

Deep Mutual Learning across Task Towers for Effective Multi-Task Recommender Learning

Yi Ren¹, Ying Du¹, Bin Wang¹ and Shenzheng Zhang¹

¹Tencent, Beijing, China

Abstract

Recommender systems usually leverage multi-task learning methods to simultaneously optimize several objectives because of the multi-faceted user behavior data. The typical way of conducting multi-task learning(MTL) is to establish appropriate parameter sharing across multiple tasks at lower layers while reserving a separate task tower for each task at upper layers. With such design, the lower layers intend to explore the structure of task relationships and mine valuable information to be used by the task towers for accurate prediction.

Since the task towers exert direct impact on the prediction results, we argue that the architecture of standalone task towers is sub-optimal for promoting positive knowledge sharing. First, for each task, attending to the input information of other task towers is beneficial. For instance, the information useful for predicting the "like" task is also valuable for the "buy" task. Furthermore, because different tasks are inter-related, the training labels of multiple tasks should obey a joint distribution. It is undesirable for the prediction results for these tasks to fall into the low density areas. Accordingly, we propose the framework of **Deep Mutual Learning** across task towers(**DML**), which is compatible with various backbone multi-task networks. At the entry layer of the task towers, the shared component of Cross Task Feature Mining(**CTFM**) is introduced to transfer input information across the task towers while still ensuring one task's loss will not impact the inputs of other task towers. Moreover, for each task, dedicated network component called Global Knowledge Distillation(**GKD**) are utilized to distill valuable knowledge from the global results of the upper layer task towers to enhance the prediction consistency. Extensive offline experiments and online A/B tests are conducted to evaluate and verify the proposed approach's effectiveness.

Keywords

Recommender Systems, Multi-Task Learning, Parameter Sharing

1. Introduction


Recently, we have seen the widespread application of recommender systems, which involve different types of user feedback signals, such as clicking, rating, commenting, etc. Moreover, no single feedback signal can accurately reflect user satisfaction. For example, over-concentrating on clicking may aggravate the click-bait issue. Therefore, it is highly desirable to be able to effectively learn and estimate multiple types of user behaviors at the same time. And Multi-Task Learning is a promising technique to address this challenge. Given several related learning tasks, the goal of multi-task learning is to enhance the overall performance of different tasks by

ORSUM@ACM RecSys 2023: 6th Workshop on Online Recommender Systems and User Modeling, jointly with the 17th ACM Conference on Recommender Systems, September 19th, 2023, Singapore

✉ yiren_bj@outlook.com (Y. Ren); yingdu@tencent.com (Y. Du); hillmwang@tencent.com (B. Wang); qjzcyzhang@tencent.com (S. Zhang)



© 2023 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

 CEUR Workshop Proceedings (CEUR-WS.org)

leveraging knowledge transfer among tasks. With the Multi-Task Learning(MTL) paradigm[1, 2], multiple tasks are learned simultaneously in a single model. Compared with single-task solutions, MTL costs much fewer machine resources and improves learning efficiency since we just need to train and deploy a single model. Moreover, MTL usually can warrant enhanced recommendation performance through appropriate parameter sharing.

Most existing methods of Multi-Task Learning(MTL) appropriately share parameters across multiple tasks at lower layers while keeping separate task towers at upper layers. These methods can be roughly classified into four categories. The first category comprises the methods of **hard parameter sharing**. Among them, embedding sharing is the most intuitive structure to share information. For instance, the ESSM model [3] shares embedding parameters between the tasks of CTR (Click-Through Rate) and CVR (Conversion Rate) for improving the prediction performance of the sparse CVR task. In addition to the embedding parameters, the Shared-Bottom structure[2] is introduced to share the parameters of lower-layer MLPs among tasks. But these methods are severely plagued by the task conflicts and negative transfer issue. Second, for the methods of **soft parameter sharing**, each task owns separate parameters, which are regularized during training to minimize the differences between the shared parameters. L2-constrained[4] is a typical algorithm belonging to this category. Third, for the methods of **customized routing**, they learn customized routing weights for each task to combine and fuse information from lower-layer networks to counteract the negative transfer issue. Cross-stitch network[5] and sluice network[6] learn separate linear weights for each task to selectively merge representations from different lower-level branches. SNR[7] modularizes the shared low-level layers into parallel sub-networks and uses a transformation matrix multiplied by a scalar coding variable to learn their connections to upper-level layers to alleviate the task conflict and negative transfer issue. MSSM[8] learns differing combinations of feature fields for each expert and designs finer-grained sharing patterns among tasks through a set of coding variables that selectively choose which cells to route for a given task. But the learned routing parameters of these methods are static for all the samples, which can hardly warrant optimal performance. Finally, the methods of **dynamic gating** learn optimized weights for each task based on the input sample to effectively combine the outputs of lower-level networks and achieve success in industrial applications. The MMoE [9] model adapts the Mixture-of-Experts (MoE)[10] structure to multi-task learning by sharing the expert sub-networks across all tasks, while also maintaining separate gating network optimized for each task. And Zhao et al. [11] extend the MMoE model [9] and apply it to learn multiple ranking objectives in Youtube video recommender systems. Moreover, PLE [12] achieves superior performance for news recommendation by assigning both shared expert sub-networks among tasks and dedicated expert sub-networks for each task.

AITM [13] is the most similar method, which also augments the architecture of task towers. Nevertheless, as a concrete implementation, it is not validated to enhance the performance of various multi-task models. Moreover, it can only work for the tasks with sequential dependence relations.

Admittedly, the aforementioned methods achieve impressive performance. However, as the task towers exert a direct effect on the prediction results, the standalone task towers tend not to be the most effective design for promoting positive knowledge transfer by exploiting the task relationships. First, for each task, the information selected by the relevant tasks is

extremely valuable. Accordingly, we introduce the shared component of Cross Task Feature Mining(*CTFM*), which utilizes delicate attention mechanisms to extract relevant information from other tasks at the entry layer of the task tower. With the common attention mechanisms, the explicit task-specific information distilled by lower-level networks are mingled together and one task’s loss will undesirably affect the inputs of other task towers, which is the task awareness missing problem and can hinder the learning of lower-level networks. In contrast to the usual attention mechanisms, our design can ensure appropriate information separation. We argue that reserving explicit task-specific knowledge has a positive effect on performance, which is validated in the experimental section. Second, because the tasks for recommender systems are related, the training labels of multiple tasks should obey a joint distribution. The prediction results for these tasks should not densely fall into the low-density areas. Therefore, a dedicated network named Global Knowledge Distillation(*GKD*) is introduced for each task to distill valuable global knowledge from the results of the upper layer task towers. For each task, the distilled global information helps to ensure consistent predictions with other tasks. We summarize our main contributions below.

- We propose the framework of Deep Mutual Learning across task towers(**DML**), which is compatible with various backbone multi-task models.
- The proposed novel sharing structure helps to enhance effective knowledge transfer across different tasks.
- We conduct offline experiments and online A/B testing to evaluate and understand the effectiveness of our method.

2. Methodology

In this section, we first introduce the problem of multi-objective ranking for recommender systems. Second, we describe the general design of DML. Finally, we elaborate on the introduced components.

2.1. Multi-Objective Ranking for Recsys

Given a set of candidates with N items $\mathcal{C} = \{i_n\}_{1 \leq n \leq N}$, the ranking model for recommender systems is to rank and recommend the top M items $\mathcal{S} = \{i_m\}_{1 \leq m \leq M} \subseteq \mathcal{C}$ for user u so as to optimize the overall utility and enhance user experience. First, for each pair of user u and item i_n , the input feature x_n is derived. Second, a multi-task learning model is utilized to estimate K objectives corresponding to multiple user feedback signals. Furthermore, to compute the overall reward, we need to merge the multiple predictions with a function Φ shown in equation (2) to derive the item’s final reward score for greedy ranking.

$$o_n^1, o_n^2, \dots, o_n^K = MTLModel(x_n; \theta) \quad (1)$$

$$r_n = \Phi(o_n^1, o_n^2, \dots, o_n^K) \quad (2)$$

where θ denotes the model parameters and Φ is usually a function manually tuned to reflect the reward of recommending item i_n to user u based on the business goals. With the estimated reward scores $\{r_n\}_{1 \leq n \leq N}$ for the items in \mathcal{C} , the ranking model can recommend the item sequence \mathcal{S} consisting of the top M items to the user u .

2.2. Overall Design of DML

With the existing MTL algorithms[1], the equation (1) can be further decomposed as below. For simplicity, we omit the subscript n in this section.

$$l^1, l^2, \dots, l^K = G(x; \theta_l) \quad (3)$$

$$o^k = F^k(l^k; \theta_h^k) \quad (4)$$

where G represents the lower level networks that encodes x to K different latent spaces with partially or fully shared parameters θ_l . And F^k is the upper-level network for task k , which accepts l^k as input to model objective k with task-specific higher level parameter θ_h^k . Multiple candidate models [3, 9, 12, 7, 11, 5, 8, 4, 6, 1] are proposed to enhance G with different parameter sharing designs.

In this research, rather than G , we focus on the enhancement of upper-level networks for improved prediction performance. First, the shared component of *CTFM* is introduced, which leverages the attention mechanism to extract relevant information from the inputs of other task towers (the results of Equation (3)) as a complement to the target task. Please note that this attention is well-designed to solve the task-awareness missing issue, for which the excessive encouragement of knowledge sharing is not conducive to the extraction of task-specific knowledge. With our design, the gradients computed from the target task's loss will not impact the inputs of other task towers.

$$\hat{l}^1, \hat{l}^2, \dots, \hat{l}^K = CTFM(l^1, l^2, \dots, l^K; \theta_s) \quad (5)$$

where shared parameters θ_s is employed across different tasks.

Moreover, a separate multi-layer network is introduced for each task to process each element of $\{\hat{l}^k\}_{1 \leq k \leq K}$ and generate the hidden representation, based on which accurate prediction can be made.

$$h^k = H^k(\hat{l}^k; \theta_{t_0}^k) \quad (6)$$

where H^k denotes the task-specific MLP for task k with separate parameters $\theta_{t_0}^k$.

Finally, for each task, a dedicated component named *GKD* is utilized to distill information from the hidden representations for both itself and other tasks to promote prediction consistency across tasks and more precisely model the corresponding objective.

$$o^k = GKD^k(h^k, \{h^j\}_{1 \leq j \leq K}; \theta_{t_1}^k) \quad (7)$$

where GKD^k is the dedicated component for task k . In contrast to *CTFM* at Equation (5), we utilize task-specific parameters $\theta_{t_1}^k$ here as specialization is usually helpful for upper layer networks. Furthermore, proper operation is implemented to ensure the prediction error of one task does not impact the hidden representations of other tasks.

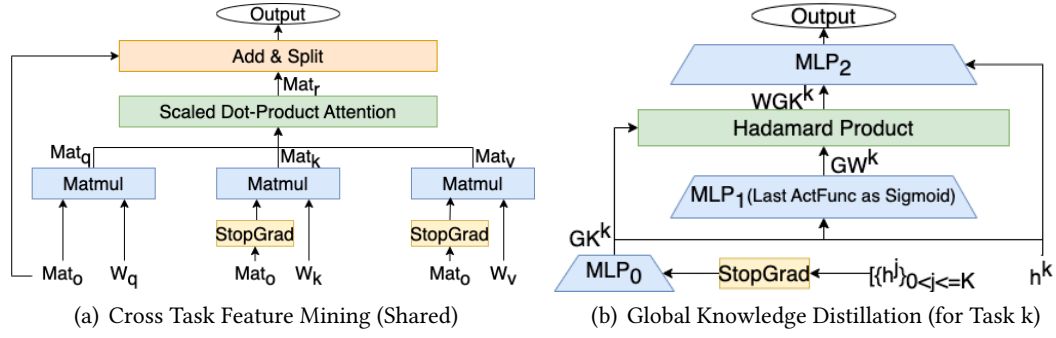


Figure 1: Introduced Components of DML

2.3. Cross Task Feature Mining

Please refer to Figure 1(a) for the detailed process of *CTFM*. From the perspective of task towers, the outputs from lower-level networks at Equation (3) can be recognized as the mined features for them. First, for each task, the features mined for related tasks can be leveraged to enhance its prediction performance. Thus, trainable task embeddings, namely $\{t^k\}_{1 \leq k \leq K}$, are derived for the tasks to facilitate the learning of general task relations. Second, the importance of features from related tasks can vary per sample. Accordingly, we stack together the items of the set $\{l^k + t^k\}_{1 \leq k \leq K}$ to derive the matrix $Mat_o \in \mathcal{R}^{K \times d_0}$ where d_0 is the size of l^k and t^k . Third, we use the projection parameters $W_q, W_k, W_v \in \mathcal{R}^{d_0 \times d_0}$ to transform Mat_o to the query, key, and value matrix of $Mat_q, Mat_k, Mat_v \in \mathcal{R}^{K \times d_0}$. And gradient backpropagation from Mat_k and Mat_v to Mat_o is forbidden. Moreover, the scaled dot-product attention [14] is performed on $Mat_q, Mat_k,$ and Mat_v to compute the result matrix $Mat_r \in \mathcal{R}^{K \times d_0}$. Finally, we add Mat_o to Mat_r for residual connection and further split based on the first axis to return $\{\hat{l}^k\}_{1 \leq k \leq K}$. Please note that the aforementioned networks are shared among the tasks to encourage generalizable modeling with parameter sharing. The usual attention mechanisms will cause the task awareness missing problem and can hinder the learning of lower-level networks. Instead, our design can ensure appropriate information separation and reserve the explicit task-specific knowledge by stopping the gradient backflow from Mat_k and Mat_v to Mat_o .

2.4. Global Knowledge Distillation

Please refer to Figure 1(b) for the detailed process of *GKD*. In contrast to *CTFM*, each task is assigned dedicated parameters for global knowledge distillation. Acting as the last step, we would like to facilitate more flexible modeling by promoting task specialization here. This module accepts the hidden representations for both the corresponding task and other tasks as input. First, a multilayer perceptron(MLP) is utilized to extract valuable global knowledge ($GK^k \in \mathcal{R}^{d_1}$) for the target task k from the concatenation of all these hidden representations ($\{h^j\}_{1 \leq j \leq K}$). Since the design goal here is to train task-specific MLPs to distill relevant global knowledge while not impacting the hidden representations, the gradient backpropagation from

Table 1

The overall performance. The bold-face font denotes the winner in that column. Moreover, the "*" symbol denotes introducing DML achieves significant ($p < 0.05$ for one-tailed t-test) gain over the corresponding baseline.

Model	ML-1M			Electronics	
	<i>AUC</i>	<i>MSE</i>	<i>Consistent Ratio</i>	<i>AUC_{rate}</i>	<i>AUC_{pos}</i>
Single Task	0.8066	0.7741	0.7154	0.7608	0.7334
SB	0.8100	0.7724	0.7530	0.7876	0.7608
SB+DML	0.8115*	0.7648*	0.7649*	0.7890*	0.7631*
MSSM	0.8128	0.7651	0.7519	0.7883	0.7627
MSSM+DML	0.8141*	0.7611*	0.7637*	0.7892*	0.7641*
MMOE	0.8105	0.7688	0.7507	0.7888	0.7628
MMOE+DML	0.8135*	0.7591*	0.7596*	0.7897*	0.7644*
PLE	0.8122	0.7606	0.7514	0.7885	0.7627
PLE+DML	0.8151*	0.7533*	0.7631*	0.7893*	0.7640*

GK^k to the MLP's input is prohibited. Second, we input GK^k and h^k to another MLP with Sigmoid as its last activation function. Then, the weights of GK^k 's different latent dimensions can be dynamically adapted for each sample. The weights are denoted by $GW^k \in \mathcal{R}^{d_1}$. Moreover, the weighted global knowledge ($WKG^k \in \mathcal{R}^{d_1}$) is computed with the hadamard product of GK^k and GW^k . Finally, WKG^k and h^k are concatenated together as the input for the last MLP to make predictions for task k .

3. Experiments

In this section, we conduct extensive offline experiments¹ and online A/B testing to prove DML's effectiveness.

3.1. Experimental Settings for Public Data

3.1.1. Datasets

We evaluate our methods on two public datasets.

- **MovieLens-1M**[15]: One of the currently released MovieLens datasets, which contains 1 million movie ratings from 6,040 users on 3,416 movies.
- **Amazon**[16]: A series of datasets consisting of product reviews from Amazon.com. We use the sub-category of "Electronics" including 1.7 million reviews from 192,403 users on 63,001 items.

¹The code can be found at: <https://github.com/renyi533/ml-consistency/tree/main>.

For ML-1M, we introduce the binary classification task of positive rating prediction (≥ 4) and the regression task of rating estimation. These two tasks are strictly correlated. For electronics, following [17], we first augment the dataset by randomly sampling un-rated items for every user. Moreover, we make sure the number of the un-rated items is the same as the number of the rated items for each user. Furthermore, we introduce two binary classification tasks, namely rating prediction (whether a rating exists) and positive rating prediction. Compared with the tasks of ML-1M, the negative transfer is more likely to occur as the task relationship here is more complex (The Pearson correlation coefficient [18] of these two labels is around 0.7). Both data are randomly split into the training set, validation set, and test set by the ratio of 8:1:1.

3.1.2. Evaluation Metrics

The merge function Φ in Equation (2) assumes that the model can estimate accurate interaction probabilities for binary classification tasks (e.g. clicking) and absolute values for regression tasks (e.g. watch time). Therefore, instead of the ranking metrics, such as NDCG [19] and MRR [20], we use the metrics of AUC [21] for classification tasks and Mean Squared Error (MSE) [22] for regression tasks. Please note that many other recommendation literature, such as [12, 11, 8], also use similar metrics. For AUC, a bigger value indicates better performance. While, for MSE, it is the smaller the better.

3.1.3. Models

As *soft parameter sharing* methods need resources to store and train multiple sets of parameters, they are not widely applied in recommender systems. Thus, we select base models to cover the other three categories. The models include Shared-Bottom(SB)[2], MSSM[8], MMOE[9], and PLE[12]. MSSM is a recent method belonging to the *Customized Routing* category and achieves better results than SNR[7] and Cross-Stitch[5]. Though with the same category of *Dynamic Gating*, both MMOE and PLE are tested owing to their popularity. For each base model, we will verify whether DML can achieve additional gains. For reference, we also provide the performance of single task models.

3.1.4. Implementation Details

For each feature, we use the embedding size 8. As suggested by the original papers, we use 1 level bottom sub-networks for MMOE, MSSM, and SB while 2 levels for PLE. For SB, a sub-network of 1 layer structure with 128 output dimensions is shared by the tasks. For other multi-task models, each bottom level includes three sub-networks, which have the same aforementioned architecture. For MSSM and PLE, task-specific and shared sub-networks are designated. For multi-task models, each task tower is of the three layers MLP structure (128,80,1) and each task is assigned equal loss weight. For the single-task model, each task utilizes the four layers MLP structure (128,128,80,1). For the first two MLPs at Figure 1(b), we utilize the one layer structure with 80 as the output dimension. For the last MLP at Figure 1(b), a one layer structure with 1 as the output size is used. If not explicitly specified, RELU [23] is used as the default activation function. All models are implemented with tensorflow [24] and optimized using the Adam [25]

Table 2
Further Analysis Results

Model	ML-1M		Electronics	
	<i>AUC</i>	<i>MSE</i>	<i>AUC_{rate}</i>	<i>AUC_{pos}</i>
<i>MSSM</i>	0.8128	0.7651	0.7883	0.7627
<i>MSSM + CTFM</i>	0.8134	0.7654	0.7891	0.7636
<i>MSSM + GKD</i>	0.8129	0.7636	0.7886	0.7633
<i>MSSM + DML_{v0}</i>	0.8138	0.7625	0.7884	0.7631
<i>MSSM + DML</i>	0.8141	0.7611	0.7892	0.7641
<i>PLE</i>	0.8122	0.7606	0.7885	0.7627
<i>PLE + CTFM</i>	0.8133	0.7603	0.7890	0.7633
<i>PLE + GKD</i>	0.8137	0.7574	0.7887	0.7634
<i>PLE + DML_{v0}</i>	0.8141	0.7536	0.7886	0.7629
<i>PLE + DML</i>	0.8151	0.7533	0.7893	0.7640

optimizer with learning rate 0.001 and mini-batch size 512. We run 20 times for each test to report the results.

3.2. Overall Performance for Public Data

Please refer to Table 1 for the overall results. First, DML achieves significant gains across all these tested multi-task models on the two public datasets, which shows DML’s wide compatibility. Second, DML-enhanced PLE and MMOE get the best performance for MovieLens and Electronics respectively. Considering their wide application in recommender systems, the results are as expected. Third, the multi-task models perform better than the single task models thanks to the knowledge transfer between tasks.

Besides AUC and MSE, DML should help to foster task consistency with *CTFM* and *GKD*. As the tasks of MovieLens are rigorously correlated, we verify whether DML really enhances task consistency on this data. First, we construct pairs of samples with different rating scores and count the pair numbers. Second, we count the number of pairs, for which the prediction scores of both tasks are in the same pair order as the rating score. The enhancement of the pair order consistency among the two prediction scores and rating score should positively contribute to the performance. Then, we can compute the metric of ‘Consistency Ratio’. The listed data in Table 1 agree with our anticipation. (For Shared-Bottom, we also observe more pairs, for which predictions of both tasks are in rating score’s reverse order. This can explain its worse performance in spite of the better consistency ratio.)

3.3. Further Analysis on Public Data

We select the two latest algorithms of PLE[12] and MSSM[8] to appraise the value of *DML*’s components, namely *CTFM* and *GKD*. Without the stop gradient operation, *CTFM* will be very similar to the common attention mechanism. To prove the benefit of *CTFM*’s design, we also add the assessment for *DML_{v0}*, which reserve the design of *GKD* while remove the

gradient blocking operation of *CTFM*. Please refer to Table 2 for the evaluation results. First, *CTFM* and *GKD* both contribute considerable gains over the base model. Second, as the integrated model, *DML* enhances the performance further. Third, *DML_{v0}* is consistently worse than *DML*, which corroborates the value of reserving task-awareness. Compared with *CTFM* and *GKD*, *DML_{v0}* performs better on MovieLens while much worse on Electronics. The task relationship of Electronics is more complex and negative transfer across tasks usually exhibits more severe impact due to task conflicts. In this case, compared with vanilla attention, *CTFM* obtains substantial gains.

3.4. Online A/B Testing

DML is applied to the ranking stage[26] of an industrial large-scale news recommender system. PLE [12] is utilized as the base model. And the main prediction tasks are the binary classification task of Click Through Rate (CTR) and the regression task of item watch time. First, after the model converge by training with billions of samples, the AUC metric for CTR consistently increases 0.12% and the MSE metric for watch time decreases 0.14%. Moreover, the most important online metrics include effective PV(count of Page Views with watch time exceeding a threshold) and total watch time. We randomly distributed online users to two buckets with the base PLE model or PLE+DML model and evaluated the performance for two weeks. DML achieves significant ($p < 0.05$) gains over the base model by 1.22% for effective PV and 0.61% for total watch time. DML has been deployed to our online environment based on the results.

4. Conclusion

In this paper, we propose the framework of **Deep Mutual Learning** across task towers(**DML**), which is compatible with various backbone multi-task networks. Extensive offline experiments help to verify DML’s effectiveness on multiple real-world datasets and across various base models. Moreover, thorough ablation studies are carried out to verify and understand the value of each newly introduced module. Finally, DML achieves significant online gains and has already been deployed to the online platform.

References

- [1] Y. Zhang, Q. Yang, A survey on multi-task learning, *IEEE Transactions on Knowledge and Data Engineering* (2021).
- [2] R. Caruana, Multitask learning, *Machine learning* 28 (1997) 41–75.
- [3] X. Ma, L. Zhao, G. Huang, Z. Wang, Z. Hu, X. Zhu, K. Gai, Entire space multi-task model: An effective approach for estimating post-click conversion rate, in: *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, 2018, pp. 1137–1140.
- [4] L. Duong, T. Cohn, S. Bird, P. Cook, Low resource dependency parsing: Cross-lingual parameter sharing in a neural network parser, in: *Proceedings of the 53rd annual meeting*

- of the Association for Computational Linguistics and the 7th international joint conference on natural language processing (volume 2: short papers), 2015, pp. 845–850.
- [5] I. Misra, A. Shrivastava, A. Gupta, M. Hebert, Cross-stitch networks for multi-task learning, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 3994–4003.
 - [6] S. Ruder, J. Bingel, I. Augenstein, A. Søgaard, Latent multi-task architecture learning, in: Proceedings of the AAAI Conference on Artificial Intelligence, volume 33, 2019, pp. 4822–4829.
 - [7] J. Ma, Z. Zhao, J. Chen, A. Li, L. Hong, E. H. Chi, Snr: Sub-network routing for flexible parameter sharing in multi-task learning, in: Proceedings of the AAAI Conference on Artificial Intelligence, volume 33, 2019, pp. 216–223.
 - [8] K. Ding, X. Dong, Y. He, L. Cheng, C. Fu, Z. Huan, H. Li, T. Yan, L. Zhang, X. Zhang, et al., Mssm: a multiple-level sparse sharing model for efficient multi-task learning, in: Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, 2021, pp. 2237–2241.
 - [9] J. Ma, Z. Zhao, X. Yi, J. Chen, L. Hong, E. H. Chi, Modeling task relationships in multi-task learning with multi-gate mixture-of-experts, in: Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, 2018, pp. 1930–1939.
 - [10] R. A. Jacobs, M. I. Jordan, S. J. Nowlan, G. E. Hinton, Adaptive mixtures of local experts, *Neural computation* 3 (1991) 79–87.
 - [11] Z. Zhao, L. Hong, L. Wei, J. Chen, A. Nath, S. Andrews, A. Kumthekar, M. Sathiamoorthy, X. Yi, E. Chi, Recommending what video to watch next: a multitask ranking system, in: Proceedings of the 13th ACM Conference on Recommender Systems, 2019, pp. 43–51.
 - [12] H. Tang, J. Liu, M. Zhao, X. Gong, Progressive layered extraction (ple): A novel multi-task learning (mtl) model for personalized recommendations, in: Fourteenth ACM Conference on Recommender Systems, 2020, pp. 269–278.
 - [13] D. Xi, Z. Chen, P. Yan, Y. Zhang, Y. Zhu, F. Zhuang, Y. Chen, Modeling the sequential dependence among audience multi-step conversions with multi-task learning in targeted display advertising, in: Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining, 2021, pp. 3745–3755.
 - [14] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, I. Polosukhin, Attention is all you need, *Advances in neural information processing systems* 30 (2017).
 - [15] F. M. Harper, J. A. Konstan, The movielens datasets: History and context, *Acm transactions on interactive intelligent systems (tiis)* 5 (2015) 1–19.
 - [16] J. McAuley, C. Targett, Q. Shi, A. Van Den Hengel, Image-based recommendations on styles and substitutes, in: Proceedings of the 38th international ACM SIGIR conference on research and development in information retrieval, 2015, pp. 43–52.
 - [17] Y. Wang, Z. Zhao, B. Dai, C. Fifty, D. Lin, L. Hong, L. Wei, E. H. Chi, Can small heads help? understanding and improving multi-task generalization, in: Proceedings of the ACM Web Conference 2022, 2022, pp. 3009–3019.
 - [18] Wikipedia, Pearson correlation coefficient — Wikipedia, the free encyclopedia, <http://en.wikipedia.org/w/index.php?title=Pearson%20correlation%20coefficient&oldid=1146097966>, 2023. [Online; accessed 15-April-2023].

- [19] Y. Wang, L. Wang, Y. Li, D. He, W. Chen, T.-Y. Liu, A theoretical analysis of ndcg ranking measures, in: Proceedings of the 26th annual conference on learning theory (COLT 2013), volume 8, 2013, p. 6.
- [20] Wikipedia contributors, Evaluation measures (information retrieval) — Wikipedia, the free encyclopedia, [https://en.wikipedia.org/w/index.php?title=Evaluation_measures_\(information_retrieval\)&oldid=1095286224](https://en.wikipedia.org/w/index.php?title=Evaluation_measures_(information_retrieval)&oldid=1095286224), 2022. [Online; accessed 9-January-2023].
- [21] P. A. Flach, J. Hernández-Orallo, C. F. Ramirez, A coherent interpretation of auc as a measure of aggregated classification performance, in: ICML, 2011.
- [22] Wikipedia contributors, Mean squared error — Wikipedia, the free encyclopedia, https://en.wikipedia.org/w/index.php?title=Mean_squared_error&oldid=1127519968, 2022. [Online; accessed 31-January-2023].
- [23] R. Arora, A. Basu, P. Mianjy, A. Mukherjee, Understanding deep neural networks with rectified linear units, in: 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings, OpenReview.net, 2018. URL: https://openreview.net/forum?id=B1J_rgWRW.
- [24] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, et al., Tensorflow: a system for large-scale machine learning., in: Osd, volume 16, Savannah, GA, USA, 2016, pp. 265–283.
- [25] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, in: Y. Bengio, Y. LeCun (Eds.), 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, 2015. URL: <http://arxiv.org/abs/1412.6980>.
- [26] P. Covington, J. Adams, E. Sargin, Deep neural networks for youtube recommendations, in: Proceedings of the 10th ACM conference on recommender systems, 2016, pp. 191–198.