# A Day in the Life of a Data Scientist

Brian Eoff (@bde)
Lead Scientist at Spring

SPRING

# What is data science?

**Jake Vanderplas**
@jakevdp

Follow

Best thing about being a data scientist: nobody actually knows what data science *is*, so you can pretty much do whatever you want.

RETWEETS
65

FAVORITES
72

9:28 AM - 17 Oct 2014

# Data Science Conference Bingo Card

| | | | | |
|---|---|---|---|---|
| In-Memory | Unstructured Data | "We're Hiring" | Predictive Analytics | Streaming |
| Iris Data Set | Machine Learning | Real Time | Datafication | Facebook and Twitter |
| NoSQL | Mobile | Free Space!! | Internet of Things | Reuters-21578 |
| Visualization | Hadoop | Social Graph | @BigDataBorat quote | Wordcount Demo |
| Sentiment Analysis | NCDC GSOD | Business Intelligence | Someone who thinks R doesn't suck | "Data Is The New Oil" |

# Vinod Khosla: In The Next 10 Years, Data Science Will Do More For Medicine Than All Biological Sciences Combined

**FREDERIC LARDINOIS** ⌄
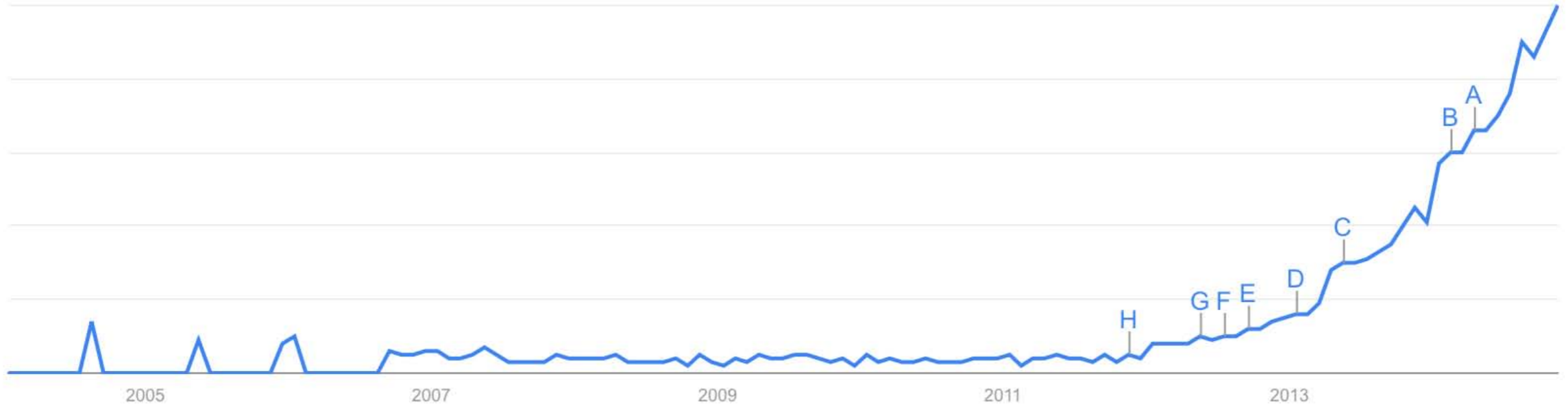
Wednesday, September 11th, 2013                    Comments

# Data Scientist: The Sexiest Job of the 21st Century

by Thomas H. Davenport and D.J. Patil

2011 May - Join Bitly as Data Scientist

2013 August - Lead Data Scientist at Bitly

2014 June - Lead Data Scientist at Spring

# bitly

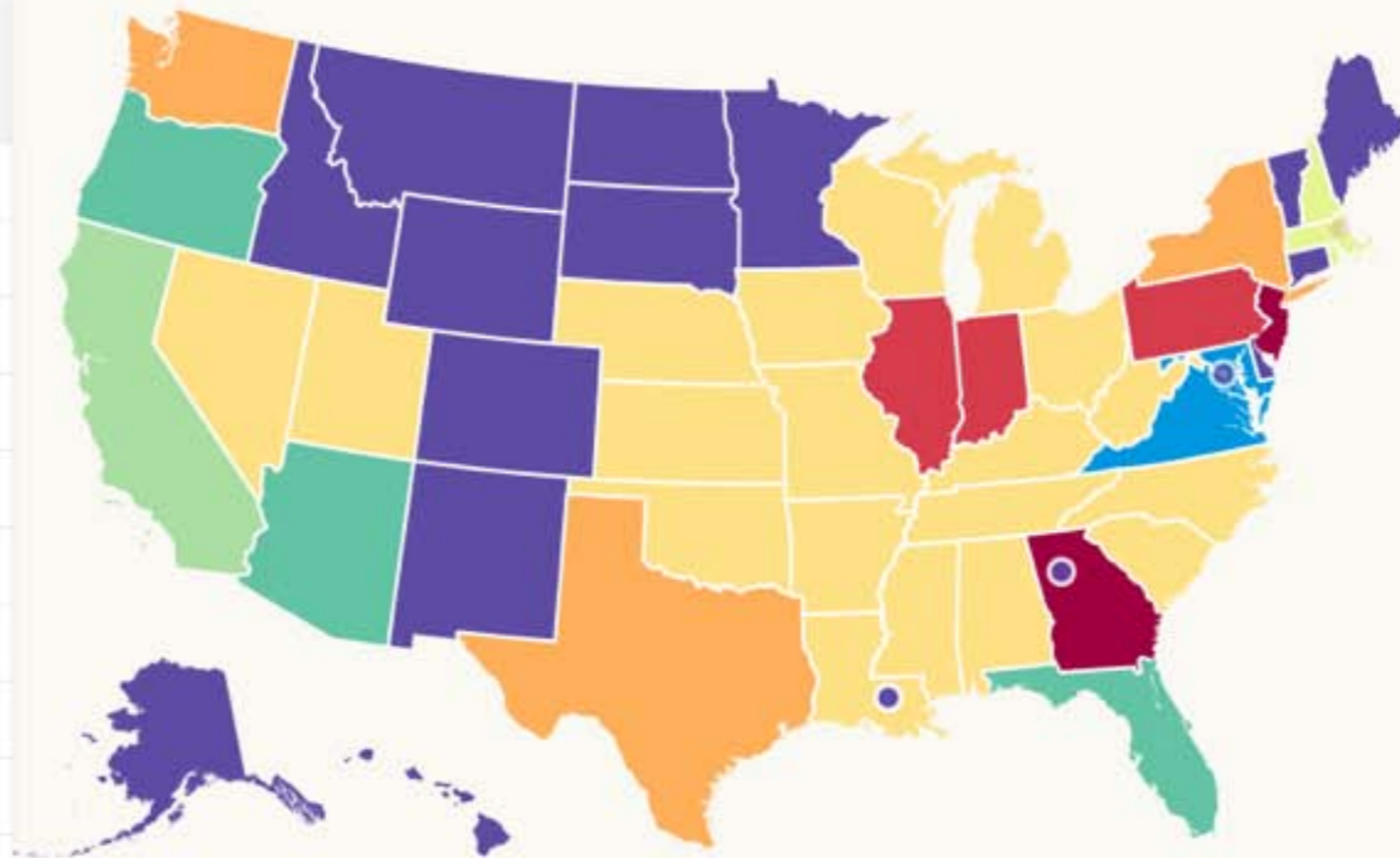## Real-Time Media Map

**1. Select a media type**

Newspapers ⌄

**2. Select your view**

⬜ Real-time traffic  ⬜ Winners by state  🔴 **Both**

### Media Properties Legend
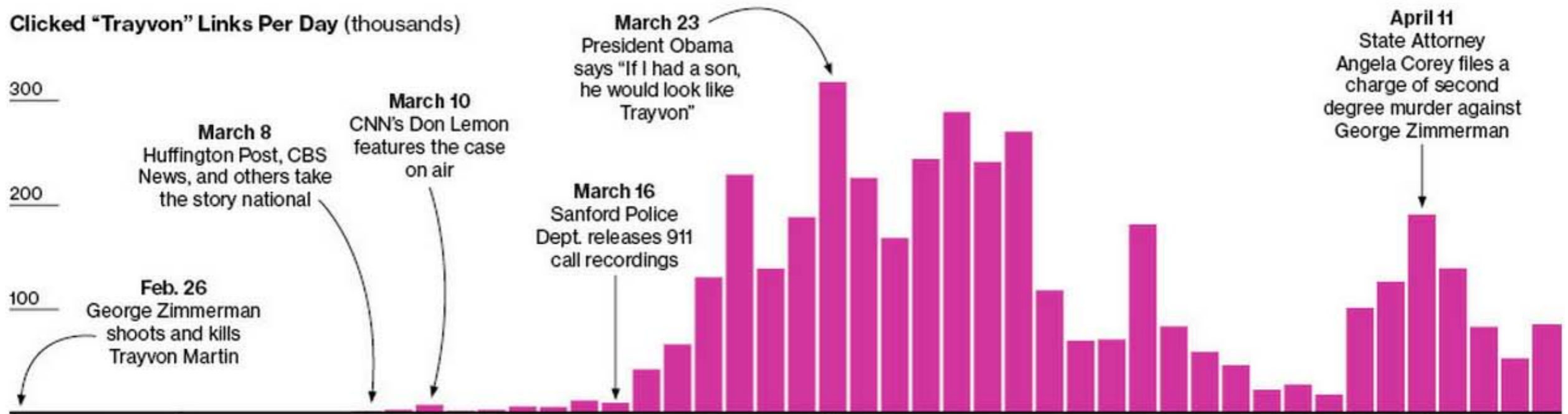(click to view by property)

- 🟣 The New York Times
- 🔵 Washington Post
- 🔴 Wall Street Journal
- 🔴 Chicago Tribune
- 🟢 Los Angeles Times
- 🟢 SF Gate
- 🟡 Boston Globe
- 🟡 USA Today
- 🟠 The Guardian
- 🔴 San Jose Mercury

Interested in what Bitly can do with your link click data? Please reach out to us at community@bitly.com
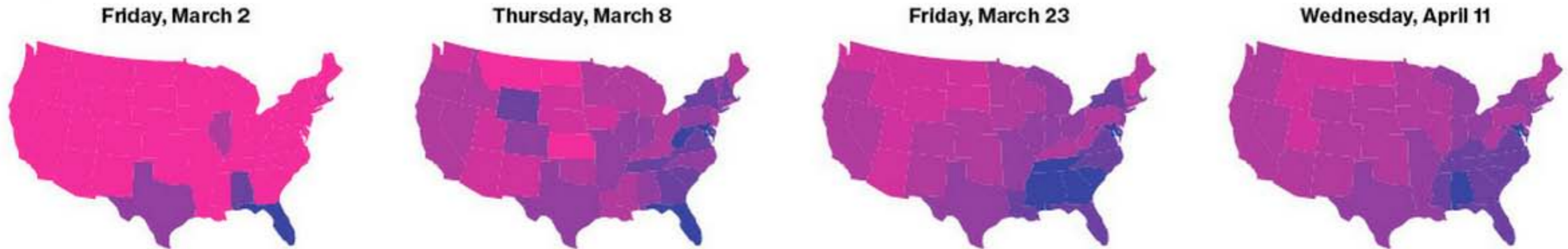
# Clicked "Trayvon" Links Per Day (thousands)

300

200

100

**March 8**
Huffington Post, CBS News, and others take the story national

**March 10**
CNN's Don Lemon features the case on air

**March 23**
President Obama says "If I had a son, he would look like Trayvon"

**March 16**
Sanford Police Dept. releases 911 call recordings

**April 11**
State Attorney Angela Corey files a charge of second degree murder against George Zimmerman
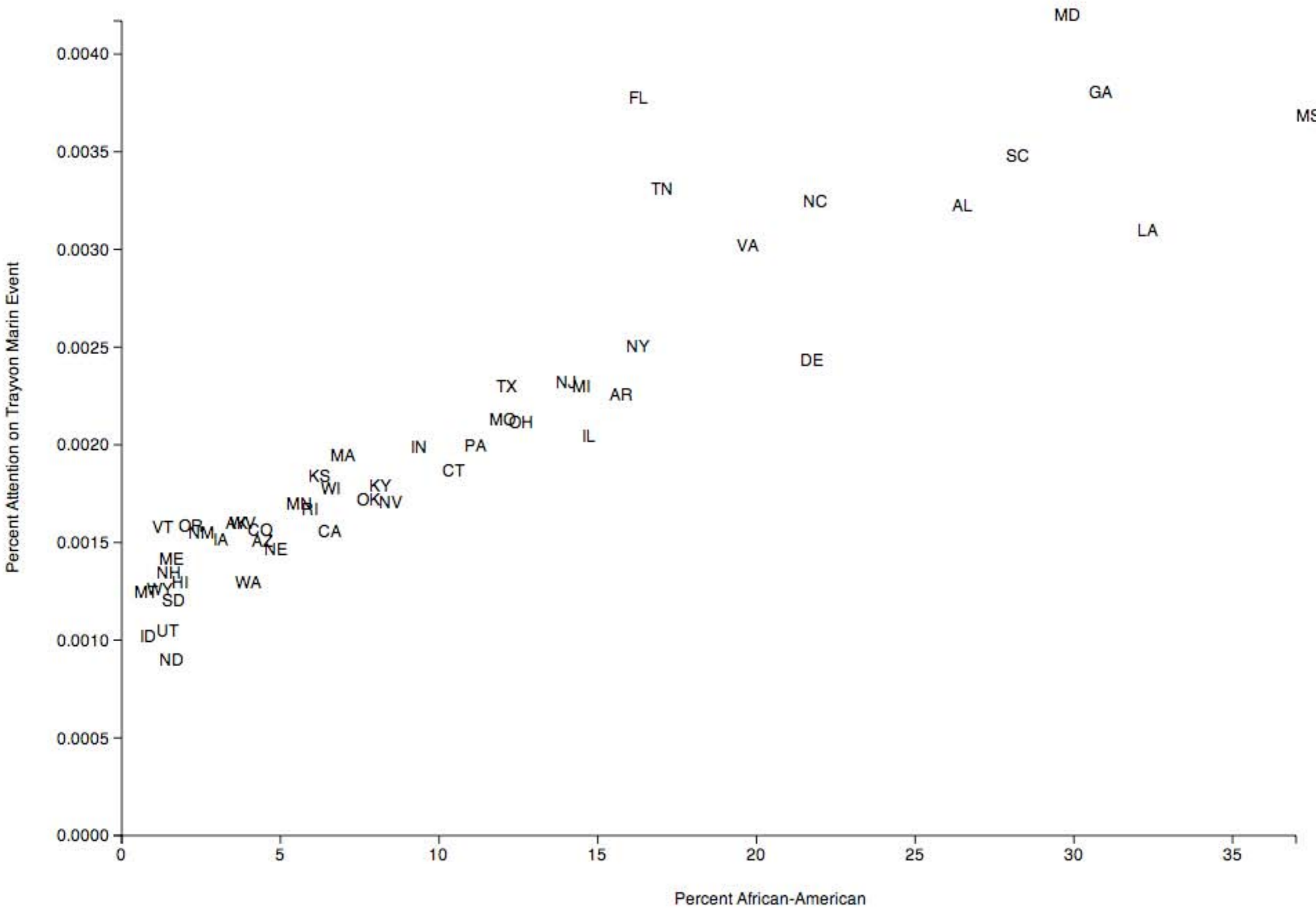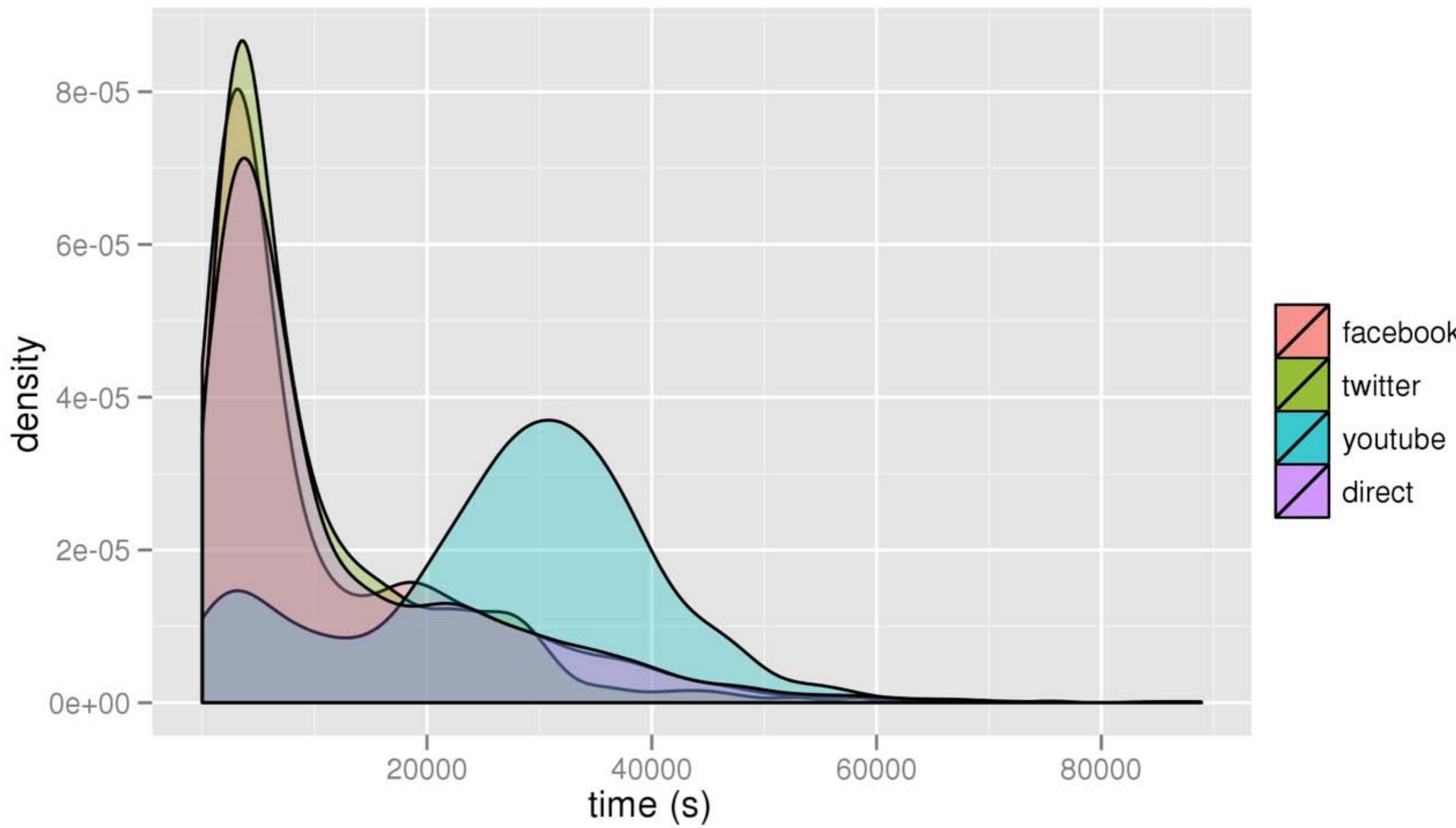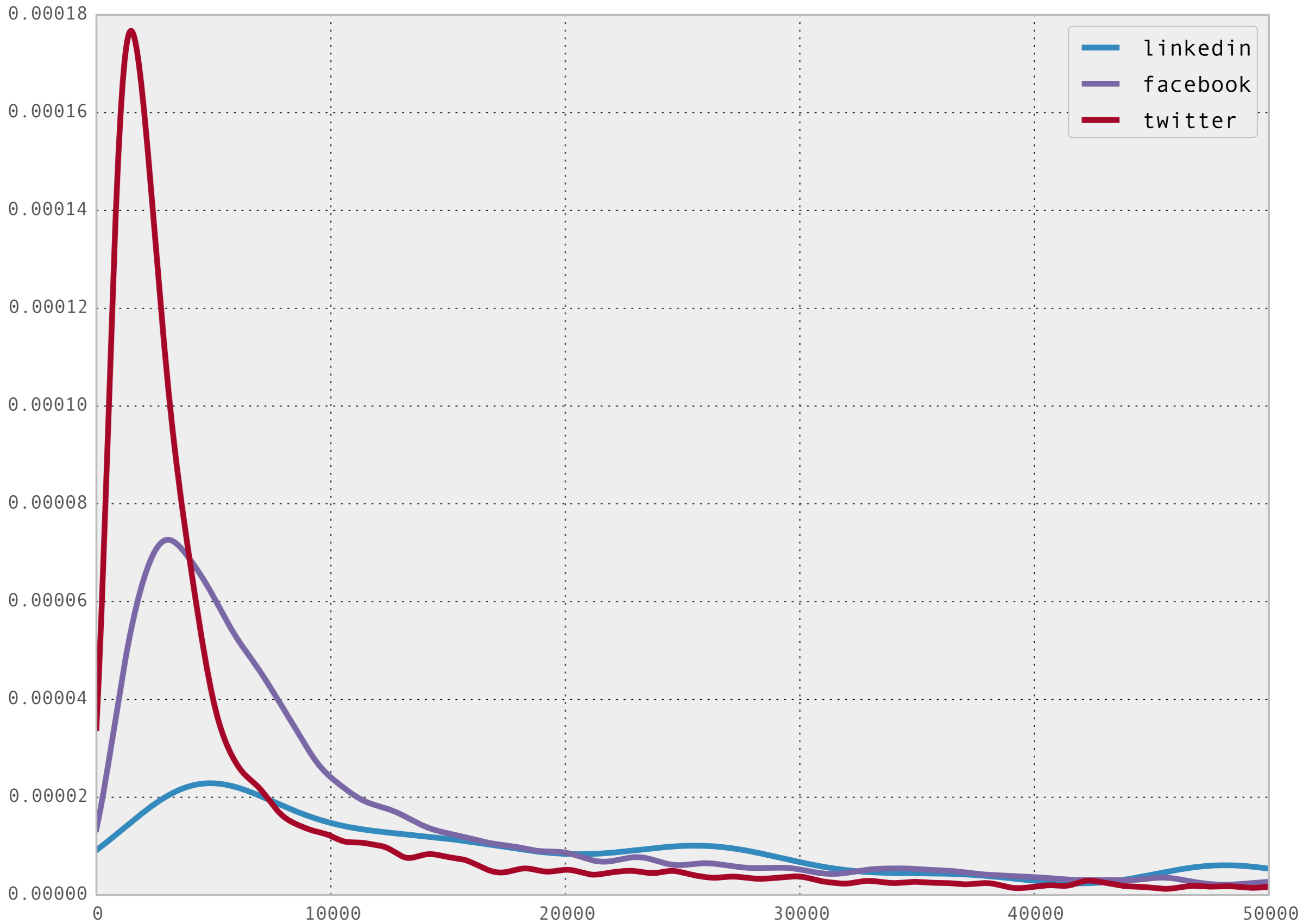
**Feb. 26**
George Zimmerman shoots and kills Trayvon Martin

## State Interest*

LESS    MORE



**Friday, March 2**

**Thursday, March 8**

**Friday, March 23**

**Wednesday, April 11**

*STATE INTEREST IS THE PORTION OF CLICKED LINKS THAT RELATE TO TRAYVON MARTIN COMPARED WITH OTHER STATES IN ONE DAY;
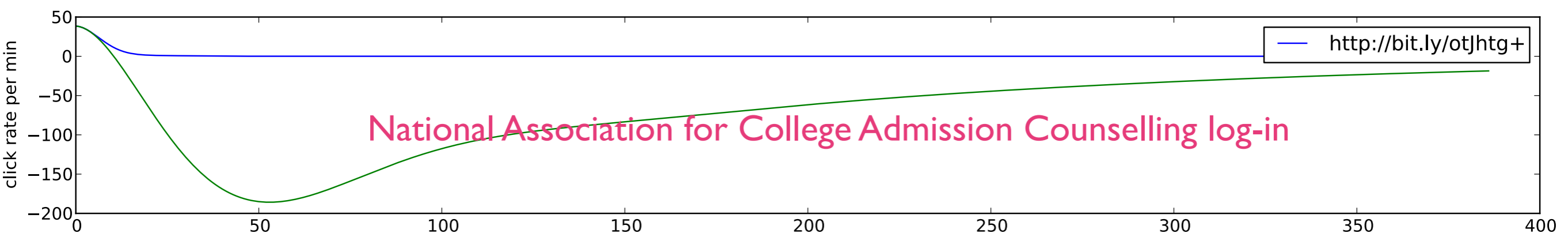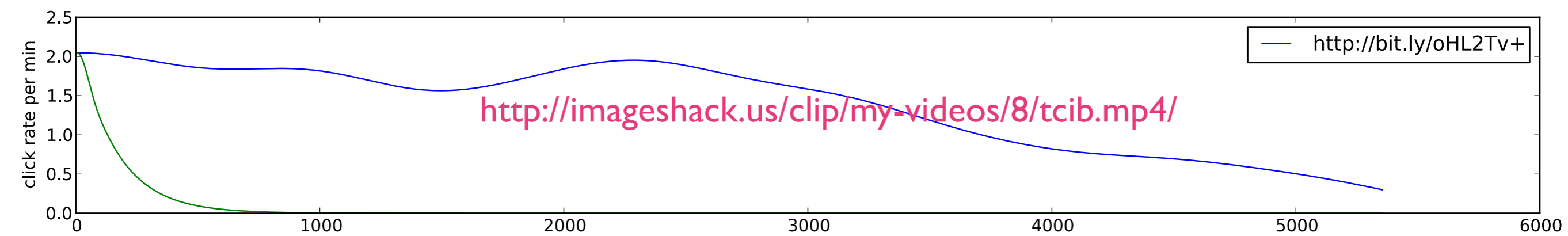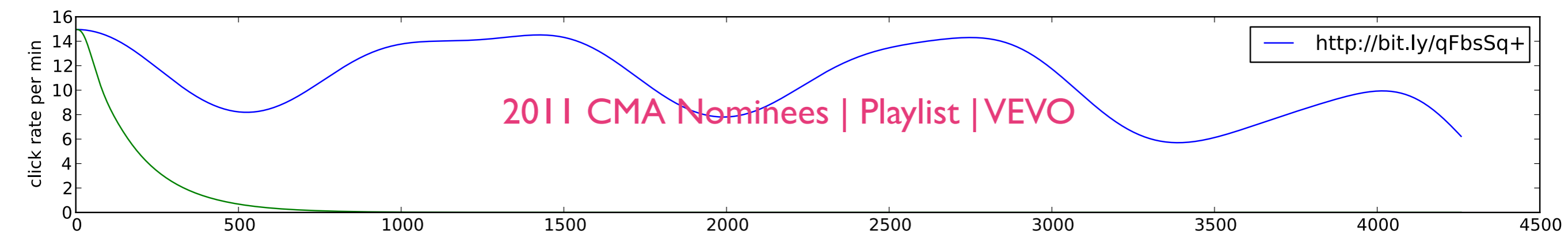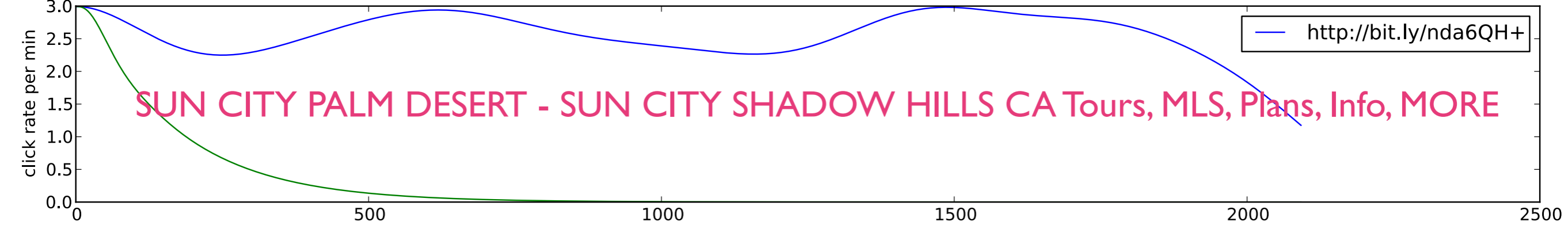GRAPHIC BY BLOOMBERG BUSINESSWEEK; DATA: BITLY

Free Game For Kids Registration - Pittsburgh Penguins - Fan Zone (2756 clicks)

Video: Carra's top five transfers - Liverpool FC (915 clicks)

Clip of the Week: Toews Nails Tot | NBC Chicago (179 clicks)

Allegro.pl nie działa (919 clicks)

DallasCowboys.com - Official Site of the Dallas Cowboys (683 clicks)

Mao's Room (2339 clicks)

Manchester United Official Web Site - Ashley Young was long term United target (526 clicks)

Shocking! Lady Gaga Poses Sans Makeup for Harper's Bazaar Cover - UsMagazine.com (12288 clicks)

way - Runway TV Collections Fashion Magazine - Nick Carter: From the Backstreet to Taking Off (3662 c

The GQ&A: Drive Director Nicolas Winding Refn (601 clicks)

SUN CITY PALM DESERT - SUN CITY SHADOW HILLS CA Tours, MLS, Plans, Info, MORE

http://bit.ly/nda6QH+

2011 CMA Nominees | Playlist | VEVO

http://bit.ly/qFbsSq+

http://imageshack.us/clip/my-videos/8/tcib.mp4/

http://bit.ly/oHL2Tv+

National Association for College Admission Counselling log-in

http://bit.ly/otJhtg+

TOP SALE Viagra from USD 0.90 per pill, Cialis from USD 1.75 per pill

http://bit.ly/ogBZAF+

sample index

SPRING

EST 2013 NYC

# Big-Data
## to
# No-Data

# Three Areas of Data Science
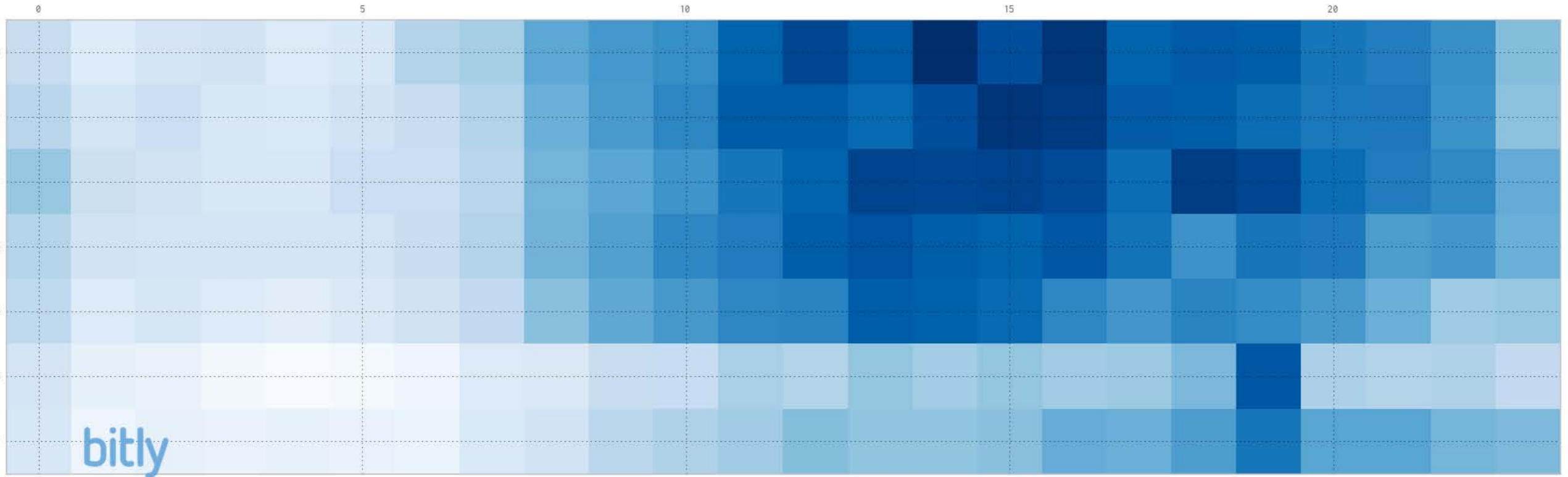
# Data Engineering

**Brian David Eoff** @bde · Jun 11

Not so much a data scientist as a data plumber. Conceiving of the flow, pipes and stores necessary to accomplish our goals.
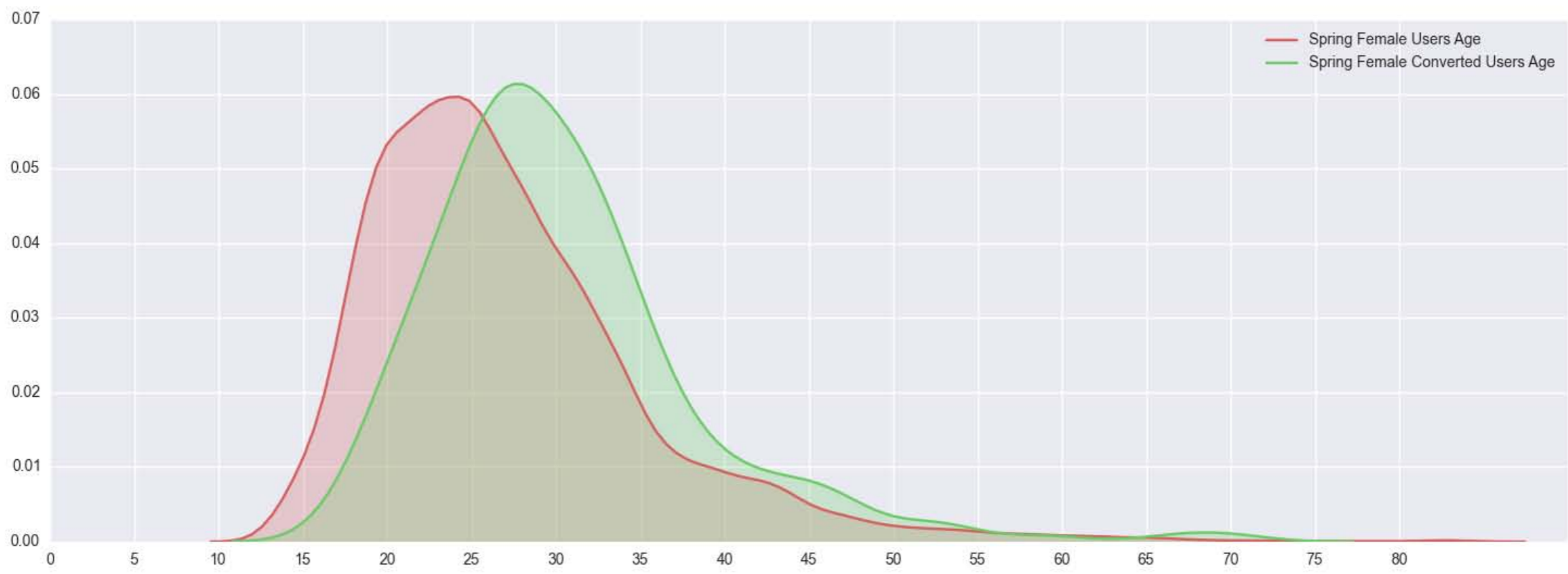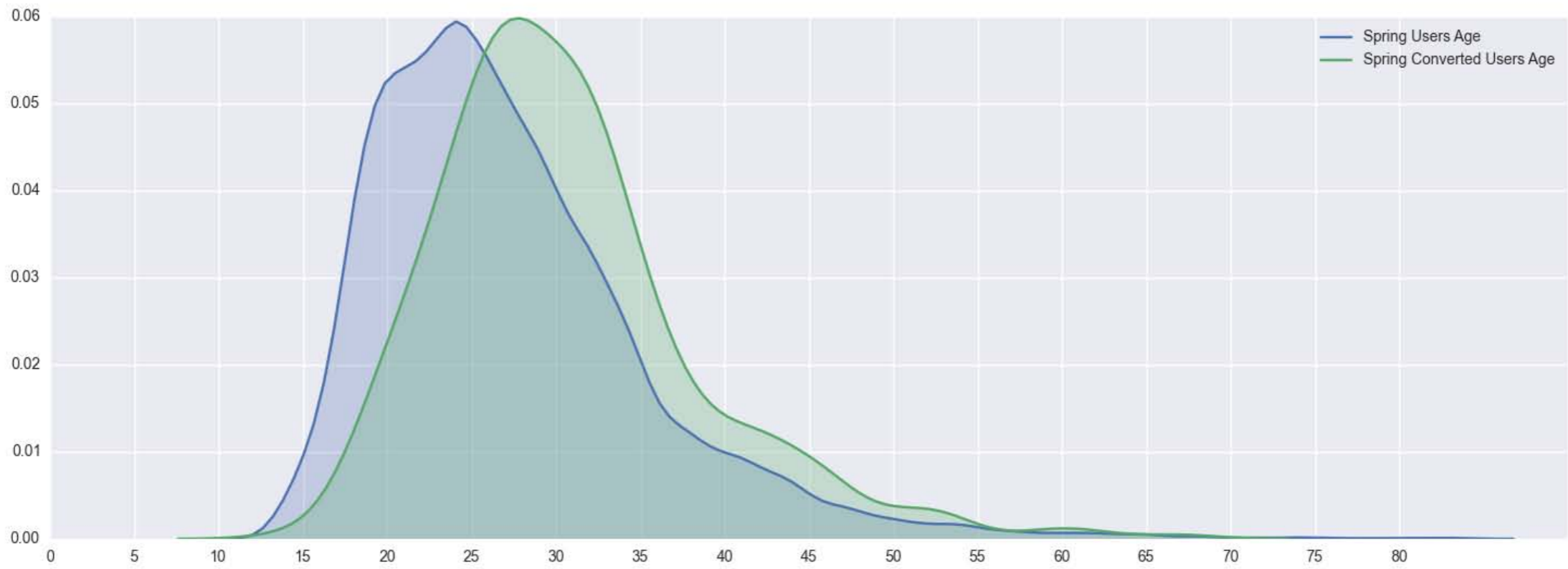
# More than…

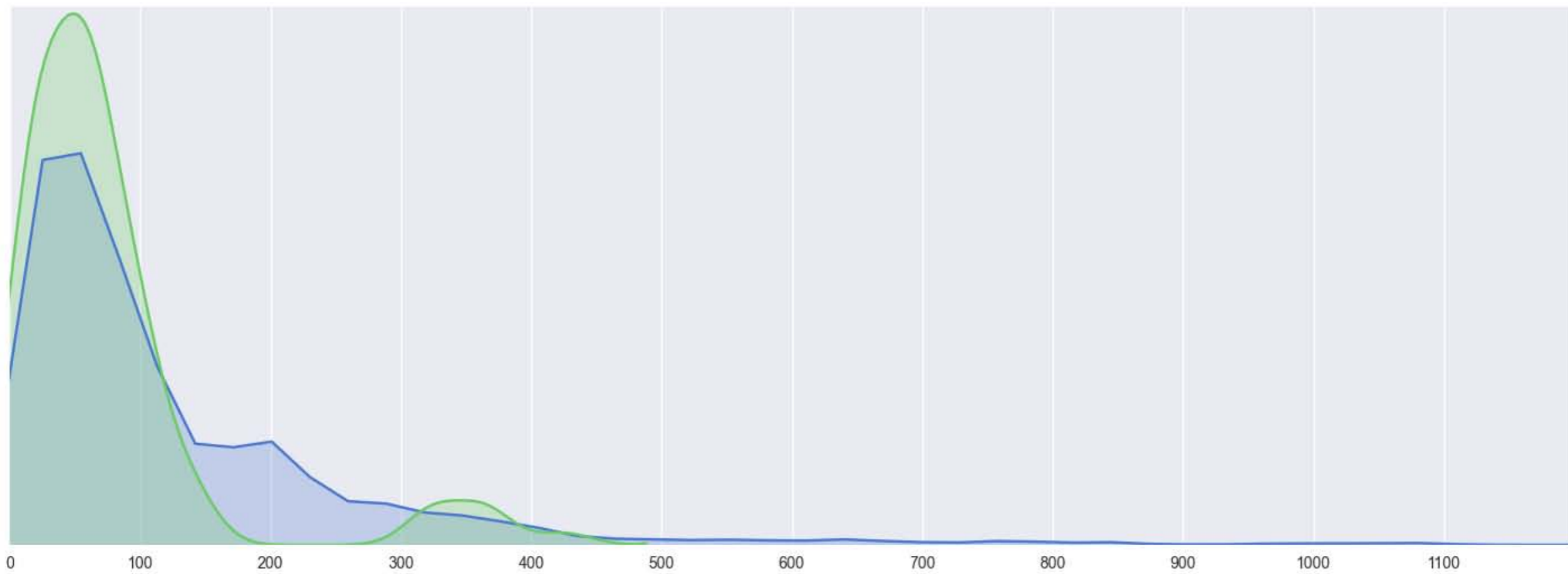# Everything in UTC

# Data Analysis

# Reproducibility, Automation, Version Control

# Data Modeling

# Discovery

# Collaborative Filtering for Implicit Feedback Datasets

Yifan Hu
AT&T Labs – Research
Florham Park, NJ 07932

Yehuda Koren*
Yahoo! Research
Haifa 31905, Israel

Chris Volinsky
AT&T Labs – Research
Florham Park, NJ 07932

## Abstract

*A common task of recommender systems is to improve customer experience through personalized recommendations based on prior* implicit feedback. *These systems passively track different sorts of user behavior, such as purchase history, watching habits and browsing activity, in order to model user preferences. Unlike the much more extensively researched* explicit feedback, *we do not have any direct input from the users regarding their preferences. In particular, we lack substantial evidence on which products consumer dislike. In this work we identify unique properties of implicit feedback datasets. We propose treating the data as indication of positive and negative preference associated with vastly varying confidence levels. This leads to a factor model which is especially tailored for implicit feedback recommenders. We also suggest a scalable optimization procedure, which scales linearly with the data size. The algorithm is used successfully within a recommender system for television shows. It compares favorably with well tuned implementations of other known methods. In addition, we offer a novel way to give explanations to recommendations given by this factor model.*

tent based approach creates a profile for each user or product to characterize its nature. As an example, a movie profile could include attributes regarding its genre, the participating actors, its box office popularity, etc. User profiles might include demographic information or answers to a suitable questionnaire. The resulting profiles allow programs to associate users with matching products. However, content based strategies require gathering external information that might not be available or easy to collect.

An alternative strategy, our focus in this work, relies only on past user behavior without requiring the creation of explicit profiles. This approach is known as *Collaborative Filtering* (CF), a term coined by the developers of the first recommender system - Tapestry [8]. CF analyzes relationships between users and interdependencies among products, in order to identify new user-item associations. For example, some CF systems identify pairs of items that tend to be rated similarly or like-minded users with similar history of rating or purchasing to deduce unknown relationships between users and items. The only required information is the past behavior of users, which might be their previous transactions or the way they rate products. A major appeal of CF is that it is domain free, yet it can address aspects of the data that are often elusive and very difficult to profile using con

# Popularity

# Tools

# Caverlee's Rule

# How to become a data scientist?

# CONFESSION

# I never took a data science class.

# What do I look for when I hire?

# Questions?

Descriptive, Prescriptive, Predictive

Data Engineering, Data Analysis, Machine Learning/Model Building/Algorithms

ETL, or never underestimate the benefit of a good plumber.

AB Test. Multi-Armed Bandit Testing. Optimization. Retention

Bitly: 200 Million Events per Day
Spring: 0

Data Science (Inductive Reasoning), obtaining data, cleaning munging data, exploring data, modeling data interpreting data

Data Engineering

ETL

Redshift

Event Streams
Kafka, Kinesis, NSQ

Tools, Process, Recruiting,
Communication, Ethics

Test- Driven Data Science

Learn from software development.

Version Control, Automated Testing

Analysis

Price Points

Conversion

Metric Aggregation

etc.

Matrix Factorization

Recommendation

Popularity