

Automatic Intent-based Classification of Citizen-to-Government Tweets

José L. Lavado*, Iván Cantador**, María E. Cortés-Cediel***, Miriam Fernández****

*Escuela Politécnica Superior, Universidad Autónoma de Madrid, Spain, jose.lavado@estudiante.uam.es

**Escuela Politécnica Superior, Universidad Autónoma de Madrid, Spain, ivan.cantador@uam.es

***Facultad de Ciencias Políticas y Sociología, Universidad Complutense de Madrid, Spain, mcorte04@ucm.es

****Knowledge Media Institute, The Open University, United Kingdom, miriam.fernandez@open.ac.uk

Abstract: Social networking technologies offer opportunities for governments to engage with citizens. However, the inability to filter relevant citizens' messages out of the vast amount of available social media content lessens their impact. In this paper, we propose a set of categories encapsulating the different citizens' intents when directing messages to public institutions, e.g., complaining, making requests, and proposing solutions to existing problems. We present a novel artificial intelligence approach, built upon natural language processing and machine learning algorithms, that enables the categorisation of citizens' messages into such intents automatically, and at scale. Through an empirical evaluation on a Twitter dataset, we show the effectiveness of our approach in terms of categorisation performance. We also discuss the value of the presented solution, as a novel tool for governments to achieve a more effective and informed communication with citizens.

Keywords: e-participation, social networks, natural language processing, machine learning

Acknowledgement: This work was conducted with financial support from the Spanish Ministry of Science and Innovation (PID2019-108965GB-I00) and the Centre of Andalusian Studies (PR137/19). José L. Lavado is partially supported by the UAM-ADIC Chair for Data Science and Machine Learning.

1. Introduction

Nowadays, the implementation of e-government models in the field of public management is mainly oriented towards the production, custody and management of large-scale data (Charalabidis et al., 2019), which are produced through technologies such as social media, IoT devices, cloud computing, and blockchain, to name a few.

Among the existing sources of information, social networks represent a prominent bidirectional communication channel between citizens and government. In them, citizens are not only content

consumers who receive the government announcements, to which they react and freely respond according to personal ideology, interests and needs, but also are content providers who generate a wide range of messages targeted to government and political stakeholders.

The amount of social media content daily generated by citizens is huge and diverse, and its processing by human actors may result too costly and overwhelming. Hence, there is increasing interest and need to use computer-assisted solutions capable of automatically gathering, processing and analysing the underlying information in the citizens' messages (a.k.a. posts) on social networks. The research literature reports extensive work on mining citizen generated content. The majority of such work has focused on i) analysing social phenomena produced through the online network structures -e.g., information spreading, fake news, and opinion polarity-, and mainly originated by particular events -e.g., natural disasters, elections, and trending news-, and (ii) extracting the most popular topics addressed by citizens' posts in social networks, as well as the general dynamics (i.e., temporal evolution) and opinions on such topics.

In this paper, we are interested in the latter case. However, differently to previous work, we go beyond the extraction of topics by attempting to automatically classify citizens' posts according to their intents or purposes. That is, we aim to determine whether a post targeting government actors expresses a question, complaint or request, presents a proposal or idea to address a particular problem, spreads an announcement or news item of interest for the general public, or reflects a personal fact or opinion. We believe this automatic classification can be very valuable for government managers and politicians in several ways. First, it would represent a mechanism to identify relevant citizen posts for which responses should be given. This may help increasing the citizens' satisfaction and engagement, who would perceive attention to their questions and requests. Hence, it may promote the openness of the public administration, and ultimately may increase the citizens' trust on a government that responds to public demands. Second, the proposed classification would allow extracting indicators about opinion on how public resources are being managed. These indicators could be used by government managers to identify problems for which new actions and public policies are needed. This may lead to increase the effectiveness and efficiency on both the management of public resources and the provision of public services, which ultimately would generate public value. Finally, the intent-oriented classification would isolate measures on current leadership perception. Taking these measures into account, political parties and leaders could make timely decisions reacting to major opinions, complaints and proposals on problematic and controversial issues.

2. Related Work

Our goal is to categorise the messages that citizens explicitly direct to public institutions in social networks. Hence, we have discarded from our literature review those papers that analyse social media content generated around particular events (e.g., elections and political uprisings), where messages are not necessarily targeted to public institutions. Among the analysed papers, we have identified two main research lines: i) works conducting topic (or thematic) analysis of the different messages that citizens direct to their public institutions, and (ii) works attempting to understand the

opinions and sentiment behind those messages. Some works use a combination of topic and sentiment analysis.

Works that focus on the analysis of sentiment or the analysis of both topics and sentiment, can be divided into (i) those that analyse messages posted by governments and politicians (Siyam et al., 2020; Zavattaro et al., 2015), and (ii) those that analyse messages posted by citizens (Guhathakurta et al., 2019; Lorenzo-Dus & Cristofaro, 2016; Masdeval & Veloso, 2015; Nummi, 2019; Picazo-Vela et al., 2012). The first set of works shows how governments and politicians that adopt a positive tone --and undertake activities like responding directly to citizens on Twitter, sharing photos, and using exclamation points-- are more likely to encourage citizen participation (Zavattaro et al., 2015). Also, they show that videos and images have a high positive impact on engagement, and tweets posted on weekdays obtain higher engagement than those posted on weekends (Siyam et al., 2020).

The works that analyse messages posted by citizens show: (i) how users present high levels of emotionality as well as a high participation rate (Lorenzo-Dus & Cristofaro, 2016) and (ii) how, within a urban context, citizens' sentiment can be used as an indicator of perceived neighbourhood quality (Guhathakurta et al., 2019), as well as an indicator to estimate urgency of urban issues, such as overflowing trash bins and broken footpaths, among others (Masdeval & Veloso, 2015). Despite the usefulness of social media data analysis, some works (Nummi, 2019; Picazo-Vela et al., 2012) argue that the integration of this knowledge in planning and decision-making has not been completely successful, and that a good implementation strategy is necessary to realise their full benefits. In this line, Garg and colleagues (2017) proposed an automatic approach to determine which of the posts that citizens direct to institutions are "actionable", i.e., can be acted upon by the government.

3. Classification of Tweets Based on Their Intent

Within the machine learning field, text classification (a.k.a. text categorization) refers to the task of automatically assigning a natural language text with one of a given set of classes (labels or categories). The classes are usually discrete values related to topics, but can also represent domain-dependent meanings, such as "spam" and "non-spam" emails, "real" and "fake" news articles, and "positive" and "negative" textual reviews. Besides, a classification problem may be binary -with two classes- or multi-class -with more than two classes.

To address this task, supervised learning assumes that a set of training data (i.e., the training set) has been provided, consisting of a set of instances (input texts) that have been labelled by hand with their correct class. On weighted feature vector representations of the training instances, a learning procedure aims to extract feature patterns and relations that allow characterizing and distinguishing instances from each class. The procedure then generates a model that attempts to meet two sometimes conflicting objectives: classifying as well as possible on the training data, and generalising as well as possible to new (test) data.

In this context, the selection and extraction of features represents a key stage for the effectiveness of the final classification process. When dealing with text documents, a typical choice is to identify

features with words, in the so-called bag of words model, and to assign each word with a weight equals to its TF-IDF (term frequency-inverse document frequency) value.

3.1 Proposed Intent-based Classes

Online social network participation can be a form of political participation that should be conceptualized, identified and measured. From a revision of the literature, he considers several forms of political participation: (i) posting (sharing) links to political stories or articles for others to read, (ii) posting own thoughts or comments on political or social issues, (iii) encouraging to take action on a political or social issue and, (iv) reposting content related to political or social issues that was originally posted by someone else. Motivated by such categorization, in this paper, we focus on identifying the intent that citizens have when posting messages to their institutions. In addition, we rely on a data-driven inspection to define categories of intent. The final ten intent-based categories extracted after this process include:

- Complaint. The intent is to state something that is unsatisfactory or unacceptable (e.g., "@MADRID after 1 week of calling, the city is yet not clean and the rats are taking over!! <http://t.co/IiIDuaPFG9>").
- Announcement. The intent is to make a public statement about a fact, occurrence or event (e.g., "The date, place and schedule of the Festival activities in La Latina have already been confirmed [@madrid @madridiario](http://t.co/U0tRwKAC)").
- News item. The intent is to objectively inform about current events. Authors of these posts are generally media news organisations and journalists (e.g., "#oladecolor #aemet @Madrid has suffered its warmest night within the latest 100 years <http://t.co/ZSjeqK6m>").
- Personal fact. The intent is to publicise self issues and experiences (e.g., "I also support the candidature from @Madrid2020ES @MADRID #aporella").
- Personal opinion. The intent is to express subjective opinions about the city, its events, activities, etc. (e.g., "The activity of #emprendeenmadrid is amazing. Congratulations @MADRID and greetings from an entrepreneur").
- Request. The intent is to explicitly ask for something specific (e.g., "Very nice but impossible to ride a bike at normal speed #MadridRio. Please @MADRID create a bike lane with cyclist priority").
- Notification. The intent is to report or give notice of urban, citizenship- or government-related issues, so that the Madrid City Council can quickly act on them and help other citizens (e.g., "@MADRID can you fix this gap in c/ San Bernardino 8-10 before someone gets hurt? <http://lockerz.com/s/117566458>").
- Question. The intent is to explicitly ask for information (e.g., "@MADRID could you please give me the telephone number of the press office of the Madrid city hall").
- Proposal. The intent is to suggest an initiative or project. Proposals indicate broader projects and ideas than the explicit and specific demands of the request category ("There is a collection of used oil in the centre of Alicante. It would be fantastic to have something similar @MADRID").

3.2 Proposed Classification Features

To automatically categorise each tweet into one of the classes categories presented in the previous subsection, it is first transformed into a vector of 37 domain- and language-independent features. From them, 27 are content-based features, including:

- Lexical features: number of characters, number of words, number of exclamation marks, number of question marks, existence of a positive emoticon, existence of a negative emoticon, and existence of a vowel (or "y") consecutively repeated 3 or more times in a word. The latter is assumed to be a signal of emphasis.
- Grammatical features (20): number of nouns, number of proper nouns, number of adjectives, number of verbs, number of adverbs, number of personal/possessive pronouns, number of time references (entities), and number of money-related references.

These content-based features were obtained by a computer program that makes use of the Stanford CoreNLP natural language processing toolkit, which, as far of March 2021, allows obtaining the syntactic parsing of sentences in English, Arabic, Chinese, French, German and Spanish. For nouns, adjectives, verbs and adverbs, we also consider the number of them which were positive/negative/neutral, according to a Spanish lexicon of word opinion polarities.

The remainder 10 features were social network-based, including:

- User features (4): number of followers, number of friends (a.k.a. followees), number of posts, and number of active days in Twitter.
- Post features (6): number of hashtags (#), number of user mentions (@), number of hyperlinks, number of multimedia, maximum hashtag length, and existence of an explicit retweet request (i.e., "RT" abbreviation).

We discarded interaction-based features, such as the number of "likes," the number of "comments," and the number of "reposts" (i.e., retweets), since our aim is to automatically categorise tweets after they are generated. Further popularity-based signals could be used in longer term processing/analysis stages. We also discarded fine-grained grammatical features, such as the number and tense of the verbs. For instance, one may expect that first-person verbs would not appear in news items, and thus may represent an informative feature to characterise that class. Similarly, imperative verbs may be much frequent in requests, whereas conditional verbs may be predominant in proposals. We did not consider these features since they depend on the language in which tweets are written. Nonetheless, they could be exploited in a language-specific solution to improve classification accuracy.

4. Experiments

4.1 Dataset

As a case study to test our approach we selected the City Council of Madrid, Spain. Its Twitter account (@Madrid) has more than 700K followers, and receives a high volume of daily posts explicitly directed to it. We aimed to categorise messages posted by citizens and directed to that public institution. To gather these messages, we first collected data for all the user accounts

following @Madrid. The Twitter API allowed us to collect the most recent 3,200 posts for each of these accounts. We then filter those messages explicitly directed to Madrid city council.

To obtain the necessary training data to build and evaluate our classification approach, we needed to categorise a subset of posts manually. For this purpose, we selected a random sample of 666 tweets. These tweets were manually annotated by four experts (each of them annotated a 500 sample), ensuring that each tweet received at least three annotations. All experts received explicit indications of the categories and their meaning before conducting the annotation process. In addition, an hour of debate was allocated for them to reflect on the categories and resolve possible doubts. The annotation process shows an agreement of Fleiss' kappa coefficient equal to 0.98, meaning almost perfect agreement. For conflicting cases, the majority class assigned to a tweet was finally selected.

4.2 Classification Algorithms

To validate the proposed method, we evaluated several machine learning algorithms on the generated dataset. The tested algorithms included:

- K-Nearest Neighbours (KNN)
- Logistic Regression (LR)
- Quadratic Discriminant Analysis (QDA)
- Decision Tree (DT), which was executed alone, and in combination with feature selection (RFECV DT) and tree pruning (AP DT) to avoid learning over-fitting
- Gaussian Process (GP)
- Support Vector Machine (SVM)
- Bagging Ensemble (BE)

4.3 Classification Results

Due to the unbalanced distribution of the instances in the 10 classes, we conducted a series of experiments where we addressed 10 binary (2-class) classification problems. Each of them aimed to distinguish the instances belonging to a particular class from the instances belonging to the other classes.

In addition to computing the accuracy (acc) metric, which measures the percentage of instances (i.e., tweets) correctly classified, we also computed the acc+ and acc- metrics, which correspond to the percentage of correctly classified instances in the minority and majority classes, respectively. As a compromise of both metrics, we considered their geometric mean $g = \sqrt{\text{acc}^+ \cdot \text{acc}^-}$. We computed average metric values from 3 independent executions of each algorithm and parameters configuration, keeping 75% of the tweets for training the machine learning models, and 25% for testing, selected randomly in each execution.

Table 1 shows the best accuracy results achieved by the evaluated algorithms on each intent-oriented classification problem. Note that the classification problems present a large unbalance between the target minority class and the majority class, ranging from $N^+ = 28\%$ of positive instances

for the complaint class to $N+ = 2\%$ for the notification, question, and proposal classes. This makes the classification problems challenging.

Despite this difficulty, by exploiting the proposed domain- and language-independent features and using generic machine learning algorithms, we were able to achieve relatively high accuracy ($acc+$) on identifying complaints, announcements, news items, personal opinions, and requests. The achieved classification performance is relatively high, as can be seen by comparing the $acc+$ and g values against the percentage of positive instances $N+$ in each class. Note that acc values alone are not informative enough, since for each intent, classifying every instance as negative, we would achieve an accuracy equals to $N-$, but we would be wrongly classifying all positive instances.

Table 1: Best accuracy (acc , $acc+$, $acc-$) values for each intent-oriented classification task. Geometric values g show the achieved accuracy balanced between minority ($N+$) and majority ($N-$) classes.

Intent	$N+$	$N-$	acc	$acc+$	$acc-$	g	Algorithm
Complaint	28%	72%	74.6%	66.6%	78.0%	72.0%	QDA/LR
Announcement	26%	74%	83.8%	75.0%	87.0%	81.0%	AP DT
News item	14%	86%	64.6%	61.0%	65.0%	63.0%	QDA
Personal fact	11%	89%	73.4%	44.0%	77.0%	59.0%	QDA
Personal opinion	8%	92%	83.6%	57.0%	86.0%	70.0%	QDA
Request	5%	95%	94.8%	56.0%	97.0%	74.0%	KNN
Notification	2%	98%	97.8%	33.0%	99.0%	57.0%	AP DT
Question	2%	98%	96.8%	33.0%	98.0%	57.0%	RFECV DT
Proposal	2%	98%	91.9%	33.0%	93.0%	56.0%	SVM/LR

5. Conclusions

As citizens are spending more time on online social networks, generating large amounts of content, there is a need for innovative methods and tools to analyse such data. In this paper, we have presented and evaluated a novel AI approach that applies natural language processing and machine learning algorithms to automatically classify citizen-to-government posts published in social networks. Differently to previous works, which have focused on topic- and opinion-based analysis, our approach aims to classify posts based on their underlying intention or purpose, distinguishing between citizens' complaints, requests, proposals and announcements, among others. This classification represents a processing stage prior to the extraction of topics and opinions, and may help filtering and prioritising citizens' messages, and further automatising processes for more efficient and effective decision and policy making.

Despite the positive classification results achieved by our approach, there is still room for improvement. For example, more sophisticated Natural Language Processing techniques, such as

language models and word embeddings, could be used to exploit the semantics of words and word sequences, e.g., "opinion is" and "really think that" could be identified as informative bigram and trigram of the personal opinion class. Furthermore, it could be possible to extend our approach with features from other sources of information, such as the user who creates a post and the users who are mentioned in a post (e.g., by considering their type: particular citizens, neighbourhood associations, organisations, or political actors), and the nature of web resources linked in the posts (e.g., articles of online news media, personal blogs, or multimedia in social networks). From a social inclusion perspective, and taking fairness concerns into account, we plan to investigate possible biases derived from the subset of the population posting these messages, as well as the possible biases that the classification algorithms may have depending on issues such as the users' posting activity and influence (i.e., number of followers), and ideological, political and popularity-based factors of the addressed topics.

References

- Charalabidis, Y., Loukis, E., Alexopoulos, C., Lachana, Z. (2019). The three generations of electronic government: From service provision to open data and to policy analytics. *Proceedings of the 18th IFIP WG 8.5 International Conference on Electronic Government*, pp. 3–17.
- Garg, H., Bansal, C., Kaushal, R., Thanaya, I. (2017). Identifying actionable information from social media for better government-public relationship. *Proceedings of the 10th International Conference on Developments in eSystems Engineering*, pp. 206–211.
- Guhathakurta, S., Zhang, G., Chen, G., Burnette, C., Sepkowitz, I. (2019). Mining social media to measure neighborhood quality in the city of Atlanta. *International Journal of E-Planning Research*, 8(1), 1–18.
- Lorenzo-Dus, N., Di Cristofaro, M. (2016). # living/minimum wage: Influential citizen talk in Twitter. *Discourse, Context & Media*, 13, 40–50.
- Masdeval, C., Veloso, A. (2015). Mining citizen emotions to estimate the urgency of urban issues. *Information Systems*, 54, 147–155.
- Nummi, P.: Social media data analysis in urban e-planning. (2019). *Smart Cities and Smart Spaces: Concepts, Methodologies, Tools, and Applications*, pp. 636–651. IGI Global.
- Picazo-Vela, S., Gutiérrez-Martínez, I., Luna-Reyes, L.F. (2012). Understanding risks, benefits, and strategic alternatives of social media applications in the public sector. *Government Information Quarterly*, 29(4), 504–511.
- Siyam, N., Alqaryouti, O., Abdallah, S. (2020). Mining government tweets to identify and predict citizens engagement. *Technology in Society*, 60, 101211.
- Zavattaro, S.M., French, P.E., Mohanty, S.D. (2015) A sentiment analysis of US local government tweets: The connection between tone and citizen involvement. *Government Information Quarterly*, 32(3), 333–341.

About the Authors

José L. Lavado

Jose Luis Lavado is a postgraduate student of Mathematics and Computer Science Universidad Autónoma de Madrid, Spain. His main research interest is functional data analysis, but he is also interested in statistical learning.

Iván Cantador

Dr. Iván Cantador is a Senior Lecturer in Computer Science at Universidad Autónoma de Madrid, Spain. His main research lines are in the Recommender Systems field, where he has investigated a wide range of issues related to user modelling, knowledge representation, and processing and mining of user-generated content.

María E. Cortés-Cediel

María Elicia Cortés Cediel is a PhD candidate and Associate Lecturer at the Faculty of Political Science and Sociology of Universidad Complutense de Madrid, Spain. She is interested in citizen engagement in decision making scenarios, focusing on new forms of citizen participation supported by electronic tools.

Miriam Fernández

Dr. Miriam Fernandez is a Senior Research Fellow at the Knowledge Media Institute, The Open University, UK. Her current research focuses on the socio-technical aspects of Artificial Intelligence, and particularly on addressing biases, inequalities, and online harm.