

Research Statement: H. Andrew Schwartz

Postdoctoral Research Fellow, University of Pennsylvania

www.seas.upenn.edu/~hansens/
hansens@seas.upenn.edu
215-746-5085

What can language analyses reveal about human health and well-being?

Upon proudly introducing her Ph.D. relative, a friend once added, “yes, a doctor, but not the kind that helps people.”

With the growth of social media, now used regularly by more than one-seventh of the world’s population, scientists are presented with an unprecedented resource of objective and quantifiable behavioral and textual data. People freely post their daily activities, feelings, and thoughts in the form of status updates, tweets, and other personal online discourse. This “big data”, unlike that from questionnaires or controlled studies, is unprompted and requires no *a priori* theories to collect, and thus opens the door for enormous data-driven discovery in the health and social sciences.

Researchers have begun leveraging natural language data for applications such as monitoring influenza outbreaks, predicting the stock market, and personally targeting advertisements. While such pursuits are exciting and worthwhile, I believe the biggest application of such natural language analyses are toward science itself: discovery of new behavioral and psychological drivers of health and well-being.

Recently, along with colleagues from psychology, medicine, and computer science, I’ve begun exploring how scalable natural language processing algorithms can be developed to mine social media for scientific discovery about people. Language use features, in the form of *words*, *phrases*, and *topics*,¹ are extracted from Facebook status updates or Tweets, sometimes dealing with billions of messages using distributed computing.² Once language use features are extracted, we perform correlational analyses and machine learning for (a) prediction of outcomes (e.g. personality, heart disease, life satisfaction) from social media messages and (b) insights from language about the behavioral and psychological factors of well-being and health.

Prediction. My research shows that *words*, *phrases*, and *topics* in Tweets and Facebook posts are predictive of many health and psychological outcomes. Language from Twitter significantly improved accuracy in a predictive model of county life satisfaction, as measured by the CDC, over and above standard demographic and socio-economic controls (age, gender, ethnicity, income, and education; see fig. 1) [1]. Similarly, atherosclerotic coronary heart disease mortality was predicted significantly higher by a twitter language model than by models based on SES, smoking, hypertension, and a combined model with ten socioeconomic demographic and health variables (see fig. 2) [2]. On the individual level, we found personality, gender, age, and life satisfaction were strongly linked with language in status updates from Facebook, and that an *open-vocabulary* analysis³ resulted in significantly higher predictive accuracies than those from traditional closed-vocabulary approaches [3, 4, 5]

Insight. Results from our predictive models of health, well-being, and personality speak to the potential of language to predict many health and psychologically-relevant outcomes. However, language can be used for more than prediction; A parallel goal is to gain a greater understanding of the nature of people. The problem changes from obtaining a single high prediction accuracy to that of describing the many links between language and outcomes. The *topics* used to predict county *life satisfaction* reveal insight into what makes a community thrive beyond broad socio-economic and demographic variables. For example, topics relating to outdoor activities, spiritual meaning, philanthropy, exercise, and work correlate with increased life satisfaction (see fig. 3) [1]. In the personality study over Facebook, topics related to laughter, college, getting a job, and having children were characteristic of the different age groups (see fig. 4). We found

¹**phrases:** words more likely to occur together than by chance (sometimes called *collocations*; e.g. “sick of”, “a beautiful day”).

topics: groups of semantically-related words automatically derived via Bayesian *Latent Dirichlet Allocation* (e.g. see fig. 3).

² Utilizing a *Hadoop*-style cluster to distribute feature extraction over many machines.

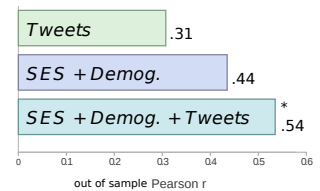


fig 1: Accuracy of predictive models on county life-satisfaction. *Tweets:* topic usage from tweets. *SES + Demog.:* age, gender, ethnicity, income, and education. *SES + Demog. + Tweets:* SES and topics from tweets. * = sig. ($p < .01$) improvement over *SES + Demog.* [1]

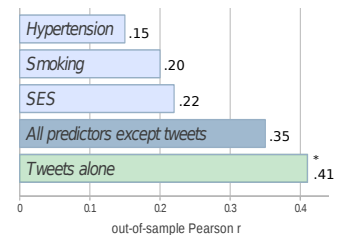


fig 2: Accuracy of predictive models on county heart disease mortality. *Hypertension, smoking:* county rates; *SES:* income and education. *All predictors except tweets:* rates of diabetes, obesity, smoking, hypertension; % black, Hispanic, female, married; income + education. *Tweets alone:* words, phrases, topics. * = sig. ($p < .02$) improvement over all predictors. [2]

³ **open-vocabulary:** data-driven approaches which discover predictive words rather than relying on *a priori* word or category judgments [3].

of causal *discourse relations* to distinguish events from their causal explanations, and sentiment analysis to find the valence of the event.

- Detecting temporal-orientation and conditional statements.** People’s ”temporal-orientation” (focus on the past, present, or future) and ”prospection” (thoughts and forecasts about the future, such as if-then thinking) have been linked with health and achievement. I am starting to use temporal orientation as revealed by social media to predict many health and well-being outcomes. Beyond the tense of verbs, people use many temporal indicators (e.g. ’tomorrow’, ’3 weeks ago’, ’next year’) that can be leveraged for identifying higher resolution temporal orientation (e.g. *4 days in the future* rather than simply *future*). Further, I seek to automatically identify conditional statements with respect to evaluation of future events (e.g. ”if I go out drinking tonight, I may not make it to my son’s soccer game tomorrow morning.”). This research will combine syntactic parsing to capture conditionals with semantic analysis for temporal orientation.

Scientific Applications. Serious behavioral, psychological, and health questions motivate my creation of novel language analyses. Over time, I hope to have solid answers to the questions outlined below.

- Can we predict individual disease recovery or longevity based on language use?** I am collaborating with researchers in medicine to write and distribute a cell-phone ”app” to consenting hospital or clinical trial patients in order to match language with individual health outcomes. I am also examining childhood writings from a cohort of 10,000 citizens of the United Kingdom with biomedical outcomes tracked across their lifetime so see which behaviors, as revealed through language, are most predictive of health outcomes.
- Can language-based prediction of personality partially replace questionnaire-based assessments?** Questionnaires have been primary tool of psychology for over a century, but they are time-consuming to administer and they are subject to self-report biases. Along with a quantitative psychologist, I am evaluating our language-based models for psychometric validity (e.g. construct, internal, and external validity and reliability).
- Which diseases are most easily predicted at the community level from social media language?** Expanding on the heart-disease work, I am now comparing which community health-outcomes are best predicted by Twitter language, and what their characteristics are. Controlling for socio-economic factors, this may give insight into those diseases that are most influenced by the psychology or behavior of a community.
- How does the temporal-orientation of individuals and communities influence their health, economics, and well-being?** Early results suggest, for example, use of more future-oriented words is predictive of lower spread of HIV among counties with larger amounts of risky behavior.

In recent decades, the biomedical sciences have been transformed based on large scale computational analyses such as the microarray and gene sequencers. Rather than testing specific hypotheses about genes, the process changed to that of *data-driven discovery*. Exploratory computational linguistics and information science over the new massive datasets of personal discourse (social media) are bringing behavioral and psychological factors into the reach of data science. Given the numerous findings revealed by language thus far, it is perhaps not far-fetched that the next major break-through in heart disease or cancer treatment will derive from the data-driven insights of computational linguistics – novel scientific discoveries of behavioral and psychological factors that ”help people”.

extraversion



introversion

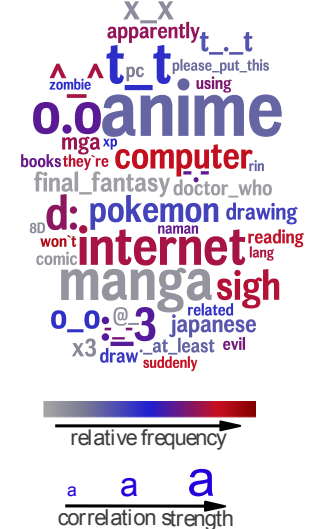


fig 5: Words and phrases most characteristic of extraversion (top) and introversion (bottom). (significantly correlated at $p < 0.001$). Word size is scaled by correlation strength. Color indicates relative frequency.) [3].

- [1] H. A. Schwartz, J. C. Eichstaedt, M. L. Kern, L. Dziurzynski, R. E. Lucas, M. Agrawal, G. J. Park, S. K. Lakshmikanth, S. Jha, M. E. P. Seligman, and L. H. Ungar, "Characterizing geographic variation in well-being using tweets," in *Proceedings of the Seventh International AAAI Conference on Weblogs and Social Media (ICWSM)*, 2013.
- [2] J. C. Eichstaedt, H. A. Schwartz, M. L. Kern, G. J. Park, S. Jha, M. Agrawal, L. Dziurzynski, M. Sap, C. Weeg, R. Merchant, D. Labarthe, M. E. P. Seligman, and L. H. Ungar, "Twitter correlates with atherosclerotic heart disease mortality across us counties.," *In Preparation*, 2013.
- [3] H. A. Schwartz, J. C. Eichstaedt, M. L. Kern, L. Dziurzynski, S. M. Ramones, M. Agrawal, A. Shah, M. Kosinski, D. Stillwell, M. E. P. Seligman, and L. H. Ungar, "Personality, gender, and age in the language of social media: The open-vocabulary approach," *PLOS ONE* 8(9), 2013.
- [4] H. A. Schwartz, J. C. Eichstaedt, L. Dziurzynski, M. L. Kern, E. Blanco, M. Kosinski, D. Stillwell, M. E. P. Seligman, and L. H. Ungar, "Toward personality insights from language exploration in social media," in *Proceedings of the AAAI Spring Symposium Series*, 2013.
- [5] H. A. Schwartz, J. C. Eichstaedt, M. L. Kern, L. Dziurzynski, M. Agrawal, A. Kapelner, M. E. P. Seligman, and L. H. Ungar, "Predicting well-being through language use: Multi-level modeling messages and users," *In Preparation*, 2013.
- [6] M. L. Kern, J. C. Eichstaedt, H. A. Schwartz, G. Park, L. H. Ungar, D. J. Stillwell, M. Kosinski, L. Dziurzynski, and M. E. P. Seligman, "From "sooo excited!!!" to "so proud": Using language to study development," *Developmental Psychology*, 2013.
- [7] H. A. Schwartz, J. C. Eichstaedt, L. Dziurzynski, M. L. Kern, E. Blanco, S. Ramones, M. E. P. Seligman, and L. H. Ungar, "Choosing the right words: Characterizing and reducing error of the word count approach," in *Proceedings of *SEM-2013: Second Joint Conference on Lexical and Computational Semantics*, 2013.
- [8] H. A. Schwartz and F. Gomez, "Acquiring knowledge from the web to be used as selectors for noun sense disambiguation," in *Proceedings of the Twelfth Conference on Computational Natural Language Learning*, (Manchester, England), pp. 105–112, August 2008.
- [9] H. A. Schwartz, F. Gomez, and L. Ungar, "Improving supervised sense disambiguation with web-scale selectors," in *COLing-2012: the 24th International Conference on Computational Linguistics*, (Mumbai, India), pp. 105–112, August 2012.

(This is a partial list of most relevant publications. The complete list is available in the curriculum vitae.)