



CS 8803 Data Analytics for Well-Being: Prediction I

Munmun De Choudhury

munmund@gatech.edu

Week 9 | March 7, 2016

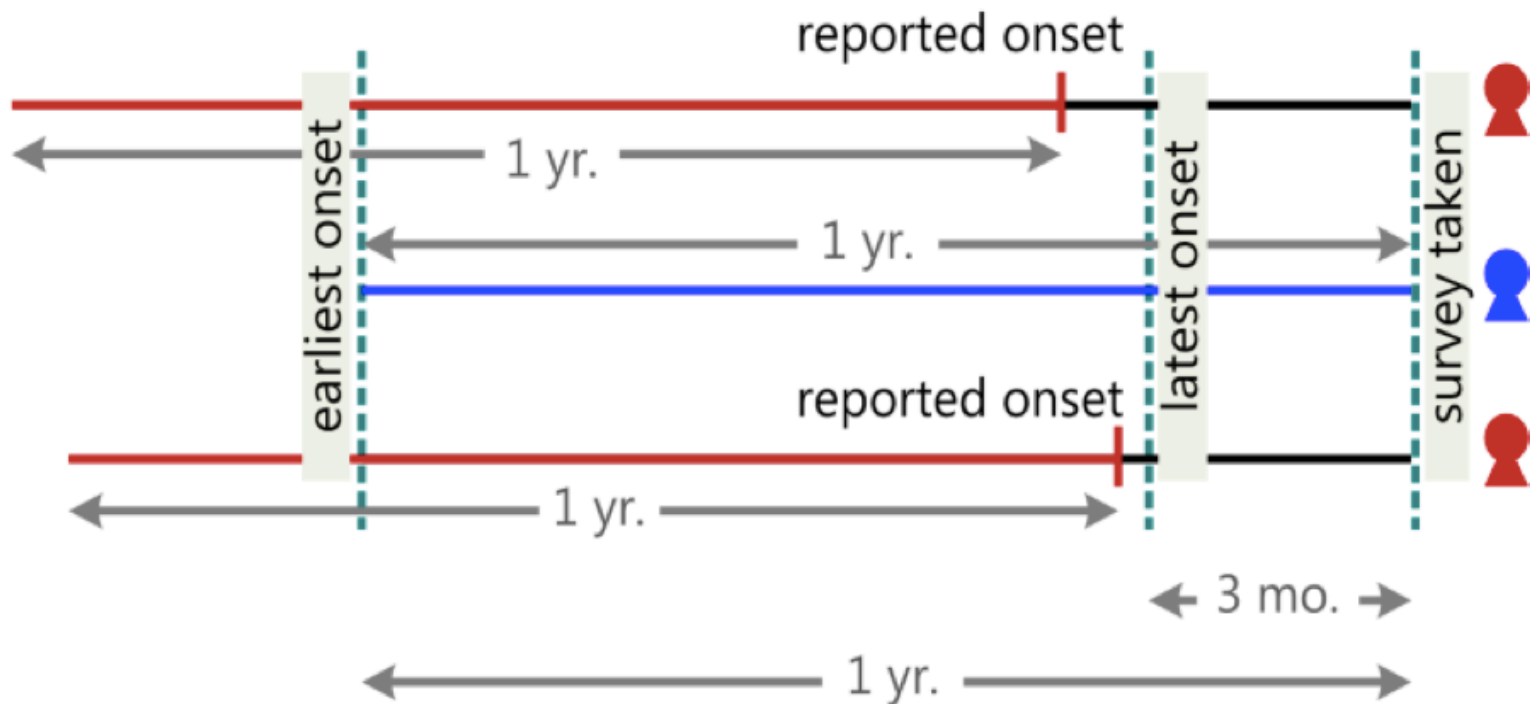
Please sign up for the class
presentations!!!!

Mid-term progress presentations due
on Mar 16

Predicting Depression via Social Media

Summary

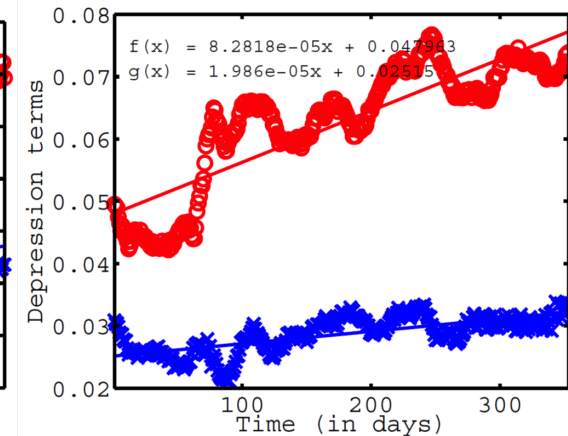
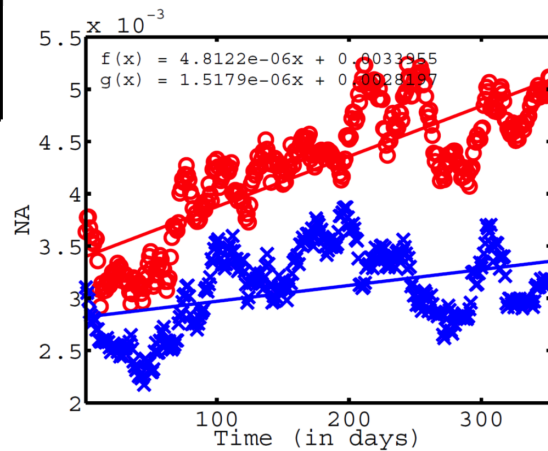
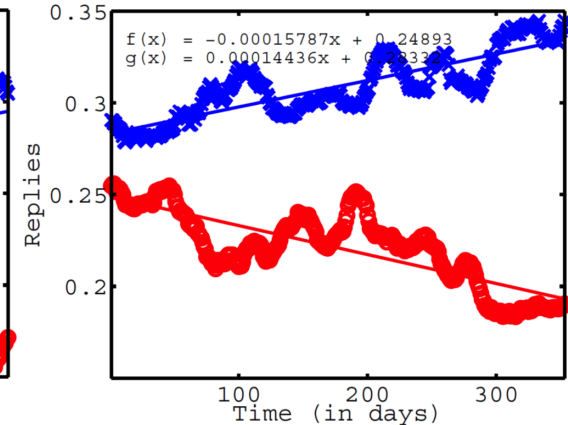
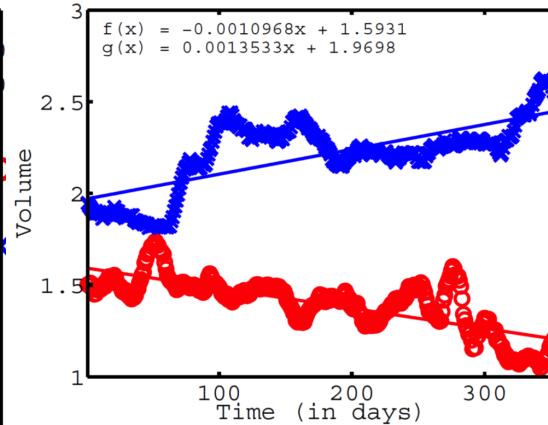
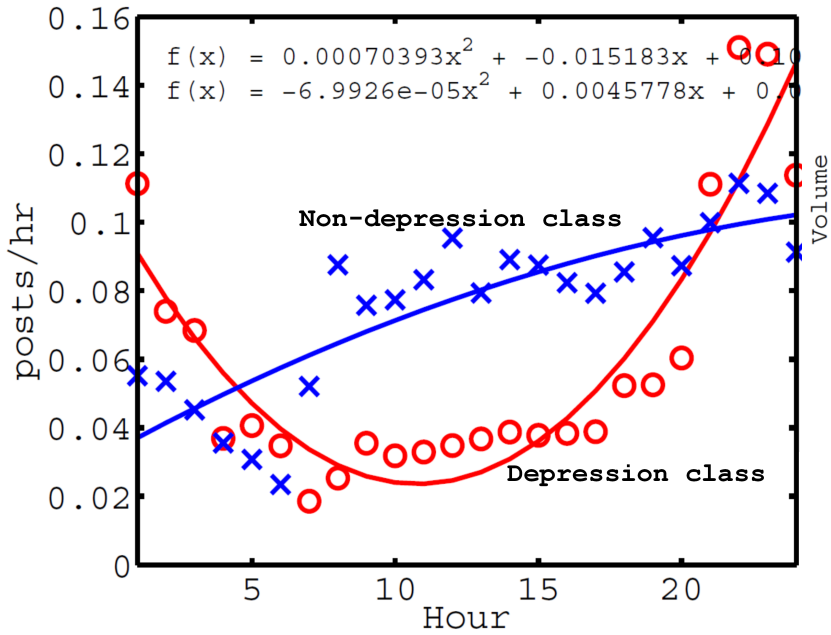
- Can social media activities and connectedness predict risk to major depressive disorder?
- Recruitment of a sample of Twitter users through a survey methodology over Amazon's Mechanical Turk
 - ~40% provided access to Twitter data



Summary

- Social engagement
- “Insomnia index” – mean z-score of an individual’s volume of Twitter activity per hour
- Ego-centric social graph – nodal properties (*inlinks, outlinks*); dyadic properties (*reciprocity, interpersonal exchange*); neighborhood properties (*density, clustering coefficient, two-hop neighborhood, embeddedness, number of ego components*)
- Language
 - Depression lexicon – top uni- and bigrams compiled from Yahoo! Answers category on mental health
 - Linguistic style

Summary



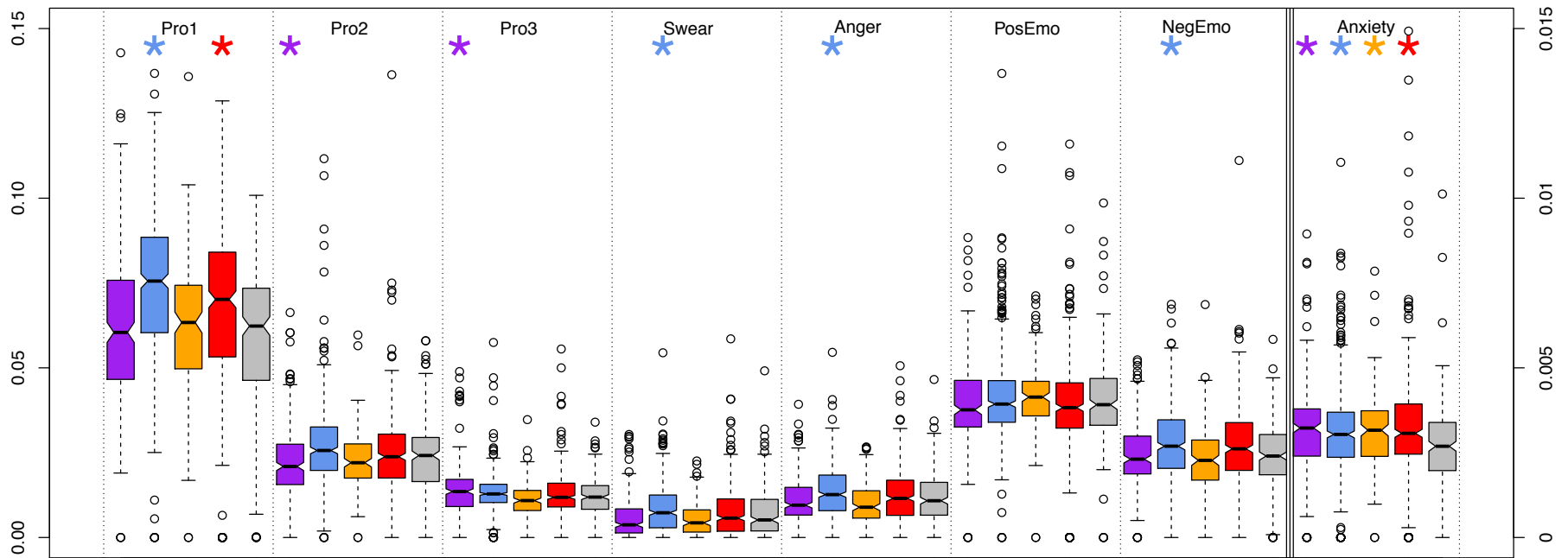
Summary

Egonetwork measures	Depres. class	Non-depres. class
#followers/inlinks	26.9 ($\sigma=78.3$)	45.32 ($\sigma=90.74$)
#followees/outlinks	19.2 ($\sigma=52.4$)	40.06 ($\sigma=63.25$)
Reciprocity	0.77 ($\sigma=0.09$)	1.364 ($\sigma=0.186$)
Prestige ratio	0.98 ($\sigma=0.13$)	0.613 ($\sigma=0.277$)
Graph density	0.01 ($\sigma=0.03$)	0.019 ($\sigma=0.051$)
Clustering coefficient	0.02 ($\sigma=0.05$)	0.011 ($\sigma=0.072$)
2-hop neighborhood	104 ($\sigma=82.42$)	198.4 ($\sigma=110.3$)
Embeddedness	0.38 ($\sigma=0.14$)	0.226 ($\sigma=0.192$)
#ego components	15.3 ($\sigma=3.25$)	7.851 ($\sigma=6.294$)

Quantifying Mental Health Signals on Twitter

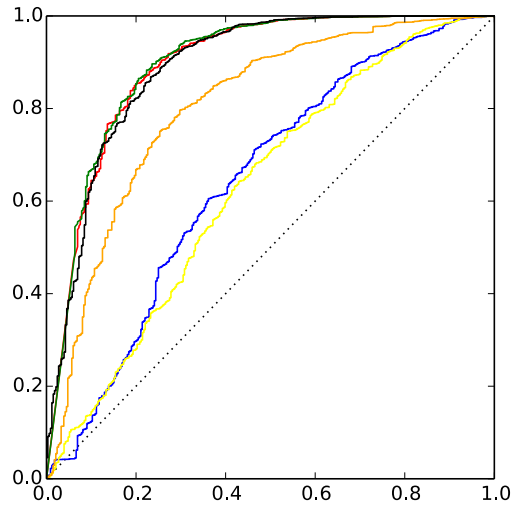
Summary

- Overcome challenges of solicited survey data – use self-identified mental health diagnoses on Twitter for obtaining ground truth.

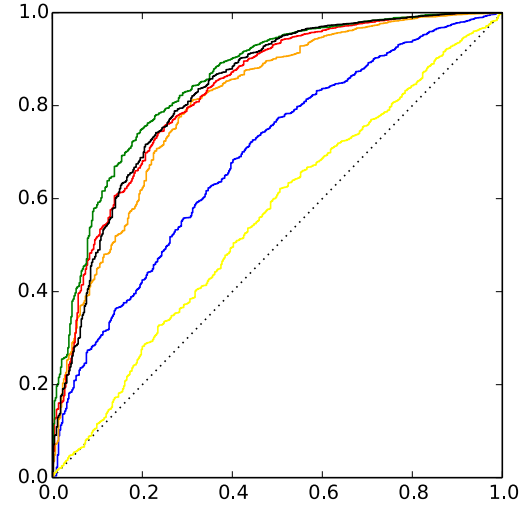


Summary

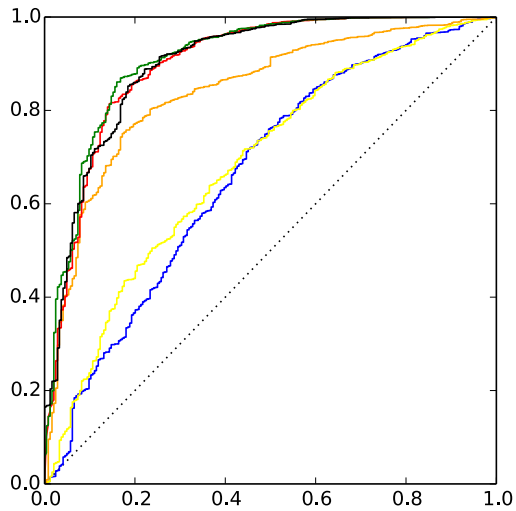
Bipolar



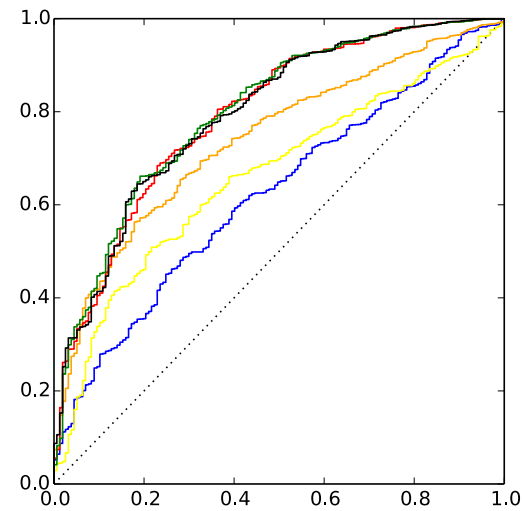
Depression



PTSD



SAD



What are the differences you observed between the two studies? What are the strengths and limitations of those differences?

A consistent challenge in many prediction tasks like these, is gathering gold standard information (or ground truth). What could be different ways to get at this problem?

In many ways, predictive models are never 100% perfect. How can social media platforms leverage the predictive methodologies outlined in the two papers?

Both papers used Twitter as the data source of study – in ways it is a nice platform where an individual can make their profile however they wish it to be. Would the same findings hold on a platform that enforces real identities, like Facebook?

Depression is not an online condition, but one that spans both the online and the offline life. The papers do not take offline attributes into their models.

Is there a way to that into account? What would be the most significant offline attributes to consider?

Anurag brings up the concern of “awareness contamination”. To what extent do you think this would impact a predictive model like the one proposed in the papers going forward? How would you combat that?

The analyses in Coppersmith et al. indicate that depression, PTSD, and bipolar disorder have high mutual correlation, but correlation is low between SAD and the others. Why do you think that is the case – is there anything in the model or the data that account for it?

In the first paper, the ground truth was obtained from Amazon mechanical turk workers. Anything odd or specific about this population that may have affected the findings? What would be alternative ways of recruiting people?

Coppersmith et al shows how n -gram language models can provide better performance over LIWC features. What is the downside of using language models though?