

Building automated vandalism detection tools for Wikidata

Amir Sarabadani
Wikimedia Deutschland
Tempelhofer Ufer 23/24
10963 Berlin, Germany
a.s.tafreshi@aut.ac.ir

Aaron Halfaker
Wikimedia Research
149 New Montgomery Street
San Francisco, USA
ahalfaker@wikimedia.org

Dario Taraborelli
Wikimedia Research
149 New Montgomery Street
San Francisco, USA
dtaraborelli@wikimedia.org

ABSTRACT

Wikidata, like Wikipedia, is a knowledge base that anyone can edit. This open collaboration model is powerful in that it reduces barriers to participation and allows a large number of people to contribute. However, it exposes the knowledge base to the risk of vandalism and low-quality contributions. In this work, we build off of past work detecting vandalism in Wikipedia to detect vandalism in Wikidata. This work is novel in that identifying damaging changes in a structured knowledge-base requires substantially different feature engineering work than in a text-based wiki. We also discuss the utility of these classifiers for reducing the overall workload of recent changes vandalism patrollers in Wikidata. We describe a machine classification strategy that is able to catch 89% of vandalism while reducing patrollers' workload by 98% by drawing lightly from contextual features of an edit and heavily from the characteristics of the user making the edit.

CCS Concepts

•Information systems → Data analytics; *Online analytical processing*;

Keywords

Wikidata; vandalism; knowledge bases; quality control

1. INTRODUCTION

Wikidata (www.wikidata.org) is a free knowledge base that everyone can edit: it is a collaborative project aiming to produce a high quality, language-independent, open-licensed, structured knowledge base. Like Wikipedia, the project is open to contributions from anyone willing to contribute productively, but also to potentially damaging/disruptive contributions. In order to combat such intentional damage, volunteer patrollers work to review changes to the database after they are saved. At a rate of about 80,000 human edits and 200,000 automated edits per day (as of February 2015),

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.

though, the task of reviewing every single edit would be daunting even for a very large pool of patrollers. Recently, substantial concerns have been raised about the quality and accuracy of Wikidata's statements[11], and therefore, the long-term viability of the project: these concerns call for the design of scalable quality control processes.

Similar concerns about quality control have been raised about Wikipedia in the past[7]. Studies of Wikipedia's quality have shown that, even at large scale and with open permissions, a high-quality information resource can be maintained[7, 17]. One of the key technologies that let Wikipedia maintain quality efficiently at scale is machine classification for detecting vandalism edits. These technologies allow the massive feed of daily changes to be filtered down to a small percentage that is most likely to actually be vandalism, substantially reducing the workload of patrollers[5, 6]. These semi-automated support systems also substantially reduce the amount of time that an article in Wikipedia remains in a vandalized state[5]. The study of vandalism detection in Wikipedia has seen substantial development as a field in the scholarly literature, to great benefit of the project [19, 8, 1, 2].

In this study, we extend and adapt methods from the Wikipedia vandalism detection literature to Wikidata's structured knowledge base. In order to do so, we develop novel techniques for extracting signal from the types of changes that editors make to Wikidata's textititems. But unlike this past literature, we focus our evaluation on the key concerns of Wikidata patrollers who are tasked with reviewing incoming edits for vandalism: reducing their workload. We show that our machine classifier can be used to reduce the amount of edits that need review by up to 98% while still maintaining a recall of 89% using an off-the-shelf implementation of a Random Forest classifier[4]¹.

1.1 Wikidata in a nutshell

Wikidata consists of mainly two types of entities: *items* and *properties*. *Items* represent define-able *things*. Since Wikidata is intended to operate in a language-independent way, each *item* is uniquely identified by a number prefixed with the letter "Q". *Properties* describe a data value of a statement that can be predicated of an item. Like items, properties are uniquely identified by a number prefixed with the letter "P".

Each *item* in Wikidata consists of five sections.

Labels a name for the item (unique per language)

¹<http://scikit-learn.org/stable/modules/generated/sklearn.ensemble.Ran>

Descriptions a short description of the item (unique per language)

Aliases alternative names that could be used as a label for the item (multiple aliases can be specified per language)

Statements *property* and data value pairs such as country of citizenship, gender, nationality, image, etc. Statements can also include *qualifiers* (which include sub-statements like the date of a census for a population count) and sources (like Wikipedia, Wikidata demands reliable sources for its data)

Site links links to Wikipedia and other Wikimedia projects (such as Wikisource²) that reference the item.

For example, the item representing the city of San Francisco (Q62) contains the following statement: (P190, Q90). P190 is a property described in English as “sister city” and Q90 is an item for the city of “Paris, France”. This statement represents the fact that San Francisco (Q62) has a sister city (P190) named Paris (Q90). Using so-called SPO triplets, standing for subject-predicate-object, as a mean to store knowledge is a common practice in knowledge bases such as Freebase.

2. RELATED WORK

The subject of quality in open production has been extensively studied in the open text editing contexts like Wikipedia, but comparatively little study has been done in open structured data editing contexts. In this section, we’ll provide an overview of some of the most relevant work exploring quality in open contexts like Wikipedia and Wikidata.

Wikidata and Wikipedia operate in a common context: they are supported by the Wikimedia Foundation³; virtually all of Wikidata users are also editor in Wikipedia and/or other Wikimedia wikis; Wikidata, like Wikipedia, is powered by MediaWiki software, but Wikidata uses the “wikibase” extension⁴ to manage structured data. Thus, damage detection in Wikipedia is closely related to vandalism detection in Wikidata projects. However, there are also open contribution structured data projects where quality and vandalism detection have been a focus of scholarly inquiry.

2.1 Quality in Wikipedia

Quality in Wikipedia has been studied so extensively that we can’t give a fair overview of all related work, so here, we provide a limited overview of the work that is related to quality prediction and editing dynamics.

Stvilia et al. built the first automated quality prediction models for Wikipedia that was able to distinguish between Featured (highest quality classification) and non-Featured articles[16]. Warncke-Wang et al. extended this work by showing that the features used in prediction could be limited *actionable* characteristics of articles in Wikipedia and maintain a high level of fitness[20] and used these predictions in task routing.

Kittur et al. explored the process by which articles improve most efficiently and found that articles with a small

group of highly active editors and a large group of less active editors were more likely to increase in quality than articles whose editors contributed more evenly[10]. They argued that this is due to the lower coordination cost when few people are primarily engaged in the construction of an article. Arazy et al. challenged the conclusions of Kittur et al. by showing a strong correlation between diversity of experience (global inequality) between editors who are active and positive changes in article quality[3]. The visibility of articles in in Wikipedia seems to be critical their development. Schneider et al., showed that hiding newly created articles from Wikipedia readers in a drafting space substantially reduced the overall productivity of editors in Wikipedia[14].

Detecting vandalism in Wikipedia using machine learning classifiers has been an active genre since 2008[15]. There are generally two types of damage detection problems discussed in the literature: *realtime* and *post-hoc*. The *realtime* framing of damage detection imagines the classifier supporting patrollers by helping them find vandalism shortly after it has happened. The *post-hoc* framing of damage detection imagines the classifier being used long after an edit has been saved (and potentially reverted by patrollers). Since the post-hoc framing allows the model to take advantage of what happens to a contribution after it is saved (e.g. that it was reverted), these classifiers are able to attain a much higher level of fitness than *realtime* classifiers that must make judgements before a human has responded to an edit. However, the utility of a post-hoc classifier is at best hypothetical while *realtime* classifiers have become a critical infrastructure for Wikipedia patrollers[6, 5]. Geiger et al. discussed the “distributed cognition” system that formed through the integration of counter-vandalism tools that use machine classification and social practices around quality control[6]. Geiger et al. showed in a follow-up work that, when systems that use automated vandalism detection go offline, vandalism is not reverted as quickly which results in twice as many views of vandalized articles[5].

Substantial effort has been put into developing high signal features for vandalism classifiers. Adler et. al was able to show substantial gains in model fitness when including user-reputation metrics as features[1]. Among several other metrics, they primarily evaluated the fitness of their classifier using the area under the receiver operating characteristic curve (ROC-AUC) and were able to attain 93.4%. Other researchers explored the use of stylometric features and were able to boost attain an ROC-AUC of 92.9%[19, 8]. West et al. explored spacio-temporal features and built and evaluated models predicting vandalism over only anonymous (not logged-in) user edits because those are the editors from which most vandalism (“offending edits” in West et al.’s terms) originate[21]. Adler et al. continues this work by comparing all of these feature extraction strategies/models and combines them to attain an ROC-AUC of 96.3%[2]. They continue to call for a focus on the area under the precision-recall curve (PR-AUC) instead of the ROC-AUC since it affords more discriminatory power between the overall fitness of models in the context of a low prevalence prediction problem (few positive examples – as is the case with vandalism in Wikipedia).

2.2 Quality in Wikidata

Like Wikipedia, Wikidata is based on the MediaWiki software which provides several means for tracking and review-

²<https://wikisource.org>

³<https://wikimediafoundation.org>

⁴<http://wikiba.se/>

ing changes to content. For example, watchlists⁵ allow editors to be notified about changes made to items and properties that they are interested in. The recentchanges feed⁶ provides an interface for reviewing all changes that have been made to the knowledge base. Wikidata also uses tools related to its own quality demands. Most notably, “Constraint violation reports”⁷ is a dynamic list of possible errors in statements that is generated using predefined rules for properties. For example, if a cat is mentioned as spouse of a human being that’s likely to need review. Other tools such as Kian⁸ also exposes possible errors in Wikidata by comparing data in Wikidata with extracted values from Wikipedia. Despite all of the efforts on quality control in Wikidata, still concerns has been raised regarding Wikidata reliability. For example Kolbe[11] calls into question whether Wikidata will ever be able to verify and source the statements in Wikidata.

Regarding vandalism detection Heindorf et al. have done studied on the demography of vandalism in Wikidata[9] showing interesting dynamics in how who vandalizes Wikidata. For example, most of the vandals in Wikidata already has vandalized at least a Wikipedia language.[9] As far as we can tell, our work is the first published about a vandalism detection classifier for Wikidata.

2.3 Quality in other structured data repositories

There have been several research projects conducted on damage detection in knowledge bases. Most notably, Tan et al. [18] worked on detecting correctness of data added to freebase⁹. They assumed that, if a statement can survive for four weeks, it’s probably a good contribution. They also showed the ratio of correct statements added by a user is not predictive in determining the correctness of future statements, but by defining the area of expertise for each user, it’s possible to make proper predictions. As we showed it doesn’t apply to Wikidata. The difference in the concept of classification in their project and ours can be the reason. They were looking into the correctness of data while our aim is detecting vandalism. Nies et al. [12] have done a research regarding vandalism in OpenStreetMap (OSM). OSM, like Wikidata, is an open structured database but unlike our work, they did not draw from the substantial history of vandalism detection in Wikipedia. Also, they did not use machine learning. Their method is poorly described “rule-based” scoring system and would be difficult to reproduce. In our work, we draw extensively from past work building high fitness vandalism detection models for Wikipedia. We use training and testing strategies that are intended to be straightforward to replicate. We’ve adopted standard metrics from the Wikipedia vandalism detection literature and supplement our own intuitive evaluation metric (filter-rate) that correspond to real effort saved for Wikidata patrollers.

3. METHODS

Using damage detection classifiers in Wikidata usable by reviewers requires two major factors. Firstly, the classifier

should be able to respond and classify edits in a timely manner i.e. within a few seconds otherwise reviewing enormous flow of edits would be unfeasible. In average, two edits are being made by humans in Wikidata every second. Post-hoc features such as time that an edit stays without being reverted, are undesirable in production environments. Secondly, There are two different use cases of classifiers. For auto-reverting by bots and for triaging edits to review by humans. At the first use case, a classifier is expected to have a high level of confidence for instance 99 or 99.9% precision. In these cases, the classifier catches obvious vandalism (e.g. blanking) but a higher recall would be helpful. While, in the second use case, the classifier is expected to have a high recall. Low precision can be tolerated however it should not be too low that in practice it classifies all edits.

In order to build a classifier usable by Wikidata users. We use Wikimedia Labs hosted by Wikimedia Foundation. The environment is called ORES, standing for Objective Revision Evaluation Service, hosts machine learning classifiers for all projects hosted by Wikimedia Foundation including Wikipedia and Wikidata. ORES accepts two methods of scoring edits. Either single edit number or batched are supported.

3.1 Building a corpus

While there has been substantial past work done on building high a quality vandalism corpus for Wikipedia[13], no such work has been done for Wikidata. The work of Heindorf et al.[9] was intended to build such a corpus, their methods (matching edit comments for the use of specific tool) leaves much to be desired as it mislabels a substantial amount of edits. They also assume that patrollers will only use tools to revert vandalism – “rollback” and “restore”. Their qualitative analysis showed that 86% of rollbacked edits and 62% of reverted edits using restore feature were vandalism. If a classifier was trained using this limited corpus was useful for predicting all cases of vandalism regardless of reverting method, then it would evaluate poorly during testing for predicting truly vandalism edits that were mis-labeled in the corpus (i.e., reverted some other way than with “rollback”). In a second scenario, the classifier may only learn how to good at classifying the type of edits that are reverted using rollback. In that case the classifier is substantially less useful in practice, but it would show high scores during evaluation for effectively ignoring the mis-labeled items in the corpus. Thus, training and testing a classifier solely based on rollbacked edits is quietly problematic.

So, rather than rely on this corpus, we applied our own strategy for identifying edits that are likely to be vandalism. First, we randomly sampled 500,000 edits saved by humans (non-bot editors) in the year 2015 in Wikidata. Next, we labeled edits that were *reverted* during an “identity revert event”. Next, we applied several filters to the dataset to examine cross sections of it. First, we split off edits that were performed by users who have attained a high status in Wikidata by receiving advanced rights (here, we include sysop, checkuser, flood, ipblock-exempt, oversight, property-creator, rollbacker, steward, sysop, translationadmin, wikidata-staff). Next, we split out edits that originated from other wikis – known as “client edits” and edits that merged together two Wikidata items. Finally, we were left with a set of regular edits by non-trusted users. Next we reviewed random samples of reverted and non-reverted edits

⁵<https://en.wikipedia.org/wiki/Special:Watchlist>

⁶<https://en.wikipedia.org/wiki/Special:RecentChanges>

⁷<https://www.wikidata.org/wiki/WD:CV>

⁸<https://github.com/Ladsgroup/Kian>

⁹Freebase, Google’s knowledge base, is currently being shut down in favor of Wikidata.

Table 1: Different types of edits in a 500,000 sample

	edits	reverted
trusted user edit	461176	1188 (0.26%)
merge edit	8241	38 (0.46%)
client edit	10099	109 (1.08%)
non-trusted regular edit	22460	622 (2.77%)

in a few key subsets to get a sense for which of these filters could be applied when identifying vandalism.

We can safely exclude client edits since, if they are vandalism, they are vandalism to the originating wiki. Edits by trusted users are reverted at an extremely low rate, but it's still worth reviewing them. Merge edits are also reverted at a low rate, but it is still worth reviewing them. Finally, non-trusted regular edits are reverted at a high rate of 2.77%, which is more in line with the rates seen for all edits in English Wikipedia[13]. To make sure that *reverts* catch most of the vandalism, we review both the reverted and non-reverted regular edits by non-trusted users.

These analyses suggest that reverted edits by non-trusted users are highly likely to be intentional vandalism (68%) or at least damaging (92%) and that non-reverted edits by users in this group are unlikely to be vandalism (1%) or damaging (4%). Further, it appears that reverted merge edits and reverted edits by trusted users are very unlikely to be vandalism (0% observed) – though many merges are good-faith mistakes that violate some Wikidata policy. Based on this analysis, we built a corpus of edits based on this 500,000 sample and labeled reverted regular edits by non-trusted users as *True* (vandalism) and all other edits as *False* (not vandalism). From this 500,000 set, we randomly split 400,000 edits for training and hyper-parameter optimization and 100,000 edits for testing. All test statistics are drawn from this 100,000 test set.

Comparing our work to that of Heindorf et al.[9] we found that only 63% (439) edits we identified as vandalism were reverted using rollback, 15% (104) were reverted using restore and 22% (155) were reverted using other methods.

3.2 Feature engineering

Before starting to build the damage detection classifier we launched a community consultation asking Wikidata users to provide common patterns in vandalism edits with examples. Subsequently, we received around thirty patterns and examples. There were two reasons for this success, Firstly one of the researchers is a prominent contributor in Wikidata. Secondly, this project is directly supported by product manager of Wikidata. Their feedback was helpful building the initial model which was launched on October 29, 2015. Then another community consultation was launched for reporting possible mistakes of the initial model and more than 20 cases of false positives and false negatives were reported which helped authors to improve the damage detection mostly by adding proper features. Also, in order to have accurately labeled data a campaign called “edit quality” also is launched by authors which asks community members to hand code 4283 edits and as the time of this writing this campaign is half-way through and its data is used in this research to examine accuracy of models and automated labels of edits that are being used in training models. Also

this classifier is accessible for everyone¹⁰ which lets users and experts comment on the algorithms and methods used.

4. FEATURES

4.1 General metrics

- Number of added/removed/changed/current site links
- Number of added/removed/changed/current labels
- Number of added/removed/changed/current descriptions
- Number of added/removed/changed/current statements
- Number of added/removed/current aliases
- Number of added/removed/current badges
- Number of added/removed/current qualifiers
- Number of added/removed/current references
- Number of changed identifiers¹¹

4.2 Wikidata vandalism inspired features

Proportion of Q-ids added It’s a common type of vandalism to add Q-id of items to contextual parts of items.

If English label has changed Changing English label is a common type of vandalism.

Proportion of language name added Wikidata user interface shows names of languages. Consequently, adding names of language such as “English” is a common pattern in vandalism.

Proportion of external links added Spamming also happens in Wikidata and this feature was helpful in catching spamming.

Is gender changed Changing gender is a common type of vandalism

Is country of citizenship changed Country of citizenship gets vandalized often.

Is member of sports team changed Changing statements regarding teams that a sportsperson has played happens to be a target for vandalism.

Is date of birth changed Date of birth is highly vandalized

Is image changed Changing image of an item to an unrelated image is a common ground for vandalism.

Is image of signature changed

¹⁰<https://github.com/wiki-ai/wb-vandalism>

¹¹Identifiers are a type of statements which holds data connecting the item to external data sources. Such as VIAF id

Table 2: Edits sampled for human review

	good	goodfaith damaging	vandalism
reverted merge edits	17	21	0
reverted trusted user edits	93	7	0
reverted nontrusted regular edits	8	24	68
nonreverted nontrusted regular edits	94	3	1

Is Wikimedia Commons identifier changed Wikidata keeps a statement as identifier to Wikimedia Commons. This property is being heavily used by Wikipedia projects and changing it directly affects Wikipedia articles. Thus, changing the statement is a common pattern in vandalism.

Is official website has changed Official website once added barely changes. Reportedly, most of changes to official websites are vandalism.

Is this item is about a human Humans are a common target of vandalism in Wikidata.

Is this item is about a living human Living humans are biggest target of vandalism in Wikidata. Wikipedia and likewise Wikidata have strict policies regarding pages related to living human being.

4.3 Wikidata’s common non-vandalism changes

Is it a client edit When a user moves a page in Wikipedia (a client of Wikidata) or deletes the page, an edit is made in wikidata to update the central repository. This type of vandalism is related to the client itself and using these features they are automatically excluded.

Is it a merge Merges, which is not enabled for new users, tends to change the item drastically hence they tend to be flagged as vandalism but merges are mostly correct and if they are wrong, they are not vandalism due to merges being disabled for new users.

Is it revert, rollback, or restore These edits are trying to undo vandalism and most of the time they are correct.

Is it creating a new item Creating new item tends to have high probability due to adding content.

4.4 Editor characteristics

Is the user is a bot Automated edits using bots is a common practice in Wikidata. 54% of Wikidata edits in January 2016 have been made by bots, consisting 59% bytes changed in the database.

Does the user has an advanced right If the user is a member of “checkuser”, “bureaucrat”, or “oversight”. There are a few users who hold such rights.

Is the user administrator Administrators are advanced users with significant contribution who are trusted by community of editors

Is the user curator If the user is a member of so-called “rollbacker”, “abusefilter”, “autopatrolled”, or “reviewer”. User with significant contribution can posses these rights.

Is the user anonymous Most of vandalism edits in Wikidata are originated from anonymous (unregistered) users.

Age of editor The time between registration of user and timestamp of the edit in seconds scaled using $\log(\text{age} + 1)$. Since old users are more likely not to vandalize Wikidata this feature is highly predictive.

5. EVALUATION

In order to evaluate the effectiveness of our prediction model, we employ three metrics.

- ROC-AUC as has been used historically in the vandalism detection literature[1, 19]
- PR-AUC as Adler et al. call for in their more recent work[2]
- Filter-rate at high recall which measures the proportion of edits that must be review by Wikidata patrollers in order to expect to catch some high percentage of all vandalism.

Our addition of “filter rates” to the discussion of vandalism classifiers is purposeful in that, our intention in designing this classifier is that it will be used by patrollers in Wikidata. As we improve fitness of the model, this filter-rate should increase and therefor the workload for patrollers should decrease. This metric directly measures theoretical changes in patroller workload.

6. RESULTS

In this section, we discuss the fitness of our model against the corpus and the realtime performance of the model via ORES, our live service for Wikidata patrollers.

6.1 Model fitness

All models were tested on the exact same set of 99,222 revisions withheld during hyper-parameter optimization and training. The table show general fitness metrics for the models. As Table 6.1 suggests, we were only able to train marginally useful prediction models when excluding user features. These models attained extremely low PR-AUC values and there was no threshold that could be set on the *True* probability that would allow for 75% of *reverted* edits to be identified to the exclusion of others – resulting in a *zero* filter rate.

Figures 6.1 and 6.1 plot precision-recall curves for the two sub-feature-sets. Figure 6.1 visually confirms the dismal results of the classifier in the case where no user features are included. At the scale of the graph it is difficult to confirm any meaningfully greater than zero precision anywhere on the recall spectrum. 6.1 shows a substantial difference. For the most part, the *all* features and *general and user* features classifiers seem to perform comparably well across the

Table 3: Model fitness for different subsets of features

features	ROC-AUC	PR-AUC	filter-rate
general	0.777	0.01	0.936 at 0.62 recall
general, context	0.803	0.013	0.937 at 0.67 recall
general, type, context	0.813	0.014	0.940 at 0.68 recall
general, user	0.927	0.387	0.985 at 0.86 reca
all	0.941	0.403	0.982 at 0.89 reca

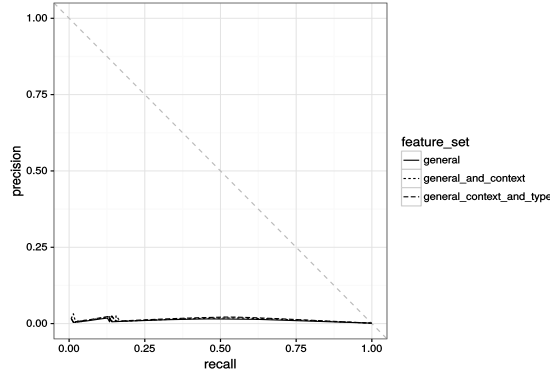


Figure 1: Precision/recall (no user). The precision/recall curve for models lacking user features is plotted.

spectrum of recall. This result implies that the inclusion of *context* and *edit type* features on top of *general* and *user*-based feature resulted in minor if any improvements.

Figures 6.1 and 6.1 show the *filter-rate* (which translates to the amount effort saved for Wikidata patrollers. See Methods.) of the classifiers across the spectrum of recall. Here, we can see that the models that lack *user* features struggle to attain 70% recall at any filter-rate while the models that include *user* features are able to attain very high filter rates up to 89% recall. This implies a theoretical reduction in patrolling workload down to 1.8% of incoming human edits assuming that it's tolerable to let 11% of potential vandalism by being caught by other means.

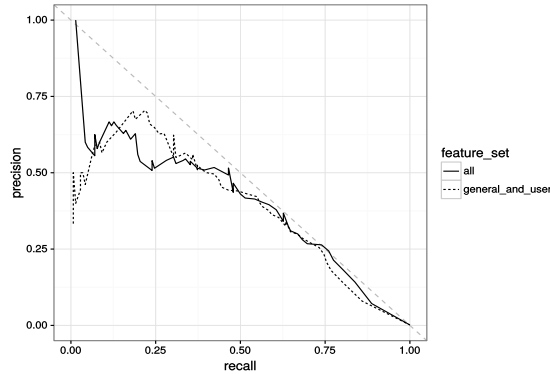


Figure 2: Precision/recall (user) The precision/recall curve for models including user features is plotted.

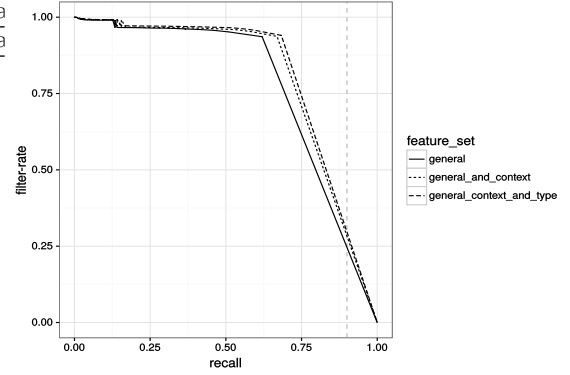


Figure 3: Filter-rate/recall (no user) The filter-rate/recall curve for models lacking user features is plotted.

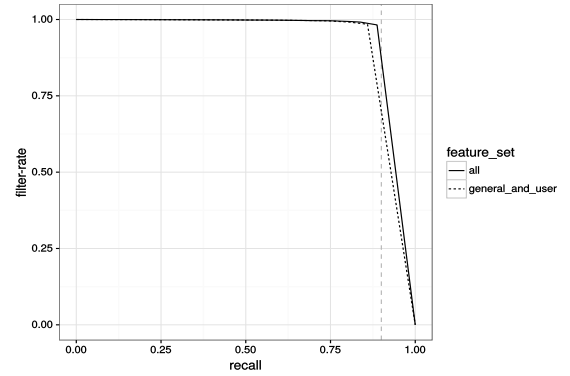


Figure 4: Filter-rate/recall (user) The filter-rate/recall curve for models including user features is plotted.

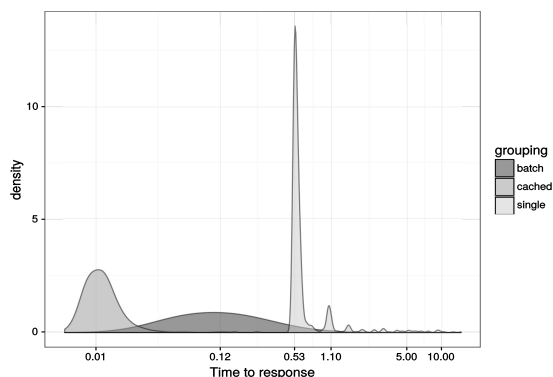


Figure 5: ORES response timing (Wikidata-reverted). The density of response timing per revision score requested is plotted for requests to ORES’ wikidatawikireverted model.

6.2 Realtime prediction speed

Requests for a single prediction will generally respond in ~ 0.5 seconds with rare cases taking up to 2-10 seconds. If the score has already been generated and cached, the system will generally respond in 0.01 seconds. This is a common use-case since we run a pre-caching service that caches scores for edits as soon as they are saved. Through ORES, we also provide the ability to request scores in batches which allows us to gather basic data for feature extraction in batch as well. When requesting predictions in 50 revision batches, we found that the service will respond in about 0.12 seconds per revision in the batch. As ORES response timing (Wikidata-reverted) suggests, this timing varies quite widely which is likely due to the rare individual edits that take a long time to score and hold back the whole batch from finishing.

7. LIMITATIONS

A major limitation of our model building and analysis exercise is our means of constructing our corpus. In our analysis of which edits and reverts are likely to represent vandalism in Wikidata, we used characteristics of the edit (e.g. is it a client edit? and is it a merge edit?) and the editor (e.g. is the user in a trusted group?) to identify vandalism. These characteristics of an edit are also included as features in our prediction model. If we had applied these filters on our vandalism corpus inappropriately, we could have simply trained our classifier to match the rules we had put in place. However, we know that this is not the case generally due to reports from Wikidata users who are using our live classification service. Reports from our users generally suggest that the classifier is generally effective at flagging edits that are vandalism. We also ran a follow-up qualitative analysis to help check whether our estimates of the filter-rate afforded by the “all features” model worked out in practice.

We randomly sampled 10,000 human edits and generated vandalism prediction scores for them. We then manually labeled (1) the highest scored edits in the dataset (100 edits at more than 93% prediction), (2) all reverted edits in the dataset, and (3) and random samples of 100 edits for each 10% strata of prediction weight (e.g. 30-40%, 40-50%, etc.) We found only 17 (0.17%) vandalism edits in the 10,000 set

and all these vandalism observations scored 93% or more. Only 100 of the 10,000 edits were scored 93% or above and by reviewing this 1% fraction of edits it’s possible to catch all damaging edits. So, in this sample set, with our classifier, it looks like we were able to attain a 99% filter-rate with 100% recall by setting the threshold at 93%. This result looks substantially better on paper than evaluation against the test set and we think that is due to the inclusion of careful human annotation. It seems likely that more of the good edits that were mistakenly labeled as vandalism in our corpus construction are probably also likely to show up as false negatives our test set – pushing down our apparent filter-rate and recall. While this analysis is not as robust and easy to replicate the formal analysis we describe above, we feel that it helps show that our classifier may be more useful than appears.

These concerns and limitations motivate the need for a PAN-like dataset for Wikidata that actually uses human judgement to identify vandalism edits rather than heuristics – at least for testing the classifiers. Otherwise, the true filter-rates and therefor reduction in workload for patrollers can only really be discovered in the context of the actual work performed by patrollers. The extremely low prominence of vandalism edits in Wikidata would mean that we would need extremely large numbers of labeled observations to attain a representative set of vandalism edits – probably on the order of 100,000 - 1,000,000. Further, asking people who are not familiar with Wikidata and the many language fields that are used to describe items and properties, would be difficult. Unlike vandalism in Wikipedia, labeling vandalism in Wikidata would require reviewers who are both familiar with the shape and form of structured data in Wikidata and able to evaluate contributions in many languages. We leave such a high quality test set to future work.

8. CONCLUSION

In this paper we described a straightforward method for classifying Wikidata edits as vandalism in realtime using a machine learned classification model. We have also shown that, using this model and our prediction service, it’s possible to reduce human labor involved in patrolling edits to Wikidata by nearly two orders of magnitude (98%). At the time of writing, there are already several tools that have adopted our service and are using the prediction model to patrol edits. Our analysis and a substantial part of our feature set are informed by the real-world experiences of patrollers who are using this classifier to do their work.

We hope that future work with focus on two key areas: (1) The development of a high quality vandalism test dataset for Wikidata and (2) The development of new features for Wikidata that draw from sources of signal other than a user’s status as “untrusted”. This high quality vandalism dataset would provide a mean to effectively compare the capabilities of prediction models without the limitations we note and the qualitative intuitions we have gained “in the wild” that are difficult to replicate. The development of high signal features beyond a users status are important from a perception point of view. Right now, our classification model is weighted strongly against edits by anonymous and new contributors to Wikidata regardless of the quality of their work. While this may be an effective way to reduce patrollers’ workload, it’s likely not fair or reasonable to these users that their edits are so carefully scrutinized. By increas-

ing the fitness of this model and adding new, strong sources of signal, a classifier could help direct patrollers attention away from good new/anonymous contributors and towards vandalism – both reducing their workload and potentially making Wikidata a more welcoming place for newcomers.

9. ACKNOWLEDGMENTS

We would like to thank Lydia Pintscher and Abraham Taherivand from Wikimedia Deutschland. Yuvaraj Pandian from Wikimedia Foundation for operational support. Adam Wight, Helder Lima, Arthur Tilley and Gediz Aksit for their help. We also want to thanks community of Wikidata editors for providing feedback and reporting mistakes.

10. REFERENCES

- [1] B. Adler, L. de Alfaro, and I. Pye. Detecting wikipedia vandalism using wikitrust. *Notebook papers of CLEF*, 1:22–23, 2010.
- [2] B. T. Adler, L. De Alfaro, S. M. Mola-Velasco, P. Rosso, and A. G. West. Wikipedia vandalism detection: Combining natural language, metadata, and reputation features. In *Computational linguistics and intelligent text processing*, pages 277–288. Springer, 2011.
- [3] O. Arazy and O. Nov. Determinants of wikipedia quality: the roles of global and local contribution inequality. In *Proceedings of the 2010 ACM conference on Computer supported cooperative work*, pages 233–236. ACM, 2010.
- [4] L. Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- [5] R. S. Geiger and A. Halfaker. When the levee breaks: without bots, what happens to wikipedia’s quality control processes? In *Proceedings of the 9th International Symposium on Open Collaboration*, page 6. ACM, 2013.
- [6] R. S. Geiger and D. Ribes. The work of sustaining order in wikipedia: the banning of a vandal. In *Proceedings of the 2010 ACM conference on Computer supported cooperative work*, pages 117–126. ACM, 2010.
- [7] J. Giles. Internet encyclopaedias go head to head. *Nature*, 438(7070):900–901, 2005.
- [8] M. Harpalani, M. Hart, S. Singh, R. Johnson, and Y. Choi. Language of vandalism: Improving wikipedia vandalism detection via stylometric analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies: short papers-Volume 2*, pages 83–88. Association for Computational Linguistics, 2011.
- [9] S. Heindorf, M. Potthast, B. Stein, and G. Engels. Towards vandalism detection in knowledge bases: Corpus construction and analysis. In *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 831–834. ACM, 2015.
- [10] A. Kittur and R. E. Kraut. Harnessing the wisdom of crowds in wikipedia: quality through coordination. In *Proceedings of the 2008 ACM conference on Computer supported cooperative work*, pages 37–46. ACM, 2008.
- [11] A. Kolbe. Whither wikidata? <https://en.wikipedia.org/wiki/Wikipedia:Wikipedia.Signpost/2015-12-02/Op-ed>, 2015. [Online; accessed 10-February-2016].
- [12] P. Neis, M. Goetz, and A. Zipf. Towards automatic vandalism detection in openstreetmap. *ISPRS International Journal of Geo-Information*, 1(3):315–332, 2012.
- [13] M. Potthast. Crowdsourcing a wikipedia vandalism corpus. In *Proceedings of the 33rd international ACM SIGIR conference on Research and development in information retrieval*, pages 789–790. ACM, 2010.
- [14] J. Schneider, B. S. Gelley, and A. Halfaker. Accept, decline, postpone: How newcomer productivity is reduced in english wikipedia by pre-publication review. In *Proceedings of the international symposium on open collaboration*, page 26. ACM, 2014.
- [15] K. Smets, B. Goethals, and B. Verdonk. Automatic vandalism detection in wikipedia: Towards a machine learning approach. In *AAAI workshop on Wikipedia and artificial intelligence: An Evolving Synergy*, pages 43–48, 2008.
- [16] B. Stvilia, M. B. Twidale, L. C. Smith, and L. Gasser. Assessing information quality of a community-based encyclopedia. In *IQ*, 2005.
- [17] B. Stvilia, M. B. Twidale, L. C. Smith, and L. Gasser. Information quality work organization in wikipedia. *Journal of the American society for information science and technology*, 59(6):983–1001, 2008.
- [18] C. H. Tan, E. Agichtein, P. Ipeirotis, and E. Gabrilovich. Trust, but verify: Predicting contribution quality for knowledge base construction and curation. In *Proceedings of the 7th ACM international conference on Web search and data mining*, pages 553–562. ACM, 2014.
- [19] W. Y. Wang and K. R. McKeown. Got you!: automatic vandalism detection in wikipedia with web-based shallow syntactic-semantic modeling. In *Proceedings of the 23rd International Conference on Computational Linguistics*, pages 1146–1154. Association for Computational Linguistics, 2010.
- [20] M. Warncke-Wang, D. Cosley, and J. Riedl. Tell me more: An actionable quality model for wikipedia. In *Proceedings of the 9th International Symposium on Open Collaboration*, page 8. ACM, 2013.
- [21] A. G. West, S. Kannan, and I. Lee. Detecting wikipedia vandalism via spatio-temporal analysis of revision metadata? In *Proceedings of the Third European Workshop on System Security*, pages 22–28. ACM, 2010.