
Auto-encoding Documents for Topic Modeling with L-2 sparsity regularization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We propose a novel yet simple neural network architecture for topic modeling.
2 The method is based on training an autoencoder structure where the bottleneck
3 represents the space of the topics distributions and the decoder output represents
4 the space of the word distribution over the topics. We also exploit an auxiliary
5 decoder to prevent mode collapsing in our model. A key feature for an effective
6 topic modeling method is having sparse topic and word distributions, where there
7 is a trade-off between the sparsity level of topics and words. This feature is im-
8 plemented in our model by an L-2 regularization and the model hyperparameters
9 take care of the trade-off. We show in our experiments that our model achieves
10 competitive results compared to the state-of-the-art deep models for topic model-
11 ing, despite its simple architecture and training procedure. The *New York Times*
12 and *20 Newsgroups* datasets are used in the experiments.

13 1 Introduction

14 Topic models are among the key models in Natural Language Processing (NLP) that aim to represent
15 a large body of text using only a few concepts or topics, on a completely unsupervised basis. Topic
16 modeling has found its application in many different areas including bioinformatics [13], computer
17 vision [5], recommendation systems [8, 15], etc. Latent semantic indexing (LSI) [4] and proba-
18 biliestic latent semantic indexing (PLSI) [6] are among the oldest algorithms, but Latent Dirichlet
19 Allocation (LDA) [3] is the most widely used algorithm for topic modeling and most of the success-
20 ful algorithms in this area are variants of LDA. The main challenge in training the LDA model is
21 its relatively complicated inference model, which makes finding the true posterior a hard task and
22 therefore LDA-based topic modeling algorithms rely on approximating methods. There are many
23 ways to approximate the inference in LDA, e.g. Mean field variational methods [2], variational in-
24 ference [2], expectation propagation [10], collapsed Gibbs sampling [11], factorization inference [1]
25 etc.

26 With the advances in deep learning, there has been some efforts to implement LDA and its variants
27 using neural networks. Recent variational autoencoding (VAE) [7, 12] model has paved the way for
28 approximating posteriors using a black box neural network model. However, VAE works best with
29 the posterior approximating distributions for which sampling can be done using reparameterization
30 trick, e.g. Gaussian distribution. Therefore, successful deep models that use VAE framework to
31 implement LDA either directly replace the Dirichlet distribution with a Gaussian distribution [9]
32 or approximate the Dirichlet using a combination of Gaussians [14]. On the other hand, imposing
33 sparsity constraint, which is an important aspect of LDA and Dirichlet distribution, is not trivial in
34 these models.

35 In this work we propose a new deep topic modeling algorithm, based on training an autoencoder,
36 which takes as its input the distribution of words in each document and represents the topics in its

37 bottleneck layer. Our model keeps the most important properties of LDA, while having a simple
38 network structure and an easy training procedure. We maintain the sparsity property of the Dirichlet
39 distribution by imposing an L-2 regularization on the softmax layers of the neural network. We also
40 resolve the mode collapsing issue of the model by adding an auxiliary decoder network that makes
41 the separation of the representations in the latent space easier.

42 2 Model Description

43 In this section we first briefly describe LDA and its properties and then explain our model and
44 how it resembles the LDA properties without having the difficulties of dealing with the intractable
45 posterior.

46 2.1 Following LDA properties

47 Let's assume we have a set of D documents with vocabulary size N . Each document is represented
48 by a vector \mathbf{x} . Also there are K topics with different distributions over the words, denoted by β_1 to
49 β_K . The LDA generative process is as follows:

```
50     For each document  $\mathbf{x}$ 
51         Draw topic distribution  $\theta \sim \text{Dirichlet}(\alpha)$ 
52     For each word in the document
53         Sample a topic  $\mathbf{z} \sim \text{Multinomial}(\theta)$ 
54         Sample a word in position  $n$ ,  $w_n \sim \text{Multinomial}(\beta_{\mathbf{z}})$ 
```

55 Two important objectives that LDA implicitly tries to achieve and they make this model suitable for
56 topic modelling are:

- 57 • Distribution of the topics for each document is sparse, therefore each document can be
58 represented by a few topics
- 59 • Distribution of the words for each topic is sparse, therefore each topic can be represented
60 by a few words.

61 There is a trade-off between these two objectives. If a document is represented using only a few
62 topics, then number of words with high probability in those topics should be large, and if topics
63 are represented using only a few words then we need a large number of topics to cover the words
64 in the document. The sparsity of the distributions is a property of the Dirichlet distribution that is
65 controlled by its concentration parameters. Also, based on LDA, the distribution of the words in a
66 document is a mixture of multinomials.

67 In our model we follow the main principals of the LDA algorithm, i.e. sparse distributions for the
68 topics and words and the final distribution of the words in a document is a mixture of multinomial.
69 On the other hand, we try to avoid the difficulties of training the LDA model. Since our downstream
70 task is finding topics in the documents, and not generating new documents, we do not need to
71 learn the true posterior probability, or find ways to approximate it. Therefore we leave the latent
72 representation unconstrained with regard to its distribution.

73 We first encode the documents to the topic space \mathcal{Z} using $f_{\text{topic}}(\mathbf{x}; \phi)$, which is implemented by
74 a neural network with parameter set ϕ . To make sure \mathcal{Z} is a probability space we use a softmax
75 layer at the last layer of this network. Also, K vectors, β_1 to β_K , with softmax activation represent
76 the words distributions in the topics, each of them is a multinomial distribution. A mixture of

77 multinomials, i.e. $\tilde{\mathbf{x}} = f_{\beta}(\mathbf{z}) = \sum_{k=1}^K z_k \beta_k$, will be a reconstruction of the input vector \mathbf{x} . We
78 intentionally do not use a matrix multiplication notation so that we can explain the constraints on
79 β_k 's in a simpler and more explicit way.

80 To make both topic and words distributions sparse, we impose an L-2 norm constraint on them.
81 Maximizing the L-2 norm over a positive, sum-to-one vector, concentrates the probability mass over
82 a few number of elements. This way we are keeping the most informative words of each topic,

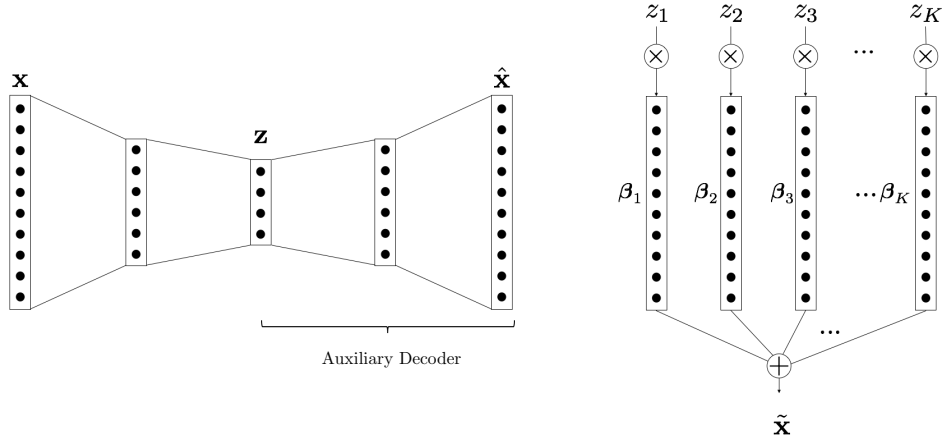


Figure 1: Networks of the model

83 because sparsity means we are minimizing the entropy of the distributions, Therefore the objective
 84 of the algorithm will be as follows:

$$\min_{\phi, \beta_k} D(\mathbf{x}, \tilde{\mathbf{x}}) - \gamma \|\mathbf{z}\|^2 - \eta \sum_{k=1}^K \|\beta_k\|^2, \quad (1)$$

85 where distance D is the cross entropy, and γ and η are hyperparameters of the model. The trade-off
 86 of the sparsity in topics and words distributions can be controlled by tuning γ and η .

87 2.2 Mode Collapsing

88 We observed that training the model using Eq. (1) causes mode collapsing, in the sense that only a
 89 very few topics will have meaningful words in them and the rest of the topics have high probability
 90 over some random words. Also, all the probability mass of the topics distribution for all of the
 91 documents are concentrated on those specific topics. In other words, all the documents are encoded
 92 to the same set of topics and the model cannot capture the variations in the documents. We believe
 93 this is due to the fact that $f_\beta(\mathbf{z})$ is not a powerful function for backpropagating the error signal from
 94 the output to the previous layers of the network. To resolve this issue and produce a richer \mathcal{Z} space,
 95 we attach an auxiliary decoder to the latent representation, which we call it $f_{\text{AUX}}(\mathbf{z}; \varphi)$ and it is a
 96 neural network with parameter set φ . The output of this decoder, denoted by $\hat{\mathbf{x}}$, also reconstructs the
 97 input document. Our observations show that by adding this decoder we can separate the documents'
 98 representations in the latent space.

99 In both topic and word level, instead of sampling, we consider \mathbf{z} and β_k 's (for all $k \in \{1, 2, \dots, K\}$)
 100 as a normalized typical set of the distribution \mathbf{z} and β_k 's. This is to avoid sampling from the multi-
 101 nomial distribution for which there is no easy way, e.g. reparameterization trick in [7] for Gaussian
 102 family, to backpropagate the error for training the neural networks. Therefore the overall objective
 103 of our model is:

$$\begin{aligned} \min_{\phi, \varphi, \beta_k} \quad & D(\mathbf{x}, \tilde{\mathbf{x}}) + \lambda D(\mathbf{x}, \hat{\mathbf{x}}) - \gamma \|\mathbf{z}\|^2 - \eta \sum_{k=1}^K \|\beta_k\|^2 \\ \text{s.t} \quad & \sum_{n=1}^N \beta_{kn} = 1 \quad \forall k \in \{1, 2, \dots, K\} \\ & \sum_{k=1}^K z_k = 1, \end{aligned} \quad (2)$$

104 where λ is another hyperparameter of the model that controls the role of the auxiliary decoder in
 105 training. Figure 1 shows the structure of the networks.

106 3 Experiments

107 In this section we compare the performance of the proposed algorithm with LDA with collapsed
 108 Gibbs, and two deep models, i.e. Neural Variational Document Model (NVDM) and ProdLDA
 109 algorithms in [14, 9]. Although comparing different topic modeling results qualitatively is a hard
 110 task, we follow standard metrics for such comparisons. The comparison is made based on *topic*
 111 *coherence* (higher is better) and *perplexity score* (lower is better) of the results.

112 3.1 New York Times

113 This dataset consists of $D = 8,447$ documents with vocabulary size $N = 3,012$ words. We down-
 114 loaded the dataset from this git repository https://github.com/moorissa/nmf_nyt.
 115 This dataset doesn't need a preprocessing phase, as the common words and stop words has already
 116 been removed from it. We try performing topic modeling using 25 and 50 topics for this dataset (CG
 117 in the tables mean Collapsed Gibbs and the best results are indicated by bold symbols).

Number of Topics	LDA with CG	ProdLDA	NVDM	Our Model
25	0.26	0.30	0.25	0.32
50	0.23	0.30	0.21	0.29

Table 1: Topic Coherence for the *New York Times* dataset

Number of Topics	LDA with CG	ProdLDA	NVDM	Our Model
25	781	910	842	762
50	770	930	892	751

Table 2: Perplexity for the *New York Times* dataset

118 In this experiment, for $K = 25$ topics the value of hyperparameters are: $\lambda = 0.1$, $\gamma = 0.1$, and
 119 $\eta = 0.001$. For $K = 50$ topics these values are: $\lambda = 0.1$, $\gamma = 0.05$, and $\eta = 0.001$.

120 3.2 20 Newsgroups

121 The 20 Newsgroups has $D = 11,000$ training documents. We follow the same preprocessing in
 122 [14], tokenization, removing some of the non UTF-8 characters and English stop word removal.
 123 These are all done using `scikit-learn` package. After this preprocessing the vocabulary size is
 124 $N = 2,000$. For this dataset we try training the models with 50 and 200 topics.

Number of Topics	LDA with CG	ProdLDA	NVDM	Our Model
50	0.18	0.23	0.10	0.25
200	0.14	0.19	0.08	0.18

Table 3: Topic Coherence for the *20 Newsgroups* dataset

Number of Topics	LDA with CG	ProdLDA	NVDM	Our Model
50	737	1180	830	795
200	690	1139	842	806

Table 4: Perplexity for the *20 Newsgroups* dataset

125 In this experiment, for $K = 50$ topics the value of hyperparameters are: $\lambda = 0.5$, $\gamma = 0.1$, and
 126 $\eta = 0.001$. For $K = 200$ topics these values are: $\lambda = 0.4$, $\gamma = 0.01$, and $\eta = 0.001$.

127 We can see that for both datasets, our algorithm achieves competitive results with some of the state-
 128 of-the-art deep models for topic modeling. ProdLDA shows better performance in term of the topic
 129 coherence when the number of topics gets large. However, for lower number of topics our model
 130 outperforms all othe algorithms. We can also see that our results is better than ProdLDA in terms of
 131 perplexity for both datasets, although LDA with collapsed Gibbs yilelds the best results for the *20*
 132 *Newsgroups* dataset. Some random topics and the 10 highest probable words in them are shows in
 133 tables 5 and 6 in appendix.

References

- 134
- 135 [1] S. Arora, R. Ge, Y. Halpern, D. Mimno, A. Moitra, D. Sontag, Y. Wu, and M. Zhu. A practical
136 algorithm for topic modeling with provable guarantees. In *International Conference on*
137 *Machine Learning*, pages 280–288, 2013.
- 138 [2] D. M. Blei, M. I. Jordan, et al. Variational inference for dirichlet process mixtures. *Bayesian*
139 *analysis*, 1(1):121–143, 2006.
- 140 [3] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent dirichlet allocation. *Journal of machine Learning*
141 *research*, 3(Jan):993–1022, 2003.
- 142 [4] S. Deerwester, S. T. Dumais, G. W. Furnas, T. K. Landauer, and R. Harshman. Indexing by
143 latent semantic analysis. *Journal of the American society for information science*, 41(6):391–
144 407, 1990.
- 145 [5] L. Fei-Fei and P. Perona. A bayesian hierarchical model for learning natural scene categories.
146 In *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Con-*
147 *ference on*, volume 2, pages 524–531. IEEE, 2005.
- 148 [6] T. Hofmann. Probabilistic latent semantic analysis. In *Proceedings of the Fifteenth conference*
149 *on Uncertainty in artificial intelligence*, pages 289–296. Morgan Kaufmann Publishers Inc.,
150 1999.
- 151 [7] D. Kingma and M. Welling. Auto-encoding variational Bayes. In *Proceedings of ICLR*, 2014.
- 152 [8] X. Li and J. She. Collaborative variational autoencoder for recommender systems. In *Proceed-*
153 *ings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data*
154 *Mining*, pages 305–314. ACM, 2017.
- 155 [9] Y. Miao, L. Yu, and P. Blunsom. Neural variational inference for text processing. In *Internat-*
156 *ional Conference on Machine Learning*, pages 1727–1736, 2016.
- 157 [10] T. P. Minka. *A family of algorithms for approximate Bayesian inference*. PhD thesis, Mas-
158 sachusetts Institute of Technology, 2001.
- 159 [11] I. Porteous, D. Newman, A. Ihler, A. Asuncion, P. Smyth, and M. Welling. Fast collapsed gibbs
160 sampling for latent dirichlet allocation. In *Proceedings of the 14th ACM SIGKDD international*
161 *conference on Knowledge discovery and data mining*, pages 569–577. ACM, 2008.
- 162 [12] D. Rezende, S. Mohamed, and D. Wierstra. Stochastic backpropagation and approximate in-
163 ference in deep generative models. In *Proceedings of the 31st International Conference on*
164 *Machine Learning*, pages 1278–1286, 2014.
- 165 [13] S. Rogers, M. Girolami, C. Campbell, and R. Breitling. The latent process decomposition of
166 cdna microarray data sets. *IEEE/ACM transactions on computational biology and bioinfor-*
167 *matics*, 2(2):143–156, 2005.
- 168 [14] A. Srivastava and C. Sutton. Autoencoding variational inference for topic models. *arXiv*
169 *preprint arXiv:1703.01488*, 2017.
- 170 [15] H. Wang, N. Wang, and D.-Y. Yeung. Collaborative deep learning for recommender systems.
171 In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery*
172 *and Data Mining*, pages 1235–1244. ACM, 2015.

173 **A Some results**

174 Here we present some randomly selected topics with their top 10 words with the highest probabilities
 175 for the two datasets.

176 **New York Times**

Topic 1	Topic 2	Topic 3	Topic 4	Topic 5
building	police	percent	preseident	military
city	man	price	executive	plane
build	kill	market	chief	flight
open	arrest	rate	director	fly
house	officer	rise	name	mission
project	pficial	fall	vice	attack
site	fire	sale	chairman	airline
street	charge	share	company	security
space	shoot	report	agency	air
home	death	increase	management	pilot

Table 5: Topc 10 words in randomly selected topics

177 **20 Newsgroups**

Topic 1	Topic 2	Topic 3	Topic 4	Topic 5
windows	jews	goverment	space	game
card	israel	people	station	team
drive	israeli	federal	moon	games
monitor	law	american	orbit	play
pc	anti	political	launch	win
drivers	jewish	civil	nasa	season
disk	arab	war	shuttle	players
mouse	religious	society	developed	hit
mac	killed	national	object	baseball
printer	muslim	majority	lunar	league

Table 6: Topc 10 words in randomly selected topics