

Analysis on Filipino Word Embeddings

CSC714M - Theories in Natural Language Processing

Paolo Edni Andryn V. Espiritu

De La Salle University

1 Abstract

The development of word embeddings throughout the years have advanced various research areas in natural language processing (NLP). These are representations of a word where its semantic meanings are also captured in the form of a vector in multi-dimensional space. The distance of a vector from another indicates the similarity between the two word representations.

In this study, an analysis between two word embedding techniques will be explored — Word2Vec and FastText. The models were downloaded from Dan John Velasco’s GitHub repository. These will be evaluated on the scores it achieved from predicting the top-10 most similar word similarities or analogies. Overall, the results show that Word2Vec provided better predictions compared to the FastText model.

2 Description of the Models

2.1 Word2Vec

Word2Vec is a word embedding technique that was developed and published by Google in 2013. It represents each distinct word in the corpus as a vector of numbers in a multi-dimensional space. Each vector contains the word’s semantic meaning based on a span of n -words surrounding it. To measure the similarity between words, cosine similarity is used which implies that the two vectors are nearby one another in the vector space. Its architecture includes the implementation of Continuous Bag of Words (CBOW) and Skip-gram. To briefly describe the two methods, CBOW aims to reduce the distance between the target and actual words by predicting a target word based on its context. Conversely, skip-gram is used to predict the context words given the target word and its goal is to minimize the difference between the actual and target context words.

2.2 FastText

On the other hand, FastText is an open-source library created by Facebook’s AI Research (FAIR) Lab that may also be used for training word embeddings and performing text classification tasks. In this paper, the term ‘FastText’ will be used to refer to the word embedding model which extends from the features of a Word2Vec model. In contrast to Word2Vec, FastText considers each word as a bag of character n -gram. For instance, the word ‘hello’ at $n = 3$ will be broken down into {hel, ell, llo}.

3 Experimental Set-up

The word embedding models used in this assignment were extracted from Dan John Velasco’s GitHub repository.

3.1 Word Embedding Corpus

The author noted that the corpora used for training the models were the following [1]:

- WikiText-TL-39
- NewsPH-NLI
- Unpublished Twitter dataset

Before using these datasets, the author also performed the following pre-processing operations:

- | | | |
|------------------------------------|------------------------------------------------------------|---------------------------------------------------|
| 1. Retain stop words | 5. Removed brackets, parenthesis, braces, and its contents | 8. Replace numbers with <i>xx_digit</i> |
| 2. Retain commas | | 9. Replace amounts with <i>xx_amount</i> |
| 3. Lowercase text | 6. Removed punctuations | |
| 4. Removed quotes and its contents | 7. Removed symbols | 10. Replace percentages with <i>xx_percentage</i> |

As a result, the following are the statistics of the corpora after pre-processing:

- 4.68 million sentences
- 14.28 average sentence length
- 66.9 million tokens
- 1.08 million unique tokens

3.2 Models Used

The models were loaded using the *gensim* library as shown in the GitHub repository [1]. The following are the two models used in this study along with the dimensions, vocabulary size, and file size:

Models	Dimensions	Vocabulary Size	File Size
Word2Vec	300	126,687	269.9 MB
FastText	300	126,687	2.34 GB

Table 1: Word2Vec and FastText models

3.3 Evaluation

Two sets of test cases were prepared. The first set will gauge the performance of the model in predicting the top-10 most similar words based on the input. Whereas, the second set will test the capabilities of the model in providing the missing word for the given analogy. In Table 3, the types of analogies are synonymy, antonymy, part-whole, superclass, and geography from numbers 1 to 5 respectively.

To evaluate the models, the function, *most_similar*, was used under the *gensim* library. For word similarity, only two arguments were supplemented to the function, *positive* and *topn*, where *positive* contributes positively towards similarity and *topn* shows the resulting top-n words. For word analogies, another argument named *negative* was also used to specify the word that negatively affects the similarity of the first word. For example, given test 2 in Table 3, *buhay* and *laki* will be passed on to the function as *positive* then *patay* will be passed as *negative* to indicate that *laki* should correspond to *liit*

Test	Sample word
1	<i>gamot</i>
2	<i>ilaw</i>
3	<i>tubig</i>
4	<i>tao</i>
5	<i>pusa</i>

Table 2: List of word test cases

4 Analysis of Results

The tables discussed in this section contain outputs that were ranked from 1 to 10 where 1 is the best result of the model on the given task. For 4.1, the headers of the table indicate the input to the model. Furthermore, for 4.2, the headers of the table indicate the analogy given to the model and the blank line will be predicted by the model.

Test	Sample analogy	Expected Output
1	<i>Pinggan is to plato, as gwapo is to _____</i>	<i>pogi</i>
2	<i>Buhay is to patay, as laki is to _____</i>	<i>liit</i>
3	<i>Papel is to libro, as mata is to _____</i>	<i>mukha</i>
4	<i>Talong is to gulay, as asul is to _____</i>	<i>kulay</i>
5	<i>Australia is to Canberra, as Thailand is to _____</i>	<i>bangkok</i>

Table 3: List of analogy test cases

4.1 Word Similarity

Tables 4 and 5 show the results of Word2Vec and FastText models on predicting the most similar word given an input. As shown in the tables, Word2Vec produced better outputs in terms of semantic similarity to the input. This shows its effectiveness in understanding semantic relationships by using word-level vectors. On the other hand, the FastText model generated similar words based on the sequence of the word. For instance, the results under the *ilaw* column in Table 5 would always contain the substring *ilaw* in each of its output. It can be observed that it does not prioritize the semantic meaning of the input when predicting the most similar word.

One key difference of Word2Vec and FastText is that the FastText model uses subword-level vectors. As previously mentioned, this implies that FastText generate word embeddings by considering character n-grams in each word. With this, the model is often used to better represent embeddings in morphologically-rich languages such as Filipino. This is evident in the column of *gamot* where various forms of the word were also predicted by the model (e.g. *panggamot*, *manggamot*, *pagamot*, *gagamot*).

Rank	<i>gamot</i>	<i>ilaw</i>	<i>tubig</i>	<i>tao</i>	<i>pusa</i>
1	antibiotic (0.6556)	<i>kuryente</i> (0.5558)	<i>kuryente</i> (0.6309)	<i>taong</i> (0.7032)	<i>aso</i> (0.8511)
2	antibiotics (0.6533)	<i>kandila</i> (0.5447)	<i>hangin</i> (0.6051)	<i>mamamayan</i> (0.5896)	<i>kuting</i> (0.6781)
3	<i>alak</i> (0.6217)	<i>aircon</i> (0.5371)	<i>gripo</i> (0.5996)	<i>babae</i> (0.5783)	<i>babae</i> (0.6219)
4	<i>pagkain</i> (0.6126)	electricfan (0.5120)	<i>yelo</i> (0.5942)	<i>bagay</i> (0.5742)	<i>langgam</i> (0.5901)
5	<i>inumin</i> (0.6039)	<i>liwanag</i> (0.5055)	<i>kumukulong</i> (0.5640)	<i>lalake</i> (0.5586)	hamster (0.6309)
6	<i>gatas</i> (0.5966)	<i>binata</i> (0.5049)	<i>tubig-baha</i> (0.5634)	<i>pilipino</i> (0.5495)	<i>langaw</i> (0.5482)
7	vitamins (0.5947)	<i>kable</i> (0.5007)	<i>maiinom</i> (0.5615)	<i>bata</i> (0.5455)	<i>bubuyog</i> (0.5479)
8	biogesic (0.5868)	<i>kawad</i> (0.4946)	<i>dugo</i> (0.5602)	<i>nilalang</i> (0.5427)	<i>lamok</i> (0.5461)
9	<i>bakuna</i> (0.5765)	<i>kurtina</i> (0.4879)	<i>inumining</i> (0.5585)	<i>lalaki</i> (0.5294)	<i>ipis</i> (0.5432)
10	<i>antibiyotiko</i> (0.5672)	<i>tubig</i> (0.4759)	<i>pagkain</i> (0.5567)	<i>indibidwal</i> (0.4961)	<i>kambing</i> (0.5425)

Table 4: Word2Vec word similarity results

4.2 Word Analogy

In tables 6 and 7, the models results for completing the analogy are illustrated. From left to right excluding the rank column, the type of analogy used in each column is synonymy, antonymy, part-whole, superclass, and geography. The target words for these types of analogies are *gwapo*, *liit*, *mukha*, *kulay*, and *bangkok* respectively.

Consistent to the findings discussed in 4.2, the Word2Vec model completed the analogy by providing its best prediction of a word that is semantically similar to the expected output. For table 6, it can be observed that the model performed well on synonymy and antonymy since the model produced the expected words which are *pogi* (synonymous with *gwapo*) and *liit* (antonym of *laki*).

Rank	<i>gamot</i>	<i>ilaw</i>	<i>tubig</i>	<i>tao</i>	<i>pusa</i>
1	<i>gamotea</i> (0.8889)	<i>pailaw</i> (0.7953)	<i>tubig*</i> (0.9671)	<i>taoo</i> (0.8426)	<i>aso</i> (0.8613)
2	<i>panggamot</i> (0.8157)	<i>ilawom</i> (0.7897)	<i>tubigg</i> (0.9416)	<i>taod*</i> (0.8217)	<i>pusakal</i> (0.7980)
3	<i>pampagamot</i> (0.7876)	<i>pilaw</i> (0.7763)	<i>tubig-dagat</i> (0.8721)	<i>taoos</i> (0.8001)	<i>pusan</i> (0.7869)
4	<i>manggamot</i> (0.7831)	<i>tilaw</i> (0.7545)	<i>catubig</i> (0.8658)	<i>taooo</i> (0.7624)	<i>pusaaa</i> (0.7694)
5	<i>paggamot</i> (0.7692)	<i>ilaw</i> (0.7525)	<i>patubig</i> (0.8652)	<i>taong</i> (0.7246)	<i>pusanggala</i> (0.7630)
6	<i>pagamot</i> (0.7635)	<i>ilawa</i> (0.7512)	<i>tubig-baha</i> (0.8589)	<i>taoooo</i> (0.7239)	<i>daga</i> (0.7319)
7	<i>gamos</i> (0.7585)	<i>madilaw</i> (0.7464)	<i>tubig-alat</i> (0.8560)	<i>taob</i> (0.7122)	<i>pusang</i> (0.7298)
8	<i>gagamot</i> (0.7569)	<i>kilaw</i> (0.7270)	<i>matubig</i> (0.8499)	<i>taoyuan</i> (0.7103)	<i>pusaaaa</i> (0.7295)
9	<i>gamo</i> (0.7512)	<i>umiilaw</i> (0.7231)	<i>tubi</i> (0.8461)	<i>tao'y</i> (0.7097)	<i>kambing</i> (0.7218)
10	<i>panggagamot</i> (0.7384)	<i>ilawn</i> (0.7193)	<i>tubigan</i> (0.8407)	<i>taong-bayan</i> (0.6994)	<i>manol</i> (0.7201)

Table 5: FastText word similarity results

Table 7 illustrates the results of the FastText model on word analogies. As anticipated, the FastText model performed poorly on word analogies since it only based on the similarity of word sequences. It may also be noticed that FastText returned outputs with typographical errors. Despite producing erroneous outputs, this occurrence showcases one of the strengths of FastText since it can generate embeddings for words that did not exist during its training phase. This attribute is formally referred to as subword information where linguistic units smaller than words may be used to generate meaningful word embeddings.

5 Conclusion

In conclusion, this paper shows the differences between two word embedding techniques — specifically, Word2Vec and FastText. Through the experiments, it was observed that Word2Vec provided better predictions when it comes to semantic relationships of words. With the use of word-level vectors, it showed how words are similar to one another based on the context of its surrounding words.

Conversely, the FastText model predicted the most similar words based on a character-level. It was also evident that the model often predicted typographical errors and most of it are words that are not present in its training data. This shows that FastText can handle out-of-vocabulary words and better represent languages that are morphologically-rich.

References

- [1] Velasco, D. (2021). Filipino Word Embeddings. *GitHub*. Retrieved from <https://github.com/danjohnvelasco/Filipino-Word-Embeddings>.

Rank	pinggan : plato gwapo : ____	buhay : patay laki : ____	papel : libro mata : ____	talong : gulay asul : ____	australia : canberra thailand : ____
1	<i>pogi</i> (0.7468)	<i>liit</i> (0.4497)	<i>paningin</i> (0.4956)	<i>kahel</i> (0.4755)	japan (0.5281)
2	<i>ampogi</i> (0.5761)	<i>anlaki</i> (0.4467)	<i>braso</i> (0.4918)	<i>berde</i> (0.4515)	india (0.5260)
3	<i>gwapooo</i> (0.5610)	<i>napakalaki</i> (0.4350)	<i>pisngi</i> (0.4549)	<i>pula</i> (0.4441)	vietnam (0.5219)
4	<i>napakagwapo</i> (0.5427)	<i>bohail</i> (0.4337)	<i>ilong</i> (0.4541)	<i>abuhing</i> (0.4284)	indonesia (0.5177)
5	<i>cute</i> (0.5410)	<i>malaki</i> (0.4160)	<i>leeg</i> (0.4502)	<i>matingkad</i> (0.4224)	taiwan (0.5177)
6	<i>gwapoo</i> (0.5349)	<i>lahat</i> (0.3835)	<i>dibdib</i> (0.4458)	<i>puting</i> (0.4200)	singapore (0.5071)
7	<i>popogi</i> (0.5269)	<i>laking</i> (0.3755)	<i>tainga</i> (0.4251)	<i>pulang</i> (0.4197)	malaysia (0.5032)
8	<i>gwapoooo</i> (0.5193)	<i>buhai</i> (0.3697)	<i>kamay</i> (0.4223)	<i>itim</i> (0.4195)	myanmar (0.4905)
9	<i>pogiii</i> (0.5109)	<i>napakaliit</i> (0.3693)	<i>matang</i> (0.4221)	<i>berdeng</i> (0.4162)	turkey (0.4795)
10	<i>gwapooooo</i> (0.5017)	<i>paglaki</i> (0.3656)	<i>ulo</i> (0.4174)	<i>matitingkad</i> (0.4131)	korea (0.4692)

Table 6: Word2Vec analogy results

Rank	pinggan : plato gwapo : ____	buhay : patay laki : ____	papel : libro mata : ____	talong : gulay asul : ____	australia : canberra thailand : ____
1	<i>gwapo-gwapo</i> (0.8186)	<i>laki-laki</i> (0.6621)	<i>paningin</i> (0.6699)	<i>rasul</i> (0.5702)	australians (0.6875)
2	<i>g-gwapo</i> (0.7936)	<i>mlaki</i> (0.6124)	<i>mapapel</i> (0.6687)	<i>sul</i> (0.5470)	australis (0.6754)
3	<i>gugwapo</i> (0.7936)	<i>anlaki</i> (0.5820)	<i>takipmata</i> (0.6683)	<i>puting</i> (0.54541)	australasia (0.6733)
4	<i>pogi</i> (0.7768)	<i>laking</i> (0.5781)	<i>pilikmata</i> (0.6572)	<i>polong</i> (0.5428)	japan (0.6641)
5	<i>gwagwapo</i> (0.7765)	<i>lakin</i> (0.5688)	<i>namamalikmata</i> (0.6522)	<i>itim</i> (0.5382)	singapore (0.6621)
6	<i>ga-gwapo</i> (0.7635)	<i>kalaki</i> (0.5620)	<i>pagkatawan</i> (0.6471)	<i>berdeng</i> (0.5318)	indonesia (0.65971)
7	<i>ggwapo</i> (0.7601)	<i>anlaking</i> (0.5489)	<i>matamlay</i> (0.6470)	<i>talolong</i> (0.5256)	australian (0.6457)
8	<i>gagwapo</i> (0.7546)	<i>lakim</i> (0.5488)	<i>kamatyan</i> (0.6458)	<i>talon-talon</i> (0.5174)	malaysia (0.6355)
9	<i>gwapoo</i> (0.7515)	<i>buhay</i> (0.5481)	<i>pamata</i> (0.6440)	<i>hasul</i> (0.5172)	singaporeans (0.6304)
10	<i>angwapo</i> (0.7433)	<i>apakalaki</i> (0.5443)	<i>pakatawa</i> (0.6436)	<i>zilong</i> (0.5166)	australoid (0.6274)

Table 7: FastText analogy results