**Discovering Psychological and Health Insights from Social Media Langugae**

H. Andrew Schwartz

andrew.schwartz@cs.stonybrook.edu
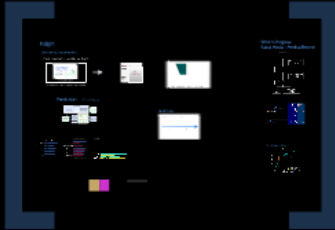
Computational Linguistics | Psychology/ Health

October 16, 2015
@ SBU-SUNYK Seminar

x 20mil.

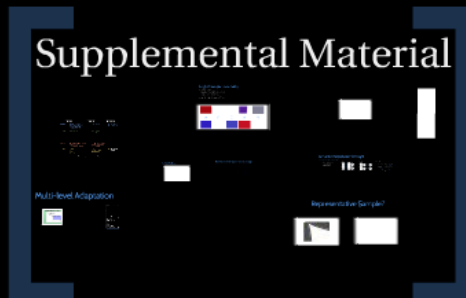**1. Individual Analyses**

x 1bil.

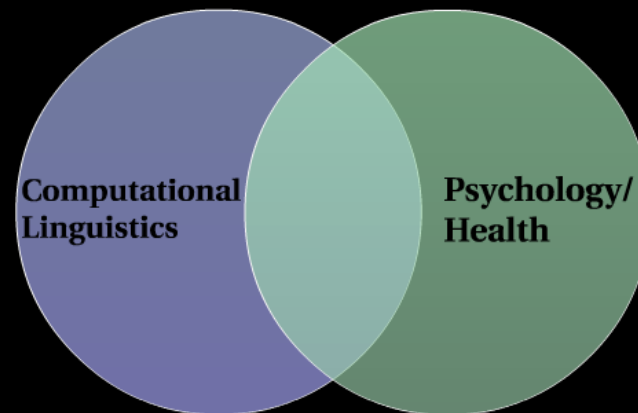**2. Community Analyses**

**Introduction**

Supplemental Material

# Discovering Psychological and Health Insights from Social Media Langugae

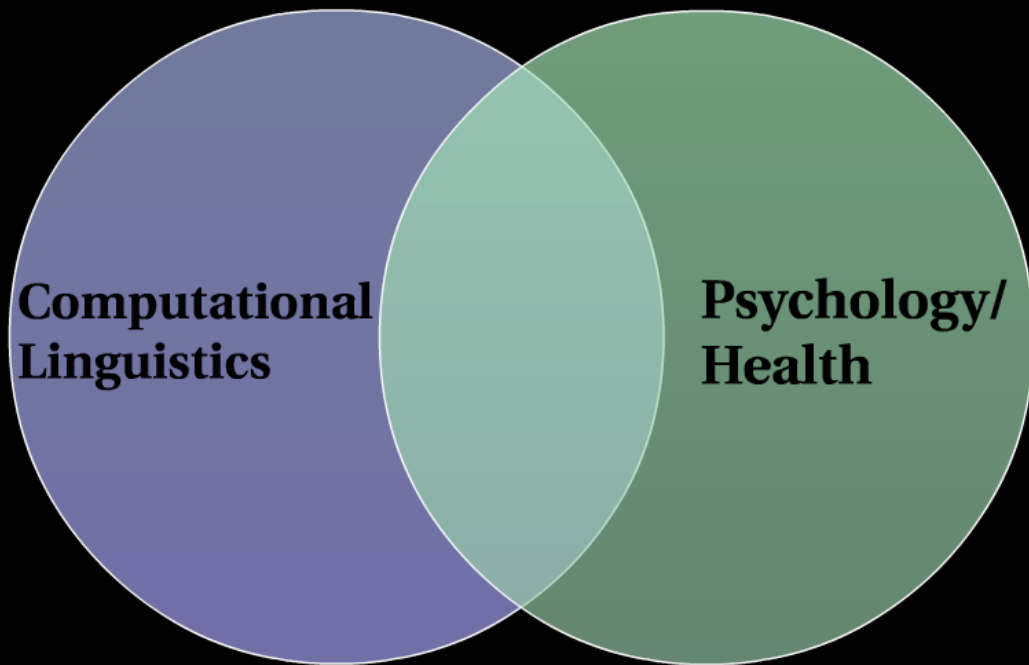## H. Andrew Schwartz

andrew.schwartz@cs.stonybrook.edu

Computational Linguistics

Psychology/ Health

October 16, 2015
@ SBU SUNYK Seminar

# H. Andrew Schwartz

andrew.schwartz@cs.stonybrook.edu

Computational Linguistics

Psychology/ Health

r

# WWBP Collaborators

Rosie Hancock  Darwin Labarthe
Richard Lucas  Luke Dziurzynski  Michal Kosinski
Yiyi Guo  Sneha Jha  Gregory Park  Shawndra Hill
Tadas Antanavicius
Chris Weeg  Lyle H. Ungar  Robert Backer
Jeanette Elstein  Margaret Kern  Rigel Swavely
Liwei Xu  Stephanie Ramones  George Wan  Megha Agrawal
Dolores Albaraccin  Maarten Sap  Johannes Eichstaedt
Achal Shah  Emily Larson  Martin E. P. Seligman
Brian Galla  Libby Benson  Eduardo Blanco  Marie Foregeard
Annie Roepke  Yoon Hwang
Arsenij Kouriatov  Raina Merchant  Saif Mohammad
Bob Stine  Winnie Cheng  David Stillwell  Molly Ireland
Evan Weingarten  Jonah Berger
Dean Foster
Daniel Peotiuc  Jordan Carpenter

UNIVERSITY OF CAMBRIDGE    HARVARD UNIVERSITY    LYMBA the power to answer

Wharton UNIVERSITY of PENNSYLVANIA    Perelman School of Medicine UNIVERSITY of PENNSYLVANIA    PENN SOCIAL MEDIA AND HEALTH INNOVATION LAB

# a friend

- measures non-objective outcomes
- many developed over decades
- a gold-standard

## a friend

- measures non-objective outcomes
- many developed over decades
- a gold-standard



## a foe

- hard to administer at scale
  - spatially
  - temporally
- not exactly "ground truth"
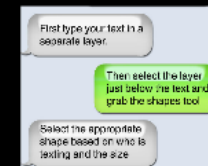- limited to preconceived theory

# Social Media

350m tweets/day

4b messages/day

in

INSTAGRAM

Insta

tumblr.

g+

WhatsApp

WeChat

First type your text in a separate layer.

Then select the layer just below the text and grab the shapes tool

Select the appropriate shape based on who is texting and the size

...

# Social Media

350m tweets/day

4b messages/day

...the largest dataset of who we are

**f** 1b **people**

150m **people**
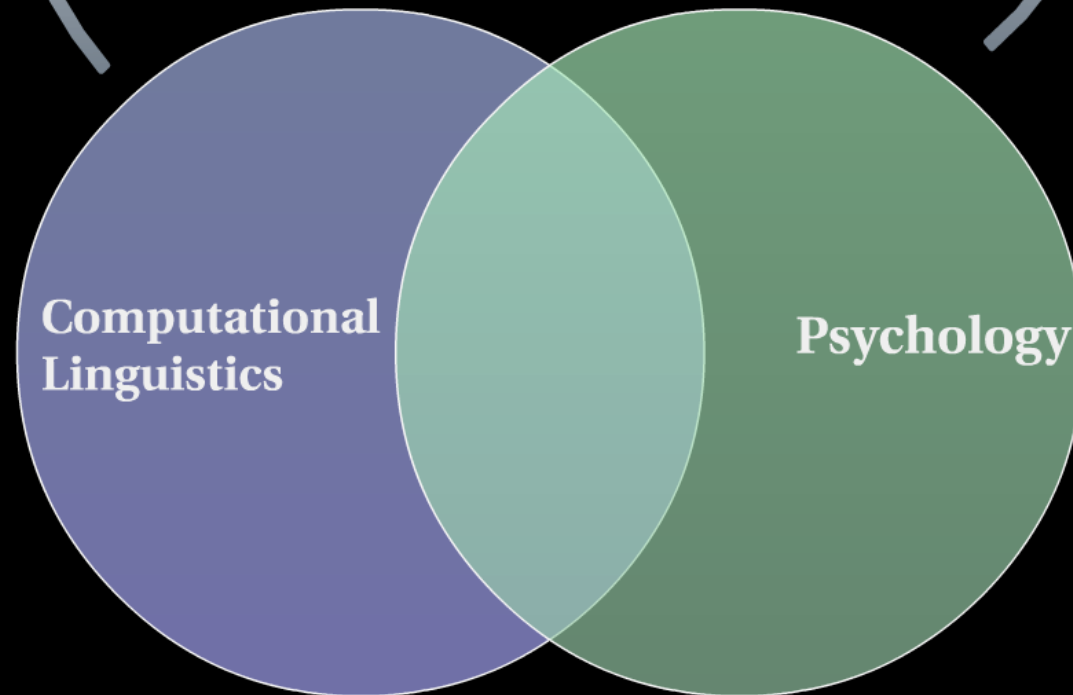
from *modeling language*
to *understanding people*

Computational Linguistics

Psychology

**data-driven human insights**

**goal**: accurate prediction

**method**: data-driven
(large data)

**goal**: human insights

**method**: closed-vocabulary
(typically "small" samples)

N=75,000 users...          ...20M statuses.

extraversion

neuroticism

⚠ explicit language warning

Schwartz, H. A., Eichstaedt, J. C., Dziurzynski, L., Kern, M. L., Blanco, E., Kosinski, M., Stillwell, D., Seligman, M. E. P., & Ungar, L. H. (2013). Toward Personality Insights from Language Exploration in Social Media. In *Proceedings of the AAAI Spring Symposium Series*. Stanford, CA.

correlation strength

relative frequency

extraversion

⚠ explicit language warning

neuroticism

Schwartz, H. A., Eichstaedt, J. C., Dziurzynski, L., Kern, M. L., Blanco, E., Kosinski, M., Stillwell, D., Seligman, M. E. P., & Ungar, L. H. (2013). Toward Personality Insights from Language Exploration in Social Media. In *Proceedings of the AAAI Spring Symposium Series*. Stanford, CA.

# explicit language warning

extraversion — neuroticism

⚠ explicit language warning

correlation strength — relative frequency

Schwartz, H. A., Eichstaedt, J. C., Dziurzynski, L., Kern, M. L., Blanco, E., Kosinski, M., Stillwell, D., Seligman, M. E. P., & Ungar, L. H. (2013). Toward Personality Insights from Language Exploration in Social Media. In *Proceedings of the AAAI Spring Symposium Series*. Stanford, CA.

extraversion — neuroticism

⚠ explicit language warning

Schwartz, H. A., Eichstaedt, J. C., Dziurzynski, L., Kern, M. L., Blanco, E., Kosinski, M., Stillwell, D., Seligman, M. E. P., & Ungar, L. H. (2013). Toward Personality Insights from Language Exploration in Social Media. In *Proceedings of the AAAI Spring Symposium Series*. Stanford, CA.

from *modeling language*
to *understanding people*

**Computational Linguistics**

**Psychology**

Can we predict disease risk and recovery from language use?

What psychological factors emerge in language as drivers of health and well-being?

To what extent can language analyses replace and extend traditional psychological asessement

# Interdisciplinary Research Questions

What psychological factors emerge in language as drivers of health and well-being?

Can we predict disease risk and recovery from language use?

To what extent can language analyses replace and extend traditional psychological asessement

**Discovering Psychological and Health Insights from Social Media Langugae**

H. Andrew Schwartz
andrew.schwartz@cs.stonybrook.edu

Computational Linguistics / Psychology/Health

October 16, 2015
@ SBU SUNYK Seminar

x 20mil.

**1. Individual Analyses**

x 1bil.

**2. Community Analyses**

**Introduction**

# Differential Language Analysis

## Goal: succinct accessible method



(black box)

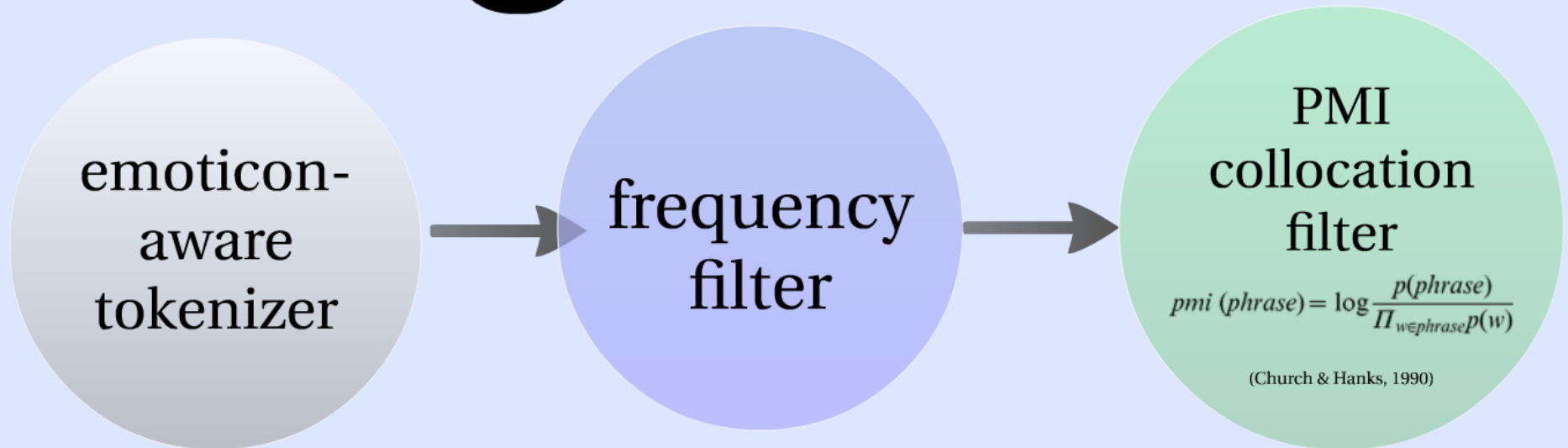"Simplicity is the ultimate sophistication" -da Vinci

# Goal: succinct accessible metho[d]



(black box)

"Simplicity is the ultimate sophistication"  -da Vinci

# ngrams

emoticon-aware tokenizer

→

frequency filter

→

PMI collocation filter

$$pmi\,(phrase) = \log \frac{p(phrase)}{\Pi_{w \in phrase}\, p(w)}$$

(Church & Hanks, 1990)

$$p(topic \mid subject) \qquad \sum_{word \in topic} p(topic \mid word) \cdot p(word$$

# topics

latent Dirichlet allocation



- 2000 social media topics
- derived over 14m status updates
- status update = document
- status updates are shorter than news articles
  => lower alpha hyperparameter

chicken cheese dinner yum
soup made rice bacon bread
yummy fried eating salad
cooking eggs sauce making
eat potatoes

on friday monday sunday saturday tuesday
thursday night wednesday weekend next
week until morning afternoon tomorrow till
working nights

snow cold weather outside warm hot its
degrees winter heat freezing snowing
ice here inside inches summer degree
storm

my mom dad with husband
wife parents son daughter
kids love ex mum brother
wonderful hubby sister
boyfriend mother

party tonight at halloween
birthday night fun club bar
dance costume bday parties
come saturday via dj sms
house

pain blood hospital teeth surgery
doctor from tooth dentist wisdom
after has brain having doctors had
doc pulled pressure

wonderful hubby sister boyfriend mother

pain blood hospital teeth surgery doctor from tooth dentist wisdom after has brain having doctors had doc pulled pressure

made h...
mmy fried eating...
oking eggs sauce making
eat potatoes

ill morn...
wo...

snow cold weather outside warm hot its degrees winter heat freezing snowing ice here inside inches summer degree storm

party tonight at halloween
birthday night fun club...
dance costu...

$$p(topic \mid subject) = \sum_{word \in topic} p(topic \mid word) * p(word \mid subject)$$

# topics

- 2000 social media topics

# analysis

adjust for other variables via ordinary least squares linear regression over standardized variables

e.g. $\beta_0 + \beta_1(ngram) + \beta_2(age) + \beta_3(gender) = extraversion$

feature 1 $\longrightarrow$ fit $\longrightarrow$ correlation

feature 2 $\longrightarrow$ fit $\longrightarrow$ correlation

... ...

feature n $\longrightarrow$ fit $\longrightarrow$ correlation

SAM BIDDLE — FACEBOOK · Tuesday 2:36pm          65,818 🔥    204 💬

# Science Shows Men and Women Are Both Awful Stereotypes on Facebook

⚠️
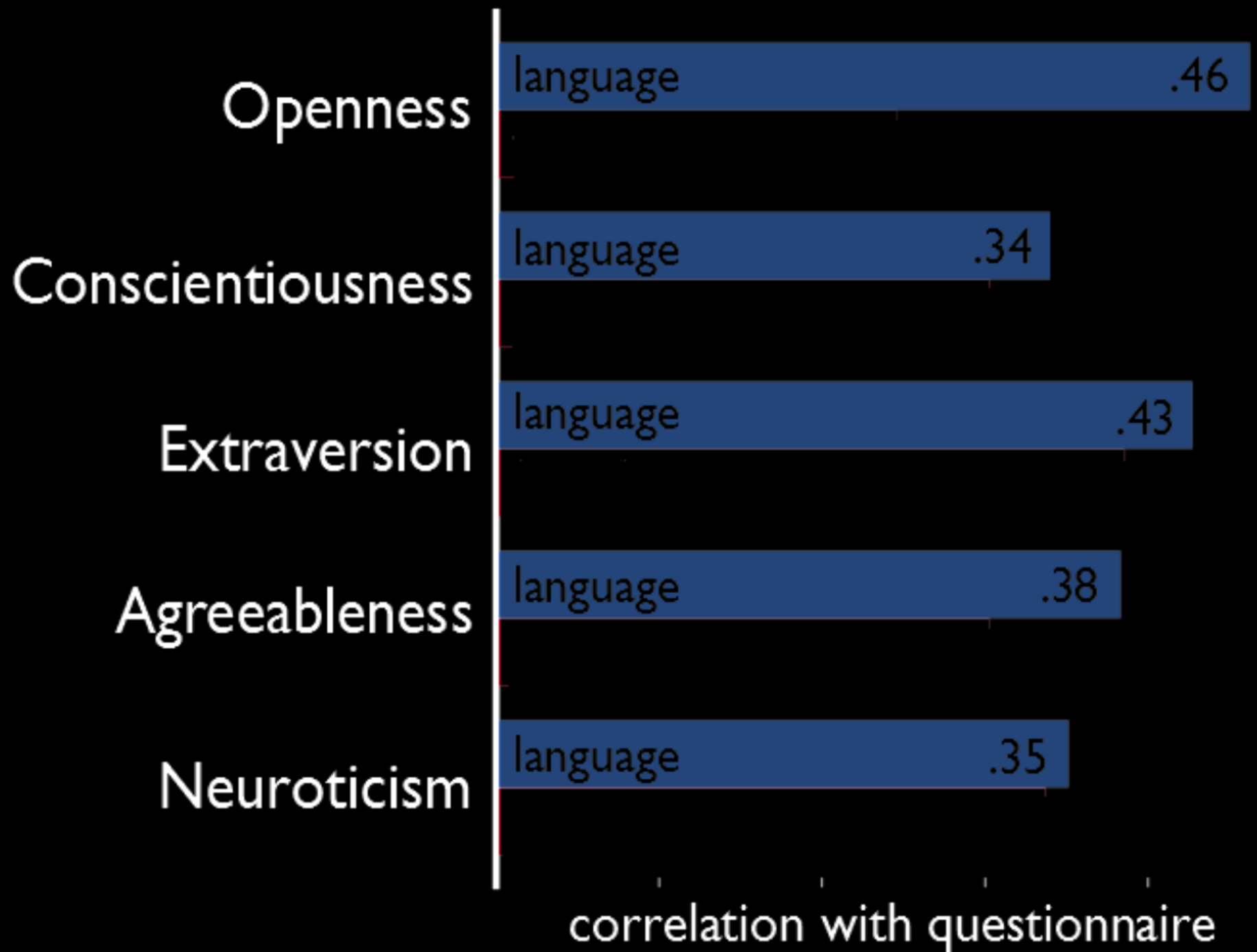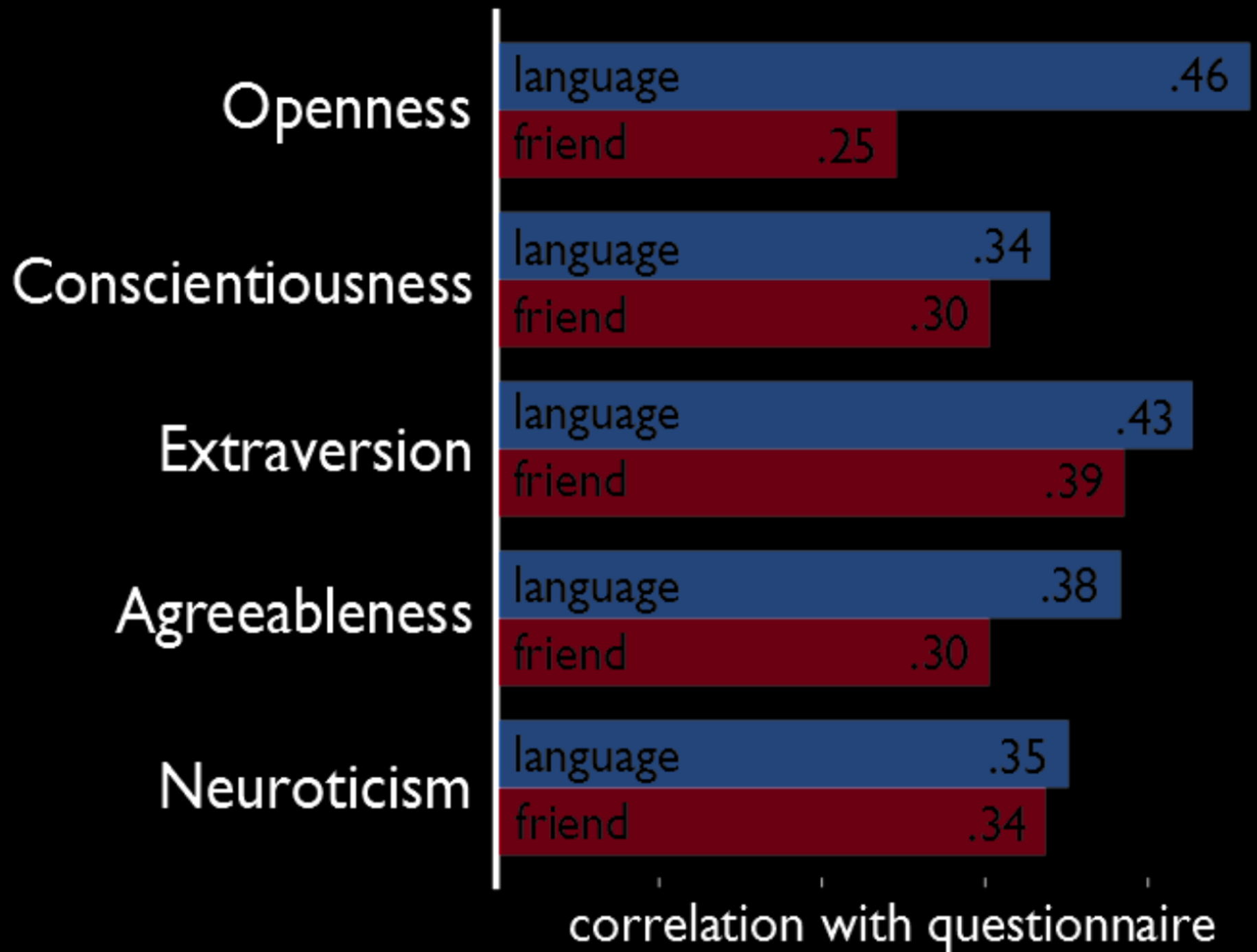
## gender cliches warning

season
nfl
football
baseball
fantasy
sports play league
playing player
team game basketball
players coach

government
state
power country
freedom thomas
nation america political
rights human democracy
civil
liberty society

vs
engineering
call_of_duty
win cod league match his_life
hit_me_up nfl fucking
creed himself team
xbox youtube.com http:/ world_cup
loves_his ps3 president
black_ops shit my_girlfriend. metal
dick fifa fans my_wife fuckin
teams modern_warfare war
championship government album
arsenal fucked wins
fucking ftw holy_shit nba shave
fuckers lebron halo
haters sake fuck thinks_he his_mind 360 live
shut bullshit fucks beard wii
bitches xbox
outta shit play ps

fight win fighting
won battles
defeat war battle
enemy sword
bands defeated meet
fought victory

online playing
gaming
fifa pc cod
games
playin
tag

pay economy
public tax state cuts
country income taxes
debt government
obama budget health
benefits

a a **a**
*correlation strength*

*relative frequency*

b b **b**
*prevalence in topic*

Schwartz, H. A., Eichstaedt, J. C., Kern, M. L., Dziurzynski, L., Ramones, S. M., Agrawal, M., Shah, A., Kosinski, M., Stillwell, D., Seligman, M. E. P., & Ungar, L. H. (2013). Personality, Gender, and Age in the Language of Social Media: The Open-Vocabulary Approach. In PLOS ONE 8(9).

# Results: Age

# Results: Age

## Predictive Accuracy

| | |
|---|---|
| Openness | language .46 |
| Conscientiousness | language .34 |
| Extraversion | language .43 |
| Agreeableness | language .38 |
| Neuroticism | language .35 |

correlation with questionnaire

Predictive Accuracy

| | language | friend |
|---|---|---|
| Openness | .46 | .25 |
| Conscientiousness | .34 | .30 |
| Extraversion | .43 | .39 |
| Agreeableness | .38 | .30 |
| Neuroticism | .35 | .34 |

correlation with questionnaire

# extraversion

| January-June Year 1 | July-December Year 1 | January-June Year 2 | July-December Year 2 |

Correlations between predictions made at different time points:

.69

.66

.61

correlation between outcome and extraversion

**Model Comparison for Gender Prediction Across Test Sets**

**Model Comparison for Age Prediction Across Test Sets**

Legend: FB_r1k, FB_stratified, blogs_r1k

| Model | FB_r1k | FB_stratified | blogs_r1k |
|-------|--------|---------------|-----------|
| baseline | .000 | .000 | .000 |
| FB | .835 | .801 | .710 |
| Blog | .664 | .657 | .768 |
| FB+Blogs | .831 | .795 | .763 |

Y-axis: Peasron r

X-axis: Model

# Work in Progress:
# Social Media + Medical Record

| Terms searched | Have dx | % with dx and used a term | % without dx and used a term |
|---|---|---|---|
| abdominal pain&stomach pain&belly pain&tummy pain&stomach hurts&belly hurts&tummy hurts&tummyache&stomachache&bellyache | 383 | 21% | 8% *** |
| nausea&vomiting&vomit&throwing up&spitting up&threw up&puke&puked&vomited | 348 | 29% | 22% * |
| headache&migraine&head hurts | 237 | 59% | 46% ** |
| leg hurts&arm hurts&finger hurts&toe hurts | 194 | 3% | 1% |
| uti&urinary tract infection | 160 | 1% | 1% |
| back pain&backache&back hurts | 190 | 15% | 11% |
| cough&coughing&coughed | 156 | 26% | 22% |
| giving birth&gave birth | 188 | 33% | 10% *** |
| anemia | 160 | 2% | 0% ** |
| dizzy&dizziness&vertigo | 127 | 22% | 15% |
| asthma&hard to breathe | 128 | 28% | 7% *** |
| caught a cold&have a cold | 101 | 7% | 4% |
| sore throat&throat hurts | 110 | 24% | 11% *** |
| depression&depressed | 92 | 38% | 30% |

# Clinically Diagnosed Depression

**Discovering Psychological and Health Insights from Social Media Langugae**

H. Andrew Schwartz
andrew.schwartz@cs.stonybrook.edu

Computational Linguistics / Psychology/Health

October 16, 2015
@ SBU SUNYK Seminar

x 20mil.

**1. Individual Analyses**

x 1bil.

**2. Community Analyses**

**Introduction**

# Well-Being

Life Satisfaction (Diener, 1987)
   *In general, how satisfied are you with your life?*

# Why?

UK Survey: *greatest happiness* or *greatest wealth*?
- =>81% - *greatest happiness*

around the world...
- **OECD** set guidelines (2013)
=> UK, France, Bhutan,
Australia, Canada, Mexico, ...

in the US...
- **Gallup** Well-Being Index
- **CDC** Life Satisfaction



OECD Guidelines on Measuring Subjective Well-being



CDC — CENTERS FOR DISEASE CONTROL AND PREVENTION

# Prediction



- tweet language featutes:
  - Lexica: LIWC, PERMA Well-Being
  - LDA Topics: 2000 (Facebook derived)

- tweet language featutes:
  - Lexica: LIWC, PERMA Well-Being
  - LDA Topics: 2000 (Facebook derived)
- controls:
  - demographics (ethnicity, gender, age)
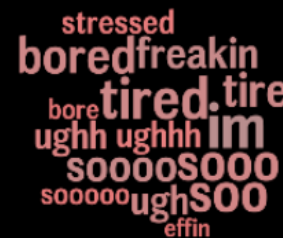  - socio-economics (income, education)

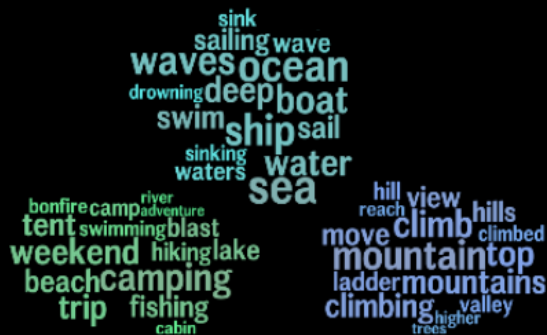- tweet language featutes:
  - Lexica: LIWC, PERMA Well-Being
  - LDA Topics: 2000 (Facebook derived)
- controls:
  - demographics (ethnicity, gender, age)
  - socio-economics (income, education)

- tweet language featutes:
  - Lexica: LIWC, PERMA Well-Being
  - LDA Topics: 2000 (Facebook derived)
- controls:
  - demographics (ethnicity, gender, age)
  - socio-economics (income, education)

# Insight

Schwartz, H. A., Eichstaedt, J. C., Kern, M. L., Dziurzynski, L., Lucas, R. E., Agrawal, M., Park, G. J., Lakshmikanth, S. K., Jha, S., Seligman, M. E. P., & Ungar, L. H. (2013). Characterizing Geographic Variation in Well-Being using Tweets. In *Proceedings of the Seventh International AAAI Conference on Weblogs and Social Media (ICWSM)*. Boston, MA.

# Insight

Schwartz, H. A., Eichstaedt, J. C., Kern, M. L., Dziurzynski, L., Lucas, R. E., Agrawal, M., Park, G. J., Lakshmikanth, S. K., Jha, S., Seligman, M. E. P., & Ungar, L. H. (2013). Characterizing Geographic Variation in Well-Being using Tweets. In *Proceedings of the Seventh International AAAI Conference on Weblogs and Social Media (ICWSM)*. Boston, MA.

**Discovering Psychological and Health Insights from Social Media Langugae**

H. Andrew Schwartz
andrew.schwartz@cs.stonybrook.edu

October 16, 2015
@ SBU SUNYK Seminar
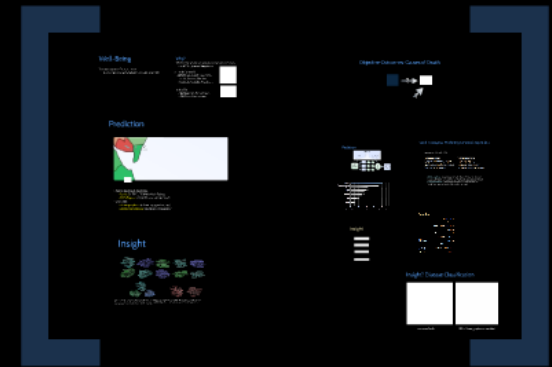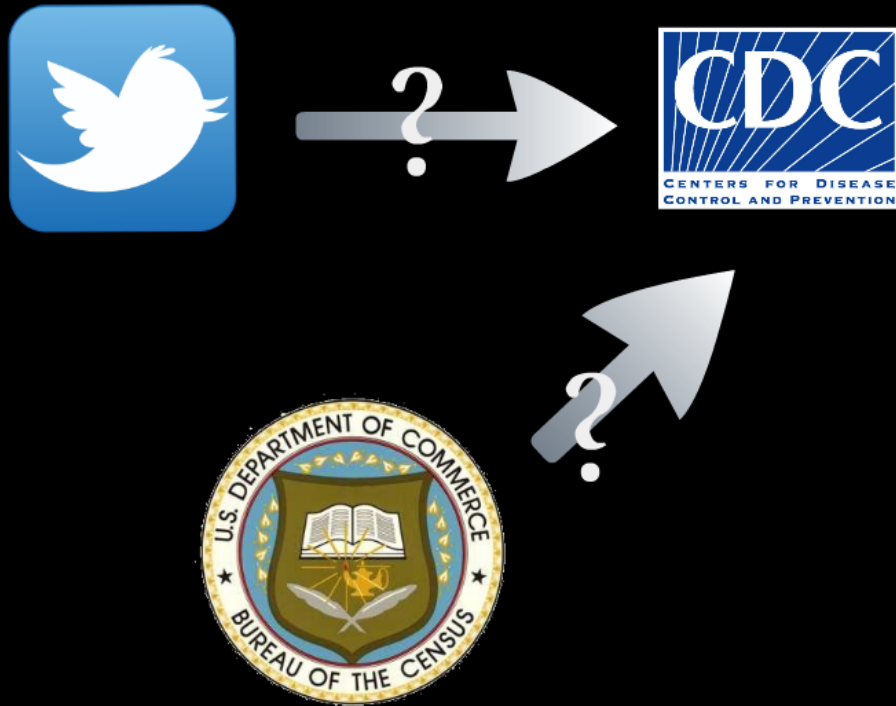
x 20mil.

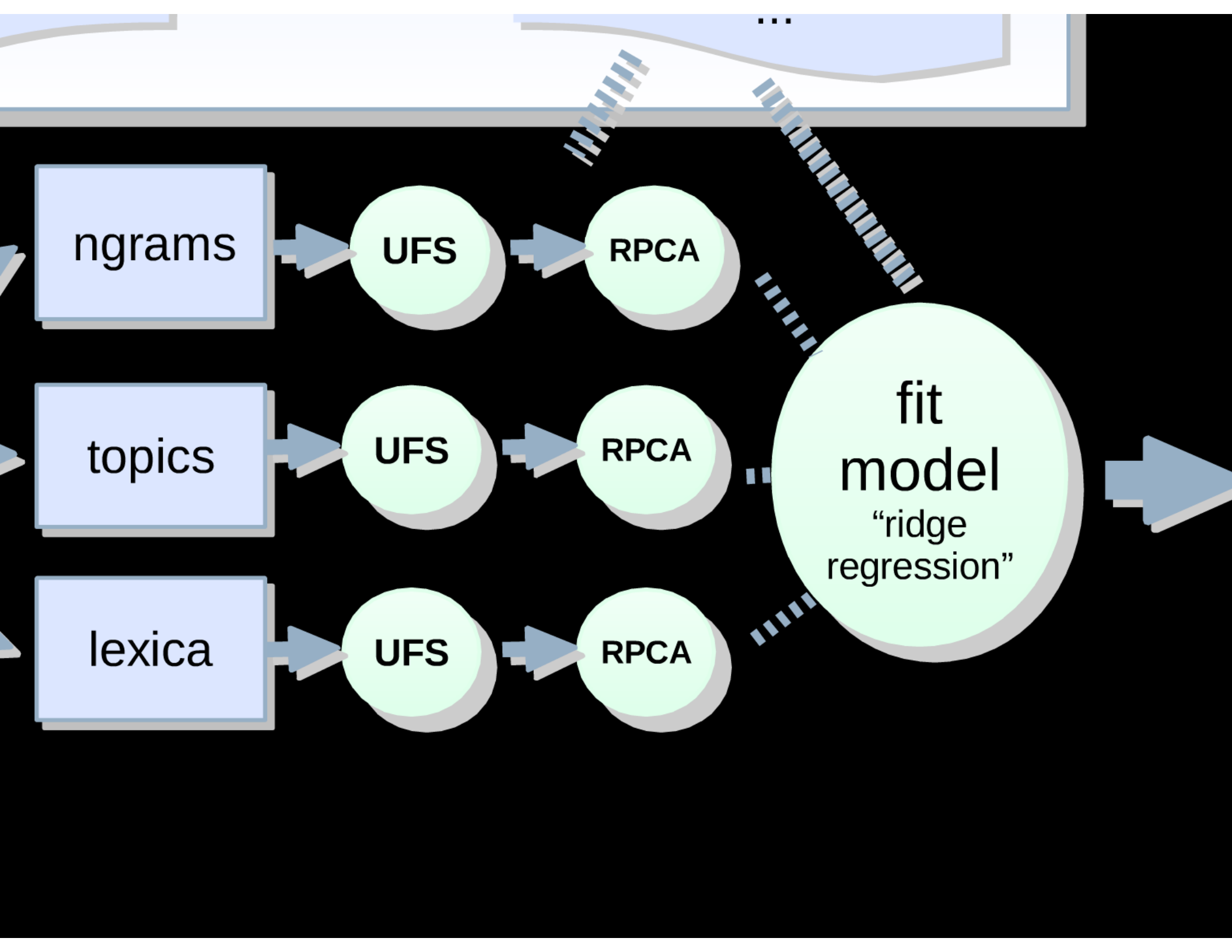**1. Individual Analyses**

x 1bil.

**2. Community Analyses**

**Introduction**

# Objective Outcomes: Causes of Death

# Prediction

**Accuracy of County-Level ACHD Predictions (Pearson r with CDC-reported ACHD)**

# Insight

Higher Status Occupations

public charity
company customer
service entertainment
center community
customers announcement
enemy
rep suggestions
provide services

skills
research engineering
learning education
communication analysis
process
management
business information
development
technology design
marketing

group
leadership attend students
youth meetings student
convention conference board
meeting staff center
council members

dumb *ssholes*
annoying

Anger, Hostility, Aggression

Negative Relationships

Disengagement

# work in progress: Predicting America's top Killers

top causes of death, 2010

1. Diseases of heart
2. Malignant neoplasms (cancers)
3. Chronic lower respiratory diseases
4. Cerebrovascular diseases (strokes)
5. Accidents, unintentional

6. Alzheimer's disease
7. Diabetes mellitus
8. Nephritis (kidney diseases)
9. Influenza and pneumonia
10. Intentional self-harm (suicide)

demographics: percentage female, black, Hispanic, foreign born, and married residents, as well as the population density

socioeconomics: percentage completed high school / a bachelors, unemployed, log median income

# Results

# Insight? Disease Classification



no controls

SES + Demographic controlled

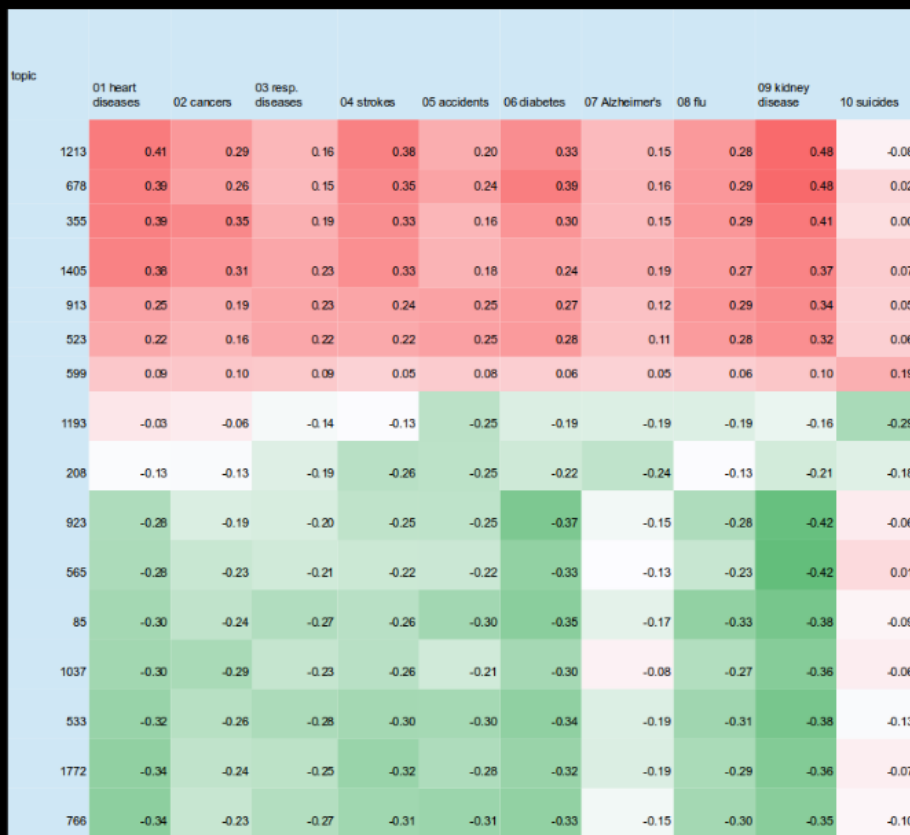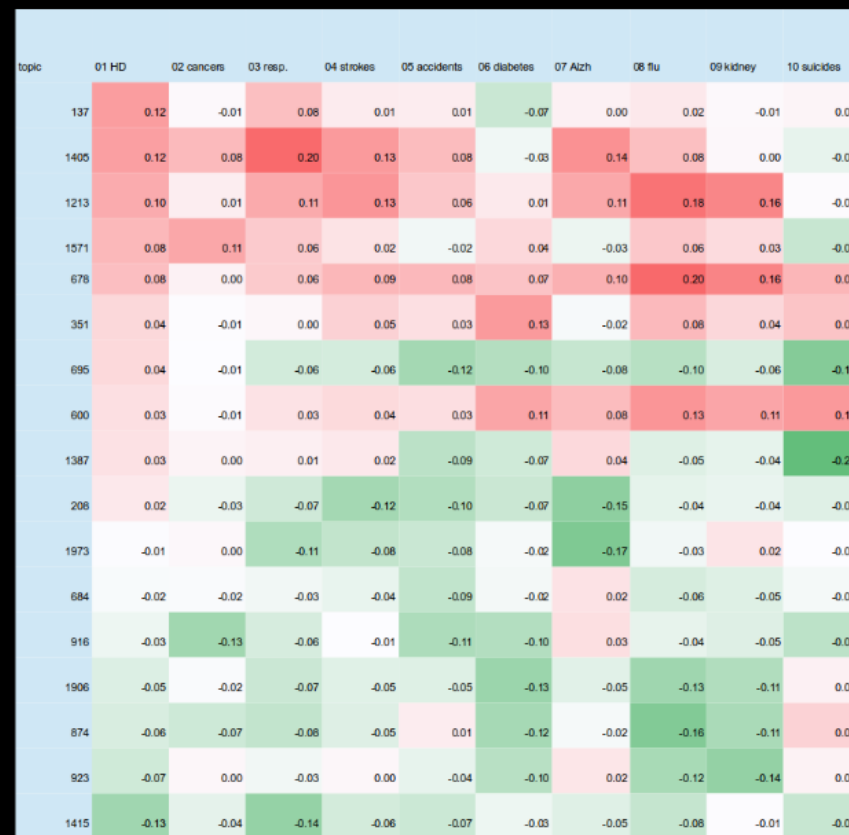| topic | 01 HD | 02 cancers | 03 resp. | 04 strokes | 05 accidents | 06 diabetes | 07 Alzh | 08 flu | 09 kidney | 10 suicides |
|---|---|---|---|---|---|---|---|---|---|---|
| 137 | 0.12 | -0.01 | 0.08 | 0.01 | 0.01 | -0.07 | 0.00 | 0.02 | -0.01 | 0.00 |
| 1405 | 0.12 | 0.08 | 0.20 | 0.13 | 0.08 | -0.03 | 0.14 | 0.08 | 0.00 | -0.04 |
| 1213 | 0.10 | 0.01 | 0.11 | 0.13 | 0.06 | 0.01 | 0.11 | 0.18 | 0.16 | -0.01 |
| 1571 | 0.08 | 0.11 | 0.06 | 0.02 | -0.02 | 0.04 | -0.03 | 0.06 | 0.03 | -0.08 |
| 678 | 0.08 | 0.00 | 0.06 | 0.09 | 0.08 | 0.07 | 0.10 | 0.20 | 0.16 | 0.09 |
| 351 | 0.04 | -0.01 | 0.00 | 0.05 | 0.03 | 0.13 | -0.02 | 0.08 | 0.04 | 0.07 |
| 695 | 0.04 | -0.01 | -0.06 | -0.06 | -0.12 | -0.10 | -0.08 | -0.10 | -0.06 | -0.16 |
| 600 | 0.03 | -0.01 | 0.03 | 0.04 | 0.03 | 0.11 | 0.08 | 0.13 | 0.11 | 0.13 |
| 1387 | 0.03 | 0.00 | 0.01 | 0.02 | -0.09 | -0.07 | 0.04 | -0.05 | -0.04 | -0.20 |
| 208 | 0.02 | -0.03 | -0.07 | -0.12 | -0.10 | -0.07 | -0.15 | -0.04 | -0.04 | -0.05 |
| 1973 | -0.01 | 0.00 | -0.11 | -0.08 | -0.08 | -0.02 | -0.17 | -0.03 | 0.02 | -0.01 |
| 684 | -0.02 | -0.02 | -0.03 | -0.04 | -0.09 | -0.02 | 0.02 | -0.06 | -0.05 | -0.02 |
| 916 | -0.03 | -0.13 | -0.06 | -0.01 | -0.11 | -0.10 | 0.03 | -0.04 | -0.05 | -0.09 |
| 1906 | -0.05 | -0.02 | -0.07 | -0.05 | -0.05 | -0.13 | -0.05 | -0.13 | -0.11 | 0.01 |
| 874 | -0.06 | -0.07 | -0.08 | -0.05 | 0.01 | -0.12 | -0.02 | -0.16 | -0.11 | 0.05 |
| 923 | -0.07 | 0.00 | -0.03 | 0.00 | -0.04 | -0.10 | 0.02 | -0.12 | -0.14 | 0.00 |
| 1415 | -0.13 | -0.04 | -0.14 | -0.06 | -0.07 | -0.03 | -0.05 | -0.08 | -0.01 | -0.08 |

Partial column at left edge:

| 10 suicides |
|---|
| -0.08 |
| 0.02 |
| 0.00 |
| 0.07 |
| 0.05 |
| 0.06 |
| 0.19 |
| -0.29 |
| -0.18 |
| -0.06 |
| 0.01 |
| -0.09 |
| -0.06 |
| -0.13 |
| -0.07 |
| -0.10 |

| topic | 01 heart diseases | 02 cancers | 03 resp. diseases | 04 strokes | 05 accidents | 06 diabetes | 07 Alzheimer's | 08 flu | 09 kidney disease | 10 suicides |
|---|---|---|---|---|---|---|---|---|---|---|
| 1213 | 0.41 | 0.29 | 0.16 | 0.38 | 0.20 | 0.33 | 0.15 | 0.28 | 0.48 | -0.08 |
| 678 | 0.39 | 0.26 | 0.15 | 0.35 | 0.24 | 0.39 | 0.16 | 0.29 | 0.48 | 0.02 |
| 355 | 0.39 | 0.35 | 0.19 | 0.33 | 0.16 | 0.30 | 0.15 | 0.29 | 0.41 | 0.00 |
| 1405 | 0.38 | 0.31 | 0.23 | 0.33 | 0.18 | 0.24 | 0.19 | 0.27 | 0.37 | 0.07 |
| 913 | 0.25 | 0.19 | 0.23 | 0.24 | 0.25 | 0.27 | 0.12 | 0.29 | 0.34 | 0.05 |
| 523 | 0.22 | 0.16 | 0.22 | 0.22 | 0.25 | 0.28 | 0.11 | 0.28 | 0.32 | 0.06 |
| 599 | 0.09 | 0.10 | 0.09 | 0.05 | 0.08 | 0.06 | 0.05 | 0.06 | 0.10 | 0.19 |
| 1193 | -0.03 | -0.06 | -0.14 | -0.13 | -0.25 | -0.19 | -0.19 | -0.19 | -0.16 | -0.29 |
| 208 | -0.13 | -0.13 | -0.19 | -0.26 | -0.25 | -0.22 | -0.24 | -0.13 | -0.21 | -0.18 |
| 923 | -0.28 | -0.19 | -0.20 | -0.25 | -0.25 | -0.37 | -0.15 | -0.28 | -0.42 | -0.06 |
| 565 | -0.28 | -0.23 | -0.21 | -0.22 | -0.22 | -0.33 | -0.13 | -0.23 | -0.42 | 0.01 |
| 85 | -0.30 | -0.24 | -0.27 | -0.26 | -0.30 | -0.35 | -0.17 | -0.33 | -0.38 | -0.09 |
| 1037 | -0.30 | -0.29 | -0.23 | -0.26 | -0.21 | -0.30 | -0.08 | -0.27 | -0.36 | -0.06 |
| 533 | -0.32 | -0.26 | -0.28 | -0.30 | -0.30 | -0.34 | -0.19 | -0.31 | -0.38 | -0.13 |
| 1772 | -0.34 | -0.24 | -0.25 | -0.32 | -0.28 | -0.32 | -0.19 | -0.29 | -0.36 | -0.07 |
| 766 | -0.34 | -0.23 | -0.27 | -0.31 | -0.31 | -0.33 | -0.15 | -0.30 | -0.35 | -0.10 |

**Discovering Psychological and Health Insights from Social Media Langugae**

H. Andrew Schwartz
andrew.schwartz@cs.stonybrook.edu

Computational Linguistics / Psychology/Health

October 16, 2015
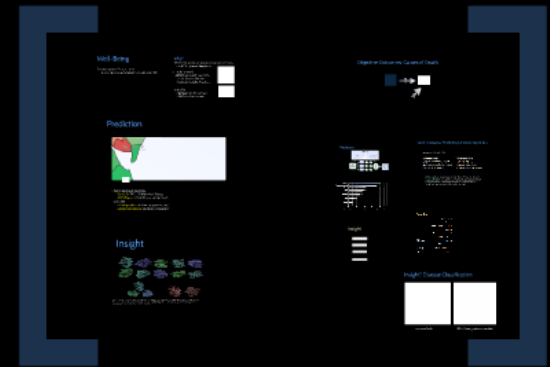@ SBU SUNYK Seminar

x 20mil.

**1. Individual Analyses**

x 1bil.

**2. Community Analyses**

**Introduction**

**Discovering Psychological and Health Insights from Social Media Langugae**

H. Andrew Schwartz
andrew.schwartz@cs.stonybrook.edu

Computational Linguistics | Psychology/Health

October 16, 2015
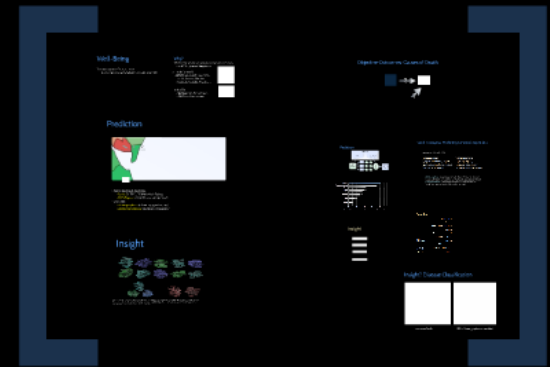@ SBU SUNYK Seminar

x 20mil.

**1. Individual Analyses**

x 1bil.

**2. Community Analyses**

**Introduction**

# Social Media

350m tweets/day

4b messages/day

...the largest dataset of who we are

# Interdisciplinary Research Questions

What psychological factors emerge in language as drivers of health and well-being?

Can we predict disease risk and recovery from language use?

To what extent can language analyses replace and extend traditional psychological asessement

# ...the largest dataset of who we are

**NLP can foster data-driven human discovery at unprecedented scale.**

# Thank You

andrew.schwartz@cs.stonybrook.edu