
The Challenge of Diacritics in Yorùbá Embeddings

Tosin P. Adewumi*, Foteini Liwicki & Marcus Liwicki

EISLAB, SRT
Luleå University of Technology, Sweden
firstname.lastname@ltu.se

Abstract

The major contributions of this work include the empirical establishment of a better performance for Yorùbá embeddings from undiacritized (normalized) dataset and provision of new analogy sets for evaluation. The Yorùbá language, being a tonal language, utilizes diacritics (tonal marks) in written form. We show that this affects embedding performance by creating embeddings from exactly the same Wikipedia dataset but with the second one normalized to be undiacritized. We further compare average intrinsic performance with two other work (using analogy test set & WordSim) and we obtain the best performance in WordSim and corresponding Spearman correlation.

1 Introduction

The Yorùbá language is spoken by about 40 million people in West Africa and around the world (Fakinlede, 2005). Of the various dialects around, the standard Yorùbá language (pioneered by Bishop Ajayi Crowther) is the focus of this paper. Standard Yorùbá orthography uses largely the Latin alphabet and is the widely spoken dialect among the educated (Bamgbose, 2000). Yorùbá has 25 letters in its alphabet, though counting the 5 nasal vowels makes it 30 (Fakinlede, 2005; Asahiah et al., 2017). Being a tonal language, 3 diacritics are used on vowels based on syllables per word: depression tone (grave), optional mid tone and elevation tone (acute) (Society, 1913). Besides these differences between the English and the Yorùbá languages, Yorùbá has no gender identification for verbs or pronouns (Nurse et al., 2010). Yorùbá verb tenses are usually determined within context and remain mostly the same in spelling and tone (Lamidi, 2010; Uwaezuoke and Ogunkeye, 2017).

The research question we address in this work is "Do diacritics affect the performance of Yorùbá embeddings and in what way?" This is because it has been observed by Asubiaro (2014) that web-search without diacritics produced more relevant results than search-words containing them, while evaluating four popular search engines. He also found out that the effectiveness of two of the search engines were adversely affected with diacritics. Thus, the objectives in this work include providing optimal Yorùbá embeddings and creating new analogy test set to evaluate the embeddings. Optimal hyper-parameter combination for the embeddings were chosen based on the work by Adewumi et al. (2020c,b). The heavily pre-processed (cleaned) Wikipedia dataset and the new analogy test set will provide valuable contributions to the natural language processing (NLP) community for the Yorùbá language, a low-resource language. The rest of this paper include the related work, the methodology, the results & discussion and the conclusion sections.

*Corresponding author

2 Related work

Initial effort by Ajayi Crowther to document Yorùbá barely had tonal marks (Bowen, 1858). In fact, early dictionary by Society (1913) had minimal diacritics compared to the modern Yorùbá dictionary by Smith and Onayemi (2005). This implies the language has been evolving and usage or discernment of diacritics between then and now is different. Revised efforts, later, standardized the diacritics and afforded others the opportunity to expand the work (Asahiah et al., 2017; Fagborun, 1989). For example, the word *abandon* in the Society (1913) dictionary is *kò-silẹ* while it is *kò-sílẹ* in the modern Lexilogos dictionary² and that by Smith and Onayemi (2005).

Absence of diacritics made contextual semantics of words, probably, more important back then than they are today, given that some words with the same spelling can have different meanings, depending on the context. Even the English language has words which are spelled the same way but pronounced differently and have different meanings (homographs), exposed by context, e.g. *lead*, *row* or *fair*. Given the relative challenge of producing Yorùbá diacritics among some users, the versions without diacritics or partial diacritics have been increasing (Asubiaro, 2014; Asahiah et al., 2017; Fagborun, 1989). This has led some to push for the normalization (restricting diacritized letters to their base versions) of the Yorùbá language, especially in electronic media (Asubiaro, 2014). This attempt may also lead to canonicalization of Yorùbá text, through the relationship between diacritized and undiacritized words that will be established.

Other researchers, like Asahiah et al. (2017) argue that diacritic restoration is a necessity. However, their own research showed the possible challenge for beginners of adding diacritics when the corpus they utilized had roughly the same percentage for the 3 diacritic marks (Asahiah et al., 2017). Yorùbá diacritic restoration is being undertaken by some researchers from word-level, syllable-level or character-level restoration and some of the methods for automatic diacritization utilize Machine Learning (ML) methods (Asahiah et al., 2017).

Word embeddings have shortcomings, such as displaying biases in the data they are trained on (Bolukbasi et al., 2016). However, they can be very useful for practical NLP applications. For example, subword representations have proven to be helpful when dealing with out-of-vocabulary (OOV) words and Thomason et al. (2020) used word embeddings to guide the parsing of OOV words in their work on meaning representation for robots. Intrinsic tests, in the form of word similarity or analogy tests, despite their weaknesses, have been shown to reveal meaningful relations among words in embeddings, given the relationship among words in context (Mikolov et al., 2013; Pennington et al., 2014). It is inappropriate to assume such intrinsic tests are sufficient in themselves, just as it is inappropriate to assume one particular extrinsic (downstream) test is sufficient to generalise the performance of embeddings on all NLP tasks (Gatt and Kraemer, 2018; Faruqui et al., 2016; Adewumi et al., 2020c,a).

3 Methodology

Three Yorùbá training datasets were used in this work. They include the cleaned 2020 Yorùbá Wikipedia dump containing diacritics to different levels across articles (Wikipedia, 2020), a normalized (undiacritized) version of it and the largest, diacritized data used by Alabi et al. (2020). The original Yorùbá Wikipedia dump has a lot of vulgar content, in addition to English, French & other language content. Manual cleaning brought the file size down to 182MB from 1.2GB, after using a Python script to remove much of the HTML tags, from the initial raw size of 1.7GB. Using the recommended script by Grave et al. (2018) to preprocess the original dataset did not work as intended, as it retained all the English & foreign content and removed characters with diacritics from the Yorùbá parts. An excerpt from the cleaned Wikipedia data, discussing about the planet Jupiter, is given below:

Awo osan ati brown inu isujo Júpítèrì wa lati iwusoke awon adapo ti won unyi awo won pada nigba ti won ba dojuko imọle [[ultraviolet]] lati ọdọ Orun. Ohun to wa ninu awon adapo wonyi ko daju, botilejepe fosforu, sulfur tabi boya [[hydro-carbon|haidrokarbon]] ni won je gbigbagbo pe won je.

²www.lexilogos.com/english/yoruba_dictionary.htm

The authors created two analogy test sets: one with diacritics and an exact copy without diacritics. However, all results reported in the next section were for the standard diacritic versions of the analogy and WordSim sets. The results based on the undiacritized WordSim set for both Wiki versions were poorer than what is reported in the next section but the undiacritized Wiki version still gave better results than the diacritized against that set. Creating the analogy sets (containing over 4,000 samples each) was challenging for some of the sections in the original Google version by Mikolov et al. (2013). For example, in the *capital-common-countries* sub-section of the semantic section, getting consistent representations of some countries, like *Germany*, is difficult, as it is translated as *Jemani* by some or *Jamani* by others. A very useful resource is Lexilogos, which translates from English to Yorùbá and, importantly, displays a number of contextual references where the translation is used in Yorùbá texts. The analogy sets are smaller versions of the original, with 5 sub-sections in the semantic section and only 2 sub-sections in the syntactic section. All datasets and relevant code used are available for reproducibility of these experiments.³ Four samples from the *gram2-opposite* of the diacritized version are given below:

wá lọ àgbà ọdọ
wá lọ òwúro irọlẹ
wá lọ ọtá ọrẹ
wá lọ nlá kékeré

Two types of embedding (word2vec and subword) per dataset were created, using the combination: skipgram-negative sampling with window size 4. The minimum and maximum values for the character ngram are 3 and 6, respectively, though the embedding by Grave et al. (2018) used ngram size of 5. Each embedding creation and evaluation was run twice to take an average, as reported in the next section. A Python-gensim (Řehůřek and Sojka, 2010) program was used to conduct the evaluations after creating the embeddings with the original C++ implementation by Grave et al. (2018). The Yorùbá WordSim by Alabi et al. (2020) was also used for intrinsic evaluation. This Yorùbá WordSim was based on the original English version by Finkelstein et al. (2001), containing a small set of 353 samples. However, the Yorùbá version had a few issues, which we corrected before applying it. For example, *television* is translated as *telifósiònù* instead of *telifíṣòn*, in one instance, and the bird *crane* is translated as *oti-bráńdi* (brandy) instead of *wádòwádò*, according to the Yorùbá dictionary.

4 Results & discussion

Tables 1 & 2 show results from the experiments while table 3 gives nearest neighbor result for the random word *iya* (*mother or affliction, depending on the context or diacritics*). Average results for embeddings from the 3 training datasets and the embedding by Grave et al. (2018) are tabulated: Wiki, U_Wiki, C3 & CC, representing embeddings from the cleaned Wikipedia dump, its undiacritized (normalized) version, the diacritized data from Alabi et al. (2020) and the Common Crawl embedding by Grave et al. (2018), respectively. Performance of the original, contaminated Wikipedia dump was poorer than the cleaned version reported here, hence, it was left out from the table. It can be observed from table 1 that the cleaned Wiki embedding have lower scores than the C3, despite the larger data size of the Wiki. This may be attributed to the remaining noise in the Wiki dataset. In spite of this noise, the exact undiacritized version (U_Wiki) outperforms C3, giving the best WordSim score & corresponding Spearman correlation. This seems to show diacritized data affects Yorùbá embeddings. The negative effect of noise in the Wiki word2vec embedding seems to reduce in the subword version in table 2.

The best analogy score is given by the embedding from Grave et al. (2018), though very small. The performance of the embeddings are much lower for analogy evaluations than their English counterparts as demonstrated by Adewumi et al. (2020c), though the comparison is not entirely justified, since different dataset sizes are involved. Other non-English work, however, show it's not unusual to get lower scores, depending, partly, on the idiosyncrasies of the languages involved (Adewumi et al., 2020b; Köper et al., 2015). NLP downstream tasks, such as named entity recognition (NER), with significance tests, will be the definitive measure for the performance of these embeddings, and this is being considered for future work.

³<https://github.com/tosingithub/ydesk>

Table 1: Yorùbá word2vec embeddings intrinsic scores (%)

Data	Vocab	Analogy	WordSim	Spearman
Wiki	275,356	0.65	26.0	24.36
U_Wiki	269,915	0.8	86.79	90
C3	31,412	0.73	37.77	37.83

Table 2: Yorùbá subword embeddings intrinsic scores (%)

Data	Vocab	Analogy	WordSim	Spearman
Wiki	275,356	0	45.95	44.79
U_Wiki	269,915	0	72.65	60
C3	31,412	0.18	39.26	38.69
CC	151,125	4.87	16.02	9.66

Table 3: Example qualitative assessment of undiacritized word2vec model

Nearest Neighbor	Result
iya	AgnEs (0.693), Arnauld (0.6798), olọlajulo (0.678), Rabiātu (0.6249), Alhaja (0.6186),...

5 Conclusion

The Yorùbá language is a tonal language and performance in NLP is affected, depending on diacritics, as shown in this work. It appears it is advantageous normalizing diacritized texts before working on them for NLP purposes, as they produce better intrinsic performance, generally. Our embeddings, based on normalized text, achieved better intrinsic performance than others tested. Future work will involve utilizing the embeddings in downstream tasks, such as NER, using state-of-the-art (SotA) architectures. Such downstream tasks will serve as the definitive measure for evaluating these embeddings. There’s ongoing effort on the sizable NER dataset to achieve this.

Broader Impact

The broader impact of this paper is the insight it provides for NLP researchers in Yorùbá language with regards to the differences in performance, based on diacritics. It provides 2 new analogy test sets for evaluating Yorùbá embeddings, depending on diacritics or the lack of it, and also provides an improved WordSim set. Furthermore, a heavily preprocessed Wikipedia dataset for training embeddings is provided, in the diacritized and undiacritized versions.

Acknowledgments and Disclosure of Funding

The work in this project is partially funded by Vinnova under the project number 2019-02996 "Språkmodeller för svenska myndigheter".

References

- Adewumi, T. P., Liwicki, F., and Liwicki, M. (2020a). Corpora compared: The case of the swedish gigaword & wikipedia corpora. *arXiv preprint arXiv:2011.03281*.
- Adewumi, T. P., Liwicki, F., and Liwicki, M. (2020b). Exploring swedish & english fasttext embeddings with the transformer. *arXiv preprint arXiv:2007.16007*.
- Adewumi, T. P., Liwicki, F., and Liwicki, M. (2020c). Word2vec: Optimal hyper-parameters and their impact on nlp downstream tasks. *arXiv preprint arXiv:2003.11645*.
- Alabi, J., Amponsah-Kaakyire, K., Adelani, D., and España-Bonet, C. (2020). Massive vs. curated embeddings for low-resourced languages: the case of yorùbá and twi. In *Proceedings of The 12th Language Resources and Evaluation Conference*, pages 2754–2762.

- Asahiah, F. O., Odejobi, O. A., and Adagunodo, E. R. (2017). Restoring tone-marks in standard yorùbá electronic text: improved model. *Computer Science*, 18.
- Asubiaro, T. V. (2014). Effects of diacritics on web search engines’ performance for retrieval of yoruba documents. *Journal of Library & Information Studies*, 12(1).
- Bamgbose, A. (2000). *A grammar of Yoruba*, volume 5. Cambridge University Press.
- Bolukbasi, T., Chang, K.-W., Zou, J. Y., Saligrama, V., and Kalai, A. T. (2016). Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Advances in neural information processing systems*, pages 4349–4357.
- Bowen, T. J. (1858). *Grammar and dictionary of the Yoruba language: with an introductory description of the country and people of Yoruba*, volume 10. Smithsonian institution.
- Fagborun, J. G. (1989). Disparities in tonal and vowel representation: some practical problems in yoruba orthography. *Journal of West African Languages*, 19(2):74–92.
- Fakinlede, K. J. (2005). *Beginner’s Yoruba*. Hippocrene Books.
- Faruqui, M., Tsvetkov, Y., Rastogi, P., and Dyer, C. (2016). Problems with evaluation of word embeddings using word similarity tasks. *arXiv preprint arXiv:1605.02276*.
- Finkelstein, L., Gabrilovich, E., Matias, Y., Rivlin, E., Solan, Z., Wolfman, G., and Ruppin, E. (2001). Placing search in context: The concept revisited. In *Proceedings of the 10th international conference on World Wide Web*, pages 406–414.
- Gatt, A. and Krahmer, E. (2018). Survey of the state of the art in natural language generation: Core tasks, applications and evaluation. *Journal of Artificial Intelligence Research*, 61:65–170.
- Grave, E., Bojanowski, P., Gupta, P., Joulin, A., and Mikolov, T. (2018). Learning word vectors for 157 languages. *arXiv preprint arXiv:1802.06893*.
- Köper, M., Scheible, C., and im Walde, S. S. (2015). Multilingual reliability and “semantic” structure of continuous word spaces. In *Proceedings of the 11th international conference on computational semantics*, pages 40–45.
- Lamidi, T. (2010). Tense & aspect in english & yoruba: Problem areas of yoruba learners of english. *Journal of the Linguistic Association of Nigeria Volume*, 13(2):349–358.
- Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013). Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- Nurse, D., Rose, S., and Hewson, J. (2010). Verbal categories in niger-congo languages.
- Pennington, J., Socher, R., and Manning, C. D. (2014). Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543.
- Řehůřek, R. and Sojka, P. (2010). Software framework for topic modelling with large corpora. In *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, pages 45–50, Valletta, Malta. ELRA. <http://is.muni.cz/publication/884893/en>.
- Smith, P. and Onayemi, A. (2005). Yoruba dictionary.
- Society, C. M. (1913). *Dictionary of the Yoruba Language*. Church Missionary Society Bookshop.
- Thomason, J., Padmakumar, A., Sinapov, J., Walker, N., Jiang, Y., Yedidsion, H., Hart, J., Stone, P., and Mooney, R. (2020). Jointly improving parsing and perception for natural language commands through human-robot dialog. *Journal of Artificial Intelligence Research*, 67:327–374.
- Uwaezuoke, A. H. and Ogunkeye, O. M. (2017). A contrastive morphological analysis of tense formation in igbo and yoruba: implication on learners and teachers. *UJAH: Unizik Journal of Arts and Humanities*, 18(3):193–219.
- Wikipedia (2020). Yoruba wikipedia multistream articles.