

Free and Open Source Software for Geospatial (FOSS4G) Conference Proceedings

Volume 17 *Boston, USA*

Article 21

2017

GIS Investigation of Crime Prediction with an Operationalized Tweet Corpus

Anthony J. Corso

California Baptist University

Abdulkareem Alsudais

Claremont Graduate University

Follow this and additional works at: <https://scholarworks.umass.edu/foss4g>

 Part of the [Business Commons](#), [Engineering Commons](#), and the [Social and Behavioral Sciences Commons](#)

Recommended Citation

Corso, Anthony J. and Alsudais, Abdulkareem (2017) "GIS Investigation of Crime Prediction with an Operationalized Tweet Corpus," *Free and Open Source Software for Geospatial (FOSS4G) Conference Proceedings*: Vol. 17 , Article 21.

DOI: <https://doi.org/10.7275/R5WM1BKZ>

Available at: <https://scholarworks.umass.edu/foss4g/vol17/iss1/21>

This Poster is brought to you for free and open access by ScholarWorks@UMass Amherst. It has been accepted for inclusion in Free and Open Source Software for Geospatial (FOSS4G) Conference Proceedings by an authorized editor of ScholarWorks@UMass Amherst. For more information, please contact scholarworks@library.umass.edu.

GIS Investigation of Crime Prediction with an Operationalized Tweet Corpus

Anthony J. Corso^{a,*}, Abdulkareem Alsudais^b

^a*California Baptist University*

^b*Claremont Graduate University*

Abstract: Social media as the de facto communication channel is being used to disseminate one's diurnal self-revelations. This profound discovery often contains double-talk, peculiar insights, or contextual information about real-world events. Natural language processing is regularly used to uncover both obvious and latent knowledge claims within disclosures published amid the complex environment. For example, a perpetrator with first-hand knowledge of their criminal incident uses social media to post critical information about it. A geographic information system (GIS) is capable of large-scale point data analysis and possesses methods that enable dataset processing, evaluation, and automatic spatial visualization. Such an artifact fused with traditional environmental criminology theory and social media erects guidelines, tools, and models for substantive construction and evaluation of GIS crime analysis solutions. Provided the social media stream is timely and correctly processed, corrective action can be taken. The construction of a natural language processing social media annotation pipe identifies latent indicators extracted from a social media corpus and is an integral part of societal mishap prediction. Spatial visualizations and regression analyses were used to describe and evaluate project artifacts. As a result, a social media corpus was operationalized, and subsequently used as a proxy for a traditional environmental criminology risk layer in construction of a social media GIS crime analysis artifact. Using such multi-domain collaboration, the artifact was able to increase the predictive crime incident outcome with an overall R-squared increase of 21.94%. This result is the state-of-the-art; there are no other results to compare it to.

Poster Download: <http://scholarworks.umass.edu/foss4g/vol17/iss1/21>

*Corresponding author
Email address: acorso@calbaptist.edu (Anthony J. Corso)

GIS Investigation of Crime Prediction with an Operationalized Tweet Corpus

Anthony J. Corso, Ph.D.

California Baptist University



Introduction

Social media (e.g., tweets) are the de facto communication channel to disseminate one's diurnal self-revelations. This profound phenomenon contains double-talk, peculiar insight, and contextual data or information about real-world events. Amid such complex and personal expose, natural language processing (NLP) techniques uncover both obvious and latent knowledge claims published within.

A geographical information system is capable of large-scale data analysis and possesses methods that enable dataset processing, evaluation, and spatial visualization. When fused with traditional research theory—such an artifact defines guidelines, algorithms, and models for substantive and predictive investigation.

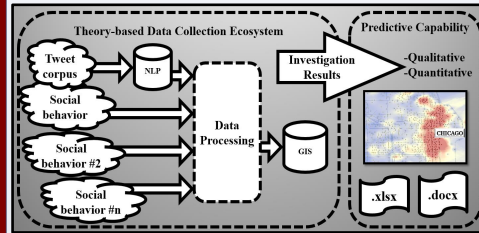
Objectives

Despite a tweet's sparse content, NLP makes their use in a predictive GIS artifact feasible. For example, subsequent to processing, useful tweets are able to:

- Predict the validity of a real-world event only recorded by observation of social media eyewitness; or
- Predict real-time trends by amalgamating social media with traditional social behavior variables.

Thus, inquiry explores GIS outcomes when consuming "useful" or "not useful" tweets as identified via NLP techniques. In addition, a research framework illustrates social media being coalesced with other behavior variables and subsequently used as a social behavior GIS proxy layer.

Research Framework



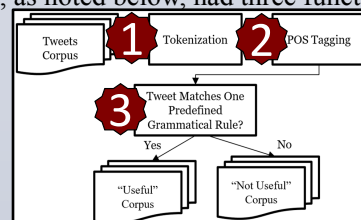
The research framework defines a process and exposes opportunity to fill the gap between sparse text social media and its representation of real-world events by examining meaningful tweet content and purging useless structures. That is, some tweets are so sparse they cannot represent the real-world context in which they exist; hence, a "Not Useful" tweet (illustrated in the table below). However, some tweets

Useful	Not Useful
I'm at Old Navy in Chicago IL https://t.co/lczpu9NLF	balloooooonsss ??? http://t.co/mjhuKyH7DM
My Phone Die So Fast David Bowie is... my favorite! @ David wie is At Mca Chicago	WAYYOHANDSIDETOOSI Funniest
I aint ga stunt on nobody...trust me Yo LoL! I can't wait to see lora tomorrow	#CuppyCoffee!!!!!! @_lorShane????

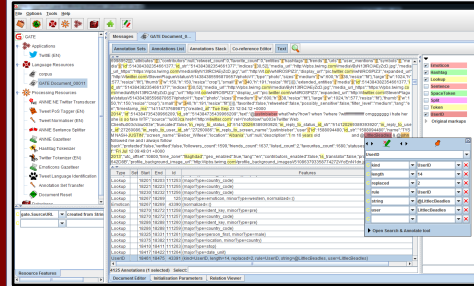
are "Useful" but require extra processing.

Tweets & Natural Language Processing

Operationalizing "useful" or "not useful" tweets was accomplished via the General Architecture for Text Engineering (GATE) NLP suite of tools. The NLP pipeline built, as noted below, had three functions.

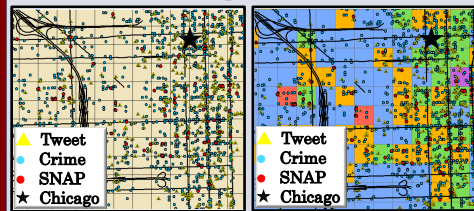


Therefore, it tokenized, part-of-speech tagged, and applied custom grammar rules to each tweet. A custom GATE NLP application (below) executed the pipeline.



GIS Analysis

Association between features of a tweet, e.g., acronym use and its grammatical structure, and its potential usefulness were operationalized via NLP preprocessing. GIS capability examined both quantifiable and meaningful qualitative results; each are required in data analysis, information dissemination, and predictive artifacts.



Data Map

Solution Map

The maps represent the area of downtown Chicago with a fishnet spacing of 750 feet. Tweets,^[2] crime,^[3] and SNAP^[4] locations are the variables displayed. The Data Map is a visualization of the data. The Solution Map represents the results of a GIS grouping analysis tool used for exploratory variable analysis; the attributes are combined and cell shading represents latent structures.

Discussion and Conclusion

With a novel NLP pipeline tweets were processed and used to measure the change in performance of an ArcGIS^[5] 10.4.1 artifact. A 1,000 tweet sample was hand tagged and compared to a baseline model, and to an innovative social media grammar applied by a rule-based social media NLP pipeline. GIS evaluation tools answer the question, prior to content analysis of a tweet, does a method exist to support identifying a tweet as "useful" for subsequent GIS processing? Indeed, "useful" tweet identification via NLP returned precision of 0.9256, recall of 0.6590, and F-measure of 0.7699; consequently, exploratory GIS processing of a social media variable increased 0.2194 over baseline.

Predictive capability potential of a GIS artifact implementing social media's latent behavior attributes is vast. Yes, preliminary results are encouraging but future research is important and needs to identify its value.

References

Alonso, O., Marshall, C. C., & Najjar, M. (2015). Are some tweets more interesting than others? A visualization. Paper presented at the Proceedings of the Symposium on Human Computer Interaction, 1-10.

Anstee, P., Bernstein, M., & Lathin, R. (2012). Who gives a tweet?: evaluating microblog content value. Paper presented at the Proceedings of the ACM 2012 conference on Computer Supported Cooperative Work.

Banerjee, J., Bhatia, A., & Veeramani, D. (2014). Crime Mapping through Geo-Spatial Social Media Activity.

Blumenthal, K., Brody, L., Fink, A., Greenwood, M. A., Maynard, D., & Awan, N. (2013). TweetIE: An Open-Source Information Extraction Pipeline for Microblog Text. Paper presented at the SIGSPATIAL.

Bramson, P., Escobar-Melians, M., Patel, A., & Shivas, R. (2011). Extracting social power relationships from natural language. Paper presented at the Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies.

Caslin, J. M., Kennedy, L. W., & Miller, J. (2011). Risk Terrain Modeling: Bridging Criminological Theory and GIS Methods for Crime Forecasting. *Justice Quarterly*, 28(2), 168-181. doi:10.1080/07418818.2010.518057

City of Chicago. (2014). Retrieved December 31, 2014, from <https://data.cityofchicago.org/>

Corso, A., & Abulafia, A. (2016). Big Social Data and GIS: Visualizing Predictive Crime. *AMIS '16 Conference 2016 Proceedings*.

Corso, A. J., & Abulafia, A. (2015). GIS, Big Data, and a Tweet Corpus Operationalized on Natural Language Processing. *AMIS '15 Conference 2015 Proceedings*.

Deerwaele, C. (2014). A Metric Comparison of Predictive Hot Spot Techniques and RTM. *Justice Quarterly*, 17(2), doi:10.1007/s10992-014-9449-0

ESRI. (2011). ArcGIS Desktop: Release 10. Redlands, CA: Environmental Systems Research Institute.

Frantz, K., Anastakis, S., & Mims, H. (2009). Automatic recognition of multi-word terms: the "naïve" N-gram method. *International Journal on Digital Libraries*, 1(2), 115-128.

Gardner, M. S. (2014). Predicting crime using Twitter and kernel density estimation. *Business Support Systems*, 4(1), 115-124.

Grice, G., & Ferguson, C. (2007). Biagrams of syntactic labels for authorship discrimination of short texts. *Literary and Linguistic Computing*, 22(4), 485-497.

Harris, N., Yamamoto, T., Jerome, M., & Tashiro, H. (2012). Effective Classification and Visualization of Photo-tagged Geotagged Tweets.

Harfield, J., & Wilson, M. L. (2011). Searching Twitter: Separating the Tweet from the Chaff. Paper presented at the ICWSM.

Jaworski, D., & Martin, J. E. (2009). Spreads and Language Processing.

Karim, L. M., Caplan, J. M., & Papp, C. (2014). Risk Factors, Strategies, and Spatial Intelligence: Risk Terrain Modeling as an Algorithm for Police Resource Allocation.

Leary, C. (2011). *Business User Statistics in Information Systems Science & Business Media*.

Pian, S., Hui, X., Li, M., & Harter, S. (2008). Learning to classify short and sparse text with hidden topics from large-scale data collections. Paper presented at the Proceedings of the 2008 International Conference on Weblogs and Social Media.

Pian, S., & White, J. (2011). A Feasibility Study on Extracting Twitter Users' Interest Using NLP Tools for Nonacademic Connections. Paper presented at the Proceedings of the 2011 International Conference on Social Networks and Social Media.

Leary, C. (2011). *Business User Statistics in Information Systems Science & Business Media*.

Quatman, C., Pomeroy, H., Ouerfelli, D., Copra, L., & Maassini, M. (2014). Who Benefits from the "Sharing" Economy of Airbnb? Paper presented at the Proceedings of the 2014 International Conference on Weblogs and Social Media.

Seaman, R., Fisher, D., Hunter, E., Forthmann, H., & Dunfee, M. (2008). Short text classification in Twitter to improve information filtering. Paper presented at the Proceedings of the 2008 ACM SIGMIS Database Conference on Data Mining and Knowledge Discovery.

Schuchman, A., Crank, A., & Rothblat, J. (2013). Retrieving and analyzing geotagged information from social media feeds. *Geoinformatics*, 7(2), 119-128. doi:10.1007/s10706-011-9424-2

Tierce-Morvan, J.-M. (2014). Three Statistical Summaries of CLEF EKN 2013 Tweet Classification Task. Paper presented at the Proceedings of the 2014 International Conference on Text Mining and Applications.

U.S. Patent. (2014). Retrieved December 31, 2014, from <http://www.uspto.gov/patents/office/>

United States Department of Agriculture. (2014). Retrieved December 31, 2014, from <http://www.usda.gov/aginfo/geoportal/>

Zhang, M., Chen, L., & Boreck, C. (2015). Statistical and Semantic Approaches for Tweet Classification. *Proceedings Computer Systems*, 61, 494-507.

Zhang, A., & Ji, H. (2014). Tweet, but verify: empirical study of information verification on Twitter. *Journal of Intelligent and Social Network Analysis and Mining*, 4(1), 1-12.

Contact

Dr. Corso holds Ph.D. in Information Systems and Technology from Claremont Graduate University. Since 2007, he is an Associate Professor in the Gordon and Jill Bourns College of Engineering at California Baptist University. E-mail: acorso@calbaptist.edu