

Efficient and Joint Hyperparameter and Architecture Search for Collaborative Filtering

Yan Wen

Department of Electronic Engineering
Beijing National Research Center for
Information Science and Technology
Tsinghua University
Beijing, China
wenyan0531@gmail.com

Chen Gao*

Department of Electronic Engineering
Beijing National Research Center for
Information Science and Technology
Tsinghua University
Beijing, China
chgao96@gmail.com

Lingling Yi

Tencent Inc.
Shenzhen, China
chrisyi@tencent.com

Liwei Qiu

Tencent Inc.
Shenzhen, China
drolcaqiu@tencent.com

Yaqing Wang

Baidu Inc.
Beijing, China
wangyaqing01@baidu.com

Yong Li

Department of Electronic Engineering
Beijing National Research Center for
Information Science and Technology
Tsinghua University
Beijing, China
liyong07@tsinghua.edu.cn

ABSTRACT

Automated Machine Learning (AutoML) techniques have recently been introduced to design Collaborative Filtering (CF) models in a data-specific manner. However, existing works either search architectures or hyperparameters while ignoring the fact they are intrinsically related and should be considered together. This motivates us to consider a joint hyperparameter and architecture search method to design CF models. However, this is not easy because of the large search space and high evaluation cost. To solve these challenges, we reduce the space by screening out usefulness hyperparameter choices through a comprehensive understanding of individual hyperparameters. Next, we propose a two-stage search algorithm to find proper configurations from the reduced space. In the first stage, we leverage knowledge from subsampled datasets to reduce evaluation costs; in the second stage, we efficiently fine-tune top candidate models on the whole dataset. Extensive experiments on real-world datasets show better performance can be achieved compared with both hand-designed and previous searched models. Besides, ablation and case studies demonstrate the effectiveness of our search framework.

CCS CONCEPTS

• **Information systems** → **Recommender systems.**

*Corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
KDD '23, August 6–10, 2023, Long Beach, CA, USA.

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 979-8-4007-0103-0/23/08...\$15.00
<https://doi.org/10.1145/3580305.3599322>

KEYWORDS

Recommendation System; Collaborative Filtering; Automated Machine Learning

ACM Reference Format:

Yan Wen, Chen Gao, Lingling Yi, Liwei Qiu, Yaqing Wang, and Yong Li. 2023. Efficient and Joint Hyperparameter and Architecture Search for Collaborative Filtering. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '23)*, August 6–10, 2023, Long Beach, CA, USA. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3580305.3599322>

1 INTRODUCTION

Collaborative Filtering (CF) is the most widely used approach for Recommender Systems [15, 19, 20, 30], aiming at calculating the similarity of users and items to recommend new items to potential users. They mainly use Neural Networks to build models for users and items, simulating the interaction procedure and predict the preferences of users for items. Recent works also built CF models based on Graph Neural Networks (GNNs) [12, 14, 36, 42]. While CF models may have different performance on different scenes [7], recent works [4, 49] have begun to apply Automated Machine Learning (AutoML) to search data-specific CF models. Previous works, including SIF [39], AutoCF [11] and [37], applied Neural Architecture Search (NAS) on CF tasks. They split the architectures of CF models into several parts and searched each part on architecture space.

However, most of these methods focus on NAS in architecture space, only considering hyperparameters as fixed settings and therefore omitting the dependencies among them. A CF model can be decided by a given architecture and a configuration of hyperparameters. Especially in the task of searching best CF models, hyperparameter choice can affect search efficiency and the evaluation performance an architecture can receive on a given dataset. Recent methods mainly focus on model search, ignoring the important role of hyperparameters. For instance, SIF focuses on interaction function, while it uses grid search on hyperparameters space. AutoCF

does not use hyperparameter tuning on each architecture, which may make the searched model sub-optimal since proper hyperparameters vary for different architectures. [37] includes GNN models in the architecture space, but the search and evaluation cost for architectures may be high.

We find that these works only focus on either hyperparameters or fixed parts in CF architecture, neglecting the relation between architecture and hyperparameters. If architecture cannot be evaluated with proper hyperparameters, a sub-optimal model may be searched, possibly causing a reduction in performance. To summarize, there exists a strong dependency between hyperparameters and architectures. That is, the hyperparameters of a model are based on the design of architectures, and the choices of hyperparameters also affect the best performance an architecture may approach. We suppose that hyperparameters can be adaptively changed when the architecture change in CF tasks, so we consider that the CF search problem can be defined on a joint space of hyperparameters and architectures. Therefore, the CF problem can be modeled as a joint search problem on architecture and hyperparameter space. While hyperparameters and architectures both influences the cost and efficiency of CF search task, there exist challenges for joint search problems: (1) Since the joint search space is designed to include both hyperparameters and architectures, the joint search problem has a large search space, which may make it more difficult to find the proper configuration of hyperparameters and architectures; (2) In a joint search problem, since getting better performance on a given architecture requires determining its hyperparameters, the evaluation cost may be more expensive.

We propose a general framework, which can optimize CF architectures and their hyperparameters at the same time in search procedure. The framework of our method is shown in Figure 1b, consisting of two stages. Prior to searching hyperparameters and architectures, we have a full understanding of the search space and reduce hyperparameter space to improve efficiency. Specifically, we reduced the hyperparameter space by their performance ranking on CF tasks from different datasets. We also propose a frequency-based sampling strategy on the user-item matrix for fast evaluation. In first stage, we search and evaluate models on reduced space and subsampled datasets, and we jointly search architecture and hyperparameters with a surrogate model. In second stage, we propose a knowledge transfer-based evaluation strategy to leverage surrogate model to larger datasets. Then we evaluate the model with transferred knowledge and jointly search the hyperparameters and architectures to find the best choice in original dataset.

Overall, we make the following important contributions:

- We propose an approach that can jointly search CF architectures and hyperparameters to get performance from different datasets.
- We propose a two-stage search algorithm to efficiently optimize the problem. The algorithm is based on a full understanding of search space and transfer ability between datasets. It can jointly update CF architectures and hyperparameters and transfer knowledge from small datasets to large datasets.
- Extensive experiments on real-world datasets demonstrate that our proposed approach can efficiently search configurations in designed space. Furthermore, results of ablation and case study show the superiority of our method.

2 RELATED WORK

2.1 Automated Machine Learning (AutoML)

Automated Machine Learning (AutoML) [18, 40] refers to a type of method that can learn models adaptively to various tasks. Recently, AutoML has achieved great success in designing the state-of-the-art model for various applications such as image classification and segmentation [25, 34, 51], natural language modeling [33], and knowledge graph embedding [45].

AutoML can be used in mainly two fields: *Neural Architecture Search (NAS)* and *Hyperparameter Optimization (HPO)*.

- NAS [23, 28, 51] splits architectures into several components and searches for each part of the architecture to achieve the whole part. DARTS [26] uses gradient descent on continuous relaxation of the architecture representation, while NASP [41] improves DARTS by including proximal gradient descent on architectures.
- HPO [1, 10, 22], usually tuning the hyperparameters of a given architecture, always plays an important role in finding the best hyperparameters for the task. Random Search is a frequently used method in HPO for finding proper hyperparameters. Algorithms for HP search based on model have been developed for acceleration [6], including Bayesian Optimization (BO) methods like Hyperopt [2], BOHB [17] and BORE [35], etc.

Recent studies on AutoML have shown that incorporating HPO into the NAS process can lead to better performance and a more effective exploration of the search space. For example, a study on ResNet [43] showed that considering the NAS process as HPO can lead to improved results. Other works, such as AutoHAS [8], have explored the idea of considering hyperparameters as a choice in the architecture space. FEATHERS [31] has focused on the joint search problem in Federated Learning. ST-NAS [3] uses weight-sharing NAS on architecture space and consider HP as part of architecture encoding. These studies demonstrate the potential of joint search on hyperparameters and architectures in improving the performance and efficiency of machine learning models.

2.2 Collaborative Filtering (CF)

2.2.1 Classical CF Models. Collaborative Filtering (CF) [30] is the most fundamental solution for Recommender Systems (RecSys). CF models are usually designed to learn user preferences based on the history of user-item interaction. Matrix Factorization (MF) [20] generates IDs of users and items, using a high-dimensional vector to represent the location of users and items' features. The inner product is used as interaction function, calculating the similarity of user/item vectors. The MF-based method has been demonstrated effective in SVD++ [20] and FISM [19]. NCF [15] applied neural networks to building CF models, using a fused model with MF and multi-layer perceptron (MLP) as an interaction function, taking user/item embeddings as input, and inferring preference scores. JNCF [5] extended NCF by using user/item history to replace user/item ID as the input encoding.

Recently, the user-item interaction matrix can also be considered as a bipartite graph, thus Graph Neural Networks (GNNs) [12] are also applied to solve CF tasks for their ability to capture the high-order relationship between users and items [36]. They consider both users and items as nodes and the interaction of users and

items as edges in bipartite graph. For example, PinSage[42] uses sampling on a graph according to node (user/item) degrees, and learn the parameters of GNN on these sampled graphs. NGCF [36] uses a message-passing function on both users themselves and their neighbor items and collects both the information to build proper user and item embeddings. LightGCN [14] generates embeddings for users and items with simple SGC layers.

2.2.2 AutoML for CF. Recently, AutoML has been frequently used in CF tasks, aiming at finding proper hyperparameters and architectures for different tasks [18, 40]. Hyperparameter Optimization (HPO) has been applied on the embedding dimension of RecSys models. For example, AutoDim [47] searches embedding dimension in different fields, aiming at assigning embedding dimension for duplicated content. PEP [27] used learnable thresholds to identify the importance of parameters in the embedding matrix, and trains the embedding matrix and thresholds by sub-gradient decent. AutoFIS [24] is designed to learn feature interactions by adding an attention gate to every potential feature interaction. Most of these works tune the embedding dimension adaptively on Recommender Systems tasks, mostly in Click-Through-Rate (CTR) tasks.

Neural Architecture Search (NAS) has also been applied on CF tasks, including SIF [39] AutoCF [11] and [37]. In detail, SIF adopts the one-shot architecture search for adaptive interaction function in the CF model. AutoCF designs an architecture space of neural network CF models, and the search space is divided into four parts: encoding function, embedding function, interaction function and prediction function. AutoCF selects an architecture and its hyperparameters in the space with a performance predictor. Hyperparameters are considered as a discrete search component in search space, neglecting its continuous characteristic. [37] designs the search space on graph-based models, which also uses random search on the reduced search space.

3 SEARCH PROBLEM

As mentioned in the introduction, finding a proper CF model should be considered as a joint search problem on hyperparameters and architectures. We propose to use AutoML to find architectures and their proper hyperparameters efficiently. The joint search problem on CF hyperparameters and architectures can be modeled as a bilevel optimization problem as follows:

DEFINITION 1 (JOINT AUTOMATED HYPERPARAMETERS AND ARCHITECTURE SEARCH FOR CF). *Let f^* denote the proper CF model, and then the joint search problem for CF can be formulated as:*

$$\alpha^*, h^* = \max_{\alpha \in \mathcal{A}, h \in \mathcal{H}} \mathcal{M}(f(\mathbf{P}^*; \alpha, h), S_{val}), \quad (1)$$

$$\text{s.t. } \mathbf{P}^* = \arg \max_{\mathbf{P}} \mathcal{M}(f(\mathbf{P}; \alpha, h), S_{tra}). \quad (2)$$

where \mathcal{H} contains all possible choices of hyperparameters h , where \mathcal{A} contains all possible choices of architectures α , S_{val} and S_{tra} denote the training and validation datasets, \mathbf{P} denotes the learnable parameters of the CF architecture α , and \mathcal{M} denotes the performance measurement, such as Recall and NDCG.

We encounter the following key challenges in effectively and efficiently solving the search problem: Firstly, the joint search space must contains a wide range of architectural operations within \mathcal{A} ,

Table 1: The operations we use for architecture space.

Architecture	Operations	
Input Features	User	ID, H
	Item	ID, H
Features Embedding	NN-based	Mat, MLP
	Graph-based	GraphSAGE, HadamardGCN, SGC
Interaction Function	multiply, minus, min, max, concat	
Prediction Function	SUM, VEC, MLP	

as well as frequently utilized hyperparameters in the learning stage within \mathcal{H} . Since this space contains various types of components, including continuous hyperparameters and categorical architectures, it is essential to appropriately encode the joint search space for effective exploration. Secondly, considering the dependency between hyperparameters and architectures, the search strategy should be robust and efficient, meeting the accuracy and efficiency requirements of real-world applications. Compared to previous AutoML works on CF tasks, our method is the first to consider a joint search on hyperparameters and architectures on CF tasks.

3.1 Architecture Space: \mathcal{A}

In this paper, the general architecture of CF models can be separated into four parts [11]: Input Features, Feature Embedding, Interaction Function, and Prediction Function. Based on the frequently used operations, we build the architecture space \mathcal{A} , illustrated in Table 1.

Input Features The input features of CF architectures come from original data: user-item rating matrix. We can apply the interactions between users and items and map them to high-dimension vectors. There are two manners for encoding users and items as input features: one-hot encoding (ID) and multi-hot encoding (H). As for one-hot encoding (ID), we consider the reorganized id of both users and items as input. Since the number of users and items is different, we should maintain two matrices when we generate these features. As for multi-hot encoding (H), we can consider the interaction of users and items. A user vector can be encoded in several places by its historical interactions with different items. The items can be encoded in the same way.

Feature Embedding The function of feature embedding maps the input encoding with high dimension into vectors with lower dimension. According to designs in Section 2.2, we can elaborate the embedding manners in two categories: Neural Networks based (NN-based) embeddings and Graph Neural Networks based (Graph-based) embeddings. The embedding function is related to the size of input features. As for NN-based methods, a frequently used method of calculation is Mat, mainly consists of a *lookup-table* in ID level, and mean pooling on both users/items sides. We can also use a multi-layer perceptron (MLP) for each user and item side, helping convert multi-hot interactions into low-dimensional vectors. As for Graph-based methods, the recent advances in GNNs use more

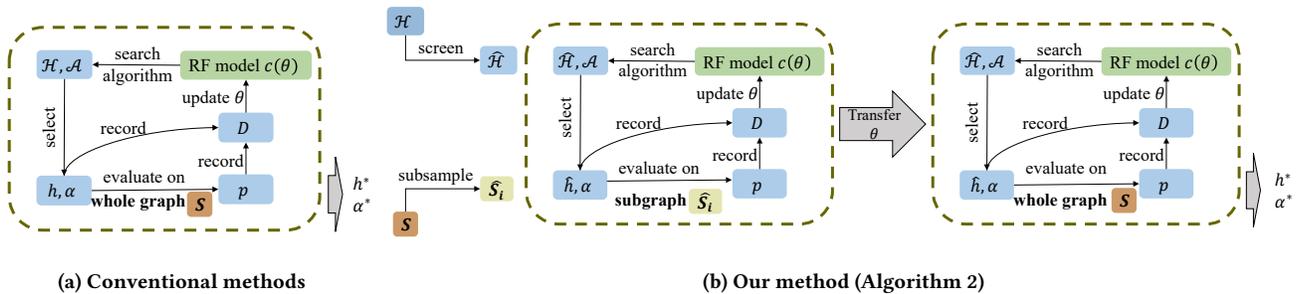


Figure 1: Joint search strategy of conventional methods and our method.

Table 2: Origin hyperparameter space, discretized values, and shrunk range after using screening method in Section 4.1.

Hyperparameter	Original range	Discrete values	Shrunk range
optimizer	Adagrad, Adam, SGD	Adagrad, Adam, SGD	Adagrad, Adam
learning rate	[1e-6, 1e-1]	$\{10^{-6}, 10^{-5}, \dots, 10^0\}$	[1e-5, 1e-2]
embedding dimension	[1, 512]	$\{2^0, 2^1, \dots, 2^9\}$	[2, 64]
weight decay	[1e-5, 1e-1]	$\{10^{-5}, 10^{-4}, \dots, 10^{-1}\}$	1e-1
batch size	[500, 5000]	$\{500, 1000, \dots, 5000\}$	{2000}

complex graph neural networks to aggregate the neighborhoods, such as GraphSAGE [13], HadamardGNN [36], and SGC [14] etc.

Interaction Function The interaction function calculates the relevance between a given user and an item. In this operation, the output is a vector affected by the embeddings of the user and item. The inner product is frequently used in many research works on CF tasks. We split the inner product and consider an element-wise product as an essential operation, which is noted as `multiply`. In coding level, we can also use `minus`, `max`, `min` and `concat`. They help us join users and items with different element-wise calculations.

Prediction Function This operation stage helps turn the output of the interaction function into an inference of similarity. As for the output vector of a specific interaction, a simple way is to use summation on the output vector. Thus, `multiply+SUM` can be considered the inner product in this way. Besides, we use a weight vector with learnable parameters, noted as `VEC`. Multi-layer perceptron (MLP) can also be used for more complex prediction on similarity.

3.2 Hyperparameter Space: \mathcal{H}

Besides the model architecture, the hyperparameter (HP) setting also plays an essential role in determining the performance of the CF model. The used components for hyperparameter space are illustrated in the first column of Table 2.

The CF model, like any machine learning model, consists of standard hyperparameters such as learning rate and batch size. Specifically, excessively high learning rates can hinder convergence, while overly low values result in slow optimization. The choice of batch size is also a trade-off between efficiency and effectiveness.

In addition, CF model has specific and essential hyperparameters, which may not be so sensitive in other machine learning models. Embedding dimension for users and items influence the representability of CF models. Besides, the embedding size determines the model’s capacity to store all information of users and items.

In general, too-large embedding dimension leads to over-fitting, and too-small embedding dimension cannot fit the complex user-item interaction data. The regularization term is always adopted to address the over-fitting problem.

4 SEARCH STRATEGY

As mentioned in Section 1, joint search on hyperparameters and architectures have two challenges in designing a search strategy effectively and efficiently: the large search space and the high evaluation costs on large network datasets. Previous research works on model design with joint search space of hyperparameters and architectures such as [8, 43] consider the search problem in the manner of Figure 1a. They considered the search component in space jointly and used different search algorithms to make the proper choice. The search procedure can be costly on a large search space and dataset. Therefore, addressing the challenges of a vast search space and the costly evaluation process is crucial.

The first challenge means the joint search space of \mathcal{H} and \mathcal{A} described in Section 3 is large for accurate search. Thus, we need to reduce the search space. In practice, we choose to screen \mathcal{H} choices by comparing relative ranking in controlled variable experiments, which is explained in Section 4.1. The second challenge means the evaluation time cost in Equation (2) is high. A significant validation time will lower the search efficiency. Thus, we apply a sampling method on datasets by interaction frequency, elaborated in Section 4.2.

In comparison to conventional joint search problem in Figure 1a, we design a two-stage search algorithm in Figure 1b. We use Random Forest (RF) Regressor, a surrogate model, to improve the efficiency of the search algorithm. To transfer the knowledge, including the relation modeling between hyperparameters, architectures, and evaluated performance, we learn the surrogate model’s parameters θ in the first stage, and we use θ as initialization of RF model in the

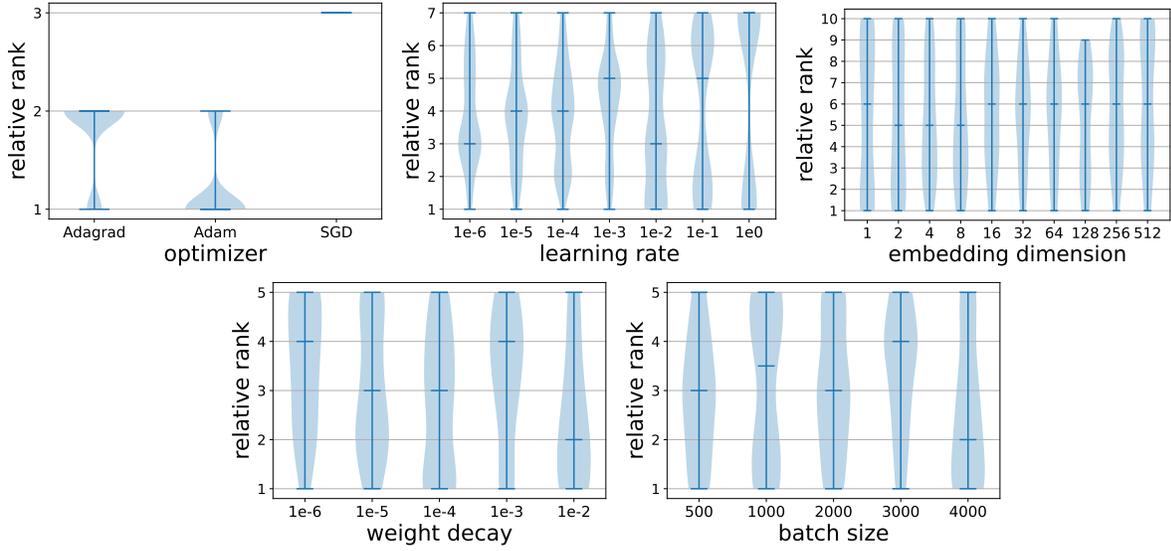


Figure 2: Ranking distribution of hyperparameters.

second stage. Our search algorithm is shown in Algorithm 2 and described in detail in Section 4.3. Besides, we have a discussion on our choices in Section A.3, elaborating how we solve the challenge in the joint search problem.

4.1 Screening Hyperparameter Choices

We screen the hyperparameter (HP) choices from \mathcal{H} to $\hat{\mathcal{H}}$ with two techniques. First, we shrink the HP space by comparing relative performance ranking of a special HP while fixing the others. After we get the ranking distribution of different HPs, we find the performance distribution among different choices of a HP, thus we may find the proper range or shrunk set of a given HP. Second, we decouple the HP space by calculating the consistency of different HP. If the consistency of a HP is high, that means the performance can change positively or negatively by only alternating this HP. Thus, this HP can be tuning separately neglecting its relation with other HPs in HP set.

4.1.1 Shrink the hyperparameter space. The screening method on hyperparameter space is based on analysis of ranking distribution on the performance with fixed value for a certain HP and random choices for other HPs and architectures. In this part, we denote a selection of hyperparameter $h \in \mathcal{H}$ as a vector, noted as $h = (h^{(1)}, h^{(2)}, \dots, h^{(n)})$. For instance, $h^{(1)}$ means optimizer, and $h^{(2)}$ means learning rate.

To obtain the ranking distribution of a certain HP $h^{(i)}$, we start with a controlled variable experiment. We vary $h^{(i)}$ in discrete values as the third column in Table 2, and we vary other HPs in original range as the second column. Specifically, given H_i as a discrete set of HP $h^{(i)}$, we choose a value $\lambda \in H_i$ and we can obtain the ranking $\text{rank}(h, \lambda)$ of the anchor HP $h \in \mathcal{H}_i$ by fixing other HPs except the i -th HP.

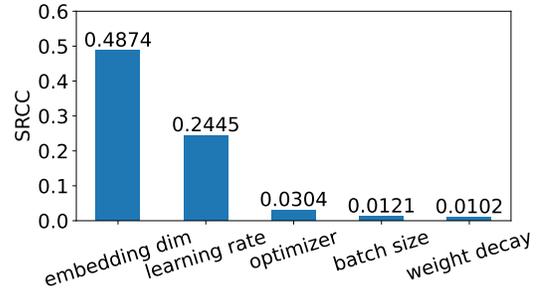


Figure 3: Consistency of each hyperparameters.

To ensure a fair evaluation of different architecture, we traverse the architecture space and calculate rank of performance with different configurations, then we can get the distribution of a type of HP. The relative performance ranking with different HPs is shown as violin plots in Figure 2. In this figure, we can get $\hat{\mathcal{H}}$ through the distribution of different HP values. We learn that the proper choice for optimizer can be shrunk to Adam and Adagrad; Proper range for learning rate is (1e-5, 1e-2); Proper range for embedding dimension can be reduced to [2, 64]; And we can fix weight decay in our experiments. We demonstrate the conclusion in the fourth column in Table 2.

4.1.2 Decouple the hyperparameter space. To decouple the search space, we consider the consistency of the ranking of hyperparameters when only alternating a given hyperparameter. For the i -th element $h^{(i)}$ of $h \in \mathcal{H}$, we can change different values for $h^{(i)}$, and then we can decouple the search procedure of the i -th hyperparameter with others. We use Spearman Rank-order Correlation Coefficient (SRCC) to show the consistency of various types of HPs, which is defined in Equation (3).

$$\text{SRCC}(\lambda_1, \lambda_2) = 1 - \frac{\sum_{h \in \mathcal{H}_i} |\text{rank}(h, \lambda_1) - \text{rank}(h, \lambda_2)|^2}{|\mathcal{H}_i| \cdot (|\mathcal{H}_i|^2 - 1)}. \quad (3)$$

where $|\mathcal{H}_i|$ means the number of anchor hyperparameters in \mathcal{H}_i . SRCC demonstrates the matching rate of rankings for anchor hyperparameters in \mathcal{H}_i with respect to $h^{(i)} = \lambda_2$ and $h^{(i)} = \lambda_1$.

The SRCC of the i -th HP is evaluated by the average of $\text{SRCC}(\lambda_1, \lambda_2)$ among different pairs of $\lambda \in H_i$, as is shown in Equation (4).

$$\text{SRCC}_i = \frac{1}{|\mathcal{H}_i|^2} \sum_{(\lambda_1, \lambda_2) \in H_i \times H_i} \text{SRCC}(\lambda_1, \lambda_2). \quad (4)$$

The SRCC results is demonstrated in Figure 3, we can directly find that the important hyperparameter with higher SRCC has a more linear relationship with its values. Since high embedding dimension model is time-costly during training and evaluation, we decide to use lower dimension in the first stage to reduce validation time, and then raise them on the original dataset to get better performance.

To summarize, we shrink the range of HP search space and find the consistency of different HPs. The shrunk space shown in Table 2 help us search more accurately, and the consistency analysis on performance ranking help us find the dependency between different HPs, thus we can tune HP with high consistency separately.

4.2 Evaluating Architecture with Sampling

To evaluate architecture more efficiently, we collect performance information evaluated from subgraphs, since the subgraph can approximate the properties of whole graph [46]. In this part, we introduce our frequency-based sampling method, and then we show the transfer ability of subgraphs by testing the consistency of architectures' performance from subsampled dataset to origin dataset.

4.2.1 Matrix sampling method. Since dataset for CF based on interaction of records, our sampling method is based on items' appearance frequency. That is, we can subsample the original dataset when we preserve part of the bipartite graph, and the relative performance on smaller datasets should have a similar consistency (i.e. ranking distribution of performance on S_{val} and \hat{S}_{val}).

The matrix subsample algorithm is demonstrated in Algorithm 1. First, we set the number of user-item interactions to be preserved first, which can be controlled by a sample ratio γ , $\gamma \in (0, 1)$. We calculate the interactions for each item, and then we preserve the item with a higher frequency. The items can be chosen in a fixed list (i.e. topk, Line 6-7 in Algorithm 1), or the interaction frequency count of items can be normalized to a probability, then different items have corresponding possibility to be preserved (i.e. distribute, Line 8-10 in Algorithm 1).

4.2.2 Transfer ability of architecture on subsampling matrix. To ensure that the relative performance ranking on subsampled datasets is similar to that on original datasets, we need to test the consistency of architecture ranking on different datasets. We evaluate the transfer ability among from subgraph to whole graph by SRCC.

For a given value of γ , we choose to select a sample set of architecture from $A_\gamma \in \mathcal{A}$. Then we evaluate them on a subsampled dataset \hat{S} and origin dataset S . The relative rank of $\alpha \in A_\gamma$ on \hat{S} and S is noted as $\text{rank}(\alpha, \hat{S})$ and $\text{rank}(\alpha, S)$.

Algorithm 1 MatrixSample Sampling Algorithm

Input: Matrix Dataset $S \in \mathbb{R}^{M \times N}$ (M : number of users; N : number of items), sample ratio $\gamma \in (0, 1)$, subsample mode (topk or distribute)

Output: Subsampled Matrix Dataset \hat{S}

- 1: Initialize item frequency set $f \leftarrow \{\}$;
 - 2: **for** $j = 1, 2, \dots, N$ **do**
 - 3: Calculate the users interacted with item j , record it as f_j ;
 - 4: $f \leftarrow f \cup f_j$;
 - 5: **end for**
 - 6: **if** subsample mode == topk **then**
 - 7: Rank f , choose top γN items in \hat{f} , neglect unrelated users;
 - 8: **else if** subsample mode == distribute **then**
 - 9: Calculate probability $\beta_j = f_j / \sum_{k=1}^N f_k$;
 - 10: Select \hat{f} with respect to $\{\beta_j\}$;
 - 11: **end if**
 - 12: Construct \hat{S} with f ;
 - 13: **return** \hat{S} .
-

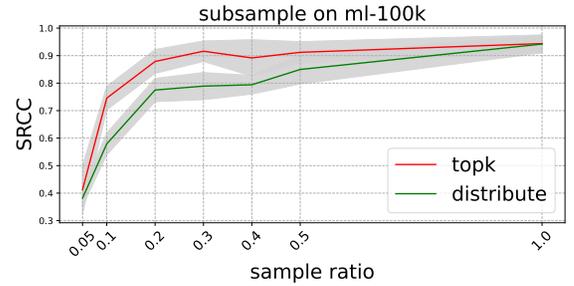


Figure 4: Consistency of different sample ratio. The light grey area is the standard deviation of SRCC on different \hat{S} .

$$\text{SRCC}_\gamma = 1 - \frac{\sum_{\alpha \in A_\gamma} |\text{rank}(\alpha, \hat{S}) - \text{rank}(\alpha, S)|^2}{|A_\gamma| \cdot (|A_\gamma|^2 - 1)}. \quad (5)$$

We can choose different subsampled dataset \hat{S} to get average consistency. As is demonstrated in Figure 4, sample ratio with higher SRCC has better transfer ability among graphs with sample mode topk, and the proper sample ratio should be in $\gamma \in [0.2, 1)$.

To summarize, through the sampling method, the evaluation cost will be reduced, and thus the search efficiency is improved. Since the ranking distribution among subgraphs is similar to that of the original dataset, we can transfer the evaluation modeling from small to large dataset.

4.3 Two-stage Joint Hyperparameter and Architecture Search

As discussed above, the evaluation cost can be highly reduced with sampling methods. Since the sampling method ensure the transfer ability from subgraphs to whole graphs, we propose a two-stage joint search algorithm, shown in Algorithm 2, and the framework is also shown in Figure 1b. The main notations in algorithm can be

found in Table A1 in Appendix. We also compare our method with conventional method in Figure A1.

To briefly summarize our algorithm: In the first stage (Lines 3-14), we sample several subgraphs with frequency-based methods, and preserve $\gamma = 0.2$ of interactions from the original rating matrix. Based our understanding of hyperparameters in Section 4.1, we select architecture and its hyperparameters in our reduced hyperparameter space.

We use a surrogate model to find proper architecture and hyperparameters to improve search efficiency, noted as $c(\cdot, \theta)$, θ is the parameter of the model. After we get evaluations of architecture and hyperparameters, we update parameters θ of c . In our framework, we choose BORE [35] as a surrogate and Random Forest (RF) as the regressor to modulate the relation between search components and performance. Details of the BORE+RF search algorithm (RFSurrogate) is shown in Algorithm A1 in Appendix.

We first transfer the parameters of c trained in the first stage by collected configurations and performance. Then, in the second stage (Lines 15-22), we select architectures and hyperparameters and evaluate them in the original dataset. Besides, we increase the embedding dimension to reach better performance. Similarly, the configurations and performance are recorded for the surrogate model (RF+BORE), which can give us the next proper configuration. Finally, after two-stage learning on the subsampled and original datasets, we get the final output of architecture and hyperparameters as the best choice of architecture and its hyperparameters.

5 EXPERIMENTS

Extensive experiments are performed to evaluate the performance of our search strategy by answering the following several research questions:

- **RQ1:** How does our algorithm work in comparison to other CF models and automated model design works?
- **RQ2:** How efficiently does our search algorithm work in comparison to other typical search algorithms?
- **RQ3:** What is the impact of every part of our design?
- **RQ4:** How specific is our model for different CF tasks?

5.1 Experimental Settings

5.1.1 Datasets. We use MovieLens-100K, MovieLens-1M, Yelp, and Amazon-Book for CF tasks. The detailed statistics and preprocess stage of datasets are shown in Table A2 in Appendix A.4.

5.1.2 Evaluation metrics. As for evaluation metrics, we choose two widely used metrics for CF tasks, Recall@ K and NDCG@ K . According to recent works [15, 36], we set the length for recommended candidates K as 20. We use Recall@20 and NDCG@20 as validation. As for loss function, we use BPR loss [29], the state-of-the-art loss function for optimizing recommendation models.

5.1.3 Baselines for Comparison. Since we encode both hyperparameters and architectures, we can use previous search algorithms on hyperparameters [1, 9, 32, 35] and extend them on a joint search space. The details of search algorithms can be found in Appendix A.5.1.

Algorithm 2 Joint Search Algorithm.

Input: Train/valid set S_{tra}, S_{val} , hyperparameter space \mathcal{H} , architecture space \mathcal{A} , sample number L , surrogate model $c(\cdot, \theta)$

Output: Best hyperparameters h^* , architecture α^*

```

1: Screen  $\mathcal{H}$  to  $\hat{\mathcal{H}}$  according to Section 4.1;
2: Initialize configuration set  $D \leftarrow \{\}, k, j \leftarrow 0$ ;
   %stage one start
3: for  $i = 1, 2, \dots, L$  do
4:   Sample matrix  $\hat{S}_i = \text{MatrixSample}(\mathcal{S}, \gamma)$ ;
5:   Split train/valid/test dataset  $\hat{S}_i = \hat{S}_{i,tra} \cup \hat{S}_{i,val} \cup \hat{S}_{i,tst}$ ;
6:   Initialize surrogate model  $c(\cdot; \theta)$ ;
7:   while not converge do
8:     Select  $\alpha_k$  and  $h_k$  by RFSurrogate;
9:     Evaluate  $p_k \leftarrow \mathcal{M}(f(\mathbf{P}^*; \alpha_k, h_k), \hat{S}_{i,val})$ ;
10:    Save to set  $D \leftarrow D \cup \{\alpha_k, h_k, p_k\}$ ;
11:    Update  $c(\cdot; \theta)$  with  $D$  by RFSurrogate;
12:     $k \leftarrow k + 1$ ;
13:  end while
14: end for
15: Transfer parameters  $\theta$  in  $c$ ;
   %stage two start
16: while not converge do
17:   Select  $\alpha_j$  and  $h_j$  by RFSurrogate;
18:   Evaluate  $p_j \leftarrow \mathcal{M}(f(\mathbf{P}^*; \alpha_j, h_j), S_{val})$ ;
19:   Save to set  $D \leftarrow D \cup \{\alpha_j, h_j, p_j\}$ ;
20:   Update  $c(\cdot; \theta)$  with  $D$  by RFSurrogate;
21:    $j \leftarrow j + 1$ ;
22: end while
23: return  $\alpha^*, h^*$ .
```

5.2 Performance Comparison (RQ1)

For CF tasks, we compare our results with NN-based CF models and Graph-based CF models. Besides, we also compare our search algorithm with other search algorithms designed for CF models. The search space is based on analysis of hyperparameter understanding in Section 4.1. We report the performance on four datasets in Table 3.

We summarize the following observation:

- We find that in our experiment, our CF model trained by searched hyperparameters and architectures can achieve better performance than the classical CF models. Some single models also perform well due to their special design. For example, LightGCN performs well on ML-100K and Yelp, while NGCF also performs well on ML-1M. Since our search algorithm has included various operations used in CF architectures, the overall architecture space can cover the architectures of classical models.
- Compared to NAS method on CF models, our method works better than SIF for considering multiple operations in different stages of architecture. Our methods also outperform AutoCF by 3.10% to 12.1%. The reason for that is we have an extended joint search space for both hyperparameters and architectures. Besides, our choice for hyperparameters is shrunk to a proper range to improve search efficiency and performance.
- We can observe in Table 3 that our searched models can achieve the best performance compared with all other baselines. Note that our proposed method can outperform the best baseline by

Table 3: Comparison of different methods on CF tasks.

Dataset	MovieLens-100K		MovieLens-1M		Yelp		Amazon-Book	
Metric	Recall@20	NDCG@20	Recall@20	NDCG@20	Recall@20	NDCG@20	Recall@20	NDCG@20
MF [21]	0.1145	0.1179	0.0896	0.0584	0.0307	0.0286	0.0291	0.0213
FISM [19]	0.1434	0.1422	0.0995	0.0621	0.0528	0.0316	0.0302	0.0211
NCF [15]	0.1980	0.1521	0.1204	0.0684	0.0534	0.0335	0.0315	0.0223
J-NCF [48]	0.2016	0.1448	0.1265	0.0629	0.0636	0.0393	0.0351	0.0252
Pinsage [42]	0.1561	0.1421	0.1354	0.0631	0.0561	0.0417	0.0321	0.0239
NGCF [36]	0.1654	0.1479	0.1678	0.0852	0.0573	0.0454	0.0349	0.0247
LightGCN [14]	0.2342	0.1755	0.1637	0.0842	0.0689	0.0471	0.0355	0.0273
SIF [39]	0.1935	0.1532	0.1294	0.0695	0.0581	0.0459	0.0350	0.0249
AutoCF [11]	0.2259	0.1703	0.1309	0.0739	0.0643	0.0467	0.0354	0.0265
Ours	0.2647	0.1913	0.1787	0.0945	0.0721	0.0482	0.0365	0.0281
Improvement	13.02%	9.00 %	6.50 %	10.92%	4.64%	2.33%	2.82%	2.93%

Table 4: Top 3 architectures for each datasets.

Dataset	Top-1	Top-2	Top-3
ML-100K	(H, H, SGC, SGC, min, VEC)	(ID, ID, SGC, SGC, multiply, SUM)	(H, H, Mat, Mat, multiply, VEC)
ML-1M	(H, H, SGC, SGC, min, VEC)	(H, H, HadamardGCN, HadamardGCN, min, VEC)	(H, H, Mat, Mat, multiply, VEC)
Yelp	(H, H, SGC, SGC, multiply, MLP)	(H, H, SGC, SGC, multiply, VEC)	(H, H, MLP, MLP, multiply, VEC)
Amazon-Book	(H, H, SGC, SGC, min, MLP)	(ID, ID, SGC, SGC, multiply, VEC)	(ID, H, Mat, MLP, max, VEC)

2.33% to 13.02%. The performance improvement on small datasets is better than that on large ones.

5.3 Algorithm Efficiency (RQ2)

We compare the different search algorithms mentioned in Section 5.1.3. The search results are shown in Figure 5a on dataset ML-100K and Figure 5b on dataset ML-1M. We plot our results by the output of the search algorithm. As is demonstrated in Figure 5b and 5a, search algorithms of BO perform better than random search, since the BO method considers a Gaussian Process surrogate model for simulating the relationship between performance output and hyperparameters and architectures. We find that BORE+RF outperforms other search strategies in efficiency. The reason is that the surrogate model RF can better classify one-hot encoding of architectures. Besides we also compare different time curves for BORE+RF in single-stage and two-stage evaluations. Since our two-stage algorithm uses subsampling method on rating matrix, and ensure the consistency between subgraph and whole graph, we can achieve better performance and higher search efficiency.

5.4 Ablation Study (RQ3)

In this subsection, we analyze how important and sensitive the various components of our framework are. We focus on the performance’s improvement and efficiency by reducing and decoupling the hyperparameter space. We also elaborate on how we choose sample ratio and effectiveness of tuning hyperparameters.

5.4.1 Plausibility of screening hyperparameters. To further validate the effectiveness of our design of screening hyperparameter choices,

we choose hyperparameters on the origin space, shrunk space, and decoupled space. The search time curve is shown in Figure 5c.

We demonstrate that screening hyperparameter choices can help improve performance on CF tasks and search efficiency. According to our results, performance on \mathcal{H} is better than the one on the origin space \mathcal{H} . The reason is that the shrunk hyperparameter space has a more significant possibility of including better HP choices for CF architectures to achieve better performance. The search efficiency is improved since the choice of hyperparameters is reduced, and proper hyperparameters can be found more quickly in a smaller search space. We also find that the search efficiency reduces when we increase the batch size and embedding dimension when we search on the decoupled search space. While increasing batch size and embedding dimension help improve final performance, the cost evaluation is higher, which may reduce search efficiency.

5.4.2 Choice of Sampling Ratio. To find the impact of changing sampling ratio, we search with different sampling ratios on different datasets in the same controlled search time. The evaluation results (Recall@20) for experiments on different sampling ratio settings are listed in Table 5.

According to the table, smaller sampling ratios have lower performance. The reason is that the user-item matrix sampled by low sample ratio may not capture the consistency in the origin matrix, as the consistency results shown in Section 4.2. While sampled dataset generated by higher sample ratio has more considerable consistency with the original dataset, the results in a limited time may not be better. The reason is that too much time on evaluation in the first stage may have fewer results for the surrogate model

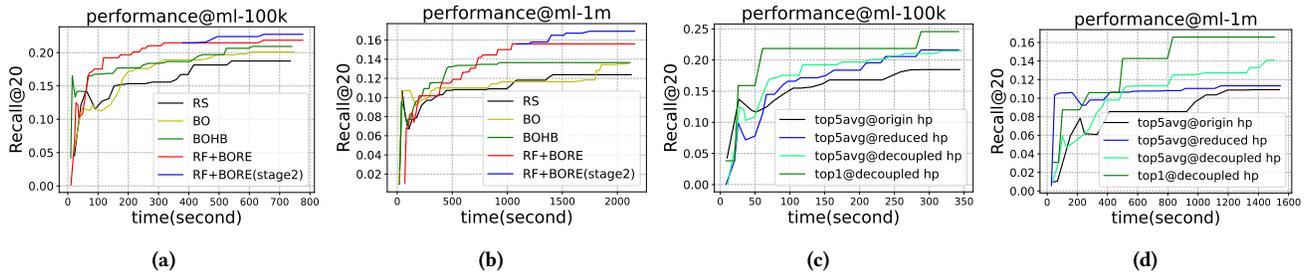


Figure 5: Comparison for different search algorithms and HP space. (a)-(b) Time curve of different search algorithms on ML-100K and ML-1M; (c)-(d) ablation studies on our search strategy in different search space on ML-100K and ML-1M.

Table 5: Performance with different sample ratio γ .

sample ratio	5%	10%	20%	50%
ML-100K	0.2113	0.2224	0.2646	0.2623
ML-1M	0.1471	0.1632	0.1787	0.1715
Yelp	0.0523	0.0654	0.0721	0.0706
Amazon-Book	0.0311	0.0342	0.0365	0.0356

Table 6: Comparison on search algorithm with/without HP tuning.

Performance	AutoCF	AutoCF (HP)	Improvement
ML-100K	0.2259	0.2335	3.36%
ML-1M	0.1309	0.1358	3.74%
Yelp	0.0643	0.0672	4.51%
Amazon-Book	0.0354	0.0359	1.41%

to learn the connection between performance and configurations of hyperparameters and architectures. Thus, the first stage has a trade-off between sample ratio and time cost, and we choose 20% as our sample ratio in our experiments.

5.4.3 Tuning on Hyperparameters. To show the effectiveness of our design on joint search, we propose to apply our joint search method to previous research work AutoCF [11]. The results are shown in Table 6. We apply our search strategy to AutoCF to include hyperparameters in our search space, noted as AutoCF(HP). We find that hyperparameter searched on shrunk space can perform better than the one that uses architecture space and random search on different datasets.

5.5 Case Study (RQ4)

In this part, we mainly focus on our search strategy’s search results of different architectures. The search results with top performance share some similarities, while different architecture operations have different performances.

According to search results, we present the top models for each task on all datasets in Table 4. It is easy to find that interaction history-based encoding features may have better performance and more powerful representative ability. Besides, the embedding function of SGC has stronger representative ability since it can capture

the high-order relationship. Both SGC and HadanardGCN collect information from different layers, simply designed SGC have stronger ability. As for the interaction function, we find the classical element-wise multiply can receive strong performance; The prediction function of learnable vector and MLP can capture more powerful and complex interaction than SUM. We can find that these top models of each task have similar implementations, but there also have some differences among different datasets. One top architecture on a given dataset may not get the best performance on another one.

In summary, the proper architectures for different datasets may not be the same, but these top results may share some same operations in architecture structures. The result and analysis can help human experts to design more powerful CF architectures.

6 CONCLUSION AND FUTURE WORK

In this work, we consider a joint search problem on hyperparameters and architectures for Collaborative Filtering models. We propose a search framework based on a search space consisting of frequently used hyperparameters and operations for architectures. We make a complete understanding of hyperparameter space to screen choices of hyperparameters, We propose a two-stage search algorithm to find proper hyperparameters and architectures configurations efficiently. We design a surrogate model that can jointly update CF architectures and hyperparameters and can be transferred from small to large datasets. We do experiments on several datasets, including comparison on different models, search algorithms and ablation study.

For future work, we find it important to model CF models based on Knowledge Graphs as a search problem. With additional entities for items and users, deep relationships can be mined for better performance. An extended search framework can be built on larger network settings. Models on extensive recommendation tasks and other data mining tasks can also be considered as a search problem.

ACKNOWLEDGMENTS

This work is partially supported by the National Key Research and Development Program of China under 2021ZD0110303, the National Natural Science Foundation of China under 62272262, 61972223, U1936217, and U20B2060, and the Fellowship of China Postdoctoral Science Foundation under 2021TQ0027 and 2022M710006.

REFERENCES

- [1] J. Bergstra and Y. Bengio. 2012. Random search for hyper-parameter optimization. *JMLR* 13, Feb (2012), 281–305.
- [2] James Bergstra, Dan Yamins, David D Cox, et al. 2013. Hyperopt: A python library for optimizing the hyperparameters of machine learning algorithms. In *Proceedings of the 12th Python in science conference*, Vol. 13. Citeseer, 20.
- [3] Jinhang Cai, Yimin Ou, Xiu Li, and Haoqian Wang. 2021. ST-NAS: Efficient Optimization of Joint Neural Architecture and Hyperparameter. In *Neural Information Processing*, Teddy Mantoro, Minh Lee, Media Anugerah Ayu, Kok Wai Wong, and Achmad Nizar Hidayanto (Eds.). Springer International Publishing, Cham, 274–281.
- [4] Bo Chen, Xiangyu Zhao, Yejing Wang, Wenqi Fan, Huifeng Guo, and Ruiming Tang. 2022. Automated Machine Learning for Deep Recommender Systems: A Survey. <https://doi.org/10.48550/ARXIV.2204.01390>
- [5] Wanyu Chen, Fei Cai, Honghui Chen, and Maarten De Rijke. 2019. Joint Neural Collaborative Filtering for Recommender Systems. *ACM Transactions on Information Systems (TOIS)* 37, 4 (2019), 1–30.
- [6] Marc Claesen and Bart De Moor. 2015. *Hyperparameter search in machine learning*. Technical Report.
- [7] Maurizio Ferrari Dacrema, Paolo Cremonesi, and Dietmar Jannach. 2019. Are we really making much progress? A worrying analysis of recent neural recommendation approaches. In *Proceedings of the 13th ACM Conference on Recommender Systems*. 101–109.
- [8] Xuanyi Dong, Mingxing Tan, Adams Wei Yu, Daiyi Peng, Bogdan Gabrys, and Quoc V Le. 2020. AutoHAS: Efficient hyperparameter and architecture search. *arXiv preprint arXiv:2006.03656* (2020).
- [9] Stefan Falkner, Aaron Klein, and Frank Hutter. 2018. BOHB: Robust and efficient hyperparameter optimization at scale. In *ICML*. PMLR, 1437–1446.
- [10] Matthias Feurer, Aaron Klein, Katharina Eggenberger, Jost Tobias Springenberg, Manuel Blum, and Frank Hutter. 2019. Auto-sklearn: efficient and robust automated machine learning. In *Automated Machine Learning*. Springer, 113–134.
- [11] Chen Gao, Quanming Yao, Depeng Jin, and Yong Li. 2021. Efficient Data-specific Model Search for Collaborative Filtering. In *KDD*. 415–425.
- [12] Chen Gao, Yu Zheng, Nian Li, Yinfeng Li, Yingrong Qin, Jinghua Piao, Yuhan Quan, Jianxin Chang, Depeng Jin, Xiangnan He, et al. 2023. A survey of graph neural networks for recommender systems: challenges, methods, and directions. *ACM Transactions on Recommender Systems* 1, 1 (2023), 1–51.
- [13] William L Hamilton, Rex Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. In *NIPS*. 1025–1035.
- [14] Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yongdong Zhang, and Meng Wang. 2020. Lightgcn: Simplifying and powering graph convolution network for recommendation. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval (SIGIR)*. 639–648.
- [15] Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. 2017. Neural collaborative filtering. In *International World Wide Web Conference (WWW)*. 173–182.
- [16] Xiangnan He, Hanwang Zhang, Min-Yen Kan, and Tat-Seng Chua. 2016. Fast matrix factorization for online recommendation with implicit feedback. In *International Conference on Research and Development in Information Retrieval (SIGIR)*. 549–558.
- [17] Frank Hutter, Holger H Hoos, and Kevin Leyton-Brown. 2011. Sequential model-based optimization for general algorithm configuration. In *LION*. Springer, 507–523.
- [18] Frank Hutter, Lars Kotthoff, and Joaquin Vanschoren. 2019. *Automated Machine Learning*. Springer.
- [19] Santosh Kabbur, Xia Ning, and George Karypis. 2013. Fism: factored item similarity models for top-n recommender systems. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*. 659–667.
- [20] Yehuda Koren. 2008. Factorization meets the neighborhood: a multifaceted collaborative filtering model. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*. 426–434.
- [21] Yehuda Koren, Robert Bell, and Chris Volinsky. 2009. Matrix factorization techniques for recommender systems. *IEEE Transactions on Computers (Computer)* 42, 8 (2009).
- [22] Lisha Li, Kevin G Jamieson, Giulia DeSalvo, Afshin Rostamizadeh, and Ameet Talwalkar. 2017. Hyperband: Bandit-based configuration evaluation for hyperparameter optimization. In *ICLR (Poster)*.
- [23] Liam Li and Ameet Talwalkar. 2019. Random search and reproducibility for neural architecture search. *arXiv preprint arXiv:1902.07638* (2019).
- [24] Bin Liu, Chenxu Zhu, Guilin Li, Weinan Zhang, Jincai Lai, Ruiming Tang, Xiquang He, Zhenguo Li, and Yong Yu. 2020. AutoFIS: Automatic Feature Interaction Selection in Factorization Models for Click-Through Rate Prediction. (2020). <https://doi.org/10.48550/ARXIV.2003.11235>
- [25] Chenxi Liu, Liang-Chieh Chen, Florian Schroff, Hartwig Adam, Wei Hua, Alan L Yuille, and Li Fei-Fei. 2019. Auto-deeplab: Hierarchical neural architecture search for semantic image segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 82–92.
- [26] H. Liu, K. Simonyan, and Y. Yang. 2019. DARTS: Differentiable architecture search. In *ICLR*.
- [27] Siyi Liu, Chen Gao, Yihong Chen, Depeng Jin, and Yong Li. 2021. Learnable Embedding Sizes for Recommender Systems. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=vQzcqQWIS0q>
- [28] H. Pham, M. Guan, B. Zoph, Q. Le, and J. Dean. 2018. Efficient Neural Architecture Search via Parameter Sharing. In *ICML*. 4092–4101.
- [29] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2009. BPR: Bayesian personalized ranking from implicit feedback. In *Conference on Uncertainty in Artificial Intelligence (UAI)*. 452–461.
- [30] Badrul Sarwar, George Karypis, Joseph Konstan, and John Riedl. 2001. Item-based collaborative filtering recommendation algorithms. In *International World Wide Web Conference (WWW)*. 285–295.
- [31] Jonas Seng, Pooja Prasad, Martin Mundt, Devendra Singh Dhami, and Kristian Kersting. 2023. FEATHERS: Federated Architecture and Hyperparameter Search. *arXiv:2206.12342* [cs.LG]
- [32] Jasper Snoek, Hugo Larochelle, and Ryan P Adams. 2012. Practical bayesian optimization of machine learning algorithms. *NIPS* 25 (2012).
- [33] David R So, Chen Liang, and Quoc V Le. 2019. The evolved transformer. *arXiv preprint arXiv:1901.11117* (2019).
- [34] Mingxing Tan and Quoc V Le. 2019. Efficientnet: Rethinking model scaling for convolutional neural networks. *arXiv preprint arXiv:1905.11946* (2019).
- [35] Louis C Tiao, Aaron Klein, Matthias Seeger, Edwin V Bonilla, Cedric Archambeau, and Fabio Ramos. 2021. BORE: Bayesian Optimization by Density-Ratio Estimation. In *ICML*.
- [36] Xiang Wang, Xiangnan He, Meng Wang, Fuli Feng, and Tat-Seng Chua. 2019. Neural Graph Collaborative Filtering. In *International Conference on Research and Development in Information Retrieval (SIGIR)*.
- [37] Zhenyi Wang, Huan Zhao, and Chuan Shi. 2022. Profiling the Design Space for Graph Neural Networks based Collaborative Filtering. (2022).
- [38] Jiancan Wu, Xiang Wang, Fuli Feng, Xiangnan He, Liang Chen, Jianxun Lian, and Xing Xie. 2021. Self-supervised graph learning for recommendation. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 726–735.
- [39] Quanming Yao, Xiangning Chen, James Kwok, and Yong Li. 2020. Efficient Neural Interaction Functions Search for Collaborative Filtering. In *WWW*. *arXiv preprint arXiv:1906.12091*.
- [40] Q. Yao and M. Wang. 2018. *Taking human out of learning applications: A survey on automated machine learning*. Technical Report. Arxiv: 1810.13306.
- [41] Quanming Yao, Ju Xu, Wei-Wei Tu, and Zhanxing Zhu. 2019. Differentiable Neural Architecture Search via Proximal Iterations. *arXiv preprint arXiv:1905.13577* (2019).
- [42] Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L Hamilton, and Jure Leskovec. 2018. Graph convolutional neural networks for web-scale recommender systems. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD)*. ACM, 974–983.
- [43] Arber Zela, Aaron Klein, Stefan Falkner, and Frank Hutter. 2018. Towards automated deep learning: Efficient joint neural architecture and hyperparameter search. *arXiv preprint arXiv:1807.06906* (2018).
- [44] Yongqi Zhang and Quanming Yao. 2022. *Knowledge Graph Reasoning with Relational Digraph*. Technical Report.
- [45] Yongqi Zhang, Quanming Yao, Wenyuan Dai, and Lei Chen. 2019. AutoSF: Searching Scoring Functions for Knowledge Graph Embedding. *ICDE*.
- [46] Yongqi Zhang, Zhanke Zhou, Quanming Yao, and Yong Li. 2022. KGtuner: Efficient Hyper-parameter Search for Knowledge Graph Learning. <https://doi.org/10.48550/ARXIV.2205.02460>
- [47] Xiangyu Zhao, Haochen Liu, Hui Liu, Jiliang Tang, Weiwei Guo, Jun Shi, Sida Wang, Huiji Gao, and Bo Long. 2021. Autodim: Field-aware embedding dimension search in recommender systems. In *Proceedings of the Web Conference 2021*. 3015–3022.
- [48] Lei Zheng, Vahid Noroozi, and Philip S Yu. 2017. Joint deep modeling of users and items using reviews for recommendation. In *International Conference on Web Search and Data Mining (WSDM)*. 425–434.
- [49] Ruiqi Zheng, Liang Qu, Bin Cui, Yuhui Shi, and Hongzhi Yin. 2022. AutoML for Deep Recommender Systems: A Survey. <https://doi.org/10.48550/ARXIV.2203.13922>
- [50] Yu Zheng, Chen Gao, Xiang Li, Xiangnan He, Yong Li, and Depeng Jin. 2021. Disentangling user interest and conformity for recommendation with causal embedding. In *Proceedings of the Web Conference 2021*. 2980–2991.
- [51] B. Zoph and Q. Le. 2017. Neural architecture search with reinforcement learning. In *ICLR*.

A APPENDIX

A.1 Search Space

The main notations in this paper are listed in Table A1, and we discuss some details for search space in this section.

A.1.1 Hyperparameter Choice. We explain how to reduce the hyperparameters by Figure 2 in this section. We shrink hyperparameter search space based on the performance ranking distribution. We can split the hyperparameters into four categories:

- **Reduction:** This kind of hyperparameter is usually categorical, like optimizer. The choices of hyperparameters can be reduced.
- **Shrunk range:** This kind of hyperparameter is usually selected in a continuous range, such as learning rate. The choices of these values can be constrained to a smaller range.
- **Monotonous related:** The performance with this kind of hyperparameter usually rises when the hyperparameter increase, such as embedding dimension. However, we can not choose these HPs with too large values for limited memory. Thus, we choose a smaller value in first stage and a larger one in second stage.
- **No obvious pattern:** We do not have to change this kind of HP in our experiment, just as weigh decay.

A.2 Search Algorithms

A.2.1 Surrogate Model Design. We demonstrate our search algorithm with surrogate model in Algorithm A1.

We design our search algorithm with BORE and Random Forest (RF) regressor. In the first stage, we train the surrogate model with D , update parameters of surrogate model. The output of BORE+RF can help give an inference between $(0, 1)$, and we choose the least one as output. After the first stage, we save the parameters of this surrogate model, and transfer it to second stage, which we do experiment on larger datasets. The configuration of hyperparameters and architectures and the performance on larger dataset can also update the parameters of surrogate model. With the knowledge we learn on the first stage, the surrogate model can better choose the proper architecture and hyperparameter for CF tasks.

A.2.2 Fair Comparison. Compared with hyperparameters, the choice of architectures can affect the performance of CF models more. Thus, to compare different experiment settings fairly, we use the average of top5 configurations instead for the search time curve, noted as top5avg.

A.2.3 Search Procedure. We show the comparison of conventional joint search method and our method in Figure A1. In Figure A1, conventional one-stage method search components separately, while our method search hyperparameters on a shrunk space. Besides, the evaluation time is lower in the first stage in our method.

A.3 Discussion

In this section, we discuss about the method we have chosen, and the difference with previous works on joint search problems.

The first is the reason for screening in the hyperparameters space rather than the architecture space. The search space of hyperparameters mainly consists of components of continuous values with infinite choices. And screening range of hyperparameters [37, 44] is proven effective on deep neural networks and graph-based models.

Algorithm A1 RFSurrogate RF+BORE surrogate model

Input: Training data S_{Tra} and evaluation data S_{val} , Reduced hyperparameter space $\hat{\mathcal{H}}$, Architecture space \mathcal{A} , percentage threshold $\tau = 0.2$, RF regressor $y = c(\alpha||h; \theta)$

Output: A proper configuration of hyperparameters h^* and architecture α^*

```

1: Initialize configuration set  $D \leftarrow \{\}$ ,  $i \leftarrow 0$ 
   % Initialize surrogate model;
2: for  $i = 1, 2, \dots, N$  do
3:   Select architecture  $\alpha_i \in \mathcal{A}$ , hyperparameter  $h_i \in \mathcal{H}$ ;
4:   Get CF evaluation performance  $p_i = \mathcal{M}(f(\mathbf{P}^*; \alpha_i, h_i), S_{val})$ ;
5:   Save to set  $D \leftarrow D \cup \{\alpha_i, h_i, p_i\}$ ;
   % BORE
6:   Find  $p_\tau$  as the  $\tau$ -quantile of the performance set  $\{p_i\}$ ;
7:   if  $p_i < \tau$  then
8:     Set  $z_i$  label 0;
9:   else if  $p_i \geq \tau$  then
10:    Set  $z_i$  label 1;
11:   end if
12:   Fit the surrogate model with  $D$ ;
13: end for
   % Update surrogate model (with a new data);
14: Get performance  $p_k = \mathcal{M}(f(\mathbf{P}^*; \alpha_k, h_k), S_{val})$ ;
15: Save to set  $D \leftarrow D \cup \{\alpha_k, h_k, p_k\}$ ;
16: Set  $z_i$ ;
17: Fit to update RF model  $c(\cdot; \theta)$ ;
   % Find next choice through  $c(\cdot; \theta)$ ;
18:  $z_{\max} \leftarrow 0$ 
19: Randomly sample configuration set  $H = \{h_l, \alpha_l\}$ ;
20: for configuration in  $H$  do
21:    $z_l = c(\alpha_l||h_l; \theta)$ ;
22:   if  $z_l > z_{\max}$  then
23:      $h^*, \alpha^* = h_l, \alpha_l$ ;
24:   end if
25: end for
26: return  $h^*, \alpha^*$ 

```

Architecture space components are typically categorized, and each one is necessary in some manner and should not be ignored. Furthermore, we believe that it is unnecessary to reduce the architecture space after conducting a fair ranking of different architectures.

Our work is the first on CF tasks, which is different from previous research studies on joint search. [43] mainly focuses on a joint search problem on typical neural networks. Another study on the joint search problem, AutoHAS [8], focuses on model search with weight sharing. FEATHERS [31] focuses on a joint search problem on Federate Learning, an aspect of Reinforcement Learning. In ST-NAS [3], the sub-models are candidates of a designed super-net, and they sample sub-ST-models from the super-net and weights are update using training loss while updating HPs with the validation loss.

A.4 Data Preprocess

A.4.1 Details of Datasets. In our experiments, we choose four real-world raw datasets and build them for evaluation on CF tasks.

Table A1: Notations

Variable	Definition
h, h^*	Hyperparameters (HP), best HP from \mathcal{H}
\hat{h}	HP from $\hat{\mathcal{H}}$
$h^{(i)}$	The i -th HP of $h, h = (h^{(1)}, h^{(2)}, h^{(i)}, \dots, h^{(n)})$
H_i	Discrete values set of $h^{(i)}$
$\lambda, \lambda_1, \lambda_2$	Some values of hyperparameter from H_i
$\text{rank}(h, \lambda)$	Relative performance ranking when $h^{(i)} = \lambda$
\mathcal{H}	Origin space of hyperparameters
\mathcal{H}_i	Subset of \mathcal{H} , the i -th HP is selected from H_i
$\hat{\mathcal{H}}$	Subset of \mathcal{H} , screened space of hyperparameters
α, α^*	Architecture/Best architecture from \mathcal{A}
\mathcal{A}	Space of architecture
S	Origin dataset (whole graph)
\hat{S}	The i -th subsampled dataset (subgraph)
γ	Subsample ratio
p, p_i	Test performance
D	Records of (α, h, p)
$c(\cdot; \theta)$	Surrogate model with parameter θ

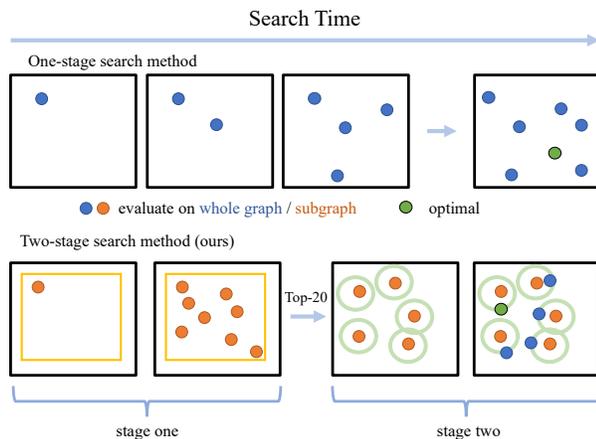


Figure A1: Comparison of search procedure using one-stage method and our method.

- **MovieLens-100K**¹ This widely used movie-rating dataset contains 100,000 ratings on movies from 1 to 5. We also convert the rating form to binary form, where each entry is marked as 0 or 1, indicating whether the user has rated the item.
- **MovieLens-1M**² This widely used movie-rating dataset contains more than 1 million ratings on movies from 1 to 5. Similarly, for MovieLens-1M, we build an implicit dataset.
- **Yelp**³ This is published officially by Yelp, a crowd-sourced review forum website where users can write comments and reviews for various POIs, such as hotels, restaurants, etc.

¹<https://grouplens.org/datasets/movielens/100k>²<https://grouplens.org/datasets/movielens/1m>³<https://www.yelp.com/dataset/download>

Table A2: Statistics of datasets

Name	#users	#items	#records	density
ML-100K	943	1,682	100,000	6.304%
ML-1M	6,040	3,952	1,000,209	4.190%
Yelp	6,102	18,599	445,576	0.393%
Amazon-Book	25,774	80,211	3,040,864	0.147%

- **Amazon-Book**⁴ This book-rating dataset is collected from users' uploaded review and rating records on Amazon.

A.4.2 Preprocess. Since the origin dataset in Table A2 is too large for training, we reduce them in frequency order. We use 10-core select on Yelp and 50-core select on Amazon-Book. N -core means we choose users and items that appear more than N times in the whole record history. After we select the dataset, we split the data into training, validation and test sets. We shuffle each set when we start a new evaluation task.

A.4.3 Dataset sampling. Some papers studying long-tail recommendation [50] or self-supervised recommendation [38] may discuss the impact of data sparsity. We can sample the subgraph according to the popularity [16], mainly user-side and item-side, and the long-tail effect on the item-side is more severe. Thus we sample the rating matrix based on the item frequency. The more the item appears in rating records, the more likely it is reserved in the subsampled matrix.

A.5 Detailed Settings of Experiments

A.5.1 Baseline algorithms.

- **Random Search** [1]. It is a simple method for finding hyperparameters and architectures with random choices on either continuous or categorical space.
- **Bayesian Optimization** [32]. Bayesian Optimization (BO) is a search algorithm based on the analysis of posterior information, and it can use Gaussian Process as a surrogate model.
- **BOHB** [9]. BOHB is a method consisting of both Bayesian Optimization (BO) and HyperBand (HB), helping modulate functions with a lower time budget.
- **BORE** [35]. BORE is a BO method considering expected improvement as a binary classification problem. It can be combined with regressors, including Random Forest (RF), Multi-Layer Perceptron (MLP), and Gaussian Process (GP). We choose RF as a regressor in our experiments.

A.5.2 Hardware Environment. We implement models with PyTorch 1.12 and run experiments on a 64-core Ubuntu 20.04 server with NVIDIA GeForce RTX 3090 GPU with 24 GB memories each. It takes 3.5-4 hours to search on a dataset with one million records.

A.5.3 Codes. We have released our implementation code for experiments in <https://github.com/overwenyan/Joint-Search>

⁴<https://nijianmo.github.io/amazon/index.html>