
Generating Text via Adversarial Training

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Generative Adversarial Networks (GANs) have achieved great success in generating
2 realistic synthetic real-valued data. However, the discrete output of language model
3 hinders the application of gradient-based GANs. In this paper we propose a generic
4 framework employing Long short-term Memory (LSTM) and convolutional neural
5 network (CNN) for adversarial training to generate realistic text. Instead of using
6 standard objective of GAN, we match the feature distribution when training the
7 generator. In addition, we use various techniques to pre-train the model and handle
8 discrete intermediate variables. We demonstrate that our model can generate
9 realistic sentence using adversarial training.

10 1 Introduction

11 Learning sentence representations is central to many natural language applications. The aim of a
12 model for such task is to learn fixed-length feature vectors that encode the semantic and syntactic
13 properties of sentences. One popular approach to learn a sentence model is by encoder-decoder
14 framework via recurrent neural network (RNN) [1]. Recently, several approaches has been proposed.
15 The skip-thought model of [2] describes an encoder-decoder model to reconstruct the surrounding
16 sentences of an input sentence, where both the encoder and decoder are modeled as RNN. The
17 sequence autoencoder of [3] is a simple variant of [2], in which the decoder is used to reconstruct the
18 input sentence itself.

19 These types of models enjoyed great success in many aspects of language modeling tasks, including
20 sentence classification and word prediction. However, autoencoder-based methods may fail when
21 generating realistic sentences from arbitrary latent representations [4]. The reason behind this is that
22 when mapping sentences to their hidden representations using an autoencoder, the representations of
23 these sentences may often occupy a small region in the hidden space. Thereby, most of regions in
24 the hidden space do not necessarily maps to a realistic sentence. Consequently, using a randomly
25 generated hidden representation from a prior distribution would usually leads to implausible sentences.
26 [4] attempt to use a variational auto-encoding framework to ameliorate this problem, however in
27 principle the posterior of the hidden variables would not cover the hidden space, rendering difficulties
28 to randomly produce sentences.

29 Another underlying challenge of generating realistic text relates to the nature of RNN. Suppose we
30 attempt to generate sentences from certain latent codes, the error will accumulate exponentially with
31 the length of the sentence. The first several words can be relatively reasonable, however the quality
32 of sentence deteriorates quickly. In addition, the lengths of sentences generated from random latent
33 representations could be difficult to control.

34 In this paper we propose a framework to generate realistic sentences with adversarial training scheme.
35 We adopted LSTM as generator and CNN as discriminator, and empirically evaluated various model
36 training techniques. Due to the nature of adversarial training, the generated text is discriminated with
37 the real text, thus the training is from a holistic perspective, rendering generated sentences to maintain

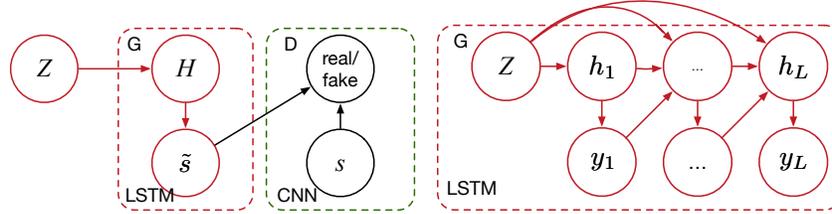


Figure 1: Left: Illustration of the textGAN model. The discriminator is a CNN, the sentence decoder is an LSTM. Right: the structure of LSTM model

38 high quality from the start to the end. As a related work, [5] proposed a sentence-level log-linear
 39 bag-of-words (BoW) model, where a BoW representation of an input sentence is used to predict
 40 adjacent sentences that are also represented as BoW. CNNs have recently achieved excellent results
 41 in various supervised natural language applications [6, 7, 8]. However, CNN-based *unsupervised*
 42 sentence modeling has previously not been explored.

43 We highlight that our model can: (i) learn a continuous hidden representation space to generate
 44 realistic text; (ii) generating high quality sentence in a holistic manner; (iii) we propose several
 45 training techniques for training such a model; (iv) be potentially applied to unsupervised disentangling
 46 learning and transferring literary styles.

47 2 Model description

48 2.1 TextGAN

49 Assume we are given a corpus $\mathcal{S} = \{s_1, \dots, s_n\}$, where n is the total number of sentence. Let w^t
 50 denote the t -th word in sentences s . Each word w^t is embedded into a k -dimensional word vector
 51 $\mathbf{x}_t = \mathbf{W}_e[w^t]$, where $\mathbf{W}_e \in \mathbb{R}^{k \times V}$ is a word embedding matrix (to be learned), V is the vocabulary
 52 size, and notation $[v]$ denotes the index for the v -th column of a matrix. Next we describe the model
 53 in three parts: CNN discriminator, LSTM generator and training strategies.

54 **CNN discriminator** The CNN architecture in [7, 9] is used for sentence encoding, which con-
 55 sists of a convolution layer and a max-pooling operation over the entire sentence for each feature
 56 map. A sentence of length T (padded where necessary) is represented as a matrix $\mathbf{X} \in \mathbb{R}^{k \times T}$, by
 57 concatenating its word embeddings as columns, *i.e.*, the t -th column of \mathbf{X} is \mathbf{x}_t .

58 A convolution operation involves a filter $\mathbf{W}_c \in \mathbb{R}^{k \times h}$, applied to a window of h words to produce
 59 a new feature. According to [9], we can induce one feature map $\mathbf{c} = f(\mathbf{X} * \mathbf{W}_c + \mathbf{b}) \in \mathbb{R}^{T-h+1}$,
 60 where $f(\cdot)$ is a nonlinear activation function such as the hyperbolic tangent used in our experiments,
 61 $\mathbf{b} \in \mathbb{R}^{T-h+1}$ is a bias vector, and $*$ denotes the convolutional operator. Convoluting the same filter
 62 with the h -gram at every position in the sentence allows the features to be extracted independently of
 63 their position in the sentence. We then apply a max-over-time pooling operation [9] to the feature
 64 map and take its maximum value, *i.e.*, $\hat{c} = \max\{\mathbf{c}\}$, as the feature corresponding to this particular
 65 filter. This pooling scheme tries to capture the most important feature, *i.e.*, the one with the highest
 66 value, for each feature map, effectively filtering out less informative compositions of words. Further,
 67 this pooling scheme also guarantees that the extracted features are independent of the length of the
 68 input sentence.

69 The above process describes how one feature is extracted from one filter. In practice, the model uses
 70 multiple filters with varying window sizes. Each filter can be considered as a linguistic feature detector
 71 that learns to recognize a specific class of n -grams (or h -grams, in the above notation). Assume we
 72 have m window sizes, and for each window size, we use d filters; then we obtain a md -dimensional
 73 vector \mathbf{f} to represent a sentence. Above this md -dimensional vector feature layer, we use a softmax
 74 layer to map the input sentence to an output $D(\mathbf{X}) \in [0, 1]$, represents the probability of \mathbf{X} is from
 75 the data distribution, rather than from adversarial generator.

76 There exist other CNN architectures in the literature [6, 8, 10]. We adopt the CNN model in [7, 9]
 77 due to its simplicity and excellent performance on classification. Empirically, we found that it can
 78 extract high-quality sentence representations in our models.

79 **LSTM generator** We now describe the LSTM decoder that translates a latent vector \mathbf{z} into the
 80 synthetic sentence \tilde{s} . The probability of a length- T sentence \tilde{s} given the encoded feature vector \mathbf{z} is
 81 defined as

$$p(\tilde{s}|\mathbf{z}) = p(w^1|\mathbf{z}) \prod_{t=2}^T p(w^t|w^{<t}, \mathbf{z}) \quad (1)$$

82 Specifically, we generate the first word w^1 from \mathbf{z} , with $p(w^1|\mathbf{z}) = \text{argmax}(\mathbf{V}\mathbf{h}_1)$, where $\mathbf{h}_1 =$
 83 $\tanh(\mathbf{C}\mathbf{z})$. Bias terms are omitted for simplicity. All other words in the sentence are then sequentially
 84 generated using the RNN, until the end-sentence symbol is generated. Each conditional $p(w^t|w^{<t}, \mathbf{z})$,
 85 where $<t = \{1, \dots, t-1\}$, is specified as $\text{argmax}(\mathbf{V}\mathbf{h}_t)$, where \mathbf{h}_t , the hidden units, are recursively
 86 updated through $\mathbf{h}_t = \mathcal{H}(\mathbf{y}_{t-1}, \mathbf{h}_{t-1}, \mathbf{z})$. The LSTM input \mathbf{y}_{t-1} for t -th step is the embedding vector
 87 of the previous word that maximize the w^{t-1}

$$\mathbf{y}_{t-1} = \mathbf{W}_e[w^{t-1}]. \quad (2)$$

88 \mathbf{V} is a weight matrix used for computing a distribution over words. The synthetic sentence is obtained
 89 by $s = [\text{argmax}(w^1), \dots, \text{argmax}(w^T)]$.

90 The transition function $\mathcal{H}(\cdot)$ is implemented with an LSTM [1]. Each LSTM unit has a cell containing
 91 a state \mathbf{c}_t at time t . Reading or writing the memory unit is controlled through sigmoid gates, namely,
 92 input gate \mathbf{i}_t , forget gate \mathbf{f}_t , and output gate \mathbf{o}_t . The hidden units \mathbf{h}_t are updated as follows:

$$\mathbf{i}_t = \sigma(\mathbf{W}_i\mathbf{y}_{t-1} + \mathbf{U}_i\mathbf{h}_{t-1} + \mathbf{C}_i\mathbf{z}) \quad \mathbf{f}_t = \sigma(\mathbf{W}_f\mathbf{y}_{t-1} + \mathbf{U}_f\mathbf{h}_{t-1} + \mathbf{C}_f\mathbf{z}) \quad (3)$$

$$\mathbf{o}_t = \sigma(\mathbf{W}_o\mathbf{y}_{t-1} + \mathbf{U}_o\mathbf{h}_{t-1} + \mathbf{C}_o\mathbf{z}) \quad \tilde{\mathbf{c}}_t = \tanh(\mathbf{W}_c\mathbf{y}_{t-1} + \mathbf{U}_c\mathbf{h}_{t-1} + \mathbf{C}_c\mathbf{z}) \quad (4)$$

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tilde{\mathbf{c}}_t \quad \mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t) \quad (5)$$

93 where $\sigma(\cdot)$ denotes the logistic sigmoid function, and \odot represents the element-wise multiply operator
 94 (Hadamard product). $\mathbf{W}_{\{i,f,o,c\}}$, $\mathbf{U}_{\{i,f,o,c\}}$, $\mathbf{C}_{\{i,f,o,c\}}$, \mathbf{V} and \mathbf{C} are the set of parameters. Note that
 95 \mathbf{z} is used as an explicit input at each time step of the LSTM to guide the generation of \tilde{s} .

96 **Training techniques** Given the sentence corpus \mathcal{S} , instead of directly minimizing the objective
 97 function from standard GAN [11], we adopted an approach similar to feature matching [12]. The
 98 iterative optimization schemes consists of two steps:

$$\text{minimizing: } L_D = -\mathbb{E}_{s \sim \mathcal{S}} \log D(s) - \mathbb{E}_{z \sim p_z(z)} \log[1 - D(G(z))] \quad (6)$$

$$\text{minimizing: } L_G = \text{tr}(\boldsymbol{\Sigma}_s^{-1}\boldsymbol{\Sigma}_r + \boldsymbol{\Sigma}_r^{-1}\boldsymbol{\Sigma}_s) + (\boldsymbol{\mu}_s - \boldsymbol{\mu}_r)^T (\boldsymbol{\Sigma}_s^{-1} + \boldsymbol{\Sigma}_r^{-1})(\boldsymbol{\mu}_s - \boldsymbol{\mu}_r) \quad (7)$$

99 where $\boldsymbol{\Sigma}_s, \boldsymbol{\Sigma}_r$ represents the covariance matrices of real and synthetic sentence feature vector $\mathbf{f}_s, \mathbf{f}_r$,
 100 respectively. $\boldsymbol{\mu}_s, \boldsymbol{\mu}_r$ denote the mean vector of $\mathbf{f}_s, \mathbf{f}_r$, respectively. $\boldsymbol{\Sigma}_s, \boldsymbol{\Sigma}_r, \boldsymbol{\mu}_s$ and $\boldsymbol{\mu}_r$ are empirically
 101 estimated on minibatch. By setting $\boldsymbol{\Sigma}_s = \boldsymbol{\Sigma}_r = \mathbf{I}$, this reduces to the feature matching technique
 102 from [12]. Note that this second loss L_G is the symmetric KL divergence between two multivariate
 103 Gaussian distribution $\mathcal{N}(\boldsymbol{\mu}_r, \boldsymbol{\Sigma}_r)$ and $\mathcal{N}(\boldsymbol{\mu}_s, \boldsymbol{\Sigma}_s)$. Instead of capturing the first moment similarity,
 104 we cast more stringent criteria of matching the feature covariance of real and synthetic data. Despite
 105 the fact that the feature vector is not necessarily Gaussian distributed, empirically this loss (7) works
 106 well. Intuitively, this technique provides stronger signal for modifying the generator to make the
 107 synthetic data more realistic.

108 To train the generator G which contains discrete variable, direct application of gradient estimation
 109 would fails. Score function based algorithms, such as REINFORCE [13] algorithm, obtains unbiased
 110 gradient estimation for discrete variables by using Monte Carlo estimation. However the variance of
 111 the gradient estimation could be large [14]. Here we consider a soft-argmax function when performing
 112 the inference as an approximation to the (2):

$$\mathbf{y}_{t-1} = \mathbf{W}_e \text{softmax}(\mathbf{V}\mathbf{h}_{t-1} \odot L). \quad (8)$$

113 where \odot represents element-wise product. When $L \rightarrow \infty$, this approximation would becomes (2).
 114 The w^{t-1} is constrained to be either 0 or have an absolute value greater than 1.

115 Previous literature [11, 12] has discussed the fundamental difficulty in training GAN model using
 116 gradient-based method. In general, gradient descent optimization scheme would fail to converges
 117 to the Nash equilibrium by moving along orbit trajectory among saddle points. Presumably, a
 118 good initialization would encourage convergence by reducing the orbit movement. To achieve good

119 initialization, we initialize the LSTM parameters for the generator by pre-training a standard auto-
 120 encoder LSTM model. For the discriminator, we use a confusion training strategy. For each sentence
 121 in the corpus, we randomly swap two words to construct a tweaked counterpart sentence. The
 122 discriminator is pre-trained to classify the tweaked sentences from the true sentence. The swapping
 123 operation, rather than adding/deleting is used because we attempt to train the CNN discriminator to
 124 learn the *sentence structure* feature, rather than absence/presence of certain words (*words presence*
 125 feature). This strategy helps avoid the generator to produce repeated "realistic" words that rewards
 126 the *words presence* feature to achieve a higher score of being classified as from real data.

127 3 Experiments

128 Our model is trained using a combination of two datasets: 1). BookCorpus dataset [15], which
 129 consists of 70 million sentences from over 7000 books; 2). ArXiv dataset, which consists of 5 million
 130 sentences from various subjects, scrawled from arXiv website. The purpose of including two different
 131 corpus is to investigate whether the model can generate sentences that integrates both scientific
 132 writing style and vocabulary with informal writing style and vocabulary. We randomly choose 1
 133 million sentences from BookCorpus and 1 million sentences from arXiv dataset to construct our
 134 training dataset.

135 We train the generator and discriminator iteratively. Given that the LSTM generator typically
 136 involves more parameters and is more difficult to train than the CNN discriminator, we perform one
 137 optimization step for the discriminator for every 5 steps for the generator. Both of the generator and
 138 discriminator are per-trained with the strategy described in Section 2.

139 For the CNN encoder, we employ filter windows (h) of sizes $\{3,4,5\}$ with 300 feature maps each,
 140 hence each sentence is represented as a 900-dimensional vector. For both, the LSTM sentence decoder
 141 and paragraph generator, we use one hidden layer of 500 units.

142 Gradients are clipped if the norm of the parameter vector exceeds 5 [16]. The Adam algorithm [17]
 143 with learning rate 1×10^{-4} for both discriminator and generator is utilized for optimization. We
 144 use mini-batches of size 128. All experiments are implemented in Theano [18], using a NVIDIA
 145 GeForce GTX TITAN X GPU with 12GB memory. The model was trained for roughly one week.

146 We first examine how well the generator can produce similar feature distributions to mimic the
 147 real data. We calculate the empirical expectation and covariance of 900 CNN top layer features
 148 over a minibatch of 128 real sentences and 128 synthetic sentences. As shown in Figure 2, the
 149 expectation of these 900 features from synthetic sentences matches well with the feature expectation
 150 from the real sentences. The covariances of 900 features against the first feature (*i.e.*, $Cov(f_i, f_1)$)
 151 also demonstrate consistency between the real sentences and synthetic sentences.

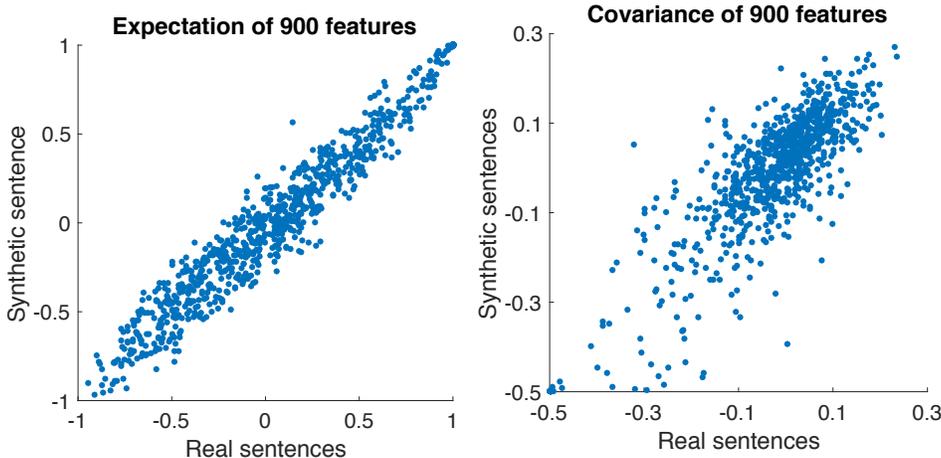


Figure 2: Left: the scatter plot of the expectations of features from real data against the expectations of features from synthetic data. Right: the scatter plot of the covariances of features (*v.s.* 1st feature, $Cov(f_i, f_1)$) from real data against the covariances of features from synthetic data.

Table 1: Intermediate sentences produced from linear transition between two random points in the latent space. Each sentence is generated from a latent point on a linear grid with equidistance

A	we show the efficacy of our new solvers , making it up to identify the optimal random vector .
-	we made the price of such new solvers and reduce its time series pairs .
-	we realized and minimize the best of such tasks in terms of multidimensional high dimensional kernels .
-	we realized and minimize the best of such tasks are invariant on such a sparse linear program .
-	we realized and minimize the price of <UNK> 's plans , by drawing up a constant .
-	we realized the price used to pay , whatever they are .
-	i realized mrs. <UNK> was surprised that most of them got safe .
B	dylan realized the distance of the walls to them .

152 We further empirically evaluate whether the latent variable space can densely encode sentences with
 153 appropriate grammar structure, and to visualize the transition from sentence to sentence. We construct
 154 a linear path between two randomly selected points in the latent space, and generate the intermediate
 155 sentences along the linear trajectory. The results are presented in Table 1. Generally the produced
 156 sentence is grammatically correct and semantically reasonable. The transition demonstrate smooth-
 157 ness and interpretability, however the wording choices and sentence structure showed dramatical
 158 changes in some regions in the latent space. Presumably the local “transition smoothness” can varies
 159 from region to region.

160 Note that the generated text demonstrates some entangling of both scientific style of writing and
 161 informal style of writing, because the training is on the combination of two corpus. The generator
 162 can conjure novel sentence by leveraging the grammatical rule and property of words.

163 We observed that the discriminator can still sufficiently distinguish the synthetic sentences from
 164 the real ones (the probability to predict synthetic data as real is around 0.08), even if the synthetic
 165 sentences seems to perserve reasonable grammatical structure and use proper wording. Probably the
 166 CNN is able to sophisticatedly characterize the semantic meaning and differentiate sentences in a
 167 holistic perspective, while the generator may stuck into a sweet point where any slight modification
 168 would render a higher loss (7) for the generator.

169 4 Conclusion

170 We presented a text generation model via adversarial training and discussed several techniques for
 171 training such a model. We demonstrate that the proposed model can produce realistic sentences by
 172 mimicking the input real sentences, and the learned latent representation space can continuously
 173 encode plausible sentences.

174 In future work, we attempt to disentangle the latent representations for different writing style in
 175 an unsupervised manner. This would enable a smooth lexical and grammatical transition between
 176 different writing styles. The model also will be fine-tuned in order to achieve more compelling results,
 177 such as improved sentences with better semantical interpretation. Furthermore, a more comprehensive
 178 quantitative comparison will be performed.

179 **References**

- 180 [1] S. Hochreiter and J. Schmidhuber. Long short-term memory. In *Neural computation*, 1997.
- 181 [2] R. Kiros, Y. Zhu, R. Salakhutdinov, R. Zemel, R. Urtasun, A. Torralba, and S. Fidler. Skip-thought vectors.
- 182 In *NIPS*, 2015.
- 183 [3] A. Dai and Q. Le. Semi-supervised sequence learning. In *NIPS*, 2015.
- 184 [4] Samuel R Bowman, Luke Vilnis, Oriol Vinyals, Andrew M Dai, Rafal Jozefowicz, and Samy Bengio. Generating sentences from a continuous space. *arXiv preprint arXiv:1511.06349*, 2015.
- 185 [5] F. Hill, K. Cho, and A. Korhonen. Learning distributed representations of sentences from unlabelled data.
- 186 In *NAACL*, 2016.
- 187 [6] N. Kalchbrenner, E. Grefenstette, and P. Blunsom. A convolutional neural network for modelling sentences.
- 188 In *ACL*, 2014.
- 189 [7] Y. Kim. Convolutional neural networks for sentence classification. In *EMNLP*, 2014.
- 190 [8] B. Hu, Z. Lu, H. Li, and Q. Chen. Convolutional neural network architectures for matching natural
- 191 language sentences. In *NIPS*, 2014.
- 192 [9] R. Collobert, J. Weston, L. Bottou, M. Karlen, K. Kavukcuoglu, and P. Kuksa. Natural language processing
- 193 (almost) from scratch. In *JMLR*, 2011.
- 194 [10] R. Johnson and T. Zhang. Effective use of word order for text categorization with convolutional neural
- 195 networks. In *NAACL HLT*, 2015.
- 196 [11] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron
- 197 Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing*
- 198 *Systems*, pages 2672–2680, 2014.
- 199 [12] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved
- 200 techniques for training gans. *arXiv preprint arXiv:1606.03498*, 2016.
- 201 [13] Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement
- 202 learning. *Machine learning*, 8(3-4):229–256, 1992.
- 203 [14] Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of
- 204 discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016.
- 205 [15] Y. Zhu, R. Kiros, R. Zemel, R. Salakhutdinov, R. Urtasun, A. Torralba, and S. Fidler. Aligning books and
- 206 movies: Towards story-like visual explanations by watching movies and reading books. In *ICCV*, 2015.
- 207 [16] I. Sutskever, O. Vinyals, and Q. Le. Sequence to sequence learning with neural networks. In *NIPS*, 2014.
- 208 [17] D. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- 209 [18] F. Bastien, P. Lamblin, R. Pascanu, J. Bergstra, I. Goodfellow, A. Bergeron, N. Bouchard, D. Warde-Farley,
- 210 and Y. Bengio. Theano: new features and speed improvements. *arXiv:1211.5590*, 2012.
- 211