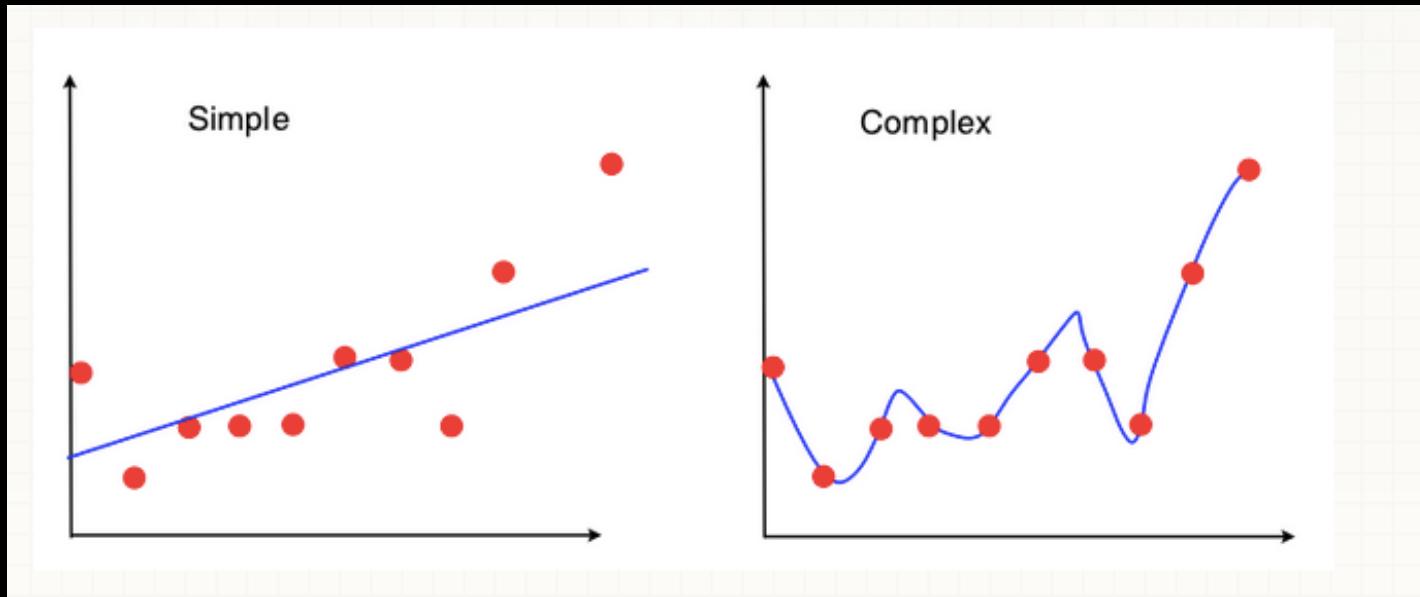


h/t michael littman

# learnings

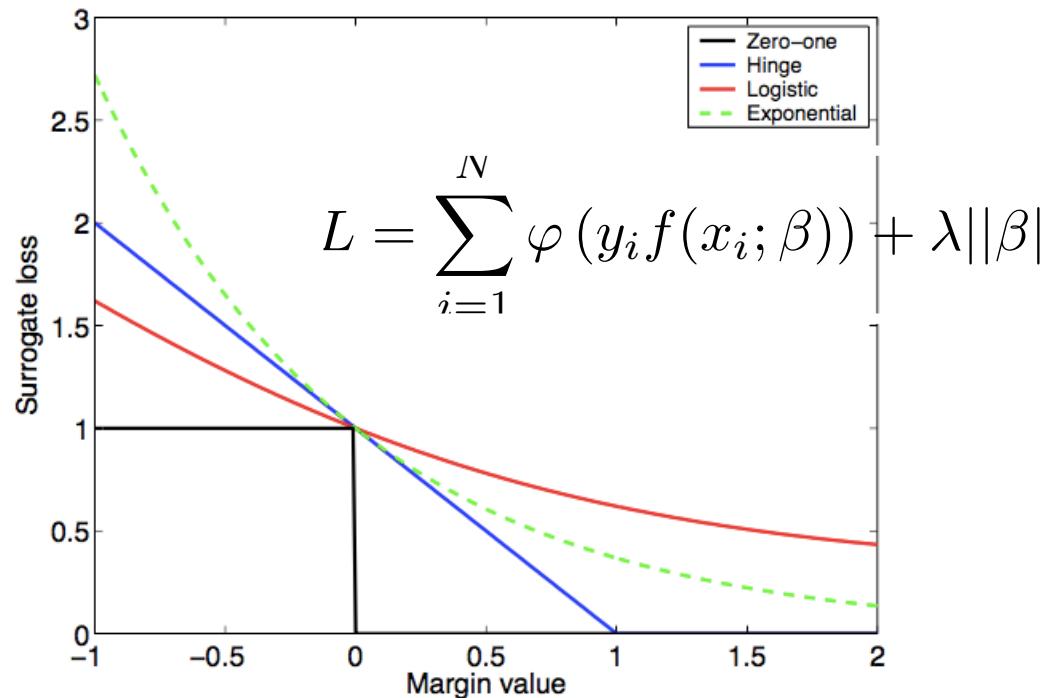
- predictive modeling
- descriptive modeling
- prescriptive modeling

cf. [modelingsocialdata.org](http://modelingsocialdata.org)



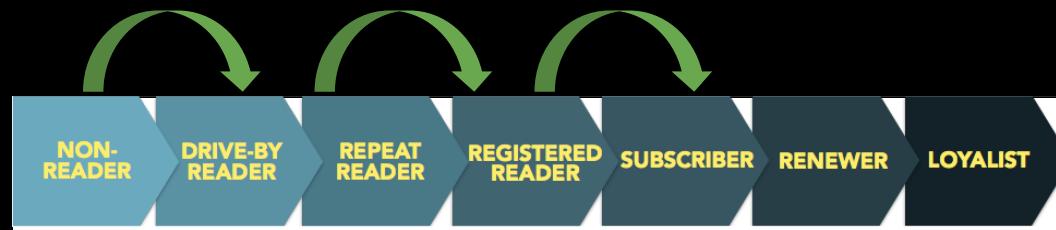
[stats.stackexchange.com](https://stats.stackexchange.com)

## Margin-Based Surrogate Loss Functions



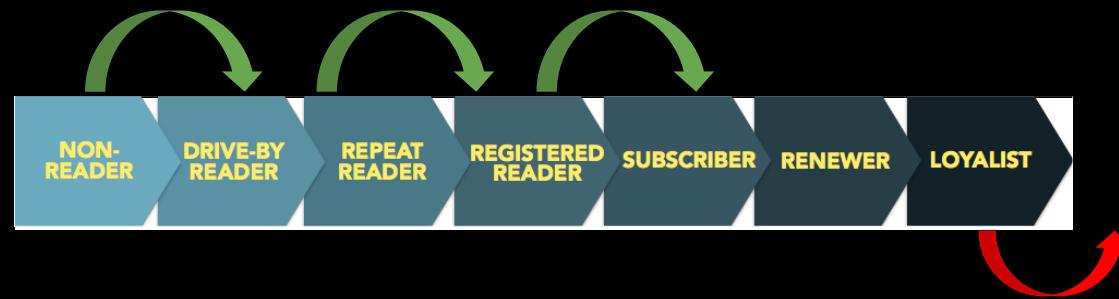
from “are you a bayesian or a frequentist”  
–michael jordan

*predictive modeling, e.g.,*



cf. [modelingsocialdata.org](http://modelingsocialdata.org)

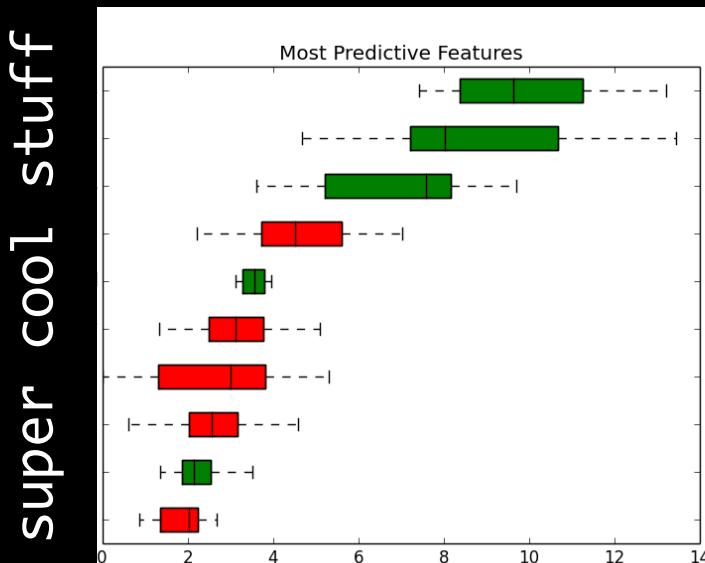
*predictive modeling, e.g.,*



“the funnel”

cf. [modelingsocialdata.org](http://modelingsocialdata.org)

# interpretable predictive modeling



cf. [modelingsocialdata.org](http://modelingsocialdata.org)

interpreting

super cool stuff

cf. me

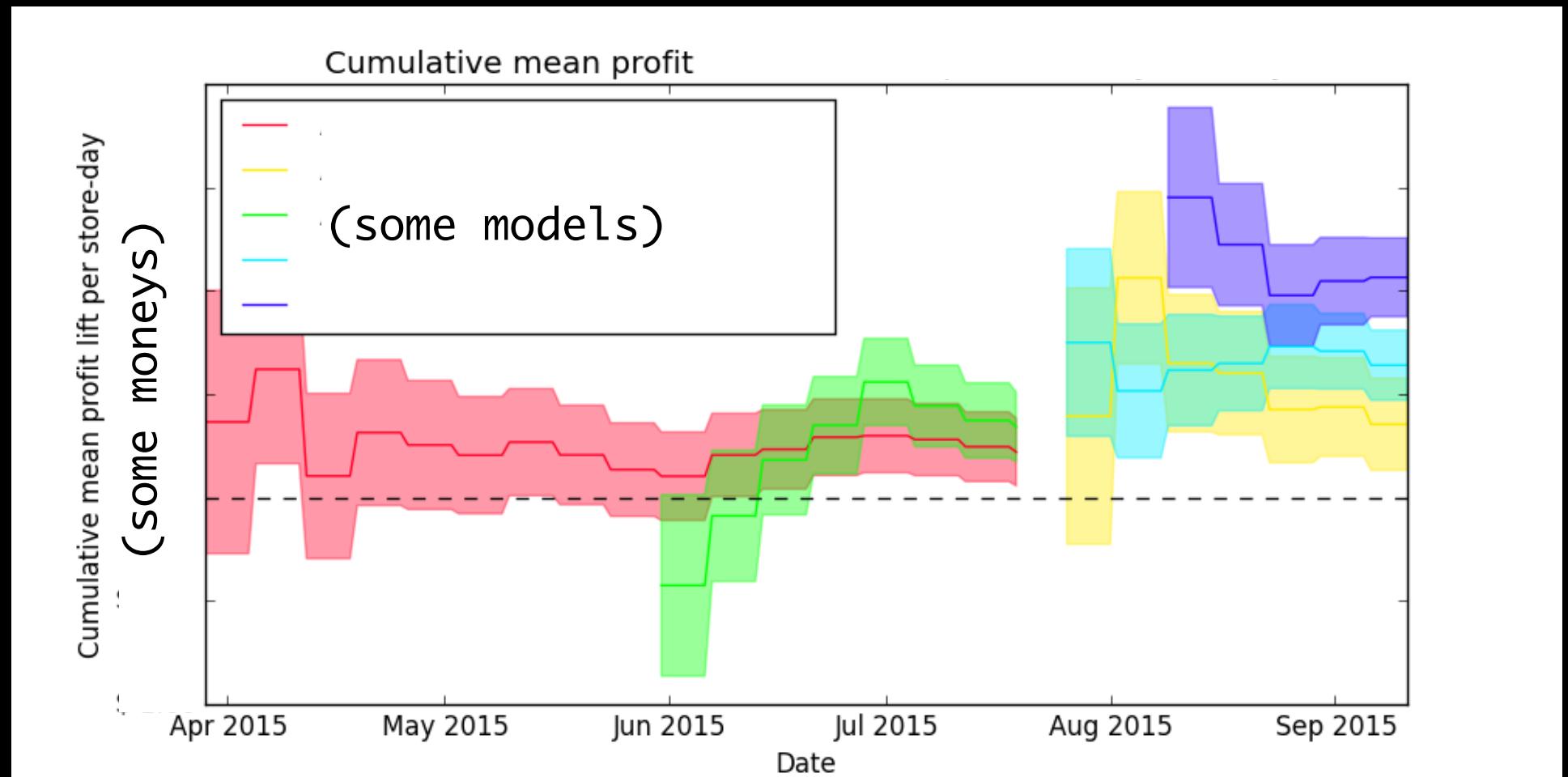
TFNAME	DB-MOTIF	MOTIF	DBNAME	d(p,q)
CBF1	CACGTG		YPD	0.032635
CGG everted repeat	CGGN*CCG		YPD	0.032821
MSN2			TRANSFAC	0.085626
HSF1	TTCNNNGAA		SCPD	0.102410
XBP1			TRANSFAC	0.140561
STE12			TRANSFAC	0.256750
GCN4			SCPD	0.292221
TBP			TRANSFAC	0.376601
HAP1	CGGNNTWNCGG		YPD	0.423004
RAP1	RMACCCA		SCPD	0.523059
mPAC			AlignACE	0.552493
mRRPE			AlignACE	0.630740
PHO4			TRANSFAC	0.672961
YAP1			TRANSFAC	0.777816
MIG1	CCCCCACAAA		YPD	0.799412
MET31,32	AAACTGTGG		YPD	0.84893
HAP2,3,4			TRANSFAC	1.070837

optimization & learning, e.g.,



“How The New York Times Works” popular mechanics, 2015

optimization & prediction, e.g.,



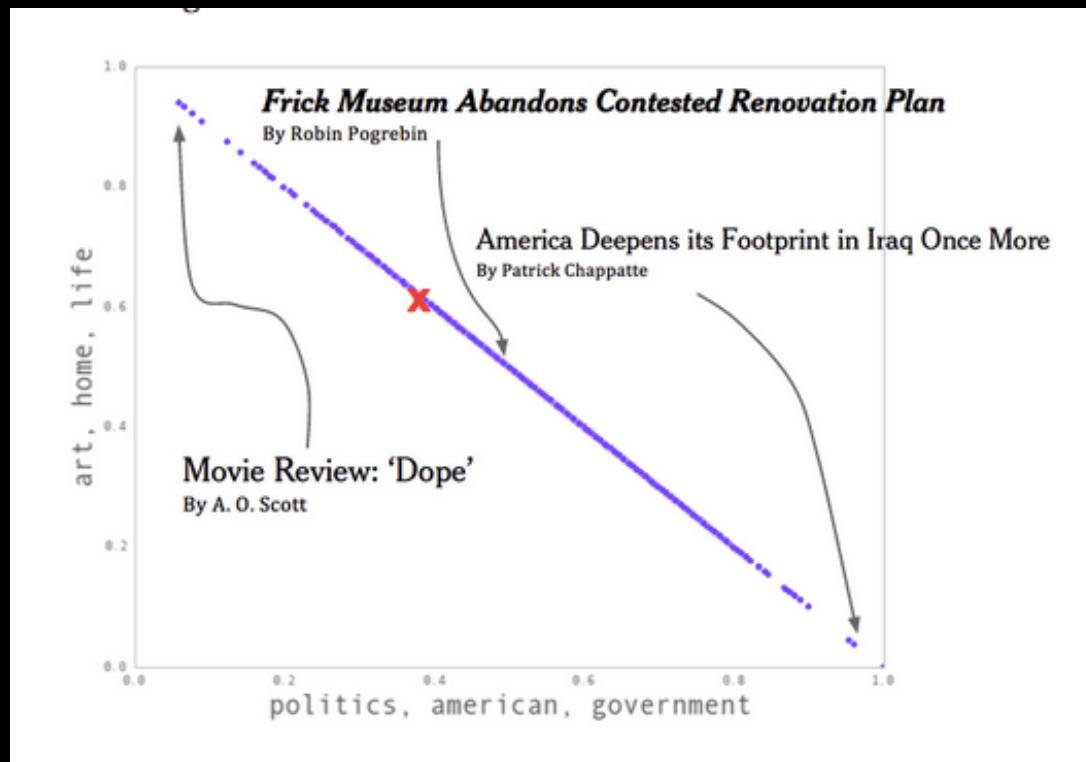
“newsvendor problem,” literally (+prediction+experiment)

# recommendation as inference

MOST EMAILED   MOST VIEWED   RECOMMENDED FOR YOU

- 1. THE OUTLAW OCEAN  
A Renegade Trawler, Hunted for 10,000 Miles by Vigilantes 
- 2. Campus Suicide and the Pressure of Perfection 
- 3. As Tech Booms, Workers Turn to Coding for Career Change 
- 4. Prison Worker Who Aided Escape Tells of Sex, Saw Blades and Deception 
- 5. Under Oath, Donald Trump Shows His Raw Side 
- 6. American Hunter Killed Cecil, Beloved Lion That Was Lured Out of Its Sanctuary 
- 7. A Creature on the Loose Puts Milwaukee Residents on Edge 
- 8. N.F.L. Upholds Tom Brady's Ban; Cellphone's Fate Helped Make the Call 
- 9. Escalator Death in China Sets Off Furor Online 
- 10. DAVID BROOKS  
The Structure of Gratitude 

## recommendation as inference



[bit.ly/AlexCTM](http://bit.ly/AlexCTM)

descriptive modeling, e.g,

“segments”

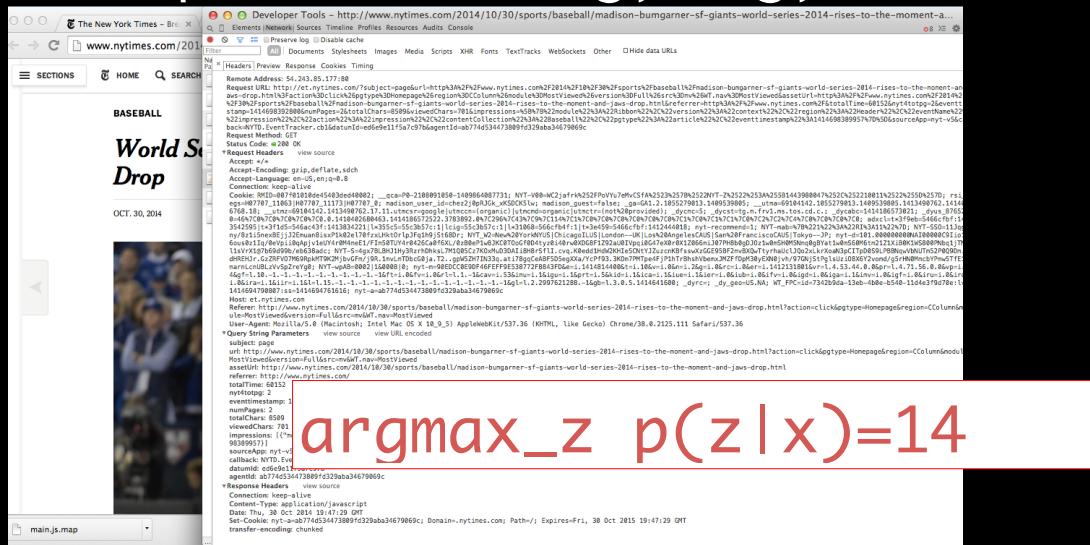
cf. [modelingsocialdata.org](http://modelingsocialdata.org)

descriptive modeling, e.g.,

# “segments”

cf. [modelingsocialdata.org](http://modelingsocialdata.org)

descriptive modeling, e.g.,

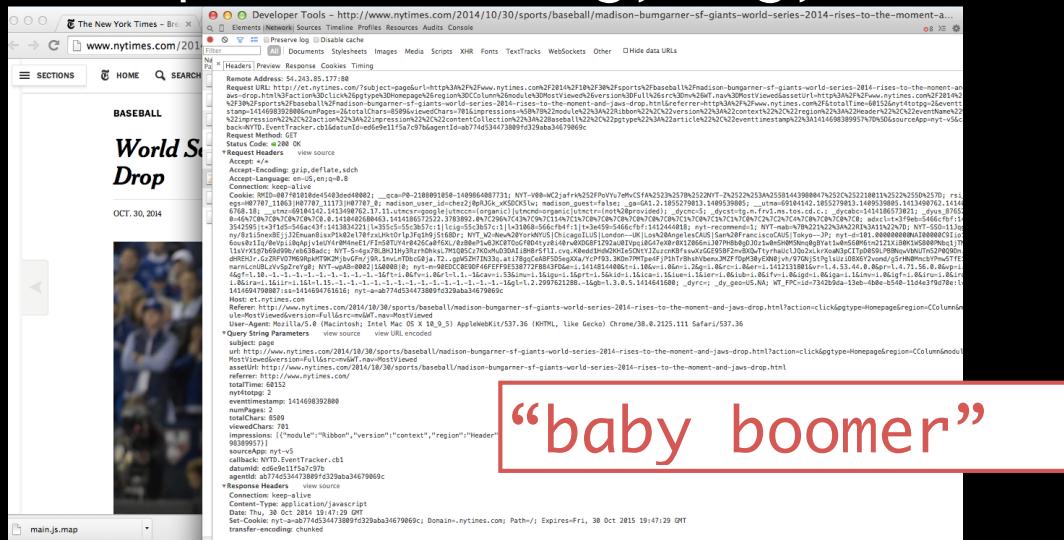


argmax\_z p(z|x)=14

# “segments”

cf. [modelingsocialdata.org](http://modelingsocialdata.org)

descriptive modeling, e.g.,



# “segments”

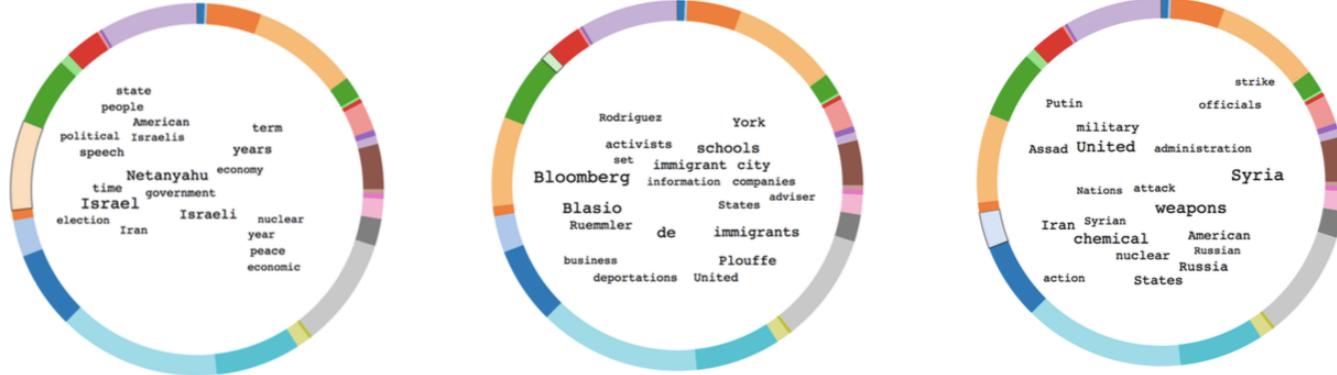
cf. [modelingsocialdata.org](http://modelingsocialdata.org)

- descriptive data product

### A Quick Refinery Demo

The New York Times Developers Article Search API v2 →  Extracting NYT articles from keyword "obama" in 2013. → 

What themes / topics defined the Obama administration during 2013?

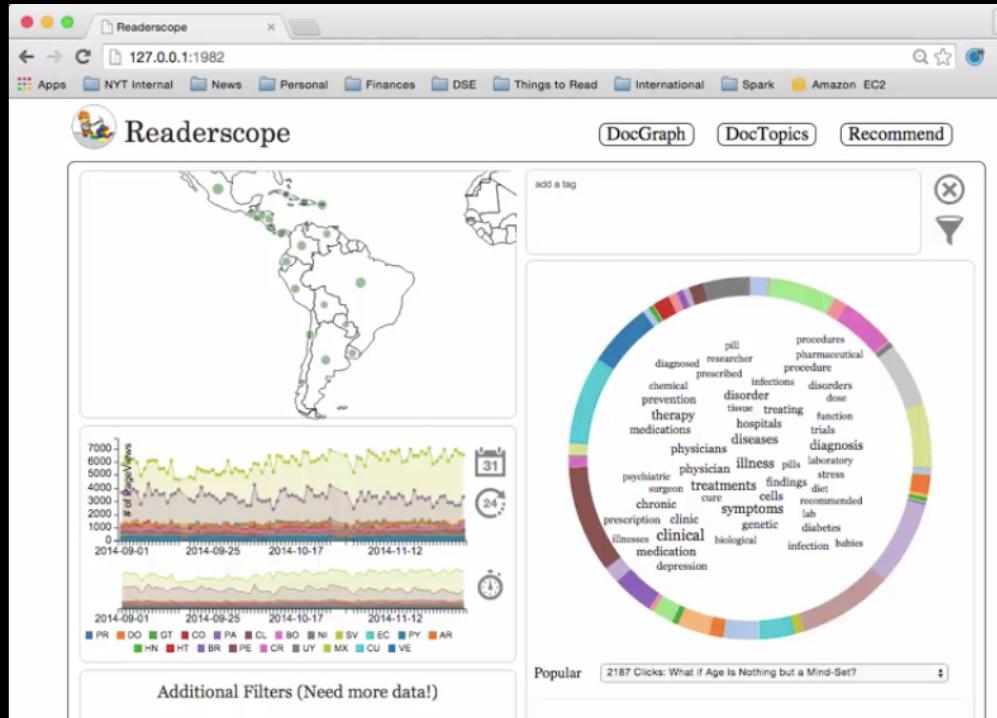


The three sunburst charts illustrate the most prominent themes in news articles about the Obama administration in 2013. The inner ring represents the primary theme, and the outer ring represents sub-themes.

- Chart 1 (Left):** Focuses on international relations and politics. Key terms include: state, people, American, political, Israelis, speech, Netanyahu, economy, time, government, Israel, election, Iran, Israeli, nuclear, year, peace, economic.
- Chart 2 (Middle):** Focuses on domestic politics and administration. Key terms include: Rodriguez, Bloomberg, Blasio, Ruemmler, de, business, deportations, Plouffe, United, activists, set, immigrant, city, information, companies, schools, York, York, immigrants, United, adviser.
- Chart 3 (Right):** Focuses on international conflicts and military actions. Key terms include: Putin, Assad, United, Syria, Nations, attack, weapons, Iranian, Syrian, chemical, nuclear, American, Russian, Russia, States, strike, officials, military, administration, action.

cf. [daeilkim.com](http://daeilkim.com)

descriptive modeling, e.g.,



cf. [daeilkim.com](http://daeilkim.com) ; import bnpy

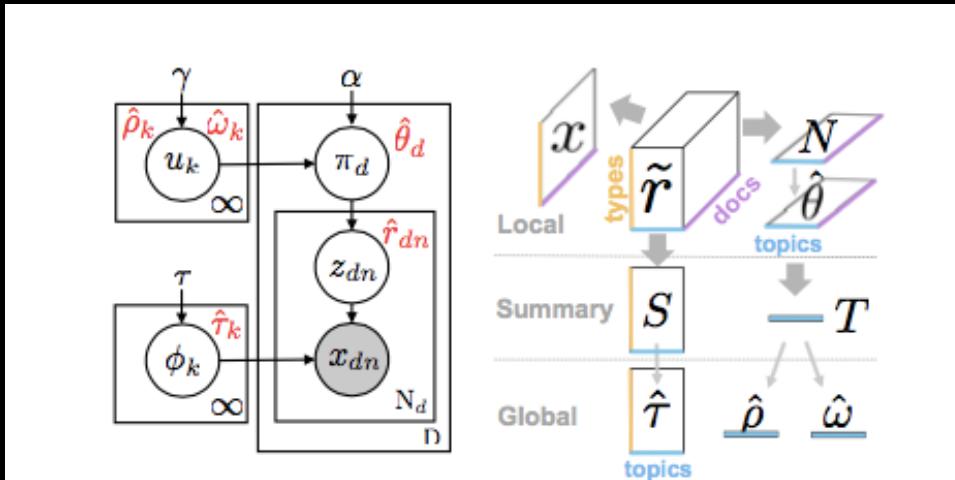


Figure 1: *Left:* Directed graphical model for the HDP admixture (Sec. 2). Free parameters for mean-field variational inference (Sec. 3) shown in red. *Right:* Flow chart for our inference algorithm, specialized for bag-of-words data, where we can use sparse type-based assignments  $\tilde{r}$  instead of per-token variables  $\hat{r}$ . We define  $\tilde{r}_{dwk}$  to be the total mass of all tokens in document  $d$  of type  $w$  assigned to  $k$ :  $\tilde{r}_{dwk} = \sum_{n=1}^{N_d} \hat{r}_{dnk} \delta_{x_{dn}, w}$ . Updates flow from  $\tilde{r}$  to global topic-type parameters  $\hat{\tau}$  and (separately) to global topic weight parameters  $\hat{\rho}, \hat{\omega}$ . Each variable's shape gives its dimensionality. Thick arrows indicate summary statistics; thin arrows show free parameter updates.

modeling your audience  
[bit.ly/Hughes-Kim-Sudderth-AISTATS15](http://bit.ly/Hughes-Kim-Sudderth-AISTATS15)

**Objective function.** Mean field methods optimize an evidence lower bound  $\log p(x|\gamma, \alpha, \tau) \geq \mathcal{L}(\cdot)$ , where

$$\mathcal{L}(\cdot) \triangleq \mathcal{L}_{data}(\cdot) + H_z(\cdot) + \mathcal{L}_{HDP}(\cdot) + \mathcal{L}_u(\cdot). \quad (4)$$

The final term  $\mathcal{L}_u(\cdot)$ , which depends only on  $q(u)$ , is discussed in the next section. The first three terms account for data generation, the assignment entropy, and the document-topic allocations. These are defined below, with expectations taken with respect to Eq. (3):

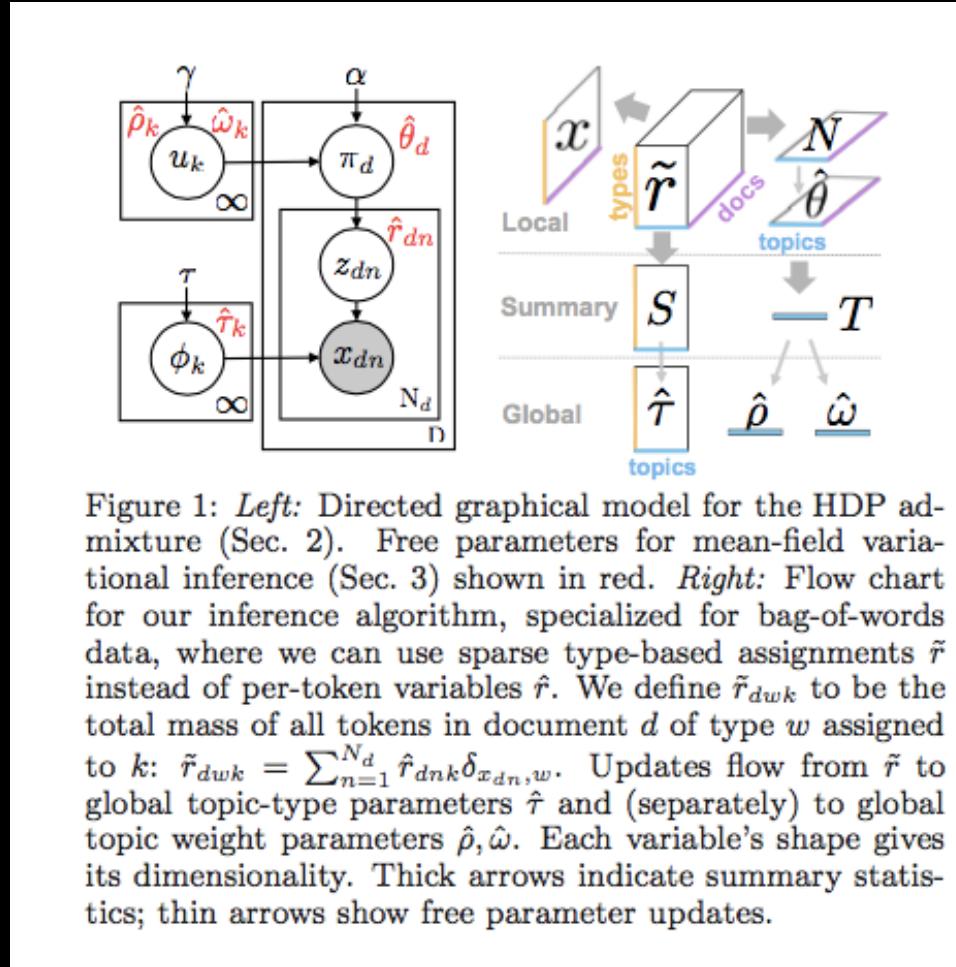
$$\mathcal{L}_{data}(\cdot) \triangleq \mathbb{E}_q[\log p(x|z, \phi) + \log \frac{p(\phi|\bar{\tau})}{q(\phi|\hat{\tau})}], \quad (5)$$

$$H_z(\cdot) \triangleq -\sum_{k=1}^K \sum_{d=1}^D \sum_{n=1}^{N_d} \hat{r}_{dnk} \log \hat{r}_{dnk},$$

$$\mathcal{L}_{HDP}(\cdot) \triangleq \mathbb{E}_q \left[ \log \frac{p(z|\pi)p(\pi|\alpha, u)}{q(\pi|\hat{\theta})} \right].$$

The forms of  $\mathcal{L}_{data}$  and  $H_z$  are unchanged from the simpler case of mean-field for DP mixtures. Closed-form expressions are in the Supplement.

modeling your audience  
(optimization, ultimately)



modeling your audience  
also allows insight+targeting as inference

prescriptive modeling

## prescriptive modeling

---

descriptive:	specify $x$ ; learn $z(x)$ or $p(z x)$ where $z$ is “simpler” than $x$
predictive:	specify $x$ and $y$ ; learn to predict $y$ from $x$
prescriptive:	specify $x, y$ , and $a$ ; learn to prescribe $a$ given $x$ to maximize $y$

---

prescriptive modeling

$$V = E_+(y) = \sum_{yax} y P_+(y, a, x)$$

“off policy value estimation”  
(cf. “causal effect estimation”)

$$\hat{V} = \frac{1}{N} \sum_{i=1}^{i=N} y_i \frac{\mathbf{1}(a_i = h(x_i))}{\hat{B}(a_i|x_i)}$$

cf. Langford `08-`16;  
Horvitz & Thompson `52;  
Holland `86

“off policy value estimation”  
(cf. “causal effect estimation”)

$$\hat{V} = \frac{1}{N} \sum_{i=1}^{i=N} y_i \frac{1(a_i = h(x_i))}{\hat{B}(a_i|x_i)}$$

Vapnik's razor

“ When solving a (learning) problem of interest,  
do not solve a more complex problem as an  
intermediate step.”

# prescriptive modeling

## Experimental evidence of massive-scale emotional contagion through social networks

Adam D. I. Kramer<sup>a,1</sup>, Jamie E. Guillory<sup>b,2</sup>, and Jeffrey T. Hancock<sup>b,c</sup>

Author Affiliations 

<sup>a</sup>Core Data Science Team, Facebook, Inc., Menlo Park, CA 94025; and

Departments of <sup>b</sup>Communication and

<sup>c</sup>Information Science, Cornell University, Ithaca, NY 14853

Edited by Susan T. Fiske, Princeton University, Princeton, NJ, and approved March 25, 2014 (received for review October 23, 2013)

A correction has been published

A correction has been published

cf. [modelingsocialdata.org](http://modelingsocialdata.org)

# prescriptive modeling

## Experimental evidence of massive-scale emotional contagion through social networks

Adam D. I. Kramer<sup>a,1</sup>, Jamie E. Guillory<sup>b,2</sup>, and Jeffrey T. Hancock<sup>b,c</sup>

Author Affiliations 

<sup>a</sup>Core Data Science Team, Facebook, Inc., Menlo Park, CA 94025; and

Departments of <sup>b</sup>Communication and

<sup>c</sup>Information Science, Cornell University, Ithaca, NY 14853

Edited by Susan T. Fiske, Princeton University, Princeton, NJ, and approved March 25, 2014 (received for review October 23, 2013)

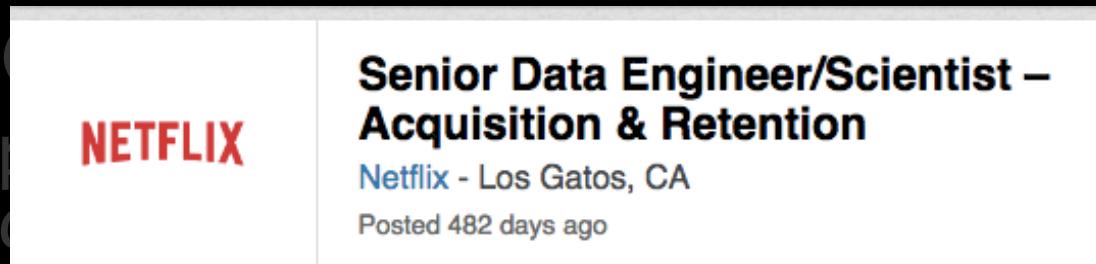
A correction has been published

A correction has been published

aka “A/B testing”;  
RCT

cf. [modelingsocialdata.org](http://modelingsocialdata.org)

prescriptive modeling: from A/B to....



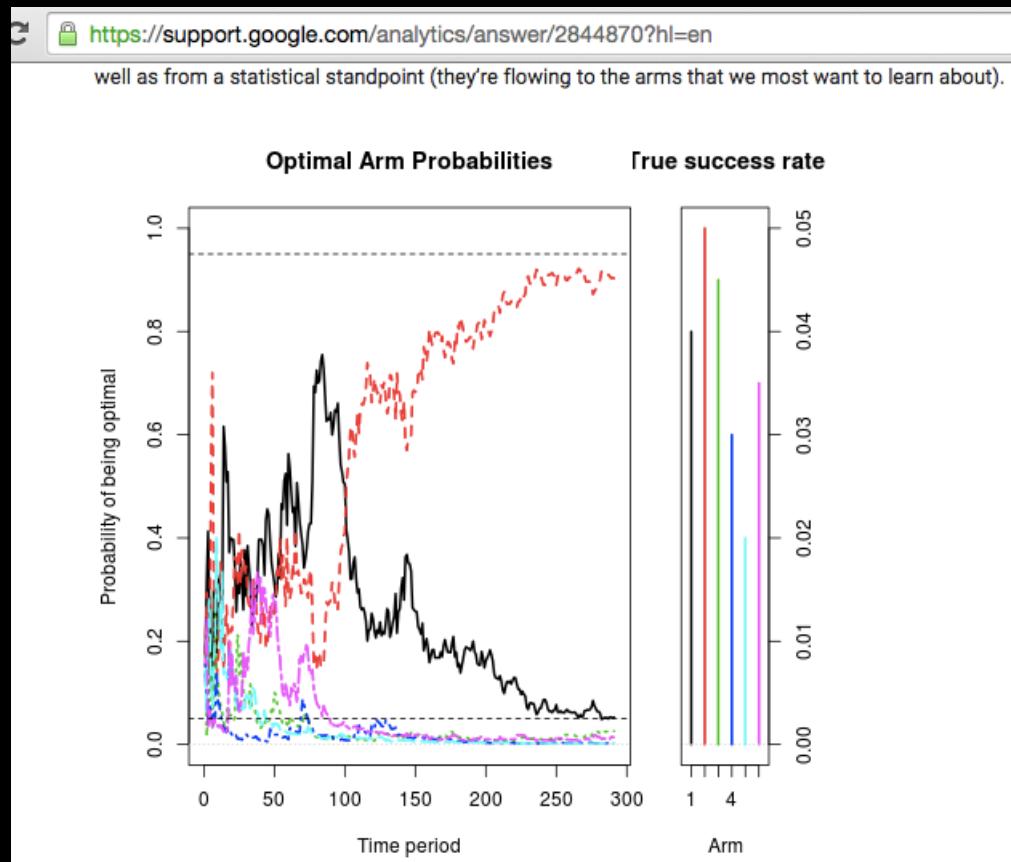
testing”; Test

Some of the most recognizable personalization in our service is the collection of “genre” rows. . . Members connect with these rows so well that we measure an **increase in member retention** by **placing the most tailored rows higher** on the page instead of lower.

business as usual Reporting

cf. [modelingsocialdata.org](http://modelingsocialdata.org)

## real-time A/B -> “bandits”



cf. [modelingsocialdata.org](http://modelingsocialdata.org)

*prescriptive modeling, e.g.,*

prescriptive modeling, e.g.,



**Colin Russel** 10:21 AM

!blossom facebook? all



**blossombot** BOT 10:21 AM ★

Blossom has the following suggestions for your next Facebook posts:

**nytimes:** <http://www.nytimes.com/2015/08/12/opinion/frank-bruni-can-we-interest-you-in-teaching.html>

**nytopinion:** <http://www.nytimes.com/2015/08/12/opinion/when-innocence-is-no-defense.html>

**nytpolitics:** <http://www.nytimes.com/2015/08/16/magazine/president-obamas-letter-to-the-editor.html>

**upshot:** <http://www.nytimes.com/interactive/2015/07/03/upshot/a-quick-puzzle-to-test-your-problem-solving.html>

**tmagazin:** Blossom currently has no suggestions

**nytimestravel:** <http://www.nytimes.com/2015/08/13/us/politics/us-jets-meet-limit-as-iraqi-ground-fight-against-isis-plods-on.html>

**nytimesscience:** <http://www.nytimes.com/interactive/2015/07/03/upshot/a-quick-puzzle-to-test-your-problem-solving.html>

**nytfood:** Blossom currently has no suggestions

**WellNYT:** <http://www.nytimes.com/2015/08/12/sports/kaci-lickteig-climbs-toward-the-top-in-ultrarunning.html>

**NewYorkTodayNYT:** <http://www.nytimes.com/2015/08/08/nyregion/she-answered-his-ad-for-a-roommate-moved-in-and-never-left.html>



prescriptive modeling, e.g.,

Colin Russel 10:21 AM  
!blossom facebook? all

blossombot BOT 10:21 AM ★  
Blossom has the following suggestions for your next Facebook posts:  
[nytimes: http://www.nytimes.com/2015/08/12/opinion/frank-bruni-can-we-interest-you-in-teaching.html](http://www.nytimes.com/2015/08/12/opinion/frank-bruni-can-we-interest-you-in-teaching.html)

blossombot BOT 3:40 PM ★  
Blossom Alert!  
This piece is predicted to go viral if posted on Facebook next:  
<http://www.nytimes.com/2015/08/12/sports/football/ikemefuna-enemkpali-the-backup-who-broke-geno-smiths-jaw.html>  
Here is its Stela Link

nytimesscience: <http://www.nytimes.com/interactive/2015/07/03/upshot/a-quick-puzzle-to-test-your-problem-solving.html>

nytfood: Blossom currently has no suggestions

WellNYT: <http://www.nytimes.com/2015/08/12/sports/kaci-lickteig-climbs-toward-the-top-in-ultrarunning.html>

NewYorkTodayNYT: <http://www.nytimes.com/2015/08/08/nyregion/she-answered-his-ad-for-a-roommate-moved-in-and-never-left.html>

+ ☺

prescriptive modeling, e.g.,

Colin Russel 10:21 AM  
!blossom facebook? all

blossombot BOT 10:21 AM ★  
Blossom has the following suggestions for your next Facebook posts:  
[nytimes: http://www.nytimes.com/2015/08/12/opinion/frank-bruni-can-we-interest-you-in-teaching.html](http://www.nytimes.com/2015/08/12/opinion/frank-bruni-can-we-interest-you-in-teaching.html)

blossombot BOT 3:40 PM ★  
Blossom Alert!  
This piece is predicted to go viral if posted on Facebook next:  
<http://www.nytimes.com/2015/08/12/sports/football/ikemefuna-enemkpali-the-backup-who-broke-geno-smiths-jaw.html>  
Here is its Stela Link

nytimesscience: <http://www.nytimes.com/interactive/2015/07/03/upshot/a-quick-puzzle-to-test-your-problem-solving.html>

nytfood: Blossom currently has no suggestions

WellNYT: <http://www.nytimes.com/2015/08/12/sports/kaci-lickteig-climbs-toward-the-top-in-ultrarunning.html>

NewYorkTodayNYT: <http://www.nytimes.com/2015/08/08/nyregion/she-answered-his-ad-for-a-roommate-moved-in-and-never-left.html>

leverage methods which are predictive yet performant

## NB: data-informed, not data-driven

 Colin Russel 10:21 AM  
!blossom facebook? all

 blossombot BOT 10:21 AM ★  
Blossom has the following suggestions for your next Facebook posts:  
[nytimes: http://www.nytimes.com/2015/08/12/opinion/frank-bruni-can-we-interest-you-in-teaching.html](http://www.nytimes.com/2015/08/12/opinion/frank-bruni-can-we-interest-you-in-teaching.html)

 blossombot BOT 3:40 PM ★  
**Blossom Alert!**  
This piece is predicted to go viral if posted on Facebook next:  
<http://www.nytimes.com/2015/08/12/sports/football/ikemefuna-enemkpali-the-backup-who-broke-geno-smiths-jaw.html>  
Here is its [Stela Link](#)

nytimesscience: <http://www.nytimes.com/interactive/2015/07/03/upshot/a-quick-puzzle-to-test-your-problem-solving.html>  
nytfood: Blossom currently has no suggestions  
WellNYT: <http://www.nytimes.com/2015/08/12/sports/kaci-lickteig-climbs-toward-the-top-in-ultrarunning.html>  
NewYorkTodayNYT: <http://www.nytimes.com/2015/08/08/nyregion/she-answered-his-ad-for-a-roommate-moved-in-and-never-left.html>

predicting views/cascades: doable?

## **Optimizing Web Traffic via the Media Scheduling Problem**

Lars Backstrom\*  
Dept. of Computer Science  
Cornell University  
Ithaca, NY 14853.  
[lars@cs.cornell.edu](mailto:lars@cs.cornell.edu)

Jon Kleinberg†  
Dept. of Computer Science  
Cornell University  
Ithaca, NY 14853.  
[kleinber@cs.cornell.edu](mailto:kleinber@cs.cornell.edu)

Ravi Kumar  
Yahoo! Research  
701 First Ave.  
Sunnyvale, CA 94089.  
[ravikumar@yahoo-inc.com](mailto:ravikumar@yahoo-inc.com)

KDD 09: how many people are online?

predicting views/cascades: doable?

## Can cascades be predicted?

Justin Cheng  
Stanford University  
[jcccf@cs.stanford.edu](mailto:jcccf@cs.stanford.edu)

Lada A. Adamic  
Facebook  
[ladamic@fb.com](mailto:ladamic@fb.com)

P. Alex Dow  
Facebook  
[adow@fb.com](mailto:adow@fb.com)

Jon Kleinberg  
Cornell University  
[kleinber@cs.cornell.edu](mailto:kleinber@cs.cornell.edu)

Jure Leskovec  
Stanford University  
[jure@cs.stanford.edu](mailto:jure@cs.stanford.edu)

WWW 14: FB shares

# predicting views/cascades: features?

Content Features	
$score_{food/nature/\dots}$	The probability of the photo having a specific feature (food, overlaid text, landmark, nature, etc.)
$is\_en$	Whether the photo was posted by an English-speaking user or page
$has\_caption$	Whether the photo was posted with a caption
$lwc_{pos/neg/soc}$	Proportion of words in the caption that expressed positive or negative emotion, or sociality, if English
Root (Original Poster) Features	
$views_{0..k}$	Number of users who saw the original photo until the $k$ th reshare was posted
$orig\_is\_page$	Whether the original poster is a page
$outdeg(v_0)$	Friend, subscriber or fan count of the original poster
$age_0$	Age of the original poster, if a user
$gender_0$	Gender of the original poster, if a user
$fb\_age_0$	Time since the original poster registered on Facebook, if a user
$activity_0$	Average number of days the original poster was active in the past month, if a user
Resharer Features	
$views_{1..k-1..k}$	Number of users who saw the first $k - 1$ reshares until the $k$ th reshare was posted
$pages_k$	Number of pages responsible for the first $k$ reshares, including the root, or $\sum_{i=0}^k \mathbb{1}\{v_i \text{ is a page}\}$
$friends_k^{avg/90p}$	Average or 90th percentile friend count of the first $k$ reshancers, or $\frac{1}{k} \sum_{i=1}^k outdeg_{friends}(v_i) \mathbb{1}\{v_i \text{ is a user}\}$
$fans_k^{avg/90p}$	Average or 90th percentile fan count of the first $k$ reshancers, or $\frac{1}{k} \sum_{i=1}^k outdeg(v_i) \mathbb{1}\{v_i \text{ is a page}\}$
$subscribers_k^{avg/90p}$	Average or 90th percentile subscriber count of the first $k$ reshancers, or $\frac{1}{k} \sum_{i=1}^k outdeg_{subscriber}(v_i) \mathbb{1}\{v_i \text{ is a user}\}$
$fb\_ages_k^{avg/90p}$	Average or 90th percentile time since the first $k$ reshancers registered on Facebook, or $\frac{1}{k} \sum_{i=1}^k fb\_age_i$
$activities_k^{avg/90p}$	Average number of days the first $k$ reshancers were active in July, or $\frac{1}{k} \sum_{i=1}^k activity_i$
$ages_k^{avg/90p}$	Average age of the first $k$ reshancers, or $\frac{1}{k} \sum_{i=1}^k age_i$
$female_k$	Number of female users among the first $k$ reshancers, or $\sum_{i=1}^k \mathbb{1}\{gender_i \text{ is female}\}$
Structural Features	
$outdeg(v_i)$	Connection count (sum of friend, subscriber and fan counts) of the $i$ th resharer (or out-degree of $v_i$ on $G = (V, E)$ )
$outdeg(v'_i)$	Out-degree of the $i$ th reshare on the induced subgraph $G' = (V', E')$ of the first $k$ reshancers and the root
$outdeg(\hat{v}_i)$	Out-degree of the $i$ th reshare on the reshare graph $\hat{G} = (\hat{V}, \hat{E})$ of the first $k$ reshares
$orig\_connections_k$	Number of first $k$ reshancers who are friends with, or fans of the root, or $ \{v_i \mid (v_0, v_i) \in E, 1 \leq i \leq k\} $
$border\_nodes_k$	Total number of users or pages reachable from the first $k$ reshancers and the root, or $ \{v_i \mid (v_i, v_j) \in E, 0 \leq i, j \leq k\} $
$border\_edges_k$	Total number of first-degree connections of the first $k$ reshancers and the root, or $ \{(v_i, v_j) \mid (v_i, v_j) \in E, 0 \leq i, j \leq k\} $
$subgraph'_k$	Number of edges on the induced subgraph of the first $k$ reshancers and the root, or $ \{(v_i, v_j) \mid (v_i, v_j) \in E', 0 \leq i, j \leq k\} $
$depth'_k$	Change in tree depth of the first $k$ reshares, or $\min_\beta \sum_{i=1}^k (depth_i - \beta i)^2$
$depths_k^{avg/90p}$	Average or 90th percentile tree depth of the first $k$ reshares, or $\frac{1}{k} \sum_{i=1}^k depth_i$
$did\_leave$	Whether any of the first $k$ reshares are not first-degree connections of the root
Temporal Features	
$time_i$	Time elapsed between the original post and the $i$ th reshare
$time'_{1..k/2}$	Average time between reshares, for the first $k/2$ reshares, or $\frac{1}{k/2-1} \sum_{i=1}^{k/2-1} (time_{i+1} - time_i)$
$time'_{k/2..k}$	Average time between reshares, for the last $k/2$ reshares, or $\frac{1}{k/2-1} \sum_{i=k/2}^{k-1} (time_{i+1} - time_i)$
$time''_{1..k}$	Change in the time between reshares of the first $k$ reshares, or $\min_\beta \sum_{i=1}^{k-1} (time_{i+1} - time_i - \beta i)^2$
$views'_{0..k}$	Number of users who saw the original photo, until the $k$ th reshare was posted, per unit time, or $\frac{views_{0..k}}{time_k}$
$views'_{1..k-1..k}$	Number of users who saw the first $k - 1$ reshares, until the $k$ th reshare was posted, per unit time, or $\frac{views_{1..k-1..k}}{time_k}$

## predicting views/cascades: features?

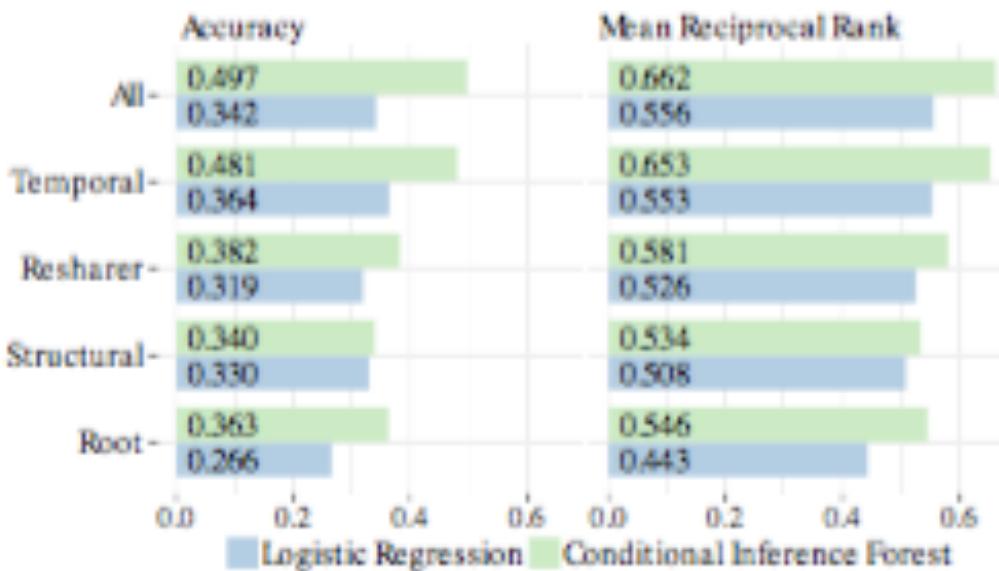


Figure 10: In predicting the largest cascade in clusters of 10 or more cascades of identical photos, we perform significantly above the baseline of 0.1.

# **predicting views/cascades: doable?**

## **Exploring limits to prediction in complex social systems**

**Travis Martin**

University of Michigan  
Dept. of Computer Science  
Ann Arbor, MI  
[travisbm@umich.edu](mailto:travisbm@umich.edu)

**Jake M. Hofman**

Microsoft Research  
641 6th Ave, Floor 7  
New York, NY  
[jmh@microsoft.com](mailto:jmh@microsoft.com)

**Amit Sharma**

Microsoft Research  
[amshar@microsoft.com](mailto:amshar@microsoft.com)

**Ashton Anderson**

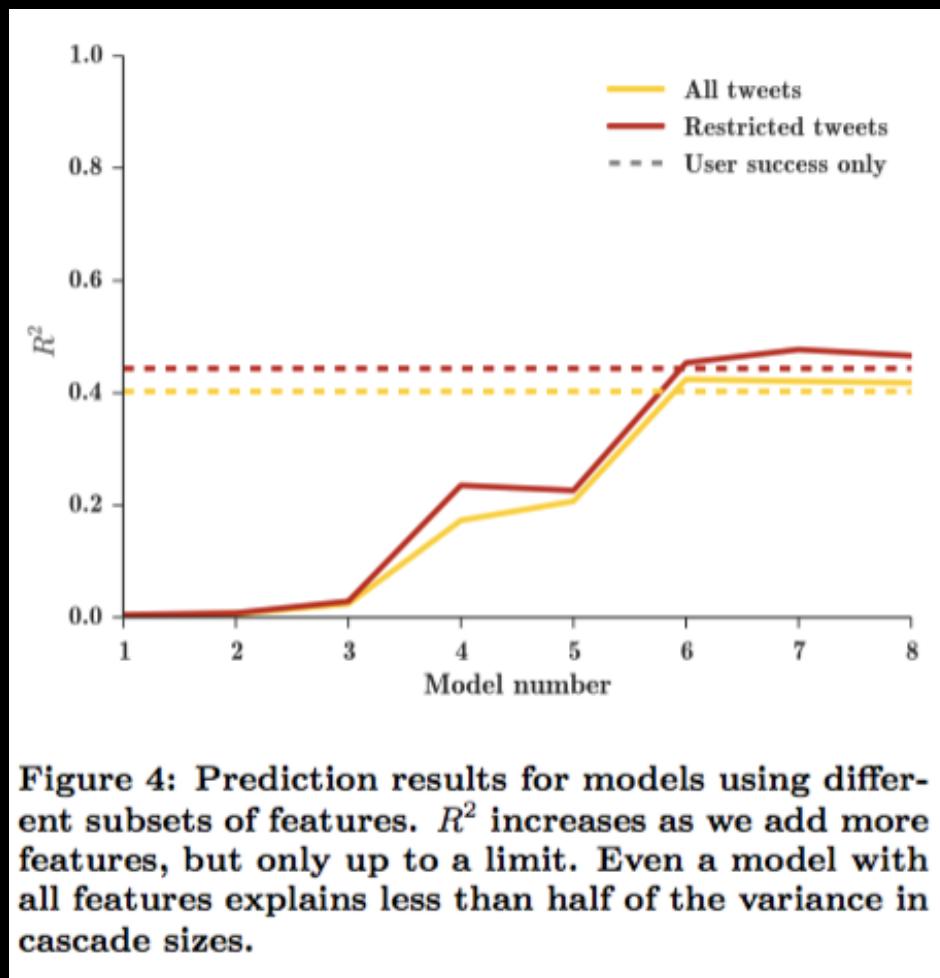
Microsoft Research  
[ashton@microsoft.com](mailto:ashton@microsoft.com)

**Duncan J. Watts**

Microsoft Research  
[duncan@microsoft.com](mailto:duncan@microsoft.com)

**WWW 16: TWIT RT's**

## predicting views/cascades: doable?



descriptive:

predictive:

prescriptive:

Explore

Learning

Test

Optimizing

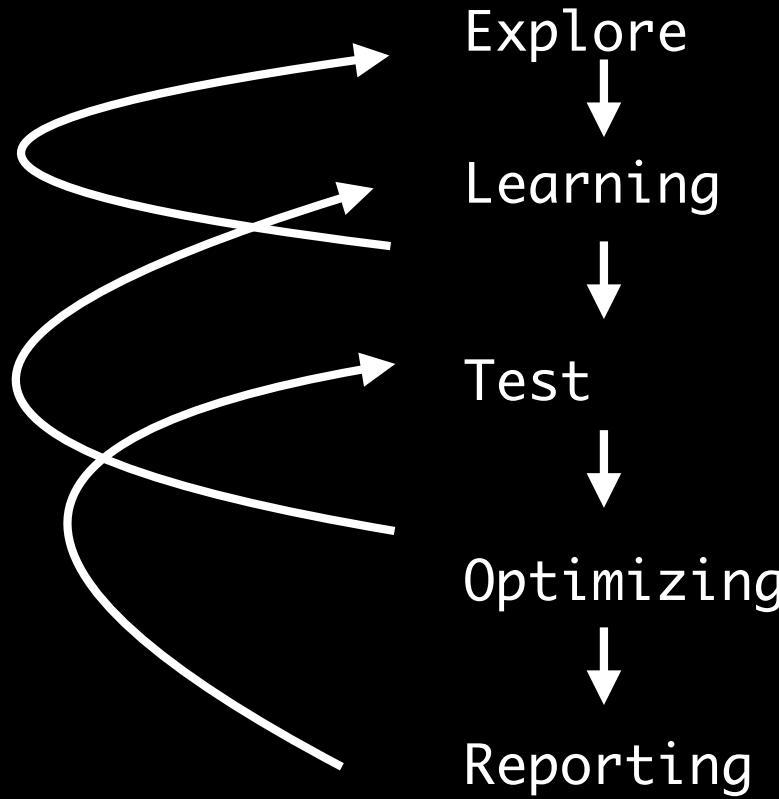
Reporting



descriptive:

predictive:

prescriptive:



things:  
what does DS team deliver?

- build data product
- build APIs
- impact roadmaps

data science @ The New York Times



chris.wiggins@columbia.edu  
chris.wiggins@nytimes.com  
@chrishwiggins

references: <http://bit.ly/stanf16>

