

# Improvement on Word Representations using Sentiment Analysis

Abhijat Shrivastava

SBUID 112584928

ashrivastava@cs.stonybrook.edu

Aishwarya Vijayakumar

SBUID 112673842

hvijayakumar@cs.stonybrook.edu

Debapriya Mukherjee

SBUID 112683344

dmukherjee@cs.stonybrook.edu

December 13, 2019

## Abstract

Word embeddings which consider the semantic and syntactic information from contexts are generally used for numerous natural language processing tasks. However the limitation is that they lack the ability to understand the underlying sentiment of the words. Analyzing the sentiment behind the words has become one of the most important research tasks in natural language processing currently. Sentiment analysis can be of huge importance and can be used in many applications such as advertising and marketing, business, politics and others. Various techniques are there for analyzing sentiments but lately, word embeddings method is widely being used in sentiment analyzing tasks. In this paper, we aim to improve the accuracy of the pre-trained word embeddings. We use a baseline model implemented with Word2Vec/GloVe methods and evaluate the accuracy of the word embeddings without Part-of-Speech (POS) tagging techniques and later with POS tagging techniques as our Idea 1, then evaluate the accuracy by implementing the bi-directional Long short-term memory (LSTM) model with POS tagging techniques. We evaluated the accuracy of the word-embeddings via different deep learning models and sentiment analyzing datasets. Our experimental evaluations show that the Improved Word Vectors (IWF) proposed in Idea 1, the implementation of Bi-LSTM model as proposed in Idea 2 are very effective for improving the accuracy of the word embeddings.

## 1 Introduction

The Word2Vec and GloVe require large corpuses for training and presenting an acceptable vector for each word. Investigators use pre-trained word vectors such as GloVe and Word2Vec because of the small size of some datasets. This would not be a best fit for the data. One more disadvantage of these word vector calculations is that they do not consider the context of a document. For example, the word vector for “beetle” as a car is equal to its word vector as an animal. These models ignore the sentiment information of the given text. The side effect of this problem is that words with opposite polarity are mapped into close vectors and it is a disaster for sentiment analysis.

Improved Word Vectors (IWF) is a method mentioned in the paper [1] which increases the accuracy of pre-trained word embeddings in sentiment analysis. This method is based on Part-of-Speech (POS) tagging techniques, lexicon-based approaches and Word2Vec/GloVe methods.

We demonstrate the advantage of specializing semantic spaces for either similarity or relatedness. Specializing for similarity is achieved by learning from both a corpus and a thesaurus, and for relatedness by learning from both a corpus and a collection of psychological association norms. We also compare the recently introduced technique of graph-based retrofitting (Faruqui et al., 2015) with a skip-gram retrofitting and a skip-gram

joint-learning approach.

All three methods yield specialized semantic spaces that capture human intuitions regarding similarity and relatedness significantly better than unspecialized spaces, in one case yielding state-of-the-art results for word similarity. More importantly, we show clear improvements in downstream tasks and applications: specialized similarity spaces improve synonym detection, while association spaces work better than both general-purpose and similarity specialized spaces for document classification.

## 2 Task

### 2.1 Baseline Model

A part of the baseline model being used for this task was available to us as a github implementation [3]. The dataset used for the training the baseline model is that of IMDB dataset.

We ran sentiment analysis on the IMDB dataset from the keras.datasets library for python. The model was trained using keras.models library for python.

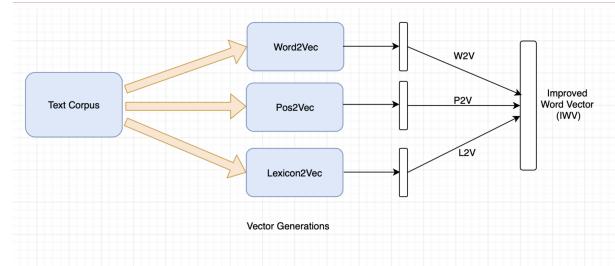


Figure 1: The main architecture of the the Improved Word Vector ( IWV)

The baseline on which we worked is based on the paper “Improving the Accuracy of Pre-trained Word Embeddings for Sentiment Analysis” by Seyed Mahdi Rezaeinia et al [1]. The paper talks about ‘Improved Word Vector’ (IWV), which takes into account word2vec, pos2vec and lexicon2vec before producing the final embedding. We have evaluated the accuracies based on word2vec and

pos2vec.

Word2Vec is a successful word embedding algorithm and is based on continuous Bag-of-Words (CBOW) and Skip-gram architectures which can provide high quality word embedding vectors. CBOW predicts a word given its context and Skip-gram can predict the context given a word. The generated vectors of words which appear in common contexts in the corpus are located close to each other in the vector space.

### 2.2 The issues

Word embedding techniques based on word2vec ignores the sentiment information of the given text . The side effect of this problem is that those words with opposite polarity are likely to be mapped into closer vectors and it might result in the loss of the true sentiments. [1]

The accuracy of the obtained model on a sentiment analysis task is 85%. This means that on a huge dataset, we'll lose a significant portion of the underlying true sentiments of the inputs. Our goal is to improve the word embeddings obtained from this technique using the three ideas mentioned below, so that true sentiments are retained while analysing a task.

## 3 Approach

### 3.1 Idea 1:

Part-of-speech (POS) tagging is an important and fundamental step in Natural Language Processing which is the process of assigning to each word of a text the proper POS tag.

Our idea 1 comprises of evaluating the improvement in accuracy of the word embeddings when computed with the Part-of-speech (POS) tagging techniques. The Part-of-speech gives large amount of information about a word and its neighbors, syntactic categories of words (nouns, verbs, adjectives, adverbs, etc.) and similarities and dissimilarities between them.

### 3.2 Idea 2:

For Idea 2, we thought of using a different method of generating the word vector embeddings instead of just using Word2Vec method. We used a Recurrent Neural Network (RNN) associated with Word2Vec network to achieve this. We applied Bidirectional LSTM in our model with the motive of improving the accuracy of the model.

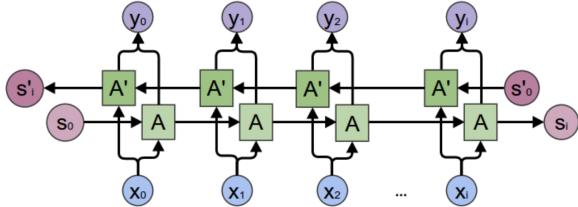


Figure 2: A Bi-directional LSTM model

### 3.3 Idea 3:

An idea that can be incorporated into our progress is that retrofitting of word vectors can be performed to further improve the training of the model. We target to minimize the euclidean distance between particular target vectors, so that the words with similar meaning get similar representations. This means that the words with same polarity will have their euclidean distance minimized. We propose this as a future scope of our project. The euclidean distance between two words can be represented by the formula:

$$\Psi(Q) = \sum_{i=1}^n \left[ \alpha_i \|q_i - \hat{q}_i\|^2 + \sum_{(i,j) \in E} \beta_{ij} \|q_i - q_j\|^2 \right]$$

Figure 3: Representative objective function. alpha beta represents relative strengths of information from original word vector and neighbours

We observe that most embedding models depend on the

correlation between the words appearing close to each other. The hypothesis behind this is that words occurring in similar contexts have similar meanings. But actually the words having similar contexts only means similar grammar, but not similar sentiment. So the vector spaces being trained from these models are suitable for analyzing word semantic relatedness(such as the relationship between keyboard and mouse), not word semantic similarity(such as the relationship between gasoline and diesel).

The idea 3 requires the usage of two thesauruses: the word similarity thesaurus from the electronic dictionary [MyThes] and the word relatedness thesaurus from the data set [USF Free Association Norms] to compare the similarity and the relatedness of the words. The paper[8] states the advantage of specializing semantic word embeddings for either similarity or relatedness. Here two variants of retrofitting and a joint-learning approach are taken care of, and capture human intuitions regarding similarity and relatedness.

### 3.4 Implementation Details

For implementing Idea 1, we used the baseline model and evaluated the model accuracy. We converted each generated POS tag to a constant vector and concatenated with Word2Vec vectors. The model was trained on 5 epochs, with a batch size of 1024. We found the accuracy to be 85%. The next step was to evaluate the model along with POS tagging technique to validate the significance of incorporating POS tags to the word embeddings.

For implementing Idea 2, we used a Bi-LSTM model associated with the Word2Vec model on the baseline model to generate the word embeddings.

## 4 Evaluation

### 4.1 Dataset Details

The datasets that have been used here [3] are as follows:

1. IMDB movie reviews dataset

```
Training done (data=1000000000)
  "she" was just brilliant, casting location, story, direction, everyone really called the
  "she" they played, and "she" could just imagine, being there, robert downey jr., is an amazing actor, and now, the
  "she" was real connection, with this film, the witty, breaks throughout the film were great, and "she" was just
  brilliant so much that "she" bought the ticket as soon as "she" was released from "down" and "she" would recommend
  it to anyone, "she" know, "she" know, "she" say, "you" cry a lot, "she" must have been good, and "this" definitely was
  a good movie, "she" know, "she" know,
  "she" children, are often left out of the "she" list, "she" think, because the stars, that play them all grow
  up, what they have done, "don't you" think, the whole story, was so lovely, because "it" was true, and "she" was
  Training a FastText model from Facebook Research
  read 1M words
  11355
  Number of labels: 2
  Number of words/sec/thread: 67299 lr: 0.000000 loss: 1.690068 ETA: 8h 0m
  Training complete
  vector model built
  oodles
```

Figure 4: IMDB Dataset

## 2. Another IMDB movie review dataset [4]

For evaluation of idea 2, we used a specific dataset which is well suited for sentiment analysis. A sentiment is either positive or negative. It has been categorized into 2 classes: 0 (for IMDB rating  $< 5$ , and 1 (for IMDB rating  $\geq 7$ ). Each movie has an upper bound of 30 reviews. The test and training datasets are different from each other. We have used the following files as our dataset:

- (a) labeledTrainData - The labeled training set. The file is tab-delimited and has a header row followed by 25,000 rows containing an id, sentiment, and text for each review.
  - (b) testData - The test set. The tab-delimited file has a header row followed by 25,000 rows containing an id and text for each review. Your task is to predict the sentiment for each one.
  - (c) unlabeledTrainData - An extra training set with no labels. The tab-delimited file has a header row followed by 50,000 rows containing an id and text for each review.

3. Two thesaurus used for checking the similarity and relatedness between words for Idea 3.

## 4.2 Evaluation Measures

Idea 1 was evaluated on dataset 1. Idea 2 was implemented and validated against dataset 2.

The metrics for the above mentioned ideas are: accuracy, loss, precision, recall and f1-loss.

For idea 3, the metrics are the learning rate and the loss of the model and pearson co-efficient while calculating the similarity and the relatedness of the word vectors.

## 4.3 Baselines

**For ideas 1 & 2:**

A part of the baseline model was available to us. We made some changes to the available baseline in order for it to work on our systems. Certain libraries required by the baseline were installed on our conda environment.

```

Epoch 2/5
33500/33500 [=====] - 45s 1ms/step - loss: 0.4068 - accuracy: 0.8218 - val_loss: 0.3882 - val_accuracy: 0.8282
Epoch 3/5
33500/33500 [=====] - 45s 1ms/step - loss: 0.3426 - accuracy: 0.8547 - val_loss: 0.3563 - val_accuracy: 0.8435
Epoch 4/5
33500/33500 [=====] - 45s 1ms/step - loss: 0.3149 - accuracy: 0.8684 - val_loss: 0.3456 - val_accuracy: 0.8563
Epoch 5/5
33500/33500 [=====] - 45s 1ms/step - loss: 0.2864 - accuracy: 0.8850 - val_loss: 0.3410 - val_accuracy: 0.8853
Classification Report:
          precision    recall   f1-score   support
0       0.84      0.87     0.85     8100
1       0.87      0.84     0.85     8391
   accuracy                           0.85      16500
   macro avg       0.85      0.85     0.85     16500
  weighted avg     0.85      0.85     0.85     16500
model built

```

Figure 5: Result after running the baseline model

**For idea 3:**

We obtain a part of another baseline model for the paper mentioned in [8] in github [9] to check the similarity and relatedness between word embeddings based on two available thesaurus [7,8]. We made certain changes in the baseline to make the model run on our systems and many requirements needed to be installed for that.

We executed the baseline model for this on our systems and observed the following:

```
[INFO] 106848 106848 0.000 secs: 1166.393
running time: 35.51 mins
(nlp-project) D:\PycharmProjects\DeepLearning\Specialized\crossdomain\multimodel\master\deepbayes\ymmler\models$ python evaluation_with_pearson_coeff.py ./models/raw_emb
relatedness evaluation: raw embedding (0.398) vs specialize for relatedness embedding (0.547)
relatedness evaluation: raw embedding (0.398) vs specialize for relatedness embedding (0.547)
```

Figure 6: Executing the baseline model for the Idea 3

- relatedness evaluation: raw embedding(0.398) vs specialize for relatedness embedding(0.547)
  - similarity evaluation: raw embedding(0.260) vs specialize for similarity embedding(0.513)

We have an executable baseline for this and we can extend the idea behind the baseline as a part of our future scope to work on for improving the word embeddings.

## 4.4 Results

### Idea 1:

On comparing the accuracies of the word embeddings with and without POS tagging technique, we see an accuracy improvement of 1% after using POS tagging on the words.

Model (Evaluated in 5 epochs)	Accuracy
With POS tags	86.0%
Without POS tags	85.0%

Figure 7: Result for evaluating word embeddings with and without POS tagging technique

### Idea 2:

In this idea, we have implemented a Bi-LSTM model to generate the word embeddings. On implementing the Bi-LSTM model, we see a significant increase in the accuracy of the model. We evaluate it on 20 epochs and observe the following over the 20 epochs:

1. We train the model on 23750 samples, and validate on 1250 samples.
2. The accuracy gradually increases from 0.6306 to 0.9805.
3. The loss gradually decreases, starting from 0.6439 at the first epoch and ending with 0.0547 at the last epoch.

```
Train on 23750 samples, validate on 1250 samples
Epoch 0/20
23750/23750 [=====] - 135s/step - loss: 0.6439 - accuracy: 0.6306 - val_loss: 0.6467 - val_accuracy: 0.6272
23750/23750 [=====] - 135s/step - loss: 0.5381 - accuracy: 0.7332 - val_loss: 0.5355 - val_accuracy: 0.6984
23750/23750 [=====] - 202s/step - loss: 0.4604 - accuracy: 0.7883 - val_loss: 0.4586 - val_accuracy: 0.7832
Epoch 4/20
23750/23750 [=====] - 164s/step - loss: 0.4056 - accuracy: 0.8187 - val_loss: 0.4088 - val_accuracy: 0.8088
Epoch 8/20
23750/23750 [=====] - 138s/step - loss: 0.3556 - accuracy: 0.8445 - val_loss: 0.4407 - val_accuracy: 0.8064
23750/23750 [=====] - 135s/step - loss: 0.3062 - accuracy: 0.8730 - val_loss: 0.4689 - val_accuracy: 0.8088
Epoch 12/20
23750/23750 [=====] - 132s/step - loss: 0.2541 - accuracy: 0.8946 - val_loss: 0.5029 - val_accuracy: 0.8032
Epoch 16/20
23750/23750 [=====] - 132s/step - loss: 0.2090 - accuracy: 0.9147 - val_loss: 0.5308 - val_accuracy: 0.7928
Epoch 20/20
23750/23750 [=====] - 131s/step - loss: 0.1718 - accuracy: 0.9313 - val_loss: 0.6232 - val_accuracy: 0.8088
Epoch 0/20
23750/23750 [=====] - 131s/step - loss: 0.1439 - accuracy: 0.9446 - val_loss: 0.6597 - val_accuracy: 0.8088
Epoch 1/20
23750/23750 [=====] - 131s/step - loss: 0.1274 - accuracy: 0.9496 - val_loss: 0.7188 - val_accuracy: 0.7944
Epoch 2/20
23750/23750 [=====] - 133s/step - loss: 0.1060 - accuracy: 0.9595 - val_loss: 0.7885 - val_accuracy: 0.7872
Epoch 3/20
23750/23750 [=====] - 133s/step - loss: 0.0935 - accuracy: 0.9654 - val_loss: 0.8185 - val_accuracy: 0.8096
Epoch 4/20
23750/23750 [=====] - 133s/step - loss: 0.0809 - accuracy: 0.9691 - val_loss: 0.7678 - val_accuracy: 0.7968
Epoch 5/20
23750/23750 [=====] - 136s/step - loss: 0.0800 - accuracy: 0.9693 - val_loss: 0.7580 - val_accuracy: 0.8016
Epoch 6/20
23750/23750 [=====] - 134s/step - loss: 0.0734 - accuracy: 0.9724 - val_loss: 0.8146 - val_accuracy: 0.8112
Epoch 7/20
23750/23750 [=====] - 133s/step - loss: 0.0669 - accuracy: 0.9750 - val_loss: 0.8168 - val_accuracy: 0.8128
Epoch 8/20
23750/23750 [=====] - 136s/step - loss: 0.0581 - accuracy: 0.9793 - val_loss: 0.8183 - val_accuracy: 0.8048
23750/23750 [=====] - 133s/step - loss: 0.0580 - accuracy: 0.9790 - val_loss: 0.8197 - val_accuracy: 0.8056
23750/23750 [=====] - 133s/step - loss: 0.0547 - accuracy: 0.9889 - val_loss: 1.0423 - val_accuracy: 0.9868
```

Figure 8: Evaluating the accuracies in 20 epochs on implementing Bi-LSTM model

The Bi-LSTM model is a bidirectional recurrent neural networks for document-level sentiment analysis with word2vec Embedding. The model was implemented with Bi-LSTM model, and the experiment shows that the Bi-LSTM model associated with word2vec word embedding outperformed the baseline model and achieved a 98.05% accuracy.

## 4.5 Analysis

### Idea 1:

We analyzed our results by comparing the accuracy of the models on a sentiment analysis task on a IMDB dataset obtained from keras.datasets library, available for python. We evaluate the system once without using the POS tagging techniques and then using it. On comparing both, the results show that the system appears to work well when Part-of-Speech tags are incorporated into the word embeddings, since the accuracy increases by 1% (for 5 epochs). This shows that incorporating POS tags into word embeddings helps them absorb more sentiment-related information, which can be put to use in sentiment analysis tasks.

### Idea 2:

On evaluating the accuracies of the word embeddings with POS tagging techniques, we found that there is a slight improvement. On implementing the Bi-LSTM

model associated with the Word2Vec model along with the POS tagging technique (idea 1) to generate the word embeddings, we observed a significant improvement in the word vector representations. This improvement is supported by the boost in accuracy of the sentiment analysis task.

## 4.6 Code

The project we've been working on has been uploaded on the following Github repository: <https://github.com/abhijat96/NLPPProject>

## 5 Conclusions

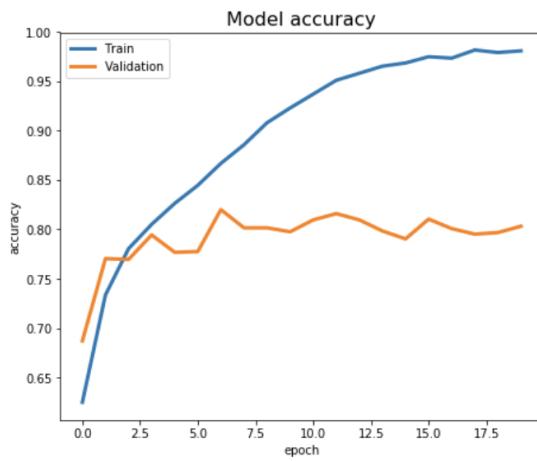


Figure 9: Graph of accuracy of model

From the graph above, we observed that the model accuracy keeps increasing gradually as the number of epochs increases on implementing the biLSTM model associated with the word2Vec model, incorporated with POS tagging techniques. Thus we can conclude that the final model shows a significant improvement on the baseline model. Therefore, we have improved the accuracy of the pre-trained word embeddings.

## References

- [1] Seyed Mahdi Rezaeinia, Ali Ghodsi, Rouhollah Rahmani. *Improving the Accuracy of Pre-trained Word Embeddings for Sentiment Analysis.*
- [2] Liang-Chih Yu, Jin Wang, K. Robert Lai and Xuejie Zhang *Refining Word Embeddings for Sentiment Analysis.*
- [3] Github Link (Baseline 1): <https://github.com/PrashantRanjan09/Improved-Word-Embeddings>.
- [4] Kaggle Link for Data: <https://www.kaggle.com/jiaofenx/imdb-review-word2vec-rnn-tutorial/data>.
- [5] Manaal Faruqui, Jesse Dodge, Sujay K. Jauhar, Chris Dyer, Eduard Hovy, Noah A. Smith *Retrofitting Word Vectors to Semantic Lexicons.*
- [6] Thesaurus: <http://www.openoffice.org/lingucomponent/thesaurus.html>.
- [7] Thesaurus: <http://w3.usf.edu/FreeAssociation/>.
- [8] Douwe Kiela, Felix Hill and Stephen Clark *Specializing Word Embeddings for Similarity or Relatedness.*
- [9] Github Link (Baseline for Idea 3): <https://github.com/huangmiaoxin/SpecializeWordEmbedding>.