

# Implementing Adversarial Trainings to Text Classification

Shaguna Awasthi

Information Technology, Maharaja Agrasen Institute of Technology

**KEYWORDS:** Text Classification , Research Papers, Adversarial Trainings, Supervised algorithm , Word Embeddings, LSTM

## I. Introduction

Adversarial examples are models that are made by making small perturbations to the information de- endorsed to fundamentally increase the loss incurred by an AI model (Szegedy et al., 2014; Goodfellow et al., 2015). A few models, including cutting edge convolutional neural networks, come up short on the capacity to classify adversarial models effectively, at times in any event, adversarial perturbations obliged to be little to the point that a humans can't see it. Adversarial trainings are the way toward preparing a model to accurately order both unmodified models and promotion adversarial models. It improves robustness but also generalization performance for original examples.

Past work has basically applied to adversarial training to image classification errands to text

Classification. Adversarial perturbations generally consist of making little adjustments to a lot of real valued input. For text classification, the info is discrete, and generally addressed as a progression of high-dimensional one-hot vectors. Since the arrangement of high-dimensional one-hot vectors doesn't concede infinitesimal perturbation, we characterize perturbation on word embeddings rather than discrete word inputs. Conventional adversarial trainings can be deciphered both as a regularization procedure (Szegedy et al., 2014; Goodfellow et al., 2015; Miyato et al., 2016) and as defense against an adversary who can supply malicious information sources (Szegedy et al., 2014; Goodfellow et al., 2015). Since the perturbed embedding doesn't guide to any word. We consequently propose this methodology solely as a method for regularizing a text classifier by

balancing out the classification function.

## II. Related work

Dropout (Srivastava et al., 2014) is a regularization strategy broadly utilized for some areas including text. There are some past works adding arbitrary noise to the data and hidden layer during preparation, to prevent overfitting (for example (Sietsma and Dow, 1991; Poole et al., 2013)). Nonetheless, in our trials and in past works (Miyato et al., 2016), preparing with adversarial perturbations outperformed the models with random perturbations.

For semi-supervised learning with neural networks, a typical methodology, particularly in the image space, is to prepare a generative model whose latent highlights might be utilized as features for classification (for example (Hinton et al., 2006; Maaløe et al., 2016)). These models presently accomplish best in class execution on the image space. Notwithstanding, these strategies require various extra hyperparameters with generative models, and the conditions under which the generative model will give great regulated learning execution are inadequately perceived. By comparison adversarial training requires only one hyperparameter.

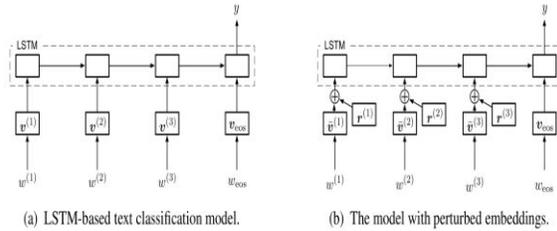
Adversarial look like some semi-supervised or transductive SVM approaches (Joachims, 1999; Chapelle and Zien, 2005; Collobert et al., 2006; Belkin et al., 2006) in that the two groups of techniques push the direction limit a long way from preparing models (or on account of transductive SVMs, test models). In any case, adversarial training strategies demand edges on the input space, while SVMs demand edges on the feature space characterized by the kernel function. This property permits adversarial training techniques to accomplish the models with a more adaptable capacity on the space where the edges are imposed. In our trials (Table 2, 4) and Miyato et al. (2016), adversarial training accomplishes preferable execution over SVM based techniques.

There have likewise been semi-supervised approaches applied to message grouping with both CNNs and RNNs. These methodologies use 'view-embeddings'(Johnson and Zhang, 2015b; 2016b) which utilize the window around a word to produce its embeddings. When these are utilized as a pretrained model for the classification model, they are found to improve generalization. These techniques and our strategy are integral as we showed that our

strategy improved from a recurrent pretrained language model.

### III. Model

We denote a grouping of  $T$  words as  $\{w(t)|t = 1, \dots, T\}$ , and a target as  $y$ . To change a discrete word contribution to a persistent vector, we characterize the word embeddings matrix  $V \in \mathbb{R}^{(K+1) \times D}$  where  $K$  is the quantity of words in the corpus and each column  $v_k$  relates to the word embeddings of the  $k$ -th word. Note that the  $(K + 1)$ -th word embeddings is utilized as an embedding of an 'end of sequence (eos)' token,  $v_{eos}$ . As a text classification model, we utilized a basic LSTM-based neural network model, which appeared in Figure. At time step  $t$ , the info is the discrete word  $w$ , and the comparing word embeddings is  $v$ . Furthermore we attempted the bidirectional LSTM design (Graves and Schmidhuber, 2005). For building the bidirectional LSTM model for text order, we add an extra LSTM on the reversed arrangement to the unidirectional LSTM model. The model then, at that point, predicts the name on the connected LSTM yields of the two closures of the succession. In adversarial training, we train the classifier to be robust to perturbations of the embeddings.



### IV. Adversarial Training

Adversarial training (Goodfellow et al., 2015) is a novel regularization technique for classifiers to improve robustness to little, roughly worst case scenario perturbations. Let's indicate  $x$  as the info and  $\theta$  as the boundaries of a classifier. When applied to a classifier, Adversarial training adds the accompanying term to the cost functions:

Where  $r$  is a perturbation on information and  $\theta$  is a steady set to the current boundaries of a classifier. The utilization of the steady duplicate  $\theta$  as opposed to  $\theta$  shows that the backpropagation calculation ought not be utilized to propagate gradients through the Adversarial model development process.

$$-\log p(y | x + r_{adv}; \theta) \text{ where } r_{adv} = \arg \min_{r, \|r\| \leq \epsilon} \log p(y | x + r; \theta)$$

This perturbation can be effectively computed utilizing backpropagation in neural organizations.

## V a) Recurrent Language Model Pre- Training

Following Dai and Le (2015), we instated the word embedding matrix and LSTM loads with a pre-prepared recurrent language model (Bengio et al., 2006; Mikolov et al., 2010) that was prepared on both labelled and unlabeled models. We utilized a unidirectional single-layer LSTM with 1024 secret units. The word embeddings dimensions D was 256 on IMDB and 512 on the other datasets. We utilized an inspected softmax loss with 1024 competitor samples for preparing. For the enhancement, we utilized the Adam Optimizer (Kingma and Ba, 2015), with bunch size 256, an underlying initial learning rate of 0.001, and a 0.9999 learning rate exponential rot factor at each preparation step. We prepared for 100,000 steps. To lessen runtime on GPU, we utilized shortened backpropagation up to 400 words from each finish of the arrangement. For regularization of the recurrent language model, we applied dropout (Srivastava et al., 2014) on the word embeddings layer with 0.5 dropout rate.

For the bidirectional LSTM model, we utilized 512 hidden units LSTM for both the standard request and switched request groupings, and we utilized 256 dimensional word embeddings which are imparted to both of the LSTMs. The other hyperparameters are equivalent to the unidirectional LSTM.

Pretraining with a repetitive language model was viable on classification performance on all the datasets we tried on thus our outcomes are with this pretraining.

## V b) Training Classification Models

After pre-training, we prepared the text classification model. Between the softmax layer for the objective y and the last yield of the LSTM, we added a hidden layer, which has measurement 30 on IMDB, Elec and Rotten Tomatoes, and 128 on DBpedia and RCV1. The initiation work on the secret layer was ReLU(Jarrett et al., 2009; Nair and Hinton, 2010; Glorot et al., 2011). For improvement, we again utilized the Adam Optimiser agent, with 0.0005 starting learning rate 0.9998 remarkable rot. Clump sizes are 64 on IMDB, Elec, RCV1, and 128 on DBpedia. For the Rotten Tomatoes

dataset, for each progression, we take a clump of size 64 for ascertaining the deficiency of the negative log-probability and 512 for figuring the deficiency of virtual ill-disposed preparation. Likewise for Rotten Tomatoes, we utilized writings with lengths  $T$  under 25 in the unlabeled dataset. We iterated 10,000 preparing steps on all datasets with the exception of IMDB and DBpedia, for which we utilized 15,000 and 20,000 preparing steps individually. We again applied angle cutting with the standard as 1.0 on every one of the boundaries aside from the word installing. We likewise utilized shortened backpropagation up to 400 words, and furthermore produced the ill-disposed and virtual ill-disposed irritation up to 400 words from each finish of the arrangement.

We tracked down the bidirectional LSTM to merge all the more gradually, so we iterated for 15,000 preparation steps when preparing the bidirectional LSTM characterization model.

For each dataset, we isolated the original training set into testing set and validation set, and we generally enhanced some hyperparameters; (model engineering, bunch size, preparing ventures) with the approval execution of the base model with installing dropout. For every

technique, we upgraded two scalar hyperparameters with the approval set. These were the dropout rate on the embeddings and the standard imperative  $\rho$  of adversarial trainings. We didn't do early halting with these strategies. The strategy with just pre-training and embedding dropout is utilized as the baseline. Because adversarial training only uses a labelled subset of the training data, it eventually overfits even the task of resisting adversarial perturbations.

## VI Results

We created a function that will save the average accuracy for each epoch in a separate automatically created text file. The final average accuracy will also be written in the text file similar to the figure.

A common misunderstanding is that adversarial training is the same as training on noisy examples.

Noise is a far weaker regularizer than adversarial perturbations because an average noise vector is approximately orthogonal to the cost gradient in high dimensional input spaces.

Adversarial perturbations are explicitly chosen to increase the cost in a consistent manner.

To demonstrate the superiority of adversarial training over noise addition, we include control experiments on each embedding in the

sequence that replaced adversarial perturbations with random perturbations from a multivariate Gaussian with scaled norm.

```
100
101 Epoch 9 of 10 -- average accuracy is 0.990 (train) -- average loss is 0.6
102 dev validation...
103 Average accuracy on dev is 0.868
104 Best dev acc till now is 0.889
105
106 Epoch 10 of 10 -- average accuracy is 0.991 (train) -- average loss is 0.
107 dev validation...
108 Average accuracy on dev is 0.870
109 Best dev acc till now is 0.889
110
111 Train Finished!
112
113 Load the best model params and test...
114 test validation...
115 Average accuracy on test is 0.871
```

## VII Conclusion

In our analyses, we found that adversarial training has good standardization performance on text classification. On all datasets, our proposed technique surpassed or was comparable to the cutting edge models. We likewise found that adversarial training improved classification as well as the nature of word embeddings. This approach could also be used for other general sequential tasks, such as for video or speech.(Sutskever et al., 2014), learning conveyed portrayals of words or paragraphs(Mikolov et al., 2013; Le and Mikolov, 2014)

## VIII References

Martin Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, et al. Tensorflow: Large-scale machine learning on heteroge-9 Published as a conference paper at ICLR 2017 heterogeneous distributed systems. arXiv preprint arXiv:1603.04467, 2016.

Mikhail Belkin, Partha Niyogi, and Vikas Sindhwani. Manifold regularization: A geometric framework for learning from labeled and unlabeled examples. *The Journal of Machine Learning Research*, 7(Nov):2399– 2434, 2006.

Yoshua Bengio, Holger Schwenk, Jean-Sébastien Senécal, Frédéric Morin, and Jean-Luc Gauvain. Neural probabilistic language models. In *Innovations in Machine Learning*, pp. 137–186. Springer, 2006.

Olivier Chapelle and Alexander Zien. Semi-supervised classification by low density separation. In *AISTATS*, 2005.

Ronan Collobert, Fabian Sinz, Jason Weston, and Léon Bottou. Large scale transductive svms. *Journal of Machine Learning Research*, 7(Aug):1687–1712, 2006.

Andrew M Dai and Quoc V Le. Semi-supervised sequence learning. In NIPS, 2015.

Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Deep sparse rectifier neural networks. In AISTATS, 2011.

Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In ICLR, 2015.

Alex Graves and Jürgen Schmidhuber. Framewise phoneme classification with bidirectional lstm and other neural network architectures. *Neural Networks*, 18(5):602–610, 2005.

Geoffrey E. Hinton, Simon Osindero, and Yee Whye Teh. A fast learning algorithm for deep belief nets. *Neural Computation*, 18:1527–1554, 2006.

Kevin Jarrett, Koray Kavukcuoglu, Marc’Aurelio Ranzato, and Yann LeCun. What is the best multi-stage architecture for object recognition? In ICCV, 2009.

Thorsten Joachims. Transductive inference for text classification using support vector machines. In ICML, 1999.

Rie Johnson and Tong Zhang. Effective use of word order for text

categorization with convolutional neural networks. NAACL HLT, 2015a.

Rie Johnson and Tong Zhang. Semi-supervised convolutional neural networks for text categorization via region embedding. In NIPS, 2015b.

Rie Johnson and Tong Zhang. Convolutional neural networks for text categorization: Shallow word-level vs. deep character-level. arXiv preprint arXiv:1609.00718, 2016a.

Rie Johnson and Tong Zhang. Supervised and semi-supervised text categorization using LSTM for region embeddings. In ICML, 2016b.

Yoon Kim. Convolutional neural networks for sentence classification. In EMNLP, 2014.

Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In ICLR, 2015.

Quoc V Le and Tomas Mikolov. Distributed representations of sentences and documents. In ICML, 2014.

Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo N Mendes, Sebastian Hellmann, Mohamed Morsey, Patrick van Kleef, Sören Auer, et al. Dbpedia—a large-scale, multilingual knowledge base

extracted from wikipedia. Semantic Web, 6(2):167–195, 2015.

David D Lewis, Yiming Yang, Tony G Rose, and Fan Li. Rcv1: A new benchmark collection for text categorization research. *The Journal of Machine Learning Research*, 5:361–397, 2004.

Lars Maaløe, Casper Kaae Sønderby, Søren Kaae Sønderby, and Ole Winther. Auxiliary deep generative models. In *ICML*, 2016.

Andrew L Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In *ACL: Human Language Technologies-Volume 1*, 2011.

Julian McAuley and Jure Leskovec. Hidden factors and hidden topics: understanding rating dimensions with review text. In *ACM conference on Recommender systems*, 2013.

Tomas Mikolov, Martin Karafiát, Lukáš Burget, Jan Cernocký, and Sanjeev Khudanpur. Recurrent neural network based language model. In *INTERSPEECH*, 2010. Published as a conference paper at *ICLR 2017*

Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words

and phrases and their compositionality. In *NIPS*, 2013.

Takeru Miyato, Shin-ichi Maeda, Masanori Koyama, Ken Nakae, and Shin Ishii. Distributional smoothing with virtual adversarial training. In *ICLR*, 2016.

Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *ICML*, 2010.

Bo Pang and Lillian Lee. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In *ACL*, 2005.

Ben Poole, Jascha Sohl-Dickstein, and Surya Ganguli. Analyzing noise in autoencoders and deep networks. In *Deep Learning Workshop on NIPS*, 2013.

J. Sietsma and R. Dow. Creating artificial neural networks that generalize. *Neural Networks*, 4(1), 1991.

Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1), 2014.



Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. In NIPS, 2014.