

SMHD: A Large-Scale Resource for Exploring Online Language Usage for Multiple Mental Health Conditions

Arman Cohan, Bart Desmet, Andrew Yates,
Luca Soldaini, **Sean MacAvaney**, Nazli Goharian

COLING 2018 “Area Chair Favorite Paper”



4.2% of Americans suffer from
Serious Mental Illnesses (SMI).

Only 65% of those with SMI received
any treatment in past year.

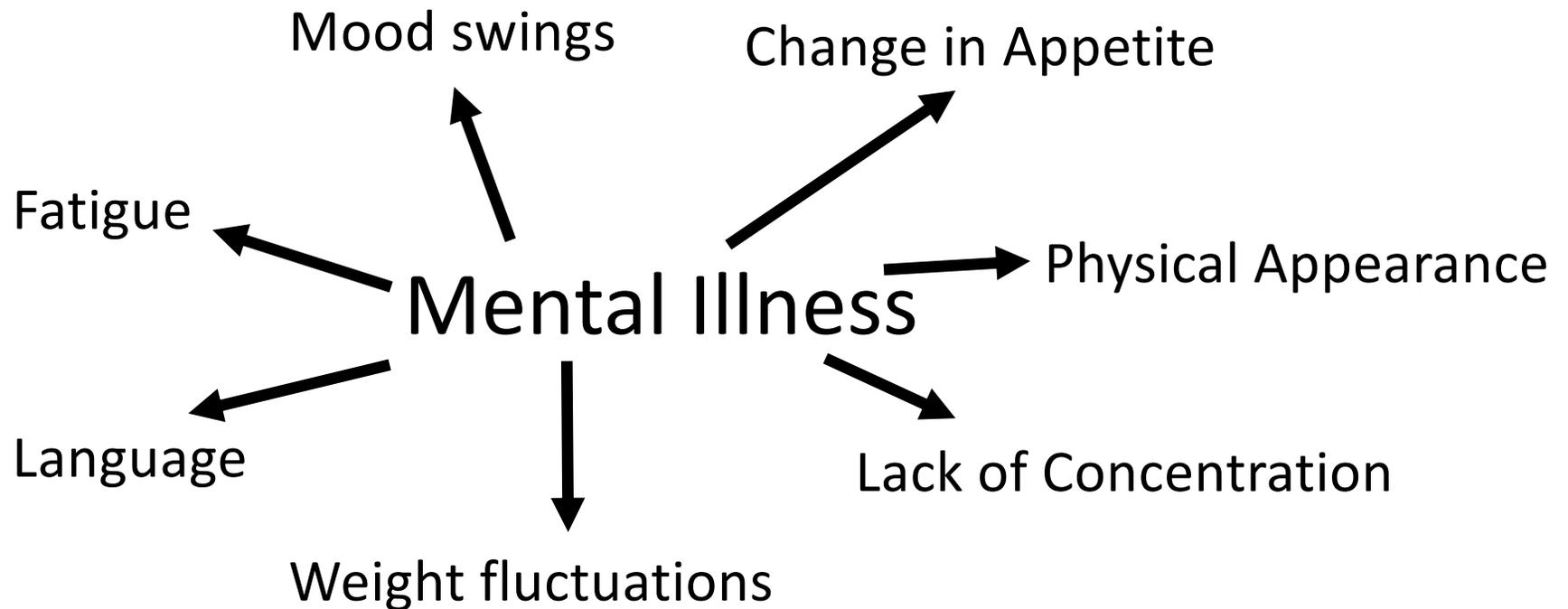
US Health and Human Services, 2016

Suicide is the 10th leading cause
of death in the US.

45,000 Americans take their own
life each year.

CDC Vital Statistics Report, 2016

Doctors and mental health providers evaluate many possible symptoms to make a diagnosis.

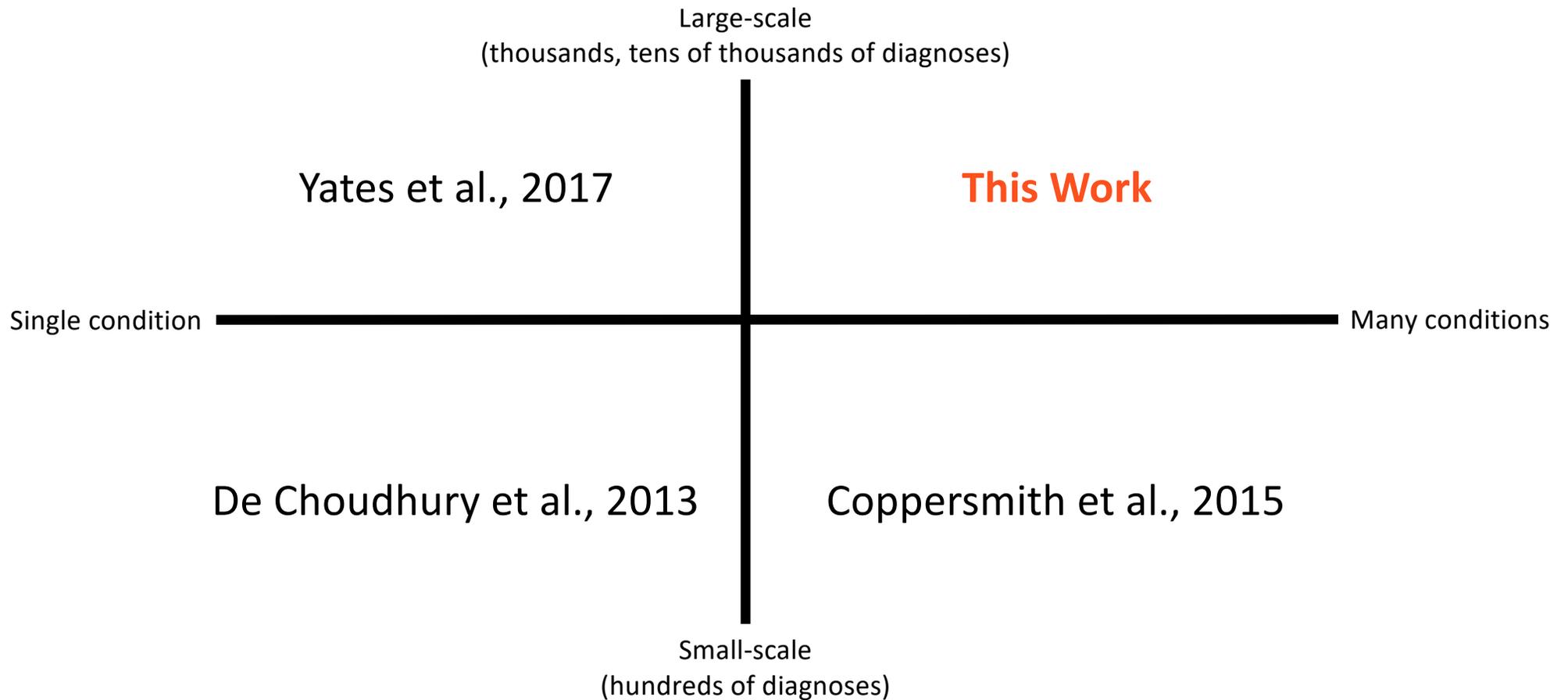


Let's try to better understand how people's language is affected by mental illness.

Goals:

- **Large-scale** – people show symptoms in a variety of ways; the more data the better
- **High quality** – a low false-positive rate for diagnoses is important
- **Multiple conditions** – conditions do not exist in isolation; the interaction between conditions is valuable to explore

Previous work has not explored this space.



We present:

SMHD: a dataset of **S**elf-reported **M**ental **H**ealth **D**iagnoses

- Large-scale labeled data from Reddit
- 37k diagnosed users from high-precision diagnosis statements
- Carefully-selected “control” users for each
- 9 mental health conditions
 - ADHD, anxiety, autism spectrum disorder, bipolar disorder, depression, eating disorder, OCD, PTSD,
- Available to researchers with signed DUA



Source of Relevance:
Self-reported diagnosis statements

was officially diagnosed with ADHD last year.

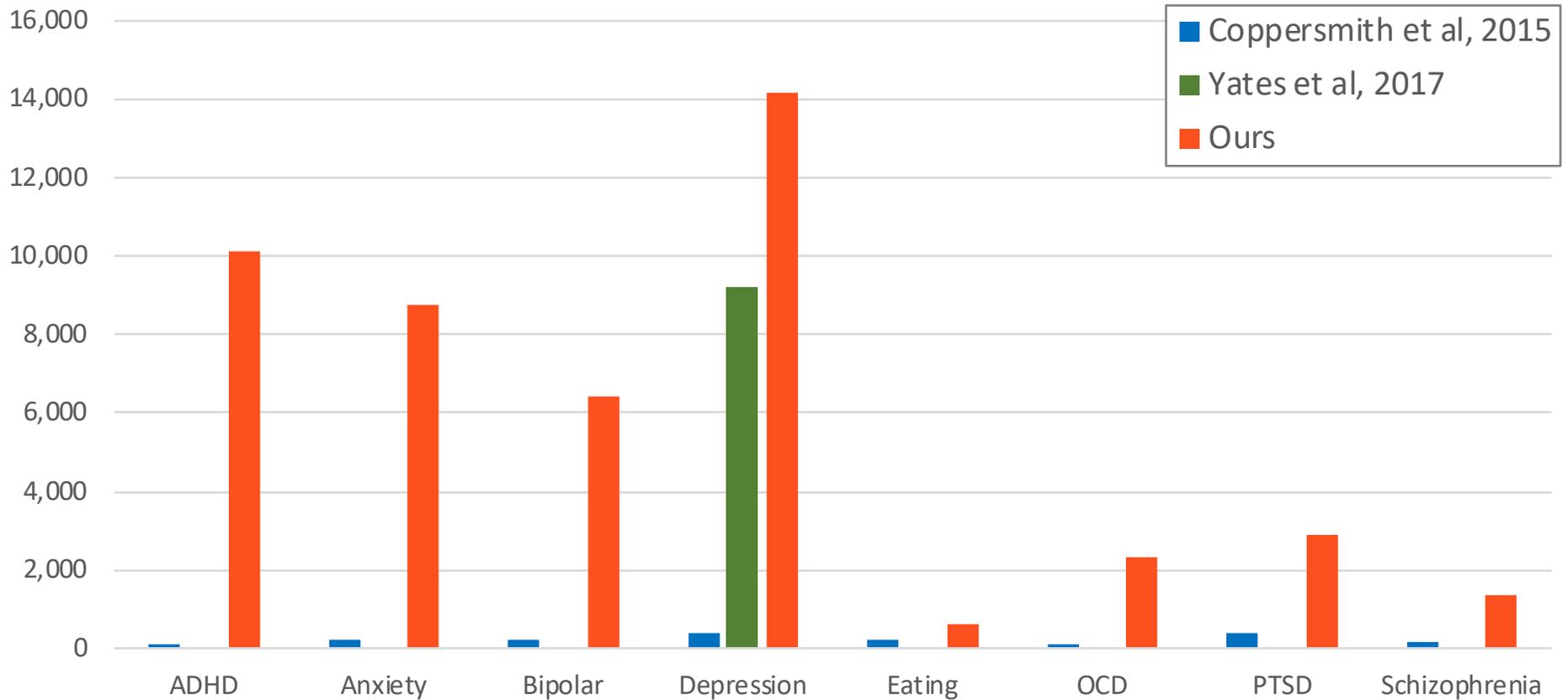
I have a diagnosed history of PTSD.

my dr just diagnosed me as schizo.

Self-reported diagnosis rules

1. Diagnosis pattern patterns, e.g. “i am diagnosed with”, “they recently diagnosed me”, “have [condition] diagnosed”...
2. Condition names matching terms found in dictionaries (MedSyn, Behavioral), manually adjusted to remove hypernyms or general terms (e.g., add), and add terms found when manually inspecting the data
3. At least 50 posts unrelated to mental health (no mental health terminology, and not in mental health sub-forums)
 - To remove bias, we only analyze these posts

Ours much larger than previous datasets.



Does not include autism (ours=2,911) borderline personality (Coppersmith=101) or Seasonal Affective (Coppersmith=100).

Control Users

- Choose 9 control users per diagnosed user
- Stringent criteria:
 - Must not have any posts that contain mental health language or be in a mental health sub-forum.
 - Must have at least 50 posts.
 - Most have posted in the same sub-forum

Dataset Characteristics

condition	posts		tokens		characters
	per user	total	per post	total	per post
control	310.0 (157.8)	115,669k	26.2 (48.3)	3,031.6M	133.9 (252.9)
depression	162.2 (84.2)	1,272k	45.1 (80.0)	57.4M	227.5 (406.9)
adhd	164.7 (83.6)	872k	46.5 (82.7)	40.5M	237.5 (433.5)
anxiety	159.7 (83.0)	795k	46.4 (83.0)	36.9M	233.9 (422.8)
bipolar	157.6 (82.4)	575k	45.5 (86.5)	26.2M	230.6 (447.0)
ptsd	160.7 (84.7)	258k	53.1 (114.0)	13.7M	267.8 (581.7)
autism	168.3 (84.5)	248k	46.5 (82.3)	11.6M	237.9 (434.0)
ocd	158.8 (81.4)	203k	46.4 (90.1)	9.4M	234.2 (459.5)
schizophrenia	157.3 (80.5)	123k	49.2 (105.6)	6.1M	253.8 (566.6)
eating	161.4 (81.0)	53k	46.3 (73.7)	2.5M	232.6 (372.8)

Textual Analysis

- Linguistic Inquiry and Word Count (LIWC) categories
- Compare scores in categories using Welch's t-test ($p < 0.001$)
- P-values adjusted with Bonferroni correction



LIWC category	depression	adhd	anxiety	bipolar	ptsd	autism	ocd	schizophrenia	eating
<i>Summary Language Variables</i>									
Clout	-0.06 [‡]	-	-0.1 [‡]	-	-	-	-	-	-
Authentic	0.2 [‡]	0.15 [‡]	0.22 [‡]	0.18 [‡]	0.21 [‡]	0.14 [‡]	0.23 [‡]	-	0.24 [*]
WPS	0.08 [‡]	0.1 [‡]	0.08 [‡]	0.06 [‡]	0.1 [‡]	0.12 [‡]	0.08 [‡]	-	-
Dictionary words	0.27 [‡]	0.22 [‡]	0.28 [‡]	0.24 [‡]	0.31 [‡]	0.2 [‡]	0.28 [‡]	0.22 [‡]	0.3 [‡]
Total function words	0.27 [‡]	0.21 [‡]	0.28 [‡]	0.24 [‡]	0.3 [‡]	0.26 [‡]	0.28 [‡]	0.23 [‡]	0.27 [‡]
Total pronouns	0.22 [‡]	0.14 [‡]	0.24 [‡]	0.2 [‡]	0.25 [‡]	0.17 [‡]	0.26 [‡]	0.18 [‡]	0.26 [‡]
Personal pronouns	0.23 [‡]	0.14 [‡]	0.26 [‡]	0.21 [‡]	0.22 [‡]	0.14 [‡]	0.23 [‡]	0.2 [‡]	0.27 [‡]
1st pers singular	0.23 [‡]	0.16 [‡]	0.28 [‡]	0.22 [‡]	0.23 [‡]	0.17 [‡]	0.26 [‡]	0.17 [‡]	0.28 [‡]
3rd pers singular	0.09 [‡]	-	0.1 [‡]	0.08 [*]	0.17 [*]	-	-	-	-
Impersonal pronouns	0.06 [‡]	0.05 [*]	0.07 [‡]	-	0.11 [‡]	0.09 [*]	0.13 [‡]	-	-
Prepositions	0.12 [‡]	0.12 [‡]	0.12 [‡]	0.12 [‡]	0.16 [‡]	0.12 [‡]	0.11 [‡]	-	-
Auxiliary verbs	0.12 [‡]	0.1 [‡]	0.14 [‡]	0.11 [‡]	0.14 [‡]	0.15 [‡]	0.13 [‡]	-	-
Common Adverbs	0.09 [‡]	0.07 [‡]	0.08 [‡]	0.06 [*]	-	-	-	-	-
Conjunctions	0.17 [‡]	0.14 [‡]	0.18 [‡]	0.16 [‡]	0.17 [‡]	0.13 [‡]	0.11 [‡]	-	0.3 [‡]
<i>Other Grammar</i>									
Common verbs	0.15 [‡]	0.11 [‡]	0.17 [‡]	0.13 [‡]	0.15 [‡]	0.13 [‡]	0.19 [‡]	-	-
Numbers	-0.1 [‡]	-0.09 [‡]	-0.11 [‡]	-0.09 [‡]	-0.11 [‡]	-	-0.1 [‡]	-	-0.13 [‡]
<i>Psychological Variables</i>									
Positive emotion	-	-	-	-	-	-0.08 [*]	-	-	-
Anxiety	0.07 [‡]	-	0.07 [*]	-	-	-	-	-	-
Social processes	0.11 [‡]	0.07 [‡]	0.11 [‡]	0.1 [‡]	0.15 [‡]	-	-	-	-
Family	0.06 [‡]	-	0.06 [*]	-	0.12 [‡]	-	-	-	-
Female references	0.13 [‡]	0.07 [‡]	0.1 [‡]	0.13 [‡]	0.22 [‡]	-	-	-	-
Cognitive processes	0.12 [‡]	0.13 [‡]	0.14 [‡]	0.09 [‡]	0.12 [‡]	0.16 [‡]	0.13 [‡]	-	-
Insight	0.09 [‡]	0.07 [‡]	0.1 [‡]	0.08 [‡]	-	0.1 [*]	0.17 [‡]	-	-
Discrepancy	-	0.06 [*]	-	-	-	-	-	-	-
Tentative	0.07 [‡]	0.08 [‡]	0.08 [‡]	0.07 [*]	-	-	-	-	-
Differentiation	0.08 [‡]	0.08 [‡]	0.1 [‡]	-	-	0.09 [*]	-	-	-
Biological processes	0.06 [‡]	-	-	-	-	-	-	-	-
Health	0.08 [‡]	0.07 [‡]	0.08 [‡]	0.11 [‡]	-	-	-	-	-
<i>Time orientation</i>									
Past focus	0.08 [‡]	0.05 [*]	0.09 [‡]	0.08 [‡]	0.09 [*]	-	0.11 [*]	-	-
Present focus	0.09 [‡]	0.06 [‡]	0.1 [‡]	0.07 [*]	-	-	-	-	-
<i>Personal concerns</i>									
Relativity	0.05 [‡]	0.06 [‡]	-	-	-	-	-	-	-
Time	0.06 [‡]	-	0.06 [*]	-	-	-	-	-	-
Work	-	0.06 [*]	-	-	-	-	-	-	-
Leisure	-0.07 [‡]	-0.07 [‡]	-0.07 [‡]	-0.09 [‡]	-0.12 [‡]	-0.09 [‡]	-	-	-
Money	-0.06 [‡]	-	-0.05 [‡]	-0.06 [‡]	-	-0.1 [‡]	-	-	-0.12 [‡]
Informal language	-0.07 [‡]	-0.07 [‡]	-0.06 [‡]	-	-0.12 [‡]	-	-0.09 [*]	-	-
Netspeak	-0.07 [‡]	-0.06 [‡]	-0.06 [‡]	-0.06 [*]	-0.1 [‡]	-0.08 [*]	-	-	-0.15 [‡]
Assent	-	-	-	-	-0.06 [*]	-	-	-	-

Control users more likely to exhibit “clout” than users with depression or anxiety.

“Authentic” language is more common among most conditions (e.g., using more personal pronouns).

Topics such as leisure and money are more common among control users.

Can we predict the condition based on posts?

FastText model
generally performs
better than others (F1)

Method	Logistic				
	Regression	XGBoost	SVM	FastText	CNN
Depression	42.9	43.3	52.8	53.6	50.5
ADHD	31.6	34.1	44.9	47.0	35.1
Anxiety	41.1	45.7	53.6	53.7	48.0
Bipolar	34.6	44.7	51.6	50.8	38.3
PTSD	32.5	53.4	57.1	57.6	54.9
Autism	14.6	34.2	42.9	49.8	38.3
OCD	9.3	39.6	41.6	44.5	32.9
Schizophrenia	4.4	38.5	39.5	45.3	37.0
Eating	0.0	24.8	23.6	41.2	44.5

SMHD: A Large-Scale Resource for Exploring Online Language Usage for Multiple Mental Health Conditions

- Large dataset of mental health language, spanning multiple conditions
<http://ir.cs.georgetown.edu/resources/>
- Observations of significant differences between diagnosed users and control users on a variety of categories
- Baseline classification results show prediction is a difficult task