# Genome-Wide Association with Digital Phenotypes

Jie Yuan

5/3/2016

# Traditional Genome-Wide Association



Questionnaires, Medical Records → Phenotype (+/- heart disease)

y

linear/logistic regression

Genomic Sequencing

```
# rsid       chromosome    position    genotype
rs4477212         1         72017          AA
rs3094315         1         742429         AG
rs3131972         1         742584         AG
rs12124819        1         766409         AA
rs11240777        1         788822         AA
rs6681049         1         789870         CC
rs4970383         1         828418         CC
rs4475691         1         836671         CC
rs7537756         1         844113         AA
rs13302982        1         851671         GG
rs1110052         1         863421         GT
...
```
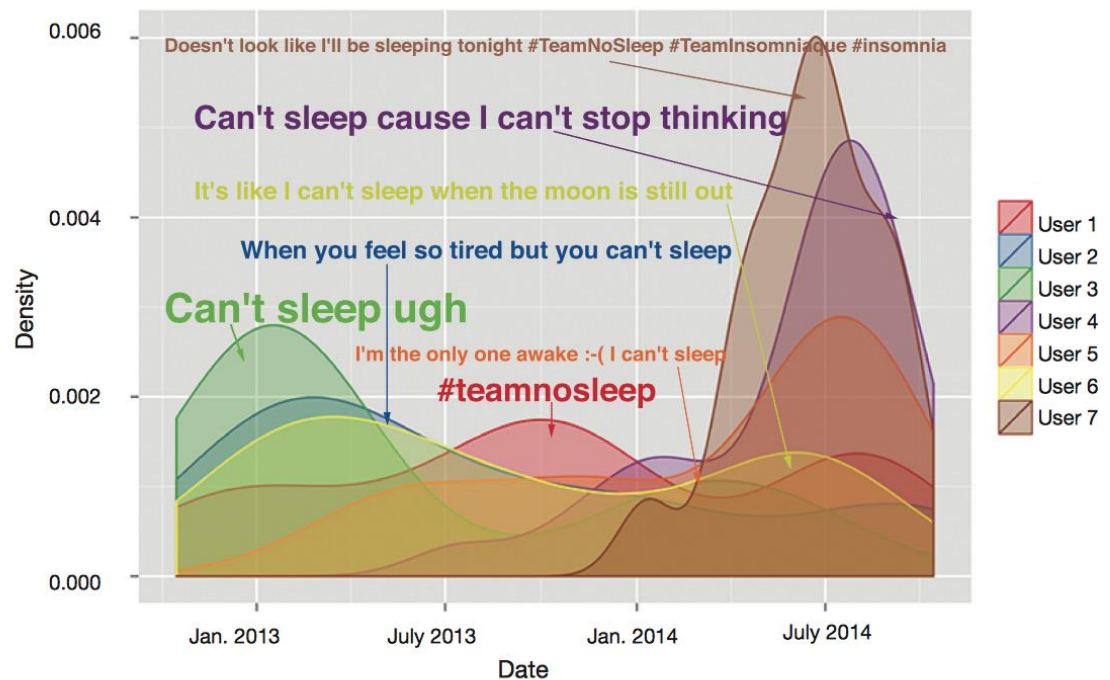
x

# The digital phenotype

Sachin H Jain, Brian W Powers, Jared B Hawkins & John S Brownstein

**In the coming years, patient phenotypes captured to enhance health and wellness will extend to human interactions with digital technology.**

In 1982, the evolutionary biologist Richard Dawkins introduced the concept of the "extended phenotype"[1], the idea that phenotypes should not be limited just to biological processes, such as protein biosynthesis or tissue growth, but extended to include all effects that a gene has on its environment inside or outside of the body of the individual organism. Dawkins stressed that many delineations of phenotypes are arbitrary. Animals and humans can modify their environments, and these modifications and associated behaviors are expressions of one's genome and, thus, part of their extended phenotype. In the animal kingdom, he cites damn building by beavers as an example of the beaver's extended phenotype[1].

As personal technology becomes increasingly embedded in human lives, we think there is an important extension of Dawkins's theory—the notion of a 'digital phenotype'. Can aspects of our interface with technology be somehow diagnostic and/or prognostic for certain conditions? Can one's clinical data be linked and analyzed together with online activity and behavior data



**Figure 1** Timeline of insomnia-related tweets from representative individuals. Density distributions (probability density functions) are shown for seven individual users over a two-year period. Density on the y axis highlights periods of relative activity for each user. A representative tweet from each user is shown as an example.

# Facebook Likes are predictive of a variety of demographic traits
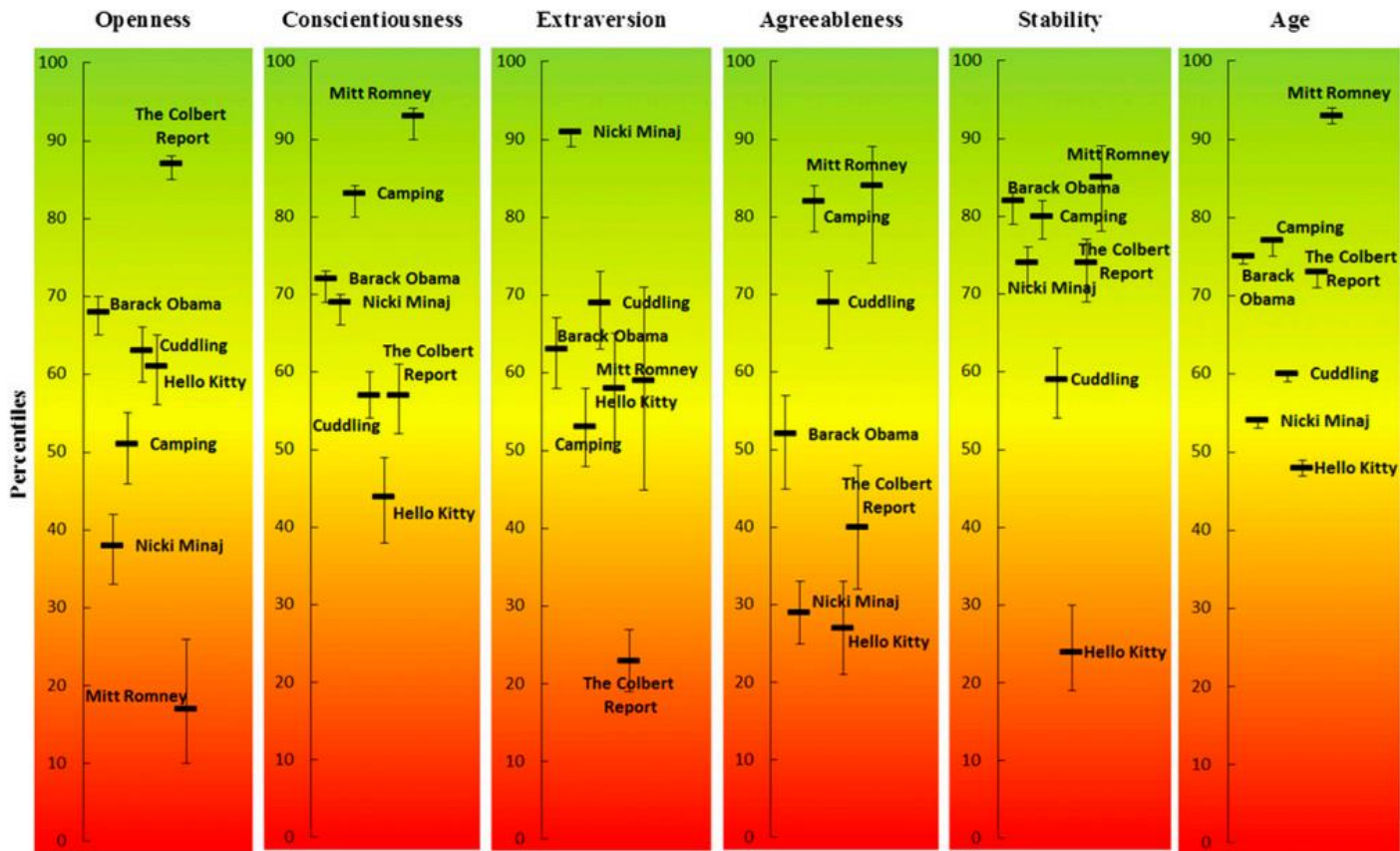


Fig. S1. Average levels of five personality traits and age of the users associated with selected Likes presented on the percentile scale. For example, the average extraversion of users associated with "The Colbert Report" was relatively low: it was lower only for 23% of other Likes in the sample. Error bars signify 95% confidence intervals of the mean.
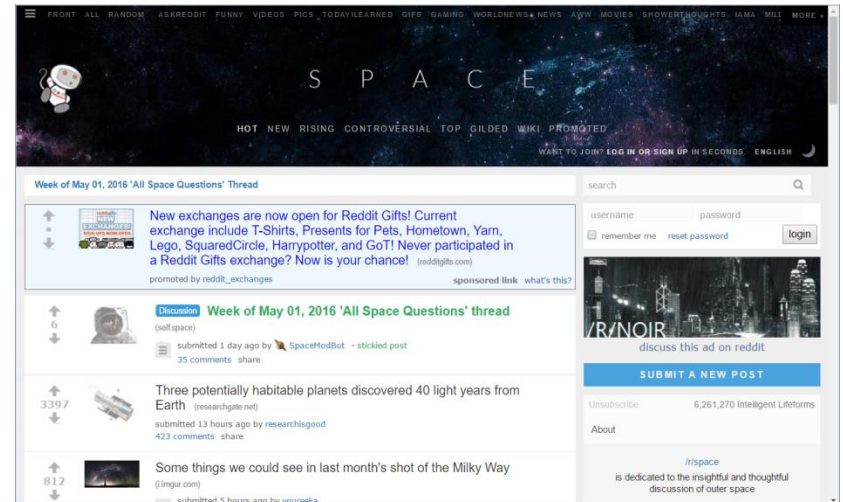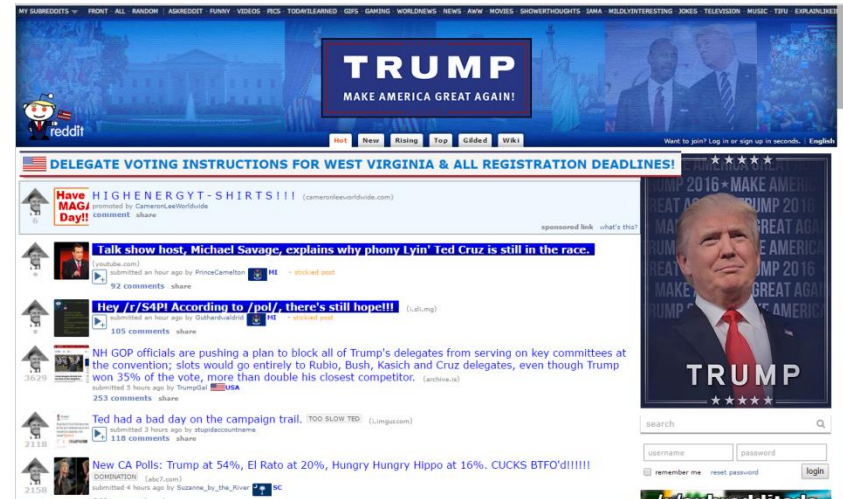
(But still requires questionnaires for labels)

Kosinski et al. (2013)

Problem:

- Can we identify traits and word-trait associations from text, without access to external labels (e.g. questionnaire results)?

  – An extension of Dr. Mark Dredze's sLDA classification problem

- Given these inferred traits, can we identify genome-trait associations?

# Reddit is divided into communities (subreddits) with highly specialized interests

# Discovering Word Associations

- Can we find overrepresented words within subreddits?

  - Calculate $\frac{f-g}{f+g}$, where $f$ and $g$ are % frequencies in the subreddit and rest of Reddit, respectively

  - Frequency cutoff filter to remove rare words (e.g. typos) appearing in the subreddit

# A sampling of medically relevant subreddits and their overrepresented words

| r/depression | r/diabetes | r/stopsmoking | r/insomnia |
|---|---|---|---|
| therapist | insulin | nicotine | insomnia |
| medication | diabetic | quitting | melatonin |
| meds | glucose | cigarettes | mg |
| suicidal | carb | smoker | sleepy |
| antidepressants | dexcom | congratulations | ambien |
| scared | diagnosed | badge | pills |
| counselor | basal | smokers | caffeine |
| mg | keto | nonsmoker | meds |
| diagnosed | endo | addicted | addictive |
| sadness | lantus | carrs | medication |

- Create a matrix of traits, and word frequencies (or sLDA topics)

|  | depression | insomnia | "sleepy" | "caffeine" | (topic 1) | (topic 2) |
|---|---|---|---|---|---|---|
| user | 0 | 1 | 15 | 10 | 0.20 | 0.05 |
| user | 1 | 1 |  |  |  |  |
| user | 0 | 0 |  |  | . |  |
| user* | ? | ? |  |  | . |  |
| user* | ? | ? |  |  | . |  |
| user* | ? | ? |  |  |  |  |

+/- phenotype (is the user a frequent poster on relevant subreddit)

frequencies for overrepresented words, and/or significant topics from sLDA

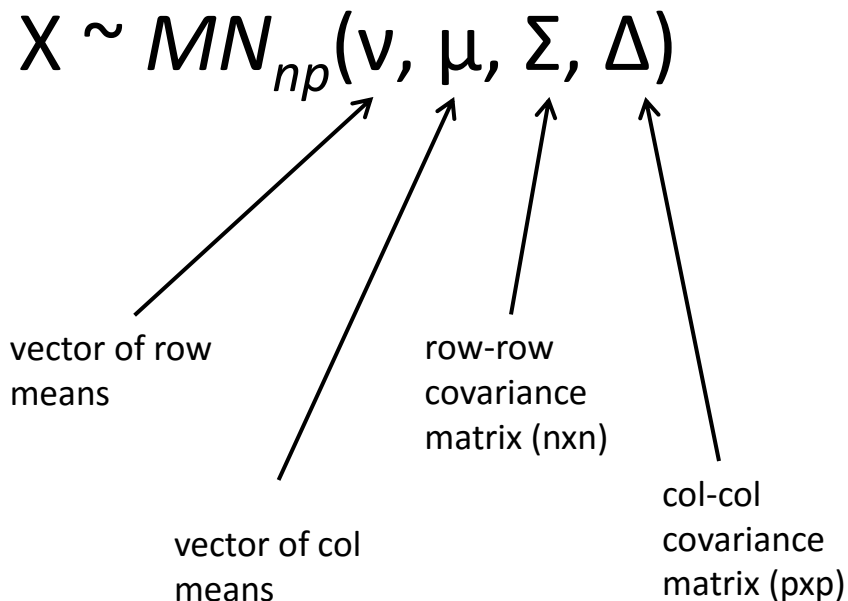- Create a matrix of traits, and word frequencies (or sLDA topics)

| | depression | insomnia | "sleepy" | "caffeine" | (topic 1) | (topic 2) |
|---|---|---|---|---|---|---|
| user | 0 | 1 | 15 | 10 | 0.20 | 0.05 |
| user | 1 | 1 | | | | |
| user | 0 | 0 | | | . | |
| user* | ? | ? | | | . | |
| user* | ? | ? | | | . | |
| user* | ? | ? | | | | |

user: a Reddit user whose digital phenotypes are the set of frequently visited subreddits

user*: a user with unlabeled social media text (e.g. from Facebook), whose digital phenotypes we want to impute
- Can we impute the "?" entries

# Matrix imputation using Transposable Regularized Covariance Model (TRCM)

$$X \sim MN_{np}(\nu, \mu, \Sigma, \Delta)$$

vector of row means

vector of col means

row-row covariance matrix (nxn)

col-col covariance matrix (pxp)

- Use EM to calculate row and col parameters

- Plug in mean estimates for imputed values

For a missing value in row $i$ and col $j$

$$x_{ij} = \nu_i + \mu_i + \varepsilon_{ij}$$

$$\varepsilon_{ij} \sim N(0, \Sigma_{ii}\Delta_{jj})$$

Allen, Tibshirani. (2010)