

How the growth of R helps data-driven organizations succeed

David Smith

Chief Community Officer

7th China R User Conference, May 2014



Agenda

- Introduction
- History of R
- Growth of R
- Applications of R
- The R Ecosystem
- Conclusion

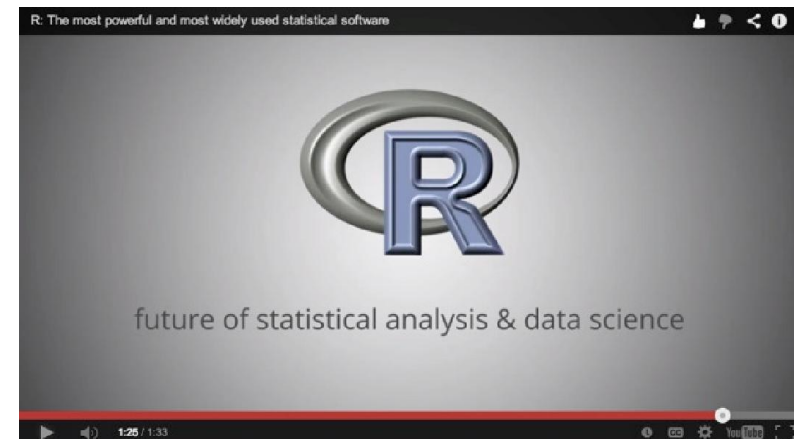
Slides will be posted to:

blog.revolutionanalytics.com



What is R?

- Most widely used data analysis software
 - Used by 2M+ data scientists, statisticians and analysts
- Most powerful statistical programming language
 - Flexible, extensible and comprehensive for productivity
- Create beautiful and unique data visualizations
 - As seen in New York Times, Twitter and Flowing Data
- Thriving open-source community
 - Leading edge of analytics research
- Fills the talent gap
 - New graduates prefer R



A brief history of R

- **1993:** Research project in Auckland, NZ
 - Ross Ihaka and Robert Gentleman
- **1995:** Released as open-source software
 - Generally compatible with the “S” language
- **1997:** R core group formed
- **2000:** R 1.0.0 released
- **2004:** R 2.0.0 released, first international user conference in Vienna
- **2009:** New York Times article on R
- **2013:** R 3.0.0 released



Search Technology

Go

Inside Technology

[Internet](#)[Start-Ups](#)[Business Computing](#)[Companies](#)Bits
Blog »

Personal Tech »

[Cellphones](#), [Cameras](#), [Computers](#) and more

Data Analysts Captivated by R's Power



Stuart Issett for The New York Times

R first appeared in 1996, when the statistics professors Robert Gentleman, left, and Ross Ihaka released the code as a free software package.

By [ASHLEE VANCE](#)

Published: January 6, 2009

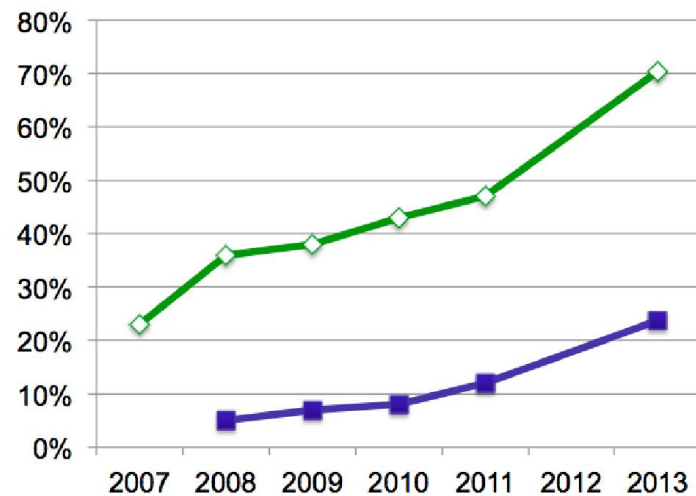
☒ RECOMMEND

New York Times:
Data Analysts Captivated by R's
Power
[6 Jan 2009](#)

R's popularity is growing rapidly

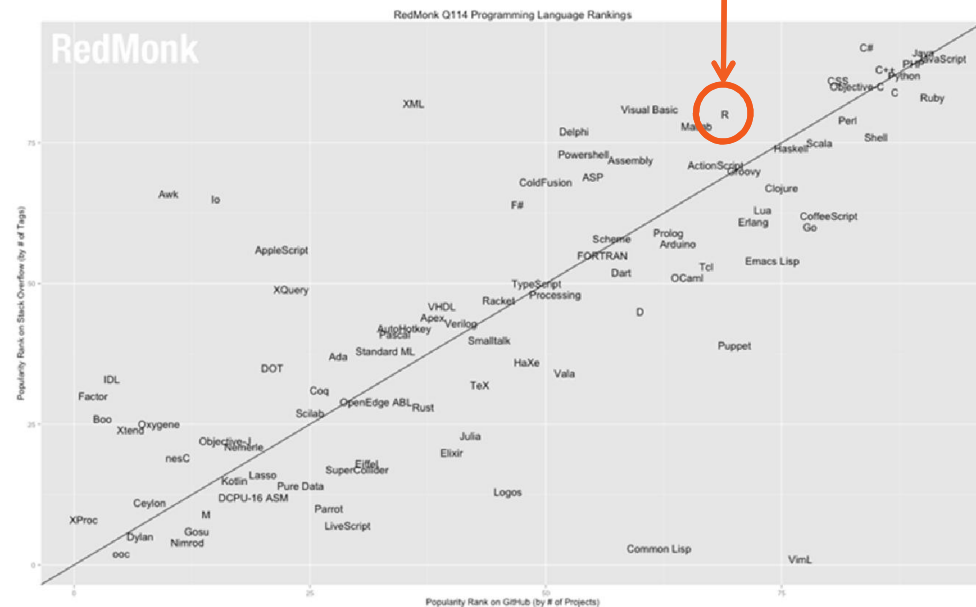
R Usage Growth

Rexer Data Miner Survey, 2007-2013



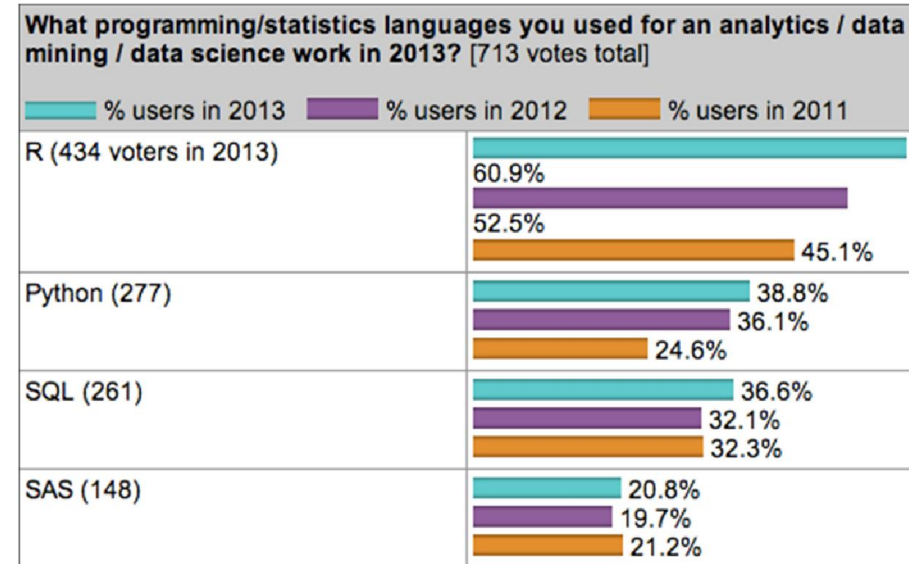
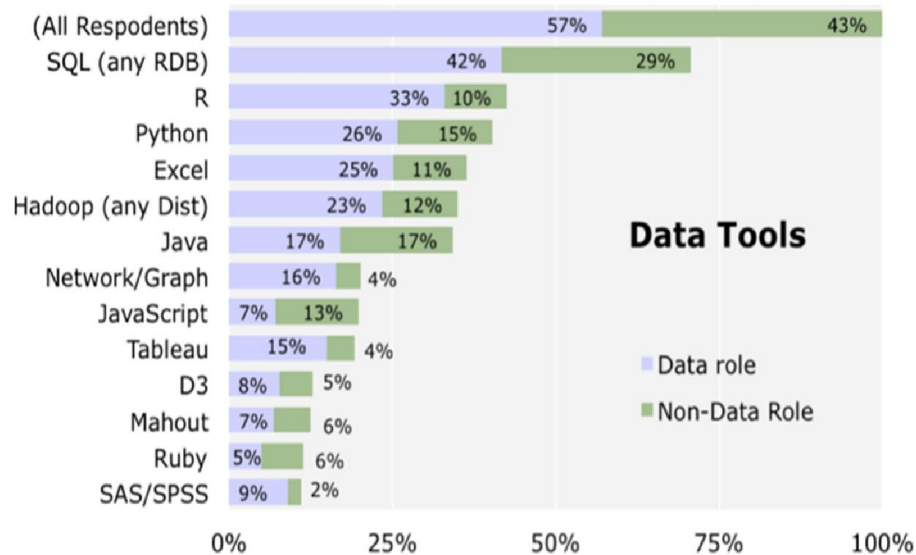
- [Rexer Data Miner Survey, 2013](#)

#15: R



- [RedMonk Programming Language Rankings, 2013](#)

R is used more than other data science tools



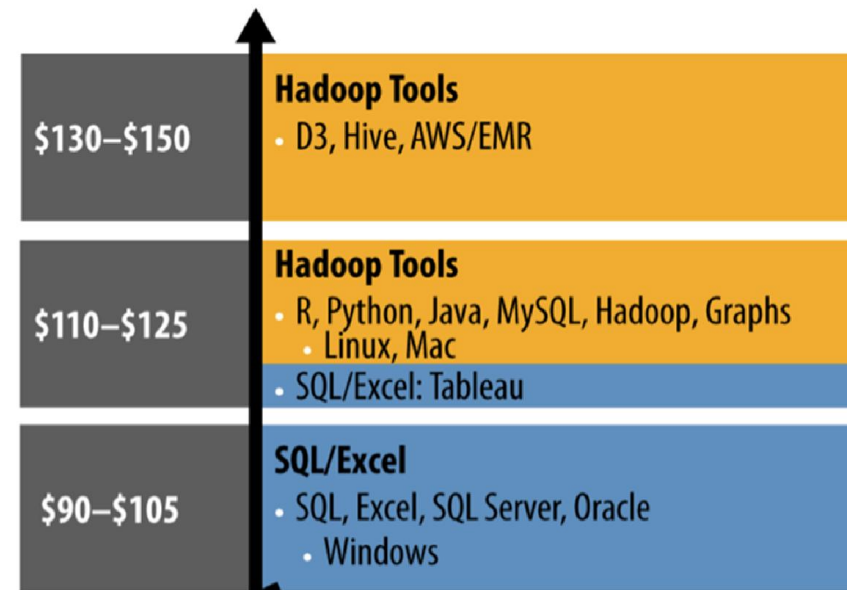
- [O'Reilly Strata 2013 Data Science Salary Survey](#)

- [KDNuggets Poll: Top Languages for analytics, data mining, data science](#)

R is among the highest-paid IT skills in the US

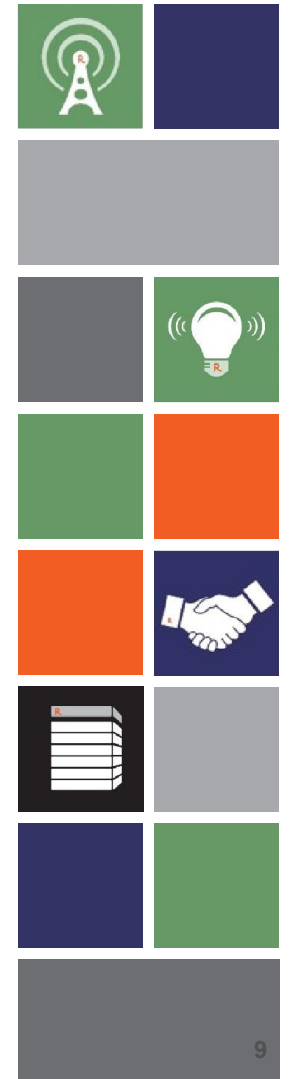
| AVERAGE SALARY FOR High Paying Skills and Experience | | |
|---|------------|--------------|
| SKILL | 2013 | YR/YR CHANGE |
| R | \$ 115,531 | n/a |
| NoSQL | \$ 114,796 | 1.6% |
| MapReduce | \$ 114,396 | n/a |
| PMBok | \$ 112,382 | 1.3% |
| Cassandra | \$ 112,382 | n/a |
| Omnigraffle | \$ 111,039 | 0.3% |
| Pig | \$ 109,561 | n/a |
| SOA (Service Oriented Architecture) | \$ 108,997 | -0.5% |
| Hadoop | \$ 108,669 | -5.6% |
| Mongo DB | \$ 107,825 | -0.4% |

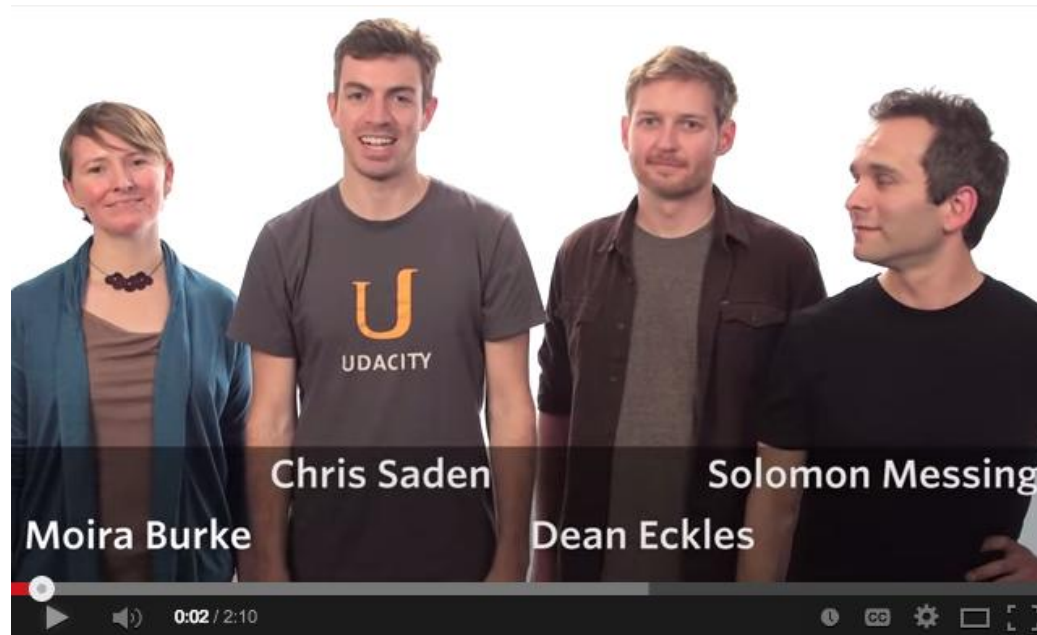
- [Dice Tech Salary Survey, January 2014](#)



- [O'Reilly Strata 2013 Data Science Salary Survey](#)

Applications of R





- [Exploratory Data Analysis](#)
- [Experimental Analysis](#)

“Generally, we use R to move fast when we get a new data set. With R, we don’ t need to develop custom tools or write a bunch of code. Instead, we can just go about cleaning and exploring the data.”

— Solomon Messing, data scientist at Facebook



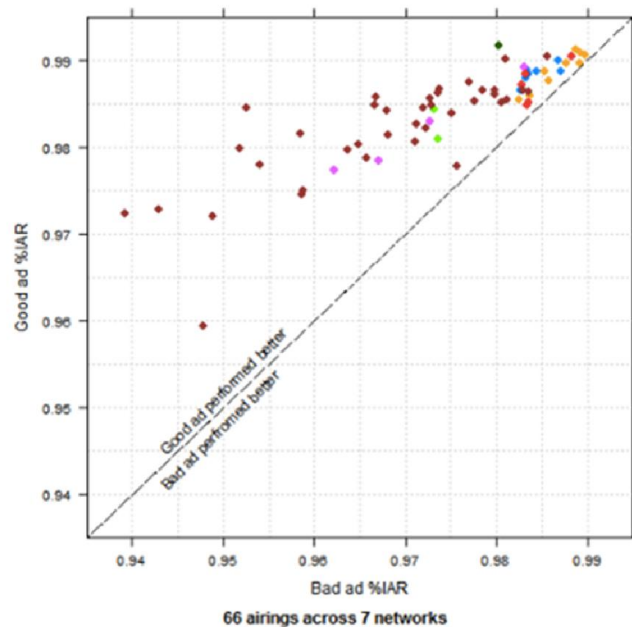
- [Big-Data Visualization](#)

“It resonated with many people. It’s not just a pretty picture, it’s a reaffirmation of the impact we have in connecting people, even across oceans and borders.” — Paul Butler, data scientist, Facebook





All live experiments, December - February



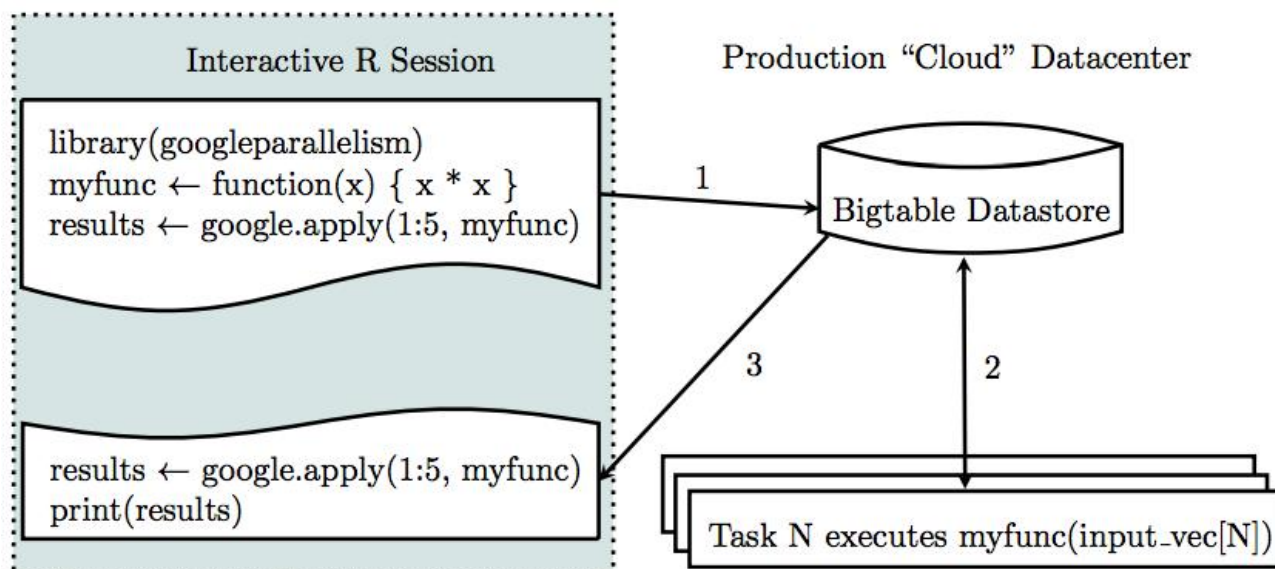
“R is really important to the point that it’s hard to overvalue it.”

— Daryl Pregibon Head of Statistics,
Google Advertising Effectiveness

“The great beauty of R is that you can modify it to do all sorts of things.”

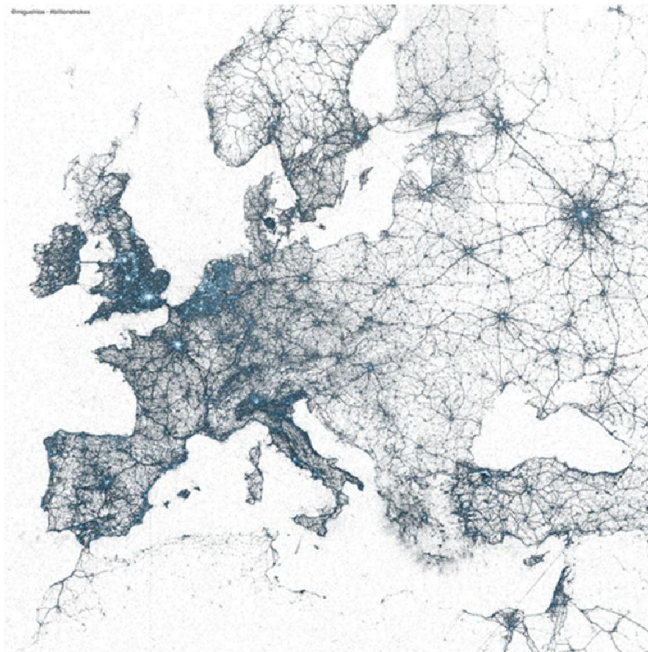
— Hal Varian
Chief Economist,
Google

- [Economic forecasting](#)



- [Big-data statistical modeling](#) (Large-Scale Parallel Statistical Forecasting Computations in R, JSM 2011)

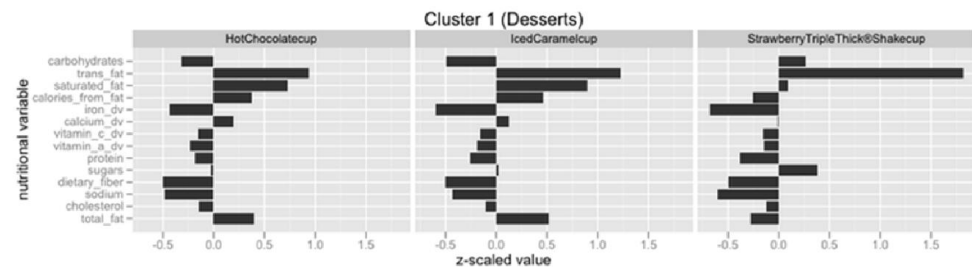
“We demonstrate the utility of massively parallel computational infrastructure for statistical computing using the MapReduce paradigm for R. This framework allows users to write computations in a high-level language that are then broken up and distributed to worker tasks in Google datacenters.”



- [Data Visualization](#)



“A common pattern for me is that I’ll code a MapReduce job in Scala, do some simple command-line munging on the results, pass the data into Python or R for further analysis, pull from a database to grab some extra fields, and so on, often integrating what I find into some machine learning models in the end” — Ed Chen, Data Scientist, Twitter



- [Semantic clustering](#)



. @cheerjoeyniz 13 Apr
Ofcourse I of all people get food poisoning the night before my last comp.
#hatelife



Foodborne Chicago
@foodbornechi

Follow

@cheerjoeyniz Sorry to hear you were sick. We can help you by clicking on this link to file a report:
foodborne.smartchicagoapps.org/32310525811084...

2:14 PM - 16 Apr 2013

Foodborne Chicago

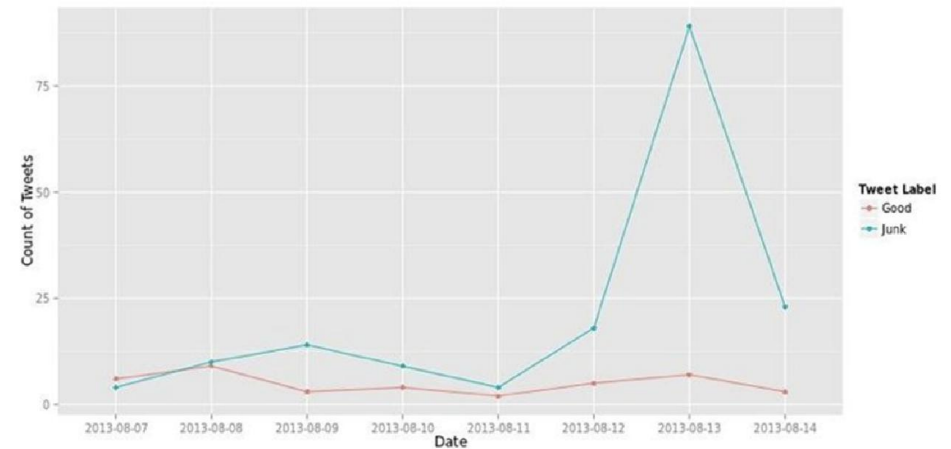
If you think you have food poisoning in Chicago, please complete this form. The info will be sent through the Chicago 311 service to the Chicago Department of Public Health so they can take any...



Foodborne Chicago @foodbornechi



Foodborne Chicago Tweets by Time



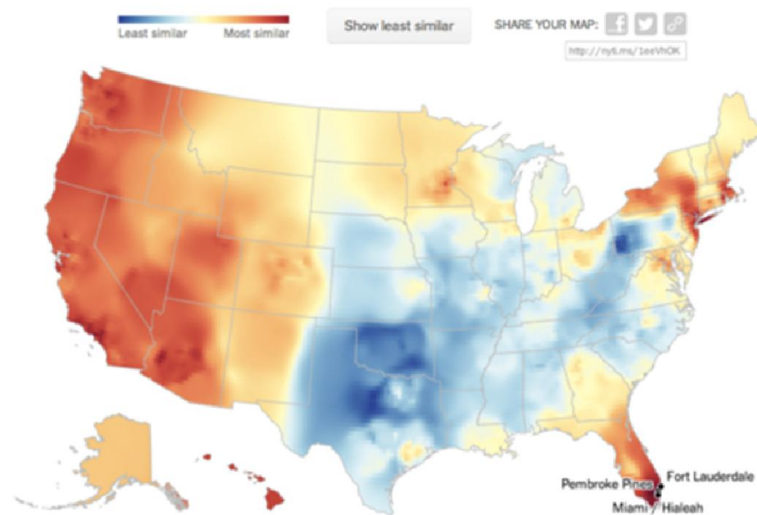
- [Food poisoning monitor](#)



REVOLUTION
ANALYTICS

The New York Times

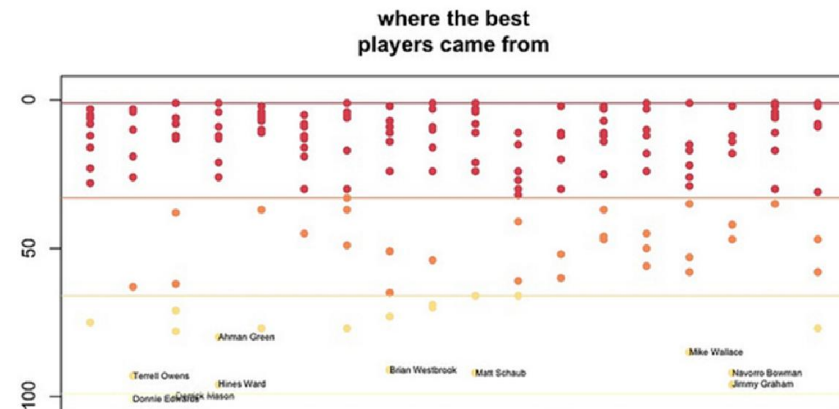
Interactive Features



- [Election Forecast](#)
- [Dialect Quiz](#)



Data Journalism



- [NFL Draft Picks](#)
- [Wealth distribution in USA](#)

The New York Times

Data Visualization

What Happens After the I.P.O.?

Since 1980, there have been about 2,400 technology, Internet and telecom I.P.O.'s. On the first day of trading, the average stock rose 32 percent above its offer price.

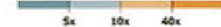
But in the three years after that, most companies had negative returns, according to statistics compiled by Jay Ritter, a professor of finance at the University of Florida. Companies with higher values compared with their revenue before the I.P.O. have fared especially poorly.

CHART KEY

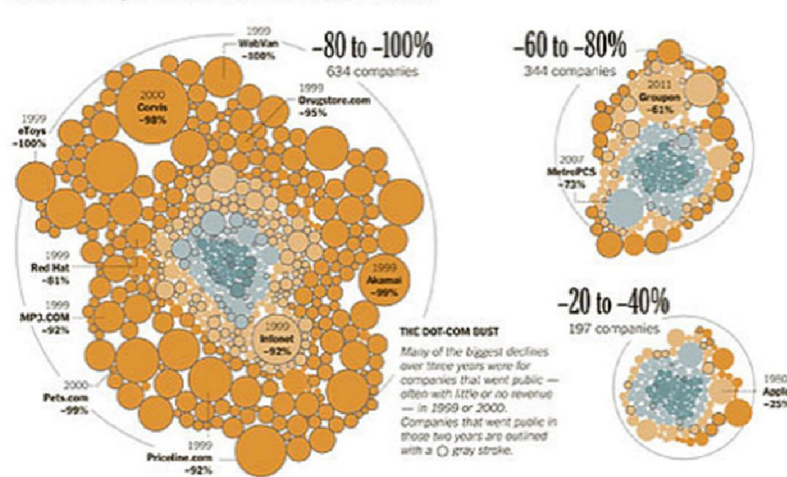
Circles are sized by value at the end of the first trading day, in today's dollars



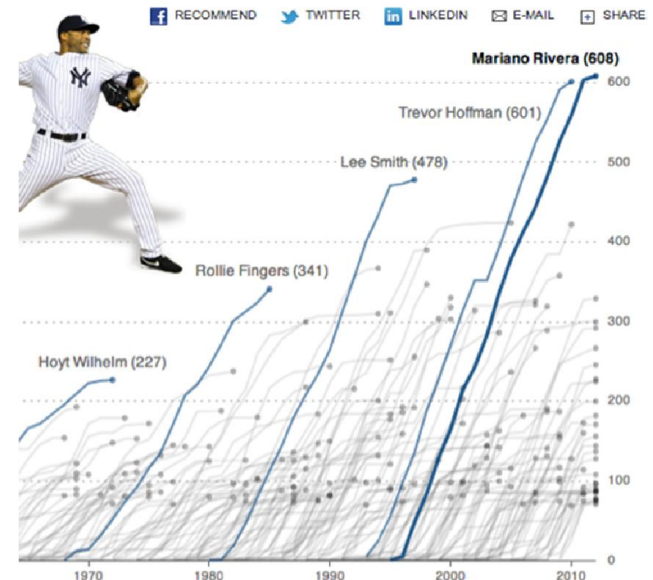
Colors show the ratio of the company's value to its revenue in the 12 months before the I.P.O.

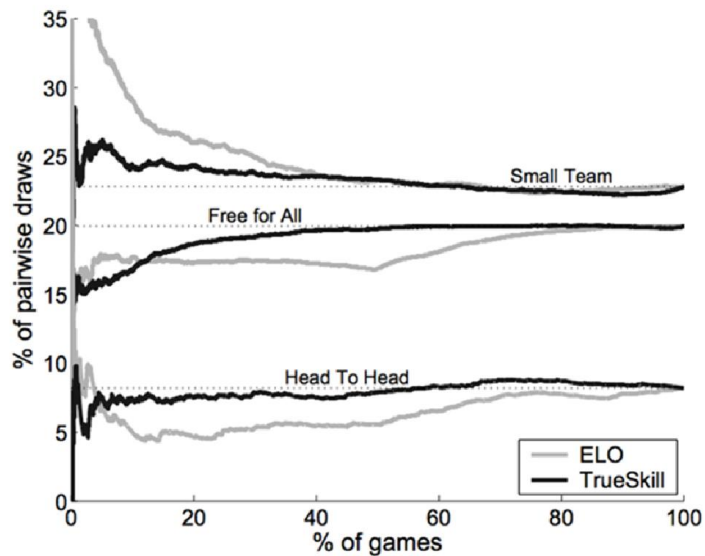


Return three years after the I.P.O.: The decliners ...



- [Facebook IPO](#)
- [Baseball legends](#)

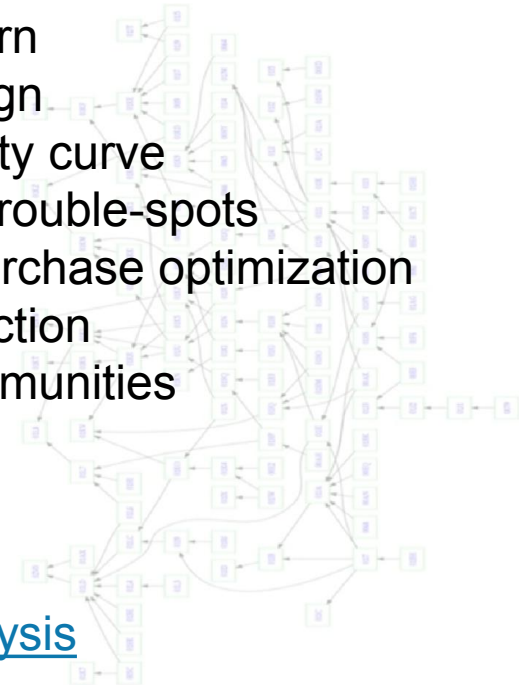




- [Multiplayer Matchmaking](#)

scientific revenue

- Player Churn
- Game design
 - Difficulty curve
 - Level trouble-spots
- In-game purchase optimization
- Fraud detection
- Player communities

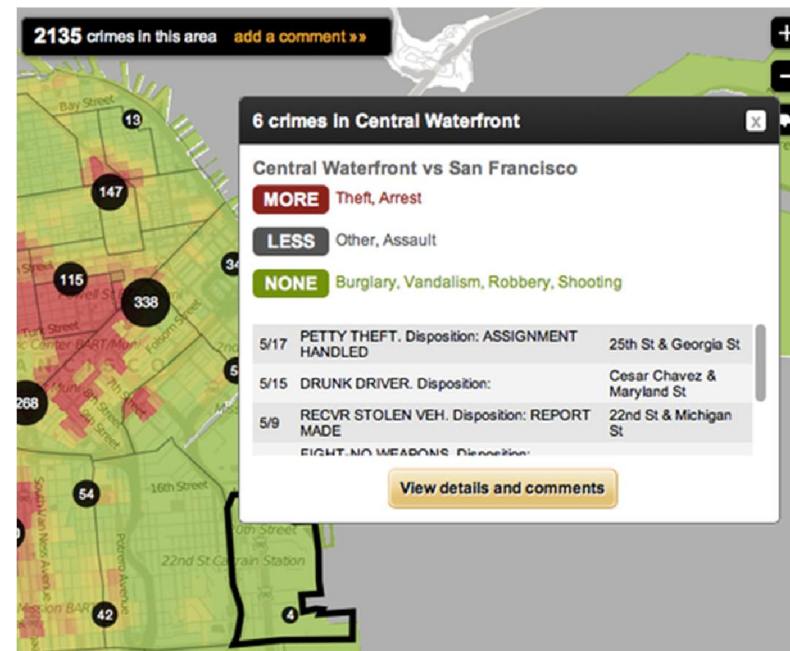


- [Game Analysis](#)



“The core innovation that Zillow offers are its advanced statistical predictive products, including the Zestimate®, the Rent Zestimate and the ZHVI® family of real estate indexes. By using R in production as well as research, Zillow maximizes flexibility and minimizes the latency in rolling out updates and new products.”

- [Statistical forecasting](#)



- [Crime mapping](#)



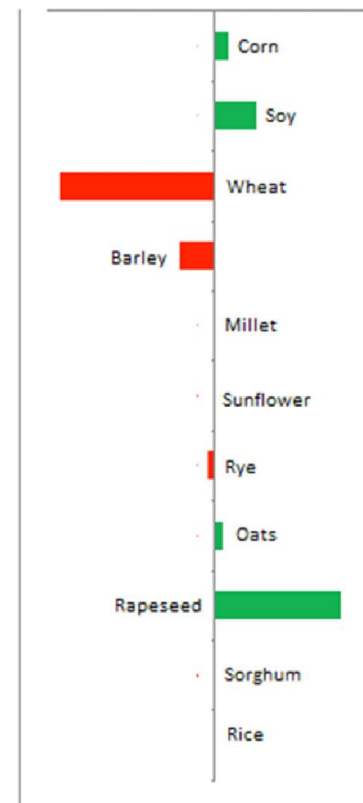
JOHN DEERE

Statistical Analysis:

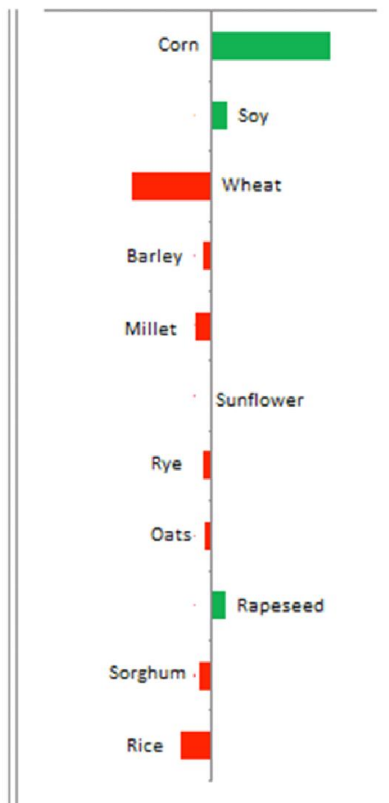
- Short Term Demand Forecasting
- Crop Forecasting
- Long Term Demand Forecasting
- Maintenance and Reliability
- Production Scheduling
- Data Coordination

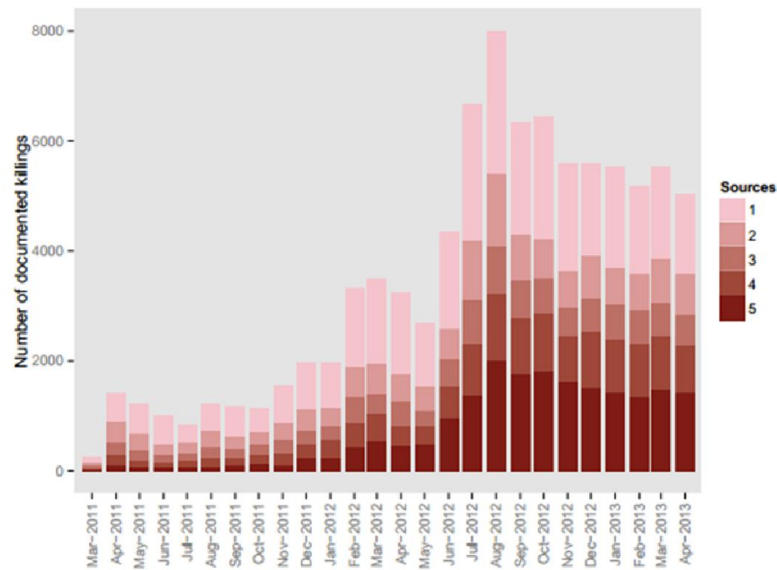


Canada



China

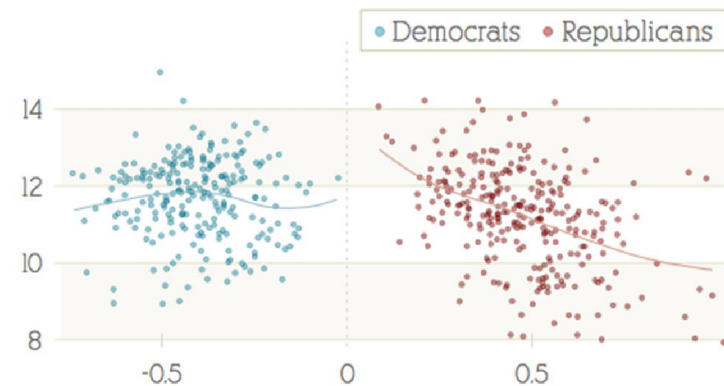




- [Casualty estimation in Warzones](#)

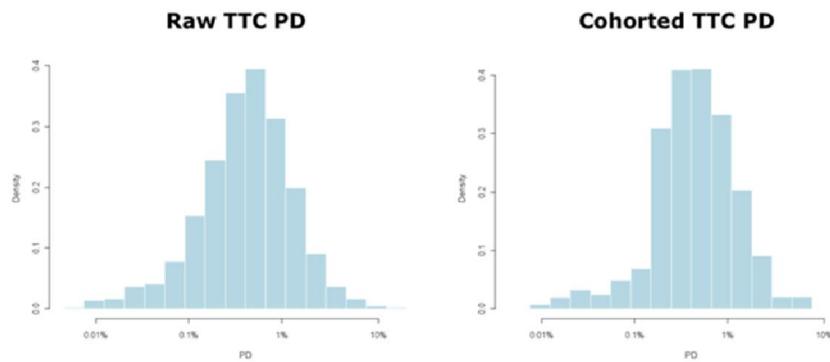


Ideology and Grade Level of Congressional Record Speeches, Current Members

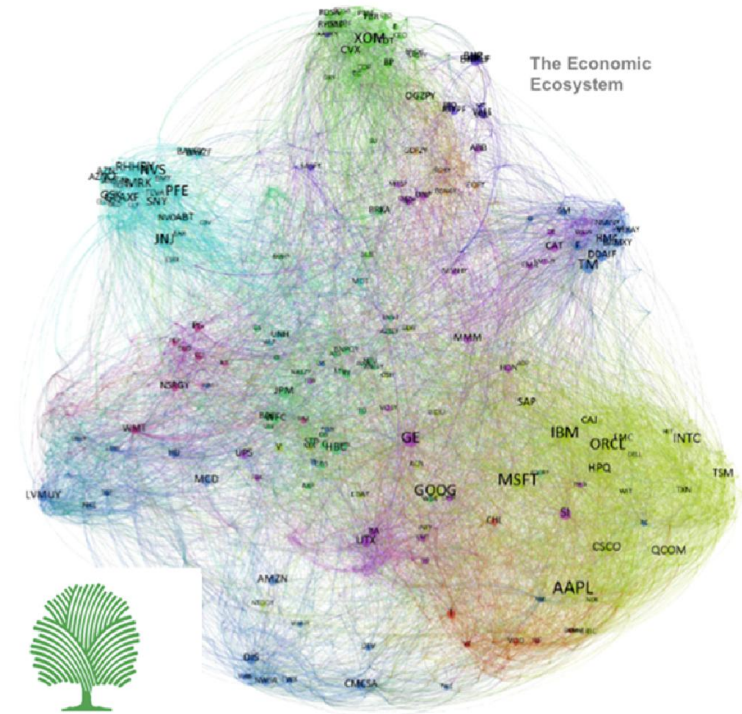


Liberal-Conservative Voting Score, 112th Congress
-1 (most liberal) to 1 (most conservative)

- [Political Analysis](#)



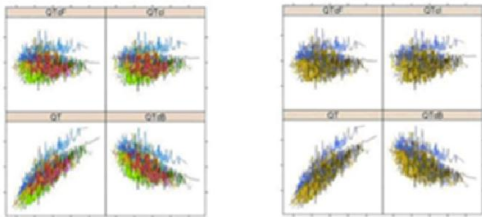
- [Credit Risk Analysis](#)



- [Financial Networks](#)



“R use at the FDA is completely acceptable and has not caused any problems.” — Dr Jae Brodsky, Office of Biostatistics, Food and Drug Administration



Example from AdminFDA report. Left: a plot made with the default rainbow color palette. This figure is not 508 compliant because some information is lost to colorblind people. Right: the red-green colorblind simulation. It is not possible to determine which lines correspond to individual subjects when the color palette is reduced.

Regulatory Drug Approvals

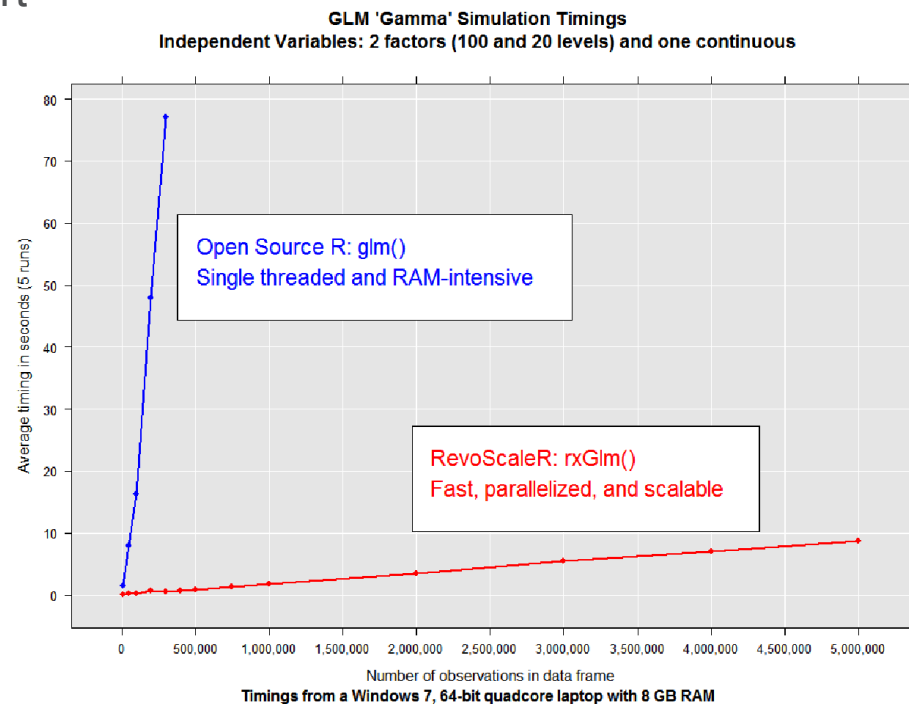
- Reproducible research
- Accurate, reliable and consistent statistical analysis
- Internal reporting (Section 508 compliance)

Software Vendors and Service Providers





- Open-Source R Technical Support
- Open Source development
 - RHadoop, ParallelR
- Community Support
 - User Group Sponsorship
 - Meetups, Events
 - Revolutions Blog
- Revolution R Enterprise
 - Big Data (ScaleR package)
 - Integration (Web Services API)
 - Enterprise Platforms
 - Cloud (Amazon AWS)



Companies Using R

Social media

Google
Facebook
Twitter
Foursquare
Kickstarter
eHarmony

Government

FDA
CPFB
City of Chicago
NOAA
NIST

Media

New York
Times
Economist
New Scientist
Xbox

Public Affairs

HRDAG
Sunlight
Foundation
Benetech
RealClimate

Software Vendors

Revolution
Analytics
Rstudio
Zementis
Alteryx
SAP
IBM
SAS
Teradata
TIBCO
Oracle
OneTick
DataCamp

Services

Mango
Accenture
Deloitte
Scientific Revenue
OpenBI
Coursera

Analytics

Zillow
Trulia
DataSong
Exelate
X+1
PredictWise

Finance

American
Century
ANZ
Credit Suisse
Nationwide
Lloyds

Manufacturing

Ford
John Deere
Monsanto
SZMF

Why are so many companies using R?

- Big Data
- Data Science
- Competition and Innovation
- Open Source
- Ecosystem
- People

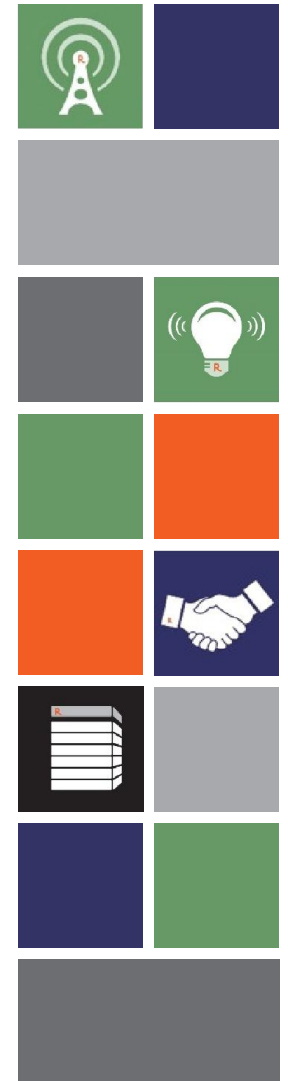


Thank You

David Smith

david@revolutionanalytics.com, @revodavid

blog.revolutionanalytics.com



Bonus Applications



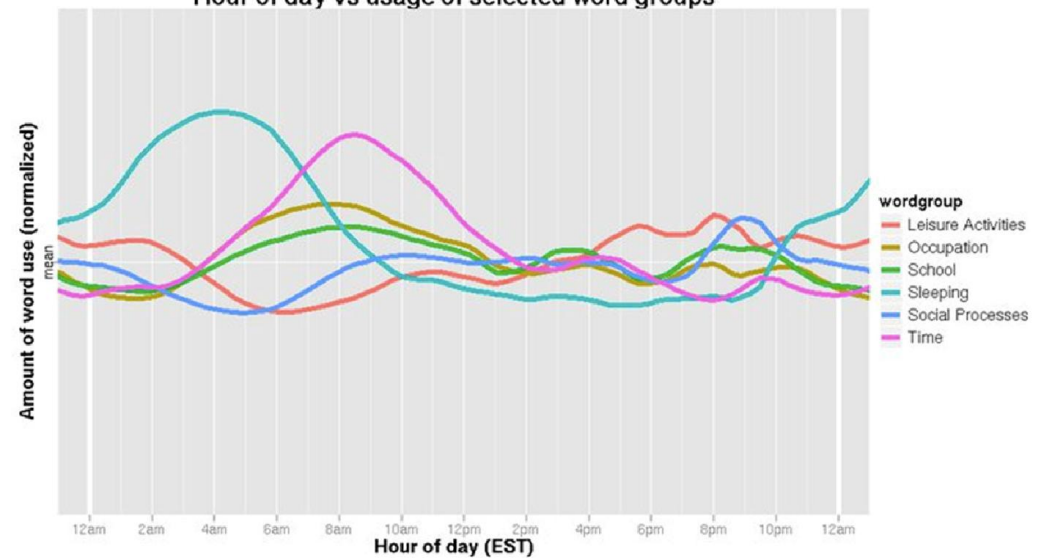


Eric Sun's Predicted Colleague Interactions

| | User | Score |
|---|--|-----------------|
| 1 |  Alexander Strehl | 9.3318008583172 |
| 2 |  Aaron Binbin Liao | 8.1162385345038 |
| 3 |  Srinivas Narayanan | 3.2938821239025 |
| 4 |  Jack Zhao | 3.0129828038274 |
| 5 |  Austin Haugen | 2.4928669082428 |

- [Human Resources](#)

Hour of day vs usage of selected word groups



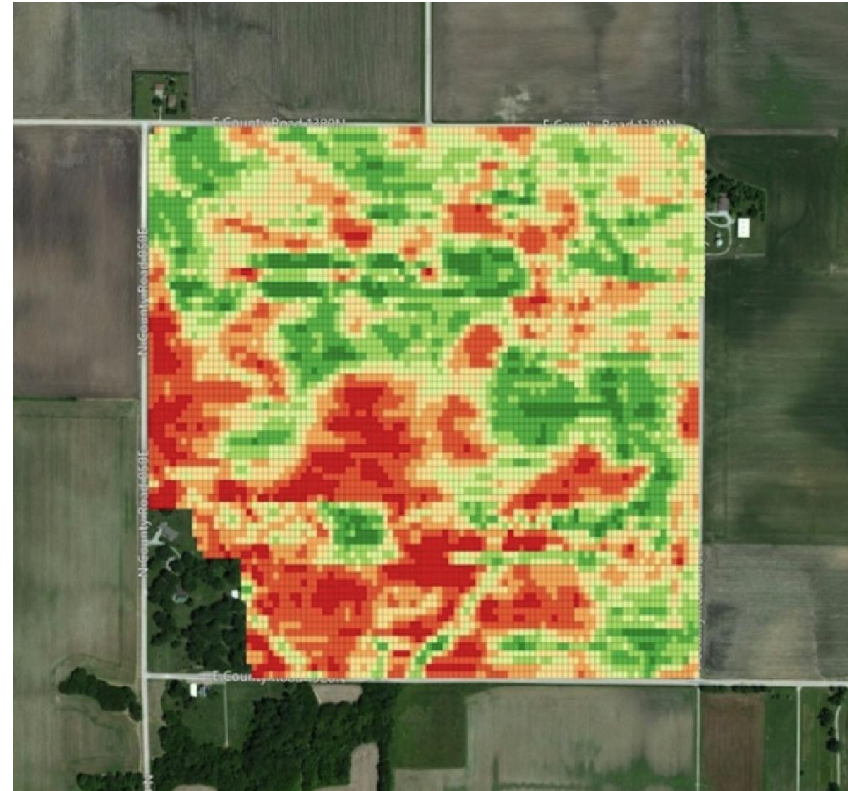
- [Status Update Trends](#)
- [User Profile Images](#)



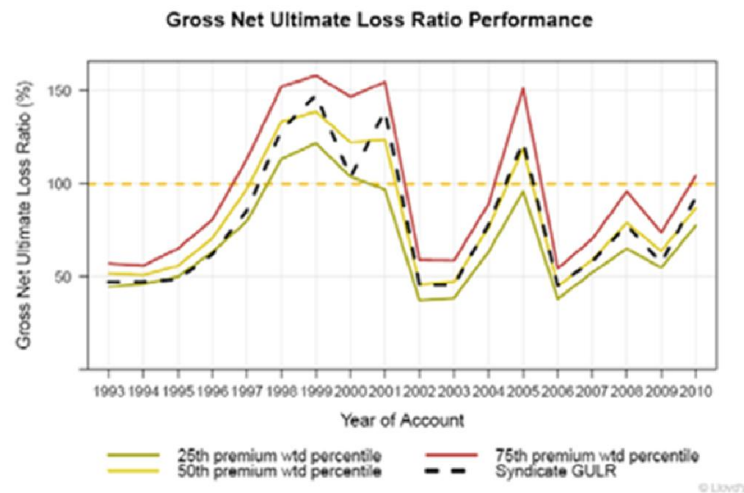


Statistical Analysis:

- Plant Breeding
- Fertility mapping
- Precision Seeding
- Disease Management
- Yield forecasting



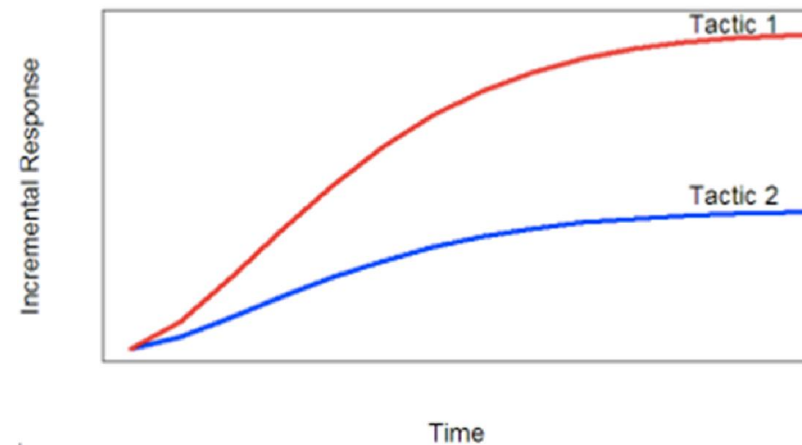
LLOYD'S LLOYD'S OF LONDON



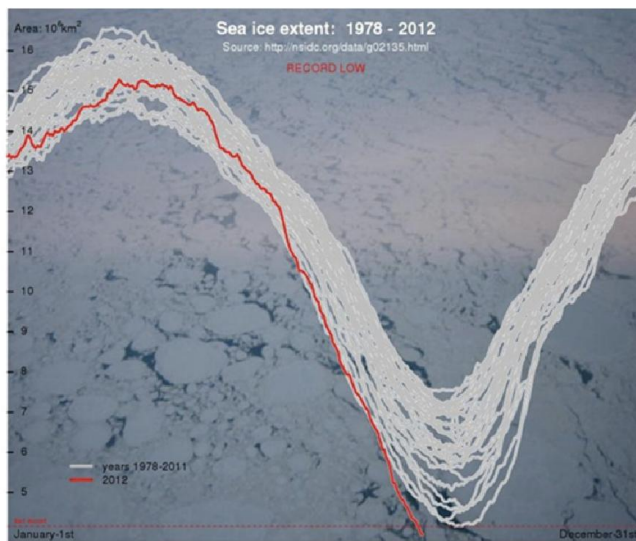
- [Risk Analysis](#)
- [Catastrophe Modeling](#)



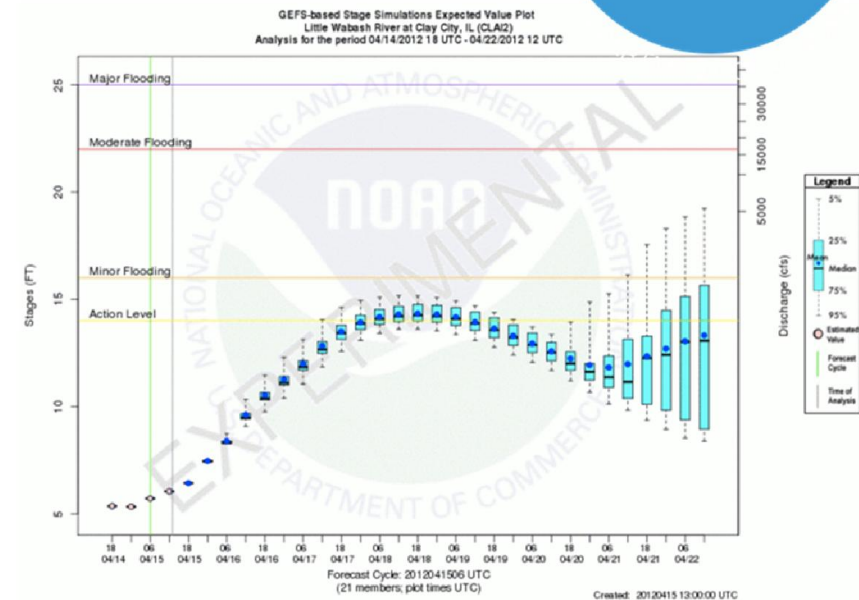
Nationwide®



- [Marketing Analytics](#)



- Climate change forecasts



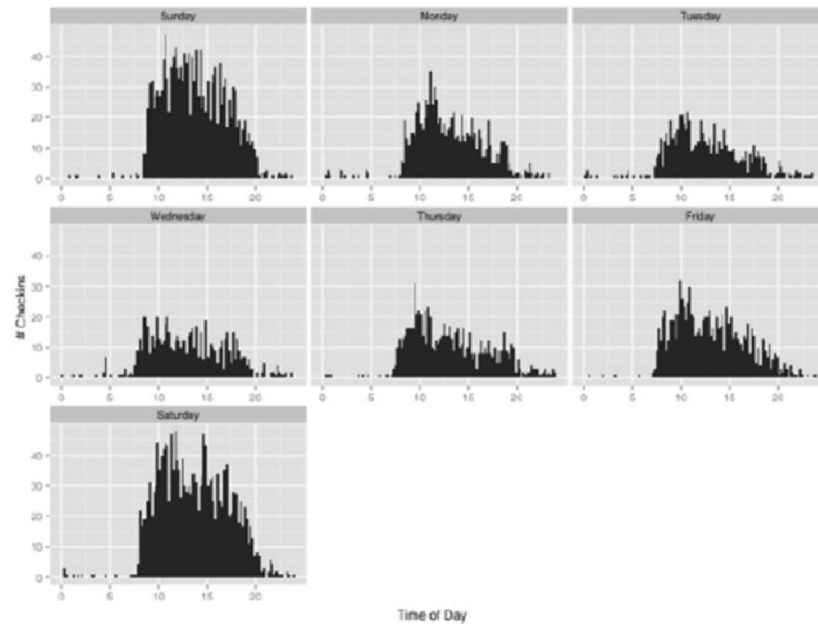
- Flood Warnings



REVOLUTION
ANALYTICS



foursquare™



- [Recommendation Engine](#)



REVOLUTION
ANALYTICS



Power

4X performance
50M records scored daily

"We've combined Revolution R Enterprise and Hadoop to build and deploy customized exploratory data analysis and GAM survival models for our marketing performance management and attribution platform. Given that our data sets are already in the terabytes and are growing rapidly, we depend on Revolution R Enterprise's **scalability and power** – we saw about a 4x performance improvement on 50 million records. **It works brilliantly.**"
- CEO, John Wallace, DataSong



Scalability

TB's data from 200+ data sources
10's thousands attributes
100's millions of scores daily

"We've been able to **scale** our solution to a problem that's so big that most companies could not address it. If we had to go with a different solution we wouldn't be as efficient as we are now."

- SVP Analytics, Kevin Lyons, eXelate



Make Every Interaction Count

Performance

2X data
2X attributes
no impact on performance

"We need a **high-performance** analytics infrastructure because marketing optimization is a lot like a financial trading. By watching the market constantly for data or market condition updates, **we can now identify opportunities for our clients that would otherwise be lost.**"

- Chief Analytics Officer, Leon Zemel,
[x+1]

