

Prompting Large Language Models for Topic Modeling

Han Wang*, Nirmalendu Prakash*, Nguyen Khoi Hoang[†], Ming Shan Hee*, Usman Naseem[‡], Roy Ka-Wei Lee*
Singapore University of Design and Technology*, Singapore

VinUniversity[†], Vietnam

James Cook University of North Queensland[‡], Australia

Email: han_wang@sutd.edu.sg, nirmalendu_prakash@sutd.edu.sg, 20nguyen.hk@vinuni.edu.vn,
mingshan_hee@mymail.sutd.edu.sg, usman.naseem@jcu.edu.au, roy_lee@sutd.edu.sg

Abstract—Topic modeling is a widely used technique for revealing underlying thematic structures within textual data. However, existing models have certain limitations, particularly when dealing with short text datasets that lack co-occurring words. Moreover, these models often neglect sentence-level semantics, focusing primarily on token-level semantics. In this paper, we propose **PromptTopic**, a novel topic modeling approach that harnesses the advanced language understanding of large language models (LLMs) to address these challenges. It involves extracting topics at the sentence level from individual documents, then aggregating and condensing these topics into a predefined quantity, ultimately providing coherent topics for texts of varying lengths. This approach eliminates the need for manual parameter tuning and improves the quality of extracted topics. We benchmark **PromptTopic** against the state-of-the-art baselines on three vastly diverse datasets, establishing its proficiency in discovering meaningful topics. Furthermore, qualitative analysis showcases **PromptTopic**'s ability to uncover relevant topics in multiple datasets.

Index Terms—topic modeling, large language models, prompt engineering

I. INTRODUCTION

Motivation. Topic modeling stands as a pivotal statistical method, focused on discerning latent thematic patterns in textual datasets [1]. It has secured a substantial footing in diverse areas, such as information retrieval and text mining. Its prowess in efficiently extracting topics from voluminous document collections bestows researchers with the capacity to delve into vast textual datasets in an efficient manner.

Over the years, topic modeling has burgeoned as an integral domain within natural language processing and machine learning. Pioneering endeavors in the field leaned towards a bag-of-words probabilistic framework to ascertain topics [1]–[4]. Contemporary strides, however, have steered towards embracing word embedding-centric [5], [6], neural frameworks [7], and transformer paradigms [8], [9]—all in a bid to encapsulate subtler intricacies within textual compositions.

Yet, the evolution hasn't rendered the field impervious to challenges. Encounters with unfamiliar words still pose significant hurdles, attributed to the fact that these models conventionally thrive on predetermined lexicons. Moreover, the fixation on word-level analysis often overshadows the deeper, contextual essence embedded within sentences. Additionally, the recurrent need to meticulously adjust hyperparameters

for superior outputs makes these models not only resource-intensive but also intricate to handle.

Research Objectives. In light of these pressing challenges, our research introduces **PromptTopic**. This avant-garde, prompt-driven strategy for topic modeling taps into the vast potential of large language models (LLMs). Explicitly, **PromptTopic** integrates capabilities of renowned LLMs like ChatGPT¹ and LLaMa [10] to seamlessly intertwine word and sentence semantics—paving the way for a more holistic topic modeling experience. Our method meticulously crafts prompts that are adept at isolating lucid and actionable topics from texts, eliminating the often laborious fine-tuning phase. The adoption of in-context learning through prompts further negates the time-consuming hyperparameter calibration, thus refining the entire topic modeling paradigm.

Contributions. The main contributions of this work are as follows: 1. We propose **PromptTopic**, a novel prompt-based model to perform topic modeling on text. To the best of our knowledge, this is the first topic modeling model that utilizes LLMs. 2. We conduct comprehensive experiments on three widely used topic modeling datasets to evaluate the performance of **PromptTopic** compared to state-of-the-art topic models. 3. We conduct a qualitative analysis of the learned topics, highlighting that our model exhibits the ability to identify meaningful and highly coherent topics.

II. RELATED WORK

A. Topic Modeling

Topic modeling, a significant area within natural language processing and information retrieval, is centered on unveiling abstract "topics" within a document collection. Topic modeling also has been applied to many other domain such as identifying topical influential users in social media [11]–[13]. Over time, a plethora of methods have evolved to improve the topic modeling performance. A standout model is the Latent Dirichlet Allocation (LDA) by [1], which has been foundational for later innovations. This model was enhanced by its successors, such as supervised LDA (sLDA) [14] and dynamic topic modeling (DTM) [15]. Other noteworthy techniques include non-

¹<https://api.openai.com/v1/chat/completions>

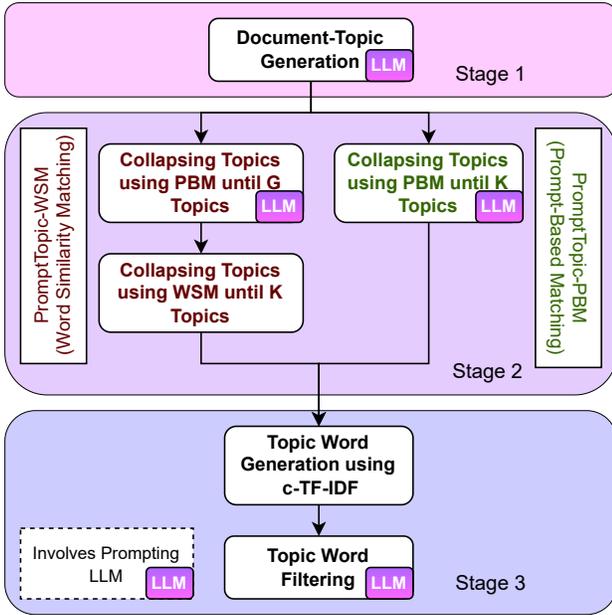


Fig. 1. Overview of PromptTopic.

negative matrix factorization (NMF) by [4] and probabilistic latent semantic analysis (pLSA) by [16].

The deep learning era brought neural models like the adapted Variational Autoencoder (VAE) and Transformer for topic modeling. ProLDA by [7] is one such adaptation of LDA. Recent models like Correlated Topic Model [5], Gaussian LDA [17], Spherical HDP [18], and Embedded Topic Model (ETM) [6], integrate word embeddings like Word2Vec [19] to grasp semantic word relationships. Strategies, including TopClus [20] and Cluster-Analysis [21], emphasize clustering word and document embeddings, offering flexibility by separating cluster creation from topic representation. Further, models like Contextualized Topic Model (CTM) [8] and BERTopic [9] leverage transformers like BERT [22] to assimilate context and shape topic representations. Similar efforts employ pretrained language embeddings for the same endeavor [23], [24].

In our paper, we introduce PromptTopic, an innovative approach that leverages Large Language Models (LLMs) for topic modeling. Unlike traditional models, it seamlessly incorporates word and sentence semantics, facilitating precise and contextually relevant topic identification. Importantly, ‘PromptTopic’ streamlines the process by avoiding extensive hyperparameter tuning, making it more accessible to researchers and users.

B. Large Language Model

In the dynamic world of large language models (LLMs), foundational works have set the stage for innovative research and applications. The groundbreaking Transformer architecture, introduced by [25], laid the groundwork for subsequent LLMs. A milestone was reached in 2020 with the release of GPT-3 by [26], a model with a staggering 175 billion param-

eters, displaying unmatched capabilities. A notable variant of GPT-3, designed for conversations, is ChatGPT.

For a long time, the consensus was that larger models equated to better performance. However, this notion was recently challenged by J. Hoffmann in [27]. They argued that superior results, within a set budget, can be achieved with smaller models trained on extensive datasets. Touvron demonstrated LLaMA in [10], a set of efficient LLMs with sizes from 7 billion to 65 billion parameters, showcasing their competitiveness against larger LLMs.

The current era of LLMs is marked by multimodality. Models like GPT-4 [28] and LLaVA [29] now encompass visual processing capabilities, enhancing their versatility.

Our research aims to utilize both accessible online LLM APIs, such as ChatGPT, and offline models like LLaMA, to assess their efficacy in topic modeling scenarios.

III. METHODOLOGY

The PromptTopic is an unsupervised approach that harnesses the robust language understanding capabilities of LLMs to generate topics. The model consists of three stages: *Topic Generation*, *Topic Collapse*, and *Topic Representation Generation*, as illustrated in Figure 1. Each stage leverages LLMs to extract, organize, and refine topics from input documents. Using prompts, the model effectively captures underlying themes and concepts, enhances topic clustering, and improves the quality of topic representations. Prompting LLMs allows PromptTopic to learn document topics comprehensively, eliminating the need for fine-tuning.

A. Topic generation

Figure 2 displays the prompt setup used for topic generation using ChatGPT. The prompt input comprises N demonstration examples, each with prompt inputs and their associated annotated answers. For LLaMA, which isn’t instruction-trained, the instructional statements are omitted from the prompt.

To find the optimal value for the demonstration parameter N , we tested values of 2, 4, 6, and 8. Our results indicated that a setting of $N = 4$ produced the best topic generation performance, notably reducing errors in the offline LLaMA model. ChatGPT, with its larger parameter size, was less sensitive to changes in N . Table I presents topics generated by LLaMA for a randomly selected document, with N ranging from 2 to 8.

B. Collapsing Overlapping Topics

LLMs often generate overlapping topics for a document. For instance, topics such as “film”, and “actor” can be merged into one topic “film”. We introduced two approaches to collapse topics: *Prompt-Based Matching (PBM)* and *Word Similarity Matching (WSM)*.

a) *Prompt-Based Matching (PBM)*: We utilize prompts to group similar topics. The process involves sorting the unique topics based on their frequency counts and selecting a specified K number of topics. Initially, we form a set T_n containing the unique topics. Then, we create a subset T_{n-1} by selecting the

TABLE I
TOPICS GENERATED BY LLAMA FOR DEMONSTRATION NUMBERS (N) FROM 2 TO 8, WITH CURRENT DOCUMENT-RELATED TOPICS HIGHLIGHTED.

Document	N (Number of Demonstrations)			
	2	4	6	8
trailer talk week movie rite mechanic week opportunity	politics, army	movies, trailers, mechanic	politics, economy, finance	social

System Instructions: *You are designated as an assistant that identify and extract high-level topics from articles. You should avoid giving specific details and provide unique topics solely.*

Please list the high-level topics in the following article.
Article: suffering bad credit history bad credit debt consolidation loan pay debt
Topics:

['finance']

○ ○ ○ N demonstrations

Please list the high-level topics in the following article.
Article: indian navy coast guard rescue thai vessel pirate joint operation indian navy coast guard intercepted neutr
Topics:

['politics', 'army']

Fig. 2. PromptTopic’s topic generation: N turns of user input (purple) and assistant’s sample answer (green). ChatGPT’s output answer is shown in red.

first $n - 1$ topics from T_n . Next, we prompt LLMs to merge each topic t_n from T_n with a topic from T_{n-1} . If no merge is possible, we merge t_n with the “Miscellaneous” topic. This iterative process continues until T_n contains only K topics.

During experimentation, we encountered a challenge with datasets that had a large number of unique topics, exceeding the maximum token length allowed by LLMs. To overcome this, we employed a sliding window approach. We selected a window of size M from the sorted unique topic set and performed the iterative topic grouping process. If LLMs successfully merged a topic with one of the M topics, we merged them and restarted the iteration. If no merge occurred, we categorized the topic as “Miscellaneous.”

b) *Word Similarity Matching (WSM)*.: This approach involves computing topic similarity and merging highly similar topics. We aggregate documents associated with each topic and compute a Class-based Term Frequency-Inverse Document Frequency (c-TF-IDF) representation, which captures word

TABLE II
DATASET STATISTICS: SIZE INDICATES THE NUMBER OF DOCUMENTS, CATEGORY INDICATES THE NUMBER OF CATEGORIES, TEXT INDICATES THE AVERAGE LENGTH OF DOCUMENTS.

Dataset	Size	Category	Text
20 NewsGroup	16,309	20	185.37
Yelp Reviews	10,000	-	131.03
Twitter Tweet	2,472	89	8.55

frequencies within topics while considering their importance across all topics. We select the top 20 words from the c-TF-IDF representation, retaining only relevant words closely related to the topic. Topic similarity is measured by counting the number of common words between the top words of each topic pair, normalized by word count. We merge the pair with the highest word similarity, create a new topic, and recalculate the similarity score based on the merged content. This process is iterated until we have K unique topics remaining.

To reduce computation time for large datasets with numerous topics, we utilize the PBM model to compress the initial n topics into a more concise set of topics, designated as G . This G value exceeds the intended number of topics, K , yet remains substantially smaller than the initial number of topics, n . Subsequently, we will proceed to further condense the set of G topics to the desired K through the utilization of the WSM technique.

C. Topic Representation Generation

In order to evaluate the performance of our PromptTopic model, we utilized well-established topic model metrics. However, evaluating the topics required representing them as word mixtures. Since our model did not generate topic-word distributions directly, we employed c-TF-IDF scores to compute the most representative words for each cluster. Initially, we obtained the top 100 c-TF-IDF words for each topic. To further refine the representation, we used LLMs to filter these words down to the top 10 most representative words. This process allowed us to assess the quality and coherence of the generated topics based on their representative word mixtures.

IV. EXPERIMENT

A. Experiment Settings

Datasets. We evaluate PromptTopic and baselines on three commonly-used topic modeling datasets: *20 NewsGroup* [30], *Yelp Reviews* [31] and *Twitter Tweet* [32]. Table II provide a statistical summary of the datasets

TABLE III
COMPARISON OF NPMI AND TOPIC DIVERSITY (TD) ACROSS THREE DATASETS. CELL COLOR INTENSITY INDICATES THE SCORE’S RANKING AMONG ALL MODELS, WITH THE BEST PERFORMING MODEL UNDERLINED.

Model	20 NewsGroup		Yelp Reviews		Twitter Tweet	
	NPMI	TD	NPMI	TD	NPMI	TD
LDA	-0.05	0.81	-0.01	0.43	-0.36	0.41
NMF	0.04	0.63	0.02	0.45	-0.25	0.60
CTM	-0.01	0.96	-0.09	0.78	0.03	0.96
TopClus	-0.13	0.92	-0.13	0.92	-0.37	0.92
Cluster-Analysis	-0.02	0.99	-0.04	0.96	-0.43	0.27
BERTopic	0.10	0.97	0.10	0.81	0.05	0.98
PromptTopic-PBM(LLaMA)	-0.12	0.97	-0.24	0.99	-0.14	0.91
PromptTopic-PBM(ChatGPT)	-0.15	0.89	-0.26	0.98	-0.04	0.95
PromptTopic-WSM(LLaMA)	0.05	0.91	0.04	0.86	0.04	0.97
PromptTopic-WSM(ChatGPT)	0.04	0.84	0.08	0.76	-0.07	0.91

TABLE IV
QUALITATIVE EVALUATION OF THE TOPIC-WORDS REPRESENTATION IN TWITTER TWEET DATASET. A SUBSET OF TOPICS THAT OCCUR FREQUENTLY ARE SELECTED. THE RELATED WORDS BELONGING TO THE CORRESPONDING TOPIC ARE HIGHLIGHTED IN BOLD.

Dataset	Topic	BASELINE MODEL				PROMPTTOPIC	
		NMF	CTM	Cluster-Analysis	BERTopic	WSM(LLaMA)	PBM(LLaMA)
Twitter Tweet	Politics	yemen	protest	king	egypt	president	president
		protest	thousand	christina	yemen	obama	protest
		president	yemen	egypt	journalist	judge	protester
		aquarium	government	yemen	president	law	government
		somali	yemeni	sundance	protest	lawsuit	judge
	Sports	superbowl	fish	superbowl	commercial	hockey	nba
		commercial	aquarium	super	superbowl	sidney	nfl
		super	bass	christina	super	crosby	football
		bowl	fishing	egypt	bowl	concussion	tennis
		ad	tackle	yemen	christina	safety	hockey

Baseline Models. PromptTopic will be evaluated against six widely used topic modeling models: LDA [1], NMF [4], CTM [8], TopClus [20], Cluster-Analysis [21] and BERTopic [9].

PromptTopic Configurations. In our experiments, we applied the PromptTopic model to two state-of-the-art LLMs: ChatGPT and LLaMA-13B [10]. However, due to the limited parameter size and absence of instruction training in LLaMA, we simplified the prompt format and reduced the number of demonstration examples. By adopting this approach, we ensured that the performance of LLaMA was comparable to that of ChatGPT. Note that for all models, we preprocess the datasets using the OCTIS package [33], which involved removing punctuation, stopwords, and performing lemmatization (except the Twitter Tweet dataset).

Setting Optimal Value of Parameter G . To determine the optimal value of parameter G across three datasets introduced in Datasets. We conducted an empirical investigation encompassing G values of 200, 400, 600, and 800. We evaluated G based on two main criteria: the time for topic collapsing and quantitative assessments using established metrics from the Quantitative Evaluation Section IV-B, including topic coherence (NPMI [34]) and topic diversity (TD [6]). After conducting our assessment, we determined that the optimal value for parameter G was 400 for both the 20 NewsGroup

and Twitter Tweet datasets, while for the Yelp Reviews dataset, a value of 200 yielded the best results. Consequently, we selected the most effective value of G for the subsequent experiments.

B. Topic Evaluation

Quantitative Evaluation. Evaluations of topic modeling often use two well-established metrics: topic coherence and topic diversity. Topic coherence gauges the extent to which the words within a topic are related, forming a coherent group. It is typically calculated using statistics and probabilities drawn from the reference corpus, focusing specifically on the context of the words. In our experiments, we employed Normalized Pointwise Mutual Information (NPMI) [34] as our measure of topic coherence. A higher NPMI score signifies better coherence, with a perfect correlation being represented by a score of 1.

Conversely, topic diversity [6] evaluates the proportion of unique words across all topic representations. The diversity score ranges from 0 to 1, where a score of 0 indicates repetitive topics, and a score of 1 indicates diverse topics. This metric is crucial for ensuring that a topic model covers a wide range of themes without overemphasizing any particular topic. Using these two metrics together provides insights into the effectiveness of topic modeling algorithms in identifying both coherent and diverse topics.

In our evaluation process, we empirically selected the number of topics (K) for each dataset. K is set to 40, 20, 20 for 20NewsGroup, Yelp Reviews, and Twitter Tweet Datasets, respectively. Table III shows NPMI and topic diversity scores for different topic models on three datasets. Our findings indicate that PromptTopic-WSM consistently outperforms the majority of baseline topic models across all datasets, as evidenced by both metrics. Notably, compared to the best-performing baseline model, BERTopic, the performance of PromptTopic remains comparable.

Remarkably, LLaMA-13b, functioning as a standalone of-line model, exhibits significantly fewer parameters while yielding results of comparable quality to ChatGPT. Their coherency levels closely align, although LLaMA-13b demonstrates a propensity for generating more diverse topics, with a higher TD score. Consequently, in the subsequent phases of qualitative and human evaluation, we shall opt for LLaMA as our preferred Language Model.

Human Evaluation. Low NPMI scores don't necessarily reflect poor topic quality, as they have been found to show a weak correlation with human ratings [35]. To gain a deeper understanding, we decided to conduct a manual evaluation.

Our assessment of topic quality relied on the word intrusion task, following [36]. In this task, we presented participants with a list of five words. Four of these words were from the prominent words of a single topic generated by the model, while the fifth, known as the 'intruder', was randomly chosen from a different topic. The word intrusion test aims to determine if the five words collectively represent a clear and distinct topic, making it easy to identify the intruder.

In our study, we compared topic words generated by PromptTopic with those from the leading baseline model, BERTopic. We then report the intrusion task accuracy results for BERTopic, PromptTopic-PBM (LLaMA), and PromptTopic-WSM (LLaMA) across all topics in the three datasets. Each topic was evaluated by two annotators.

Figure 3 presents word intrusion task results for three models across three diverse datasets. We observed a high level of consistency in word intrusion accuracy, with an average score of around 65% across all datasets. This suggests that the generated topic words maintain consistent high quality. However, it's essential to note that model performance is significantly influenced by dataset characteristics. Regarding PromptTopic-WSM and BERTopic, their performance is particularly strong when dealing with lengthy textual datasets like 20 NewsGroup and Yelp Reviews. However, they exhibit reduced performance when handling short text data, as seen in Twitter Tweets. In contrast, PromptTopic-PBM presents a unique scenario. It performs suboptimally in the Yelp Reviews dataset due to the dataset's strong focus on food-related content. This concentration results in PromptTopic-PBM generating highly specific topics, leading to a topic diversity score of 0.99. However, the disjointed occurrence of these specific food-related words within the same document affects coherence negatively. On the other hand, PromptTopic-PBM excels in the Twitter Tweet dataset, achieving an impressive word

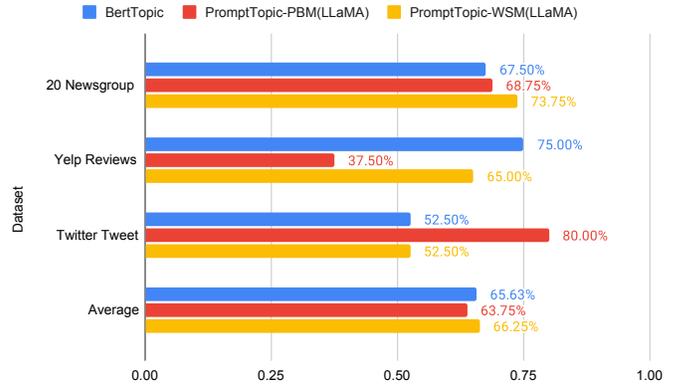


Fig. 3. Word Intrusion Study Results for 40, 20, and 20 Topics Across 20 NewsGroup, Yelp Reviews, and Twitter Tweet Datasets, generated by BERTopic, PromptTopic-PBM(LLaMA) and PromptTopic-WSM(LLaMA) Models. Average reflects the overall accuracy across all three datasets.

intrusion task accuracy of 0.80, surpassing the performance of the other models. This highlights PromptTopic-PBM's effectiveness in handling short text data, leveraging the power of Language Model Models (LLMs) to extract relevant topics.

Qualitative Evaluation. We randomly selected commonly occurring topics from all datasets and performed a manual matching procedure to determine the most relevant topic generated by each model. Table IV illustrates the top five words associated with each topic, generated by employing the LLaMA method for the PromptTopic. To enhance page space efficiency without compromising clarity, we have selectively showcased the Twitter Tweet dataset and restricted the presentation to four fundamental models, notable for their strong topical coherence across the majority of datasets, as delineated in the Quantitative Evaluation Section IV-B. Despite BERTopic's higher NPMI score, the manual assessment reveals that PromptTopic-PBM demonstrates comparable topic representation. BERTopic falls short in the short text such as the Sports topic in Twitter Tweet dataset, providing only three relevant words, with only 'superbowl' being informative. In contrast, PromptTopic-PBM generates informative words like 'nba', 'nfl', 'football', 'tennis', and 'hockey' without any overlap. The observed enhancement of PromptTopic-PBM performance in short text datasets aligns with human evaluation findings, which can be attributed to LLMs' robust language comprehension and vast knowledge base.

V. LIMITATIONS

When dealing with large datasets, the usage of LLMs for topic generation can be resource-intensive. LLMs like LLaMA require GPU devices with significant memory capacity. In our experiments, we employed the PBM method to collapse topics by prompting the LLM to merge a topic with a list of topics solely based on topic names. However, this approach lacks context and may result in the merging of unrelated topics. Furthermore, in datasets characterized by a substantial number

of topics, the need for batch-wise merging in PBM and the assistance of PBM in WSM becomes imperative.

VI. CONCLUSION

In this paper, we introduce PromptTopic, a groundbreaking approach to topic modeling utilizing LLMs. Our innovative method stands out by harnessing the power of LLMs to discern semantic structures both at the token and sentence levels. This ensures the generation of not just coherent, but also diverse topics. Through rigorous evaluations of diverse datasets, we not only validate the robustness of our approach but also highlight its superior capabilities. When compared to leading contemporary methods, PromptTopic not only matches them in terms of automatic metrics but notably surpasses them in producing more meaningful topics upon qualitative assessment. This underscores the significant contribution of our work to the field. In future projects, we'll investigate ways to enhance batch-wise merging in PBM and utilize prompt-engineering methods for better topic modeling.

REFERENCES

- [1] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent dirichlet allocation," *Journal of machine Learning research*, vol. 3, no. Jan, pp. 993–1022, 2003.
- [2] M. Rosen-Zvi, T. Griffiths, M. Steyvers, and P. Smyth, "The author-topic model for authors and documents," in *Proceedings of the 20th Conference on Uncertainty in Artificial Intelligence*, ser. UAI '04. Arlington, Virginia, USA: AUAI Press, 2004, p. 487–494.
- [3] D. Blei and J. Lafferty, "Correlated topic models," *Advances in neural information processing systems*, vol. 18, p. 147, 2006.
- [4] C. Févotte and J. Idier, "Algorithms for nonnegative matrix factorization with the β -divergence," *Neural computation*, vol. 23, no. 9, pp. 2421–2456, 2011.
- [5] G. Xun, Y. Li, W. X. Zhao, J. Gao, and A. Zhang, "A correlated topic model using word embeddings," in *International Joint Conference on Artificial Intelligence (IJCAI)*, vol. 17, 2017, pp. 4207–4213.
- [6] A. B. Dieng, F. J. R. Ruiz, and D. M. Blei, "Topic modeling in embedding spaces," *Trans. Assoc. Comput. Linguist.*, vol. 8, pp. 439–453, Dec. 2020.
- [7] A. Srivastava and C. Sutton, "Autoencoding variational inference for topic models," *International Conference on Learning Representations (ICLR)*, 2017.
- [8] F. Bianchi, S. Terragni, D. Hovy, D. Nozza, and E. Fersini, "Cross-lingual contextualized topic models with zero-shot learning," *Association for Computational Linguistics (ACL)*, 2020.
- [9] M. Grootendorst, "BERTopic: Neural topic modeling with a class-based TF-IDF procedure," *arXiv preprint arXiv:2203.05794*, Mar. 2022.
- [10] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, A. Rodriguez, A. Joulin, E. Grave, and G. Lample, "LLAMA: Open and efficient foundation language models," *arXiv preprint arXiv:2302.13971*, 2023, 1, 2, 3.
- [11] R. K.-W. Lee, T.-A. Hoang, and E.-P. Lim, "On analyzing user topic-specific platform preferences across multiple social media sites," in *Proceedings of the 26th International Conference on World Wide Web*, 2017, pp. 1351–1359.
- [12] —, "Discovering hidden topical hubs and authorities in online social networks," in *Proceedings of the 2018 SIAM International Conference on Data Mining*. SIAM, 2018, pp. 378–386.
- [13] —, "Discovering hidden topical hubs and authorities across multiple online social networks," *IEEE Transactions on Knowledge and Data Engineering*, vol. 33, no. 1, pp. 70–84, 2019.
- [14] J. McAuliffe and D. Blei, "Supervised topic models," in *Advances in Neural Information Processing Systems*, vol. 20, 2007.
- [15] D. M. Blei and J. D. Lafferty, "Dynamic topic models," in *Proceedings of the 23rd international conference on Machine learning*, June 2006, pp. 113–120.
- [16] T. Hofmann, "Probabilistic latent semantic analysis," *arXiv preprint arXiv:1301.6705*, 2013.
- [17] R. Das, M. Zaheer, and C. Dyer, "Gaussian Ilda for topic models with word embeddings," in *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 2015, pp. 795–804.
- [18] K. Batmanghelich, A. Saeedi, K. Narasimhan, and S. Gershman, "Nonparametric spherical topic modeling with word embeddings," in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Berlin, Germany: Association for Computational Linguistics, Aug. 2016, pp. 537–542. [Online]. Available: <https://aclanthology.org/P16-2087>
- [19] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," *International Conference on Learning Representations (ICLR)*, Jan. 2013.
- [20] Y. Meng, Y. Zhang, J. Huang, Y. Zhang, and J. Han, "Topic discovery via latent space clustering of pretrained language model representations," in *Proceedings of the ACM Web Conference 2022*, April 2022, pp. 3143–3152.
- [21] S. Sia, A. Dalmia, and S. J. Mielke, "Tired of topic models? clusters of pretrained word embeddings make for fast and good topics too!" *arXiv preprint arXiv:2004.14914*, 2020.
- [22] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019*, J. Burstein, C. Doran, and T. Solorio, Eds. Association for Computational Linguistics, 2019, pp. 4171–4186.
- [23] A. M. Hoyle, P. Goel, and P. Resnik, "Improving Neural Topic Models using Knowledge Distillation," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Online: Association for Computational Linguistics (ACL), Nov. 2020, pp. 1752–1771. [Online]. Available: <https://aclanthology.org/2020.emnlp-main.137>
- [24] P. Gupta, Y. Chaudhary, and H. Schütze, "Multi-source neural topic modeling in multi-view embedding spaces," in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Online: Association for Computational Linguistics, Jun. 2021, pp. 4205–4217. [Online]. Available: <https://aclanthology.org/2021.naacl-main.332>
- [25] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, ..., and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [26] Z. Chen, M. M. Balan, and K. Brown, "Language models are few-shot learners for prognostic prediction," *arXiv preprint arXiv:2302.12692*, Feb. 2023.
- [27] J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai, E. Rutherford, and L. Sifre, "Training compute-optimal large language models," *arXiv preprint*, 2022.
- [28] OpenAI, "Gpt-4 technical report," Tech. Rep., 2023, version b.
- [29] H. Liu, C. Li, Q. Wu, and Y. J. Lee, "Visual instruction tuning," *arXiv preprint*, 2023.
- [30] K. Lang, "Newsweeder: Learning to filter netnews," in *Machine Learning Proceedings 1995*. Elsevier, 1995, pp. 331–339.
- [31] "Yelp dataset challenge," Available from: http://www.yelp.com/dataset_challenge, 2015.
- [32] J. Qiang, Z. Qian, Y. Li, Y. Yuan, and X. Wu, "Short text topic modeling techniques, applications, and performance: A survey," *IEEE Trans. Knowl. Data Eng.*, vol. 34, no. 3, pp. 1427–1445, Mar. 2022.
- [33] S. Terragni, E. Fersini, B. G. Galuzzi, P. Tropeano, and A. Candelieri, "OCTIS: Comparing and optimizing topic models is simple!" in *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*. Association for Computational Linguistics, Apr. 2021, pp. 263–270. [Online]. Available: <https://www.aclweb.org/anthology/2021.eacl-demos.31>
- [34] G. Bouma, "Normalized (pointwise) mutual information in collocation extraction," *Proceedings of GSCL*, vol. 30, pp. 31–40, 2009.
- [35] A. Hoyle, P. Goel, A. Hian-Cheong, D. Peskov, J. Boyd-Graber, and P. Resnik, "Is automated topic model evaluation broken? the incoherence of coherence," *Advances in Neural Information Processing Systems*, vol. 34, pp. 2018–2033, 2021.
- [36] J. Chang, S. Gerrish, C. Wang, J. Boyd-Graber, and D. Blei, "Reading tea leaves: How humans interpret topic models," *Advances in Neural Information Processing Systems*, vol. 22, 2009.