

CASE STUDIES

Determining the Accuracy of Crowdsourced Tweet Verification for Auroral Research

Nathan A. Case^{*,†,‡}, Elizabeth A. MacDonald^{*,†}, Sean McCloat^{†,§}, Nick Lalone^{||} and Andrea H. Tapia^{||}

The Aurorasaurus project harnesses volunteer crowdsourcing to identify sightings of an aurora (the “northern/southern lights”) posted by citizen scientists on Twitter. Previous studies have demonstrated that aurora sightings can be mined from Twitter with the caveat that there is a large background level of non-sighting tweets, especially during periods of low auroral activity. Aurorasaurus attempts to mitigate this, and thus increase the quality of its Twitter sighting data, by using volunteers to sift through a pre-filtered list of geolocated tweets to verify real-time aurora sightings. In this study, the current implementation of this crowdsourced verification system, including the process of geolocating tweets, is described and its accuracy (which, overall, is found to be 68.4%) is determined. The findings suggest that citizen science volunteers are able to accurately filter out unrelated, spam-like, Twitter data but struggle when filtering out somewhat related, yet undesired, data. The citizen scientists particularly struggle with determining the real-time nature of the sightings, so care must be taken when relying on crowdsourced identification.

Keywords: twitter; crowdsourcing; aurora; sightings; citizen science

Introduction

The citizen science project Aurorasaurus (MacDonald et al. 2015) has two main goals: Improving the “nowcasting” of a visible aurora (commonly known as the “northern/southern lights”) and the ability to accurately model both the size and strength of an aurora. To do this, the project collects observations of the aurora made by the general public. These observations can be submitted directly to the project, via its website (<http://aurorasaurus.org>) and mobile apps, and are found by searching Twitter for possible sightings.

Twitter can be a useful source of data for many citizen science projects because information is freely shared by millions of users distributed around the globe. Indeed, previous studies have shown that Twitter users, who post short updates (of a maximum 140 characters in length) known as “tweets,” will often share details about the conditions around them. This is especially true for large-scale events such as earthquakes (Earle et al. 2010; Crooks et al. 2013), influenza outbreaks (Culotta 2010; Lamos

et al. 2010), and service outages (Motoyama et al. 2010). Case et al. (2015a) showed that Twitter can also be a useful source of data for studying the aurora by comparing the number of tweets relating to an aurora with auroral activity (or, more specifically, to common auroral activity indices). However, these authors also noted that Twitter data are particularly noisy and that many tweets containing aurora-related keywords (e.g., “aurora” and “northern lights”) are not actually sightings. Often such tweets are about a person or place or the desire to witness an aurora.

The Aurorasaurus project enlists volunteers, both registered and anonymous, to sort through pre-filtered, aurora-related tweets to identify and positively verify real-time aurora sightings. While combining Twitter data with other citizen science data may be a new form of crowdsourcing, many previous studies have demonstrated that crowdsourcing can be used for data classification, often using Amazon’s Mechanical Turk (Kittur et al. 2008; Ipeirotis et al. 2010). In fact, studies have shown that the crowd is sometimes more accurate than experts at identification tasks (Alonso and Mizzaro 2009).

Once a tweet has been verified as a positive sighting by the Aurorasaurus volunteers, it is treated in the same way as a direct report via the project’s website or apps. The combined observations, both direct reports and positively verified tweets, are displayed on the project home page on a real-time map alongside a modeled auroral oval (i.e., the extent to which an aurora is visible directly overhead). These observations serve several different functions,

* NASA Goddard Space Flight Center, Greenbelt, MD, USA

† New Mexico Consortium, Los Alamos, NM, USA

‡ Department of Physics, Lancaster University, Lancaster, UK

§ University of North Dakota, Grand Forks, ND, USA

|| College of Information Sciences and Technology, Pennsylvania State University, State College, PA, USA

Corresponding author: Nathan A. Case (n.case@lancaster.ac.uk)

including demonstrating where the aurora is currently being observed (Priedhorsky et al. 2012), providing data points for scientific investigation (Case et al. 2016), and providing the basis for a hybrid alert system (Lalone et al. 2015) that is analogous to disaster early warning systems (Tapia et al. 2014).

This study investigates the accuracy of volunteers in filtering useful data from a stream of tweets in an existing citizen science project. The results provide insights into the accuracy of volunteers in analysing Twitter data that may be applied to other citizen science projects.

Tweet Verification

Aurorasaurus exploits the Twitter Search API to identify publicly accessible tweets that contain any one of several different aurora-related keywords (e.g., “aurora;” “northern lights.”) The returned tweets are then filtered further on the Aurorasaurus servers to exclude most retweets, tweets from Twitter users with “aurora” in their username (although a whitelist is maintained to allow tweets from some users to go through), and tweets containing profanity or other common “spam” terms.

A location extraction process is then undertaken on the filtered tweets. Location is determined either by using the embedded GPS metadata, if the Twitter user has opted to share their location, or through the geo-parsing software CLAVIN (<https://clavin.bericotechnologies.com>), which attempts to extract a location for a tweet based upon its text (D’Ignazio et al. 2014). Using these processes, approximately 15% of the tweets can be associated with a location (with extraction through CLAVIN accounting for approximately 80% of the associations). Further filtering takes place to remove tweets whose location is determined to be anywhere containing the term “Aurora” (e.g., Aurora, CO, USA).

These “unverified tweets” are then presented to the Aurorasaurus community for verification as pins on

the main map and as a list on the “Verify Tweets” page (see **Figure 1**). The community is asked “Did they just see the aurora?” (where “they” refers to the tweet’s author) and are provided only two choices for a vote (“yes” or “no”). This subjective task allows automatic aggregation of the votes into a score and a classification based upon that score (Iren and Bilgen 2014).

For every “yes” vote a tweet receives, a value of 1 is added to its score. Conversely, for every “no” vote a tweet receives, a value of 1 is subtracted from the score. Votes from both registered and anonymous users are treated equally (i.e., there is no weighting applied to the vote based upon the user or their credentials). Once the tweet’s score reaches a certain positive threshold (currently set to +3), it is categorized as a “positively verified tweet;” its marker is updated on the map to show this new status; and votes are no longer taken on it. Similarly, once a tweet reaches a certain negative threshold (currently set to -3), the vote is categorized as a “negatively verified tweet;” the marker is removed from the map; and the tweet is no longer presented to the community for verification.

To reduce the barriers of entry for users to start verifying tweets, no compulsory training is required. However, help in verifying tweets is provided by a pop-out help menu, which opens if the user clicks on the question mark in the tweet window (see **Figure 1**). Additionally, a blog post and quiz are available, both of which guide the voter through examples of tweets and how they should be voted upon. Approximately half the respondents to a recent Aurorasaurus survey indicated that they had read at least some of this guidance (Lalone pers. comm., 2015).

Results

This study analyzes the verified tweets posted during March and April, 2015. This two-month period represents a subset of the larger Aurorasaurus data set (which spans from November 2014 to present) and includes several

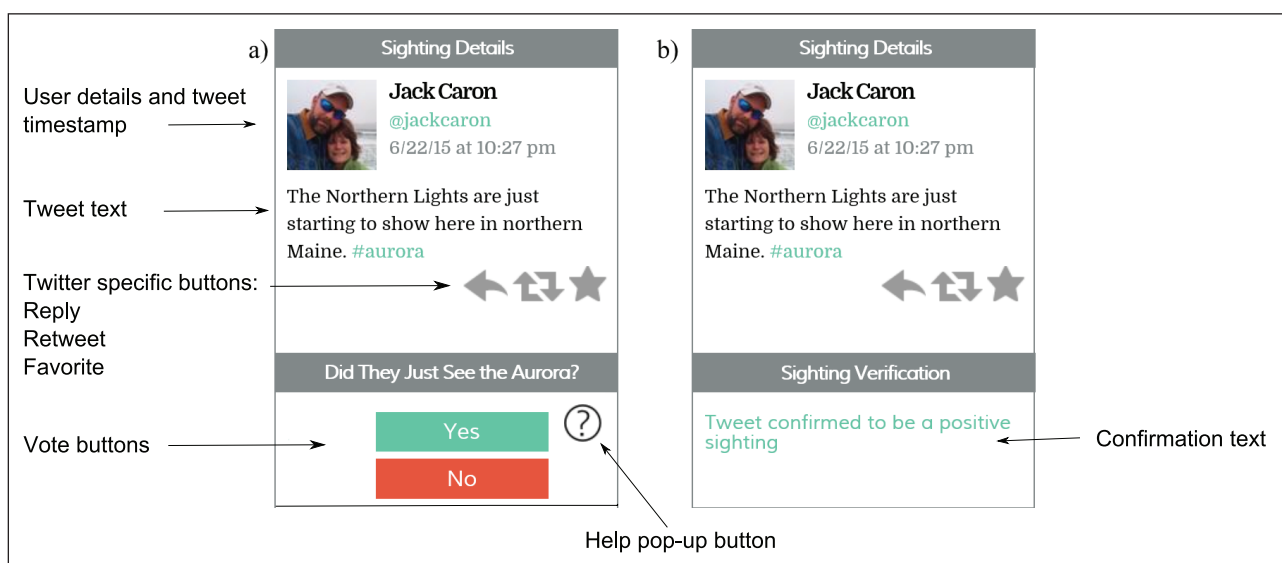


Figure 1: a) An example tweet as presented to the Aurorasaurus community for verification. The volunteers are asked “Did they just see the aurora?” and are given the two simple options of “yes” (for a positive, real-time, aurora sighting) or “no.” **b)** Once a threshold positive score is reached, the tweet is confirmed as a “positive sighting” and becomes known as a “positively verified tweet.” It is then no longer available for further voting.

large auroral events, including the largest event this decade (Case et al. 2015). It is important to note that large auroral events, where an aurora can be seen from the mid-United States and central Europe, are relatively infrequent and are dependent upon several factors including solar activity, time of day/year, and local conditions (e.g., cloud cover). Additionally, an aurora can be a widespread phenomenon, with sightings of the same event spanning multiple continents (Case et al. 2015).

During March and April, 2015, 227,280 aurora-related tweets were collected with 39,636 (17.4%) having an associated location and thus available for the Aurorasaurus community to vote on. Of these, the community verified 4,547 (11.5%) tweets: 475 positively (10.4%) and 4,072 negatively (89.6%). There were 70,331 votes cast: 49,495 by logged-in users (70.4%) and 20,836 by anonymous users (29.6%).

The distribution of the tweets and their verified status is shown in **Figure 2**. The number of each type of tweet (“total,” “with location,” “positively verified,” “negatively verified,” and “unverified”) is shown by the filled bars. Note the logarithmic scale on the y-axis.

Each of the positively verified tweets was then independently manually inspected by two members of the Aurorasaurus team. This inspection involved analyzing the text of the tweets in detail to identify any signs of non-originality and to compare the location and time of the supposed sighting with auroral models and other citizen science observations.

The verified tweets were categorized primarily into “valid” (where the tweet was indeed a real-time aurora sighting made by the tweet’s author) or “invalid” (where the tweet was incorrectly positively verified by the users). Using an open-coding method, the following categories for the invalid positively verified tweets were created:

- “Not real-time”: a sighting of an aurora by the tweet’s author, however, the tweet was posted at least several hours after the sighting took place (often the next morning).
- “Not original”: the sighting was not made by the tweet’s author (usually “retweets” or “mentions” of someone else’s tweet).
- “Overlap”: the sighting was not real-time nor was it made by the tweet’s author. This would often be the retweeting of someone else’s aurora photograph.
- “Wrong location”: the location extraction algorithm (CLAVIN) failed to determine the location correctly. These failures are particularly difficult for voters to spot, because the location of the tweet is not shown on the tweet (see **Figure 1**).
- “Not positive sighting”: the tweet did not contain a sighting of an aurora but may have been related to one (e.g., “Seeing an aurora is on my bucket list”).
- “Junk”: these tweets had nothing to do with an aurora (e.g., “Went to Aurora last night”).

The distribution of these categories is shown in **Figure 3**.

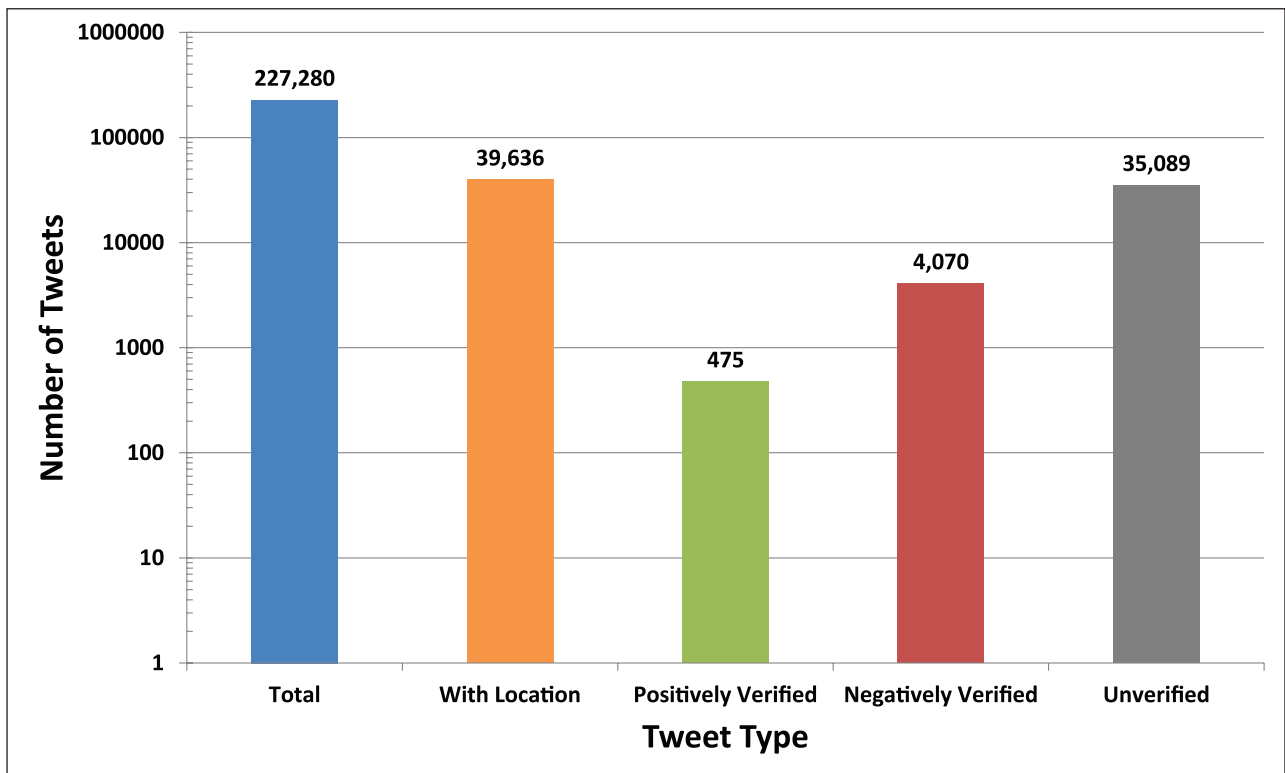


Figure 2: The distribution of tweets collected during March and April 2015. The first (blue) bar indicates the total number of tweets collected. The second (orange) shows the number of tweets with an associated location and thus available for the Aurorasaurus community to vote on. The third (green) bar shows the number of positively verified tweets, while the fourth (red) shows the number of negatively verified tweets. The final (gray) column is the number of tweets that were not verified (i.e., “unverified”).

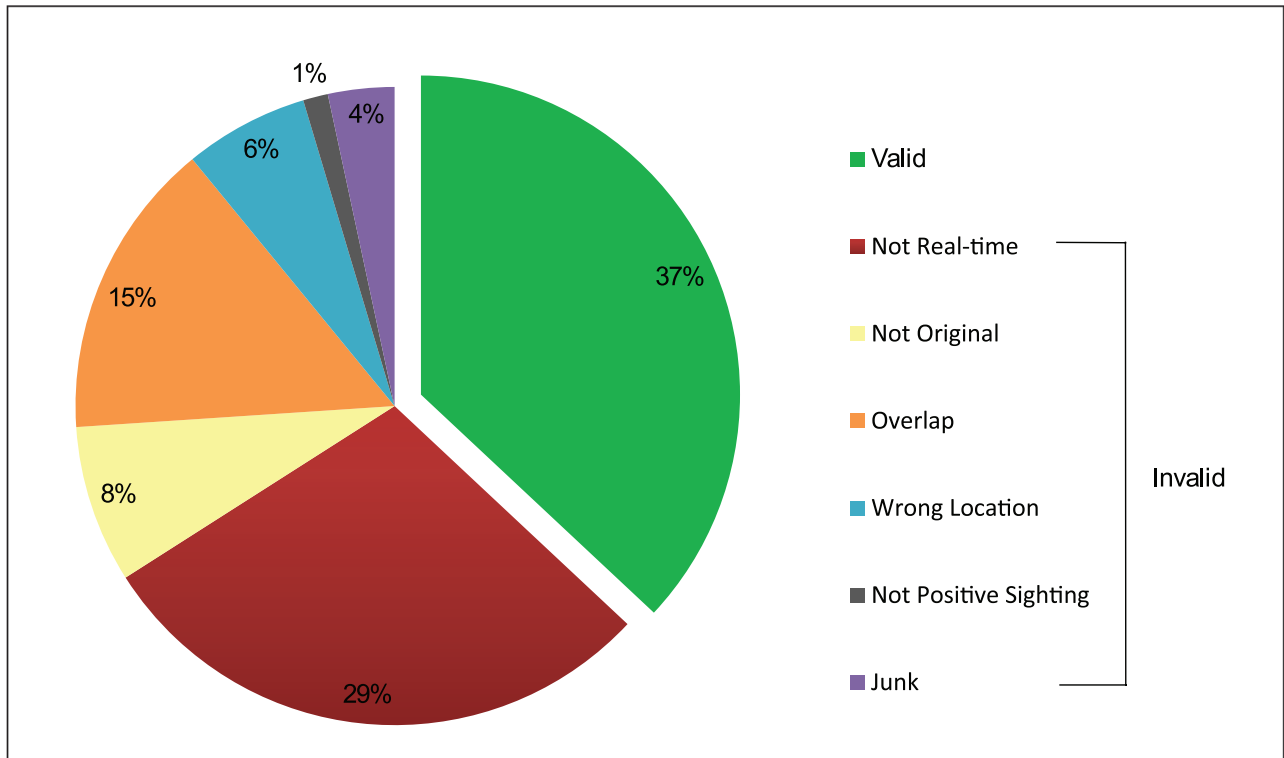


Figure 3: The distribution of positively verified tweets collected during March and April 2015. The tweets are grouped by the previous categories: valid (green), not real-time (red), not original (yellow), overlap (orange), wrong location (blue), not a positive sighting (black), and junk (purple).

Of the 475 positively verified tweets, 176 (37%) are valid. The precision, or positive predictive value (PPV), as calculated using Equation 1, of the positively verified tweets is therefore 37.1%.

$$PPV = \frac{\Sigma TP}{\Sigma TP + \Sigma FP} \quad (1)$$

where ΣTP is the number of true positives (i.e., positively verified tweets that are valid) and ΣFP is the number of false positives (i.e., positively verified tweets that are invalid).

The process was then repeated for a sample of the negatively verified tweets. This randomly selected sample included 475 negatively verified tweets (chosen to match the number of positively verified tweets). All but two of the tweets in the sample were correctly identified as negatively verified tweets. Thus, the “negative precision,” or negative predictive value (NPV), as calculated using Equation 2, was 99.6%.

$$NPV = \frac{\Sigma TN}{\Sigma TN + \Sigma FN} \quad (2)$$

where ΣTN is the number of true negatives (i.e., negatively verified tweets that are not valid sightings) and ΣFN is the number of false negatives (i.e., negatively verified tweets that are actually valid sightings).

The overall accuracy of the verified tweets, in which all of the positively verified tweets and a same-sized sample of negatively verified tweets are included, can now

be determined. Using Equation 3, the overall accuracy is found to be 68.4%.

$$ACC = \frac{\Sigma TP + \Sigma TN}{N} \quad (3)$$

where N is the total number of verified tweets in this sample (i.e., $N = 950$).

Furthermore, these results can be broken up based upon periods of when auroral activity was particularly elevated (which is when most sightings would be expected to occur). Three such events occurred during this time period: March 01–03, March 17–19, and April 10–12. The distributions of the previous categories are shown, for each of these periods, along with the distribution of “non-elevated” periods, in **Figure 4**.

The negatively verified tweets also were split by storm period. Both of the invalid negatively verified tweets occurred during the March 17–19 storm (which is not particularly surprising due to the majority of tweets occurring during this time). The PPV, NPV, and ACC are calculated for each of these storm periods and are presented in **Table 1**.

Discussion

Approximately 17.4% of the 227,280 tweets collected during this case study had a location associated with them, which is consistent with other studies (e.g., Vieweg et al., 2010). Thus, nearly 40,000 tweets were available for the Aurorasaurus community to vote on. Approximately 75% of the locations obtained were determined using the CLAVIN geo-location extraction algorithm,

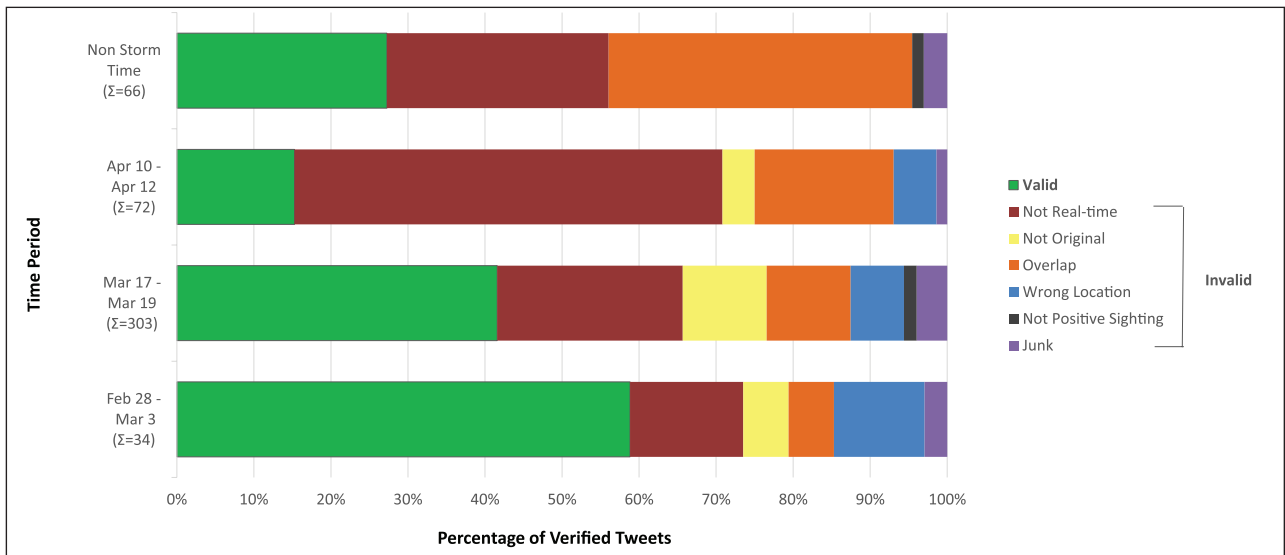


Figure 4: The positively verified tweets have been split into three active auroral time periods and one non-storm period. For each period, the percentage share of each category listed earlier is shown.

| Period | N | N _{pos} | N _{neg} | PPV (%) | NPV (%) | ACC (%) |
|----------------|-----|------------------|------------------|---------|---------|---------|
| March 01–03 | 44 | 34 | 10 | 58.8 | 100.0 | 79.4 |
| March 17–19 | 461 | 303 | 158 | 41.6 | 98.7 | 70.2 |
| April 10–12 | 117 | 72 | 45 | 16.7 | 100.0 | 58.4 |
| Non-Storm Time | 328 | 66 | 262 | 27.3 | 100.0 | 63.7 |
| Overall | 950 | 475 | 475 | 37.1 | 99.6 | 68.4 |

Table 1: Tweet numbers and verification accuracy, split by periods of auroral activity.

therefore, only a small percentage of the total tweets contained an embedded GPS location. Again, this result is consistent with other studies (e.g., Cheng et al. 2010, Lee et al. 2013).

The community cast more than 70,000 votes and verified over 4,500 tweets. The majority, around 80%, of verified tweets were negatively verified, i.e., the Aurorasaurus community voted that the tweet was not a real-time sighting of an aurora made by the tweet’s author. This result is perhaps unsurprising, because it is generally only when auroral activity is high (which occurred three times during this case study) that increased numbers of people tweet sightings of an aurora (Case et al., 2015a). Indeed, the percentage of positively verified tweets (i.e. N_{pos}/N) rises from around 20% during non-storm times to around 70% during active times (Table 1).

Notably, nearly 90% of tweets with locations went unverified (i.e., they were not positively or negatively verified). These tweets are most likely not aurora sightings; rather, they are tweets that contain aurora-related keywords. However, we cannot be certain that this set of tweets contains sightings that have simply been overlooked. While this does not affect the accuracy of the verification system, it does mean that some scientifically useful observations, such as rare sightings during low auroral activity, might be missed. Further investigation into the exact nature of the unverified tweets, and what effect the number

of unverified tweets may have on citizen science data collection on Twitter, should therefore be undertaken.

Verification Accuracy

The Aurorasaurus community was able to negatively verify tweets with extremely high accuracy. In fact, of the 475 negatively verified tweets analyzed, only two were incorrectly classified, resulting in an overall NPV of nearly 100%. The community was, however, much less accurate when positively verifying tweets. The overall PPV (or precision) was 37%, though significant variance occurred in the PPVs when splitting by event (with the highest PPV of 59% occurring during the March 01–03 storm and the lowest PPV of 27% occurring during the April 10–12 storm). At this time no reason is known for this variance unless it is attributable to differences in the sample sizes.

The overall accuracy of the verification system in this case study was 68%. Had all of the negatively verified tweets been analysed, and subsequently used in the accuracy calculation, the overall accuracy would probably have been much higher. However, because the number of negatively verified tweets was so much greater than the number of positively verified tweets, a representative sample was chosen instead. Note that the positively verified tweets (i.e., actual sightings) hold the most scientific value, so the PPV may be more important than the NPV or overall accuracy.

What affected the community's precision?

Spotting spam-like tweets that have nothing to do with sightings of an aurora is relatively easy. Much harder is differentiating between tweets that are real-time aurora sightings from those that are just related to the aurora or are true sightings that occurred several hours previous. Indeed, our analysis showed that the primary reason the community positively verified tweets incorrectly was that the community incorrectly identified the tweets as being real-time.

Identifying whether a sighting posted in a tweet is real-time can be a complex task, even for the Aurorasaurus team members. The tweet has a timestamp associated with it, but the tweet's author may be posting about a sighting that occurred several hours ago or perhaps even the day before. Unless the author explicitly uses words or phrases that chronologically identify when the sighting occurred, e.g., "just seen" or "spotted 10 mins ago," knowing exactly when the sighting occurred is difficult. In fact, even if the author includes a time, e.g., "aurora seen at 21:30," the verifier would need to know the offset between their current time zone and the time zone of the tweet's author to determine how long ago the aurora was sighted. Such detailed investigation is probably too much for most of the community to engage in, especially when they are voting on many tweets at once.

The second most common reason for incorrectly positively verifying a tweet was that the sighting was "not original." From this category we identified two themes: The tweet was of someone else's aurora photograph (85%) or the tweet was a retweet of somebody else's sighting (15%). Both of these errors likely stem from unfamiliarity with Twitter's nomenclature. For example, most of the "not original" tweets contained signs of the non-originality, i.e., the text "RT" (an acronym for retweet) or tagging of other users (which will always start with the @ symbol). We note, however, that many original real-time sightings may also tag other users, often as a way of alerting them, so this method to determine originality cannot be used on its own.

Improving the voting system

When the community incorrectly positively verifies a tweet we assume an "honest mistake" rather than a "cheater" (i.e., someone with malicious intent) because there is no gain to poor verification (Hirth et al. 2013, Iren and Bilgen 2014). Therefore, a primary way to improve the accuracy of the crowd is to improve the information provided about the task and the desired outcome (Iren and Bilgen 2014). Aurorasaurus currently provides its community with instructions/guidance via a help page, blog post, and a quiz (where members of the community can test their voting skill and receive feedback on their choices). These are all "hidden elements," however, as a user may not have seen them before beginning to vote. Indeed, a recent survey of Aurorasaurus users showed that 40% did not know that instructions on how to verify tweets were available (Lalone pers. comm. 2015).

Enforcing training upon community members before they are able to vote has been shown to improve the quality of voting (e.g., Le et al. 2010). In some implementations,

training results in a pass/fail that screens out untrustworthy or inaccurate users (Downs et al. 2010, Le et al. 2010). In others, the score attributed to each user's vote is weighted based upon how well they perform during the training (Sheng et al. 2014). We note, however, that these studies often employ contributors through Amazon's Mechanical Turk rather than volunteers in citizen science projects.

Because the Aurorasaurus project, like all citizen science projects, is reliant on volunteers, adding such compulsory activities might reduce the number of people who are willing to participate. Therefore, training that is not compulsory but that could be used to better inform the voting system on a user's trustworthiness might be desirable. For example, votes from anonymous users might be weighted to score 1, votes from registered users who have not taken the training might be weighted to score 2, votes cast by those who have taken the quiz but did not score highly might be weighted to 3, and votes from users who scored highly in the quiz might be weighted to 5. Project staff, or trusted super-users, might then have an even higher voting weight. This approach has the benefit of determining a pseudo-confidence level for each vote without erecting barriers to participation.

Vuurens et al. (2011) demonstrated that a "combined consensus algorithm," which generally used a majority vote but then took into account the voters' trustworthiness in a tie situation, consistently provided the most accurate results. A tied result, with respect to the Aurorasaurus crowdsourcing system, would be where the number of votes is over the verification threshold, however, the score has not exceeded that threshold (i.e., 10 users vote—five yes and five no—resulting in a score of 0).

The training, and subsequent vote weighting, is likely to be a one-time effort (although, in practice, users could be allowed to complete it more than once). One-time training could lead to situations where users forget what they have been taught or their voting is affected by other factors (e.g., fatigue or lack of concentration). To help mitigate the effect of "bad votes" from a trained user, an adaption of the "majority decision" cheat-detection method (Hirth et al. 2013) could be employed. If a member of the community votes against the current majority decision or the decision of a trusted voter (e.g., staff or super-user), they are advised in real-time and offered training/guidance on how they should vote. The frequency to which a user matches or does not match the majority can be stored, allowing a hybrid voting reputation to be built (Voyer et al. 2010). Based on this reputation, voting weights could again be applied.

In addition to improving the voting mechanism itself, another way to increase the quality of the verification process could be to improve the chance of a tweet being a valid sighting before presenting it to the community for validation. The current system simply uses a set of keywords for searching and another set for filtering. Machine learning, based on either a gold standard set or the community's voting, might improve the quality of the tweets being served to the community (Wang 2010, Becker et al. 2011, Truong et al. 2014). This approach was tested early in the Aurorasaurus project, however, it failed

to yield any noticeable improvements (MacDonald, pers. comm. 2015), indicating that further refinement may be needed on such an approach before it could be applied to this task successfully.

Conclusion

Like many citizen science projects, Aurorasaurus is heavily reliant upon a community of volunteers for providing data and for validating/classifying data. To complement the aurora sightings reported directly to the project, Aurorasaurus also systematically searches for observations of an aurora posted on Twitter, using the Twitter Search API and several rudimentary filters. A location is required for all sightings, so those tweets that do not contain an embedded location are passed through a location extraction algorithm that attempts to resolve a location for the tweet based upon its text. This process, while not always accurate, increases the number of usable tweets four-fold. Using a similar location extraction process is therefore recommended for other citizen science projects needing location data from tweets. Including Twitter as a data source has increased the number of observations for the Aurorasaurus project by nearly 100%. Exploiting Twitter as an available data source is therefore recommended for other citizen science projects that collect observational data.

Twitter observations are noisier than traditional citizen science reports, however, so they need more curation by both the volunteers and project staff. The Aurorasaurus community is therefore encouraged to verify these potential sightings using a simple crowdsourcing scoring system. The community is rewarded for its participation by a leader board, where each vote earns the volunteer 5 points, and by increased accuracy in localized auroral visibility alerts.

This Aurorasaurus case study has shown that volunteer citizen scientists are extremely adept at filtering out spam-like tweets and other non-aurora sightings. These tweets tend to form the majority of tweets presented to the Aurorasaurus community, especially during times with little auroral activity. For the random sample studied, the NPV of the “negatively verified” tweets was almost 100%. A good NPV is perhaps unsurprising, as filtering spam is a relatively easy task, though such a high score was somewhat unexpected. The volunteer community proved to be less accurate when identifying the true aurora sightings. The PPV, or precision, of the positively verified sightings was somewhat poor at 37%. The most common reason for the community incorrectly positively verifying a tweet was that the tweet was not real-time, followed by the tweet not being an original sighting.

While positively verifying tweets requires more detailed investigation than filtering out spam-like tweets, the PPV achieved certainly could be improved. As discussed, incorrect identifications were likely the result of honest mistakes, so the primary way to reduce them is to provide training for the community. Aurorasaurus does provide some training, although it is not compulsory. The “verifying tweets quiz,” which is the only interactive training offered, is detached from the verification process in that it is a completely separate entity and is not linked

in the “help” pop-up text (see **Figure 1**) when verifying tweets. Making any training compulsory will likely reduce the number of users who then participate in the verification process (Lintott, pers. comms. 2015). This is a quality-control cost that many projects must deal with (Iren and Bilgen 2014). However, small improvements, such as providing a link to the quiz during the verification process, are likely to increase the community’s accuracy, even if just a little, without affecting the number who are willing to participate.

Larger, systematic improvements, such as implementing vote weighting algorithms or the adaption of a real-time majority decision cheat-detection system, are likely to significantly improve the quality (particularly the PPV) of the community’s verification efforts. Such improvements will take time and resources to implement but should be on the future road map for the project.

The results of this case study suggest that other citizen science projects that plan to use volunteer crowdsourcing for data validation, especially for “noisy” data (e.g., tweets), should consider using some of the training or quality-control methods that we describe here. The information provided on Twitter by citizen scientists, and then verified by other volunteers, can be extremely useful. However, consideration must be given to training those volunteers who validate the data or else the accuracy of the crowd may be poor.

Acknowledgements

This material is based upon work supported, in part, by the National Science Foundation (NSF) under Grant #1344296. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of NSF.

Funding Information

Funding for SM was kindly provided by the North Dakota Space Grant Consortium.

Competing Interests

The authors have no competing interests to declare.

References

- Alonso, O. and Mizzaro, S., 2009. Can we get rid of TREC assessors? Using Mechanical Turk for relevance assessment. In Proceedings of the SIGIR 2009 Workshop on the Future of IR Evaluation, pp. 15–16.
- Becker, H., Naaman, M. and Gravano, L., 2011. Beyond trending topics: Real-world event identification on Twitter. *ICWSM*, 11: 438–441.
- Case, N.A., MacDonald, E.A., Heavner, M., Tapia, A.H. and Lalone, N., 2015. Mapping auroral activity with Twitter. *Geophys. Res. Lett.*, 4: 3668–3676. DOI: <http://dx.doi.org/10.1002/2015GL063709>
- Case, N.A., MacDonald, E.A. and Patel, K.G., 2015. Aurorasaurus and the St Patrick’s Day storm. *Astronomy & Geophysics*, 56(3): 13–14. DOI: <http://dx.doi.org/10.1093/astrogeo/atv089>
- Case, N.A., MacDonald, E.A. and Viereck, R., 2016. Using citizen science reports to define the equatorial extent

- of auroral visibility. *Space Weather*, 14: 198–209. DOI: <http://dx.doi.org/10.1002/2015SW001320>
- Cheng, Z., Caverlee, J. and Lee, K., 2010. You are where you tweet: A content-based approach to geo-locating twitter users. In Proceedings of the 19th ACM international conference on Information and knowledge management. ACM, pp. 759–768. DOI: <https://doi.org/10.1145/1871437.1871535>
- Crooks, A., Croitoru, A., Stefanidis, A. and Radzikowski, J., 2013. # Earthquake: Twitter as a distributed sensor system. *Transactions in GIS*, 17(1): 124–147. DOI: <http://dx.doi.org/10.1111/j.1467-9671.2012.01359.x>
- Culotta, A., 2010. Detecting influenza outbreaks by analyzing Twitter messages. *CoRR*. arXiv: 1007.4748.
- D'Ignazio, C., Bhargava, R., Zuckerman, E. and Beck, L., 2014. Cliff-clavin: Determining geographic focus for news. NewsKDD: Data Science for News Publishing, at KDD 2014.
- Downs, J.S., Holbrook, M.B., Sheng, S. and Cranor, L.F., 2010. Are your participants gaming the system? Screening mechanical turk workers. In Proceedings of the SIGCHI Conference on Human Factors in Computing Systems. ACM, pp. 2399–2402.
- Earle, P., Guy, M., Buckmaster, R., Ostrum, C., Horvath, S. and Vaughan, A., 2010. OMG earthquake! Can Twitter improve earthquake response? *Seismological Res. Lett.*, 81(2): 246–251. DOI: <http://dx.doi.org/10.1002/2015GL063709>
- Hirth, M., Hoßfeld, T. and Tran-Gia, P., 2013. Analyzing costs and accuracy of validation mechanisms for crowdsourcing platforms. *Mathematical and Computer Modelling*, 57: 11–12, 2918–2932. DOI: <http://dx.doi.org/10.1016/j.mcm.2012.01.006>
- Ipeirotis, P.G., Provost, F. and Wang, J., 2010. Quality management on amazon mechanical turk. In Proceedings of the ACM SIGKDD workshop on human computation. ACM, pp. 64–67.
- Iren, D. and Bilgen, S., 2014. Cost of quality in crowdsourcing. *Human Computation*, 1(2): 283–314. DOI: <http://dx.doi.org/10.15346/hc.v1i2.14>
- Kittur, A., Chi, E.H. and Suh, B., 2008. Crowdsourcing user studies with Mechanical Turk. In Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI '08). ACM, New York, NY, USA, pp. 453–456. DOI: <http://dx.doi.org/10.1145/1357054.1357127>
- LaLone, N., Tapia, A.H., Case, N.A., MacDonald, E.A.M. and Heavner, M., 2015. Hybrid community participation in crowdsourced early warning systems. In Proceedings of the ISCRAM 2015 Conference.
- Lamos, V., De Bie, T. and Cristianini, N., 2010. Flu Detector—Tracking epidemics on Twitter. In Balcázar, J., Bonchi, F., Gionis, A. and Sebag, M. (Eds.), *Machine Learning and Knowledge Discovery in Databases. Lecture Notes in Computer Science*, Vol. 6323. Springer Berlin Heidelberg, pp. 599–602. DOI: https://doi.org/10.1007/978-3-642-15939-8_42
- Le, J., Edmonds, A., Hester, V. and Biewald, L., 2010. Ensuring quality in crowdsourced search relevance evaluation: The effects of training question distribution. In SIGIR 2010 workshop on crowdsourcing for search evaluation, pp. 21–26.
- Lee, K., Ganti, R., Srivatsa, M. and Mohapatra, P., 2013. Spatio-temporal provenance: Identifying location information from unstructured text. In Pervasive Computing and Communications Workshops (PERCOM Workshops). 2013 IEEE International Conference. IEEE, pp. 499–504.
- MacDonald, E.A., Case, N.A., Clayton, J.H., Hall, M.K., Heavner, M., Lalone, N., Patel, K.G. and Tapia, A.H., 2015. Aurorasaurus: A citizen science platform for viewing and reporting the Aurora. *Space Weather*, 13: 548–559. DOI: <https://doi.org/10.1002/2015SW001214>
- Motoyama, M., Meeder, B., Levchenko, K., Voelker, G.M. and Savage, S., 2010. Measuring online service availability using twitter. WOSN'10, pp. 13–13.
- Priedhorsky, R., MacDonald, E. and Cao, Y., 2012. First solar maximum with social media: Can space weather forecasting be improved? In AGU Fall Meeting Abstracts, Vol. 1, P. 2324.
- Sheng, K., Gu, Z., Mao, X., Tian, X., Gan, X. and Wang, X., 2014. Answer inference for crowdsourcing based scoring. In Global Communications Conference (GLOBECOM), 2014 IEEE. IEEE, pp. 2733–2738. DOI: <https://doi.org/10.1109/glocom.2014.7037221>
- Tapia, A.H., Lalone, N., MacDonald, E.A., Hall, M., Case, N.A. and Heavner, M., 2014. AURORASURUS: Citizen science, early warning systems and space weather. In Second AAAI Conference on Human Computation and Crowdsourcing.
- Truong, B., Caragea, C., Squicciarini, A. and Tapia, A.H., 2014. Identifying valuable information from Twitter during natural disasters. Proceedings of the American Society for Information Science and Technology, 51(1): 1–4. DOI: <https://doi.org/10.1002/meet.2014.14505101162>
- Vieweg, S., Hughes, A.L., Starbird, K. and Palen, L., 2010. Microblogging during two natural hazards events: What Twitter may contribute to situational awareness. In Proceedings of the SIGCHI conference on human factors in computing systems. ACM, pp. 1079–1088. DOI: <https://doi.org/10.1145/1753326.1753486>
- Voyer, R., Nygaard, V., Fitzgerald, W. and Copperman, H., 2010. A hybrid model for annotating named entity training corpora. In Proceedings of the fourth linguistic annotation workshop. Association for Computational Linguistics, pp. 243–246.
- Vuurens, J., de Vries, A.P. and Eickhoff, C., 2011. How much spam can you take? An analysis of crowdsourcing results to increase accuracy. In Proc. ACM SIGIR Workshop on Crowdsourcing for Information Retrieval (CIR'11). pp. 21–26.
- Wang, A. 2010. Detecting spam bots in online social networking sites: A machine learning approach. Data and Applications Security and Privacy XXIV. In Foresti, S. and Jajodia, S. (Eds.), *Lecture Notes in Computer Science*, Vol. 6166. Springer Berlin Heidelberg, pp. 335–342. DOI: https://doi.org/10.1007/978-3-642-13739-6_25

How to cite this article: Case, N A, MacDonald, E A, McCloat, S, Lalone, N and Tapia, A H 2016 Determining the Accuracy of Crowdsourced Tweet Verification for Auroral Research. *Citizen Science: Theory and Practice*, 1(2): 13, pp. 1–9, DOI: <http://dx.doi.org/10.5334/cstp.52>

Submitted: 30 December 2015 **Accepted:** 24 May 2016 **Published:** 21 December 2016

Copyright: © 2016 The Author(s). This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International License (CC-BY 4.0), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited. See <http://creativecommons.org/licenses/by/4.0/>.

]u[*Citizen Science: Theory and Practice* is a peer-reviewed open access journal published by Ubiquity Press.

OPEN ACCESS 