Learning from Web Activity

Jake Hofman

Joint work with Irmak Sirer & Sharad Goel Yahoo! Research

June 27, 2011

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 1 / 39



ICSA, 2011.06.27 2 / 39

э



"... applies social sciences techniques, including algorithmic game theory, economics, network analysis, psychology, ethnography, and mechanism design, to online situations."

http://research.yahoo.com

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 3 / 39

▲ロト ▲帰 ト ▲ 臣 ト ▲ 臣 ト 一 臣 … の へ ()

SOCIAL SCIENCE

Computational Social Science

David Lazer,¹ Alex Pentland,² Lada Adamic,² Sinan Aral,²⁴ Albert-László Barabási,⁵ Devon Brewer,⁶ Nicholas Christakis, ¹ Noshir Contractor,⁷ James Fowler,⁸ Myron Gutmann,³ Tony Jebara,⁸ Gary King,¹ Michael Macy,¹⁰ Deb Roy² Marshall Van Alstyne^{2,11} A field is emerging that leverages the capacity to collect and analyze data at a scale that may reveal patterns of individual and group behaviors.

"... a computational social science is emerging that leverages the capacity to collect and analyze data with an unprecedented breadth and depth and scale ..."

http://sciencemag.org/content/323/5915/721

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 4 / 39

イロト 不得下 イヨト イヨト 二日

SOCIAL SCIENCE

Computational Social Science

David Lazer,¹ Alex Pentland,² Lada Adamic,² Sinan Aral,²⁴ Albert-László Barabási,⁵ Devon Brewer,⁶ Nicholas Christakis, ¹ Noshir Contractor,⁷ James Fowler,⁸ Myron Gutmann,³ Tony Jebara,⁸ Gary King,¹ Michael Macy,¹⁰ Deb Roy² Marshall Van Alstyne^{2,11} A field is emerging that leverages the capacity to collect and analyze data at a scale that may reveal patterns of individual and group behaviors.

"... shares with other nascent interdisciplinary fields (e.g., sustainability science) the need to develop a paradigm for training new scholars ..."

http://sciencemag.org/content/323/5915/721

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 4 / 39

イロト 不得下 イヨト イヨト 二日

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 Lecture¹, 1965

イロト イポト イヨト イヨト

ICSA, 2011.06.27 5 / 39

¹http://bit.ly/feynmannobel

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

Outline

- The clean story
- The real story
- Lessons learned

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 6 / 39

э

Demographic diversity on the Web

with Irmak Sirer and Sharad Goel

The clean story (covering our tracks)

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 7 / 39

イロト 不得下 イヨト イヨト 二日

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 and group-level differences in online access

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 8 / 39

▲□▶ ▲□▶ ▲三▶ ▲三▶ - 三 - つくや

"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)

イロト 不得下 イヨト イヨト



Figure 6: Relative saturation of ethnicities on Facebook. As the lines converge towards 100% (center), the makeup of U.S. Facebook converges towards that of the addressable Internet population.

Chang, et. al., ICWSM (2010)

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

< □ > < 同 >

10.0

ICSA, 2011.06.27

9 / 39

Focus on activity instead of access





イロト イポト イヨト イヨト

How diverse is the Web?

To what extent do online experiences vary across demographic groups?

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 10 / 39

nielsen MegaPanel

- Representative sample of 265,000 individuals in the US, paid via the Nielsen MegaPanel²
- 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.)

²Special thanks to Mainak Mazumdar

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

イロト イポト イヨト イヨト

ICSA, 2011.06.27 11 / 39

• Transform all demographic attributes to binary variables e.g., Age \rightarrow Over/Under 25, Race \rightarrow White/Non-White, Sex \rightarrow Female/Male

イロト 不得 トイヨト イヨト 二日

- Transform all demographic attributes to binary variables e.g., Age \rightarrow Over/Under 25, Race \rightarrow White/Non-White, Sex \rightarrow Female/Male
- 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

イロト 不得下 イヨト イヨト 二日

- Transform all demographic attributes to binary variables e.g., Age \rightarrow Over/Under 25, Race \rightarrow White/Non-White, Sex \rightarrow Female/Male
- 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)

イロト 不得 トイヨト イヨト 二日

- Transform all demographic attributes to binary variables e.g., Age \rightarrow Over/Under 25, Race \rightarrow White/Non-White, Sex \rightarrow Female/Male
- 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

▲□▶ ▲□▶ ▲三▶ ▲三▶ - 三 - つくや

How diverse are site audiences?

- For each site and attribute, calculate the skew in visitors (e.g., 93% of pageviews on foxnews.com are by White users)
- For each attribute, plot the distribution of visitor skew across all sites



イロト イポト イヨト イヨト

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 13 / 39



Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 14 / 39

Many sites have skew close the overall mean, but there also popular, highly-skewed sites

	Greater Than 90%	Less Than 10%
Female	youravon.com	coveritlive.com
	collectionsetc.com	needlive.com
White	foxnews.com	blackplanet.com
	wunderground.com	mediatakeout.com
College Educated	news.google.com	slumz.boxden.com
	nytimes.com	sythe.com
Over 25 Years Old	mail.yahoo.com	nanowrimo.org
	apps.facebook.com	cbox.ws
Household Income	scarleteen.com	opentable.com
Under \$50,000	boards.adultswim.com	marketwatch.com

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 15 / 39

This skew persists even when we restrict attention to the top 10k or 1k sites



Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 16 / 39

Sites vs. ZIPs

How do diversity of the online and offline worlds compare?



Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 17 / 39

3

Sites vs. ZIPs

How do diversity of the online and offline worlds compare?



As expected, neighborhoods are more gender-balanced than sites

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 17 / 39

3

Sites vs. ZIPs

How do diversity of the online and offline worlds compare?



But sites typically have more racially diverse audiences than neighborhoods have residents

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 17 / 39

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)

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 18 / 39

All groups spend more than a third of their time on a handful of email, search, and social networking sites



Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 19 / 39

But different groups distribute their time differently, both on universally popular and on more niche sites





There is both reasonable overlap and variation amongst the most popular sites within groups.



Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 21 / 39

3

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

- Represent each user by the set of sites visited
- Fit linear models³ 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

³Using SVM-perf

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 22 / 39

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



Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 23 / 39

Highly-weighted sites under the fitted models

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

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 24 / 39

▲□▶ ▲□▶ ▲臣▶ ▲臣▶ 三臣 - のへ⊙

Similar performance even when restricted to top 1k sites



Jake Hofman (hofman@yahoo-inc.com) Learning from Web Activity

 $\exists \rightarrow$ ICSA, 2011.06.27 25 / 39

A (1) > A (2) >

Substantially better performance when restricted to "stereotypical" users (~80-90%)



Learning from Web Activity

 $\exists \rightarrow$ ICSA, 2011.06.27 26 / 39

< 1 > <

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

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 27 / 39

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Demographics diversity on the Web

with Irmak Sirer and Sharad Goel

The real story (what we actually did)

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 28 / 39

イロト 不得下 イヨト イヨト 二日

0. Got several hundred GBs of MegaPanel data from Nielsen⁴, looked at a small sample

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

⁴Special thanks to Mainak Mazumdar

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

≣ ► < ≣ ► ≡ ∽ ९ ज ICSA, 2011.06.27 29 / 39

イロト 不得下 イヨト イヨト

1. Discussed (many) possible projects

- Infer the number of individuals using the same browser or behind the same ip?
- Determine number of actual uniques advertisers are receiving?
- Predict user demographics from a few minutes of browsing activity for ad-targeting?

2. Modeled real-valued age as a function of site visits

- Worked on this for an embarassingly long time
- Tried various options for cleaning and normalizing data (100GB $\rightarrow \sim 5 \text{GB})$
- Investigated several methods for feature selection (e.g., naive Bayes, mutual information, *popularity*)

▲□▶ ▲□▶ ▲三▶ ▲三▶ - 三 - つくで

Hadoop + Pig (+ awk)

```
-- define streaming command for normalizing urls
DEFINE add star domains `awk 'n=split($2,a,","): (n==3) {print $1"\t*,"a[2]","a[3]"\t"$3}'`:

    flatten user histories

user_pageviews = FOREACH users GENERATE
        uid.
        FLATTEN(history) AS (did, pageviews, weight);
-- join user pageviews against top domains
user_pageviews = JOIN user_pageviews BY did, top_domains BY did USING 'replicated';
user pageviews = FOREACH user pageviews GENERATE
        user pageviews::uid AS uid.
        top domains::domain AS domain,
        user pageviews::pageviews AS pageviews;
-- stream through awk to extract two-level domains (e.g. *.yahoo.com)
user pageviews = STREAM user pageviews THROUGH add star domains AS (uid:long. domain:chararray. pageviews:int):
-- regroup and count pageviews by normalized domains
— (userid, normalized_domain, num_pageviews)
user pageviews = GROUP user pageviews BY (uid. domain) PARALLEL 10:
user pageviews = FOREACH user pageviews GENERATE
        group.uid AS uid,
        group.domain AS domain,
        SUM(user pageviews.pageviews) AS pageviews:
```

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 32 / 39

▲□▶ ▲□▶ ▲ヨ▶ ▲ヨ▶ - ヨ - のの⊙

3. Settled for classification of binary outcomes (e.g. adult/non-adult)

- 265,000 users (examples) and 100,000 sites (features)
- Logistic regression in R, e.g.

model <- glm(is.adult ~ ., data=nielsen, family=binomial)</pre>

???

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 33 / 39

3. Settled for classification of binary outcomes (e.g. adult/non-adult)

$$\hat{y}(x_i) = w \cdot x_i + b$$

Support Vector Machine : $L(y, \hat{y}) = C \sum_{i} [1 - y_i \hat{y}(x_i)]_+ + ||w||^2$

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 34 / 39

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

4. Investigated why classification worked reasonably well

- Generated descriptive statistics across all attributes at the site and group levels
- Compared site statistics to ZIP code data from the US Census
- Compared time distribution across groups

5. Realized that we now had one of the most comprehensive studies available on demographic diversity of the web

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 36 / 39

イロト 不得下 イヨト イヨト 二日

Conclusion (lessons learned)

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 37 / 39

<ロト < 回 > < 回 > < 回 > < 回 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < 三 > < < = < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Data jeopardy

Regardless of scale, it's difficult to find the right questions to ask of the data

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 38 / 39

イロト 不得下 イヨト イヨト 二日

Rapid iteration

The ability to iterate quickly, asking and answering many questions, is crucial

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 38 / 39

Data cleaning

Cleaning and normalizing data is a substantial amount of the work

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 38 / 39

Modeling

Simple methods (e.g., linear models) work surprisingly well, especially with lots of data

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 38 / 39

▲□▶ ▲□▶ ▲三▶ ▲三▶ - 三 - つくで

It's easy to cover your tracks—things are often much more complicated than they appear

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 38 / 39

イロト 不得下 イヨト イヨト 二日

Thanks. Questions?

http://messymatters.com/webdemo

http://jakehofman.com hofman@yahoo-inc.com

Jake Hofman (hofman@yahoo-inc.com)

Learning from Web Activity

ICSA, 2011.06.27 39 / 39

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○○