

Interactive comment on “TAGGS: Grouping Tweets to Improve Global Geotagging for Disaster Response” by Jens de Bruijn et al.

Anonymous Referee #2

Received and published: 20 June 2017

Geoparsing and Geotagging is a well studied problem. The authors have missed a lot of important work in related work section (see end citations).

Although the problem is based on disaster response actually the work presented is standard geoparsing.

The location identification approach proposed by the authors is not state of the art (standard tokenization, geonames gazetteer lookup).

The location disambiguation is also quite simplistic (language/time zone/population filters, location name subsumption). Modern geoparsing approaches use extra information to boost precision, either from contextual text (e.g. linguistic patterns used and other location mentions in text), geospatial source (e.g. OpenStreetMap database lo-

C1

cations for spatial proximity) or social context (e.g. social media tags from a training corpus). The approach used by the authors is not advancing the state of the art.

The idea of clustering (i.e. grouping) posts with the same location mention to provide additional context is potentially a novel idea that does advance the state of the art.

However the choice to filter out all locations where the population density is small greatly weakens this works applied value. Geoparsing popular locations, usually cities with a high population density, is much easier than geoparsing obscure remote locations. For populous areas social media tagging has been proven to be very effective. Such work (see Kordopatis-Zilos) has not been cited at all and is not compared.

There are publicly available benchmark geoparse datasets (see PyPI geoparsepy <https://pypi.python.org/pypi/geoparsepy>) which could and should have been used to directly compare results to published work (see Middleton and Gelernter). The authors instead created a smaller dataset manually and thus the results are not directly comparable.

The evaluation itself is weak. No statistical significance is reported. A breakdown of results by event & language is not provided (known factors in performance). The authors appear to be geoparsing at the region level (the paper is not precise about this so its impossible to tell). The example of 'A10 near Amsterdam' should geoparse both the A10 (road) and Amsterdam (region) - but appears to only score on Amsterdam. The whole approach lacks rigour and depth.

The results themselves look weak (i.e precision scores), which is not surprising considering the lack of state of the art location disambiguation techniques. I would expect a precision for region level geoparsing to be 0.9+ (F1 0.8+). The authors manage around 0.7 precision.

Overall this work is too immature at this stage for a journal publication. The lack of a real comparison to the extensive related work already published in this area means

C2

this work has little current value to a serious researcher.

– missing citations –

Judith Gelernter and Nikolai Mushegian, 2011. Geo-parsing Messages from Microtext. Transactions in GIS, Vol 15, Issue 6, 753–773

Middleton, S.E. Middleton, L. Modafferi, S. "Real-time Crisis Mapping of Natural Disasters using Social Media", Intelligent Systems, IEEE , vol.29, no.2, pp.9,17, Mar.-Apr. 2014

Middleton, S.E. Krivcovs, V. "Geoparsing and Geosemantics for Social Media: Spatio-Temporal Grounding of Content Propagating Rumours to support Trust and Veracity Analysis during Breaking News", ACM Transactions on Information Systems (TOIS), 34, 3, Article 16 (April 2016), 26 pages.

Marieke van Erp, Giuseppe Rizzo and Raphaël Troncy, 2013. Learning with the Web: Spotting Named Entities on the intersection of NERD and Machine Learning. In Proceedings of the 3rd Workshop on Making Sense of Microposts (#MSM2013). Rio de Janeiro, Brazil

Alan Ritter, Sam Clark, Mausam and Oren Etzioni, 2011. Named entity recognition in tweets: An experimental study. In Proceedings of Empirical Methods for Natural Language Processing (EMNLP), Edinburgh, UK

Alexandre Davis, Adriano Veloso, Altigran S. da Silva, Wagner Meira Jr., Alberto H. F. Laender, 2012. Named entity disambiguation in streaming data. In Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics, 815 - 824, Jeju, Republic of Korea

Giorgos Kordopatis-Zilos, Symeon Papadopoulos, and Yiannis Kompatsiaris. 2015. Geotagging social media content with a refined language modelling approach. In Pacific-Asia Workshop on Intelligence and Security Informatics. Ho Chi Minh City, Vietnam. 19 May 2015

C3

Interactive comment on Nat. Hazards Earth Syst. Sci. Discuss., <https://doi.org/10.5194/nhess-2017-203>, 2017.

C4