

TOWARDS EXPLAINABLE RECOMMENDATION VIA BERT-GUIDED EXPLANATION GENERATOR

Huijing Zhan^{†1}, Ling Li^{†2}, Shaohua Li¹, Weide Liu¹, Manas Gupta¹, Alex C. Kot²

¹A*STAR, Singapore

²Nanyang Technological University, Singapore

ABSTRACT

Explainable recommender system has recently drawn increasing attention due to its capability of providing justification to recommendation. Rather than focusing on certain topics or specific item features, the explanation generated by existing works are too general without the guidance of aspects. However, such information is not given in the practical scenario. To address this issue, we propose a novel Explainable recommender system with BERT-guided explanation generator, named ExBERT to generate reliable explanation with finer granularity. More specifically, a multi-head self-attention based encoder is employed to incorporate pseudo user and item profiles into semantic representation. Moreover, we propose a novel matched explanation prediction task with discriminative ability to enable personalization of the generated sentence. Extensive experiments conducted on two real-world explainable recommendation datasets significantly outperform the state-of-the-art in generation.

Index Terms— Explanation generation, BERT, matched explanation prediction.

1. INTRODUCTION

Recommendation system has become essential to alleviate the problem of information explosion and considerable research efforts are devoted to enhance the recommendation performance, ranging from collaborative filtering [1], latent factor model [2, 3] to deep neural network [4, 5] and graph neural network [6, 7, 8], etc. However, recommender system still remains a black-box. To increase its persuasiveness, explainable recommendation is an emerging topic which sheds light on the uncertainties with convincing explanation.

Current works on generating explanation for recommendation are categorized into two types: 1) template-based approaches [12, 13, 14] with predefined templates, which suffer from the flexibility and domain knowledge is required; and 2) natural language generation approaches [10, 9, 15] with sequence-to-sequence models, which suffer from the

Table 1. An example of generated explanations by the state-of-the-art methods and the proposed ExBERT. Highlighted green and red words denote the correctly and wrongly generated fine-grained item details, respectively.

	Ratings	Reviews
Ground Truth	5	I'm a 27 in jeans and the petite small size fits me perfectly .
NRT [9]	4.09	It is a <i>nice</i> shirt .
Att2Seq [10]	-	It is a very nice watch .
PETER [11]	4.17	It is a nice shirt .
ExBERT	4.29	I bought a small and it fits perfectly .

long-range dependency issue (i.e., Recurrent Neural Networks) with long sentences as the input and the generated sentence lacks diversity in both content and styles.

Transformer [16] and BERT [17] have recently achieved excellent performances in a variety of natural language processing tasks. However, integrating transformer-based models into the explainable recommendation is relatively unexplored and existing few attempts [11, 18] have their respective limitations. As shown in Table 1, the generated explanation of PETER [11] fails to justify the user's preference towards the items by synthesizing irrelevant description. To generate reliable and accurate explanations, [18] proposes to feed the aspects as the guided input to BERT in order to derive the controllable justification. However, the fine-grained aspects need to be pre-extracted which is subject to the performance of the off-the-shelf toolkit and the expressiveness of the sentiment lexicon [12]. Moreover, in real-world scenarios, such auxiliary knowledge is not always available.

To address the above mentioned drawbacks, we propose a framework named ExBERT which jointly generates rating prediction and explanation sentences. More specifically, a multi-head self-attention based encoder is leveraged to capture the relevance between the historical explanations written by the user and that belonging to the target user-item pair. To improve the personalization of generated explanation, we propose a novel matched explanation prediction task, adapted from the next sentence prediction

[†] denotes equal contribution. This work was supported in part by AHSF project No. C211118010 and Rapid-Rich Object Search (ROSE) Lab and the NTU-PKU Joint Research Institute.

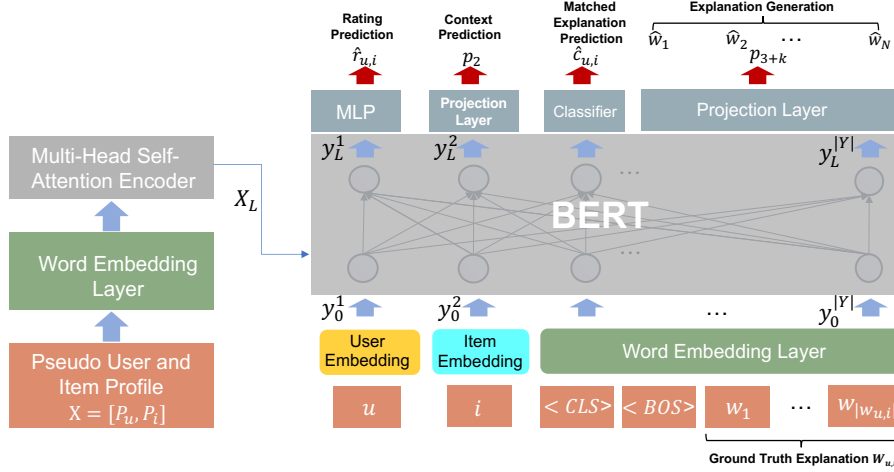


Fig. 1. The overview of the proposed ExBERT framework. Given the user u and item i , we aim to jointly generate the rating prediction \hat{r} and provide explanation $(\hat{w}_1, \hat{w}_2, \dots, \hat{w}_N)$.

(NSP) by regarding the concatenation of user and item ID as the unique identity to discriminate different user’s preferences. The explanation sentences written by the particular user towards the item pair are considered as the positive samples otherwise negatives. Experiments on real-world explainable recommendation datasets validate the superiority of our approach in both improving the rating prediction accuracy and generating fine-grained explanation without pre-extracted aspects as the guidance. The code is available at <https://github.com/zhanhuijing/ExBERT>.

2. METHODOLOGY

2.1. Target-Aware Self-Attention based Encoder

2.1.1. Pseudo User and Item Profile

To make the best of historical explanations, the pseudo user profile \mathcal{P}_u and item profile \mathcal{P}_i are constructed. For the target user-item pair (u, i) , not all the reviews written by the user u are of equal importance to encode the his/her preferences towards the item i . We introduce the semantic ranking procedure on the reviews and extract the highly-ranked ones as the input to the encoder. Firstly we utilize Sentence-BERT [19] to compute the embeddings of the historical review sentences written by u , which carry semantic meanings. Then the similarity scores between the target ground truth explanation $\mathcal{W}_{u,i}$ and the historical review are calculated. The top- K reviews are chosen as the pseudo profile for the user. Finally, these top- K scored sentences are arranged in descending order and concatenated into $\mathcal{P}_u = [P_u^1, \dots, P_u^K]$ as the representation of pseudo user profile. Due to space limitation, we will skip the description of the pseudo item profile $\mathcal{P}_i = [P_i^1, \dots, P_i^K]$, which is constructed in the similar manner.

2.1.2. Multi-Head Self-Attention based Encoder

Self-attention based encoder employs multi-head mechanism to calculate the importance of each word in the sentence. Results from each head are concatenated and a parameterized linear transformation is applied. The word embedding layer maps the input sequence $\mathcal{X} = [P_u^1, \dots, P_u^K, P_i^1, \dots, P_i^K]$ to $\mathbf{X}_0 = [\mathbf{x}_u^1, \dots, \mathbf{x}_u^K, \mathbf{x}_i^1, \dots, \mathbf{x}_i^K]$, where $\mathbf{X} \in \mathbb{R}^{|\mathcal{X}| \times d_x}$ is the d_x -dimension word embedding. To model the mutual dependency relationship between each word in the sequence, L layer of the multi-head self-attention based encoder is employed as below, where $l \in [1, L]$:

$$\mathbf{Z}_l^h = \text{softmax} \left(\frac{(\mathbf{X}_{l-1} \mathbf{W}^Q) (\mathbf{X}_{l-1} \mathbf{W}^K)^T}{\sqrt{d_x}} \right) (\mathbf{X}_{l-1} \mathbf{W}^V)$$

$$\mathbf{X}_l = \text{Concat}(\mathbf{Z}_l^1, \dots, \mathbf{Z}_l^h, \dots, \mathbf{Z}_l^{|H|}) \mathbf{W}^O, \quad (1)$$

\mathbf{X}_{l-1} is the output from the $l-1$ layer. $\mathbf{W}^Q, \mathbf{W}^K, \mathbf{W}^V \in \mathbb{R}^{d_x \times d_z}$ and $\mathbf{W}^O \in \mathbb{R}^{d_z |H| \times d_z}$ are transformation matrices and $|H|$ indicates the number of attention heads in our multi-head self-attention based encoder.

2.2. BERT-based Decoder

Three types of input with user and item IDs as well as the explanation pass through embedding layers. After the position embeddings is further incorporated, the content representation is obtained, denoted as $\mathbf{Y}_0 = [\mathbf{y}_0^1, \dots, \mathbf{y}_0^{|Y|}]$, where $|Y|$ is the length of the sequence and $|Y| = |\mathcal{W}_{u,i}| + 4$. Note that here $\langle BOS \rangle$ and $\langle CLS \rangle$ are in the input sequence with the length of 1. It is subsequently fed forward into L -layer BERT calculated as below:

$$\mathbf{A}_l = \text{MultiHead}(\mathbf{Y}_{l-1}, \mathbf{Y}_{l-1}, \mathbf{Y}_{l-1}), \quad (2)$$

where \mathbf{A}_l is the output of l -th multi-head attention layer. To integrate the output representation \mathbf{X}_L from the encoder, we adopt a cross-attention module to model its impact on the hidden states of \mathbf{A}_l :

$$\mathbf{Y}_l = \text{MultiHead}(\mathbf{A}_l, \mathbf{X}_L, \mathbf{X}_L), \quad (3)$$

where output \mathbf{Y}_L after L layers is utilized to perform the tasks mentioned below.

2.2.1. Matched Explanation Prediction

The next sentence prediction (NSP) task in the original BERT model is to predict whether sentence B follows the sentence A in the same context. Its effectiveness in modeling the sentence-level coherence has been empirically validated in a variety of downstream tasks [17]. Inspired by this, we propose matched explanation prediction task and adapt the definition of ‘‘sentence’’ to make it suitable in our explanation recommendation scenario. The newly designed task treats the concatenation of the user-item ID pair as the identity sentence A and the ground truth explanation $\mathcal{W}_{u,i}$ is sentence B , labeled as *IsMatched* (1). While the negatives refer to the explanations randomly selected from the corpus, labeled as *NotMatched* (0). Here a special classification token $\langle CLS \rangle$ is introduced to indicate the label. The prediction probability is calculated as follows:

$$\hat{c}_{u,i} = \mathbf{w}^c \tanh(\mathbf{W}^c \mathbf{y}_L^3 + \mathbf{b}^c) + b^c, \quad (4)$$

where $\mathbf{W}^c \in \mathbb{R}^{d \times d}$, $\mathbf{b}^c \in \mathbb{R}^d$, $\mathbf{w}^c \in \mathbb{R}^{1 \times d}$ and $b^c \in \mathbb{R}$ are weight parameters, and $\tanh(\cdot)$ is the hyperbolic tangent function. The output representation is offset by 2 since $\langle CLS \rangle$ follows the user and item. The cross entropy loss is leveraged to measure the matched explanation prediction:

$$\mathcal{L}_{cls} = -\frac{1}{|\mathcal{T}|} \sum_{(u,i) \in \mathcal{T}} c_{u,i} \log \hat{c}_{u,i}, \quad (5)$$

where $c_{u,i}$ is ground truth $\langle CLS \rangle$ label and $\mathcal{T} = [\mathcal{T}_{pos}, \mathcal{T}_{neg}]$ is the entire training set.

2.2.2. Rating Prediction

The first position of the representation $\mathbf{y}_{L,1}$ is corresponding to rating prediction task. A multi-layer perceptron (MLP) is applied after \mathbf{y}_L^1 as follows:

$$\hat{r}_{u,i} = \mathbf{w}^r \sigma(\mathbf{W}^r \mathbf{y}_L^1 + \mathbf{b}^r) + b^r, \quad (6)$$

where $\mathbf{W}^r \in \mathbb{R}^{d \times d}$, $\mathbf{b}^r \in \mathbb{R}^d$, $\mathbf{w}^r \in \mathbb{R}^{1 \times d}$ and $b^r \in \mathbb{R}$ are weight parameters, and $\sigma(\cdot)$ is the sigmoid function. Mean Square Error (MSE) is utilized as rating loss function:

$$\mathcal{L}_r = \frac{1}{|\mathcal{T}_{pos}|} \sum_{(u,i) \in \mathcal{T}_{pos}} (r_{u,i} - \hat{r}_{u,i})^2, \quad (7)$$

where $r_{u,i}$ is the ground-truth ratings.

Table 2. Statistics of the datasets.

	TripAdvisor	Amazon-CSJ
#users	9765	38764
#items	6280	22919
#records / user	32.77	4.62
#records / item	50.96	7.82
#words / exp	13.01	10.48

2.2.3. Explanation Generation and Context Prediction

To enable the personalization of explanation generation, the context prediction task is incorporated to guide the explanation generation similar as [11]. The word probability distribution \mathbf{p}_k over vocabulary \mathcal{V} is computed as follows. Here the second position of \mathbf{y}_L^k is utilized for context prediction and k is in the range of $[4, |Y|]$ for explanation generation.

$$\mathbf{p}_k = \text{softmax}(\mathbf{W}^v \mathbf{y}_L^k + \mathbf{b}^v), \quad (8)$$

where $\mathbf{W}^v \in \mathbb{R}^{|\mathcal{V}| \times d}$, $\mathbf{b}^v \in \mathbb{R}^{|\mathcal{V}|}$ are weight parameters, $|\mathcal{V}|$ is the size of vocabulary \mathcal{V} . The Negative Log Likelihood Loss (NLL) loss is employed as the loss function, with the probability \mathbf{p}_k offset by three positions for calculating \mathcal{L}_w and \mathbf{p}_2 utilized for computing \mathcal{L}_{ctx} , defined as below:

$$\mathcal{L}_{w(ctx)} = \frac{1}{|\mathcal{T}_{pos}|} \sum_{(u,i) \in \mathcal{T}_{pos}} \frac{1}{|\mathcal{W}_{u,i}|} \sum_{k=1}^{|\mathcal{W}_{u,i}|} -\log \mathbf{p}_{3+k(2)}, \quad (9)$$

where \mathcal{T}_{pos} denotes the positive training set and $|\mathcal{W}_{u,i}|$ is the length of ground truth explanation. Note that the context prediction aims to map the user and item IDs onto the words in the explanation so as to make the personalized explanation.

2.3. Multi-task Learning

Finally, the above-mentioned sub-tasks are integrated into the multi-task learning framework, with the overall objective loss shown as follows:

$$\mathcal{L} = \min_{\theta} (\lambda_{cls} \mathcal{L}_{cls} + \lambda_r \mathcal{L}_r + \lambda_w \mathcal{L}_w + \lambda_{ctx} \mathcal{L}_{ctx}), \quad (10)$$

where λ_{cls} , λ_r , λ_w and λ_{ctx} denote trade-off weights and θ refer to the trainable parameters of ExBERT.

3. EXPERIMENTS

3.1. Datasets

The proposed model is evaluated on two publicly available datasets from different domains released by NETE [15]. They are Amazon-Clothing Shoes & Jewellery (ecommerce) and TripAdvisor (hotel). In each piece of record, it contains user ID and item ID, rating score (scaled from 1 to 5) and the real review text. The statistical features of the datasets are shown in Table 2.

Table 3. Performance comparison on Amazon-CSJ and TripAdvisor.

Metrics	Recommendation		Explainability			Text Quality								
	RMSE↓	MAE↓	FMR↑	FCR↑	DIV↓	USR↑	B1↑	B4↑	R1-P↑	R1-R↑	R1-F↑	R2-P↑	R2-R↑	R2-F↑
Amazon-CSJ														
NRT [9]	1.06	0.76	0.03	0.01	0.42	0.01	12.41	0.80	15.02	12.91	12.90	1.66	1.58	1.46
Att2Seq [10]	-	-	0.05	0.04	0.14	0.03	12.83	0.85	15.57	13.41	13.37	1.85	1.70	1.60
PETER [11]	1.05	0.83	0.10	0.14	0.14	0.11	13.32	0.96	16.36	14.46	14.25	2.20	1.93	1.84
ExBERT	1.03	0.80	0.30	0.32	0.05	0.43	18.77	2.74	25.90	21.94	22.16	6.64	5.58	5.54
TripAdvisor														
NRT [9]	0.79	0.61	0.06	0.09	4.27	0.08	15.05	0.99	18.22	14.39	15.40	2.29	1.98	2.01
Att2Seq [10]	-	-	0.06	0.15	4.32	0.17	15.27	1.03	18.97	14.72	15.92	2.40	2.03	2.09
PETER [11]	0.81	0.63	0.07	0.13	2.95	0.08	15.96	1.11	19.07	16.09	16.48	2.33	2.17	2.09
ExBERT	0.85	0.66	0.37	0.45	1.61	0.75	25.71	4.83	34.21	28.66	29.66	10.62	9.22	9.21

Table 4. Ablation studies on Amazon-CSJ. MEP refers to the matched explanation prediction task.

	Recommendation		Explainability			Text Quality								
	RMSE↓	MAE↓	FMR↑	FCR↑	DIV↓	USR↑	B1↑	B4↑	R1-P↑	R1-R↑	R1-F↑	R2-P↑	R2-R↑	R2-F↑
ExBERT (1:0.2)	1.03	0.81	0.29	0.30	0.06	0.40	18.65	2.63	25.40	21.62	21.79	6.42	5.36	5.32
ExBERT (1:0.5)	1.03	0.80	0.30	0.32	0.05	0.43	18.77	2.74	25.90	21.94	22.16	6.64	5.58	5.54
ExBERT (1:0.8)	1.03	0.80	0.30	0.33	0.05	0.44	18.73	2.80	25.86	21.85	22.13	6.53	5.55	5.51
ExBERT w/o Re-Ranking	1.06	0.85	0.08	0.07	0.35	0.05	14.42	0.83	15.53	15.19	14.41	1.67	1.69	1.52
ExBERT w/o MEP	1.05	0.85	0.28	0.21	0.06	0.35	17.63	2.46	25.77	20.88	21.55	6.41	5.11	5.21

3.2. Baselines

To evaluate the effectiveness of our model, we compare it with three state-of-the-art baselines.

- **NRT** [9] simultaneously predicts the rating and generates abstractive tips with a multi-layer perceptron network (MLP) and Gated Recurrent Unit (GRU).
- **Att2Seq** [10] encodes a set of attributes via MLP into latent factor, which is further decoded by stacked multiple layers of RNN for review generation. It is worth noting that Att2Seq only generates review sentences.
- **PETER** [11] is built on a small unpretrained transformer backbone. It bridges the personalization factor (i.e., user and item IDs) and review words via context prediction module.

3.3. Implementation Details

Following the experimental settings of [11], we split both datasets into training, validation, and testing sets with the ratio as 8:1:1. For fair comparison, the average results of five repetitive experiments are reported. In particular, low frequency words in the review texts are removed and a vocabulary dictionary \mathcal{V} of 20,000 common words is built for each dataset. The baselines and our method are implemented on Pytorch. The proposed ExBERT trains from scratch with a 2-layer light-weight BERT model [17] and stochastic gradient descent (SGD) is utilized as the optimizer [20]. The batch size is fixed to 128. The initial learning rate is set as 1.0 and is multiplied by 0.8 when the loss stops decreasing. And the training is terminated if the number of loss drop reaches 10.

3.4. Performance Comparison and Ablation Study

The experimental results are demonstrated in Table 3. It can be seen the proposed approach improved over the state-of-the-art with a significant margin in terms of all the metrics utilized. To validate the essence of the designed modules in our framework, we provide an ablation study conducted on Amazon-Clothing Shoes & Jewellery (Amazon-CSJ), as shown in Table 4. Different ratios of positive-negatives (see Sec 2.2.1) are shown: [1:0.2, 1:0.5, 1:0.8]. The results reveal that the best performance is achieved when ratio set as 1:0.5. We replace the top- K scored ranked review history (see Sec 2.1.1) with randomly ranked ones, denoted as “ExBERT w/o Re-Ranking”. We find the explainability and text quality of generated reviews as well as the recommendation performance drop dramatically. When removing the matched explanation prediction task (see Sec 2.2.1), denoted as w/o MEP, the performance over all metrics decreases consistently, which proves the effectiveness of classification task.

4. CONCLUSIONS AND FUTURE WORK

In this paper, we develop a novel BERT-guided explanation recommender system, ExBERT to simultaneously perform the rating prediction and explanation generation. Without the guidance of aspects, the proposed framework is capable of generating high-quality explanation sentence with finer-grained details. The historical reviews are explored to extract the user’s interests and item’s characteristics. A new matched explanation prediction task is proposed to distinguish positive and negative samples with the user and item ID as the personalization factor. Experimental results on two public datasets demonstrate the superiority of ExBERT in both recommendation and explanation.

5. REFERENCES

- [1] Hao Wang, Naiyan Wang, and Dit-Yan Yeung, “Collaborative deep learning for recommender systems,” in *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*, 2015, pp. 1235–1244.
- [2] Wei Zhang, Jianyong Wang, and Wei Feng, “Combining latent factor model with location features for event-based group recommendation,” in *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2013, pp. 910–918.
- [3] Christopher C Johnson, “Logistic matrix factorization for implicit feedback data,” *Advances in Neural Information Processing Systems*, vol. 27, no. 78, pp. 1–9, 2014.
- [4] Paul Covington, Jay Adams, and Emre Sargin, “Deep neural networks for youtube recommendations,” in *Proceedings of the 10th ACM conference on recommender systems*, 2016, pp. 191–198.
- [5] Hanxin Wang, Daichi Amagata, Takuya Makeawa, Takahiro Hara, Niu Hao, Kei Yonekawa, and Mori Kurokawa, “A dnn-based cross-domain recommender system for alleviating cold-start problem in e-commerce,” *IEEE Open Journal of the Industrial Electronics Society*, vol. 1, pp. 194–206, 2020.
- [6] Huijing Zhan, Jie Lin, Kenan Emir Ak, Boxin Shi, Ling-Yu Duan, and Alex C Kot, “ α 3-fkg: Attentive attribute-aware fashion knowledge graph for outfit preference prediction,” *IEEE Transactions on Multimedia*, vol. 24, pp. 819–831, 2021.
- [7] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio, “Graph attention networks,” *arXiv preprint arXiv:1710.10903*, 2017.
- [8] Huijing Zhan, Ling Li, Xue Geng, Jie Lin, and AC Kot, “Rule-guided knowledge-graph based negative sampling for outfit recommendation,” 2022.
- [9] Piji Li, Zihao Wang, Zhaochun Ren, Lidong Bing, and Wai Lam, “Neural rating regression with abstractive tips generation for recommendation,” in *Proceedings of the 40th International ACM SIGIR conference on Research and Development in Information Retrieval*, 2017, pp. 345–354.
- [10] Li Dong, Shaohan Huang, Furu Wei, Mirella Lapata, Ming Zhou, and Ke Xu, “Learning to generate product reviews from attributes,” in *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, 2017, pp. 623–632.
- [11] Lei Li, Yongfeng Zhang, and Li Chen, “Personalized transformer for explainable recommendation,” *arXiv preprint arXiv:2105.11601*, 2021.
- [12] Yongfeng Zhang, Guokun Lai, Min Zhang, Yi Zhang, Yiqun Liu, and Shaoping Ma, “Explicit factor models for explainable recommendation based on phrase-level sentiment analysis,” in *Proceedings of the 37th international ACM SIGIR conference on Research & development in information retrieval*, 2014, pp. 83–92.
- [13] Lei Li, Li Chen, and Ruihai Dong, “Caesar: context-aware explanation based on supervised attention for service recommendations,” *Journal of Intelligent Information Systems*, vol. 57, no. 1, pp. 147–170, 2021.
- [14] Juntao Tan, Shuyuan Xu, Yingqiang Ge, Yunqi Li, Xu Chen, and Yongfeng Zhang, “Counterfactual explainable recommendation,” in *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, 2021, pp. 1784–1793.
- [15] Lei Li, Yongfeng Zhang, and Li Chen, “Generate neutral template explanations for recommendation,” in *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, 2020, pp. 755–764.
- [16] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [17] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” *arXiv preprint arXiv:1810.04805*, 2018.
- [18] Jianmo Ni, Jiacheng Li, and Julian McAuley, “Justifying recommendations using distantly-labeled reviews and fine-grained aspects,” in *Proceedings of the 2019 conference on empirical methods in natural language processing and the 9th international joint conference on natural language processing (EMNLP-IJCNLP)*, 2019, pp. 188–197.
- [19] Nils Reimers and Iryna Gurevych, “Sentence-bert: Sentence embeddings using siamese bert-networks,” *arXiv preprint arXiv:1908.10084*, 2019.
- [20] Herbert Robbins and Sutton Monroe, “A stochastic approximation method,” *The annals of mathematical statistics*, pp. 400–407, 1951.