


# Topic Modelling

 ... reviewing machine learning literature about topic modelling

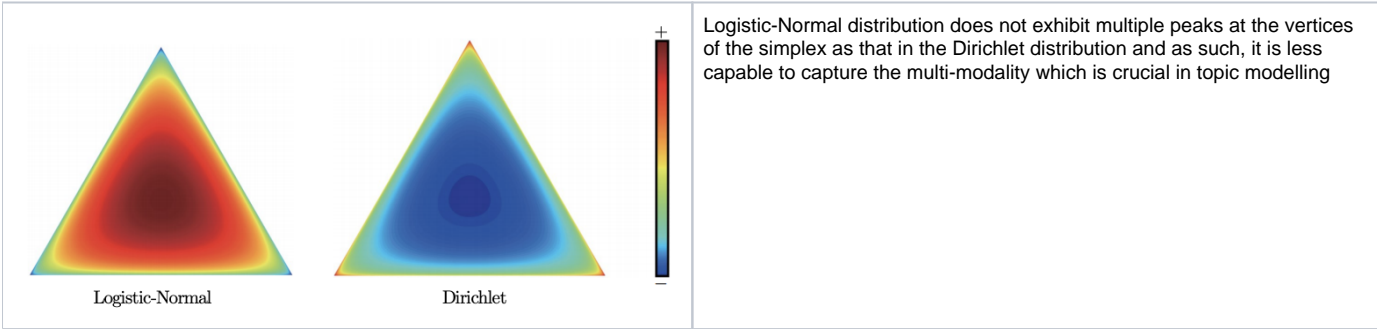
topic modelling	topic classification
<ul style="list-style-type: none"><li>unsupervised machine learning</li><li>automatically analyses text data to determine cluster words for a set of documents</li><li>If you don't have a lot of time to analyse texts, or you're not looking for a fine-grained analysis and just want to figure out what topics a bunch of texts are talking about</li></ul>	<ul style="list-style-type: none"><li>supervised machine learning techniques</li><li>automatically label a review with predefined topic tags rather than inferring what similarity cluster the review belongs to</li><li>if you have a list of predefined topics for a set of texts and want to label them automatically without having to read each one, as well as gain accurate insights</li></ul>
<p>example)</p> <p><i>"The nice thing about Eventbrite is that it's <b>free to use</b> as long as you're not <b>charging</b> for the event. There is a <b>fee</b> if you are <b>charging</b> for the event – <b>2.5% plus a \$0.99 transaction fee</b>."</i></p> <p>By identifying words and expressions such as <i>free to use, fee, charging, 2.5% plus 99 cents transaction fee</i>, topic modelling can group this review with other reviews that talk about similar things (these may or may not be about pricing).</p>	<p>example)</p> <p><i>"We have the <b>gold level plan</b> and use it for everything, <b>love the features!</b> It is one of the <b>best bang for buck</b> possible."</i></p> <p>A topic classification model that's been trained to understand these expressions (gold level plan, love the features, and best bang for buck) would be able to tag this review as topics <i>Features</i> and <i>Price</i>.</p>
topic modelling algorithms churn out collections of expressions and words that it thinks are related, leaving you to figure out what these relations mean, while topic classification delivers neatly packaged topics, with labels such as <i>Price</i> , and <i>Features</i> , eliminating any guesswork.	

## topic modelling method

- Latent Semantic Analysis (LSA)
  - one of the NLP techniques for analysis of semantics and introduced in 2005
  - trying to dig out some meaning out of a corpus of text
  - unsupervised approach
  - helpful technique in the reduction of dimensions of the topic modelling
  - main concept is to group together all the words that have a similar meaning
  - tf-idf
    - term-frequency (TF) is a number of times keyword appears in a single document divided by the total number of words in that document
$$TF(t,d) = \frac{\text{Number of times term "t" appears in a document}}{\text{Total Number of terms in a document "d"}}$$
    - inverse-document-frequency (IDF) shows how important the term is to be in the collection of documents. calculates the weight of rare term of the text in a collection of documents
$$IDF = \log \epsilon \left( \frac{\text{Total number of documents}}{\text{Number of documents that contains term "t"}} \right)$$
    - provide each word count and the frequency of rare words in order to provide them weights on the basis of their rarity
    - better than conventional counting of occurrence of the word as it only counts the frequency without classification
- Latent Dirichlet Allocation (LDA)
  - uses dirichlet distribution
  - effective technique in the next word prediction
  - unsupervised approach
  - assumption is that each document mix with various topics and every topic mix with various words
  - commonly used in topic modelling
  - <https://www.analyticssteps.com/blogs/introduction-latent-semantic-analysis-lsa-and-latent-dirichlet-allocation-lda>

## dirichlet distribution

- probability distribution
- sum of probabilities unlike the normal distribution
- continuous multivariate probability distribution
- commonly used for extraction task



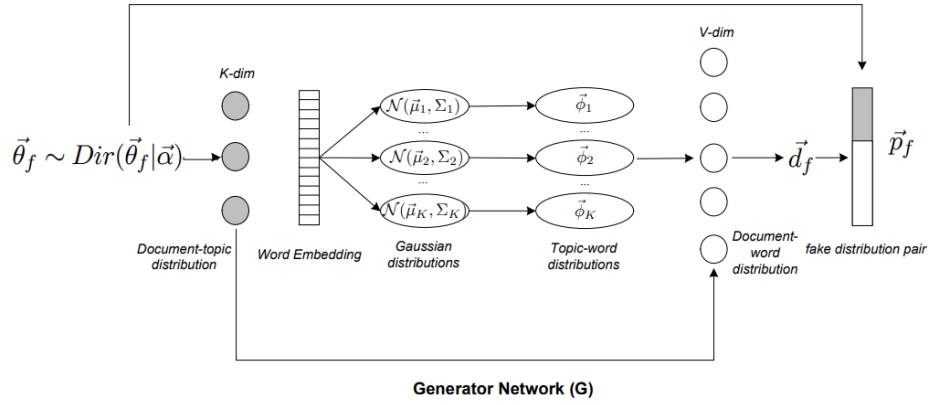
# Neural Topic Modelling with Bidirectional Adversarial Training

- Goal : propose a neural topic modelling approach, called Bidirectional Adversarial Topic (BAT) model, which represents the first attempt of applying bidirectional adversarial training for neural topic modelling
- Link : <https://aclanthology.org/2020.acl-main.32.pdf> (07/2020)
- Code : [https://github.com/zll17/Neural\\_Topic\\_Models](https://github.com/zll17/Neural_Topic_Models)
- Credible source: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics

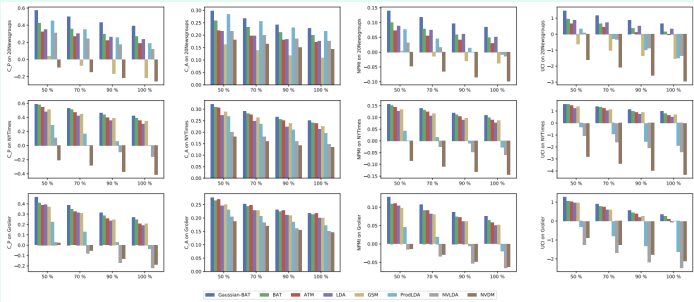
## Model

approaches	description
<ul style="list-style-type: none"> <li>• due to the difficulty of exact inference, LDA variants require approximate inference methods - small changes to the modelling assumptions result in a re-derivation of the inference algorithm</li> <li>• Bidirectional Adversarial Topic model (BAT) <ul style="list-style-type: none"> <li>• employs a generator network to learn the projection function from randomly-sampled document-topic distribution to document-word distribution</li> <li>• encoder network is used to learn the inverse projection, transforming a document-word distribution into a document-topic distribution</li> <li>• employs a discriminator which aims to discriminate between real distribution pair and fake distribution pair, thereby helps the networks (generator and encoder) to learn the two-way projections better</li> <li>• during the adversarial training phase, the supervision signal provided by the discriminator will guide the generator to construct a more realistic document and thus better capture the semantic patterns in text</li> <li>• the encoder network is also guided to generate a more reasonable topic distribution conditioned on specific document-word distributions</li> <li>• to incorporate the word relatedness information captured by word embeddings, we extend the BAT by modelling each topic with a multivariate Gaussian in the generator and propose the Bidirectional Adversarial Topic model with Gaussian (Gaussian-BAT)</li> </ul> </li> </ul>	<p>The diagram illustrates the BAT model architecture, which consists of three main components: the Encoder Network (E), the Generator Network (G), and the Discriminator Network (D).</p> <p><b>Encoder Network (E):</b> This network takes a document-word distribution <math>\vec{d}_r</math> as input. It consists of three layers: a Document-word distribution layer (V-dim), a Representation layer (S-dim), and a Document-topic distribution layer (K-dim). The output is a real distribution pair <math>\vec{p}_r</math>.</p> <p><b>Generator Network (G):</b> This network takes a random topic distribution <math>\vec{\theta}_f \sim \text{Dir}(\vec{\theta}_f   \vec{\alpha})</math> as input. It consists of three layers: a Document-topic distribution layer (K-dim), a Representation layer (S-dim), and a Document-word distribution layer (V-dim). The output is a fake distribution pair <math>\vec{p}_f</math>.</p> <p><b>Discriminator Network (D):</b> This network takes the real distribution pair <math>\vec{p}_r</math> and the fake distribution pair <math>\vec{p}_f</math> as input. It consists of a Joint distributions layer ((V+K)-dim) and a Representation layer (S-dim). The output is a discrimination signal <math>D_{out}</math>.</p> <p>BAT consists of three components</p> <ol style="list-style-type: none"> <li>1. encoder <ul style="list-style-type: none"> <li>• learns a mapping function to transform document-word distribution to document-topic distribution</li> <li>• takes V dimensional document representation <math>d</math> sampled from text corpus C as input</li> <li>• transform it into the corresponding K dimensional topic distribution</li> </ul> </li> <li>2. generator <ul style="list-style-type: none"> <li>a. provides an inverse projection from document-topic distribution to document-word distribution</li> <li>b. takes a random topic distribution drawn from a dirichlet prior as input</li> <li>c. generates a V dimensional fake word distribution</li> </ul> </li> <li>3. discriminator <ul style="list-style-type: none"> <li>a. takes the real distribution pair and fake distribution pair as input</li> <li>b. discriminates the real distribution pairs from the fake ones</li> <li>c. the outputs of the discriminator are used as supervision signals to learn during adversarial training</li> </ul> </li> </ol>

- BAT with Gaussian
  - in BAT, the generator models topics based on the bag-of-words assumption
  - to incorporate the word relatedness information captured in word embeddings into the inference process, propose Gaussian-BAT
  - models each topic with a multivariate Gaussian



## Experiment

training data	experiment	evaluation	output																																																																																																																																	
three datasets used	baseline model	topic coherence measures - average or median of pairwise word similarities formed by top words of a given topic	1.																																																																																																																																	
<div>1. 20Newsgruops<ul style="list-style-type: none"><li>collection of approximately 20,000 news group articles</li></ul></div> <div>2. Grolier<ul style="list-style-type: none"><li>built from Grolier Multimedia Encyclopedia</li><li>covers various fields</li></ul></div> <div>3. NYTimes<ul style="list-style-type: none"><li>collection of news articles published between 1987 and 2007</li><li>contains a wide range of topics</li></ul></div> <div><ul style="list-style-type: none"><li>use the full datasets of 20Newsgroups and Grolier</li><li>for the NYTimes dataset, randomly select 100,000 articles and remove the low frequency words</li></ul></div>	<div><ul style="list-style-type: none"><li>LDA - extracts topics based on word co-occurrence patterns from documents</li><li>NVDM - an unsupervised text modelling approach based on VAE</li><li>GSM - an enhanced topic model based on NVDM</li><li>NVLDA - also built on VAE but with the logistic-normal prior</li><li>ProdLDA - a variant of NVLDA, in which the distribution over individual words is a product of experts</li><li>ATM - a neural topic modelling approach based on adversarial training</li></ul></div> <div>advanced model</div> <div><ul style="list-style-type: none"><li>BAT and Gaussian-BAT</li><li>Gaussian-BAT generator : pre-trained 3000-dimensional Glove embedding is used</li></ul></div>	<div><ul style="list-style-type: none"><li>employ four topic coherence metrics (C_P, C_A, NPMI and UCI) to evaluate the topics generated by various models (Roder et al. , 2015)</li><li>each topic is represented by the top 10 words according to the topic-word probabilities</li><li>all the topic coherence values are calculated using the Palmetto library (<a href="https://github.com/dice-group/Palmetto">https://github.com/dice-group/Palmetto</a>)</li></ul></div> <div>evaluation</div> <div><div>1. compare topic coherence vs. different topic proportions<ul style="list-style-type: none"><li>calculate the average topic coherence values among topics whose coherence values are ranked at the top 50%, 70%, 90%, 100% position</li></ul></div><div>2. compare the average topic coherence values numerically to show the effectiveness of proposed BAT and Gaussian-BAT</div></div>	<div></div> <div><ul style="list-style-type: none"><li>Gaussian-BAT outperforms all the baselines except for Grolier dataset on C_A when considering 100% topics</li></ul></div> <div>2.</div> <div><table><tr><th>Dataset</th><th>Model</th><th>C_P</th><th>C_A</th><th>NPMI</th><th>UCI</th></tr><tr><td rowspan="8">20Newsgroups</td><td>NVDM</td><td>-0.2558</td><td>0.1286</td><td>-0.0984</td><td>-2.9496</td></tr><tr><td>GSM</td><td>-0.2318</td><td>0.1067</td><td>-0.0400</td><td>-1.6083</td></tr><tr><td>NVLDA</td><td>0.1205</td><td>0.1763</td><td>-0.0207</td><td>-1.3466</td></tr><tr><td>ProdLDA</td><td>0.1858</td><td>0.2155</td><td>-0.0083</td><td>-1.5044</td></tr><tr><td>LDA</td><td>0.2361</td><td>0.1769</td><td>0.0523</td><td>0.3399</td></tr><tr><td>ATM</td><td>0.1914</td><td>0.1720</td><td>0.0207</td><td>-0.3871</td></tr><tr><td>BAT</td><td>0.2597</td><td>0.1976</td><td>0.0472</td><td>0.0969</td></tr><tr><td>Gaussian-BAT</td><td><b>0.3758</b></td><td><b>0.2251</b></td><td><b>0.0819</b></td><td><b>0.5925</b></td></tr><tr><td rowspan="8">Grolier</td><td>NVDM</td><td>-0.1877</td><td>0.1456</td><td>-0.0619</td><td>-2.1149</td></tr><tr><td>GSM</td><td>0.1974</td><td>0.1966</td><td>0.0491</td><td>-0.0410</td></tr><tr><td>NVLDA</td><td>-0.2205</td><td>0.1504</td><td>-0.0653</td><td>-2.4797</td></tr><tr><td>ProdLDA</td><td>-0.0374</td><td>0.1733</td><td>-0.0193</td><td>-1.6398</td></tr><tr><td>LDA</td><td>0.1908</td><td>0.2009</td><td>0.0497</td><td>-0.0503</td></tr><tr><td>ATM</td><td>0.2105</td><td><b>0.2188</b></td><td>0.0582</td><td>0.1051</td></tr><tr><td>BAT</td><td>0.2312</td><td>0.2108</td><td>0.0608</td><td>0.1709</td></tr><tr><td>Gaussian-BAT</td><td><b>0.2606</b></td><td>0.2142</td><td><b>0.0724</b></td><td><b>0.2836</b></td></tr><tr><td rowspan="8">NYTimes</td><td>NVDM</td><td>-0.4130</td><td>0.1341</td><td>-0.1437</td><td>-4.3072</td></tr><tr><td>GSM</td><td>0.3426</td><td>0.2232</td><td>0.0848</td><td>0.6224</td></tr><tr><td>NVLDA</td><td>-0.1575</td><td>0.1482</td><td>-0.0614</td><td>-2.4208</td></tr><tr><td>ProdLDA</td><td>-0.0034</td><td>0.1963</td><td>-0.0282</td><td>-1.9173</td></tr><tr><td>LDA</td><td>0.3083</td><td>0.2127</td><td>0.0772</td><td>0.5165</td></tr><tr><td>ATM</td><td>0.3568</td><td>0.2375</td><td>0.0899</td><td>0.6582</td></tr><tr><td>BAT</td><td>0.3749</td><td>0.2355</td><td>0.0951</td><td>0.7073</td></tr><tr><td>Gaussian-BAT</td><td><b>0.4163</b></td><td><b>0.2479</b></td><td><b>0.1079</b></td><td><b>0.9215</b></td></tr></table></div> <div><ul style="list-style-type: none"><li>Gaussian-BAT gives the best overall results across all metrics and on all the datasets except for Grolier dataset on C_A</li></ul></div>	Dataset	Model	C_P	C_A	NPMI	UCI	20Newsgroups	NVDM	-0.2558	0.1286	-0.0984	-2.9496	GSM	-0.2318	0.1067	-0.0400	-1.6083	NVLDA	0.1205	0.1763	-0.0207	-1.3466	ProdLDA	0.1858	0.2155	-0.0083	-1.5044	LDA	0.2361	0.1769	0.0523	0.3399	ATM	0.1914	0.1720	0.0207	-0.3871	BAT	0.2597	0.1976	0.0472	0.0969	Gaussian-BAT	<b>0.3758</b>	<b>0.2251</b>	<b>0.0819</b>	<b>0.5925</b>	Grolier	NVDM	-0.1877	0.1456	-0.0619	-2.1149	GSM	0.1974	0.1966	0.0491	-0.0410	NVLDA	-0.2205	0.1504	-0.0653	-2.4797	ProdLDA	-0.0374	0.1733	-0.0193	-1.6398	LDA	0.1908	0.2009	0.0497	-0.0503	ATM	0.2105	<b>0.2188</b>	0.0582	0.1051	BAT	0.2312	0.2108	0.0608	0.1709	Gaussian-BAT	<b>0.2606</b>	0.2142	<b>0.0724</b>	<b>0.2836</b>	NYTimes	NVDM	-0.4130	0.1341	-0.1437	-4.3072	GSM	0.3426	0.2232	0.0848	0.6224	NVLDA	-0.1575	0.1482	-0.0614	-2.4208	ProdLDA	-0.0034	0.1963	-0.0282	-1.9173	LDA	0.3083	0.2127	0.0772	0.5165	ATM	0.3568	0.2375	0.0899	0.6582	BAT	0.3749	0.2355	0.0951	0.7073	Gaussian-BAT	<b>0.4163</b>	<b>0.2479</b>	<b>0.1079</b>	<b>0.9215</b>
Dataset	Model	C_P	C_A	NPMI	UCI																																																																																																																															
20Newsgroups	NVDM	-0.2558	0.1286	-0.0984	-2.9496																																																																																																																															
	GSM	-0.2318	0.1067	-0.0400	-1.6083																																																																																																																															
	NVLDA	0.1205	0.1763	-0.0207	-1.3466																																																																																																																															
	ProdLDA	0.1858	0.2155	-0.0083	-1.5044																																																																																																																															
	LDA	0.2361	0.1769	0.0523	0.3399																																																																																																																															
	ATM	0.1914	0.1720	0.0207	-0.3871																																																																																																																															
	BAT	0.2597	0.1976	0.0472	0.0969																																																																																																																															
	Gaussian-BAT	<b>0.3758</b>	<b>0.2251</b>	<b>0.0819</b>	<b>0.5925</b>																																																																																																																															
Grolier	NVDM	-0.1877	0.1456	-0.0619	-2.1149																																																																																																																															
	GSM	0.1974	0.1966	0.0491	-0.0410																																																																																																																															
	NVLDA	-0.2205	0.1504	-0.0653	-2.4797																																																																																																																															
	ProdLDA	-0.0374	0.1733	-0.0193	-1.6398																																																																																																																															
	LDA	0.1908	0.2009	0.0497	-0.0503																																																																																																																															
	ATM	0.2105	<b>0.2188</b>	0.0582	0.1051																																																																																																																															
	BAT	0.2312	0.2108	0.0608	0.1709																																																																																																																															
	Gaussian-BAT	<b>0.2606</b>	0.2142	<b>0.0724</b>	<b>0.2836</b>																																																																																																																															
NYTimes	NVDM	-0.4130	0.1341	-0.1437	-4.3072																																																																																																																															
	GSM	0.3426	0.2232	0.0848	0.6224																																																																																																																															
	NVLDA	-0.1575	0.1482	-0.0614	-2.4208																																																																																																																															
	ProdLDA	-0.0034	0.1963	-0.0282	-1.9173																																																																																																																															
	LDA	0.3083	0.2127	0.0772	0.5165																																																																																																																															
	ATM	0.3568	0.2375	0.0899	0.6582																																																																																																																															
	BAT	0.3749	0.2355	0.0951	0.7073																																																																																																																															
	Gaussian-BAT	<b>0.4163</b>	<b>0.2479</b>	<b>0.1079</b>	<b>0.9215</b>																																																																																																																															

## Repository

## Topic subject creation using unsupervised learning for topic modelling

- Goal : compare Non-Negative Matrix Factorization (NMF) and Latent Dirichlet Allocation (LDA) algorithms in the topic mining performance and propose methods to assign topic subject labels in an automated way
- Link : <https://arxiv.org/abs/1912.08868> (02/2019)
- Code : -
- Credible source: -

## Topic Modelling Meets Deep Neural Networks: A Survey

- Goal : provide a focused yet comprehensive overview of neural topic models for interested researchers in the AI community, so as to facilitate them to navigate and innovate in this fast-growing research area
- Link : <https://arxiv.org/abs/2103.00498> (02/2021)
- Code : -
- Credible source: -