

# Machine Learning on Social Media Textual Data for Predicting Psychological and Health outcomes

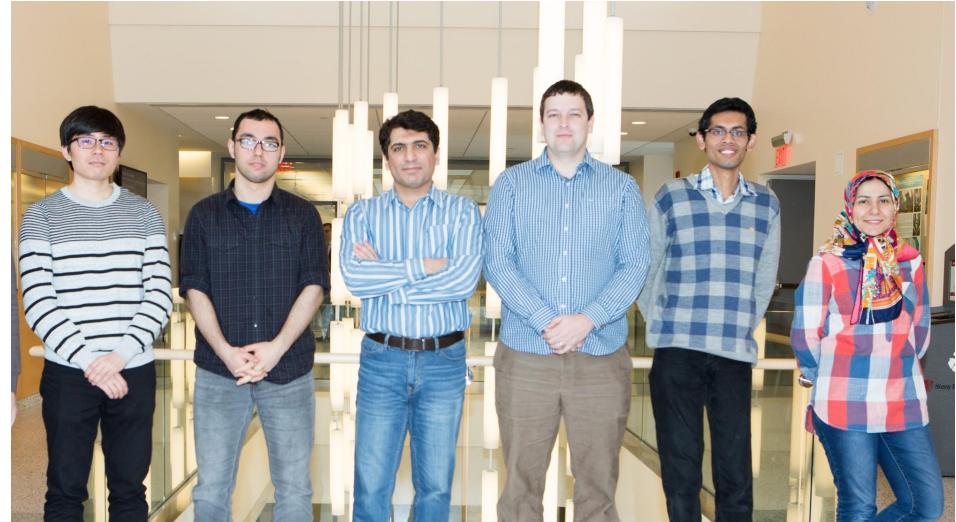
H. Andrew Schwartz



Stony Brook University  
Human Language Analysis Lab

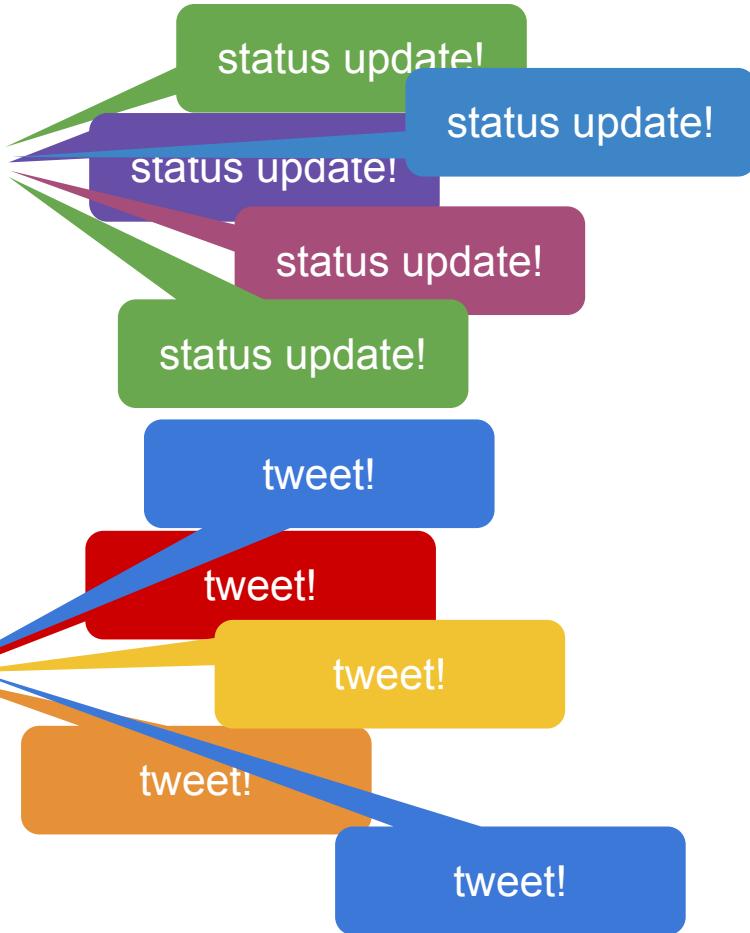


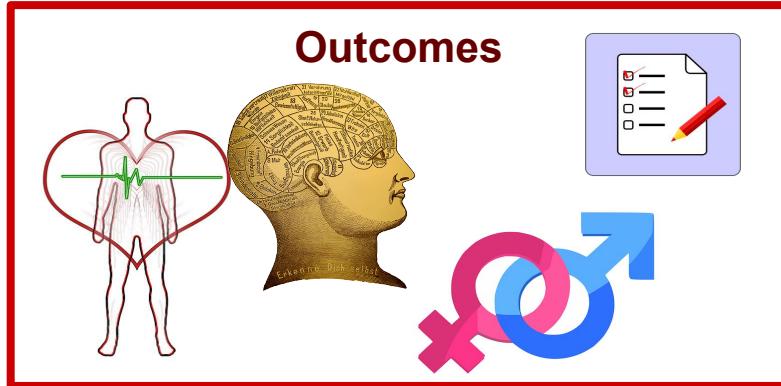
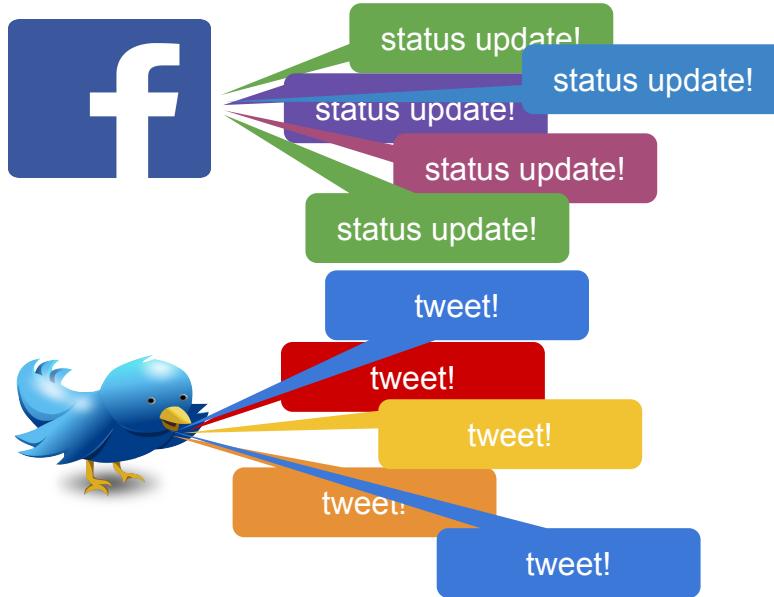
Stony Brook University  
Human Language Analysis Lab

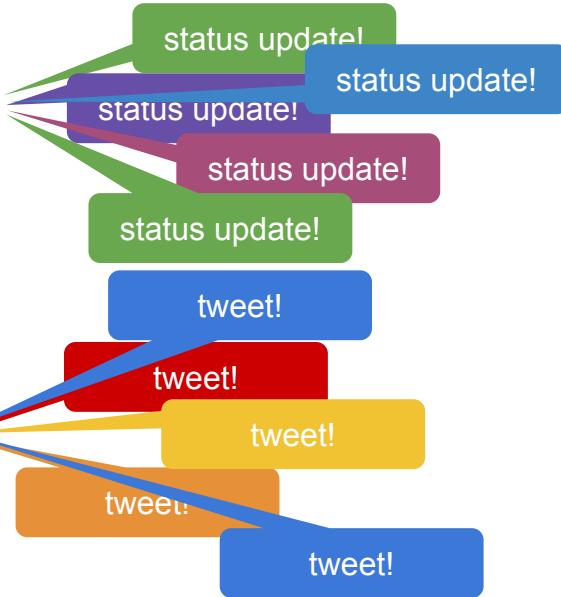


## Penn World Well-Being Project

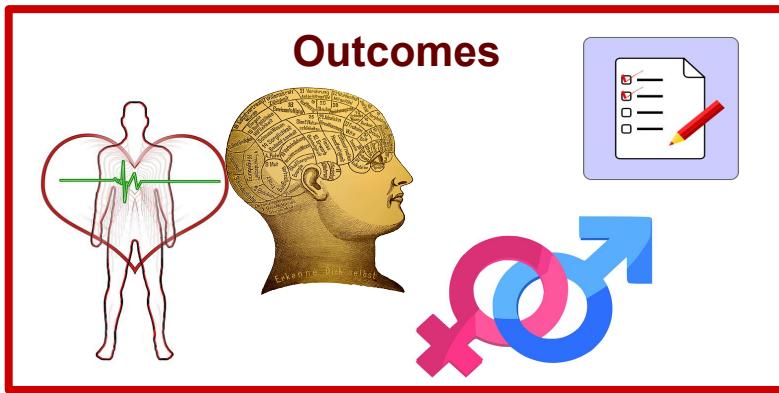


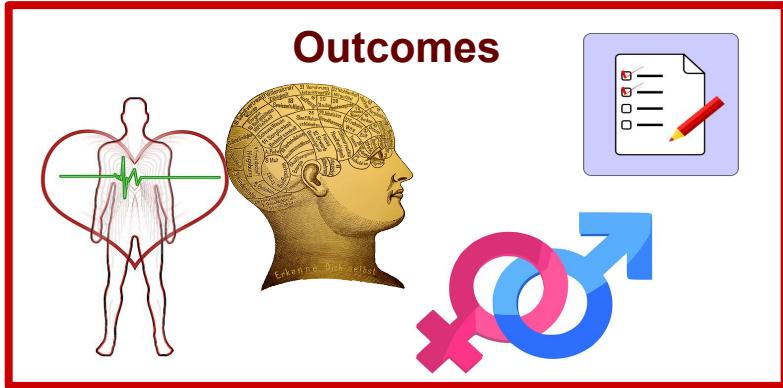
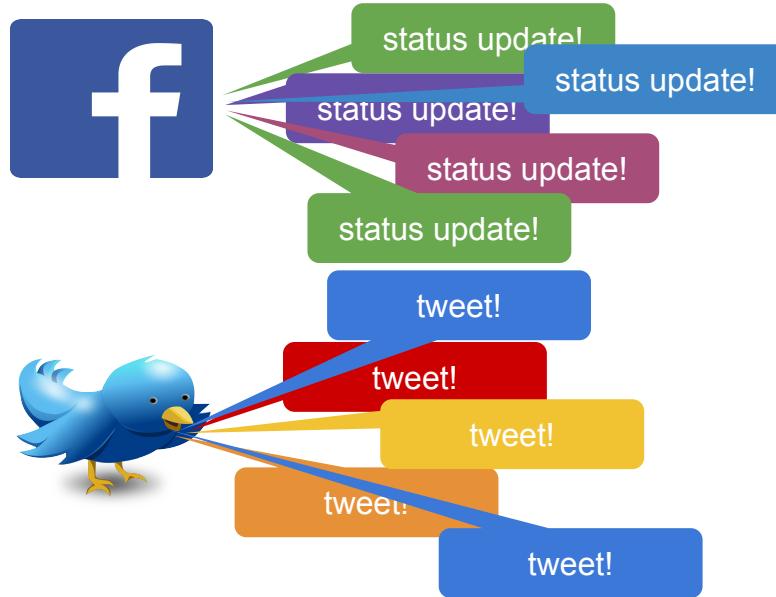






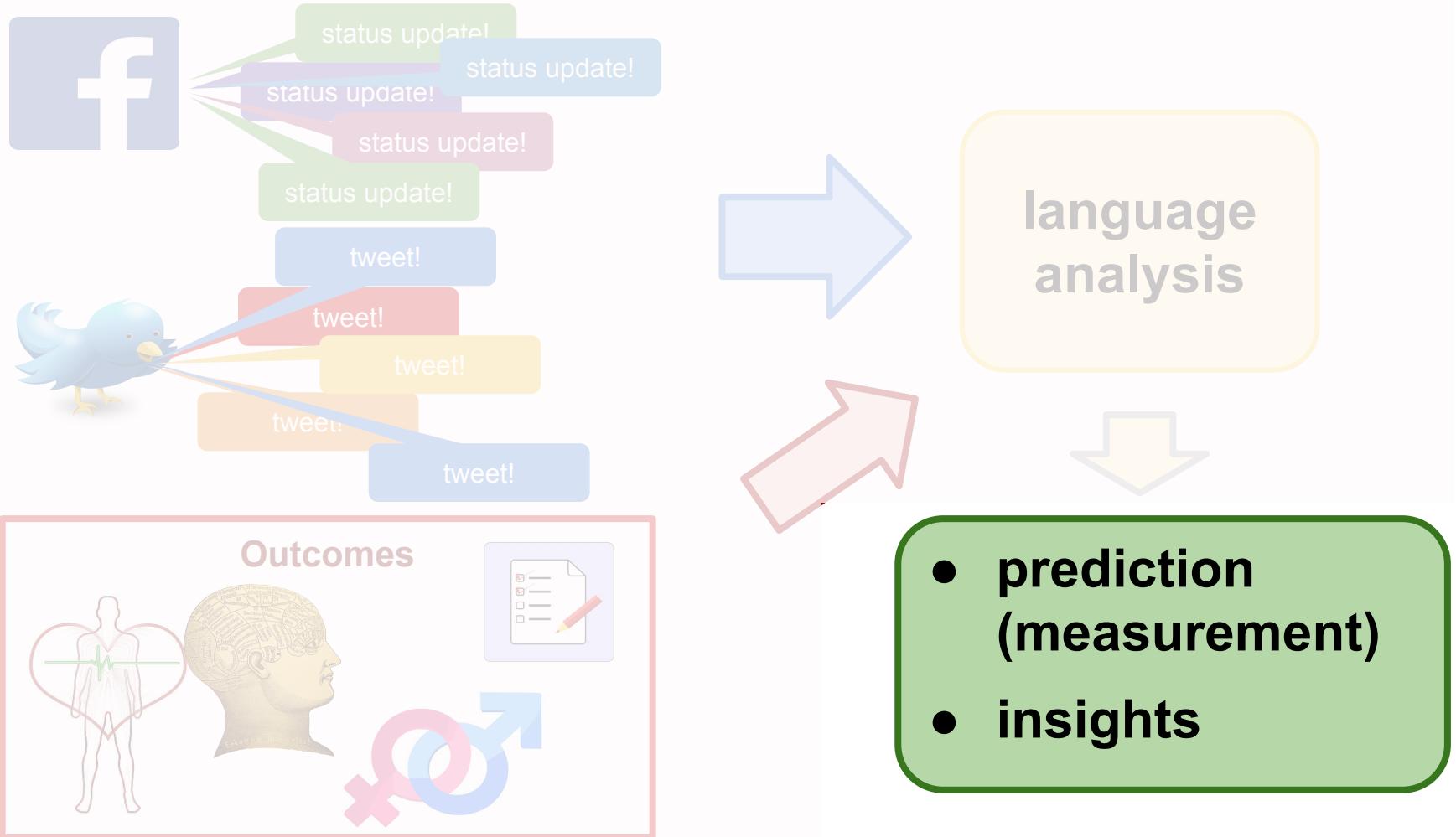
**language  
analysis**



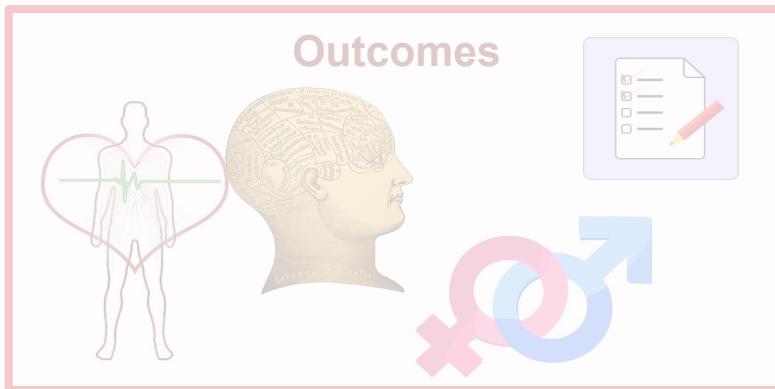
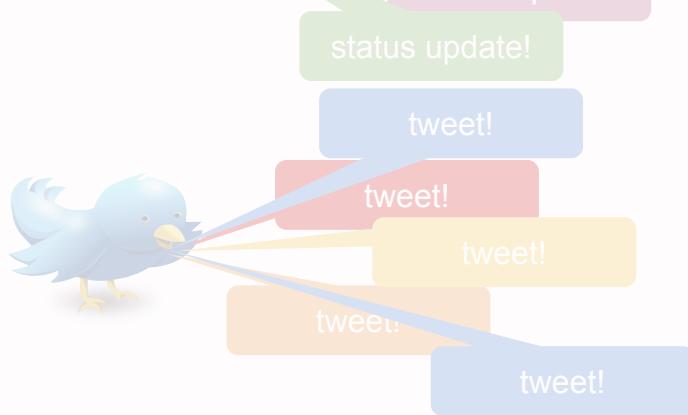


- 
- A green rounded rectangle contains three bullet points:
- prediction (measurement)
  - insights

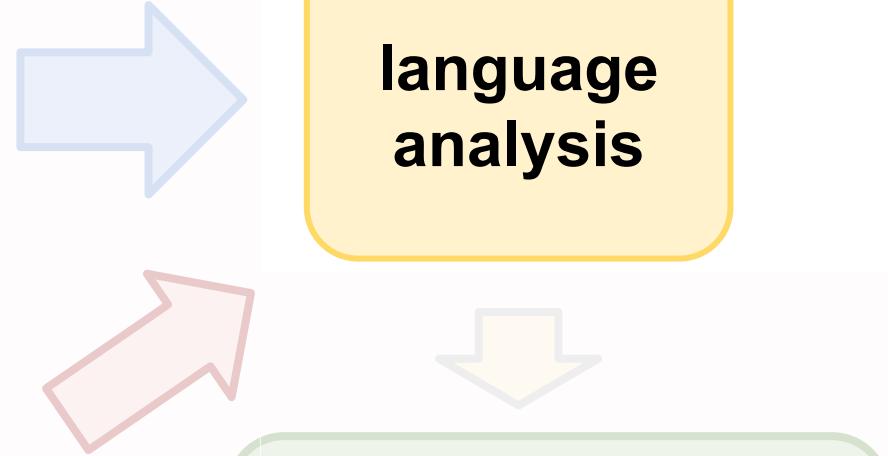
# PART I: What are the possibilities?



# PART I: What?



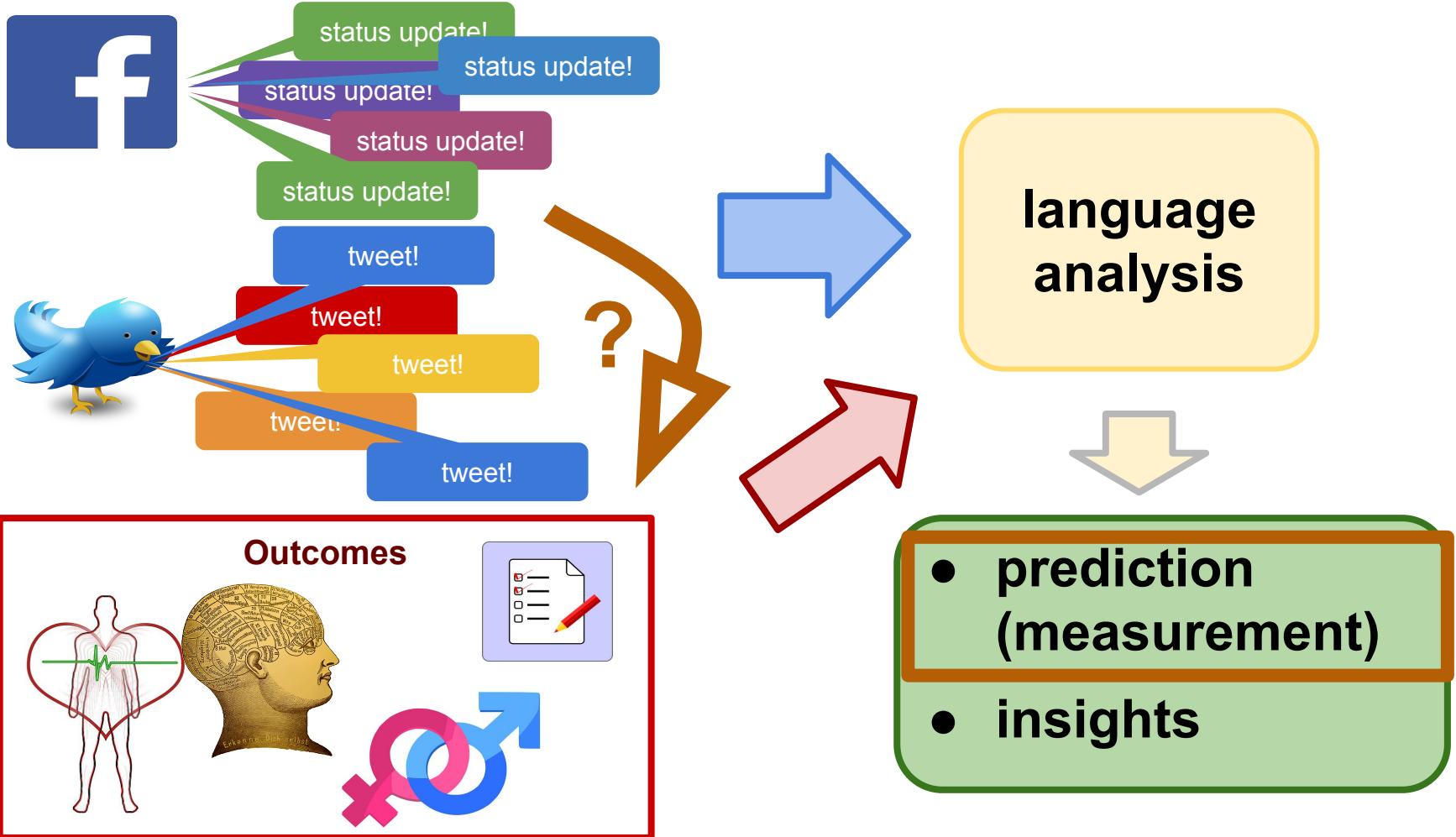
# PART II: How?



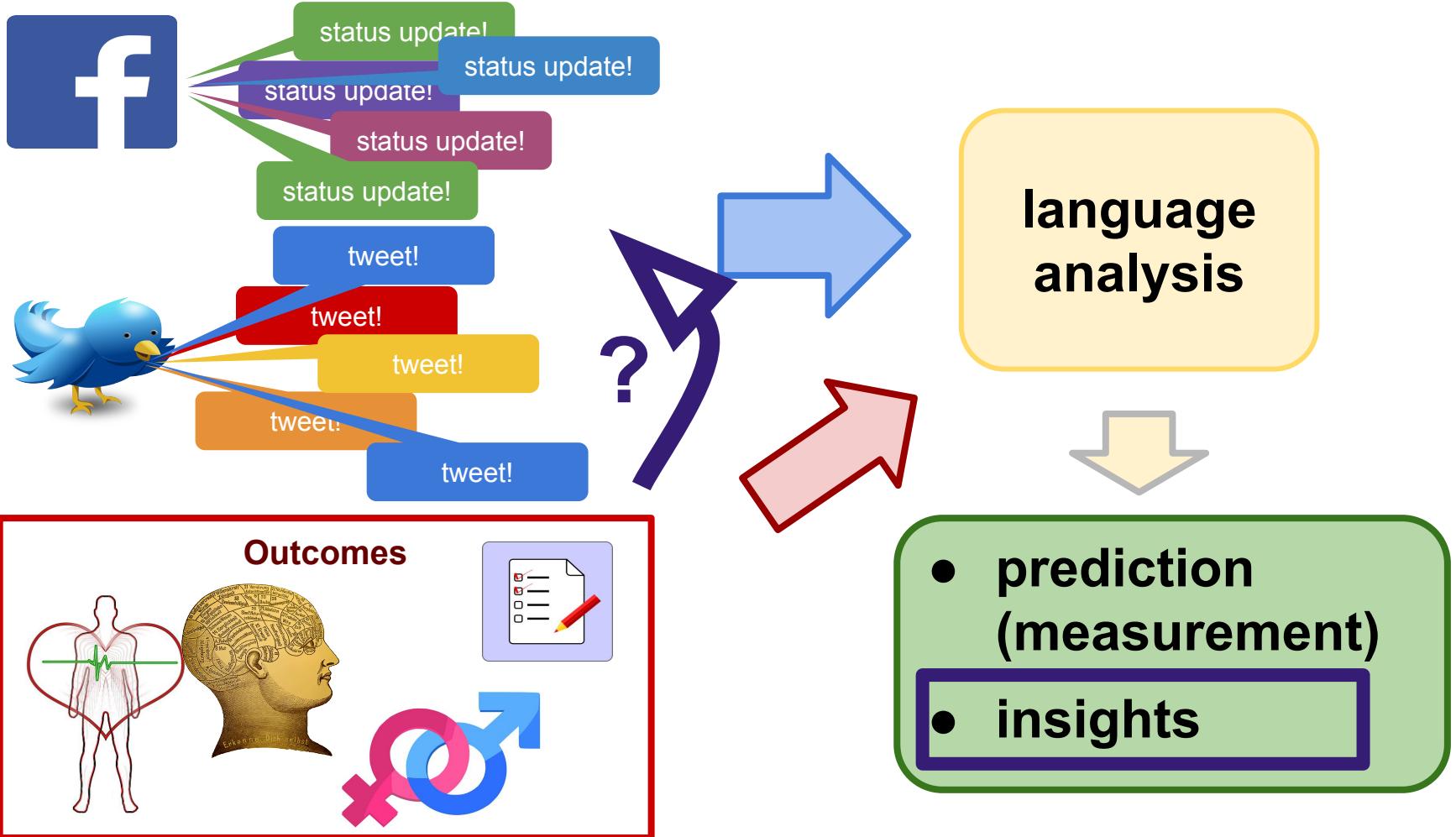
**language  
analysis**

- prediction (measurement)
- insights

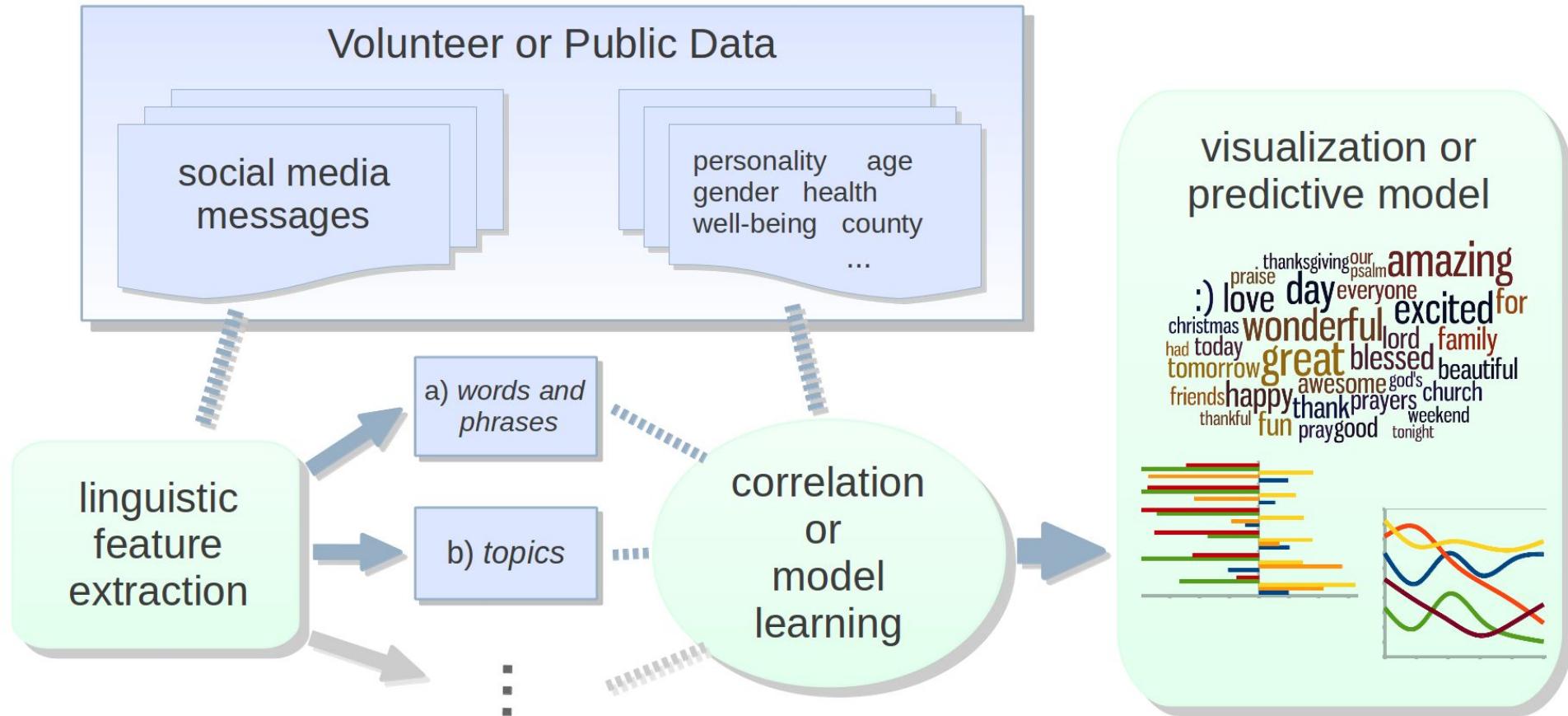
# PART I: What?



# PART I: What?



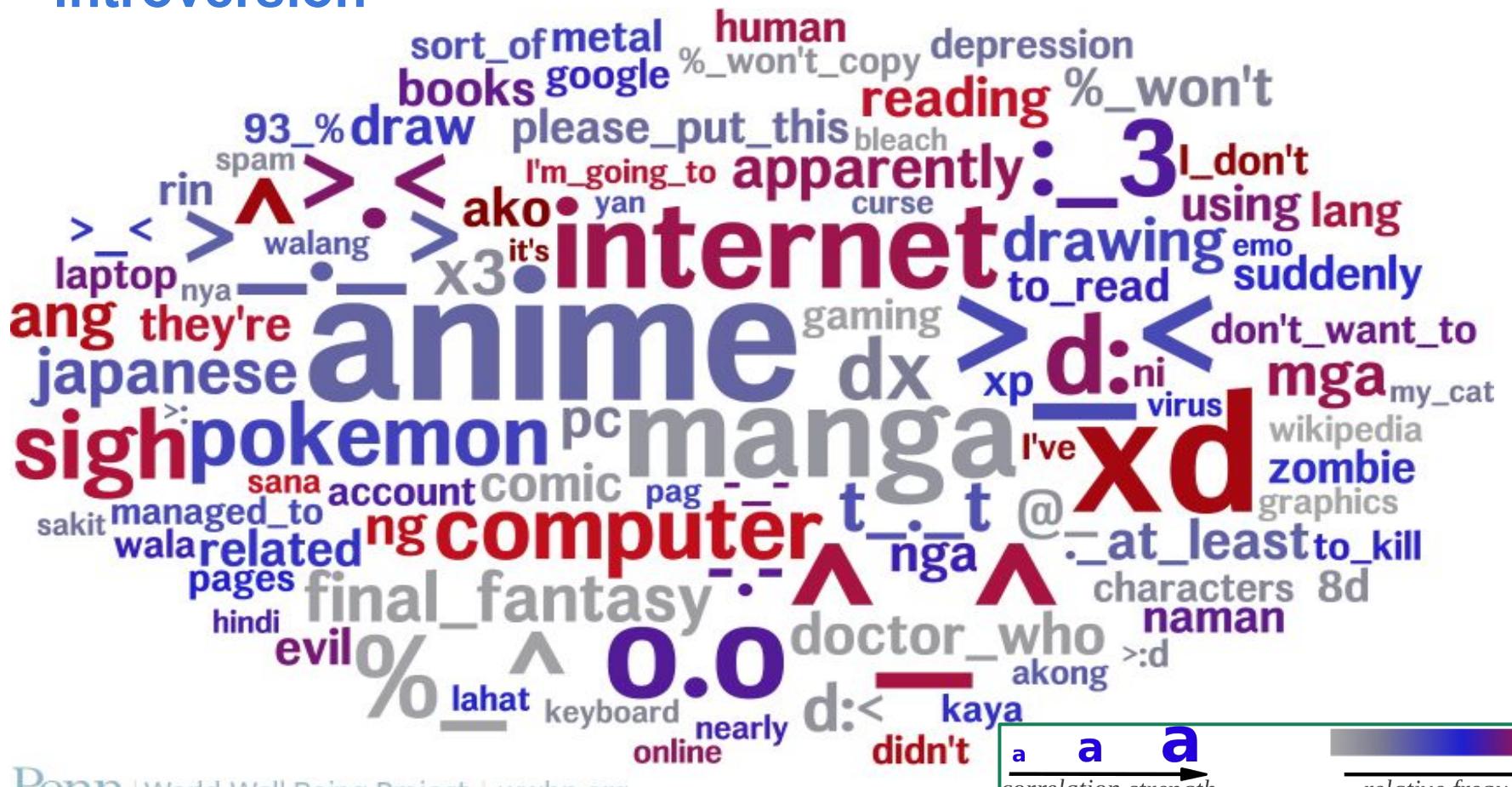
# General and Flexible Framework



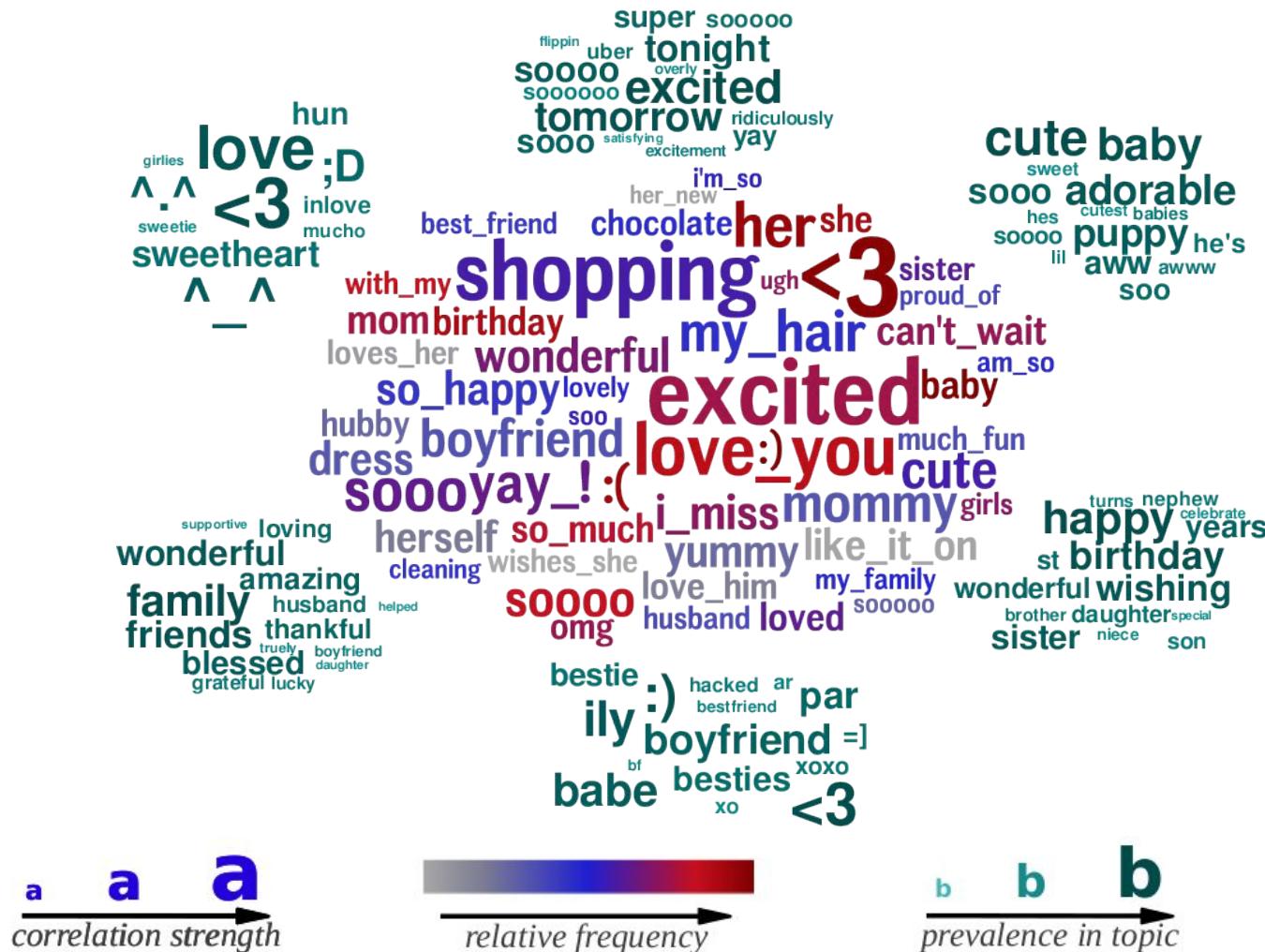
**extraversion** -- sociable, assertive, active, energetic, talkative, outgoing



# introversion



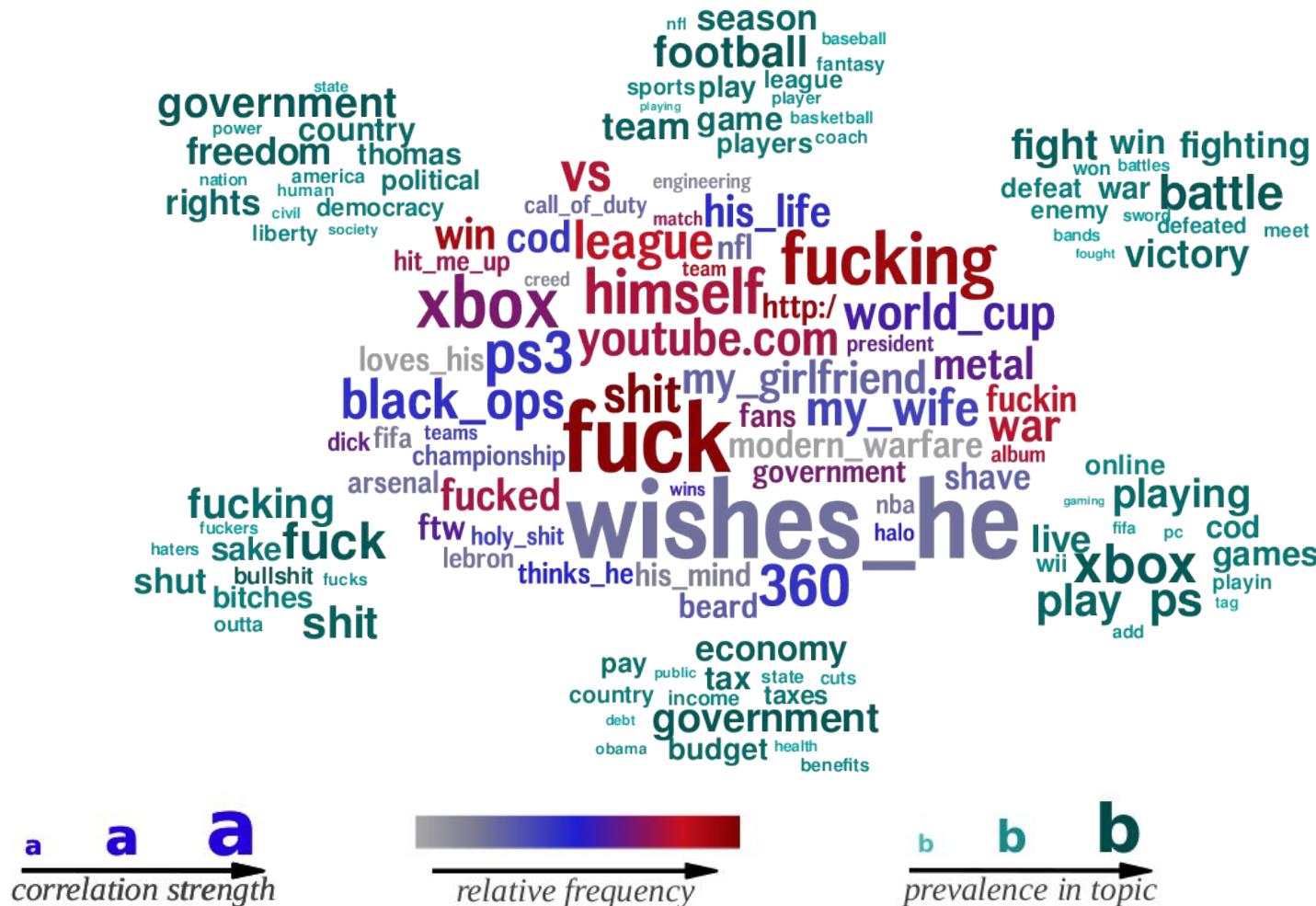
# Gender



# Gender

**Explicit Language Warning...**

# Gender

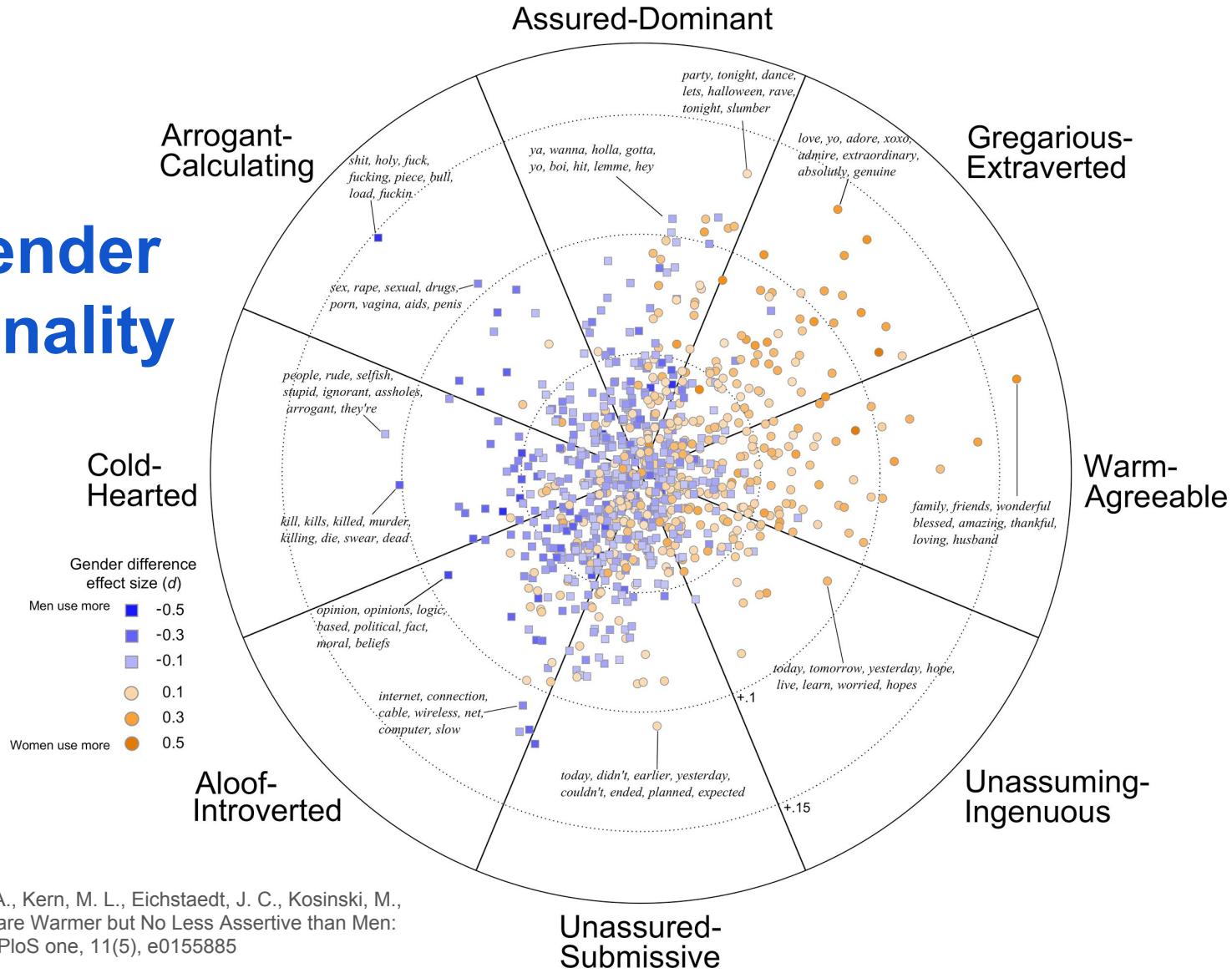


# Neuroticism

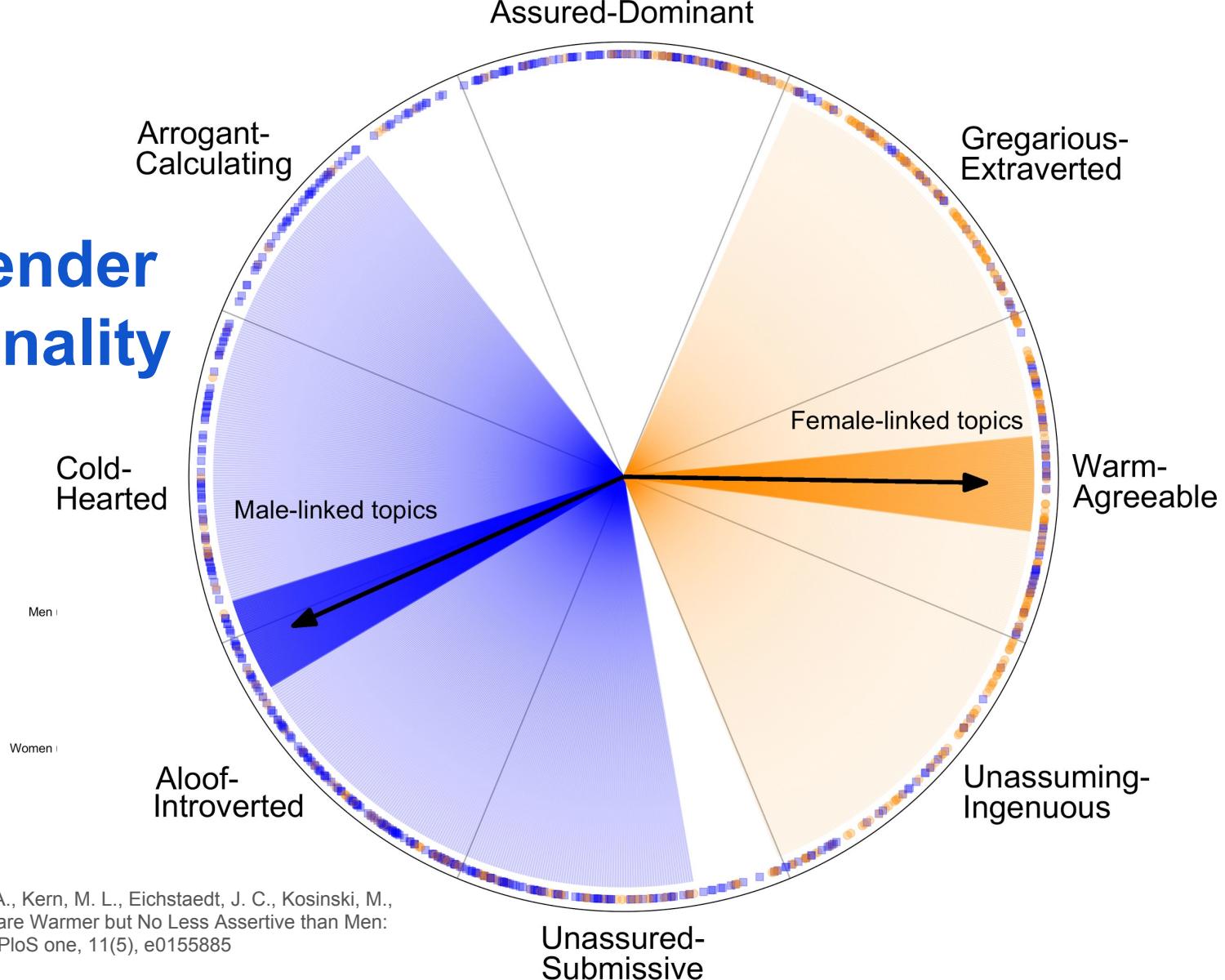
Neuroticism

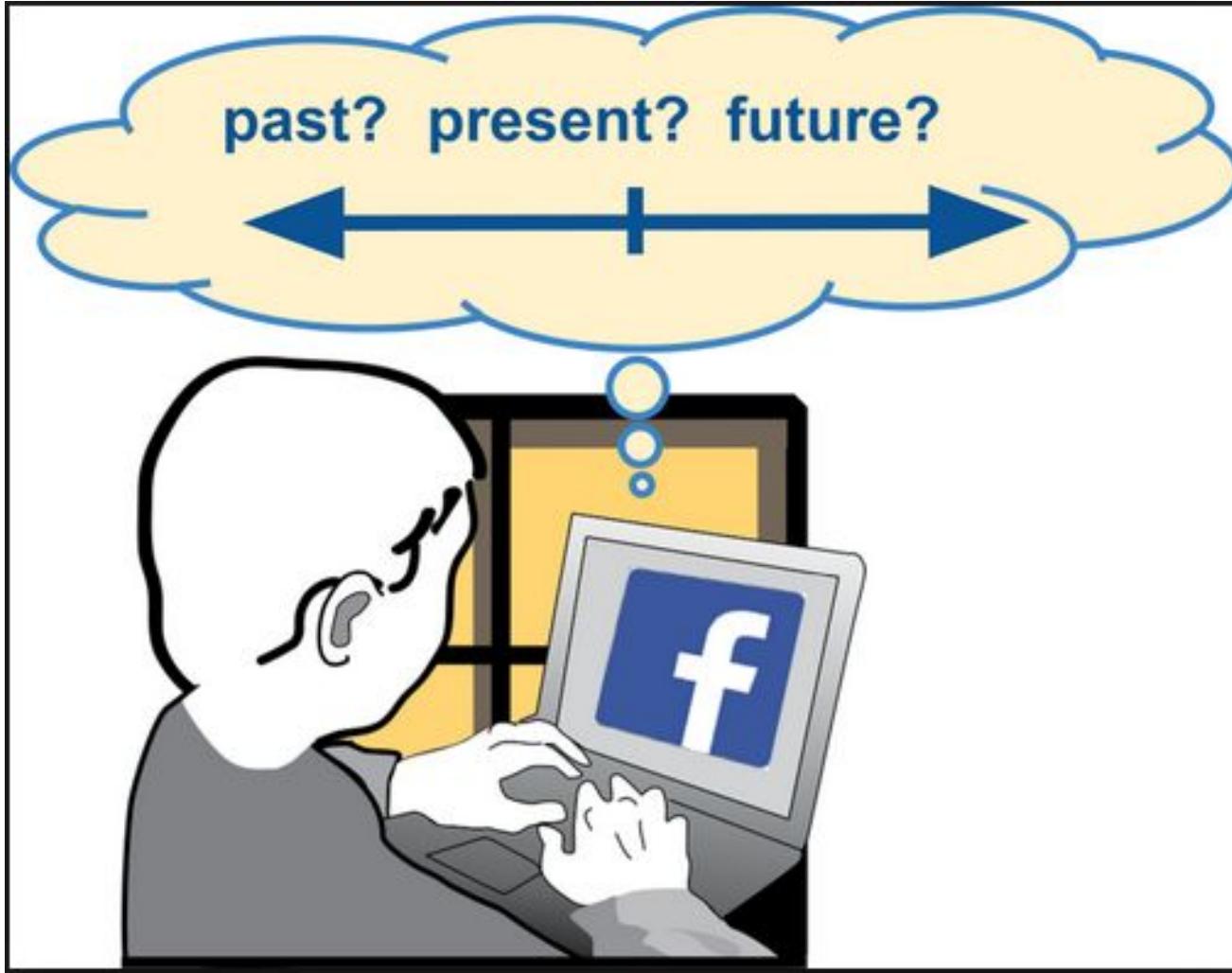


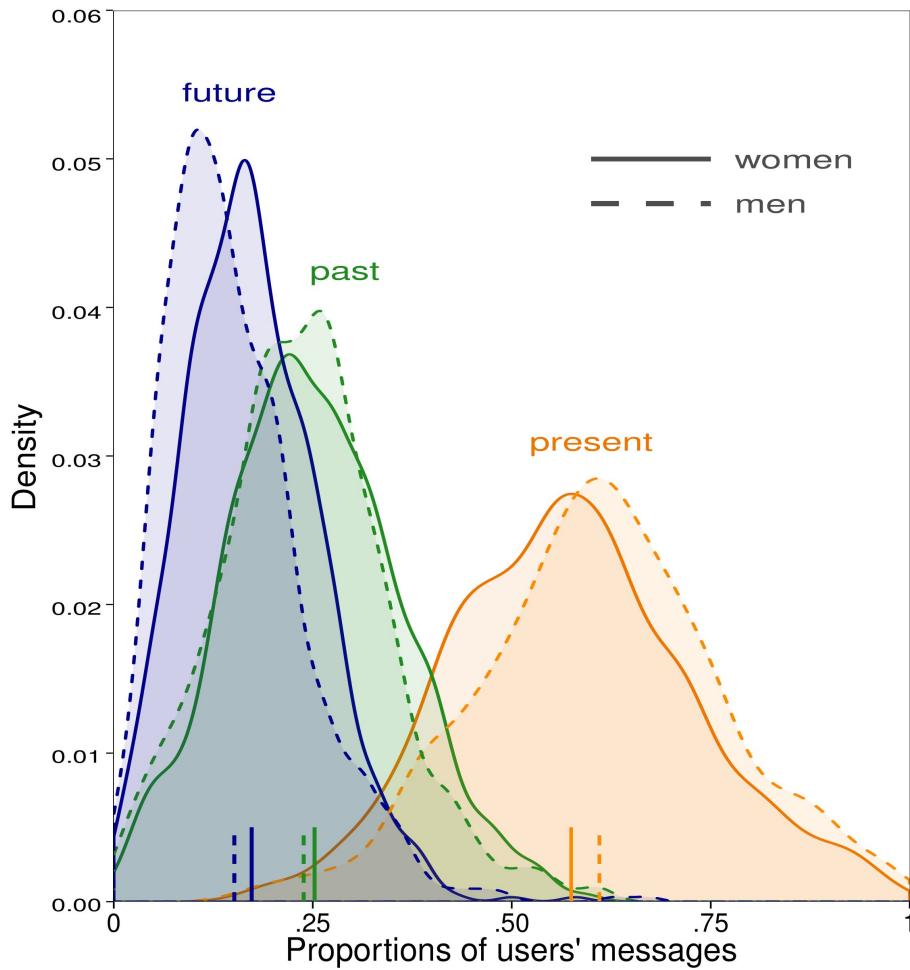
# Putting Gender and Personality Together



# Putting Gender and Personality Together







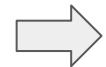
1.3m status updates  
from  
4,833 Participants  
(3,240 AG Stratified)

Schwartz, H. A., Park, G., Sap, M., Weingarten, E., Eichstaedt, J., Kern, M., Stillwell, D., Kosinski, M., Berger, J., Seligman, M., & Ungar, L. (2015). Extracting Human Temporal Orientation from Facebook Language. *NAACL-2015: Conference of the North American Chapter of the Association for Computational Linguistics*.

past

present

future → Language



**Questionnaire**  
IPIP 100 item domains



past



present



future

- complete tasks successfully
- avoid philosophical discussions
- carry out my plans
- finish what I start
- make plans and stick to them

r

past

present

future

complete tasks successfully  
avoid philosophical discussions  
carry out my plans  
finish what I start  
make plans and stick to them



past



present



future

do not like art

would describe my experiences as somewhat dull

am easy to satisfy

rarely lose my composure

don't like to draw attention to myself

past

present

future

do not like art

would describe my experiences as somewhat dull

am easy to satisfy

rarely lose my composure

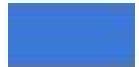
don't like to draw attention to myself



past



present



future

cut others to pieces

can say things beautifully

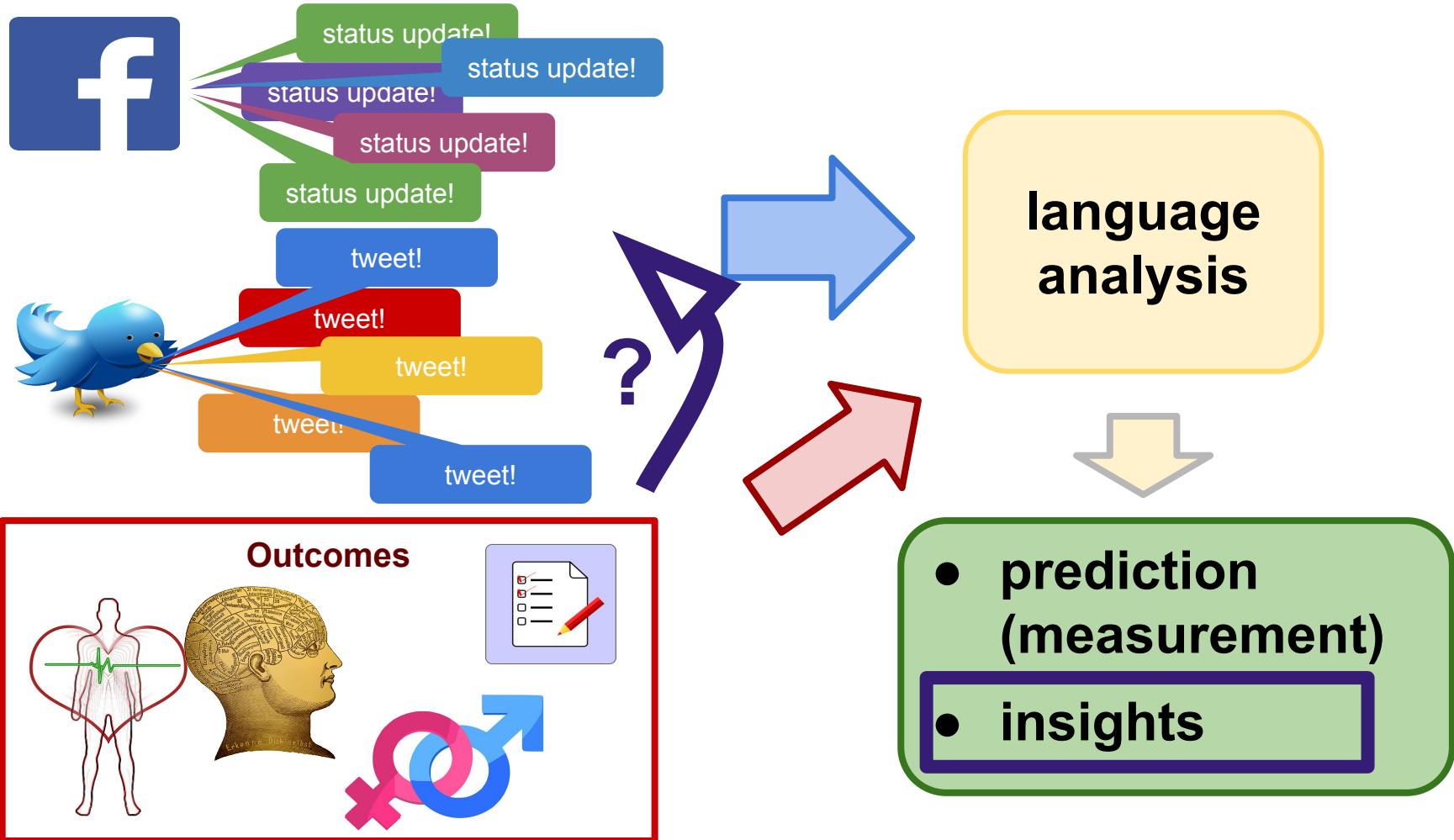
don't put my mind on the task at hand

have frequent mood swings

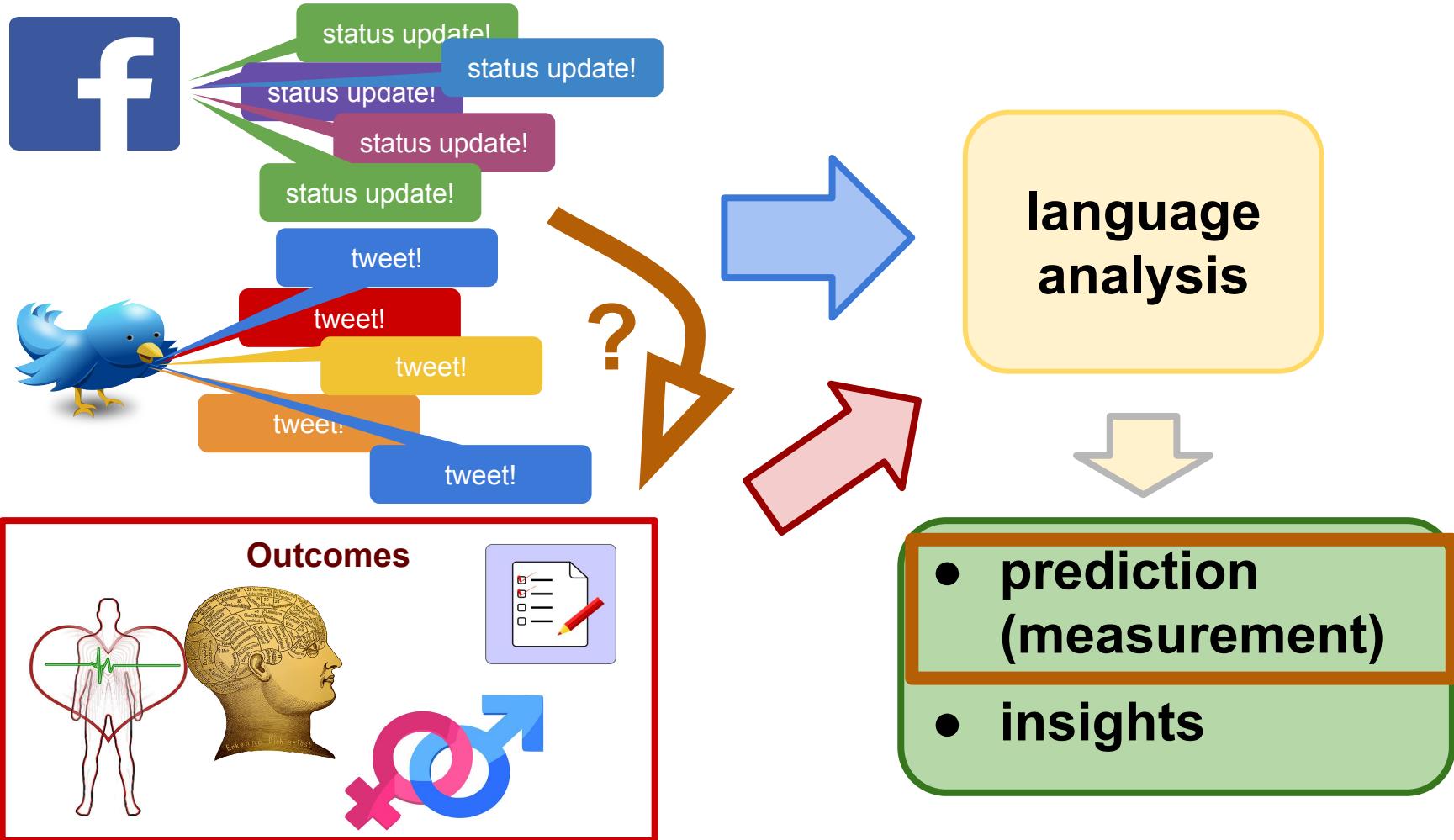
am hard to get to know

Park, G., Schwartz, H.A., Sap, M., Kern, M.L., Weingarten, E., Eichstaedt, J.C., Berger, J., Stillwell, D.J., Kosinski, M., Ungar, L.H. & Seligman, M.E. (2015). Living in the Past, Present, and Future: Measuring Temporal Orientation with Language. *Journal of personality*.

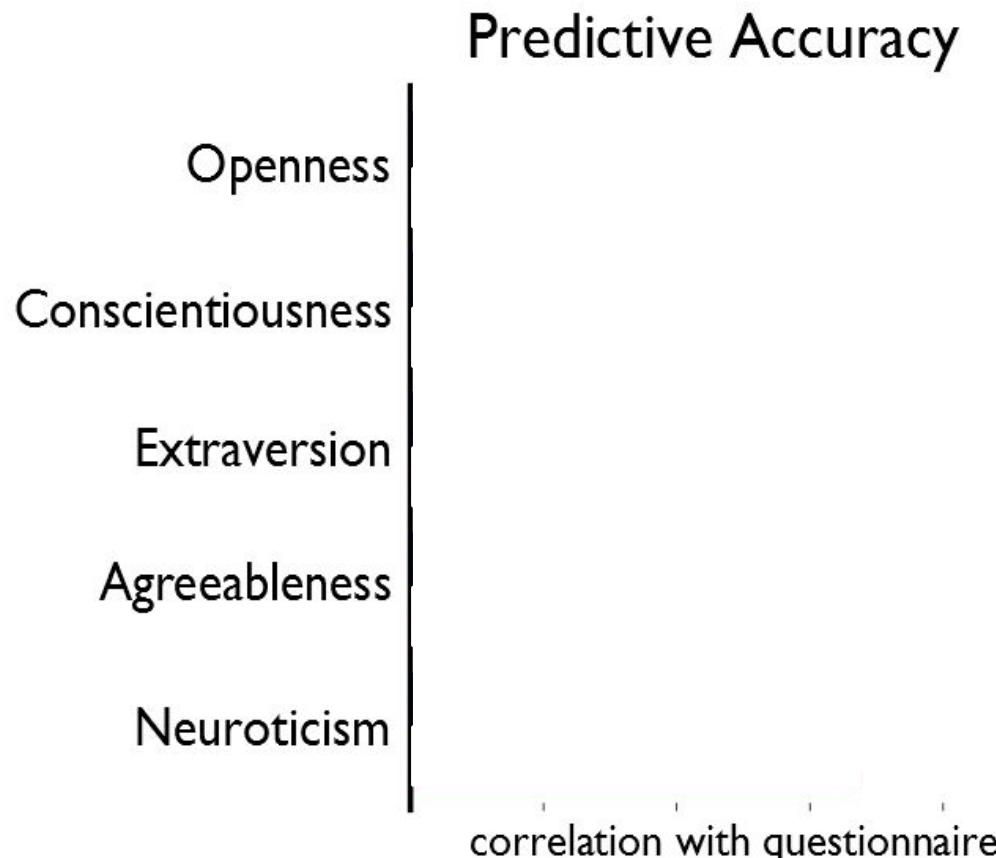
# PART I: What are the possibilities?



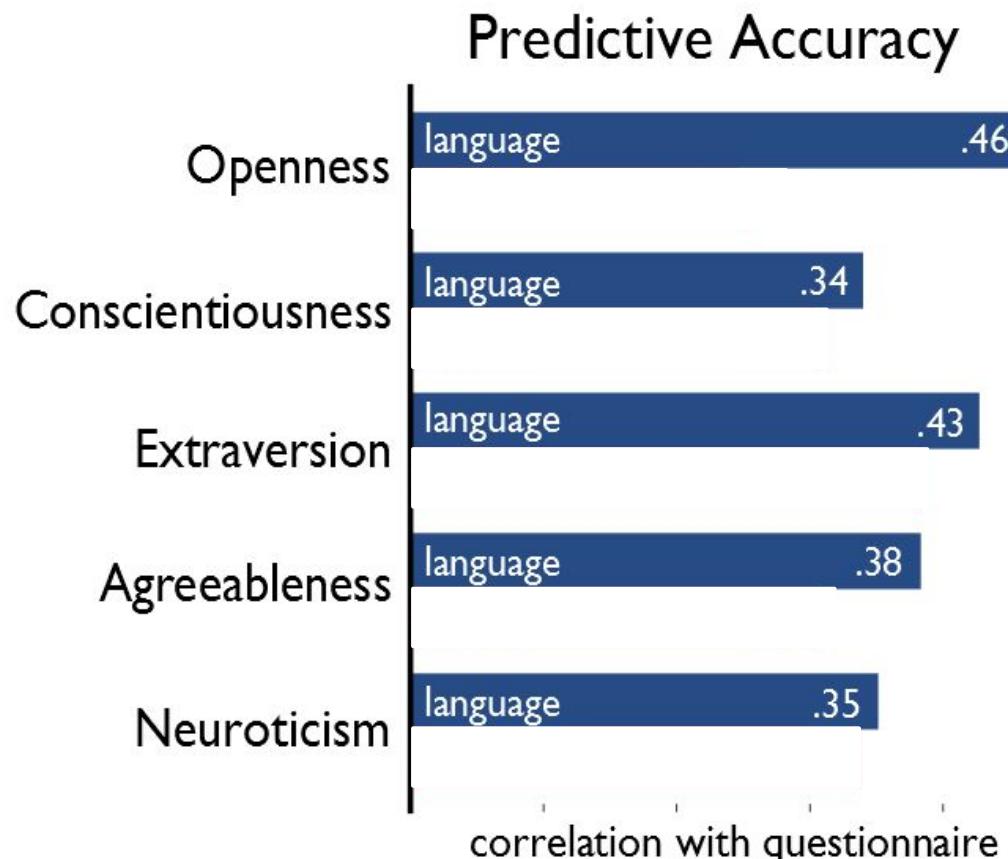
# PART I: What are the possibilities?



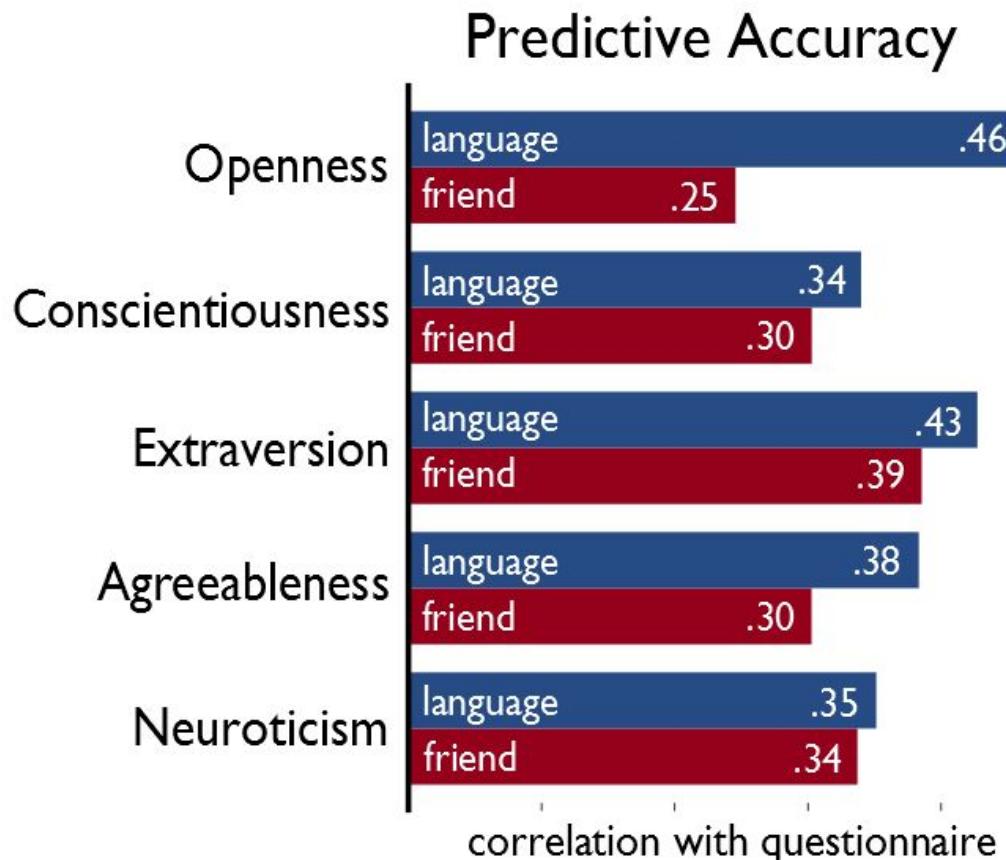
# Is it accurate? (Personality)



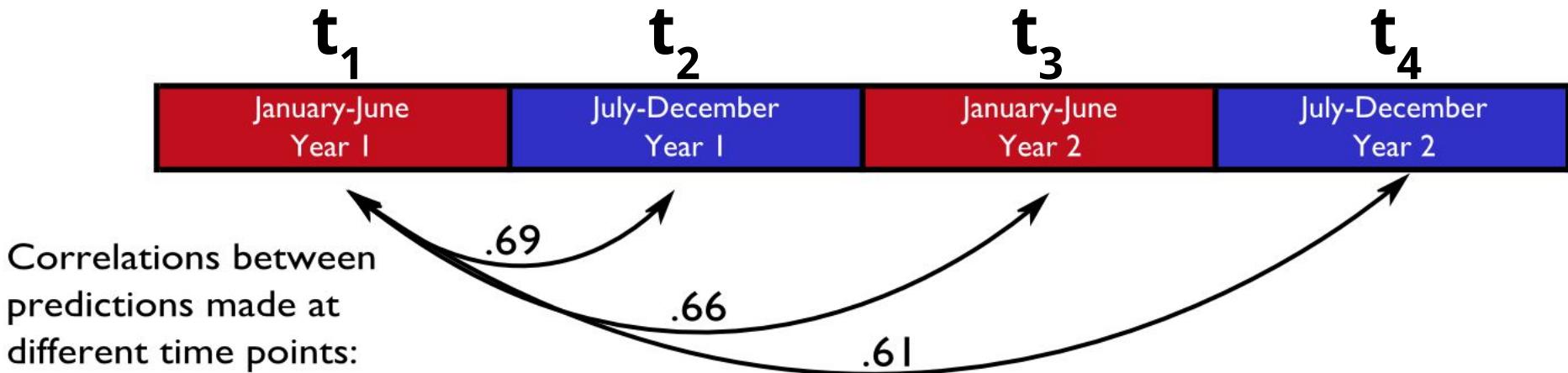
# Is it accurate?



# Is it accurate?

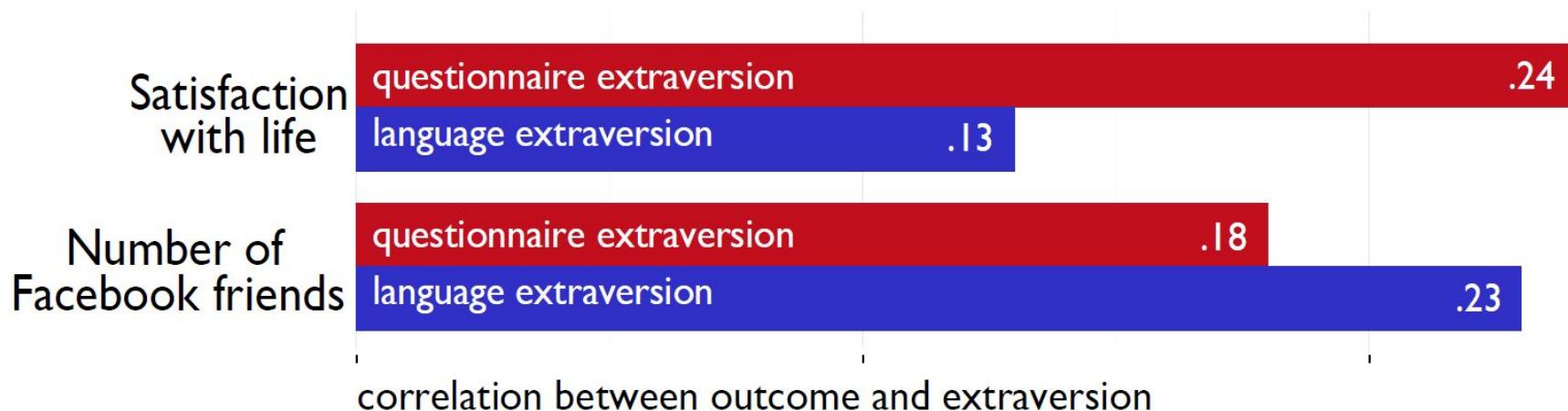


# Is it reliable and stable over time?

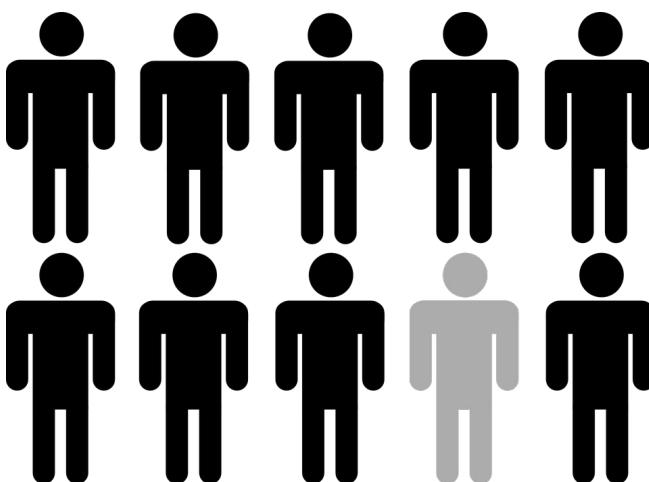
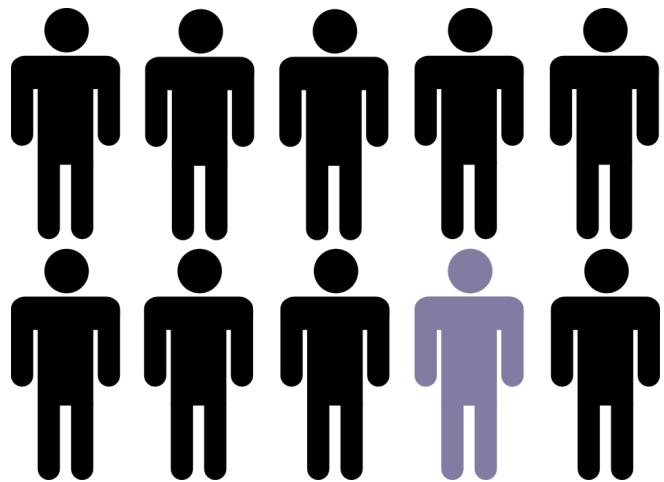
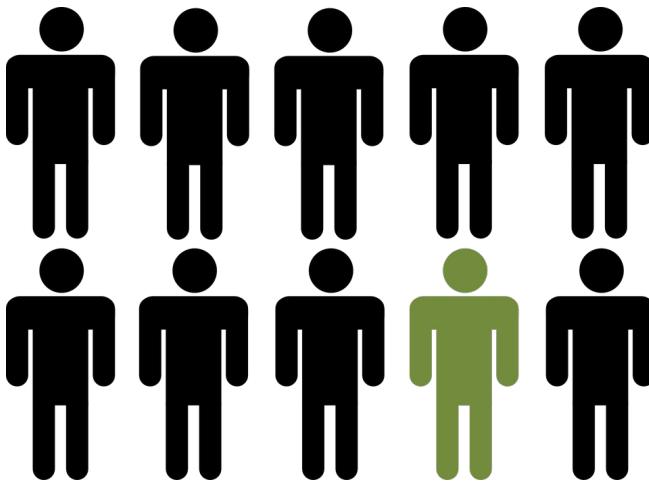
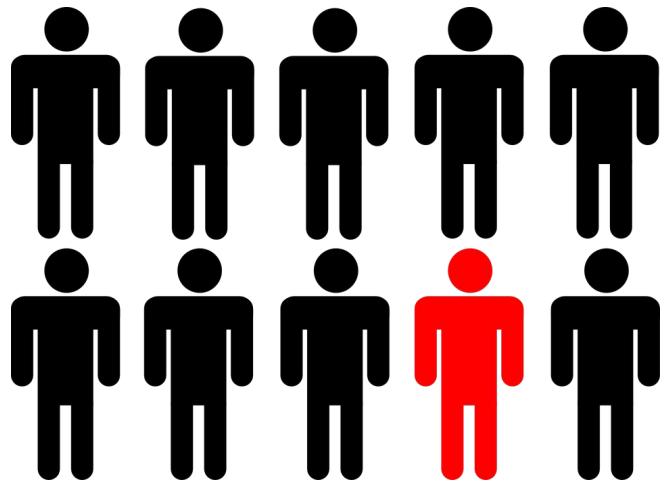


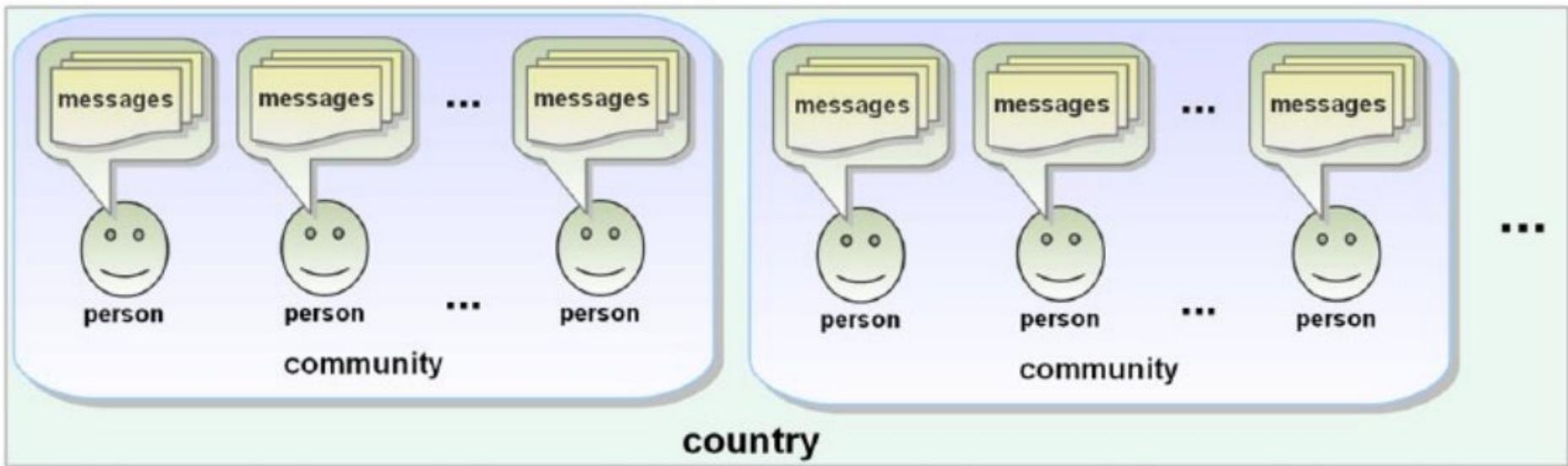
# Does it generalize to other outcomes?

## Example: extraversion

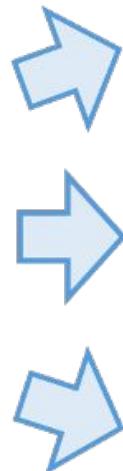
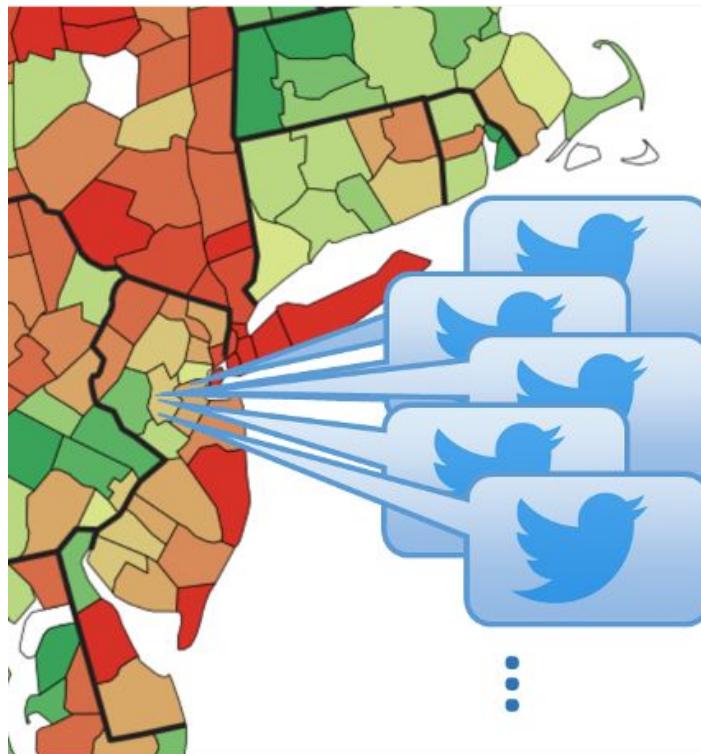




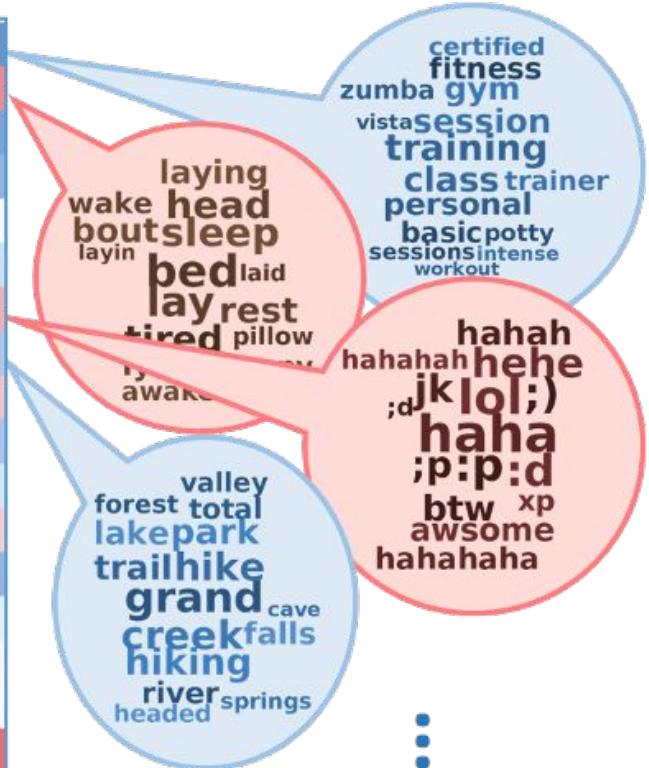




# Community Feature Extraction

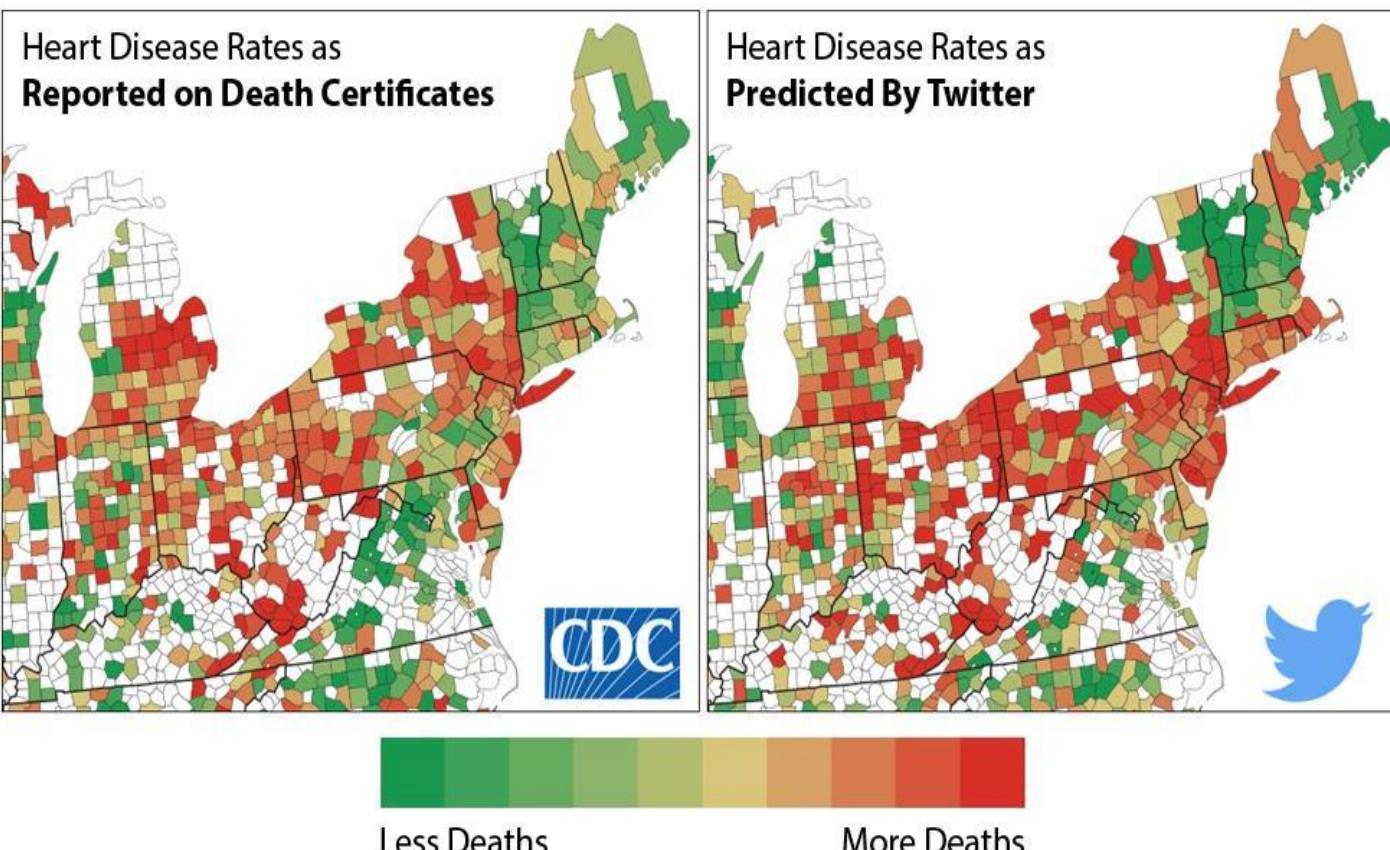


0.0852
0.8794
0.1415
0.1996
0.4561
0.3556
0.7532
0.2703
0.6872
0.2623
0.3795
0.6451
0.2032
0.4075
0.5010
0.4783
0.9845
0.6314





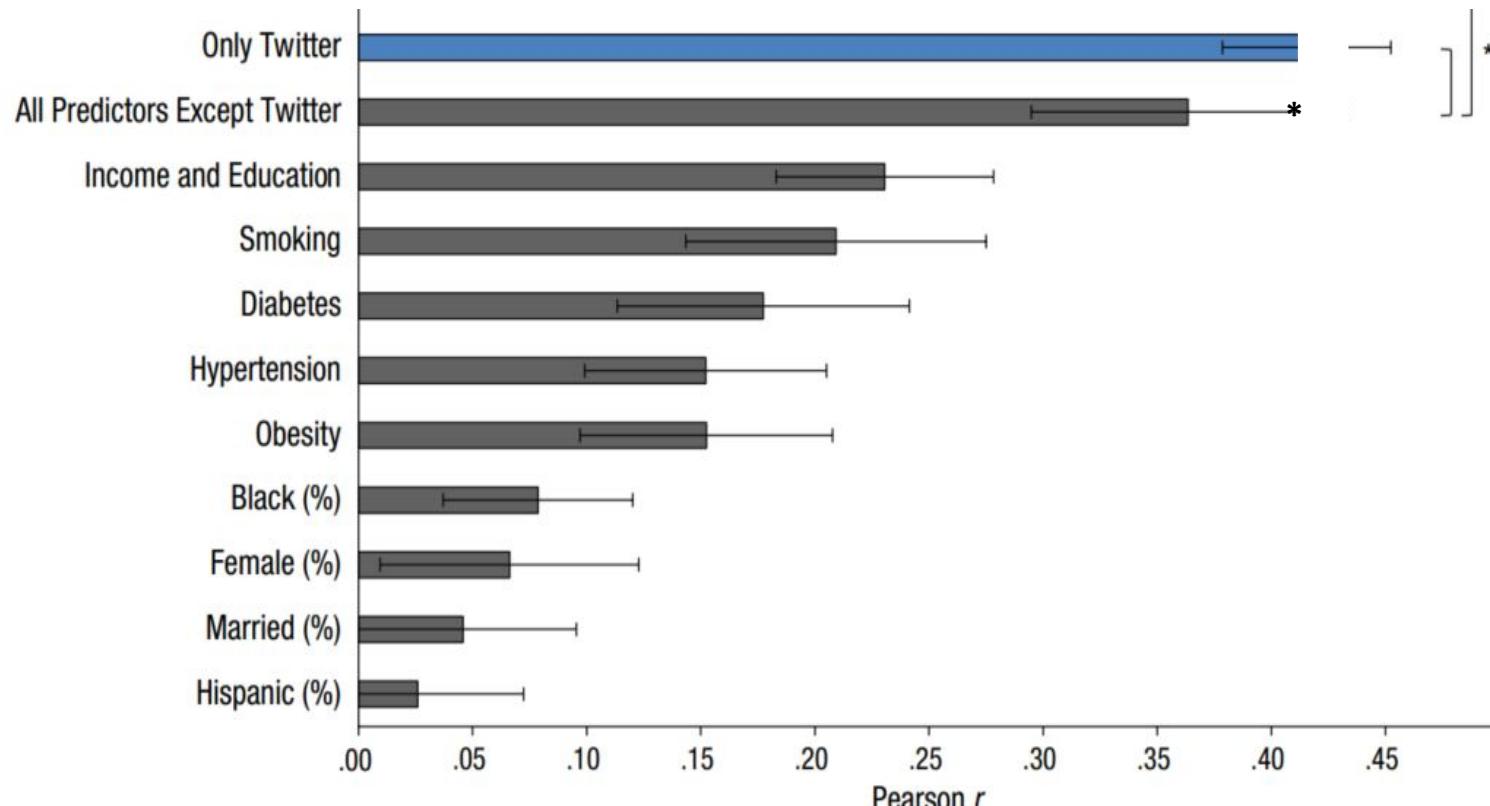
# Predicting heart disease: Accuracy



Eichstaedt, Schwartz, et al., 2015, *Psych Science*



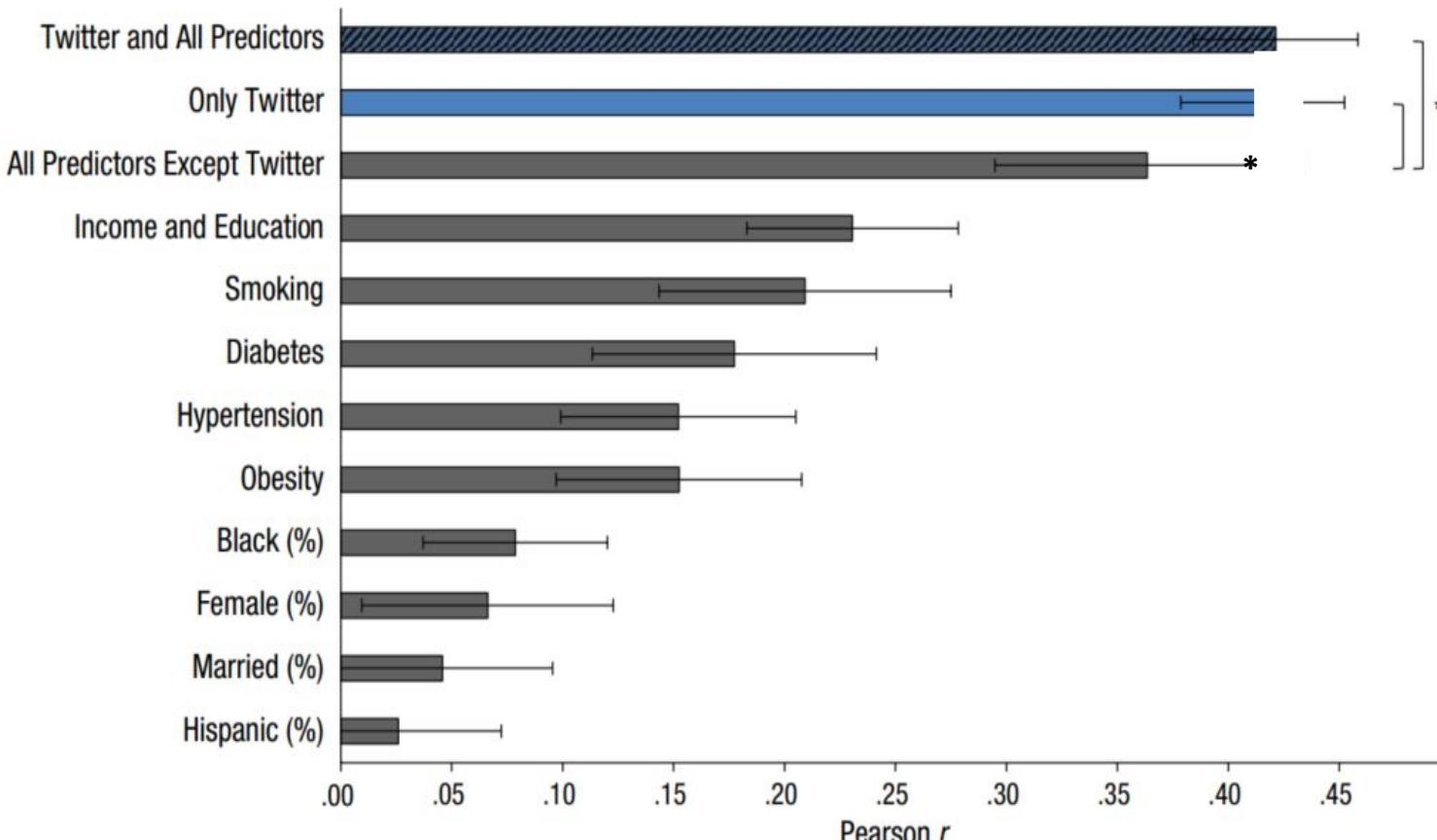
# Predicting heart disease: Accuracy



Eichstaedt, Schwartz, et al., 2015, *Psych Science*



# Predicting heart disease: Accuracy



Eichstaedt, Schwartz, et al., 2015, *Psych Science*

# Associated with Heart Disease

Hate,  
Interpersonal  
Tension

jealousy mad  
bitches  
envy hate jealous  
hating haters  
lovers famous hatin  
hater phase  
hated ya'll

nasty pieces allergic  
games faced bs head  
fake bullshit shit  
drama bull queens  
liars sneeze

grr passion  
grrr pit absolutely  
officially hate mondays  
burning hates grrrr  
despise mentioned  
fucking hating

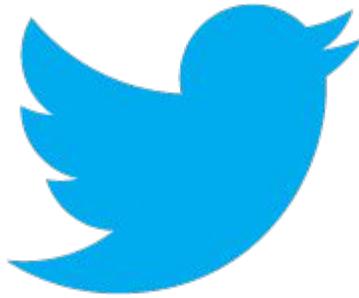
Optimism

opportunity  
possibilities  
talents  
opportunities  
discover possibility ); challenge improve  
create endless experience  
potential ability  
explore

reached dreams reaching  
perfection accomplish  
achieve goals  
greatness strive goal  
achieved set potential reach  
strive success

power strong  
overcome struggles  
strength courage  
struggle greater  
challenges faith  
peace obstacles trials  
stronger endure

# Community Real Estate



VS



	socioeconomics		demographics		socioeconomics + demographics	
	Fc	Ip	Fc	Ip	Fc	Ip
no lang	0.34	0.42	0.24	0.44	0.37	0.50
with lang	<b>0.41</b>	<b>0.56</b>	<b>0.39</b>	<b>0.57</b>	<b>0.42</b>	<b>0.59</b>

# Community Real Estate



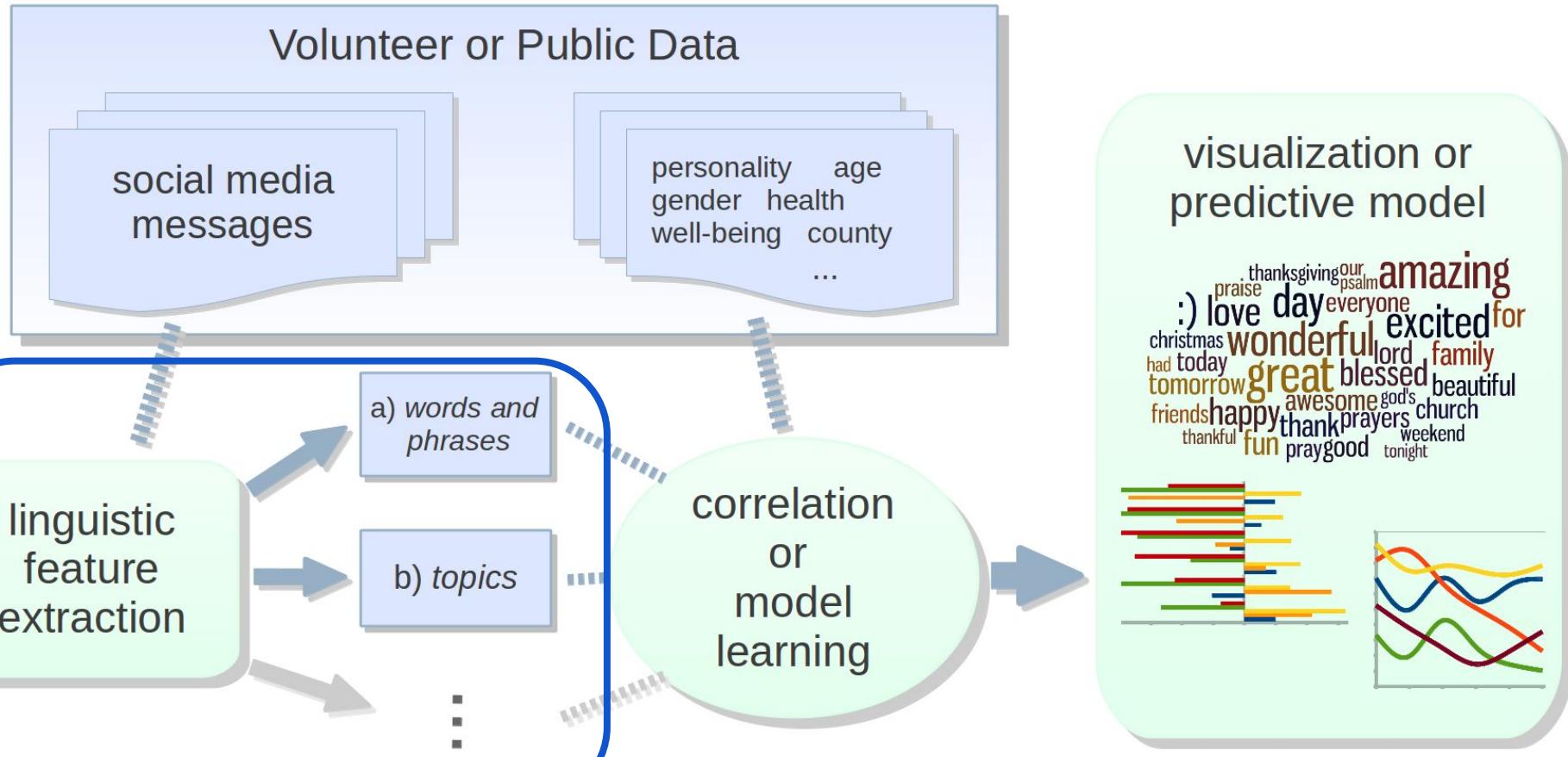
# **PART II: How?**

# PART II: How?

Two Examples:

1. Community Real Estate
2. Message Temporal Orientation

# PART II: How?



# Language Feature Space

*automatic content analysis*



# Language Feature Space

*automatic content analysis*

*closed-vocabulary*

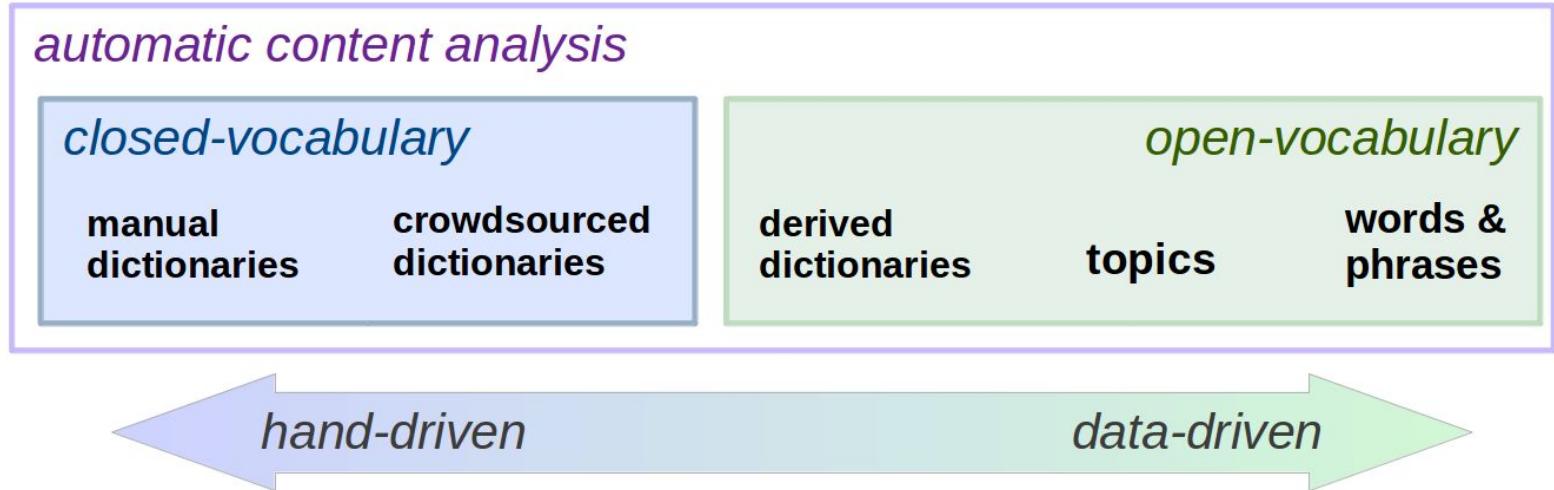
manual  
dictionaries

crowdsourced  
dictionaries

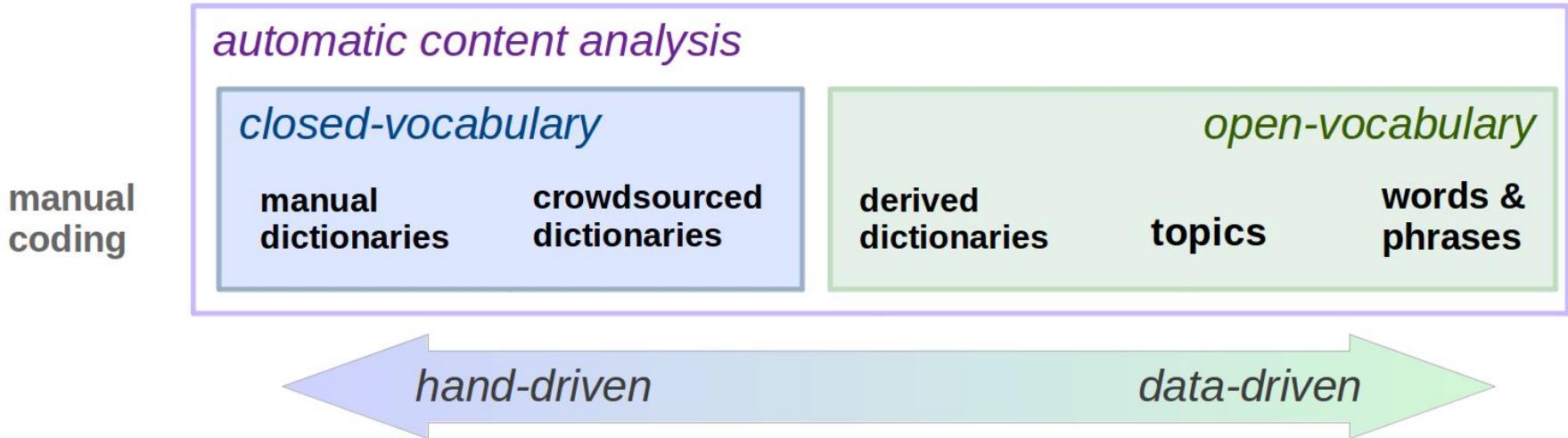
*hand-driven*

*data-driven*

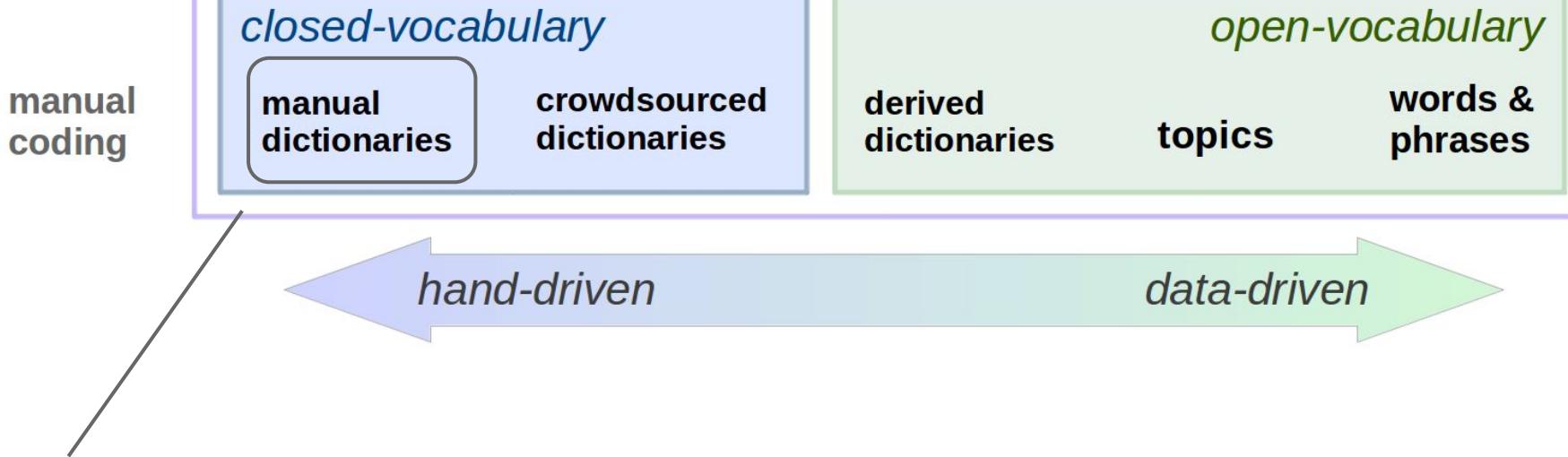
# Language Feature Space



# Language Feature Space



## *automatic content analysis*



Example:  
Linguistic Inquiry and  
Word Count  
(LIWC; Pennebaker et al., 2007)

## *automatic content analysis*

manual  
coding

### *closed-vocabulary*

manual  
dictionaries

crowdsourced  
dictionaries

### *open-vocabulary*

derived  
dictionaries

topics

words &  
phrases



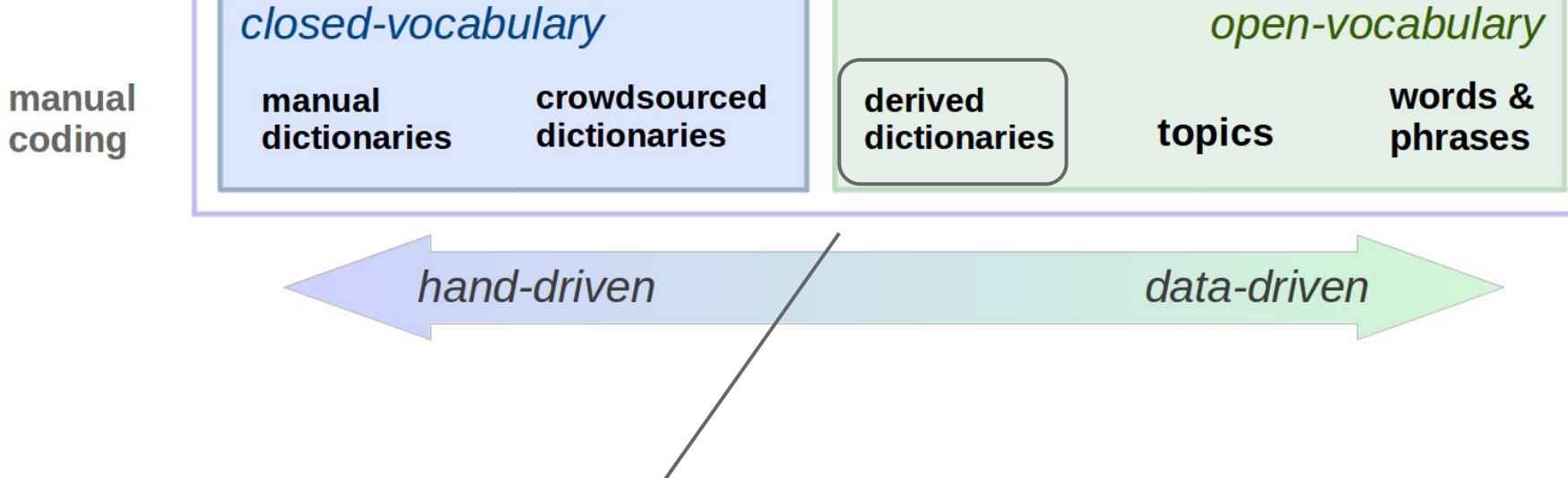
Examples:

- ANEW  
(Bradley & Lang, 1996)
- Hedonometer  
(Dodds et al., 2011)

+ often weighted

+ fuller coverage

## *automatic content analysis*

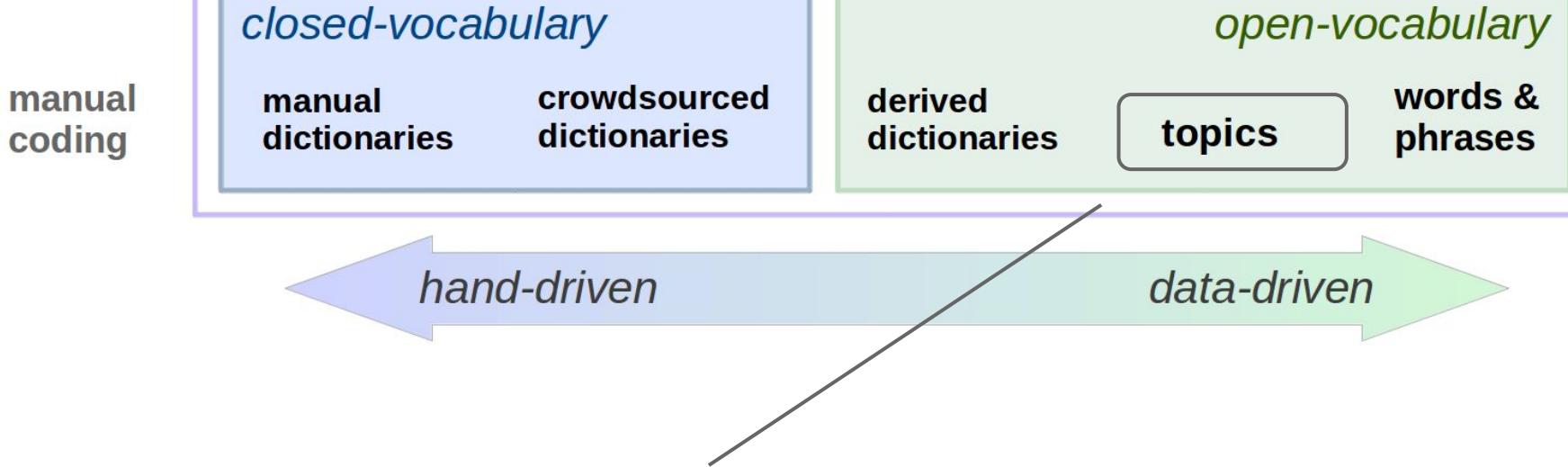


Examples:

- Sentiment  
(Pang & Lee, 2002)
- Affect & Intensity  
(Preotiuc et al., 2016)

+ real world distributions  
(still hypothesis driven)

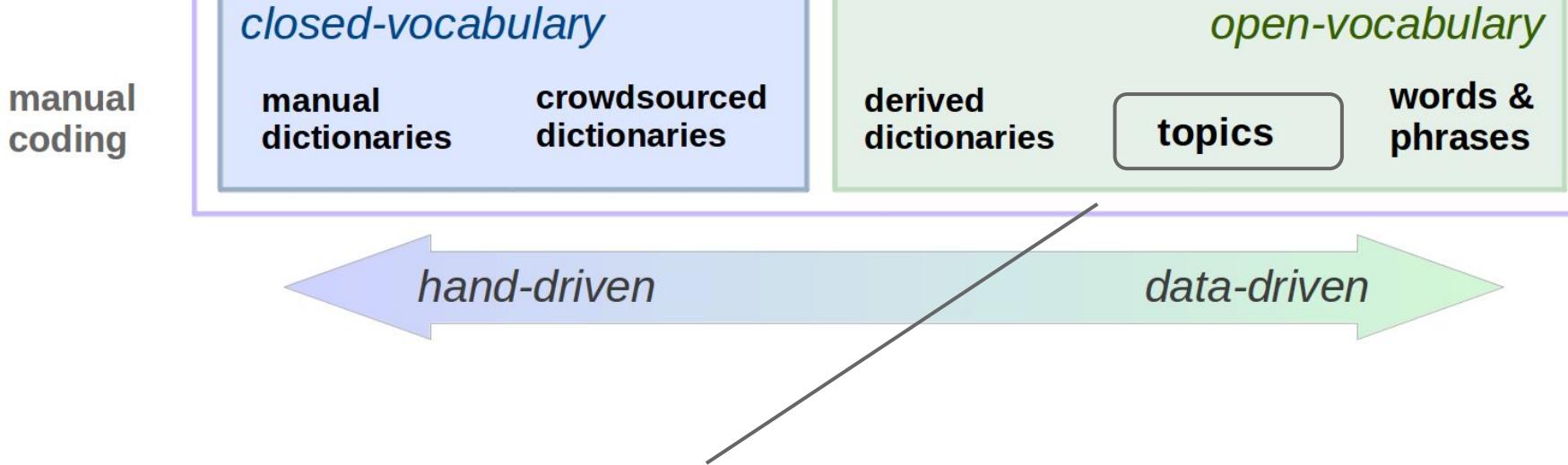
## *automatic content analysis*



Topic Modeling:

- + completely data-driven
- + “digestable”  
(still losing some information)

## *automatic content analysis*

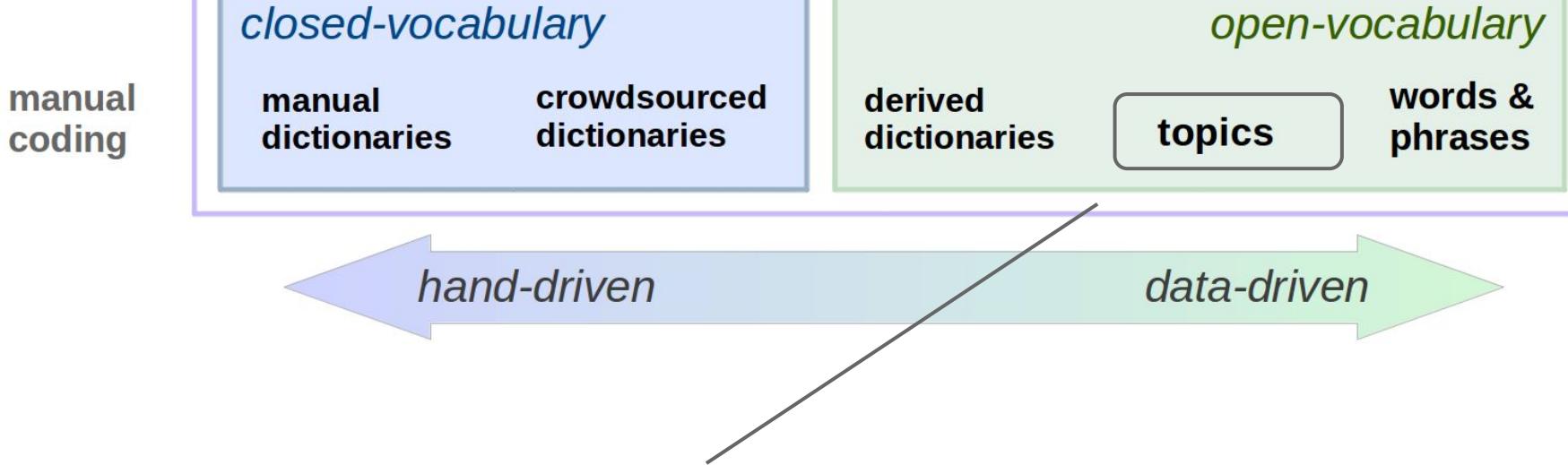


Topic Modeling:

- + completely data-driven
- + “digestable”  
(still losing some information)

$$p(\text{topic}|\text{tweet}) = \sum_{\text{term} \in \text{topic}} p(\text{topic}|\text{term})p(\text{term}|\text{tweet})$$

## *automatic content analysis*

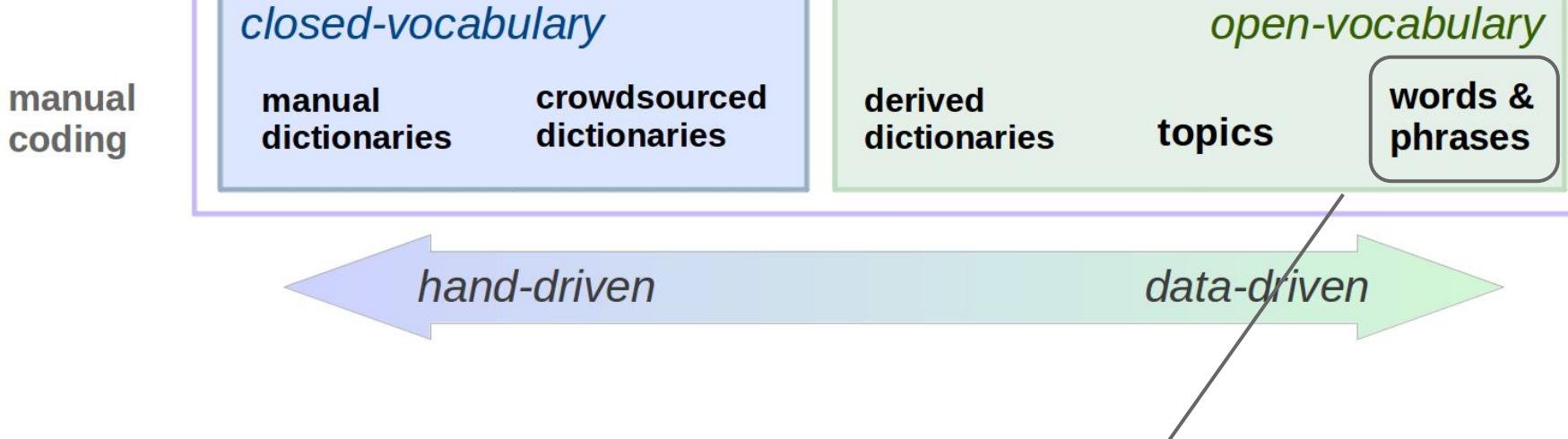


Topic Modeling:

- + completely data-driven
- + “digestable”  
(still losing some information)

$$p(\text{topic}|\text{user}) = \sum_{\text{term} \in \text{topic}} p(\text{topic}|\text{term})p(\text{term}|\text{user})$$

## *automatic content analysis*

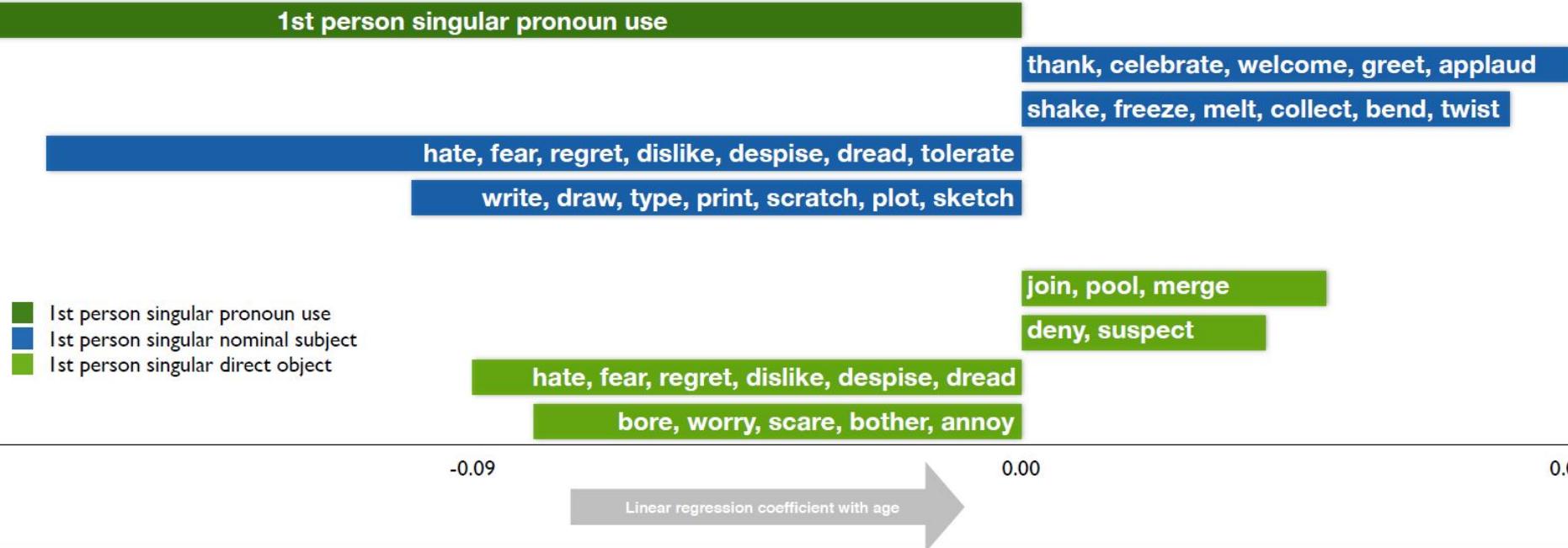


+ wide coverage  
+ fine-grained information  
phrases: capture some context

(not as “digestable”)

# **What about context?**

# What about context?

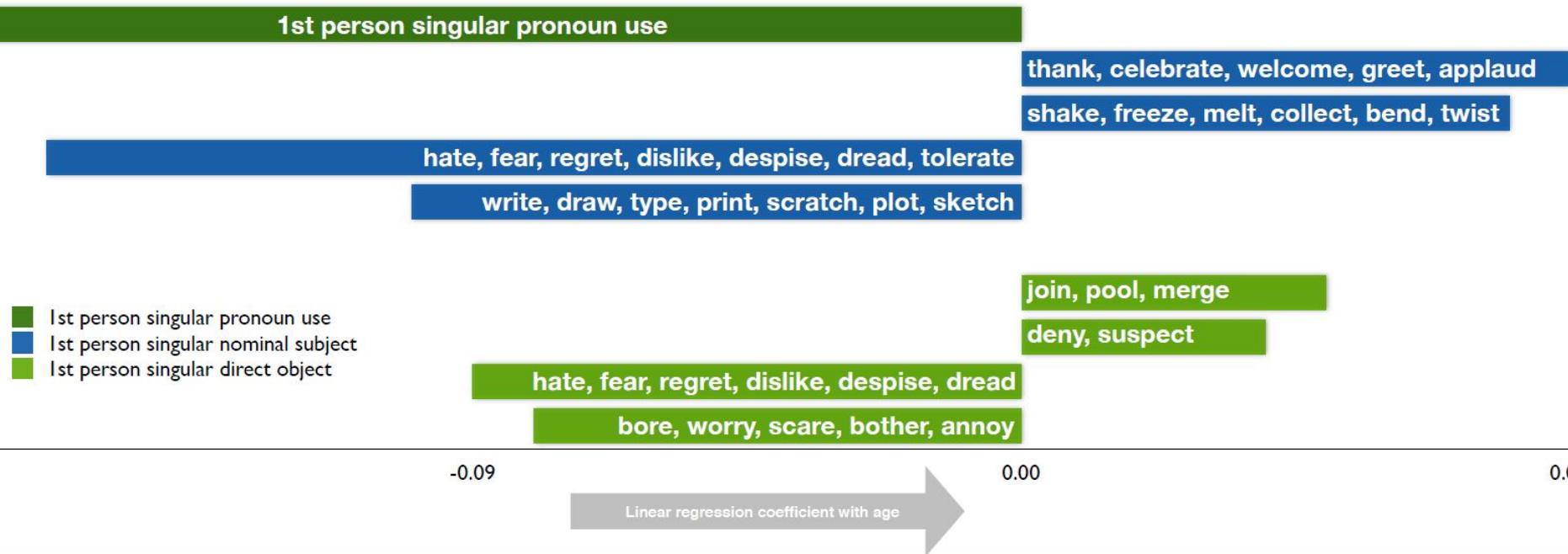


/ as subject of <verb>

me as direct object of <verb>

Rouhizadeh, M., Ungar, L., Buffone, A. & Schwartz, H.A (2016). Using Syntactic and Semantic Context to Explore Psychodemographic Differences in Self-reference. Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing

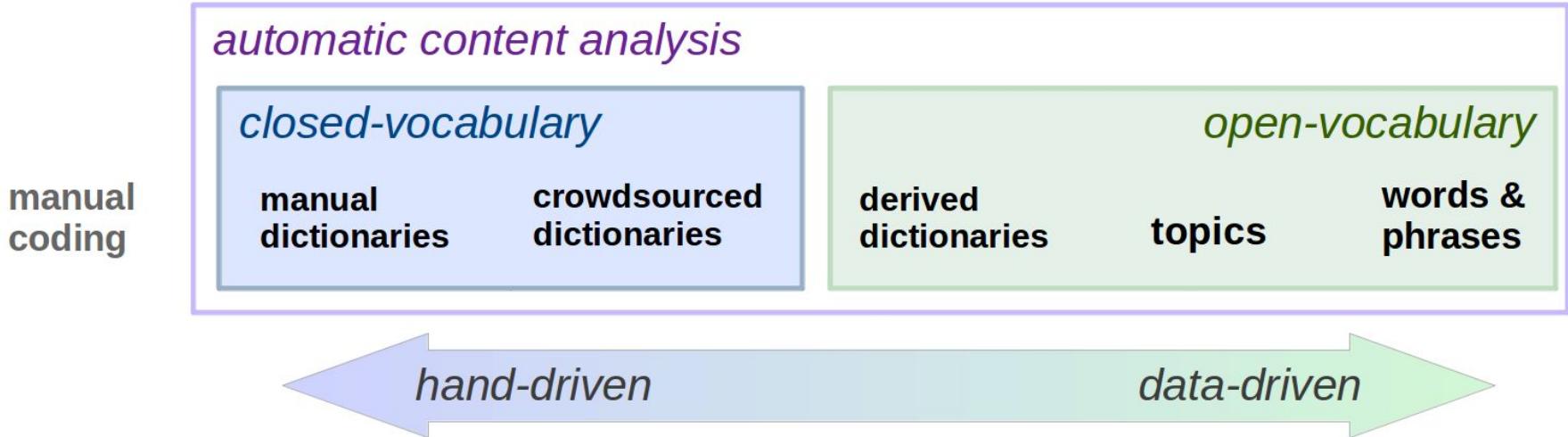
# What about context?



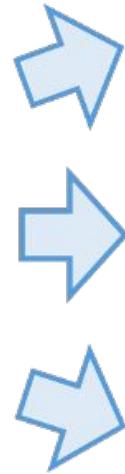
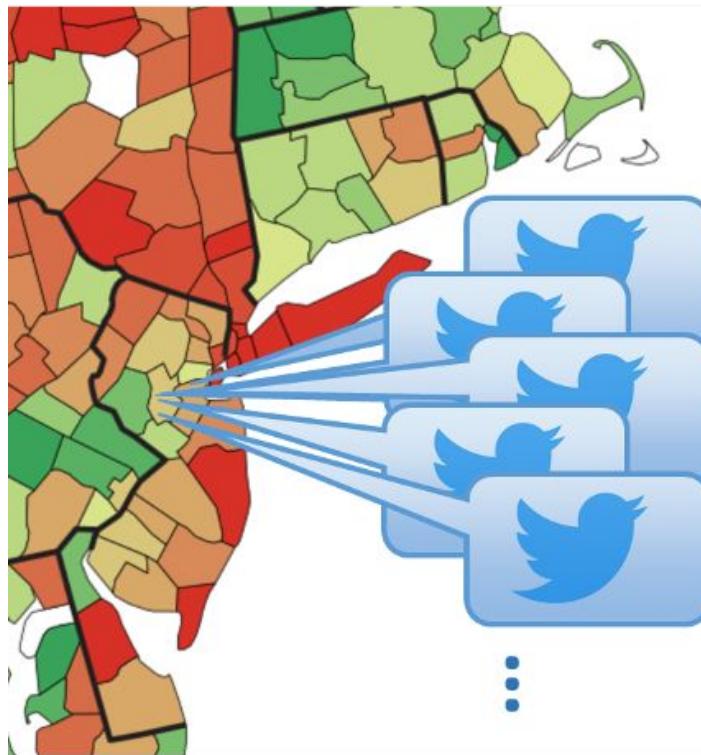
For human-level prediction, extent of benefit is open question.

Rouhizadeh, M., Ungar, L., Buffone, A. & Schwartz, H.A (2016). Using Syntactic and Semantic Context to Explore Psychodemographic Differences in Self-reference. Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing

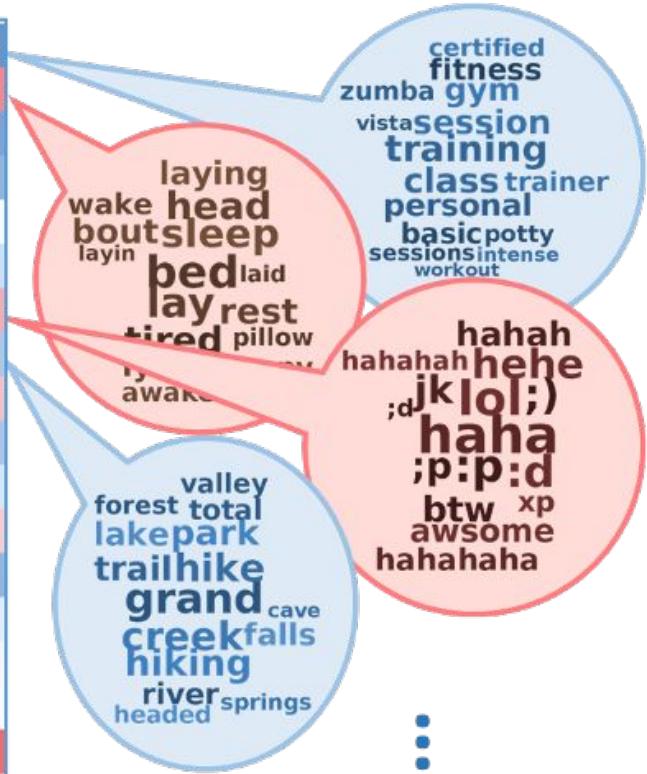
# Language Feature Space



# county\_id, feature1, value

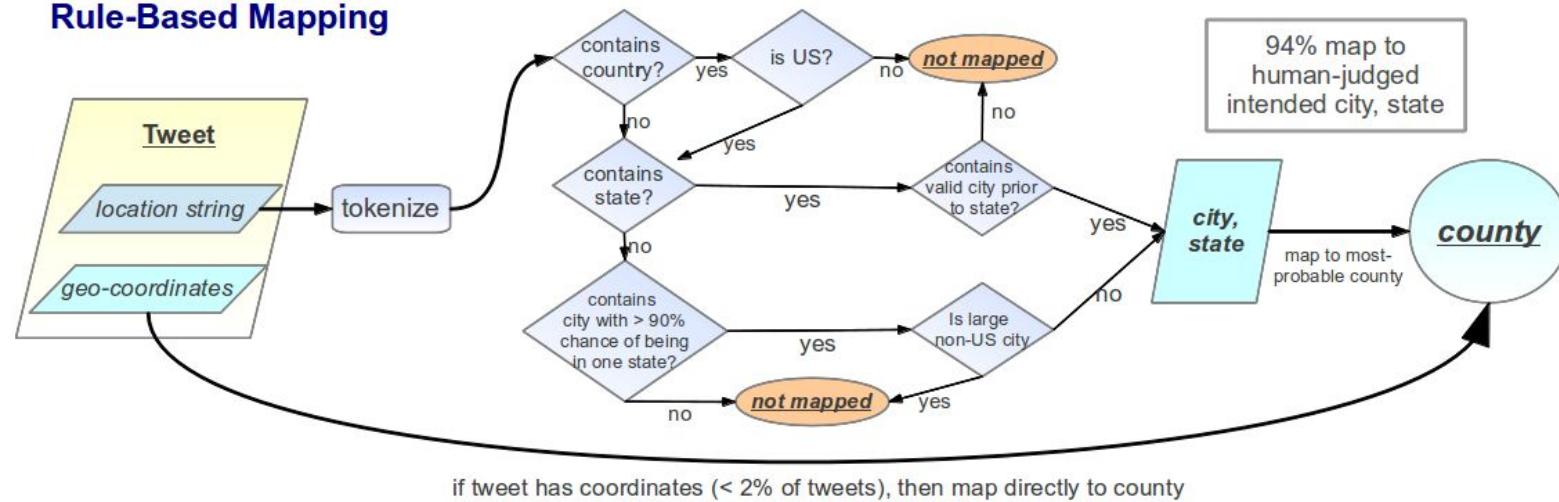


0.0852
0.8794
0.1415
0.1996
0.4561
0.3556
0.7532
0.2703
0.6872
0.2623
0.3795
0.6451
0.2032
0.4075
0.5010
0.4783
0.9845
0.6314



# Method: County-Mapping

## Rule-Based Mapping

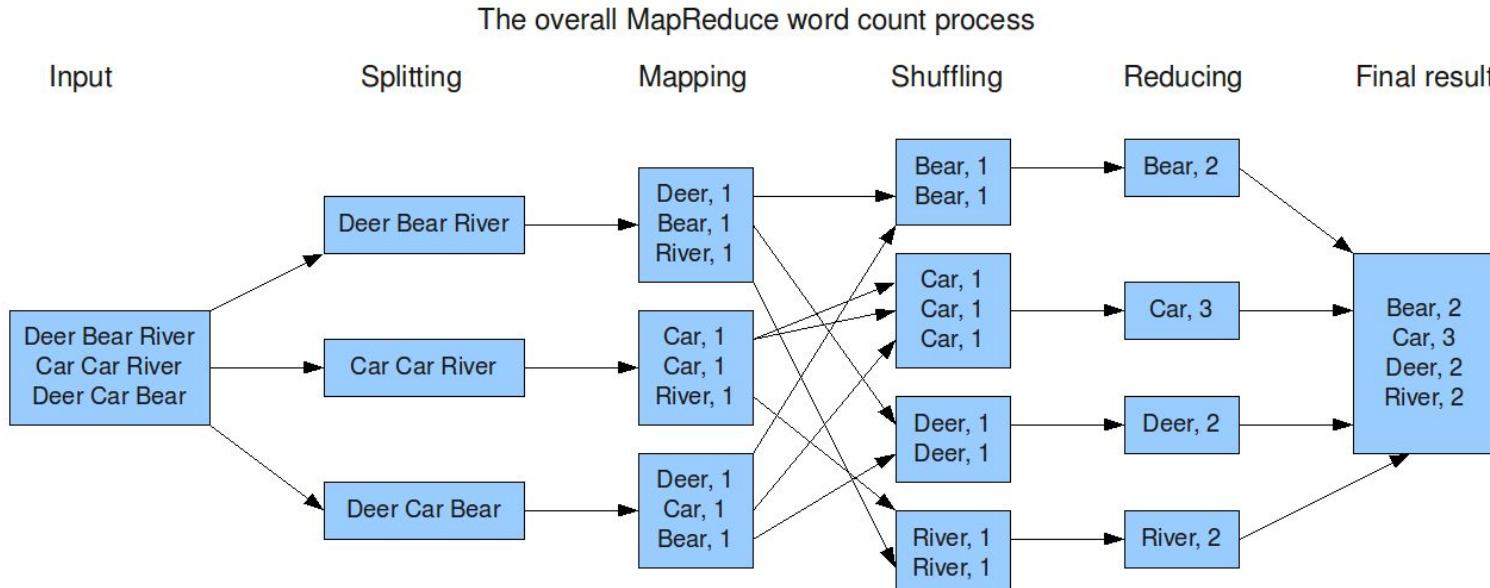


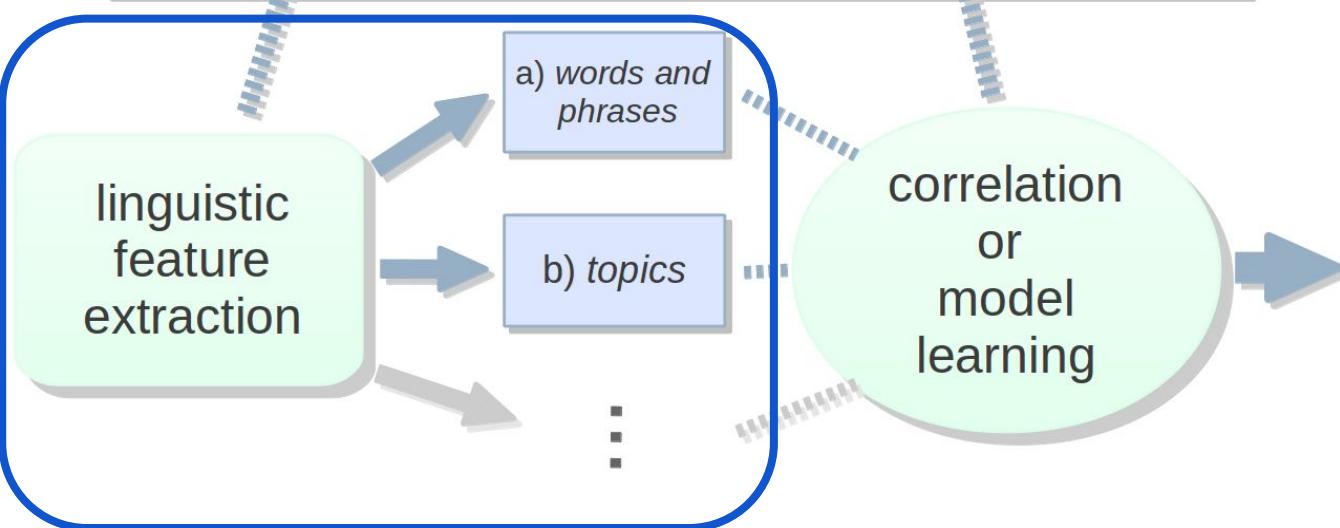
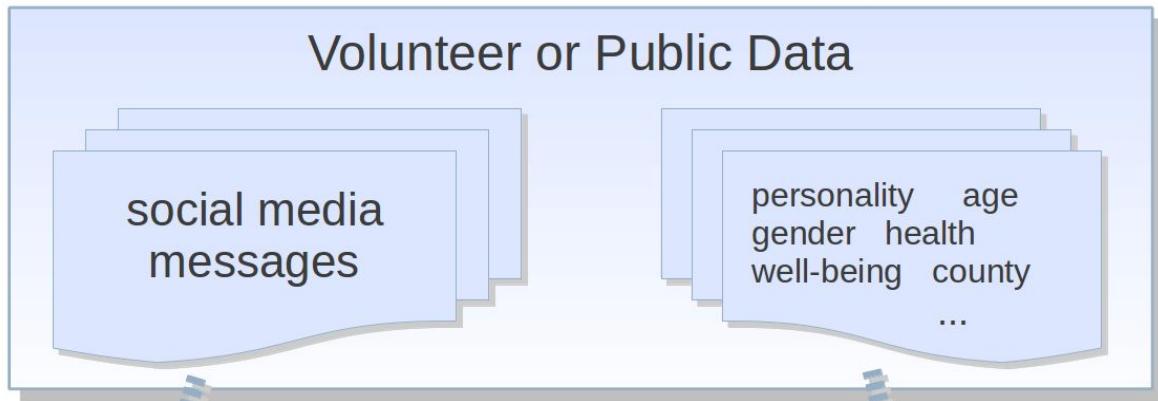
94% accurate map to human-judged intended city, state pair.

# Large Data

## Distributed Feature Extraction

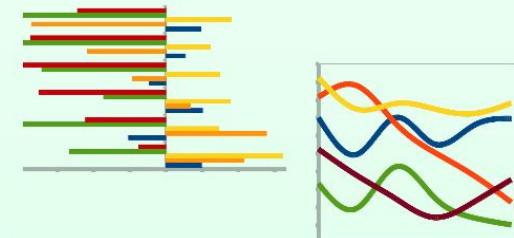
- approximate 1 billion tweets
  - for a single computer system to read the data is very time-consuming
- Utilize map-reduce in a “Hadoop” style cluster:

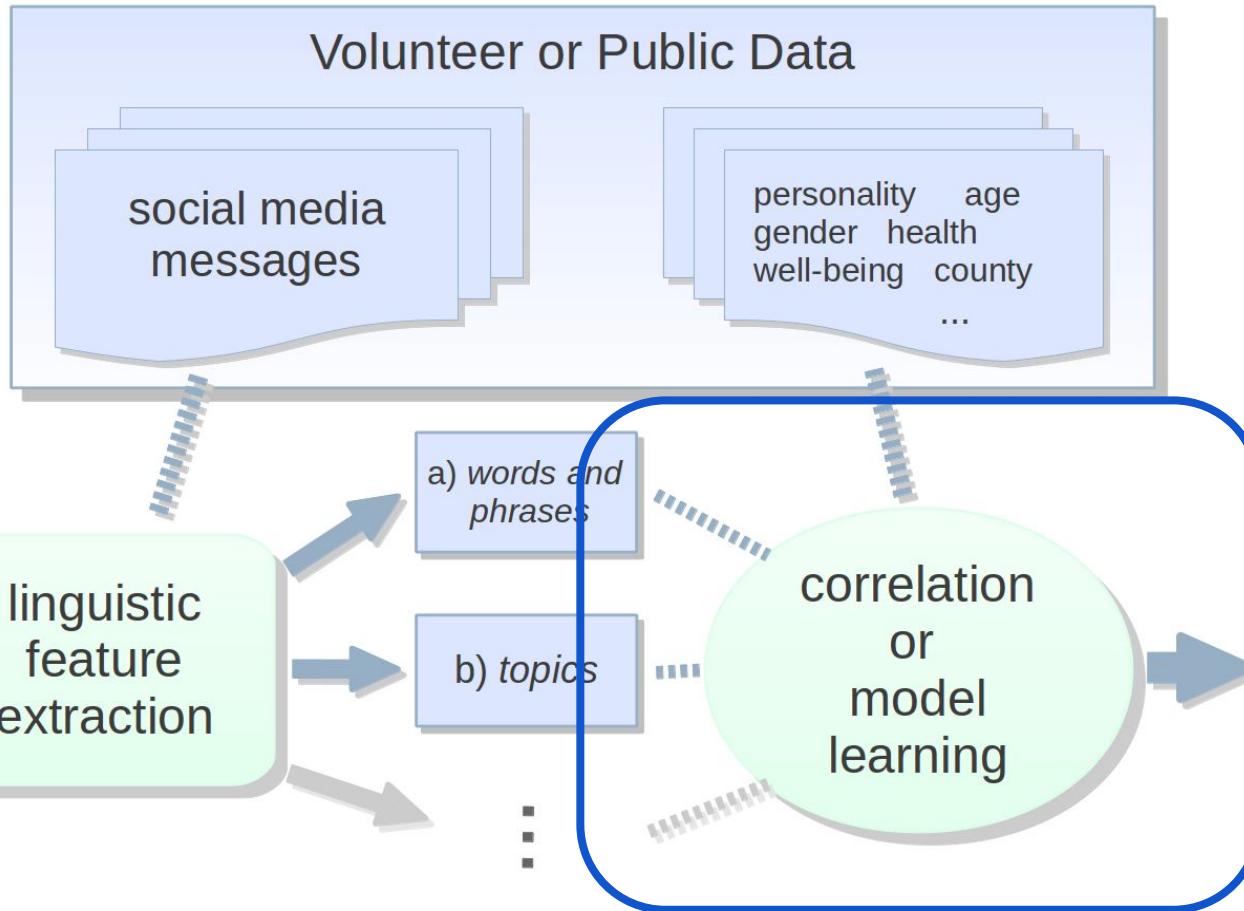




**visualization or predictive model**

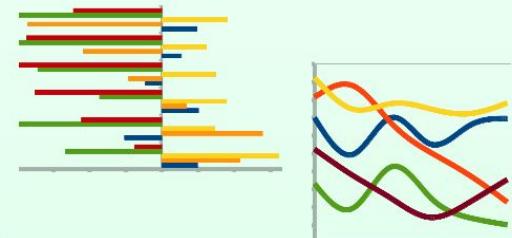
thanksgiving our psalm amazing  
praise everyone  
:) love day excited for  
christmas wonderful lord family  
had today tomorrow great blessed beautiful  
friends happy awesome god's church  
thankful thank prayers weekend  
fun pray good tonight

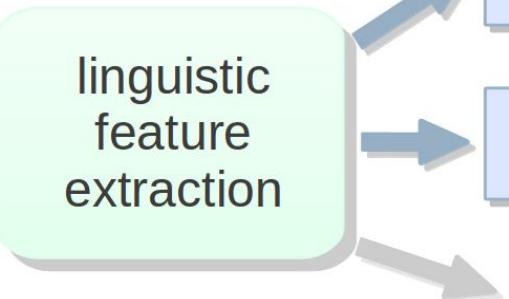
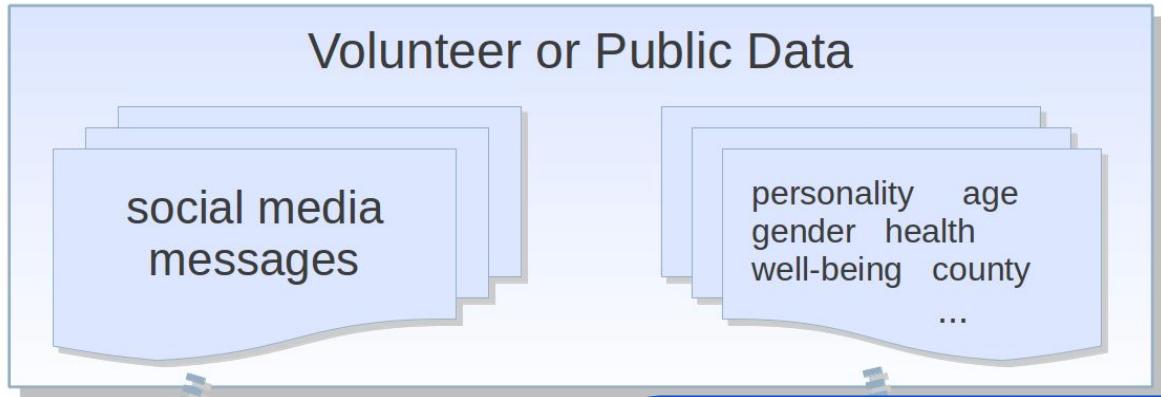




visualization or  
predictive model

thanksgiving  
praise  
:) love day everyone  
christmas wonderful for  
had today lord family  
tomorrow great blessed beautiful  
friends happy awesome god's church  
thankful thank prayers weekend  
fun pray good tonight

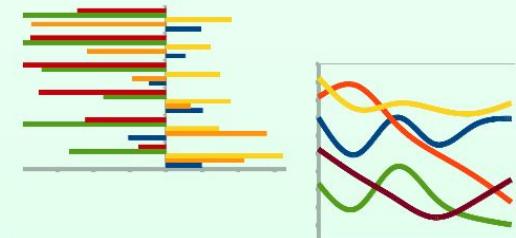




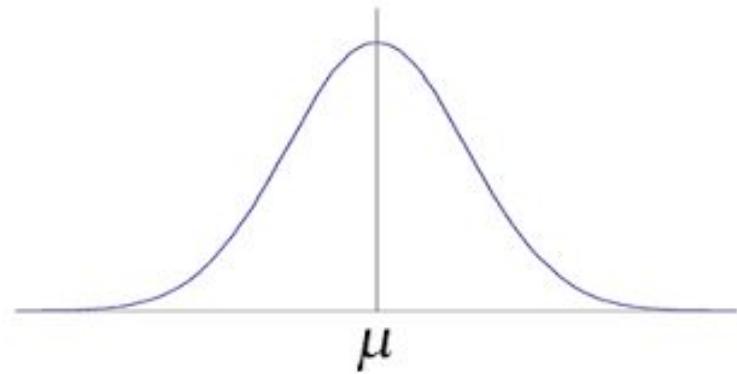
Penalized  
Regression or  
Classification  
for example,  
*Ridge Regression*  
*Support Vector Machine*

visualization or  
predictive model

thanksgiving our psalm amazing  
praise day everyone  
:) love wonderful excited for  
christmas had today lord family  
tomorrow great blessed beautiful  
friends happy awesome god's church  
thankful thank prayers weekend  
fun pray good tonight

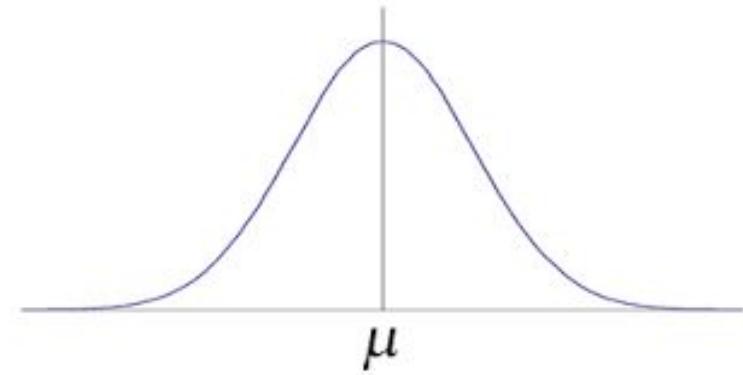
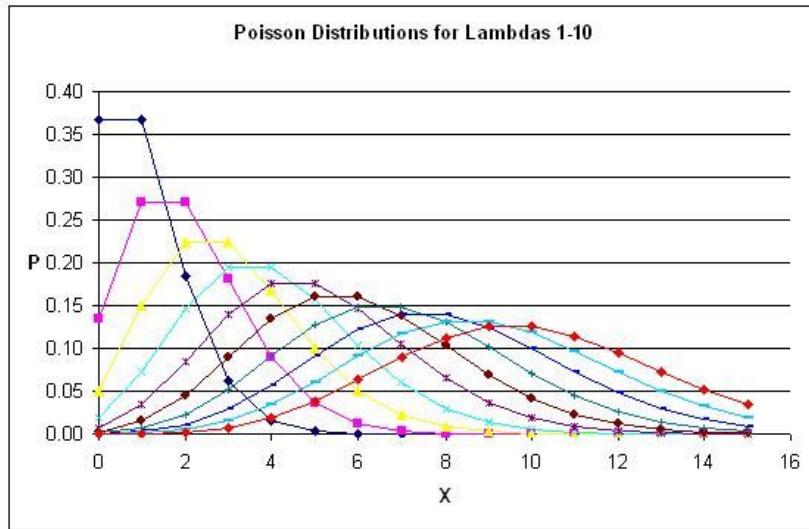


# Challenge



Topics  
~2k

# Challenge



**NGrams**  
~20k after frequency filter

**Topics**  
~2k

# Challenge

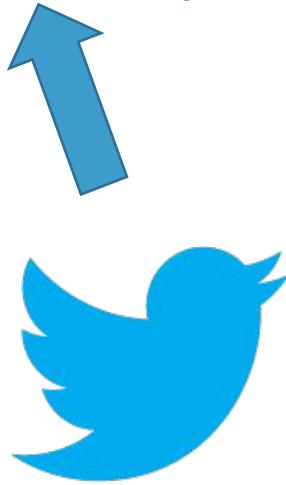


**NGrams**  
~20k after frequency filter

**Topics**  
~2k

# Challenge: Even Worse

Language:  
High-dimensional,  
sparse, and noisy.

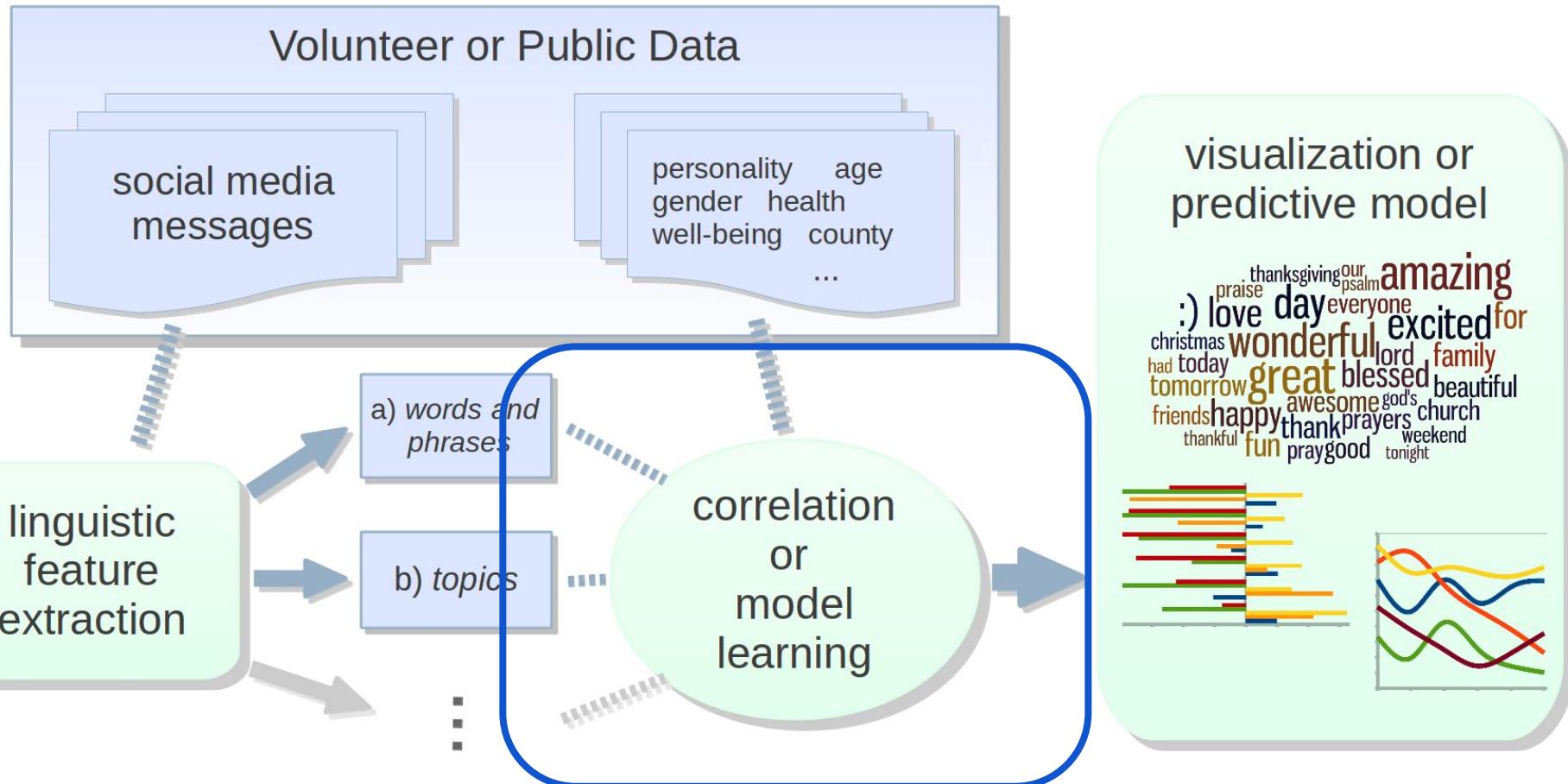


VS

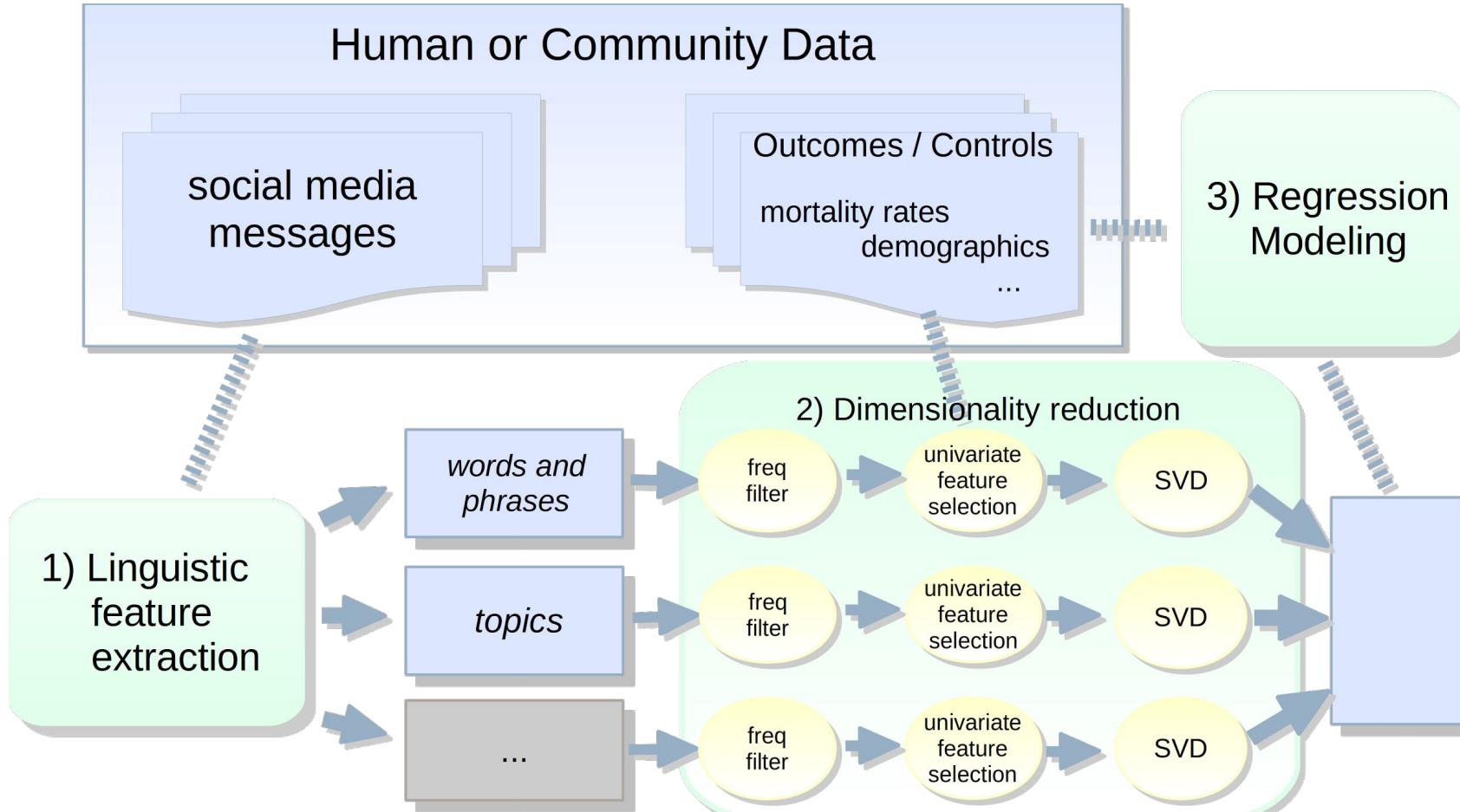
Socio-economics:  
few and well estimated



# PART II: How?



# Prediction



# Marginal gain from Socioeconomics:

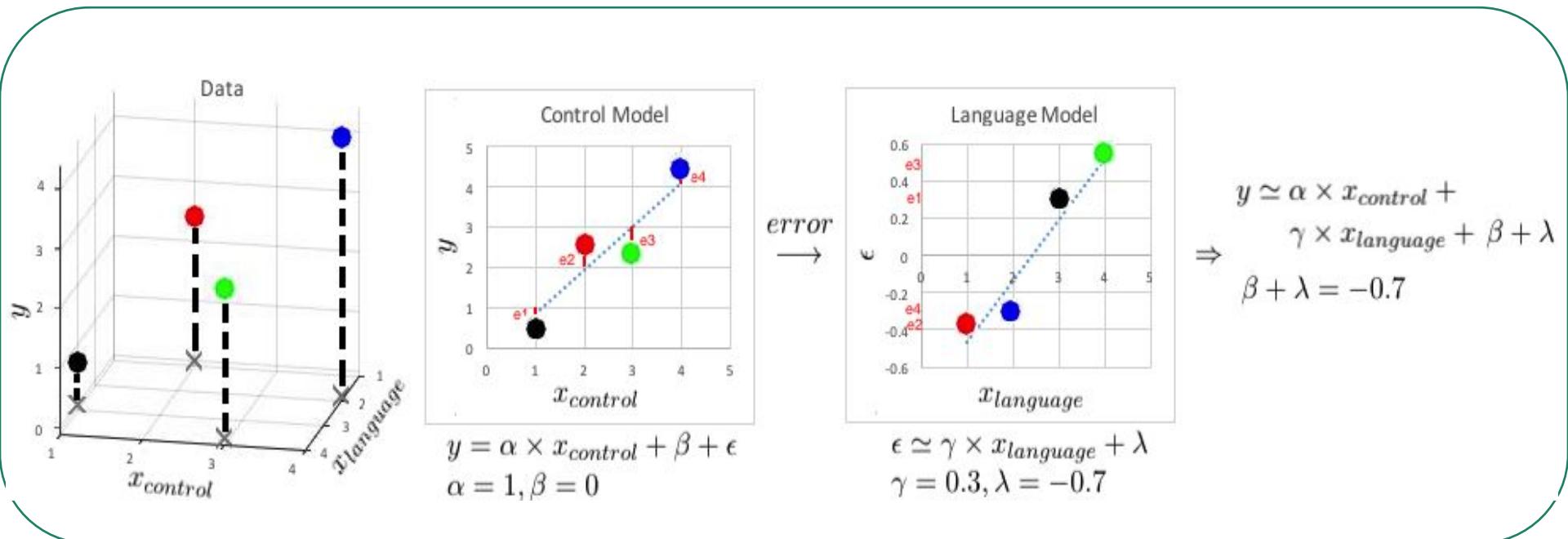
	Foreclosure	Increased-price
language	0.38	0.48
combined	0.40	0.49

Out of sample Pearson r

# *Residualized Control Model*

- Effectively use both low dimensional control features and high-dimensional, noisy language features:
  - Train a control model using the control values
  - Calculate the residual error and consider it as the new label
  - Train a language model over the new labels

# Residualized Control Model



# *Residualized Control vs. Combined Model*

$$Y = \alpha x_1 + \beta x_2 + \epsilon$$

Both learn same linear model.

Both learn the same linear model above, but

- Different learning algorithm per variable.
- As if different penalization.

# Marginal gain from Socioeconomics:

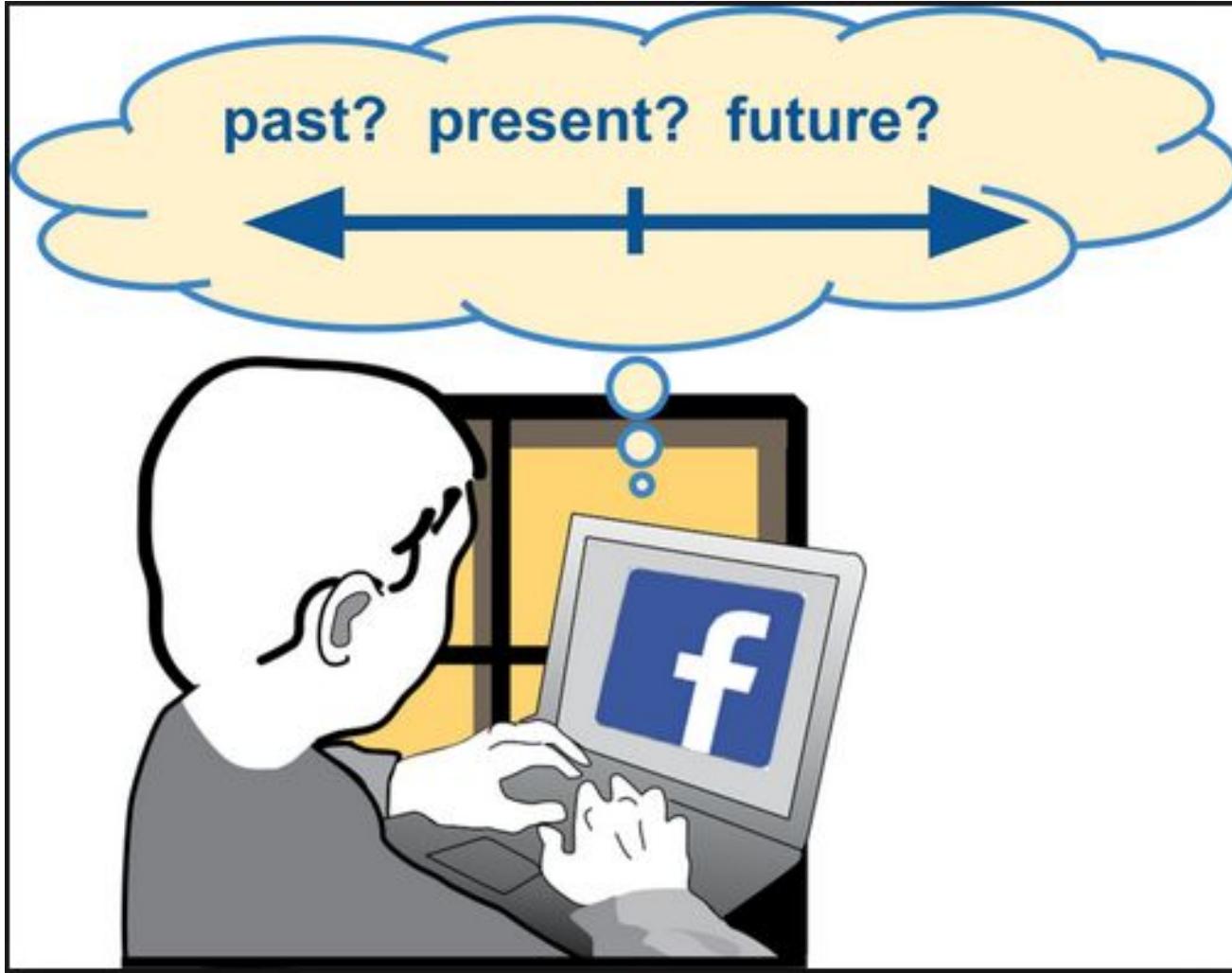
	Foreclosure	Increased-price
language	0.38	0.48
combined	0.40	0.49

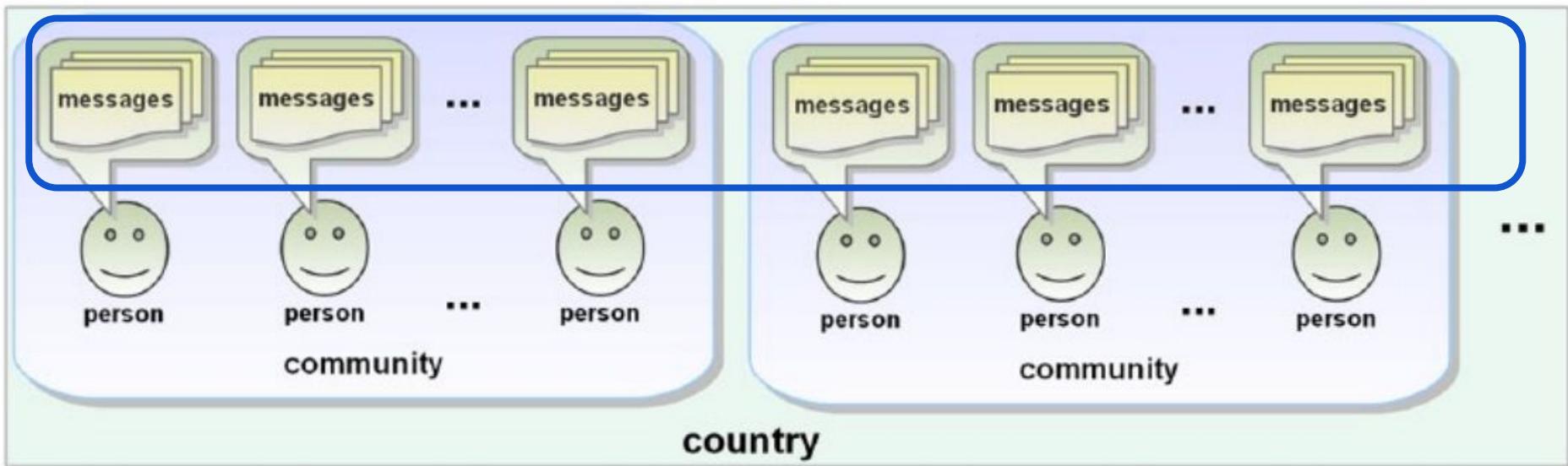
Out of sample Pearson r

# Marginal gain from Residualized Control

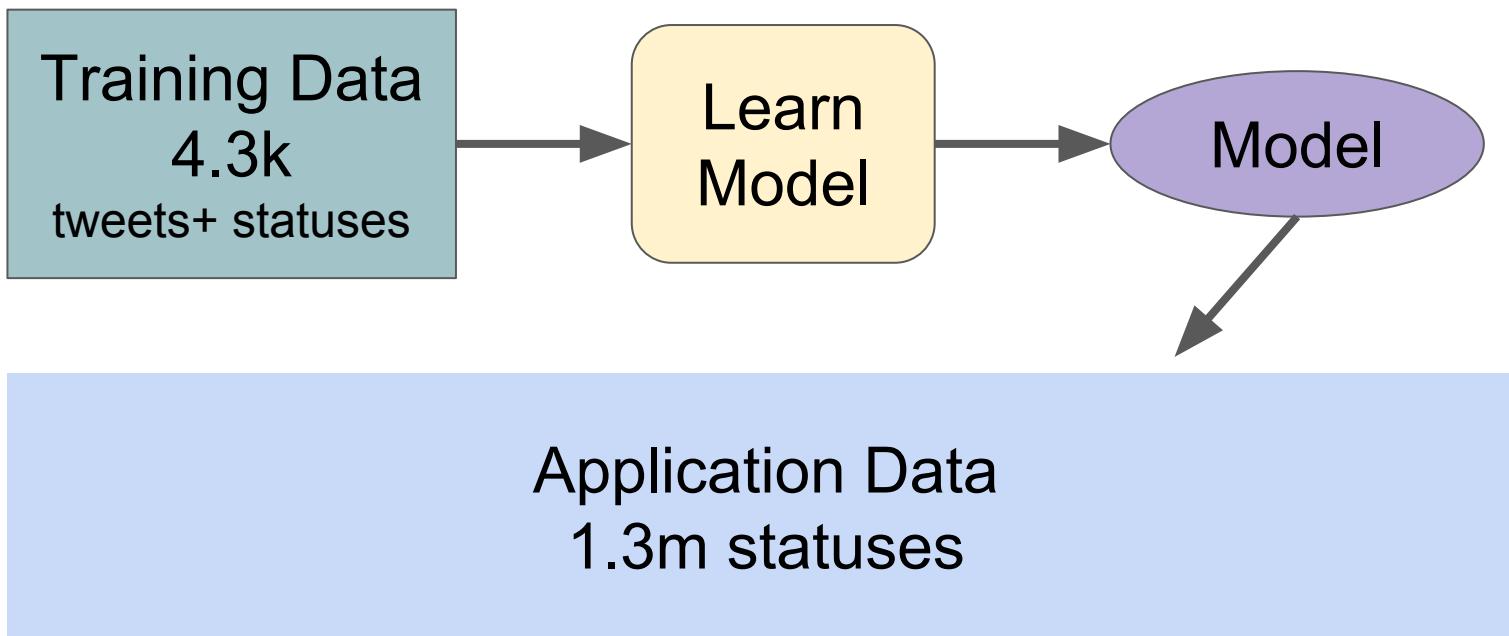
	Foreclosure	Increased-price
language	0.38	0.48
combined	0.40	0.49
residualized control	<b>0.42</b>	<b>0.59</b>

Out of sample Pearson r



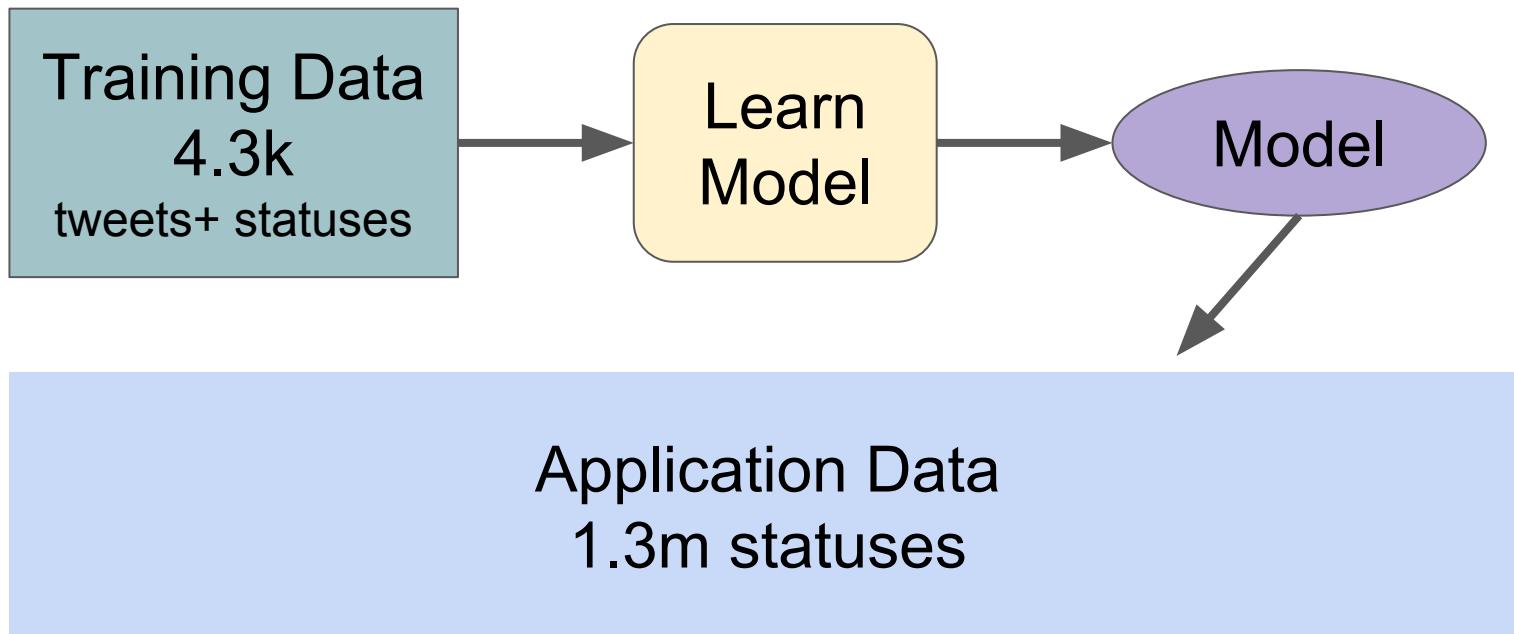


# Building a model



# Building a model

message	R1	R2	R3	m	class
<i>did nothing this morning but watch TV and it was fantastic =)</i>	-.67	-.50	-.50	-.55	past
<i>dislikes being sick.... and misses her bf</i>	0	0	0	0	present
<i>pancake day tomorrow pancake day tomorrow xxxxx</i>	.50	.50	1	.67	future



# Building a model

message	R1	R2	R3	m	class
<i>did nothing this morning but watch TV and it was fantastic =)</i>	-.67	-.50	-.50	-.55	past
<i>dislikes being sick.... and misses her bf</i>	0	0	0	0	present
<i>pancake day tomorrow pancake day tomorrow xxxxx</i>	.50	.50	1	.67	future



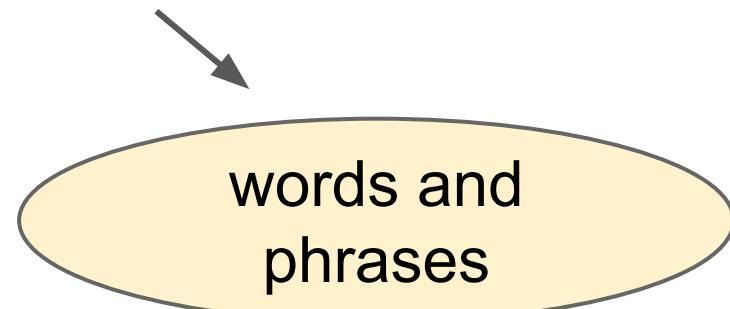
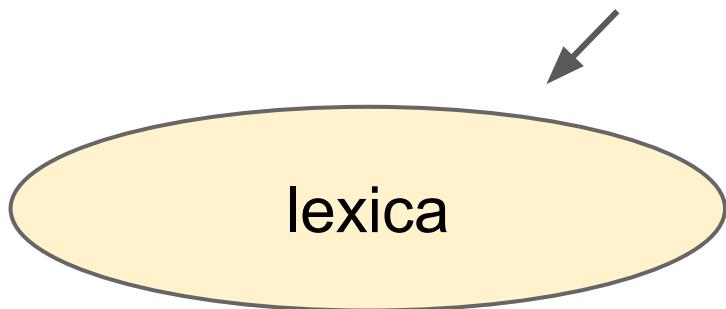
Linguistic Feature Extraction

# Building a model

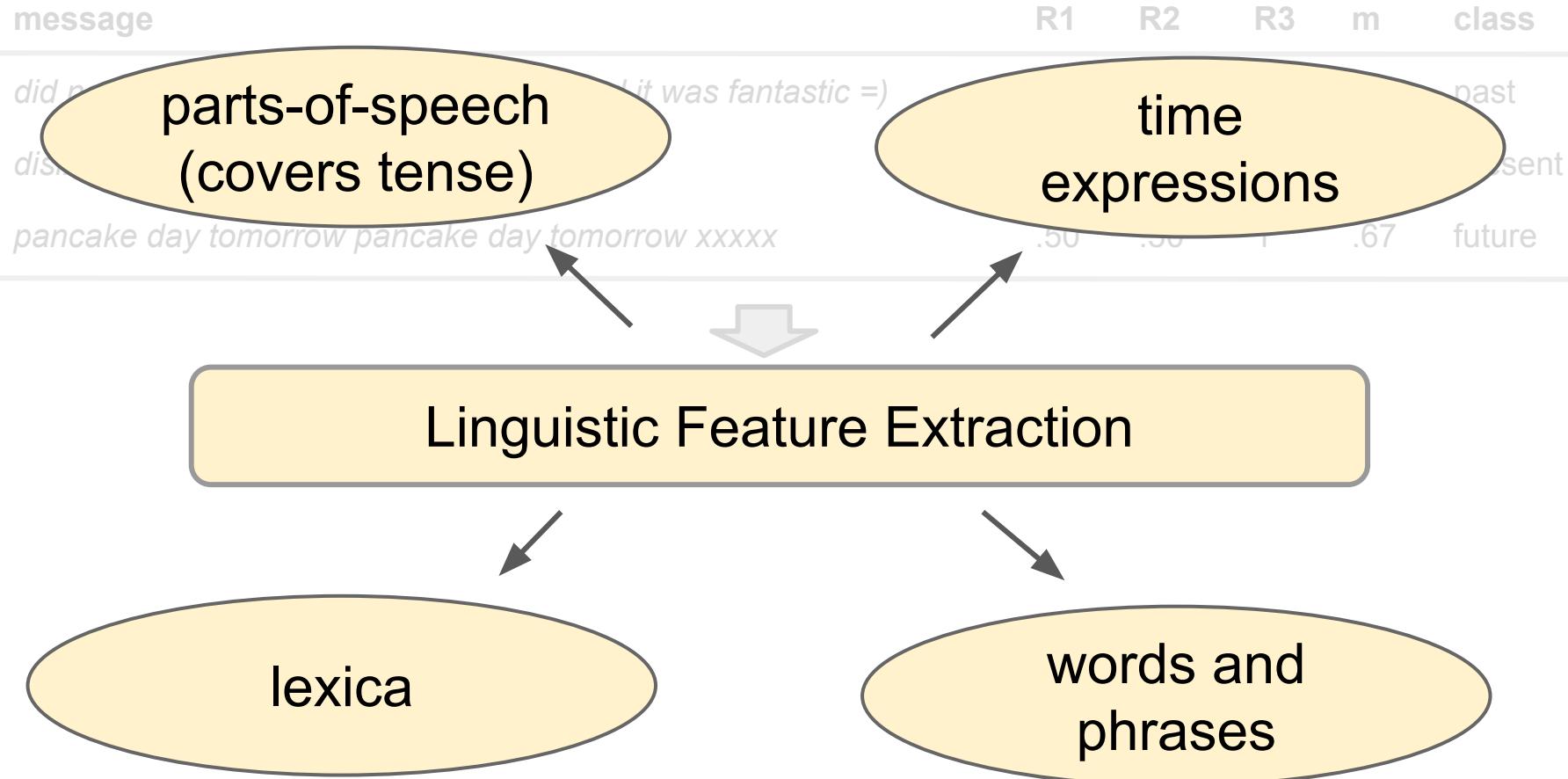
message	R1	R2	R3	m	class
<i>did nothing this morning but watch TV and it was fantastic =)</i>	-.67	-.50	-.50	-.55	past
<i>dislikes being sick.... and misses her bf</i>	0	0	0	0	present
<i>pancake day tomorrow pancake day tomorrow xxxxx</i>	.50	.50	1	.67	future



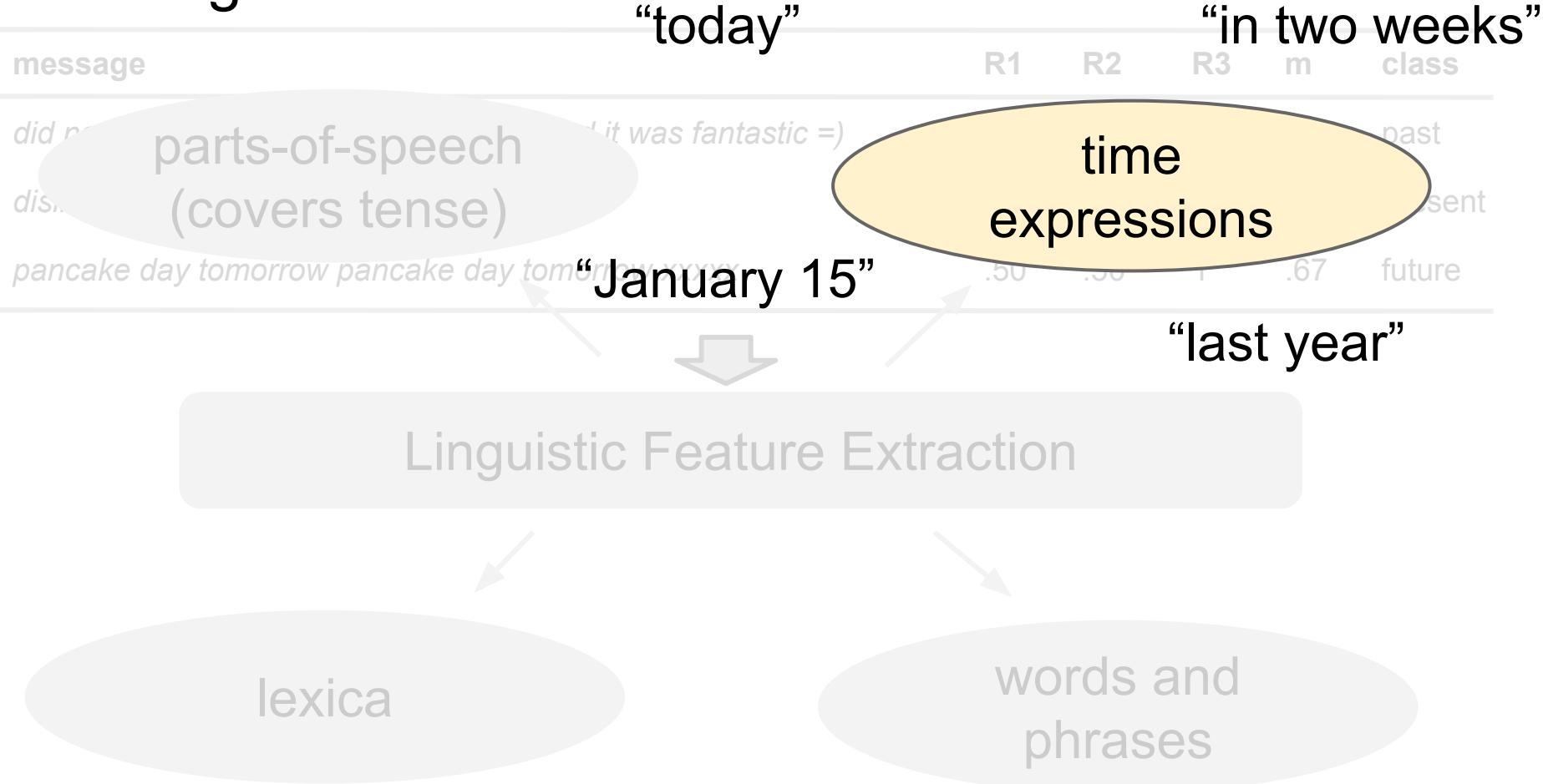
## Linguistic Feature Extraction



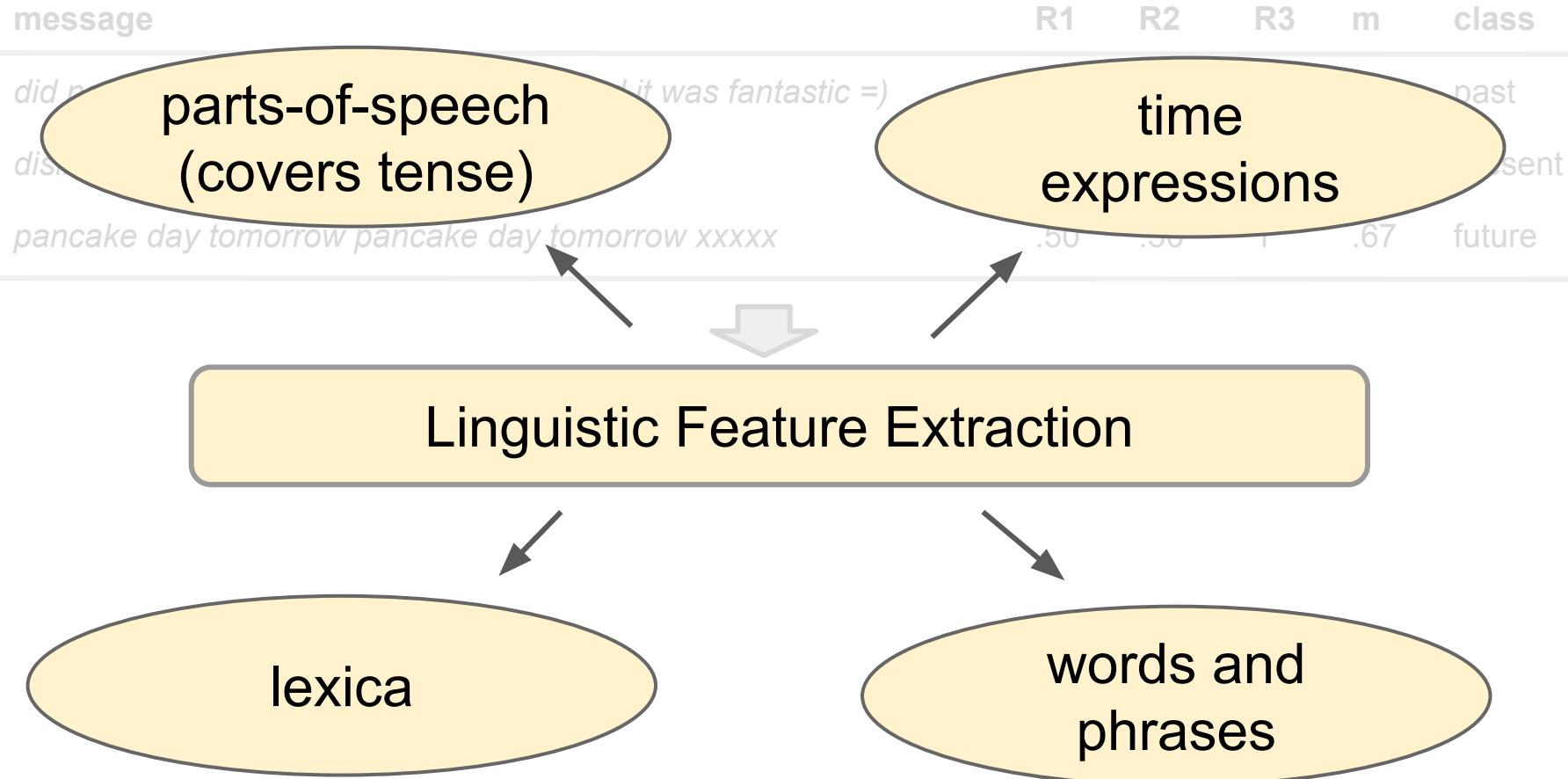
# Building a model



# Building a model



# Building a model



# Building a model

message	R1	R2	R3	m	class
<i>did nothing this morning but watch TV and it was fantastic =)</i>	-.67	-.50	-.50	-.55	past
<i>dislikes being sick.... and misses her bf</i>	0	0	0	0	present
<i>pancake day tomorrow pancake day tomorrow xxxxx</i>	.50	.50	1	.67	future



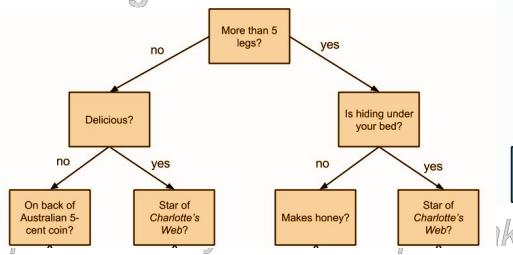
Linguistic Feature Extraction



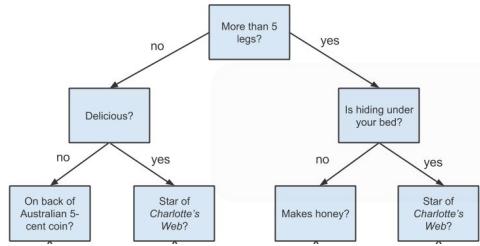
Learn Message-Level Model

# Building a model

message



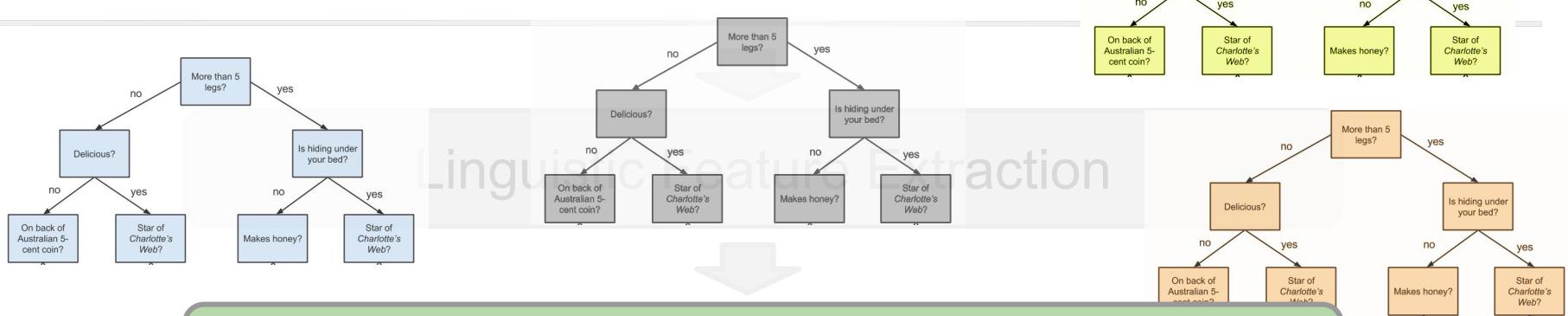
like day tomorrow xxxxx



linguistic features extraction



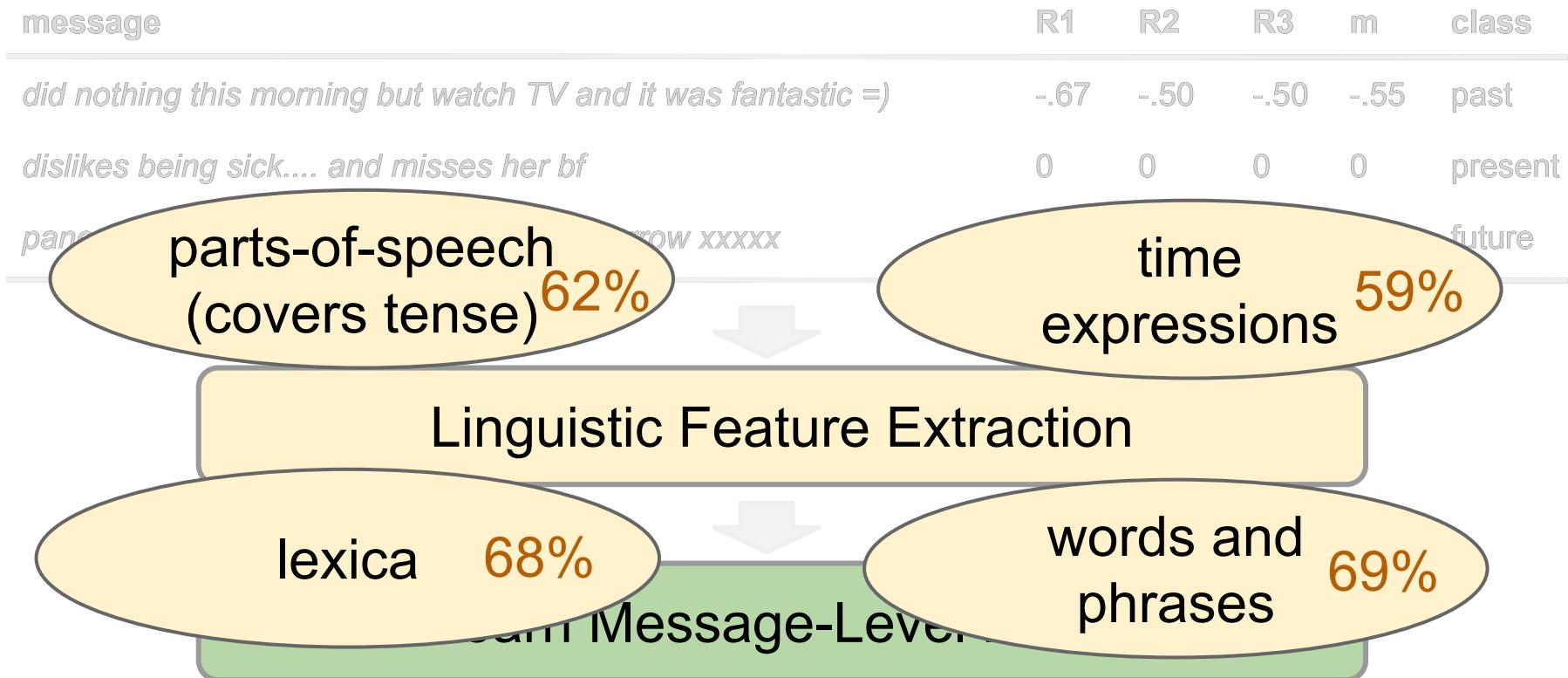
m class  
- .55 past  
it



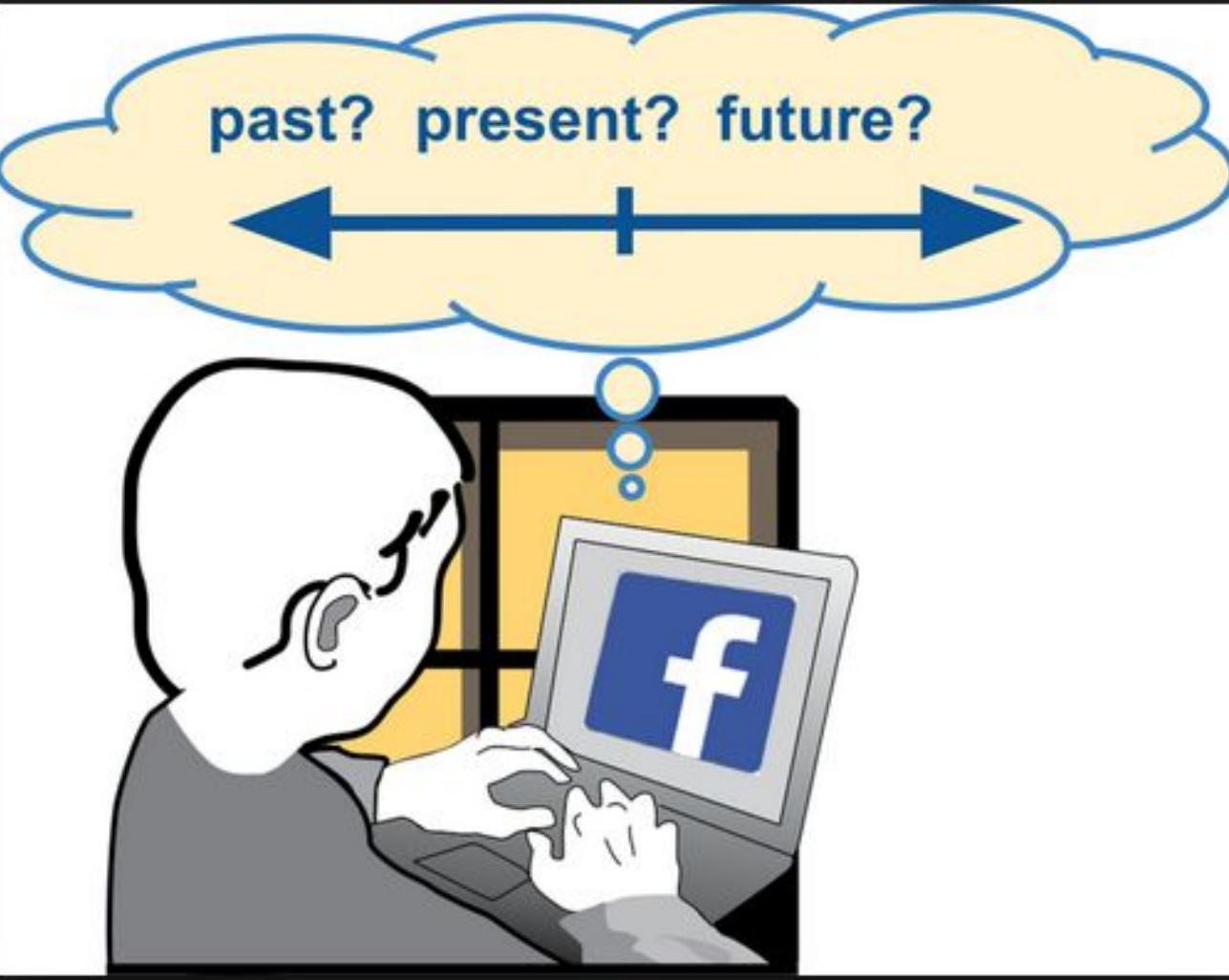
Learn Message-Level Model

Accuracy over a held-out set: 72%; baseline: 53%

# Building a model



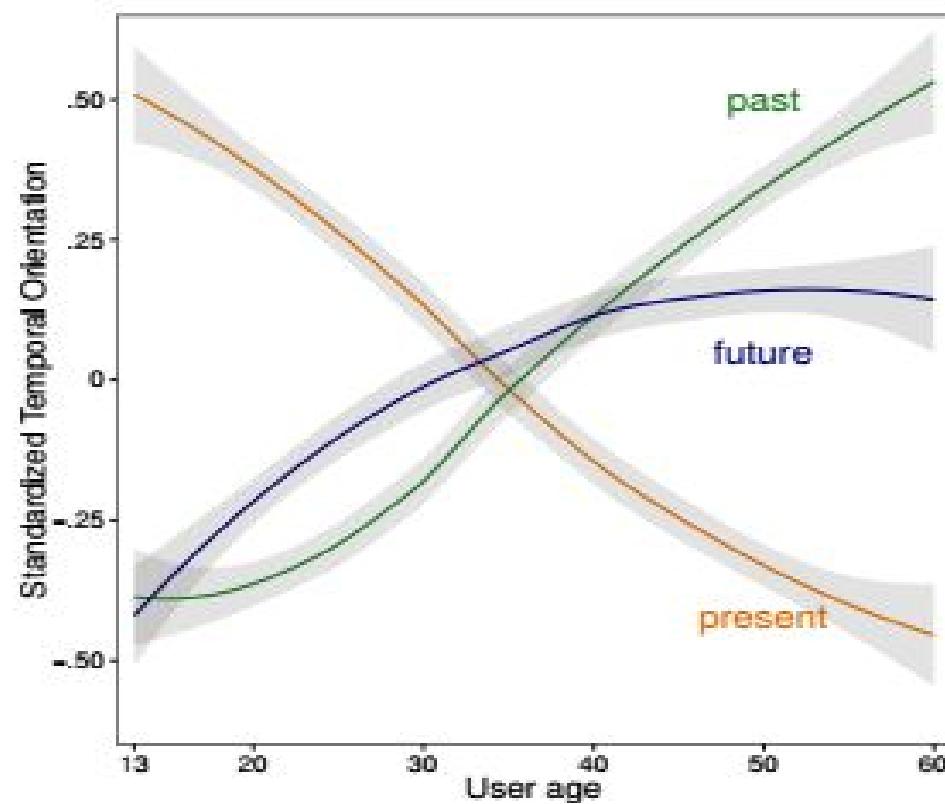
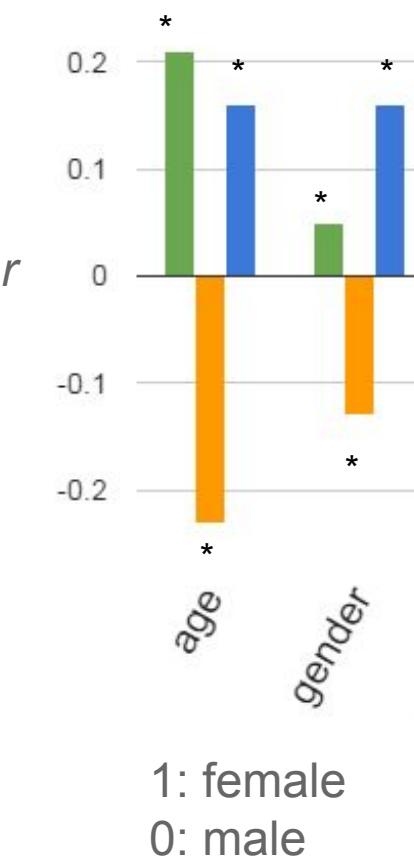
Accuracy over a held-out set: 72%; baseline: 53%



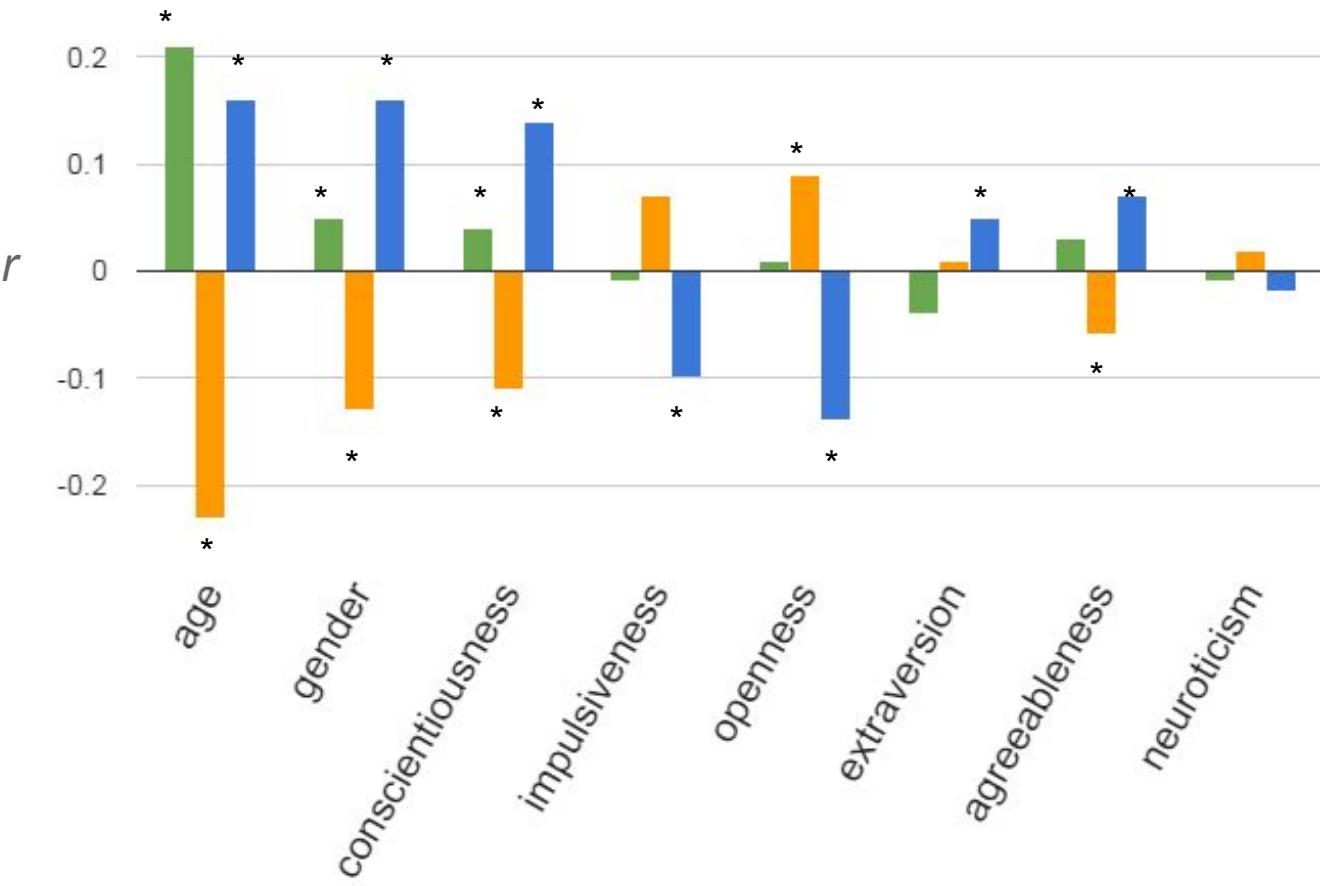
past? present? future?

1.3m statuses  
from  
4,833 Participants  
(3,240 AG Stratified)

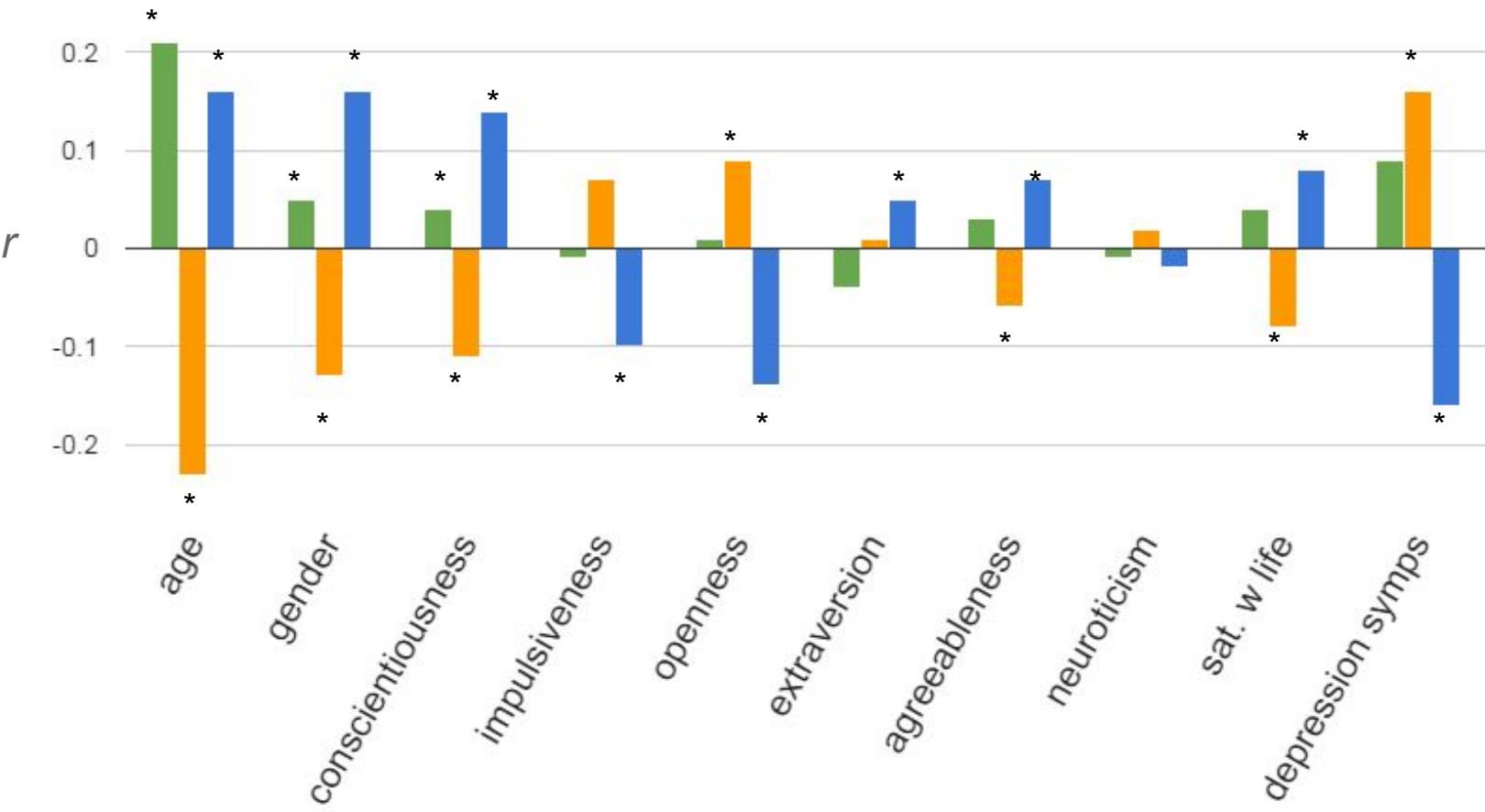
past      present      future



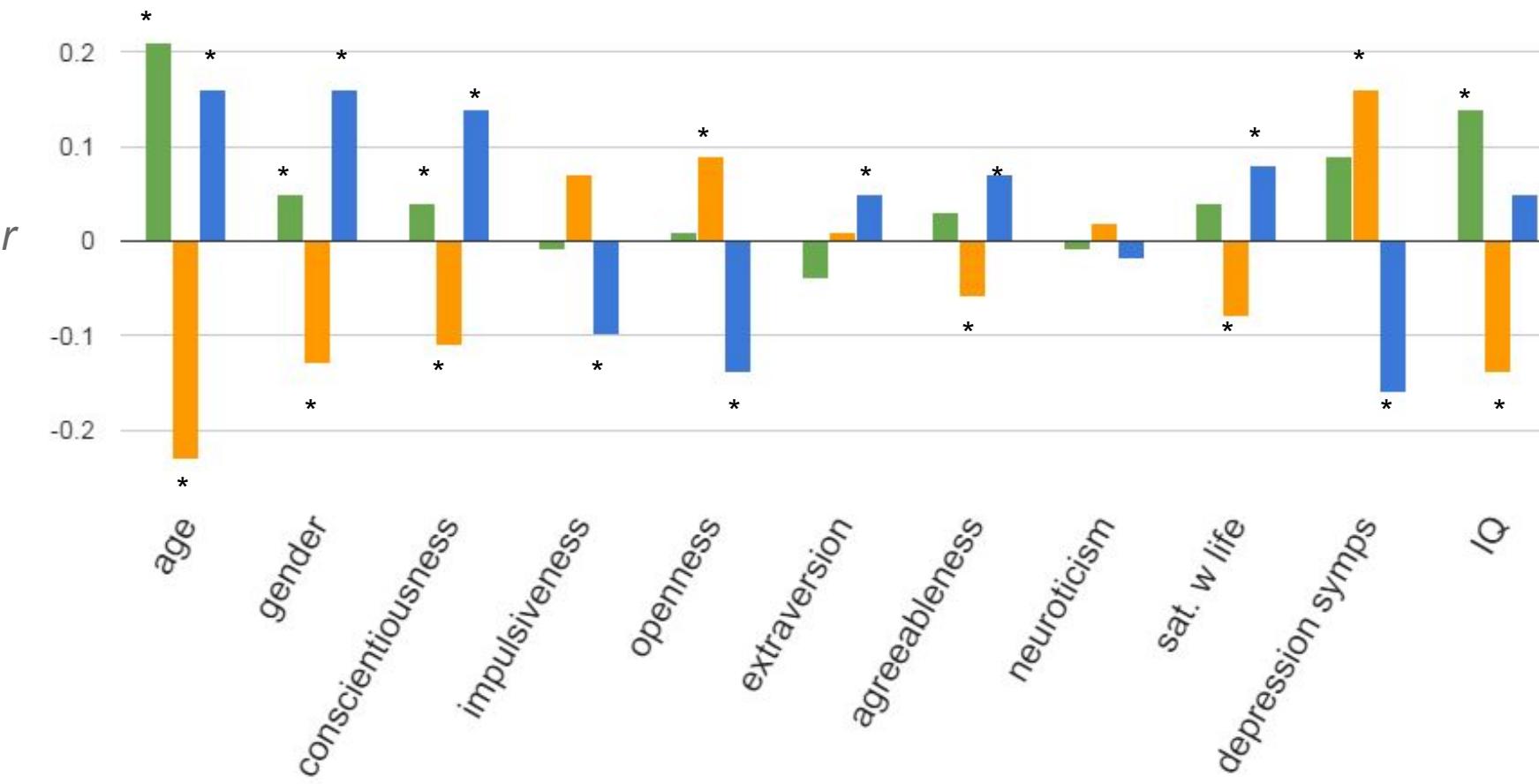
past      present      future

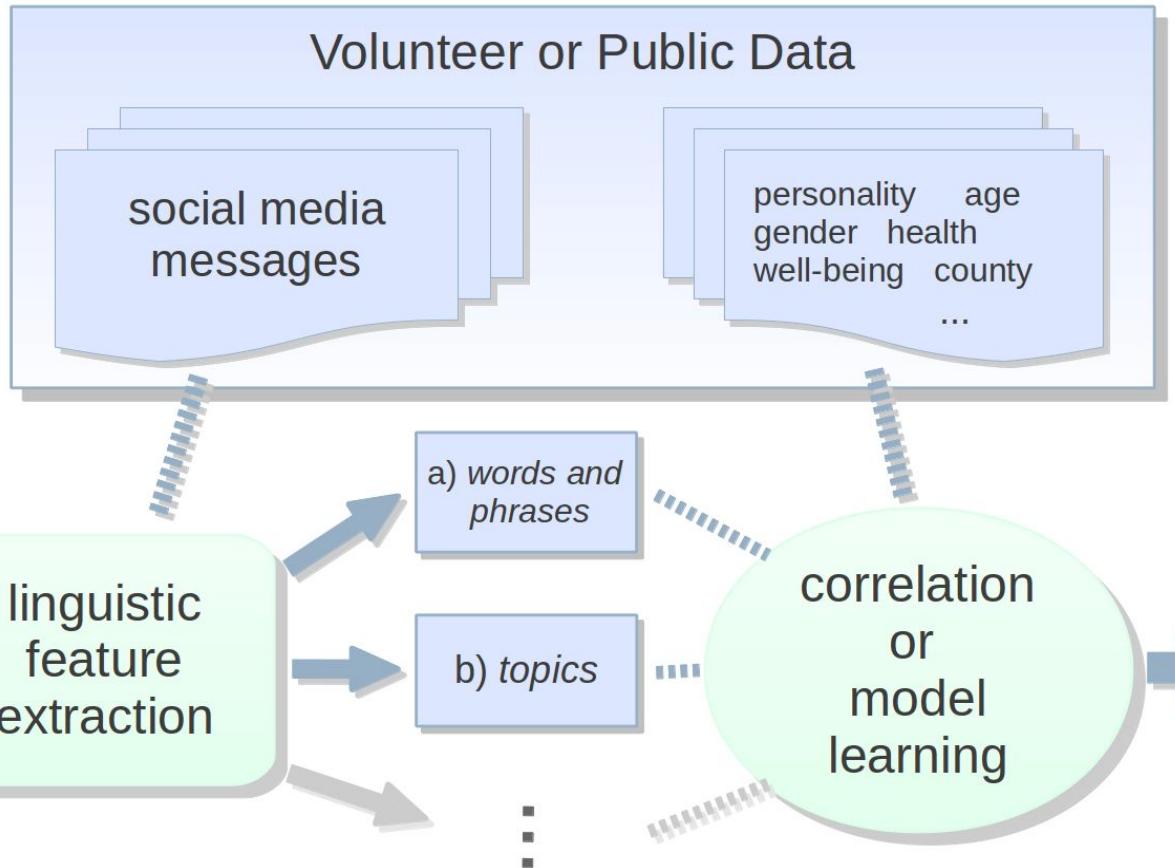


past      present      future



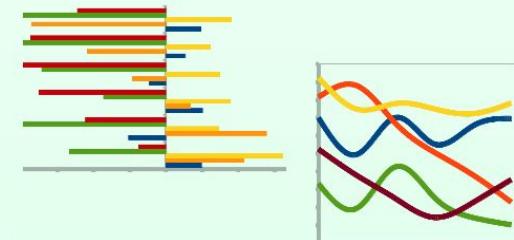
past      present      future





visualization or  
predictive model

:) love day everyone  
christmas wonderful lord family  
had today tomorrow great blessed beautiful  
friends happy awesome god's church  
thankful fun thank prayers weekend  
praise thanksgiving our psalm amazing  
for



# Resources



[dlatk.wwbp.org](http://dlatk.wwbp.org)

# Resources



[dlatk.wwbp.org](http://dlatk.wwbp.org)

```
#feature extraction:  
dlatkInterface.py  
-d <data> -t <corpus>  
-c <group_data_column>  
--add_ngrams -n 1 2 3  
--combine_feat_tables  
1to3gram
```

# Resources



[dlatk.wwbp.org](http://dlatk.wwbp.org)

```
#feature extraction:  
dlatkInterface.py  
-d <data> -t <corpus>  
-c <group_data_column>  
--add_ngrams -n 1 2 3  
--combine_feat_tables  
1to3gram
```

```
#run differential language analysis (linear regression)  
dlatkInterface.py -d <data> -t <corpus> -c user_id  
-f 'feat$cat_LIWC2007$msgs_xxx$user_id$16to16' --outcome_table  
blog_outcomes --group_freq_thresh 1000 --outcomes extraversion  
--controls age gender --output_name xxx_output --make_wordcloud
```

**extraversion** -- sociable, assertive, active, energetic, talkative, outgoing



Penn | World Well-Being Project | [wwbp.org](http://wwbp.org)

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*

# Resources



[dlatk.wwbp.org](http://dlatk.wwbp.org)

```
#feature extraction:  
dlatkInterface.py  
-d <data> -t <corpus>  
-c <group_data_column>  
--add_ngrams -n 1 2 3  
--combine_feat_tables  
1to3gram
```

```
#run differential language analysis (linear regression)  
dlatkInterface.py -d <data> -t <corpus> -c user_id  
-f 'feat$cat_LIWC2007$msgs_xxx$user_id$16to16' --outcome_table  
blog_outcomes --group_freq_thresh 1000 --outcomes extraversion  
--controls age gender --output_name xxx_output --make_wordcloud  
--[logistic_reg|AUC|IDP]
```

# Resources



[dlatk.wwbp.org](http://dlatk.wwbp.org)

```
#feature extraction:  
dlatkInterface.py  
-d <data> -t <corpus>  
-c <group_data_column>  
--add_ngrams -n 1 2 3  
--combine_feat_tables  
1to3gram
```

```
#run differential language analysis (linear regression)  
dlatkInterface.py -d <data> -t <corpus> -c user_id  
-f 'feat$cat_LIWC2007$msgs_xxx$user_id$16to16' --outcome_table  
blog_outcomes --group_freq_thresh 1000 --outcomes extraversion  
--controls age gender --output_name xxx_output --make_wordcloud
```

# Resources



[dlatk.wwbp.org](http://dlatk.wwbp.org)

```
#feature extraction:  
dlatkInterface.py  
-d <data> -t <corpus>  
-c <group_data_column>  
--add_ngrams -n 1 2 3  
--combine_feat_tables  
1to3gram
```

```
#create and cross validate a predictive model  
dlatkInterface.py -d <data> -t <corpus> -c user_id  
-f 'feat$cat_LIWC2007$msgs_xxx$user_id$16to16' --outcome_table  
blog_outcomes --group_freq_thresh 1000 --outcomes extraversion  
--controls age gender --model ridge --combo_test_regression
```

# Resources



*Intended for anyone  
(still under development)* →



**lexhub.org**



# Thank You!



...the largest data set of **who we are.**

**The most interesting results may never be hypothesized.**

has@cs.stonybrook.edu



Stony Brook University



# Thank You!



The most

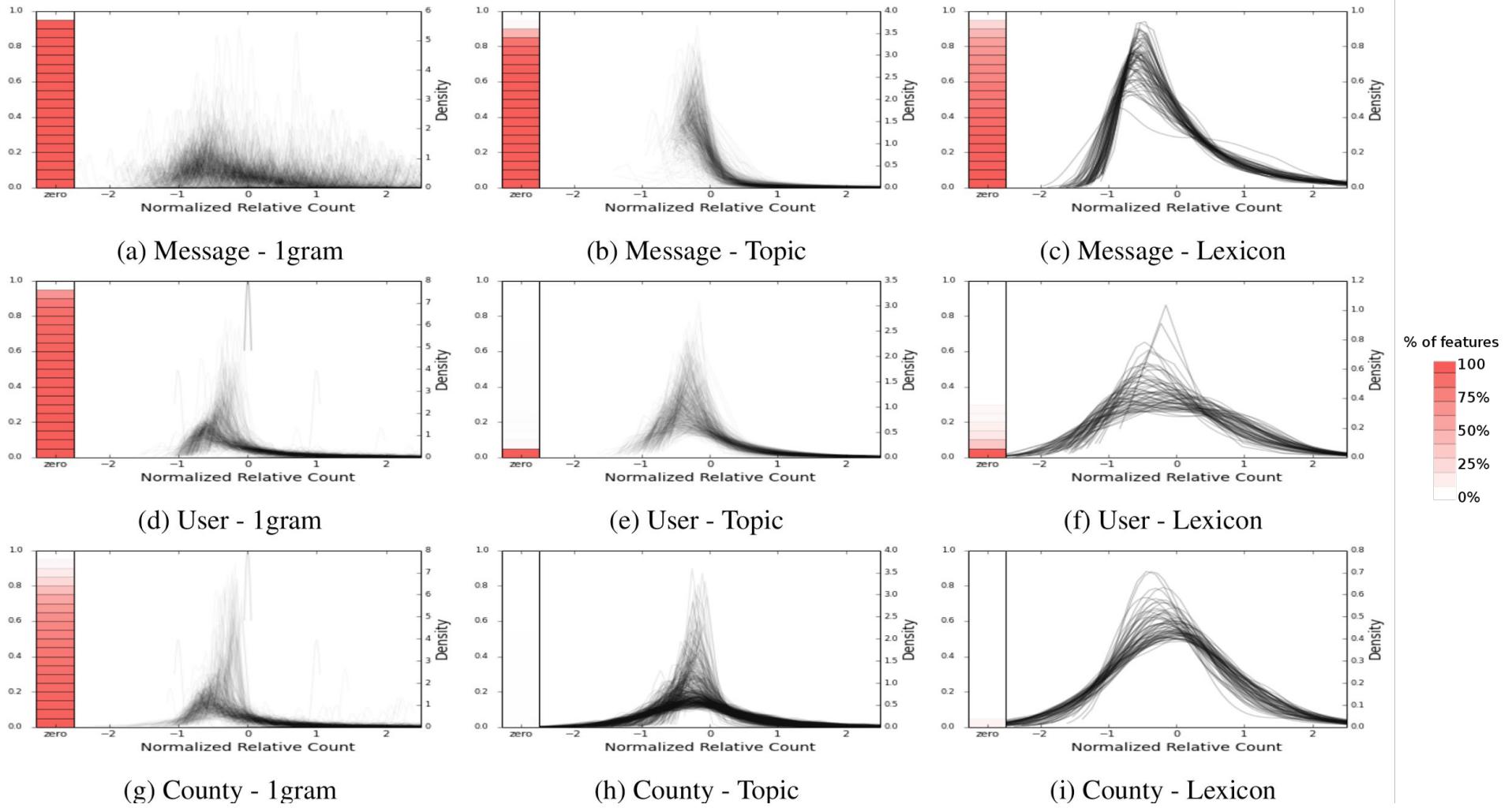


thesized.



Stony Brook University

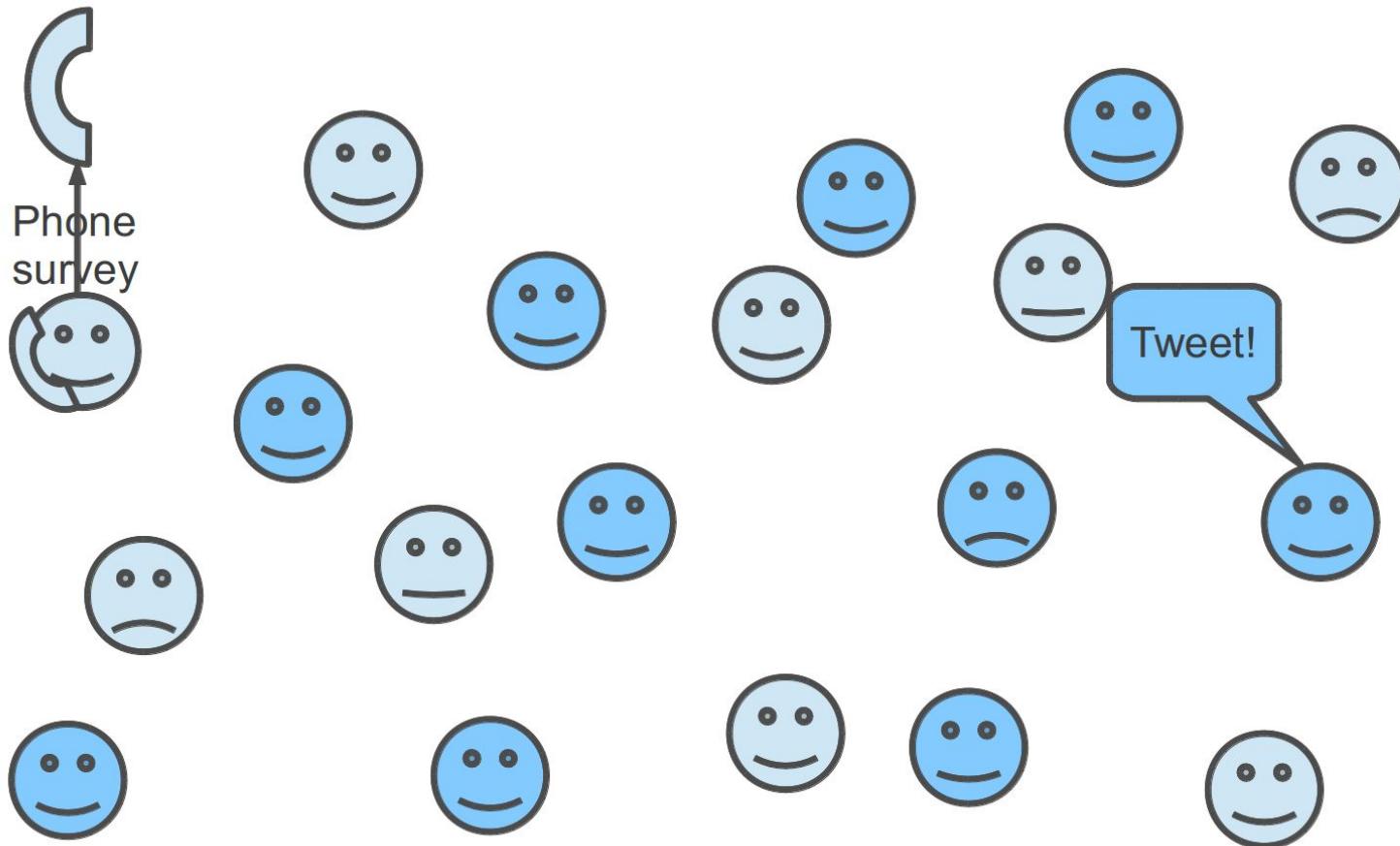
# **Supplemental Material**



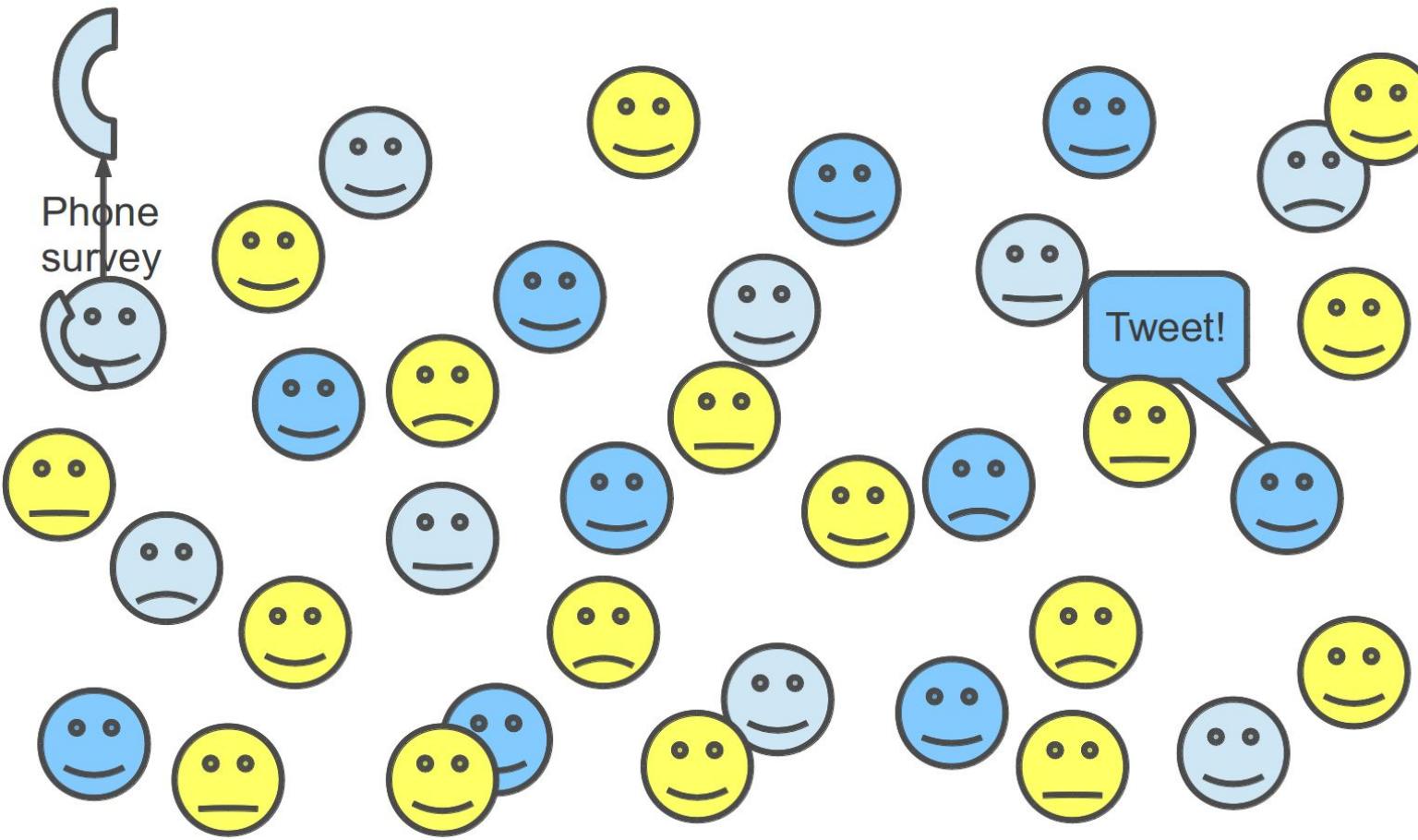
# Predicting based on a different sample



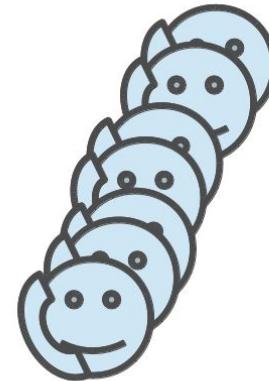
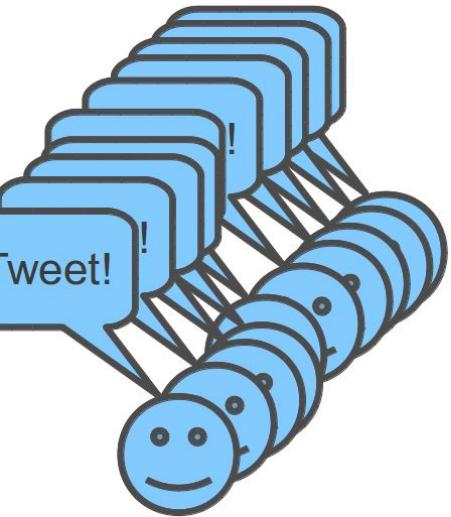
# Predicting based on a different sample



# Predicting based on a different sample

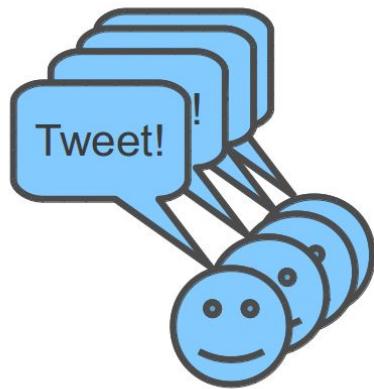
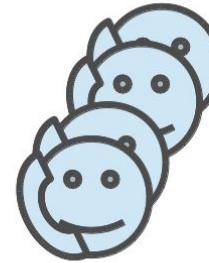
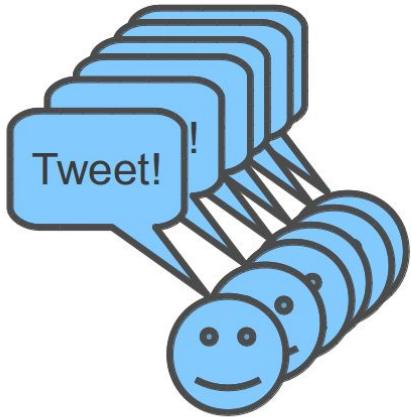


# Representative Sample?



Surveyed well-being  
from  
representative sample.

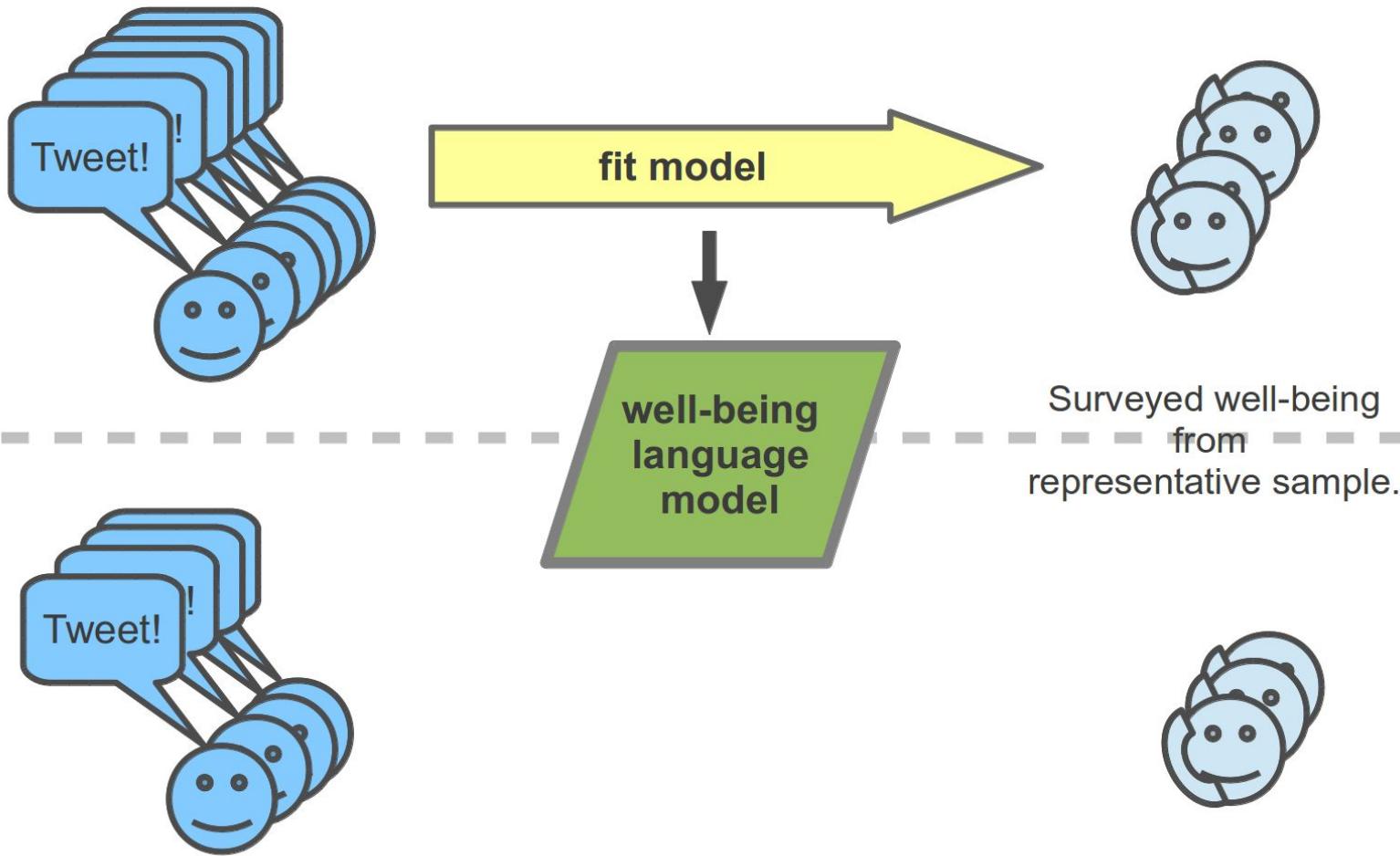
# Representative Sample?



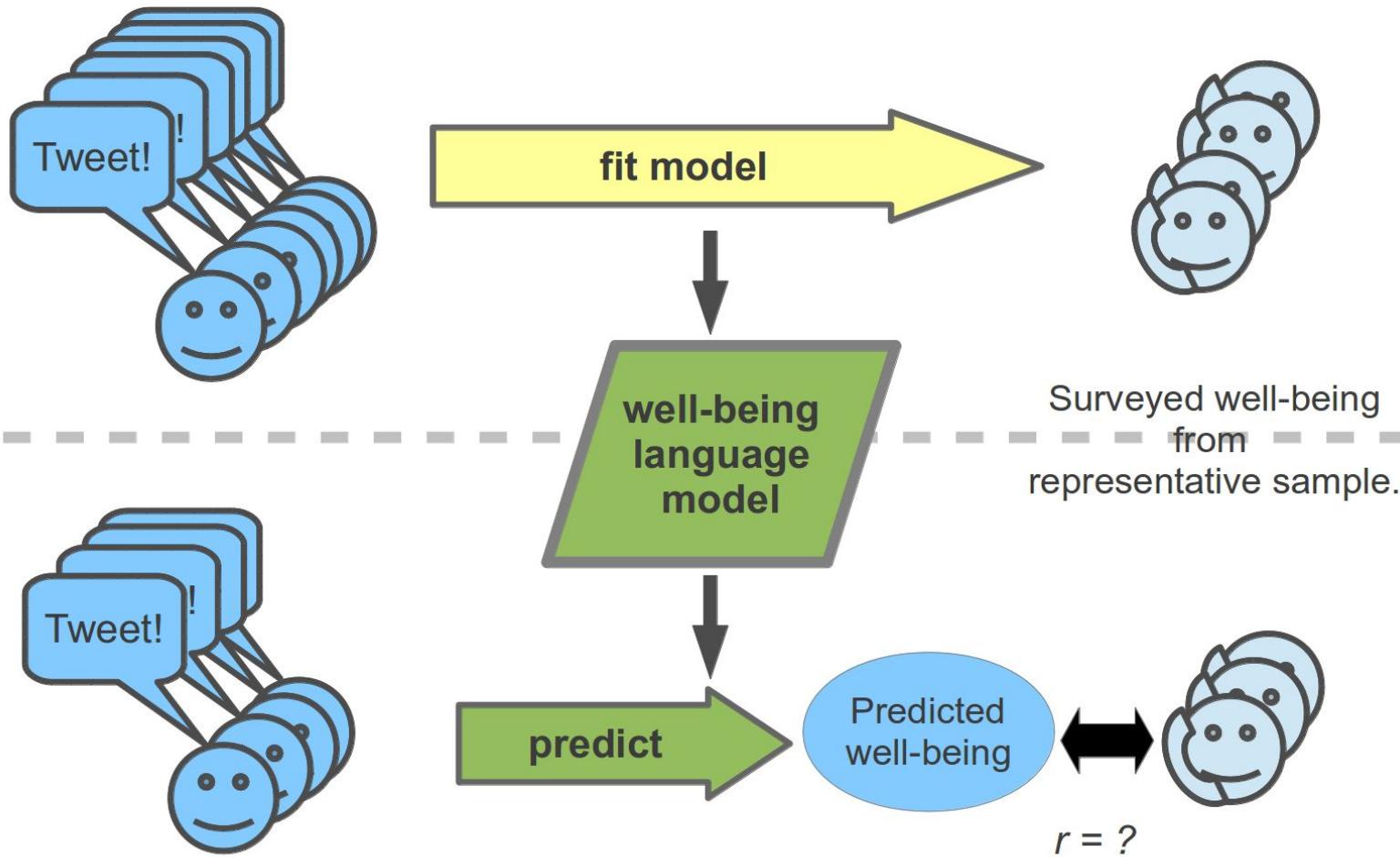
Surveyed well-being  
from  
representative sample.



# Representative Sample?

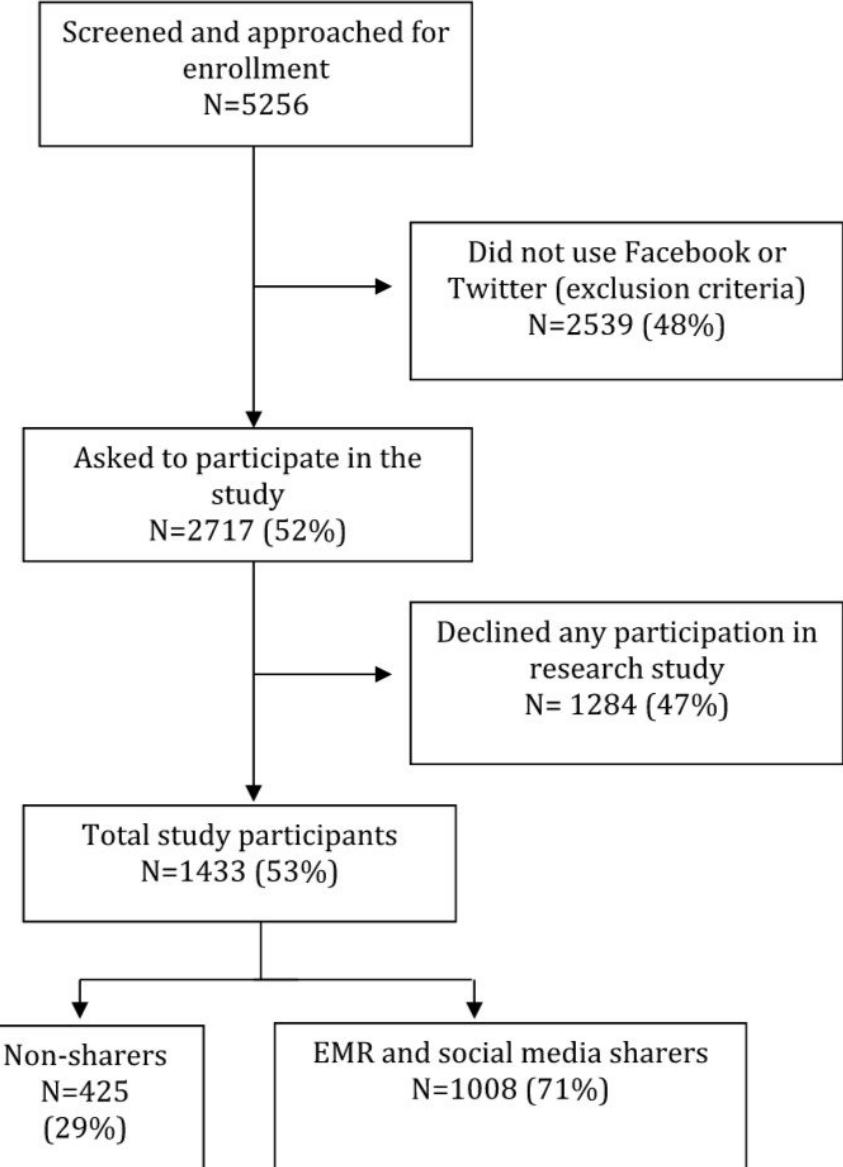


# Representative Sample?





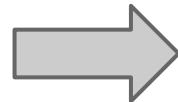
Kevin A Padrez, L Ungar, HA Schwartz, RJ Smith, S Hill, T Antanavicius, DM Brown, P Crutchley, D Asch, Merchant. Linking social media and medical record data: a study of adults presenting to an academic, urban emergency department. *BMJ Quality & Safety* | 2015



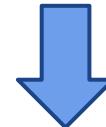
ICD-9 Diagnosis Code	Terms searched	n (%) of patients with diagnosis who used term	n (%) of patients without diagnosis who used term
Abdominal Pain (789)	abdominal pain, stomach pain, belly pain, tummy pain, stomach hurts, belly hurts, tummy hurts, tummyache, stomachache, bellyache	81 (21)	21 (8) <sup>a</sup>
Nausea/Vomiting (787)	nausea, vomiting, vomit, throwing up, spitting up, threw up, puke, puked, vomited	101 (29)	69 (22) <sup>c</sup>
Headache (339&784.0)	headache, migraine, head hurts	141 (59)	192 (46) <sup>b</sup>
Pain in limb (729.5)	leg hurts, arm hurts, finger hurts, toe hurts	5 (3)	5 (1)
UTI (599.0)	uti, urinary tract infection	1 (1)	4 (1)
Back pain (724.1&724.2&724.3&724.4&724.5)	back pain, backache, back hurts	29 (15)	51 (11)
Cough (786.2)	cough, coughing, coughed	40 (26)	109 (22)
Normal delivery (V27.0)	giving birth, gave birth	62 (33)	148 (10) <sup>a</sup>
Anemia (280&281&282&283&284&285)	anemia	3 (2)	2 (0) <sup>c</sup>
Dizziness (780.4)	dizzy, dizziness, vertigo	28 (22)	79 (15)
Asthma (493)	asthma	36 (28)	36 (7) <sup>a</sup>
Acute URI (465)	caught a cold, have a cold	7 (7)	24 (4)
Throat Pain (784.1)	sore throat, throat hurts	26 (24)	62 (11) <sup>a</sup>
Depression (311)	depression, depressed	35 (38)	169 (30)

# **Big Data** and Emergency Care: Social Media as a Source for Data-Driven Health Prediction and Insight

The Big Data Promise: enabling links with unprecedented amount of **new information**.



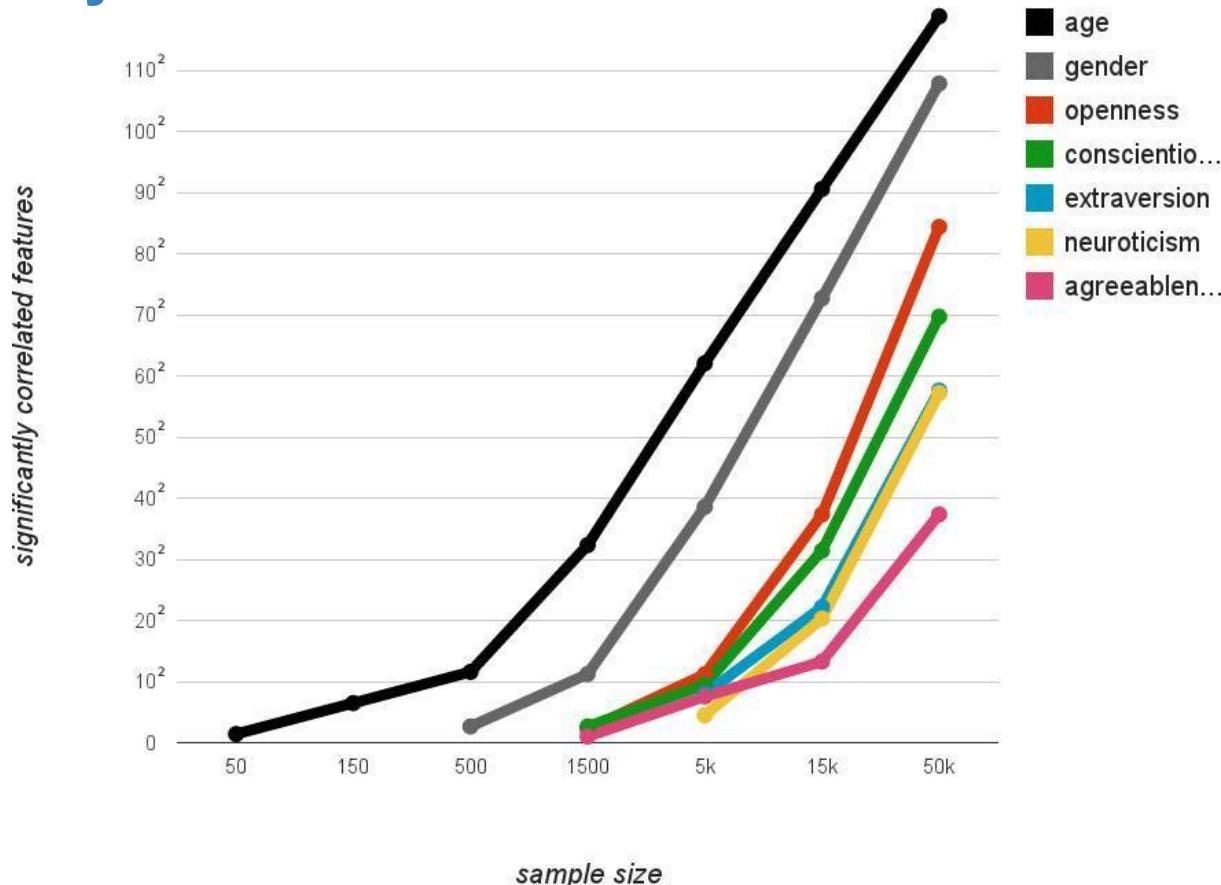
Better health care



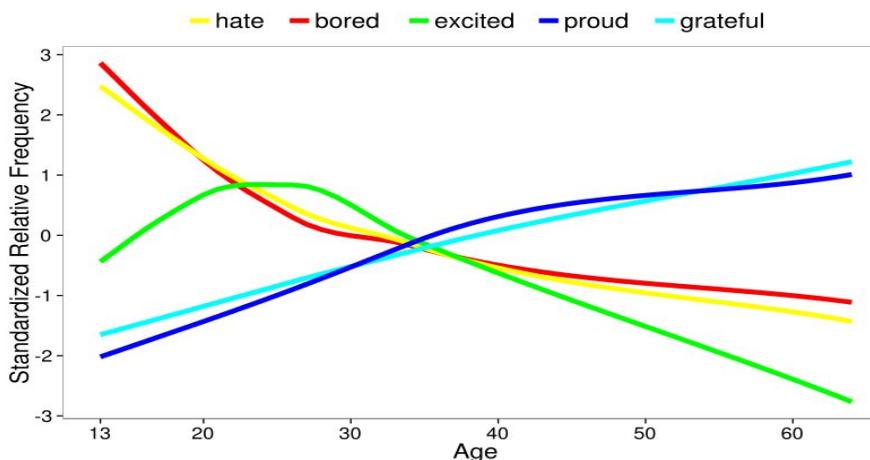
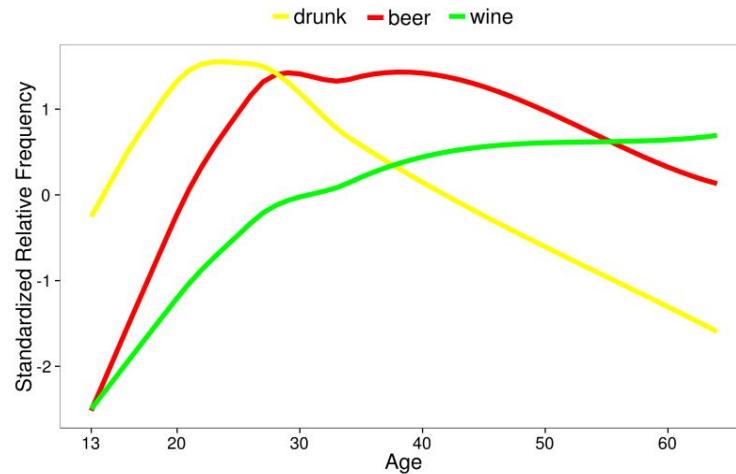
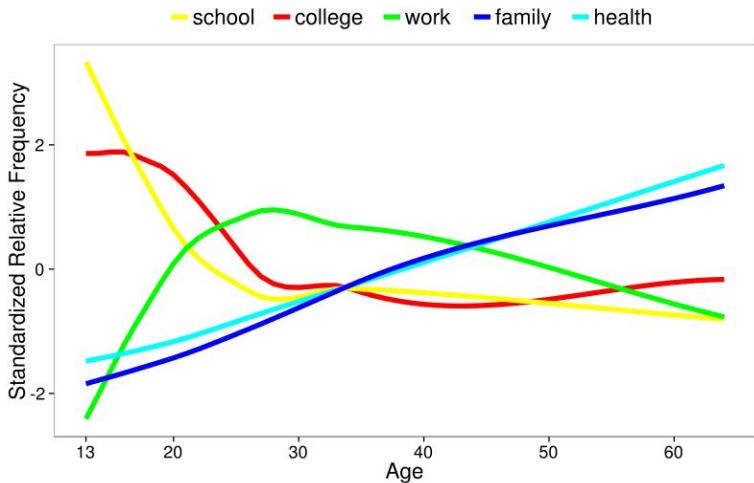
Personalized Health: **Patients' Everyday Lives**

# Individual Traits in Facebook

## Power Analyses



# Development



# Individual Well-Being

tonight tomorrow  
excited woohooo  
super pumped  
stoked soooo upcoming  
psyched bummed prom

thankful truely  
wonderful boyfriend  
helped amazing grateful  
family lucky blessed daughter  
friends loving supportive  
husband

pissed pissing wtf fucked  
bullshit >: shit fuck  
bitch asshole pisses shitty  
goddamn fucking  
piss

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

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

bored text  
bored boring bore  
entertainment insanely stiff entertained  
extremely entertain  
boredom yawn hmu  
incredibly sooooooo

satisfaction with life

# Community Well Being

sailing  
ship  
sail  
ocean  
water  
sea  
deep  
waves  
boat  
sink  
drowning  
wave  
swim  
sinking

$r' = .15$

brick  
construction  
destroy  
snowman  
blocks  
noah  
walls  
takes  
foundation  
built  
build  
bridge  
ground  
castle

$r' = .13$

valley  
falls  
trail  
creek  
grand  
park  
total  
springs  
headed  
hike  
lake  
hiking  
forest  
river

$r' = .12$

boat  
fireworks  
city  
camping  
river  
swimming  
headed  
lake  
skiing  
blast  
tube  
michigan  
fishing  
water  
salt

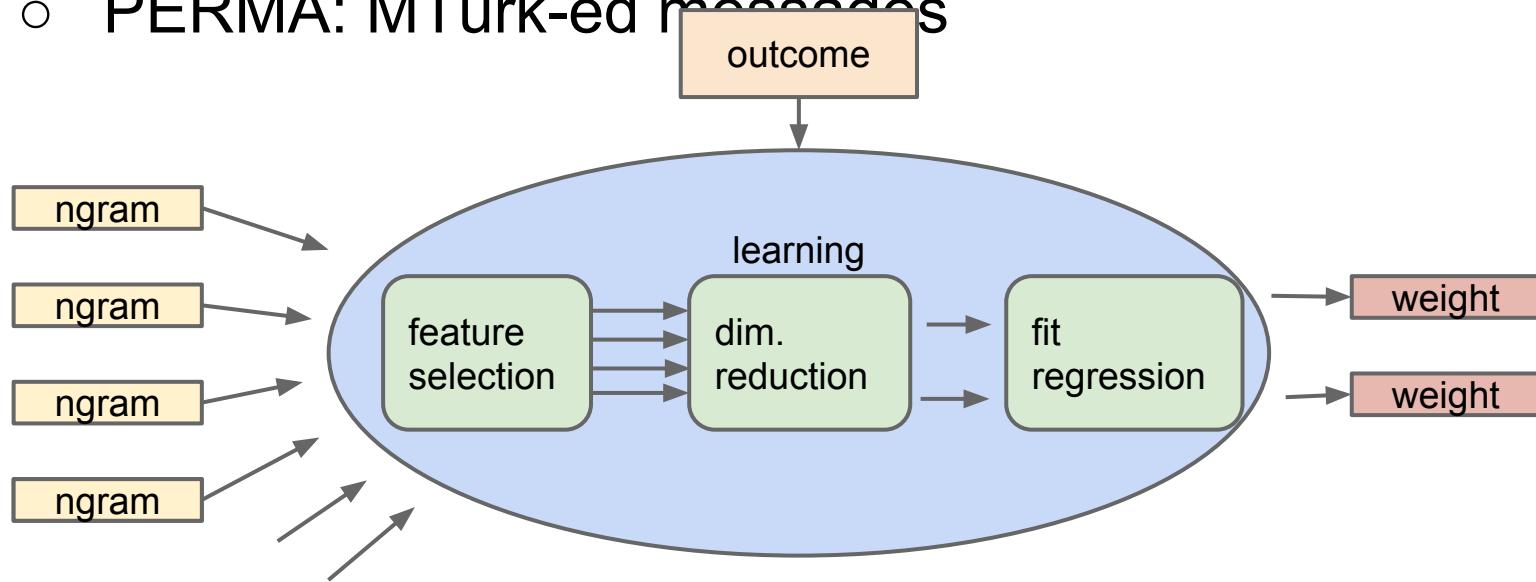
$r' = .12$

peeps  
awsome  
fabulous  
holiday  
hope  
weekend  
tgif  
great  
wonderful  
fantastic  
enjoy  
safe  
hopes  
enjoyed  
fab

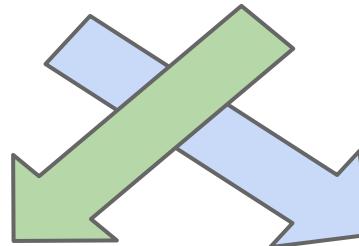
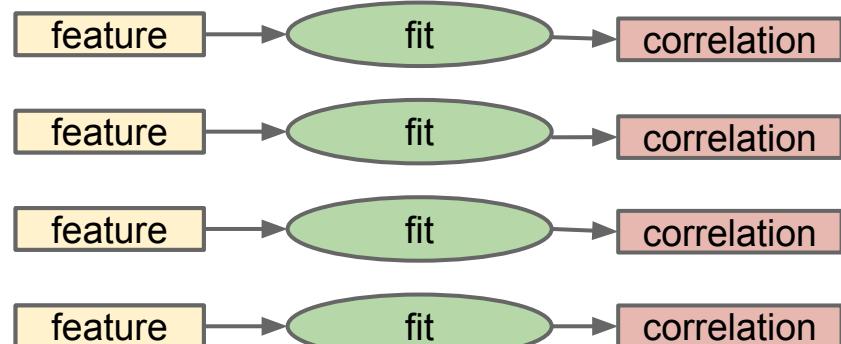
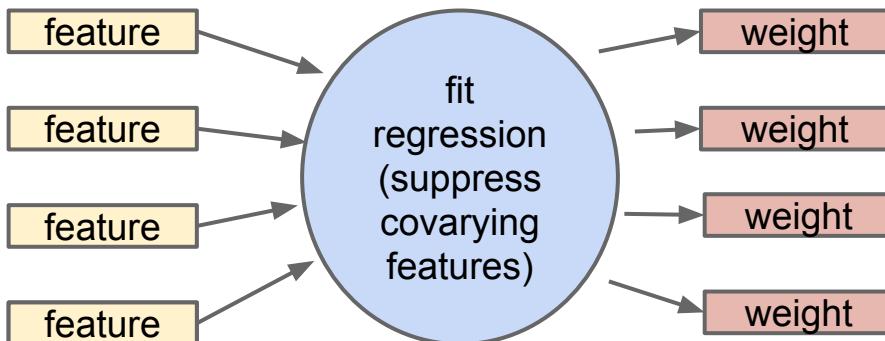
$r' = .12$

# Generating Lexica from **\*most\*** Supervised n-gram Models

- Generalize multi-variate regression model into lexica.
- Works at multiple levels:
  - Hand annotated messages or users
    - OCEAN: User-level Cambridge data set
    - PERMA: MTurk-ed messages



# Multivariate or Univariate for Insights?



A word cloud composed of various words in different sizes and shades of gray. The most prominent words include "lord", "happy", "beautiful", "great", "amazing", "wonderful", "family", "excited", "love", "thank", "god", "blessed", "for", "prayers", and "friends".

A word cloud composed of various words in different sizes and shades of gray. The most prominent words include "sara", "job", "kinds", "climb", "awesomeness", "but", "m", "spontaneous", "amazing", "hands", "sticked", "highest", "scary", "make-up", "bags", "sunrise", "fluffy", "awesome", "banana", "flaws", "everyone", "ashley", "dick", and "wow".

# Individual Well-Being: message to user-level

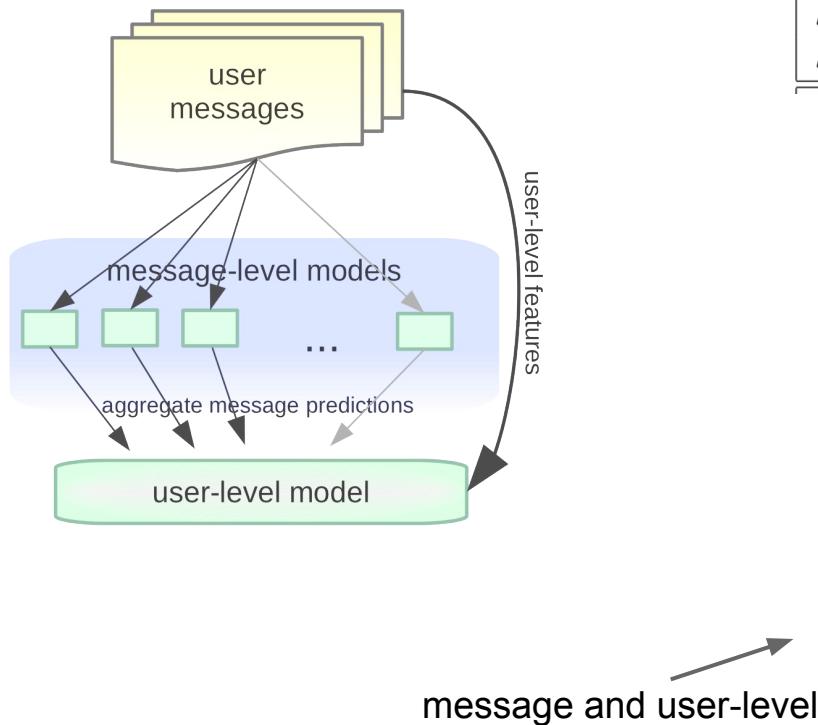
tonight tomorrow  
excited woohooo  
super pumped  
stoked soooo upcoming  
psyched bummed prom

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

thankful truely  
wonderful boyfriend  
helped amazing grateful  
family lucky blessed daughter  
friends loving supportive  
husband

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

bored text  
entertainment boring bore  
insanelytiff entertained  
extremelyentertain boredom yawn  
incredibly hmu sooooooo  
pissed wtf fucked  
bullshit >: shit fuck  
bitch asshole pisses shitty  
goddamn fucking  
piss



baselines	<i>r</i>
(mean)	.000
lexica: <i>GNH</i>	.210
lexica: <i>Hedonometer</i>	.108

writing size

sample size / populations  
(Gosling 2004; 2010)

self-descriptive variables

## *Why Social Media and Language?*

unobtrusive

longitudinal / look back in time

potential for real-time

often personal /  
everyday concerns