

A crowdsourcing approach to advance collective awareness and social good practices

Ilias Dimitriadis

School of Informatics

Aristotle University, Thessaloniki,
Greece, idimitriad@csd.auth.gr

Vasileios G. Psomiadis

School of Informatics

Aristotle University, Thessaloniki,
Greece, vpsomiadis@csd.auth.gr

Athena Vakali

School of Informatics

Aristotle University, Thessaloniki,
Greece, avakali@csd.auth.gr

ABSTRACT

Contemporary societies are challenged by many social issues, lots of which relate with environmental threats. Plastic waste's impact on societal well being, economy, and environment has already triggered the need for its revaluation and its transformation into an asset which can inspire and mobilize innovative solutions. This work envisions trusted “wisdom of the crowds” analytics and a grassroots-driven approach to facilitate citizens and stakeholders participatory and engagement practices, under an inspiring crowdsourcing framework and model. A detailed evaluation of the filtering methodology is employed. This three-layer approach for the Data Filtering enables a Crowdsourcing Component implementation to be used as a plastic-topics insightful barometer, able to detect new trends, identify interesting content, spot influential users and generally raise awareness regarding the problem of plastic overuse.

CCS CONCEPTS

• Information systems~Crowdsourcing • Information systems~RESTful web services • Computing methodologies~Topic modeling

KEYWORDS

Web 2.0 analytics, crowdsourcing tools and models, social innovation, social media innovation

ACM Reference format:

Ilias Dimitriadis, Vasileios G. Psomiadis and Athena Vakali. 2019. A crowdsourcing approach to advance collective awareness and social good practices. In *Proceedings of WI '19: IEEE/WIC/ACM International Conference on Web Intelligence (WI '19 Companion), October 14-17, 2019, Thessaloniki, Greece*. ACM, New York, NY, USA, 8 pages. <https://doi.org/10.1145/3358695.3361104>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

WI '19 Companion, October 14-17, 2019, Thessaloniki, Greece

© 2019 Association for Computing Machinery.

ACM ISBN 978-1-4503-6988-6/19/10...\$15.00

<https://doi.org/10.1145/3358695.3361104>

1 Introduction

Contemporary societies are challenged by many social issues, lots of which relate with environmental threats. Plastic waste poses such a huge threat to society and the environment, from overflowing landfills to the ever-growing ocean gyres since it's damaging fragile ecosystems and affecting the food chain. A huge waste of plastic resources, threats environment but also lives of organisms and ever humans. Plastics drawbacks become a large-scale problem which needs a global response. Plastics waste impact on societal wellbeing, economy, and environment has already triggered the need for a novel and innovative Plastics Economy¹ movement.

Novel solutions are now needed for enabling social innovation in revaluing plastics. This work is motivated by the capacity of social Web technologies to make social change and innovation under a plastics-as-an-asset paradigm. Since social Web technologies have evolved with capturing the so called “wisdom of the crowd”, in the proposed work Citizens to Communities (C2C) associations are detected and shared by harnessing crowdsourcing via a social media observatory with favors plastics economy participatory motivation by active user's engagement and plastics social innovation co-sharing. This work envisions trusted “wisdom of the crowds” analytics and a grassroots-driven approach to facilitate citizens and stakeholders participatory and engagement practices, under an inspiring crowdsourcing framework and model.

Data produced in abundance in online social media sources offer a fertile ground for harvesting users' feedback and views concerning their reactions and opinions on evolving plastic recycling and reusing key issues. The exploitation of such resources can lead to the delivery of innovative and more human-centered services for the public, leveraging the collective intelligence of the crowd. Social media data mining as a collective intelligence approach, has tremendously evolved during the last decade. It involves extracting latent information and insights from the: (i) unstructured social media User Generated Content (UGC) and (ii) users' interaction in social media, such as communities or groups of similar behaviour. These facts (massive data sizes and low-quality content, e.g. poorly formatted, short textual entities, etc.) pose significant challenges to typical data mining algorithms.

By a proper design and implementation of a social media ideas sharing terminology and an application release with crowdsourcing mechanisms plastics reuse actions can be

¹ The New Plastics Economy: Catalysing action; World Economic Forum, Jan 2017. <https://www.weforum.org/reports/the-new-plastics-economy-catalysing-action>

that will be able to identify qualitative content, and further classify topics in a set of context-specific categories. Analysis of the metadata and information around the collected UGC via topic modelling approaches is accompanied by textual-based classification approaches to extract context along with features that will enable categorization in appropriate, predefined classes (such as *reuse ideas*, *initiatives*, etc.). An advanced machine learning methodology is used to provide a much richer representation of text for more robust and accurate social data analysis. The proposed machine learning approach exploits quantitative measures (number of comments, 'likes', etc.), as well as qualitative analytics building on analytics approaches [10] and other plastics contextual feedback such in case of recycling [11] and green social behaviour [12]. In the proposed approach, a methodology for automatically collecting and identifying content relevant to specific topic areas (in this case *plastic-as-an asset* and relevant mitigation approaches) is used inline with a preliminary analysis (outlined above in Section 2). The proposed methodology steps include:

- acquiring coarse-grained streams of content (such as posts and annotated multimedia and text content, as well as user reactions to them) via appropriate data collection mechanisms;
- the identification and characterization of qualitative relevant information and the fine-tuning of the data collection parameters;
- exploiting and advancing the state of the art in the areas of qualitative crowdsourcing and on dynamic topic categorization.

The proposed data filtering methodology during the process of the Crowdsourced Data Collection and Analysis involves three main filtering types:

1. Primitive Filtering, to identify the Social media sources which will be used for the collection of data;
2. Dynamic content filtering, to identify terms related to plastic and plastics reuse;
3. Data streams & User filtering, for
 - a. Identifying and detecting influential – expert - users regarding plastic re-use thematology.
 - b. Filtering incoming data in order to produce high quality content.
 - c. Updating the keywords used to collect data.

Apart from the whole filtering process, a detailed evaluation of the filtering methodology is employed. This three-layer approach for the Data Filtering enables a Crowdsourcing Component implementation to be used as a plastic-topics insightful barometer, able to detect new trends, identify interesting content, spot influential users and generally raise awareness regarding the problem of plastic overuse.

3.1 Primitive Filtering

Data collection from social media is rapidly getting progressively limited since many concerns are raised regarding privacy issues and the exploitation of user generated data through the OSNs while recent legislation frameworks are becoming increasingly more rigid towards this sensitive subject. Thus, very few of the popular social networks allow developers to access data, even if these are meant as public posts by their

owners. The social networks that are used as sources are the following:

1. Twitter²: Twitter is the main data source for the crowdsourcing tool.
2. Flickr³: Provides access to images that have been described with certain tags relevant to plastic related terms.
3. Thingiverse⁴: Acts as the main source of the open 3D printer designs repository and is used as an interlink, where the users are able to upload their own designs.
4. Instagram⁵: Used as a secondary source, because of the access limitations to its data.

3.2 Dynamic Content Filtering

Collecting information from various sources, implies the use of specific filters so that the content of the collected data remains as close to the intended thematic as possible. For the purpose of this work the primary goal was the identification and detection of simple and sophisticated content about plastic pollution, reuse and recycling along with their correlations, trends and phenomena reflected in the society.

To identify content in Twitter a group of human experts provided an extensive list of specific keywords covering a very wide scope of the plastic thematic. These keywords were then classified using a taxonomy provided by the same expert group, forming an initial set of different categories and topics that include:

- General terms related to plastic;
- Plastic innovations;
- Plastic machines (for re-valuating plastic waste);
- Products made of plastic);
- Plastic processes (for re-valuating plastic waste);
- Plastic pollution.

Along with English this list was further translated to German, Dutch and Greek in order to capture trends and interesting topics in local level. Twitter's API was queried for all tweets that contain any on the keywords in the four languages.

Since a significant percentage of the tweets may be noisy, spam or offensive, a profanity filter was used in order to ensure the quality of the final delivered content. In our initial evaluation of the collected data, except some swearing language, no abusive content was spotted. This is mainly due to the fact that the plastic thematic does not offer a fertile ground for adult content. However, a typical text filter has been applied to filter tweets with inappropriate words included in their text. After detecting such tweets, the next step is identifying users as possibly "malicious" and blacklisting them so they are ignored by the Twitter crawler.

Apart from collecting streaming tweets, data has been collected for a number of specific user accounts that are considered to be experts in the particular field of plastic waste reuse and recycling.

Flickr is used to deliver related image content by exploiting the intelligence that has been captured from the Twitter posts

² <https://twitter.com>

³ <https://www.flickr.com>

⁴ <https://www.thingiverse.com>

⁵ <https://www.instagram.com>

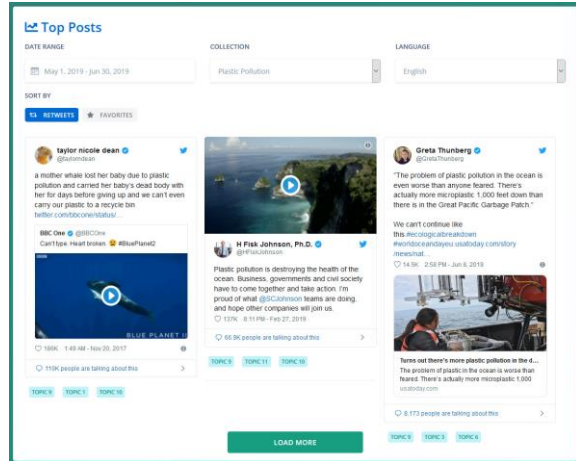


Figure 5: Top posts by retweets in the plastic pollution collection (May - June 2019)

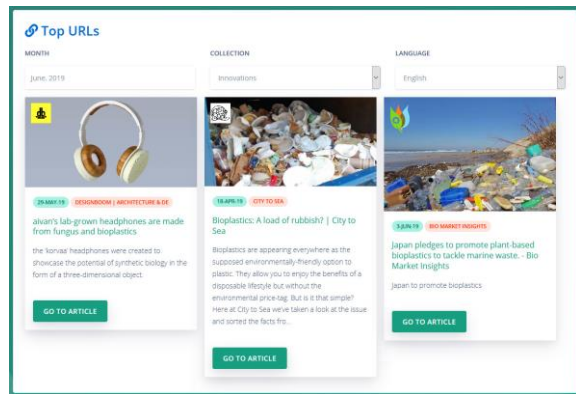


Figure 6: Top URLs in the plastic innovations collection (June 2019)

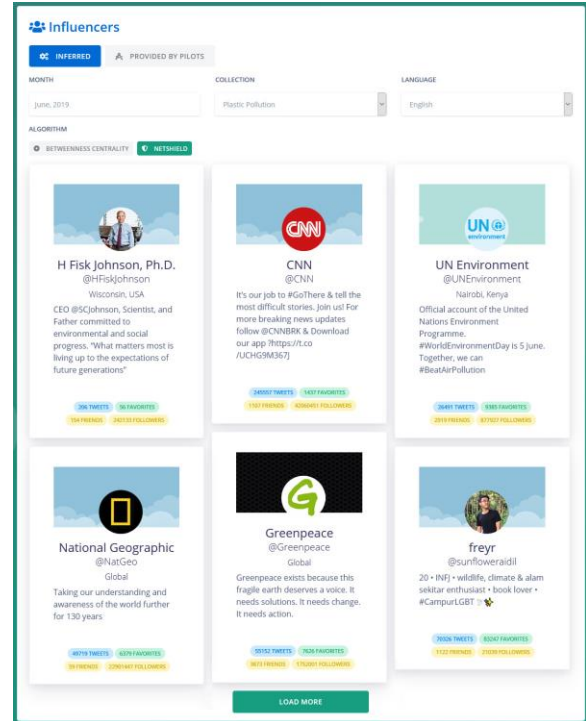


Figure 7: Top influencers regarding plastic pollution by NetShield (June 2019)

The proposed crowdsourcing approach offers an exploitable implementation since it is generically designed under a data scheme capable of integrating social media data streams with plastics classification terms hierarchies. Results of crowdsourcing are visible to the open public by an applications level observatory interface. This open Crowdsourcing API has been developed with Flask, a Python framework for building RESTful APIs. Under this interoperability scheme the research community can integrate open crowdsourced data and wisdom of the crowd's knowledge to their current research and experimentation, especially in social media analytics. At the same time, authorities can gain insights of most trendy plastics relevant topics discussed in social media. This insight impacts decision making with respect to crowd's pulse detection. The following Figure 8 and Figure 9, contain examples of related external content that is discovered and presented to the user in real time.

A crowdsourcing approach to advance collective awareness and social good practices

WI '19 Companion, October 14-17, 2019, Thessaloniki, Greece

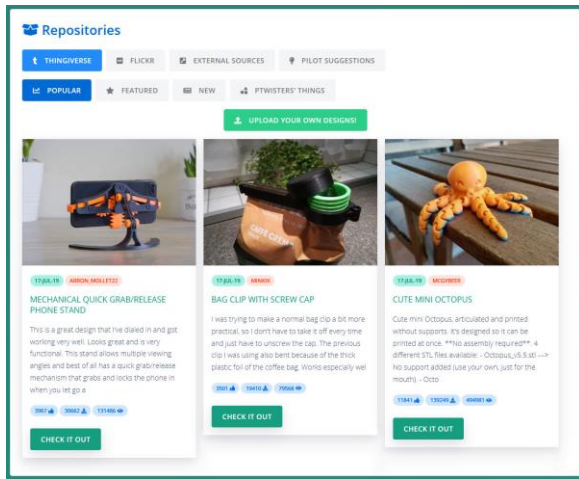


Figure 8: Popular open 3d printer designs from Thingiverse

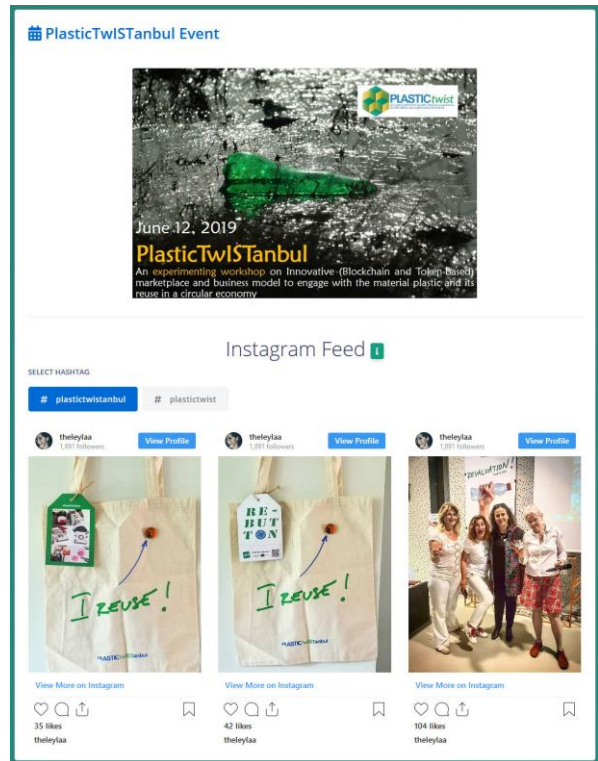


Figure 10: Social media posts of a specific event fetched from Instagram

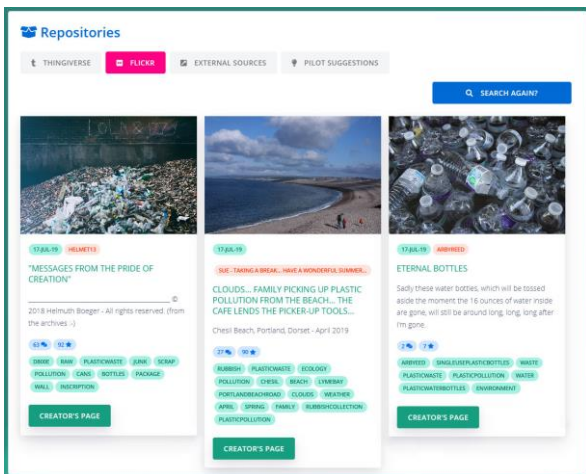


Figure 9: Photos related with the “PLASTICWASTE” term fetched from Flickr

Figure 10 contains social media posts from Instagram collected using a pre-selected hashtag, connected to a specific activity regarding plastic waste revaluation, that was chosen to spread word in social media.

5 Conclusions and Future Work

This work has addressed the potential of crowdsourcing and social media analysis for cases of societal and environmental threats, such as in the case of plastics revaluing. The publicly-available developed components along with the open API managed to serve a tri-fold purpose successfully:

- reflect public opinion and raise awareness around a critical issue that concerns today’s society;
- inform and educate on the thematic of plastic waste pollution and reuse; and
- disseminate novel solutions and provide inspiration for further grass-rooted innovation.

The followed approach, has focused on a multi-type filtering process which has been used during the data collection and analysis phase. Extracting high quality content from the users that have been identified as influencers by our system and use it to train an LDA model is a future approach, which can be used to classify other users, and extract topics using topic modelling per location. An iterative methodology can also be applied in the future work and it can be built upon the intelligence extracted by the already available high-quality content (top tweets – top URLs) to identify new trends and dynamically update the keywords used to track tweets of specific content.

ACKNOWLEDGMENTS

This work has been supported by the European Union's Horizon 2020 Research and Innovation Programme under Grant Agreement No. 780121.

REFERENCES

- [1] A. Zubiaga, D. Spina, R. Martínez, V. Fresno (2015). Real-Time Classification of Twitter Trends. *Journal of the Association for Information Science and Technology*, 66(3), 462–473.
- [2] D. Ramage, S. Dumais, D. Liebling (2010). Characterizing microblogs with topic models. In *Proceedings of the Fourth International AAAI Conference on Weblogs and Social Media*.
- [3] Ya. Duan, Z. Chen, F. Wei, M. Zhou, H.-Y. Shum (2012). Twitter Topic Summarization by Ranking Tweets using Social Influence and Content Quality. *Proceedings of COLING 2012*: 763-780.
- [4] T.-A. Hoang, W.W. Cohen, E.-P. Lim, D. Pierce, and D.P. Redlawsk (2013). Politics, sharing and emotion in microblogs. In *Proceedings of the 2013 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM '13)*. ACM, New York, NY, USA, 282-289.
- [5] A. Kanavos, I. Perikos, P. Vikatos, I. Hatzilygeroudis, C. Makris and A. Tsakalidis (2014). Modeling ReTweet Diffusion using Emotional Content, *IFIP International Conference on Artificial Intelligence Applications and Innovations (ALAI 2014)*, Rodos, Greece.
- [6] J. Zhang (2015). Voluntary information disclosure on social media. *Decision Support Systems*, 73, 28–36.
- [7] Y. Li, A.F. Smeaton (2014). From Smart Cities to Smart Neighborhoods: Detecting Local Events from Social Media. In: *Information Access in Smart Cities (i-ASC) 2014 Workshop in conjunction with ECIR 2014*, Amsterdam, Netherlands.
- [8] B. Debatin, J.P. Lovejoy, A.-K. Horn, M.A., B.N. Hughes (2009). Facebook and Online Privacy: Attitudes, Behaviors, and Unintended Consequences. *Journal of Computer-Mediated Communication*, 15(1), 83–108.
- [9] A. Gattani, D.S. Lamba, N. Garera, M. Tiwari, X. Chai, S. Das, S. Subramaniam, A. Rajaraman, V. Harinarayan, and A. Doan (2013). Entity extraction, linking, classification, and tagging for social media: a wikipedia-based approach. *Proc. VLDB* 6(11), 1126–1137.
- [10] D. Chatzakou, N. Passalis, A. Vakali (2015). Multispot: Spotting sentiments with semantic aware multilevel cascaded analysis. *International Conference on Big Data Analytics and Knowledge Discovery*. Springer International Publishing, 337–350.
- [11] M. Dupré, S. Meineri (2016). Increasing recycling through displaying feedback and social comparative feedback. *Journal of Environmental Psychology*, 48, 101–107.
- [12] A. Biswas, R. Mousumi (2016). Impact of Social Media Usage Factors on Green Choice Behavior Based on Technology Acceptance Model. *Journal of Advanced Management Science*, 4(2), 92–97.
- [13] D. M. Blei, A. Y. Ng and M. I. Jordan (2003). Latent dirichlet allocation. *Journal of Machine Learning Research*, 3, 993–1022.
- [14] D. Jurgens, T. Finethy, J. McCorriston, Y. T. Xu and D. Ruths (2015). Geolocation Prediction in Twitter Using Social Networks: A Critical Analysis and Review of Current Practice. *ICWSM*, 15, 188–197.
- [15] H. Tong, B. A. Prakash, C. Tsourakakis, T. Eliassi-Rad, C. Faloutsos and D. H. Chau (2010). On the vulnerability of large graphs. In *Proceedings - IEEE International Conference on Data Mining, ICDM*, 1091-1096.
- [16] M.Kitsak, L.K. Gallos, S. Havlin, F. Liljeros, L. Muchnik, H.E. Stanley, H.A. Makse (2010). Identification of influential spreaders in complex networks. *Nature Physics*, 6(11), 888–893.