

Arabizi sentiment analysis based on transliteration and automatic corpus annotation

Imane Guellil^{1,2}, Ahsan Adeel³, Faical Azouaou², Fodil Benali², ala-eddine Hachani², Amir Hussain³

1 Ecole Supérieure des Sciences Appliquées d'Alger ESSA-alger

2 Laboratoire des Méthodes de Conception des Systèmes (LMCS),

Ecole nationale Supérieure d'Informatique, BP 68M, 16309, Oued-Smar, Alger, Algérie

3 Institute of Computing science and Mathematics, School of Natural Sciences

University of Stirling Stirling UK

i.guellil@essa-alger.dz

{i_guellil, f_azouaou, df_benali, da_hachani}@esi.dz

{ahsan.adeel, ahu}@cs.stir.ac.uk

Abstract

Arabizi is a form of writing Arabic text which relies on Latin letters, numerals and punctuation rather than Arabic letters. In the literature, the difficulties associated with Arabizi sentiment analysis have been underestimated, principally due to the complexity of Arabizi. In this paper, we present an approach to automatically classify sentiments of Arabizi messages into positives or negatives. In the proposed approach, Arabizi messages are first transliterated into Arabic. Afterwards, we automatically classify the sentiment of the transliterated corpus using an automatically annotated corpus. For corpus validation, shallow machine learning algorithms such as Support Vectors Machine (SVM) and Naive Bays (NB) are used. Simulations results demonstrate the outperformance of NB algorithm over all others. The highest achieved F1-score is up to 78% and 76% for manually and automatically transliterated dataset respectively. Ongoing work is aimed at improving the transliterator module and annotated sentiment dataset.

1 Introduction

Sentiment analysis (SA), also called opinion mining, is the field of study that analyzes people's opinions, sentiments, evaluations, appraisals, attitudes, and emotions towards entities such as products, services, organizations, individuals, issues, events, topics, and their attributes. It represents a large problem space (Liu, 2012). To determine whether a document or a sentence expresses a positive or negative sentiment, three main approaches are commonly used, the lexicon based approach (Taboada et al., 2011), machine learning (ML) based approach (Maas et al., 2011) and a hybrid approach (Khan et al., 2015). English has

the greatest number of sentiment analysis studies, while research is more limited for other languages including Arabic and its dialects (Alayba et al., 2017; Guellil and Boukhalfa, 2015).

ML based sentiment analysis is a more dominant approach in the literature but it requires annotated training data. One of the major problems related to the treatment of Arabic and its dialect is the lack of resources. Other dominant problems include the non standard romanization (called Arabizi) that Arabic speakers often use in social media. Arabizi uses Latin alphabet, numbers, punctuation for writing an Arabic word (For example the word "mli7", combined with Latin letters and numbers, becomes the romanized form of the Arabic word "مليح" meaning "good"). To the best of our knowledge, limited work has been conducted on sentiment analysis of Arabizi ((Duwairi et al., 2016; Guellil et al., 2018)). The reason behind the lack of contribution is the complexity of Arabizi. Most researches are therefore moving towards the transformation of Arabizi into Arabic. This transformation or passage is recognized by the transliteration. Therefore, transliteration is only a process of passing from a written text in a given script or alphabet to another (Guellil et al., 2017c; Kaur and Singh, 2014). To bridge the gap, this paper proposes an approach determining the sentiment of Arabizi messages after transliterating them. This paper is organized as follows, Section 2 presents an overview of Arabizi. Section 3 presents the related work on SA and machine transliteration (MT). Section 4 presents the proposed approach and related components. Section 5 presents the simulation and experimentation. Finally, Section 6 presents the conclusion with some future directions.

2 Arabizi: An overview

Arabic speakers on social media, discussion forums, Short Messaging System (SMS), and on line chat applications often use a non standard romanization called "Arabizi" (Darwish, 2013; Bies et al., 2014). For example, the sentence: "rani fer7ana" (which means I am happy and correspond to the arabic sentence: رَاني فرحانة) is written in Arabizi. Hence, Arabizi is an Arabic text written using Latin characters, numerals and some punctuations (Darwish, 2013). The challenge behind Arabizi is the presence of many forms of the same word. For example the authors in (Ryan et al., 2014) argued that the word ان شاء الله (meaning if the god willing) could be written in 69 different manners.

3 Related work

3.1 Machine learning Arabic sentiment analysis

ML based sentiment analysis requires annotated data. Among the corpora presented in the literature and focused on MSA, we cite: LABR (Aly and Atiya, 2013), AWATIF (Abdul-Mageed and Diab, 2012), ASTD (Nabil et al., 2015) and ArTwitter (Abdulla et al., 2013). LABR contains 63,257 comments annotated with stars ranging from 1 to 5. AWATIF is a multi-genre corpus containing 10,723 sentences manually annotated in objective and subjective sentences. ASTD contains 10,000 Arab Tweets classified into objective, subjective positive, subjective negative or subjective mixed. ArTwitter contains 2,000 tweets manually annotated into positive and negative. However, most of the aforementioned works suffer from manual annotation and almost all resources are not publicly available. In addition, constructed corpora are dedicated to some dialects, neglecting others (specially Maghrebi dialect such as Moroccan or Algerian dialect).

3.2 Arabizi Transliteration

The proposed approach is inspired by the work presented in (van der Wees et al., 2016), where the authors used a table extracted from Wikipedia¹ for the passage from Arabizi to Arabic. The originality of our transliteration approach compared to this work is the treatment of ambiguities related to Arabizi transliteration such as: (a) Am-

¹https://en.wikipedia.org/wiki/Arabic_chat_alphabet

biguity of the vowels, where each vowels can be replaced by different letters or by NULL character (b) Ambiguity of the characters having the same sound or whose sounds are close, for example, the letters 's' and 'c' which can be replaced by the two letters س and ص (c) Ambiguity related to the transliteration direction, unlike the different works in (Guellil et al., 2017c,b), the rules of passage that we defined are from Arabizi to Arabic. The reverse passage may cause several ambiguities. The proposed approach is also inspired by the works presented in (Guellil et al., 2017c,b; Nouvel et al.) that uses a language model to determine the best possible candidate for a word in Arabizi. However, their work relies on a parallel corpus corresponding to the transliteration of a set of messages from Arabizi to Arabic. The realization of this corpus is usually done manually, which is a very time and effort consuming work. Hence, we avoid using a parallel corpus between Arabizi and Arabic and applied a language model (based on large corpus extracted from social media) to extract the best candidate.

3.3 Arabizi Sentiment Analysis

Different works have been proposed for handling Arabizi (Darwish (2013); Guellil and Faical (2017); Azouaou and Guellil (2017); Guellil and Azouaou (2016)). However, to the best of our knowledge, limited work has been conducted on sentiment analysis of Arabizi (Duwairi et al., 2016; Guellil et al., 2018). In (Duwairi et al., 2016), the authors presents a transliteration step before proceeding to the sentiment classification. However their approach present two majors drawbacks: (1) They rely on a very basic table for the passage from Arabizi to Arabic which cannot handle Arabizi ambiguities. (2) They construct a small annotated corpus manually (containing 3026 messages). This corpus contains Arabizi messages which therefore transliterated into Arabic. In (Guellil et al., 2018), the authors automatically construct an annotated sentiment Arabizi corpus and directly applied sentiment classification without calling the transliteration process. However, the authors confronted several ambiguity problems which resulted low F1-score of 66%. In contrast, the purpose of our paper is to present an approach dedicated to Arabizi sentiment analysis by calling transliteration process. The sentiment analysis corpus (training corpus) contains

4.3 Arabizi Transliteration

The proposed transliteration approach includes four important steps: (1) pretreatment of the Arabic corpus and the Arabizi message. (2) Proposal and application of the rules for the Algerian Arabizi. (3) Generating the different candidates. (4) Extraction of the best candidate. This part receives input, a set of messages written in Arabizi and a voluminous corpus written in DA extracted from Facebook. All these messages are pretreated (i.e. deleting exaggeration, etc). Afterwards, a set of passages rules are proposed (i.e. the letter 'a' could be replaced by 'ع، أ، ي، ة، '، etc. It could also be replaced by '، none letters when it represents a diacritic). By applying different replacements, as well as different rules developed, each Arabizi word gives birth to several words in Arabic. For example the word "kraht" generates 32 possible candidates, such as: 'كرهت', 'قرهت', 'كرهت' etc. The correctly transliterated word is 'كرهت'. The word "7iati" has 16 candidates such as: 'حياتي', 'حيطي', 'حيطي'. The correctly transliterated word is 'حياتي'. To extract the best candidate for the transliteration of a given Arabizi word into Arabic, a language model is constructed and applied.

4.4 Sentiment classification of Arabic messages

In this paper, different classification models are compared. The document embedding vectorization (Doc2vec algorithm presented within (Le and Mikolov, 2014)) is used (with default parameters). For Doc2vec, the two methods presented in (Le and Mikolov, 2014) were applied: (1) Distributed Memory Version Of Paragraph Vector (PV-DM) and (2) Distributed Bag of Words Version of Paragraph Vector (PV-DBOW). Moreover, the implementation merging these two methods is used. For the classification part, five different classifiers are used: (1) Support Vector Machine (SVM) (2) Naive Bayes (NB) (3) Logistic regression (LR) (4) Decision Tree (DT) and 5) Random Forest (RF).

5 Experimentations and results

5.1 Experimental Setup

The proposed approach is applied on a Maghrebi dialect (i.e. Algerian Arabizi) which suffers from limited available tools and other handling

resources required for automatic sentiment analysis. Algerian dialect (DALG) is largely presented in (Meftouh et al., 2012). However, the resources dedicated to the treatment of MSA cannot be directly applied to DALG. In this context, two large corpora were extracted from Facebook using RestFB³. The first one was extracted on September, 2017 which contains 8,673,285 messages with 3,668,575 written in Arabic letters. The second one was extracted on November, 2017 that contains 15,407,910 messages with 7,926,504 written in Arabic letters. The first one was used for transliteration task where the second one was used in sentiment annotation task. For testing our transliteration approach, we used Corpus_50 which is a part of Cottrell's corpus (Cottrell and Callison-Burch, 2014) used in (Guellil et al., 2017c,b,a). For testing our sentiment analysis approach, we used Corpus_500 (an Algerian Arabizi annotated corpus in (Guellil et al., 2018), containing 250 positives and negatives messages)

5.2 Experimental results

The first experiment evaluates the transliteration module. The transliteration of Corpus_50 achieves an accuracy up to 74.76% (as compared to 45.35% in (Guellil et al., 2017c)). This results shows the efficacy of the proposed transliteration approach. For sentiment analysis, we used Corpus_500. This dataset was transliterated automatically with the transliterator module. To validate the quality of the automatic transliteration, this dataset was also transliterated manually by Algerian dialect's natives. The transliteration of this dataset achieves an accuracy up to 72.05%. Afterwards, we carried out two types of experiments: (1) SA on test corpus transliterated automatically (2) SA on test corpus transliterated manually. Table 1 presents the performance of different shallow classification algorithms in terms of Precision (P), Recall (R) and F1-score (F1) for Doc2vec methods (PV-DBOW, PV-DM and PV-DBOW + PV-DM) and for Tr-automatic and Tr-manual dataset (respectively referring to the dataset transliterated automatically and manually).

5.3 Results and errors analysis

Based on the simulations and analysis, three major observations are: (1) The results with

³<http://restfb.com/>

Vectorization	classifier	Tr_automatic			Tr_manual		
		P	R	F1	P	R	F1
PV_DBOW	SVM	0.67	0.82	0.74	0.68	0.84	0.75
	NB	0.73	0.80	0.76	0.74	0.83	0.78
	LR	0.66	0.82	0.73	0.68	0.84	0.75
	RF	0.70	0.79	0.75	0.72	0.82	0.77
	DT	0.63	0.71	0.68	0.64	0.69	0.67
PV_DM	SVM	0.68	0.82	0.74	0.66	0.79	0.72
	NB	0.66	0.77	0.71	0.68	0.76	0.72
	LR	0.66	0.82	0.73	0.66	0.8	0.72
	RF	0.69	0.78	0.73	0.72	0.79	0.75
	DT	0.60	0.68	0.64	0.60	0.63	0.61
PV_DBOW + PV_DM	SVM	0.64	0.79	0.71	0.67	0.83	0.74
	NB	0.68	0.78	0.72	0.69	0.80	0.75
	LR	0.63	0.80	0.70	0.67	0.84	0.75
	RF	0.68	0.74	0.71	0.72	0.84	0.77
	DT	0.61	0.69	0.65	0.62	0.70	0.65

Table 1: Classification results with shallow machine learning

Tr_manual are slightly better than Tr_automatic (because the mistake on transliteration generally appears on only one letter), (2) The implementation PV_DBOW of Doc2vec achieved best results, (3) For classification, NB performed the best. (4) The results presented in Table 1 largely outperform the results presented in (Guellil et al., 2018) (which are up to 66%). However, we were not able to compare our results to those presented in (Duwairi et al., 2016) because their data are not available. However, the most observed errors are as follow:

- The principal error appears in transliteration process is related to technique of choosing the best candidate. The idea of language model is to extract the candidate having the most important number of occurrence. However, in some cases, this technique returns an incorrect candidate. For example the word "rakom" meaning "you are" is transliterated as "رقم" meaning "a number" rather than "راكم" (which is the correct transliteration). The solution to this problem is to integrate other parameters for determining the best candidate such as distance.
- Some sentiment classification errors are due to transliteration errors. For example, "khlwiya" meaning good and quiet is wrongly transliterated to "خليا" (meaning

empty) rather than خلوية. Improving transliteration will improve sentiment classification.

- Other sentiment classification errors are due to some errors occurred in the automatic annotated corpus (so the training corpus). For example, the messages جابوقامة الاسم تكفي meaning *Djabou the excellency of the name is sufficient* was annotated negative (where it is positive). Manually reviewing the automatic annotation will definitely improve the results.

6 Conclusion

In this paper, we present an approach to automatically classify sentiments of Arabizi messages (extracted from Facebook). The proposed approach constitutes an automatic annotation and transliteration. An Arabic sentiment lexicon is automatically constructed followed by automatic annotation and transliteration (Arabizi to Arabic). The developed dataset is validated using shallow machine learning, where the highest achieved precision is up to 78% and 76% for manual and automatic transliteration respectively with NB classifiers and PV_DBOW vectorization method. In the future, we intend to further enhance the proposed approach by improving the transliteration module focusing the annotated corpus (i.e manually reviewing the automatic annotation).

Acknowledgment

Imane Guellil and Faical Azouaou are respectively supported by Ecole Supérieure des Sciences Appliquées d'Alger ESSA-alger and Ecole nationale Supérieure d'Informatique. Amir Hussain and Ahsan Adeel were supported by the UK Engineering and Physical Sciences Research Council (EPSRC) grant No.EP/M026981/1.

References

- Muhammad Abdul-Mageed and Mona T Diab. 2012. Awatif: A multi-genre corpus for modern standard arabic subjectivity and sentiment analysis. In *LREC*, pages 3907–3914. Citeseer.
- Nawaf A Abdulla, Nizar A Ahmed, Mohammed A Shehab, and Mahmoud Al-Ayyoub. 2013. Arabic sentiment analysis: Lexicon-based and corpus-based. In *Applied Electrical Engineering and Computing Technologies (AEECT), 2013 IEEE Jordan Conference on*, pages 1–6. IEEE.
- Abdulaziz M Alayba, Vasile Palade, Matthew England, and Rahat Iqbal. 2017. Arabic language sentiment analysis on health services. In *Arabic Script Analysis and Recognition (ASAR), 2017 1st International Workshop on*, pages 114–118. IEEE.
- Mohamed Aly and Amir Atiya. 2013. Labr: A large scale arabic book reviews dataset. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, volume 2, pages 494–498.
- Faical Azouaou and Imane Guellil. 2017. Alg/fr: A step by step construction of a lexicon between algerian dialect and french. In *The 31st Pacific Asia Conference on Language, Information and Computation PACLIC 31 (2017)*.
- Stefano Baccianella, Andrea Esuli, and Fabrizio Sebastiani. 2010. Sentiwordnet 3.0: an enhanced lexical resource for sentiment analysis and opinion mining. In *LREC*.
- Ann Bies, Zhiyi Song, Mohamed Maamouri, Stephen Grimes, Haejoong Lee, Jonathan Wright, Stephanie Strassel, Nizar Habash, Ramy Eskander, and Owen Rambow. 2014. Transliteration of arabizi into arabic orthography: Developing a parallel annotated arabizi-arabic script sms/chat corpus. In *Proceedings of the EMNLP 2014 Workshop on Arabic Natural Language Processing (ANLP)*, pages 93–103.
- Ryan Cotterell and Chris Callison-Burch. 2014. A multi-dialect, multi-genre corpus of informal written arabic.
- Kareem Darwish. 2013. Arabizi detection and conversion to arabic. *arXiv preprint arXiv:1306.6755*.
- Rehab M Duwairi, Mosab Alfaqeh, Mohammad Wardat, and Areen Alrabadi. 2016. Sentiment analysis for arabizi text. In *Information and Communication Systems (ICICS), 2016 7th International Conference on*, pages 127–132. IEEE.
- Imane Guellil, Ahsan Adeel, Faical Azouaou, and Amir Hussain. 2018. Sentialg: Automated corpus annotation for algerian sentiment analysis. In *9th International Conference on Brain Inspired Cognitive Systems (BICS 2018)*.
- Imane Guellil and Faical Azouaou. 2016. Arabic dialect identification with an unsupervised learning (based on a lexicon). application case: Algerian dialect. In *Computational Science and Engineering (CSE) and IEEE Intl Conference on Embedded and Ubiquitous Computing (EUC) and 15th Intl Symposium on Distributed Computing and Applications for Business Engineering (DCABES), 2016 IEEE Intl Conference on*, pages 724–731. IEEE.
- Imane Guellil, Faical Azouaou, and Mourad Abbas. 2017a. Comparison between neural and statistical translation after transliteration of algerian arabic dialect. In *WinNLP: Women & Underrepresented Minorities in Natural Language Processing (collocated with ACL 2017)*.
- Imane Guellil, Faical Azouaou, and Mourad Abbas. 2017b. Neural vs statistical translation of algerian arabic dialect written with arabizi and arabic letter. In *The 31st Pacific Asia Conference on Language, Information and Computation PACLIC 31 (2017)*.
- Imane Guellil, Faical Azouaou, Mourad Abbas, and Sadat Fatiha. 2017c. Arabizi transliteration of algerian arabic dialect into modern standard arabic. In *Social MT 2017/First workshop on Social Media and User Generated Content Machine Translation*.
- Imane Guellil and Azouaou Faical. 2017. Bilingual lexicon for algerian arabic dialect treatment in social media. In *WinNLP: Women & Underrepresented Minorities in Natural Language Processing (collocated with ACL 2017)*. http://www.winlp.org/wp-content/uploads/2017/final_papers_2017/92_Paper.pdf.
- Imene Guellil and Kamel Boukhalfa. 2015. Social big data mining: A survey focused on opinion mining and sentiments analysis. In *Programming and Systems (ISPS), 2015 12th International Symposium on*, pages 1–10. IEEE.
- Kamaljeet Kaur and Parminder Singh. 2014. Review of machine transliteration techniques. *International Journal of Computer Applications*, 107(20).
- Aamera ZH Khan, Mohammad Atique, and VM Thakare. 2015. Combining lexicon-based and learning-based methods for twitter sentiment analysis. *International Journal of Electronics, Communication and Soft Computing Science & Engineering (IJECSCE)*, page 89.

- Quoc Le and Tomas Mikolov. 2014. Distributed representations of sentences and documents. In *International Conference on Machine Learning*, pages 1188–1196.
- Bing Liu. 2012. Sentiment analysis and opinion mining. *Synthesis lectures on human language technologies*, 5(1):1–167.
- Andrew L Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis. In *Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies-volume 1*, pages 142–150. Association for Computational Linguistics.
- Karima Meftouh, Najette Bouchemal, and Kamel Smaïli. 2012. A study of a non-resourced language: The case of one of the algerian dialects. In *The third International Workshop on Spoken Languages Technologies for Under-resourced Languages-SLTU'12*.
- Mahmoud Nabil, Mohamed Aly, and Amir Atiya. 2015. Astd: Arabic sentiment tweets dataset. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 2515–2519.
- Houda Saâdane1 Damien Nouvel, Hosni Seffih, and Christian Fluhr. Une approche linguistique pour la détection des dialectes arabes. In *24e Conférence sur le Traitement Automatique des Langues Naturelles (TALN)*, page 242.
- Cotterell Ryan, Adithya Renduchintala, Naomi Saphra, and Chris Callison-Burch. 2014. An algerian arabic-french code-switched corpus. In *Workshop on Free/Open-Source Arabic Corpora and Corpora Processing Tools Workshop Programme*, page 34.
- Maitte Taboada, Julian Brooke, Milan Tofiloski, Kimberly Voll, and Manfred Stede. 2011. Lexicon-based methods for sentiment analysis. *Computational linguistics*, 37(2):267–307.
- Mike Thelwall, Kevan Buckley, Georgios Paltoglou, Di Cai, and Arvid Kappas. 2010. Sentiment strength detection in short informal text. *Journal of the American Society for Information Science and Technology*, 61(12):2544–2558.
- Marlies van der Wees, Arianna Bisazza, and Christof Monz. 2016. A simple but effective approach to improve arabizi-to-english statistical machine translation. In *Proceedings of the 2nd Workshop on Noisy User-generated Text (WNUT)*, pages 43–50.