
Appendix for “REASONER: An Explainable Recommendation Dataset with Multi-aspect Ground Truths”

Xu Chen^{1,2}, Jingsen Zhang^{1,2,*}, Lei Wang^{1,2,*}, Quanyu Dai³, Zhenhua Dong³,
Ruiming Tang³, Rui Zhang⁴, Li Chen⁵, Wayne Xin Zhao^{1,2}, Ji-Rong Wen^{1,2}

¹Beijing Key Laboratory of Big Data Management and Analysis Methods, Beijing, China

²Gaoling School of Artificial Intelligence, Renmin University of China, Beijing, China

³Huawei Noah’s Ark Lab, China ⁴www.ruizhang.info

⁵Department of Computer Science, Hong Kong Baptist University, Hong Kong

1	Contents	
2	A Accessibility	2
3	B Library	2
4	B.1 The Structure of the Library	2
5	B.2 The Implemented Models	3
6	B.3 Examples for Using Our Library	4
7	C Benchmark	5
8	C.1 Experiment Setup	5
9	C.2 Experiment Results	5

*Equal contribution.

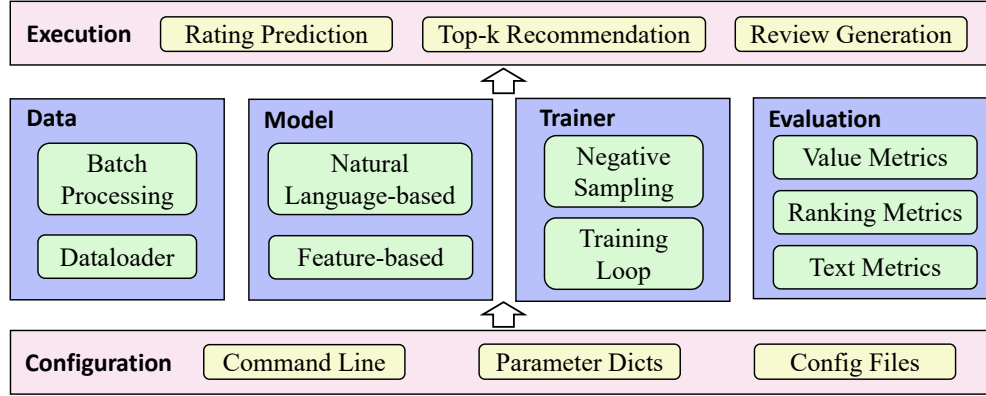


Figure 1: The structure of our library. There are six parts including the configuration, data, model, trainer, evaluation and execution modules.

10 A Accessibility

11 All the resources in our work are accessible at <https://reasoner2023.github.io>, providing comprehensive information about the dataset, library, related documents and the repository with long-term availability. Our licensing for the dataset is under a CC BY-NC 4.0 (Creative Commons Attribution-NonCommercial 4.0). See official instructions here². We will consistently maintain and update the resources to ensure the long-term usability.

16 B Library

17 B.1 The Structure of the Library

18 We show the structure of our library in Figure 1. The configuration module is the base part of the library and responsible for initializing all the parameters. We support three methods to specify the parameters, that is, the command line, parameter dictionary and configuration file. Based on the configuration module, there are four higher layer modules, that is,

22 **Data module:** this module converts the raw data into the model inputs. There are two components: the first one is responsible for loading the data and building vocabularies for the user reviews. The second part aims to process the data into the formats required by the model inputs, and generate the sample batches for model optimization.

26 **Model module:** this module implements the explainable recommender models. There are two types of methods in our library. The first one includes the feature-based explainable recommender models, and the second one contains the models with natural language explanations. We delay the detailed introduction of these models in the next section.

30 **Trainer module:** this module is leveraged to implement the training losses, such as the Bayesian Personalized Ranking (BPR) and Binary Cross Entropy (BCE). In addition, this module can also record the complete model training process.

33 **Evaluation module:** this module is designed to evaluate different models, and there are three types of evaluation tasks, that is, rating prediction, top-k recommendation and review generation.

35 Upon the above four modules, there is an execution module on the upper-most layer. It is responsible for optimizing the recommender model for different tasks, such as rating prediction, tag prediction and review generation. For more detailed introduction on our library architecture, we refer the readers to our project at <https://reasoner2023.github.io/>.

²<https://creativecommons.org/licenses/by-nc/4.0/>

39 B.2 The Implemented Models

40 In our library, we implement two types of explainable recommender models, which are widely studied
41 in the research community. The first one are feature-based explainable recommender models, where
42 the features can be the tags, item aspects and so on. The second one are the models with natural
43 language explanations.

44 More specifically, we implement the following representative feature-based explainable recommender
45 models:

46 **EFM** [15] predicts the user preferences and generates explainable recommendations based on explicit
47 product features and user opinions from the review information.

48 **TriRank** [5] models the user-item-aspect ternary relation as a heterogeneous tripartite graph based
49 on user ratings and reviews, and it devises a vertex ranking algorithm for recommendation.

50 **LRPPM** [3] is a tensor-matrix factorization algorithm which captures the user preferences using
51 ranking-based optimization objective over various item aspects.

52 **SULM** [1] enhances recommendations by recommending not only item but also the specific aspects
53 by using aspect-level sentiment analysis.

54 **MTER** [14] is a tensor factorization method which models the task of item recommendation using
55 a three-way tensor over the users, items and features. We omit the modeling of the opinions in the
56 original implementation for adapting our data.

57 **AMF** [6] improves the recommendation accuracy by using the auxiliary information extracted from
58 the user review aspects.

59 **TRDM** [18] introduces a two-stage approach to generate accurate item recommendations and effective
60 tag-based potential features simultaneously for enhancing recommendation accuracy and diversity.

61 **TRAL** [17] proposes attention-based learning to capture diverse tag-based features, and compress
62 these features with an attention pooling layer to enhance recommendation accuracy.

63 **HPTR** [16] employs hyperbolic distance to measure semantic relevance between entities, which
64 better captures hierarchical structures presented in tag information.

65 **AIRec** [2] enhances tag-aware recommender system by employing a hierarchical attention network to
66 capture multi-aspect preferences and leveraging tag intersection to improve conjunct feature learning.

67 **HAN-TR** [13] captures distinct user preferences and informative elements by employing separate
68 attention networks for element-level influence and information-level attentiveness.

69 **TNAM** [7] addresses the issues of tag weight assignment in recommender systems by introducing a
70 tag-based neural attention network that captures users' specific tag attention.

71 **BPR-T** [8] addresses high dimension and sparsity issues of tagging information by integrating tag
72 mapping into a Bayesian personalized ranking collaborative filtering model.

73 In addition to the above shallow models based on matrix factorization, we also implement the
74 following deep feature-based explainable recommender models (called **DERM** for short):

75 **DERM-MLP** is a deep recommender model for jointly predicting the ratings and tags. The two tasks
76 share the set of user/item/tag embeddings. The hidden states as well as the tag embeddings are put
77 into different layers corresponding to the different tasks.

78 **DERM-MF** firstly obtains a hidden state based on the user/item embeddings using matrix factoriza-
79 tion, and then the outputs are computed by a neural network.

80 **DERM-C** combines matrix factorization and Multi-Layer Perceptron (MLP) to derive the hidden
81 states, and the outputs are merged in a concatenated manner.

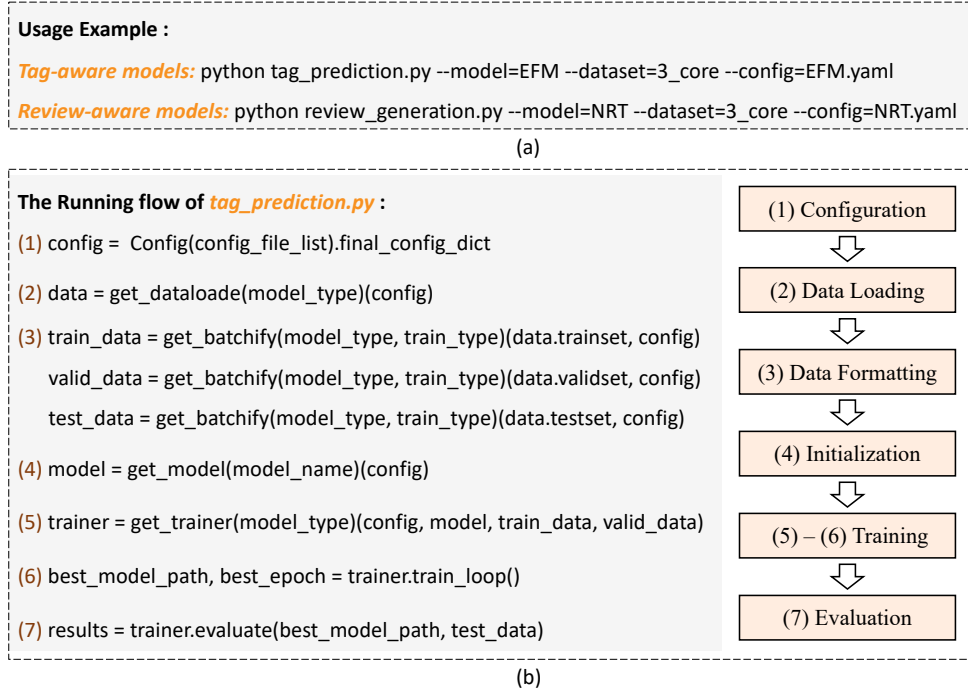


Figure 2: (a) Examples of the feature- and review-based models. (b) The running flow of our library.

82 **DERM-H** leverages the tags to profile the users and items, and then use the same architecture as
 83 **DERM-MLP** for predicting the ratings and tags.

84 For the models with natural language explanations, we implement the following representative
 85 methods:

86 **Att2Seq** [4] is a review generation model which uses LSTM as the decoder, and output the texts
 87 directly based on the user/item IDs and rating information.

88 **NRT** [10] simultaneously predicts the reviews and ratings based on the input user-item pair, where
 89 the two tasks share the same embedding and hidden layers.

90 **PETER** [9] leverages Transformer to generate the user reviews, which is a state-of-the-art review
 91 generation model.

92 **B.3 Examples for Using Our Library**

93 In this section, we introduce how to use our library. We present a simple example in Figure 2(a),
 94 where one can directly execute *tag_prediction.py* or *review_generate.py* to run a feature-based or
 95 review-based model, respectively. In each of these commands, one needs to specify three parameters
 96 to indicate the names of the model, dataset and configuration file, respectively.

97 Take *tag_prediction.py* for example, it sequentially executes the following steps: (1) Configuration.
 98 In this step, the parameters related to the model architecture and optimization process from different
 99 sources (command line, configuration dictionary and files) are integrated into a dictionary. (2) Data
 100 loading. In this step, the data loader is selected according to model type. For the review-aware models,
 101 this step reads all the records and build the vocabulary. (3) Data Formatting. The training, validation
 102 and test sets are processed into the formats required by the model input in a sample batch manner.
 103 (4) Initialization. The corresponding model class will be defined and initialized according to the
 104 parameter values in configuration. (5)-(6) Training. Selecting the optimization approach to train the
 105 model. (7) Evaluation. Measuring the model performance on different tasks.

Table 1: **Statistics of the datasets.**

Dataset	REASONER
# Users	2,997
# Items	4,672
# Tags	6,115
# Interactions	58,497
Avg. # words / review	17

106 Our library is highly extensible, and there are three steps to realize a new model: (1) implementing
 107 the basic functions of the model, including the model architecture, preference score prediction etc.
 108 (2) Customizing the training approaches in *train.py*. (3) Indicating the parameters in the config file.

109 C Benchmark

110 C.1 Experiment Setup

111 Considering that we have three types of ground truths for the explanations, we evaluate the model
 112 performance by predicting the tags for Q1, Q2, Q3, Q1+Q2, Q2+Q3, Q1+Q3 and Q1+Q2+Q3,
 113 respectively. When we have to predict multiple types of ground truths, we extend the original models
 114 to their multi-task versions by sharing the embedding parameters. We randomly split the dataset into
 115 the training, validation and testing sets according to the ratio of 8:1:1. For the review generation task,
 116 we use the most 20,000 frequently mentioned words to construct the vocabulary, and the maximum
 117 length of the generated sentences is set to 17, which is equal to the average length of reviews in
 118 dataset. **The dataset statistics are presented in Table 1.** For all the models, the batch size is set as
 119 256. We tune the other key hyper-parameters by grid search. In specific, we tune the learning rate
 120 and the weight of L2 regularization in the range of [0.1, 0.01, 0.001, 0.0001] and [0.001, 0.0001, 0]
 121 respectively. For the deep models, we tune the hidden size and the layer number in the range of [32,
 122 64, 128, 256] and [1, 2, 3, 4] respectively. More details of the experiment setting are shown in our
 123 project, which has been released at <https://reasoner2023.github.io/>. We use RMSE and MAE as the
 124 metrics to evaluate the performance of the rating prediction task. For the task of tag prediction, F1
 125 and NDCG are selected to evaluate the model performance. To evaluate the quality of the generated
 126 reviews, we leverage the metrics including BLEU [12] and ROUGE [11] for model comparison.

127 C.2 Experiment Results

128 The comparison results of the feature-based explainable recommender models on the tasks of tag
 129 and rating predictions are presented in Table 2-8. The comparison results of the models with natural
 130 language explanations on the task of review generation and rating prediction are presented in Table 9.
 131 We use the tags with top 10 prediction scores to calculate F1 and NDCG, and the results are percentage
 132 values with "%" omitted. For RMSE and MAE, a lower value indicates better performance. For
 133 each evaluation metric, we use bold fonts to label the best performance. Since the TriRank we
 134 implemented does not support to predict multiple types of tags simultaneously, we omit it in the
 135 corresponding tables.

Table 2: The benchmarking results of the feature-based explainable recommender models on predicting the tags for persuasiveness and ratings.

Metrics	Persuasiveness		Rating Prediction	
	F1	NDCG	RMSE	MAE
EFM	26.99 \pm 0.35	18.89 \pm 0.23	1.68 \pm 0.00	1.24 \pm 0.01
TriRank	18.36 \pm 0.07	13.98 \pm 0.06	2.90 \pm 0.00	2.58 \pm 0.00
LRPPM	37.31 \pm 0.23	23.25 \pm 0.10	1.22 \pm 0.00	0.96 \pm 0.00
SULM	41.68 \pm 0.63	25.77 \pm 0.22	1.65 \pm 0.08	1.30 \pm 0.06
MTER	5.66 \pm 1.74	2.65 \pm 1.13	2.27 \pm 0.64	1.96 \pm 0.62
AMF	27.93 \pm 0.08	17.62 \pm 0.17	2.28 \pm 0.00	1.86 \pm 0.00
DERM-MLP	37.74 \pm 0.12	23.45 \pm 0.02	1.30 \pm 0.01	1.06 \pm 0.01
DERM-MF	36.57 \pm 0.14	21.49 \pm 0.17	1.32 \pm 0.00	1.14 \pm 0.00
DERM-C	37.17 \pm 0.19	23.21 \pm 0.01	1.30 \pm 0.00	1.07 \pm 0.01
DERM-H	35.58 \pm 0.07	22.02 \pm 0.18	1.27 \pm 0.01	1.04 \pm 0.01
TRDM	30.24 \pm 1.57	13.44 \pm 1.11	1.19 \pm 0.00	0.95 \pm 0.01
TRAL	5.93 \pm 0.05	2.61 \pm 0.02	1.29 \pm 0.00	1.08 \pm 0.00
HPTR	38.61 \pm 0.20	23.05 \pm 1.28	2.16 \pm 0.71	1.88 \pm 0.68
AIRec	38.05 \pm 0.07	23.07 \pm 0.03	1.31 \pm 0.00	1.08 \pm 0.01
HAN-TR	34.96 \pm 0.60	18.58 \pm 3.37	2.10 \pm 0.80	1.83 \pm 0.75
TNAM	5.97 \pm 0.01	2.60 \pm 0.01	1.37 \pm 0.01	1.17 \pm 0.01
BPR-T	32.70 \pm 0.62	18.32 \pm 0.82	1.23 \pm 0.00	0.98 \pm 0.01

Table 3: The benchmarking results of the feature-based explainable recommender models on predicting the tags for informativeness and ratings.

Metrics	Informativeness		Rating Prediction	
	F1	NDCG	RMSE	MAE
EFM	5.38 \pm 0.28	3.97 \pm 0.19	1.68 \pm 0.00	1.24 \pm 0.01
TriRank	18.78 \pm 0.10	14.50 \pm 0.09	2.90 \pm 0.00	2.58 \pm 0.00
LRPPM	37.85 \pm 0.22	38.35 \pm 0.15	1.22 \pm 0.00	0.96 \pm 0.00
SULM	43.25 \pm 0.59	42.97 \pm 0.40	1.65 \pm 0.08	1.30 \pm 0.06
MTER	8.40 \pm 0.86	6.13 \pm 0.68	2.04 \pm 0.68	1.74 \pm 0.66
AMF	28.63 \pm 0.24	28.95 \pm 0.29	2.28 \pm 0.00	1.86 \pm 0.00
DERM-MLP	38.60 \pm 0.06	38.99 \pm 0.03	1.30 \pm 0.01	1.06 \pm 0.01
DERM-MF	37.10 \pm 0.09	35.37 \pm 0.14	1.32 \pm 0.00	1.14 \pm 0.00
DERM-C	37.96 \pm 0.10	38.43 \pm 0.05	1.30 \pm 0.00	1.07 \pm 0.01
DERM-H	36.36 \pm 0.55	35.85 \pm 0.41	1.29 \pm 0.01	1.06 \pm 0.01
TRDM	31.49 \pm 2.84	24.09 \pm 2.80	1.19 \pm 0.00	0.95 \pm 0.01
TRAL	6.04 \pm 0.06	4.59 \pm 0.01	1.28 \pm 0.00	1.08 \pm 0.00
HPTR	39.55 \pm 0.09	37.73 \pm 1.88	1.27 \pm 0.07	1.05 \pm 0.09
AIRec	38.90 \pm 0.06	39.33 \pm 0.06	1.31 \pm 0.00	1.08 \pm 0.01
HAN-TR	34.95 \pm 1.70	30.69 \pm 5.89	2.10 \pm 0.79	1.83 \pm 0.75
TNAM	37.77 \pm 0.35	37.97 \pm 0.20	1.36 \pm 0.00	1.16 \pm 0.00
BPR-T	33.83 \pm 0.43	31.15 \pm 0.33	1.23 \pm 0.01	0.98 \pm 0.01

Table 4: The benchmarking results of the feature-based explainable recommender models on predicting the tags for satisfaction and ratings.

Metrics	Satisfaction		Rating Prediction	
	F1	NDCG	RMSE	MAE
EFM	4.57 \pm 0.46	1.79 \pm 0.19	1.68 \pm 0.00	1.24 \pm 0.01
TriRank	16.82 \pm 0.03	13.16 \pm 0.02	2.90 \pm 0.00	2.58 \pm 0.00
LRPPM	36.04 \pm 0.16	22.35 \pm 0.08	1.22 \pm 0.00	0.96 \pm 0.00
SULM	40.46 \pm 0.62	24.80 \pm 0.19	1.64 \pm 0.09	1.29 \pm 0.07
MTER	5.97 \pm 1.92	2.85 \pm 1.07	2.26 \pm 0.65	1.96 \pm 0.62
AMF	27.16 \pm 0.19	17.05 \pm 0.21	2.28 \pm 0.00	1.86 \pm 0.00
DERM-MLP	36.76 \pm 0.07	22.53 \pm 0.12	1.30 \pm 0.01	1.06 \pm 0.01
DERM-MF	35.40 \pm 0.23	20.59 \pm 0.35	1.32 \pm 0.00	1.14 \pm 0.00
DERM-C	36.20 \pm 0.28	22.28 \pm 0.22	1.29 \pm 0.01	1.07 \pm 0.01
DERM-H	34.65 \pm 0.43	21.33 \pm 0.47	1.28 \pm 0.01	1.05 \pm 0.02
TRDM	31.29 \pm 0.63	14.74 \pm 0.69	1.19 \pm 0.00	0.95 \pm 0.01
TRAL	5.89 \pm 0.05	2.52 \pm 0.01	1.29 \pm 0.00	1.08 \pm 0.00
HPTR	35.35 \pm 2.90	17.78 \pm 3.92	1.81 \pm 0.77	1.56 \pm 0.72
AIRec	37.10 \pm 0.09	22.86 \pm 0.08	1.30 \pm 0.01	1.08 \pm 0.01
HAN-TR	33.95 \pm 1.80	18.06 \pm 4.15	2.10 \pm 0.80	1.83 \pm 0.75
TNAM	5.89 \pm 0.05	2.52 \pm 0.01	1.37 \pm 0.00	1.17 \pm 0.00
BPR-T	33.82 \pm 0.29	19.52 \pm 0.25	1.23 \pm 0.00	0.98 \pm 0.00

Table 5: The benchmarking results of the feature-based explainable recommender models on jointly predicting the tags for persuasiveness, informativeness and ratings.

Metrics	Persuasiveness		Informativeness		Rating	
	F1	NDCG	F1	NDCG	RMSE	MAE
EFM	15.69 \pm 0.03	12.74 \pm 0.07	5.38 \pm 0.73	3.94 \pm 0.42	1.66 \pm 0.00	1.23 \pm 0.00
LRPPM	37.32 \pm 0.21	23.26 \pm 0.09	37.89 \pm 0.19	38.37 \pm 0.13	1.22 \pm 0.00	0.96 \pm 0.00
SULM	41.34 \pm 0.53	25.68 \pm 0.20	42.70 \pm 0.50	42.82 \pm 0.35	1.67 \pm 0.06	1.31 \pm 0.05
MTER	35.53 \pm 0.21	21.88 \pm 0.35	36.22 \pm 0.30	36.09 \pm 0.61	1.36 \pm 0.01	1.09 \pm 0.01
AMF	27.67 \pm 0.16	17.45 \pm 0.15	28.23 \pm 0.33	28.57 \pm 0.34	2.28 \pm 0.00	1.86 \pm 0.00
DERM-MLP	38.49 \pm 0.15	23.80 \pm 0.10	39.14 \pm 0.10	39.38 \pm 0.10	1.30 \pm 0.01	1.06 \pm 0.01
DERM-MF	36.84 \pm 0.04	22.55 \pm 0.05	37.58 \pm 0.11	37.44 \pm 0.08	1.32 \pm 0.00	1.15 \pm 0.00
DERM-C	37.85 \pm 0.26	23.43 \pm 0.21	38.82 \pm 0.06	39.10 \pm 0.09	1.30 \pm 0.01	1.08 \pm 0.01
DERM-H	37.47 \pm 0.26	23.23 \pm 0.22	38.17 \pm 0.22	38.32 \pm 0.36	1.28 \pm 0.00	1.04 \pm 0.01
TRDM	32.91 \pm 1.06	15.51 \pm 0.72	33.01 \pm 0.50	25.02 \pm 0.37	1.19 \pm 0.00	0.94 \pm 0.01
TRAL	5.93 \pm 0.05	2.61 \pm 0.02	6.06 \pm 0.06	4.60 \pm 0.02	1.29 \pm 0.00	1.08 \pm 0.00
HPTR	38.68 \pm 0.25	23.11 \pm 1.20	36.20 \pm 4.79	31.40 \pm 9.27	2.16 \pm 0.71	1.88 \pm 0.68
AIRec	38.73 \pm 0.10	23.97 \pm 0.07	39.43 \pm 0.08	39.66 \pm 0.09	1.31 \pm 0.00	1.08 \pm 0.01
HAN-TR	33.32 \pm 0.75	16.83 \pm 0.53	33.31 \pm 0.93	27.72 \pm 0.46	1.29 \pm 0.02	1.06 \pm 0.03
TNAM	5.91 \pm 0.05	2.63 \pm 0.02	37.47 \pm 0.46	37.35 \pm 0.58	1.36 \pm 0.01	1.17 \pm 0.01
BPR-T	33.15 \pm 0.22	18.90 \pm 0.44	34.05 \pm 0.20	31.51 \pm 0.26	1.24 \pm 0.01	0.99 \pm 0.01

Table 6: The benchmarking results of the feature-based explainable recommender models on jointly predicting the tags for persuasiveness, satisfaction and ratings.

Metrics	Persuasiveness		Satisfaction		Rating	
	F1	NDCG	F1	NDCG	RMSE	MAE
EFM	15.58 \pm 0.03	12.84 \pm 0.07	4.58 \pm 0.45	1.77 \pm 0.18	1.66 \pm 0.00	1.23 \pm 0.00
LRPPM	37.32 \pm 0.21	23.26 \pm 0.09	36.06 \pm 0.14	22.37 \pm 0.07	1.22 \pm 0.00	0.96 \pm 0.00
SULM	41.36 \pm 0.56	35.71 \pm 0.22	40.15 \pm 0.54	24.70 \pm 0.18	1.67 \pm 0.06	1.31 \pm 0.05
MTER	5.83 \pm 0.52	2.56 \pm 0.18	5.36 \pm 0.17	2.23 \pm 0.12	2.04 \pm 0.68	1.74 \pm 0.66
AMF	27.71 \pm 0.21	17.52 \pm 0.21	26.90 \pm 0.18	16.94 \pm 0.16	2.28 \pm 0.00	1.86 \pm 0.00
DERM-MLP	38.37 \pm 0.09	23.71 \pm 0.12	37.32 \pm 0.02	22.87 \pm 0.03	1.30 \pm 0.01	1.06 \pm 0.01
DERM-MF	36.90 \pm 0.12	22.63 \pm 0.12	35.78 \pm 0.10	21.74 \pm 0.12	1.32 \pm 0.00	1.15 \pm 0.00
DERM-C	38.03 \pm 0.11	23.60 \pm 0.06	36.95 \pm 0.10	22.72 \pm 0.06	1.31 \pm 0.01	1.08 \pm 0.00
DERM-H	37.49 \pm 0.24	23.32 \pm 0.18	36.11 \pm 0.29	22.34 \pm 0.17	1.29 \pm 0.01	1.04 \pm 0.01
TRDM	32.57 \pm 1.92	15.11 \pm 1.21	30.91 \pm 1.77	13.86 \pm 1.40	1.19 \pm 0.00	0.94 \pm 0.00
TRAL	5.91 \pm 0.05	2.63 \pm 0.02	5.89 \pm 0.05	2.52 \pm 0.01	1.29 \pm 0.00	1.08 \pm 0.00
HPTR	38.64 \pm 0.20	22.93 \pm 1.46	32.96 \pm 4.03	14.78 \pm 4.06	2.16 \pm 0.71	1.88 \pm 0.68
AIRec	38.69 \pm 0.06	23.94 \pm 0.05	37.63 \pm 0.05	23.08 \pm 0.02	1.30 \pm 0.01	1.07 \pm 0.01
HAN-TR	35.95 \pm 2.15	19.75 \pm 3.93	34.95 \pm 2.02	19.13 \pm 3.74	2.09 \pm 0.80	1.83 \pm 0.75
TNAM	5.91 \pm 0.05	2.63 \pm 0.02	5.89 \pm 0.05	2.52 \pm 0.01	1.39 \pm 0.02	1.17 \pm 0.02
BPR-T	33.07 \pm 0.18	18.93 \pm 0.48	33.75 \pm 0.26	19.57 \pm 0.28	1.23 \pm 0.00	0.99 \pm 0.01

Table 7: The benchmarking results of the feature-based explainable recommender models on jointly predicting the tags for informativeness, satisfaction and ratings.

Metrics	Informativeness		Satisfaction		Rating	
	F1	NDCG	F1	NDCG	RMSE	MAE
EFM	5.15 \pm 0.79	3.75 \pm 0.38	4.57 \pm 0.54	1.75 \pm 0.21	1.66 \pm 0.00	1.23 \pm 0.00
LRPPM	37.89 \pm 0.19	38.37 \pm 0.14	36.06 \pm 0.14	22.37 \pm 0.07	1.22 \pm 0.00	0.96 \pm 0.00
SULM	42.70 \pm 0.49	42.84 \pm 0.37	40.12 \pm 0.51	24.70 \pm 0.17	1.67 \pm 0.06	1.31 \pm 0.05
MTER	6.13 \pm 0.03	4.56 \pm 0.18	5.64 \pm 0.66	2.37 \pm 0.30	2.04 \pm 0.68	1.74 \pm 0.66
AMF	28.17 \pm 0.28	28.53 \pm 0.32	26.79 \pm 0.06	16.89 \pm 0.09	2.28 \pm 0.00	1.86 \pm 0.00
DERM-MLP	39.23 \pm 0.08	39.45 \pm 0.02	37.40 \pm 0.07	22.93 \pm 0.06	1.30 \pm 0.01	1.06 \pm 0.01
DERM-MF	37.60 \pm 0.13	37.49 \pm 0.17	35.77 \pm 0.16	21.76 \pm 0.16	1.32 \pm 0.00	1.15 \pm 0.00
DERM-C	38.77 \pm 0.13	39.08 \pm 0.18	36.84 \pm 0.23	22.59 \pm 0.16	1.30 \pm 0.01	1.07 \pm 0.01
DERM-H	38.13 \pm 0.50	38.45 \pm 0.47	36.44 \pm 0.28	22.55 \pm 0.21	1.27 \pm 0.01	1.04 \pm 0.01
TRDM	33.15 \pm 0.98	25.24 \pm 1.92	30.42 \pm 1.50	13.83 \pm 1.22	1.19 \pm 0.00	0.94 \pm 0.01
TRAL	6.04 \pm 0.06	4.56 \pm 0.02	5.84 \pm 0.01	2.43 \pm 0.08	1.29 \pm 0.00	1.08 \pm 0.00
HPTR	37.36 \pm 3.15	32.22 \pm 7.55	35.35 \pm 2.93	17.79 \pm 3.92	1.81 \pm 0.77	1.56 \pm 0.72
AIRec	39.46 \pm 0.11	39.72 \pm 0.07	37.65 \pm 0.11	23.13 \pm 0.07	1.30 \pm 0.00	1.08 \pm 0.01
HAN-TR	35.79 \pm 3.09	32.11 \pm 6.85	34.04 \pm 2.94	18.48 \pm 4.06	1.31 \pm 0.01	1.10 \pm 0.02
TNAM	37.01 \pm 1.35	37.01 \pm 1.33	5.89 \pm 0.05	2.52 \pm 0.01	1.36 \pm 0.01	1.15 \pm 0.01
BPR-T	34.11 \pm 0.32	31.57 \pm 0.50	33.94 \pm 0.18	19.75 \pm 0.28	1.23 \pm 0.00	0.98 \pm 0.01

Table 8: The benchmarking results of the feature-based explainable recommender models on jointly predicting the tags for persuasiveness, informativeness and satisfaction and ratings.

Metrics	Persuasiveness		Informativeness		Satisfaction		Rating Prediction	
	F1	NDCG	F1	NDCG	F1	NDCG	RMSE	MAE
EFM	11.66 \pm 0.15	8.52 \pm 0.10	4.97 \pm 0.61	3.70 \pm 0.48	5.33 \pm 1.10	2.24 \pm 0.62	1.66 \pm 0.01	1.23 \pm 0.01
LRPPM	37.32 \pm 0.24	23.26 \pm 0.11	37.94 \pm 0.23	38.46 \pm 0.18	36.08 \pm 0.17	22.35 \pm 0.10	1.22 \pm 0.00	0.96 \pm 0.00
SULM	41.12 \pm 0.50	25.35 \pm 0.21	42.35 \pm 0.45	42.66 \pm 0.34	40.00 \pm 0.49	24.66 \pm 0.16	1.69 \pm 0.09	1.33 \pm 0.08
MTER	36.16 \pm 0.06	22.38 \pm 0.15	36.75 \pm 0.09	36.94 \pm 0.24	34.84 \pm 0.09	21.52 \pm 0.03	1.34 \pm 0.04	1.08 \pm 0.03
AMF	27.83 \pm 0.37	17.47 \pm 0.19	28.35 \pm 0.33	28.66 \pm 0.34	27.09 \pm 0.34	17.03 \pm 0.16	2.28 \pm 0.00	1.86 \pm 0.00
DERM-MLP	38.60 \pm 0.08	23.81 \pm 0.07	39.33 \pm 0.09	39.57 \pm 0.05	37.52 \pm 0.09	22.97 \pm 0.09	1.31 \pm 0.02	1.07 \pm 0.02
DERM-MF	37.42 \pm 0.21	23.16 \pm 0.08	38.26 \pm 0.13	38.46 \pm 0.16	36.60 \pm 1.01	22.18 \pm 0.19	1.33 \pm 0.00	1.16 \pm 0.00
DERM-C	38.05 \pm 0.22	23.53 \pm 0.07	39.03 \pm 0.15	39.29 \pm 0.11	37.19 \pm 0.15	22.79 \pm 0.08	1.30 \pm 0.01	1.08 \pm 0.01
DERM-H	37.64 \pm 0.24	23.36 \pm 0.18	38.52 \pm 0.44	38.83 \pm 0.39	36.70 \pm 0.40	22.60 \pm 0.16	1.28 \pm 0.01	1.05 \pm 0.02
TRDM	33.50 \pm 1.85	15.64 \pm 1.83	33.94 \pm 1.06	26.47 \pm 1.89	31.79 \pm 1.17	14.77 \pm 1.09	1.19 \pm 0.00	0.94 \pm 0.01
TRAL	5.88 \pm 0.14	2.56 \pm 0.04	5.91 \pm 0.09	4.48 \pm 0.17	35.42 \pm 0.20	21.09 \pm 0.08	1.29 \pm 0.00	1.07 \pm 0.00
HPTR	38.82 \pm 0.38	22.14 \pm 1.49	39.31 \pm 0.50	38.01 \pm 2.21	37.64 \pm 0.26	21.46 \pm 1.51	1.77 \pm 0.98	1.50 \pm 0.93
AIRec	38.89 \pm 0.09	23.98 \pm 0.05	39.61 \pm 0.11	39.82 \pm 0.09	37.85 \pm 0.05	23.16 \pm 0.06	1.30 \pm 0.00	1.07 \pm 0.01
HAN-TR	37.95 \pm 0.63	23.46 \pm 0.02	38.65 \pm 0.70	39.01 \pm 0.39	37.07 \pm 0.82	22.74 \pm 0.04	1.80 \pm 0.67	1.57 \pm 0.62
TNAM	37.96 \pm 0.13	23.51 \pm 0.02	38.74 \pm 0.23	38.81 \pm 0.17	5.89 \pm 0.08	2.44 \pm 0.12	1.37 \pm 0.01	1.17 \pm 0.01
BPR-T	33.41 \pm 0.44	19.39 \pm 0.48	34.37 \pm 0.28	32.15 \pm 0.36	34.29 \pm 0.15	20.11 \pm 0.22	1.24 \pm 0.01	1.00 \pm 0.02

Table 9: The benchmarking results of the models with natural language explanations in our library. For BLEU and ROUGE, the results are percentage values with "%" omitted. "-" means the evaluation metric is not available for the model.

Metrics	BLEU (%)		ROUGE-1 (%)			ROUGE-2 (%)		
	B-1	B-4	F1	R	P	F1	R	P
Att2Seq	19.96 \pm 0.27	3.25 \pm 0.19	22.13 \pm 0.27	19.73 \pm 0.44	26.40 \pm 0.78	5.56 \pm 0.08	5.19 \pm 0.16	6.26 \pm 0.24
NRT	17.67 \pm 1.10	2.92 \pm 0.65	20.60 \pm 0.57	16.04 \pm 1.27	30.02 \pm 2.40	5.23 \pm 0.56	4.20 \pm 0.77	7.33 \pm 0.42
PETER	17.65 \pm 1.18	2.35 \pm 0.35	20.00 \pm 1.07	15.68 \pm 1.62	28.59 \pm 1.42	4.99 \pm 0.48	3.95 \pm 0.58	7.00 \pm 0.57

References

- [1] Konstantin Bauman, Bing Liu, and Alexander Tuzhilin. Aspect based recommendations: Recommending items with the most valuable aspects based on user reviews. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 717–725, 2017.
- [2] Bo Chen, Yue Ding, Xin Xin, Yunzhe Li, Yule Wang, and Dong Wang. Airec: Attentive intersection model for tag-aware recommendation. *Neurocomputing*, 421:105–114, 2021.
- [3] Xu Chen, Zheng Qin, Yongfeng Zhang, and Tao Xu. Learning to rank features for recommendation over multiple categories. In *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*, pages 305–314, 2016.
- [4] 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*, pages 623–632, 2017.
- [5] Xiangnan He, Tao Chen, Min-Yen Kan, and Xiao Chen. Trirank: Review-aware explainable recommendation by modeling aspects. In *Proceedings of the 24th ACM international on conference on information and knowledge management*, pages 1661–1670, 2015.
- [6] Yunfeng Hou, Ning Yang, Yi Wu, and Philip S Yu. Explainable recommendation with fusion of aspect information. *World Wide Web*, 22:221–240, 2019.
- [7] Ruoran Huang, Nian Wang, Chuanqi Han, Fang Yu, and Li Cui. Tnam: A tag-aware neural attention model for top-n recommendation. *Neurocomputing*, 385:1–12, 2020.
- [8] Hongmei Li, Xingchun Diao, Jianjun Cao, Lei Zhang, and Qin Feng. Tag-aware recommendation based on bayesian personalized ranking and feature mapping. *Intelligent data analysis*, 23(3):641–659, 2019.
- [9] Lei Li, Yongfeng Zhang, and Li Chen. Personalized transformer for explainable recommendation. *arXiv preprint arXiv:2105.11601*, 2021.
- [10] 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*, pages 345–354, 2017.
- [11] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81, 2004.
- [12] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318, 2002.
- [13] Jianshan Sun, Mingyue Zhu, Yuanchun Jiang, Yezheng Liu, and Le Wu. Hierarchical attention model for personalized tag recommendation. *Journal of the Association for Information Science and Technology*, 72(2):173–189, 2021.
- [14] Nan Wang, Hongning Wang, Yiling Jia, and Yue Yin. Explainable recommendation via multi-task learning in opinionated text data. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, pages 165–174, 2018.
- [15] 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*, pages 83–92, 2014.

- 181 [16] Weibin Zhao, Aoran Zhang, Lin Shang, Yonghong Yu, Li Zhang, Can Wang, Jiajun Chen, and
182 Hongzhi Yin. Hyperbolic personalized tag recommendation. In *International Conference on*
183 *Database Systems for Advanced Applications*, pages 216–231. Springer, 2022.
- 184 [17] Yi Zuo, Shengzong Liu, Yun Zhou, and Huanhua Liu. Tral: A tag-aware recommendation
185 algorithm based on attention learning. *Applied Sciences*, 13(2):814, 2023.
- 186 [18] Yi Zuo, Yun Zhou, Shengzong Liu, and Yupeng Liu. A tag-aware recommendation algorithm
187 based on deep learning and multi-objective optimization. In *2023 International Conference on*
188 *Pattern Recognition, Machine Vision and Intelligent Algorithms (PRMVIA)*, pages 42–46. IEEE,
189 2023.