

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

Improving Paragraph2Vec

Seokho Hong
seokho@stanford.edu

Abstract

Paragraph vectors were proposed as a powerful unsupervised method of learning representations of arbitrary lengths of text. Although paragraph vectors had the advantage of being versatile, being unsupervised and unconstrained by lengths of text, the concept has not been further developed since its first publication. We propose two extensions upon the initial formulation of the paragraph vector, and test its performance on two separate semantic-based tasks. Although the results are limited by the fact that our attempt to reproduce the original paragraph vectors was not successful, we can still show that the extended models outperform the original paragraph vectors.

1 Introduction

The Paragraph Vector (Le & Mikolov 2014) was proposed with the objective of unsupervised semantic learning of arbitrary lengths of texts. While the PV-DM and PV-DBOW models proposed perform very well, their formulations leave room for improvement. In particular, the performance differences between PV-DM and PV-DBOW suggest that improvements to PV-DM could yield be significant. Also, the paper finds that PV-DM and PV-DBOW vectors concatenated achieve the best performance. Even if different training methods do not necessarily produce better vectors, concatenating them with the original paragraph vectors may offer better results than the PV-DM and PV-DBOW concatenated alone.

The main mathematical limitation of the PV-DM and PV-DBOW models is that they do not allow the paragraph vector to interact in complex, non-linear ways with the word vectors. While this design is not entirely unjustified, since different sections of text do not radically change the distribution and sequence of English words, there is nonetheless a limitation on the extent to which a paragraph vector can assist in the word-prediction task the paragraph vector uses.

We propose two different formulations of training unsupervised paragraph vectors, which we will call the hidden layer model and the tensor model. While both models offer improvements upon the original paragraph models, we cannot make definitive conclusions because we were unable to replicate the same level of performance reported in (Le & Mikolov 2014).

2 Background

Recent works have proposed various deep methods for learning the semantics of sentences. Recursive neural networks with various substructures (Socher et al. 2013), (Tai et al. 2015) and convolutional neural networks (Kalchbrenner et al. 2014), offer excellent performance that approximately matches or exceeds that of paragraph vectors across different tasks. These models, however, tend to be far more complex and deeper than the paragraph vector method and therefore take significantly longer to train. They also have other limitations that could potentially limit their range of applications. Recursive models work only for sentences and need a parser, which is not

054 only an extra requirement, but also a source of errors given that parsers are imperfect. Convolutional
055 neural networks as suggested in (Kalchbrenner et al. 2014) do not need a parser, but remain untested
056 for longer pieces of text. The max-pooling layer the paper proposes appears as though it will be less
057 effective as the text gets much longer.

058
059 The paragraph vector has the advantage of simplicity and versatility. Perhaps it is a bit too simple.
060 The models suggested here were inspired by the recent gains from the increasingly complex models.
061 It does not appear unreasonable to train modestly more complex models along the lines of the
062 paragraph vector framework if it will result in performance gains.

063
064

065 **3 Approach**

066 **Paragraph Vector Framework**

067
068 The paragraph vector framework is the general approach to training unsupervised paragraph vectors,
069 and is common to the models discussed here as well as the original PV-DM. The framework is a
070 word prediction task. Given a set of words w_1, w_2, \dots, w_n , the model trains by predicting one of
071 the words' vector, w_j given the other $n - 1$ words' vectors. The model also takes a paragraph vector
072 p_i where i identifies which body of text w_1, w_2, \dots, w_n come from.

073
074

075 The original authors proposed both hierarchical softmax and negative sampling as a replacement
076 for the traditional but expensive softmax for the objective function, but we will use only negative
077 sampling here. The cost function is therefore:

078

$$079 \log(\sigma(r \cdot v_{w_j})) - \sum_i \log(\sigma(r \cdot v_{w_i})) \quad (1)$$

080
081

082 Where r is the predicted vector, w_i is the vector of a random word, for which there are k random
083 words. For the models trained here, k is fixed at 10.

084
085

086 The training is done via backpropagation through structure using Adagrad.

087

088 In training, both the word vectors and paragraph vectors are initialized randomly with values in the
089 range of -0.01 to 0.01 .

090
091

092 **Hidden Layer Model**

093
094 The original PV-DM model has the following equation for r :

095

$$096 r = \sigma(W \cdot [c; p_i]) \quad (2)$$

097
098

099 where c is the concatenation of the input word vectors for the prediction task.

100
101

102 The hidden layer model we propose simply adds another layer between r and the two input vectors.

103
104

$$105 z = \sigma(W \cdot [c; p_i]) \quad (4)$$

106

$$107 r = \sigma(U \cdot z) \quad (5)$$

108

108 While simple, it allows the components of p_i and c to interact in more complex ways. The original
109 equation for r allows only a single non-linear transformation of a linear function of p_i and c .
110

112 **Tensor Model**

114 To allow even more interaction between the components of p_i and c , we propose the following:

$$116 r = \sigma(c \cdot T \cdot p_i + W \cdot [c; p_i]) \quad (7)$$

$$117 \quad (8)$$

119 where T is a tensor.
120

122 **3.1 Training**

124 Given a collection of text that can be divided into n documents, or paragraphs, each paragraph is
125 assigned a vector. Training involves sliding the window of context words across each word of each
126 paragraph, for every paragraph. At each data point, the input is the window of context words, and
127 the paragraph vector corresponding to the origin of the words, and the target output is a particular
128 word within or adjacent to the context words (depending on the hyper parameters). The target
129 output is obviously not given as input. At either end of the paragraph, a special "NULL" word is
130 applied as needed to fill the necessary window space. The "NULL" word is trained like every other
131 word.

132 Backpropagation is used to train both the paragraph vectors and the word vectors simultaneously.
133

135 **3.2 Testing**

137 Testing involves running gradient descent to train new paragraph vectors for each new paragraph.
138 At test time, all parameters of the model are frozen, including the word vectors, and the backpropa-
139 gation is only applied to the paragraph vectors.
140

142 **4 Experiments**

144 We tested the two new models on two fairly standard tasks, one on which other models have been
145 benchmarked: sentiment analysis and semantic textual similarity.
146

147 **4.1 Sentiment Analysis - Stanford Sentiment Treebank**

149 Each phrase in the treebank was treated as a separate paragraph, and training and testing was done ac-
150 cording to the dataset's specifications. In an attempt to reproduce the results from (Le and Mikolov,
151 2014), we tried to match the hyperparameters as closely as possible.
152

153 **4.1.1 Hyperparameters for PV-DM**

154 Word vector dimensions: 100; paragraph vector dimension: 400; context window size: 7 words,
155 predict 8th; trained using Adagrad with 0.01 learning rate, minibatches of size 300, L2 regularization
156 of $1e-4$. Trained for approximately 10 hours.
157

158 **4.1.2 Hyperparameters for Hidden Layer Model**

159 Word vector dimensions: 100; paragraph vector dimension: 400; context window size: 7 words,
160 predict 8th; trained using Adagrad with 0.01 learning rate, minibatches of size 300, L2 regularization
161 of $1e-4$. Trained for approximately 16 hours.

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

4.1.3 Hyperparameters for Tensor Model

Word vector dimensions: 100; paragraph vector dimension: 200; context window size: 7 words, predict 8th; trained using Adagrad with 0.005 learning rate, minibatches of size 300, L2 regularization of $1e-3$. Trained for approximately 48 hours.

The tensor model training becomes very expensive with large word or paragraph vectors so the paragraph vector dimensionality was reduced.

Final Classification

After pre-training the paragraph vectors (both for training and testing), a random forest classifier was used to predict the final sentiment label for each paragraph. On the fine-grained task, each paragraph was classified in categories from 1 to 5, where the original sentiment is scaled up from 0.0 to 1.0. On the binary task, each paragraph rated in categories 2 and 4 were used. We used the random forest classifier from the sci-kit learn package, with 100 trees as the only modification from the default hyperparameters.

Table 1: SST Results

Method	Fine-Grained	Binary
PV-DM	41.6	81.2
PV-Hidden	42.7	82.2
PV-Tensor	44.1	84.3
PV-DM (Le and Mikolov, 2014)	48.7	87.8

5 Attempted Optimization

We made many attempts to fully reproduce the results from (Le and Mikolov, 2014). The only major difference is that we reduced the word vector dimensionality to 100 from 400, but this reduction was necessary to allow the model to train in a reasonable amount of time. Training a model with 200 dimensions did improve the results, but not enough to suggest that the reduction from 400 is responsible for the disparity in results.

Training duration was determined by convergence on the development set (10% randomly sampled from the training set). If there was no improvement over 10 epochs, or passes through the entire training set, then the training was halted.

Dropping out for the Hidden Layer model ($p = 0.5$ on W) and Tensor model ($p = 0.2$ on T) was used.

5.1 Semantic Textual Similarity - SemEval 2014 Task 1

While (Le and Mikolov, 2014) did not benchmark their PV-DM on this task, many other papers benchmark their neural networks on this dataset, suggesting that the dataset is a reliable one.

The dataset here is the SICK dataset, 10000 pairs of English sentences, each labeled with a semantic similarity score from 1 to 5. There is a Train, Test, and Trial division of the dataset, and each were used respectively for training, testing, and validating the models.

5.1.1 Hyperparameters for PV-DM

Word vector dimensions: 200; paragraph vector dimension: 400; context window size: 9 words, predict 10th; trained using Adagrad with 0.01 learning rate, minibatches of size 300, L2 regularization of $1e-4$. Trained for approximately 6 hours.

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

5.1.2 Hyperparameters for Hidden Layer Model

Word vector dimensions: 200; paragraph vector dimension: 400; context window size: 9 words, predict 10th; trained using Adagrad with 0.01 learning rate, minibatches of size 300, L2 regularization of $1e-4$. Trained for approximately 10 hours.

5.1.3 Hyperparameters for Tensor Model

Word vector dimensions: 200; paragraph vector dimension: 200; context window size: 9 words, predict 10th; trained using Adagrad with 0.005 learning rate, minibatches of size 300, L2 regularization of $1e-3$. Trained for approximately 26 hours.

Final Classification

After pre-training the paragraph vectors (both for training and testing), a random forest regressor was used to predict the final sentiment label for each paragraph. We used the random forest regressor from the sci-kit learn package, with 100 trees as the only modification from the default hyperparameters.

Table 2: Semantic Textual Similarity Results

Method	MSE
Mean Vectors (Tai et al, 2015)	0.455
LSTM (Tai et al, 2015)	0.281
PV-DM	0.392
PV-Hidden	0.388
PV-Tensor	0.365

The Mean vectors baseline in (Tai et al, 2015) computes a semantic relatedness score from the average of the word vectors of the words in the sentence. While the LSTM is not the focus of the paper, it is one of the most effective models for the task, and puts the paragraph vector models in perspective, at least our implementations.

6 Attempted Optimization

Since this example used the same implementation as the one used for the Stanford Sentiment Treebank, no doubt there are some flaws holding back the scores here.

The SICK dataset is smaller than the previous one, which allowed slightly larger models to be trained in a reasonable time.

Training duration was determined by convergence on the development set, as specified in the dataset. If there was no improvement over 10 epochs, or passes through the entire training set, then the training was halted.

Dropping out for the Hidden Layer model ($p = 0.5$ on W) and Tensor model ($p = 0.2$ on T) was used.

7 Conclusions

Despite the subpar results on the PV-DM implementation, it is safe to say that additional hidden layer and the tensor models improve upon the original PV-DM formulation. The hidden layer model

270 probably is not worth the additional training time required for the slight performance gains. The
271 tensor model is much better, but it takes even longer to converge.
272

273 The performance of the original PV-DM showed that modeling the distribution and ordering of a
274 paragraph's words, is an effective method of modeling the semantics of a sentence or any piece
275 of text. The performance of the hidden layer and tensor models shows that there is room for
276 improvement by allowing the paragraph vector to learn a more complex function for how the
277 paragraph influences the distribution of words.
278

279 Regarding direct improvement of the model, perhaps a deeper, more complex model would further
280 improve the performance of paragraph vectors, but at that point it would no longer be the fairly
281 simple model that trains quickly compared to other deep models.
282

283 For improving performance on particular tasks such as sentiment analysis, models that train
284 directly supervised representations of text are outperforming the unsupervised paragraph vector. It
285 intuitively seems obvious that deep methods that optimize paragraph representations for a particular
286 task will outperform a shallow classifier on an unsupervised paragraph representation. It is also not
287 possible to fine tune the unsupervised paragraph vectors on a specific task unless all the paragraphs
288 are available during training time, which does not prove generalization. Thus improving the
289 paragraph vector is probably not the best method for attempting to improve upon state of the art
290 performance on separate tasks.
291

292

293 **References**

294

295 [1] Kalchbrenner, Nel, Edward Grefenstette, and Phil Blunsom. 2014. A Convolutional Neural Network for
296 Modelling Sentences. *ACL14*

297 [2] Le, Quoc V and Tomas Mikolov. 2014. Dis-tributed representations of sentences and doc-
298 *ument arXiv preprint arXiv*

299 [3] Tai, Kai Sheng, Richard Socher, Christopher D. Manning 2015. Improved Semantic Representations From
300 Tree-Structured Long Short-Term Memory Networks *uments arXiv preprint arXiv*

301 [4] Socher, Richard, Alex Perelygin, Jean Y. Wu, Jason Chuang, Christopher D. Manning, Andrew Y. Ng
302 and Christopher Potts. Recursive Deep Models for Semantic Compositionality Over a Sentiment Treebank
303 *uments EMNLP 2013*

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323