

## 百度2023-ColdNAS: 突破用户冷启动难题, 论文解析创新架构搜索框架!



SmartMindAI

专注搜索、广告、推荐、大模型和人工智能最新技术, 欢迎关注我

已关注

13 人赞同了该文章

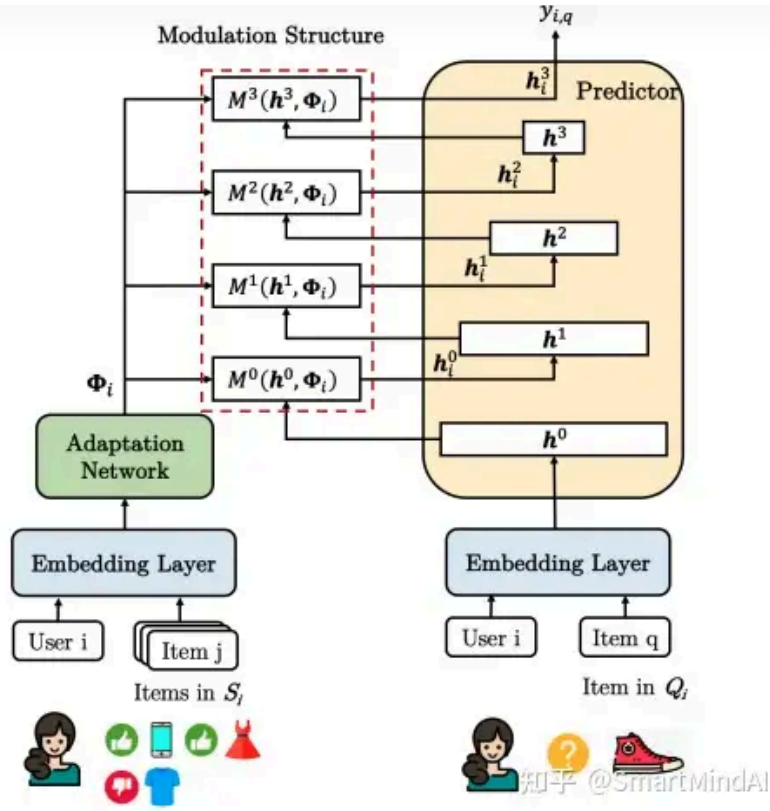
### Introduction

推荐系统<sup>\*</sup>旨在为用户提供相关的内容建议, 如电影和书籍。尽管现有的在线推荐<sup>+</sup>系统为大量用户提供了便捷的购物体验, 但冷启动推荐问题仍然是其中的一个挑战。

这个问题可以被建模为一种少样本学习问题, 其目标是通过少量训练样本快速泛化到新任务。许多论文采用了与模型无关元学习(MAML)策略, 它从一组任务中学习良好的初始化参数, 并通过在有限数量的带标签样本上进行几步骤梯度下降更新来适应新任务。然而, 这种策略需要专业知识来调整优化过程以避免过拟合<sup>+</sup>, 而且耗时可能很长。

- 本文提出了一个名为ColdNAS的架构搜索框架, 用于解决用户的冷启动问题。我们使用超网络将每个用户的历史交互映射到用户特定的参数, 然后使用这些参数来调整预测器, 并将如何调整和在哪里调整表述为一个NAS问题。
- 本文设计了一个架构搜索结构的搜索空间<sup>+</sup>, 该空间不仅可以涵盖现有的基于现有的用户冷启动模型, 而且可以包含更丰富的表达结构。由于搜索空间可能很大, 因此需要进行搜索空间转换<sup>+</sup>, 将原始空间转换为一个等效但更小的空间。我们提供了理论分析以验证其正确性。在转换后的空间上, 我们可以使用可微架构搜索算法进行高效和稳健的搜索。
- 本文在用户冷启动问题的基准数据集<sup>+</sup>进行了广泛的实验, 并观察到模型始终获得了最先进的表现。我们还验证了搜索空间和算法的设计考虑, 证明了ColdNAS的强大和合理性。

### Proposed Method



### Problem Formulation

我们设用户集为  $\mathcal{U} = \{u_i\}$ , 物品集为  $\mathcal{V} = \{v_j\}$ , 用户  $u_i$  对物品  $v_j$  的评分为  $y_{i,j}$ . 用户冷启动推荐问题的目标是为只评分过少数物品的用户  $u_i$  提供个性化推荐。这个问题可以被视为少样本学习问题, 我们需要学习一个模型, 该模型能从一组训练用户冷启动任务  $\mathcal{T}^{\text{train}}$  中学习, 并推广到新任务。每个任务  $T_i$  对应一个用户  $u_i$ , 并有一个包含现有交互历史的支持集  $\mathcal{S}_i = \{(v_j, y_{i,j})\}_{j=1}^N$  和一个查询集  $\mathcal{Q}_i = \{(v_j, y_{i,j})\}_{j=1}^M$ , 其中  $M$  表示查询集中的交互数量。

### Search Space

我们的用户冷启动方案包括嵌入层、适应网络和预测器。嵌入层将用户行为转化为向量表示, 适应网络关联用户历史行为与目标行为, 预测器预测用户的目标行为。

- 嵌入层  $E$  使用参数  $\theta_E$  将特征嵌入到向量中, 即

$$(u_i, v_j) = E(u_i, v_j; \theta_E)$$

- 具有参数  $\theta_A$  的适应网络接受特定用户的支持集  $\mathcal{S}$  作为输入, 并生成用户特定的自适应参数, 即

$$\Phi_i = \{\phi_i^k\}_{k=1}^C = A(\mathcal{S}_i; \theta_A),$$

- 预测器  $P$  使用参数  $\theta_P$ , 接收用户特定参数  $\phi_i$  和  $v_q \in \mathcal{Q}_i$  作为输入

$$\hat{y}_{i,j} = P((u_i, v_q), \Phi_i; \theta_P).$$

我们引入了一个额外的自适应网络来处理冷启动用户。对于每个  $u_i$ , 我们将支持集  $\mathcal{S}_i$  映射到用户特定参数  $\phi_i$ 。然后, 我们使用目标物品  $v_q \in \mathcal{Q}_i$  的特征、用户特征  $u_i$  和  $\phi_i$  进行预测。

最近, TaNP 提出了一种直接将式 (1) 中的  $M^l$  采用 FiLM 的形式应用于所有的  $L$  个 MLP 层的方法。这种方法可以控制每个用户在第  $l$  层中对  $h^l$  的个性化程度, 同时可学习权重  $W^l$  和  $b^l$  也在第  $l$  层中进行更新。

$$h_i^l = h^l \odot \phi_i^1 + \phi_i^2.$$

该空间的大小为  $6^{C \times L}$ 。更大的  $C$  会导致更大的搜索空间，虽然有更高的潜力包含适当的调整函数，但更难以有效地进行搜索。

### Search Strategy

我们的目标是设计一个高效的搜索算法<sup>+</sup>，用于在大型搜索空间中进行搜索。为了解决这个问题，我们提出将原始搜索空间转换为等价的但更小的空间。在这个新的空间中，我们设计了一个超级网络结构，以进行高效且稳健的可微搜索。

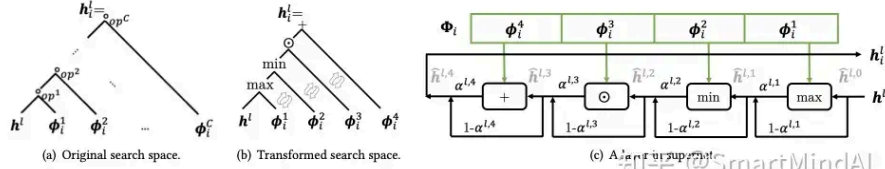


Figure 2: Illustration of our proposition 3.1, and the structure of the supernet to search on the reduced space.

### Search Space Transformation

尽管搜索空间可能很大，但我们可以将其转换为大小为  $2^{4 \times L}$  的等效空间，该空间与  $C$  无关，如命题证明所示（证明见附录）。

假设自适应网络  $A$  足够表达。任何形式为的  $M^l$ ，其中  $o_{op^k} \in \mathcal{O}$ ， $C$  是任何非负整数，和  $\phi_i^k \in \Phi_i = A(S_i, \theta_A)$ ，可以被表示为  $h_i^l = \min(\max(h^l, \hat{\phi}_i^1), \hat{\phi}_i^2) \odot \hat{\phi}_i^3 + \hat{\phi}_i^4$ ，并且上述四个操作是排列不变的。

为了更好地理解命题，让我们检查两个例子。

$$1. \min(\max(h^l, \phi_i^1) + \phi_i^2 - \phi_i^3, \phi_i^4) \odot \phi_i^5$$

$$\text{其中 } \hat{\phi}_i^1 = \phi_i^1, \hat{\phi}_i^2 = \phi_i^4 - \phi_i^2 + \phi_i^3, \hat{\phi}_i^3 = \phi_i^5, \hat{\phi}_i^4 = (\phi_i^2 - \phi_i^3) \odot \phi_i^5;$$

$$1. \max(\min(h^l + \phi_i^1, \phi_i^2), \phi_i^3) \odot \phi_i^4$$

$$\text{其中 } \hat{\phi}_i^1 = \phi_i^3 - \phi_i^1, \hat{\phi}_i^2 = \phi_i^2 - \phi_i^1, \hat{\phi}_i^3 = \phi_i^4, \hat{\phi}_i^4 = \phi_i^1 \odot \phi_i^4.$$

原始空间： $X \subseteq R^n$  变换后的空间： $Y \subseteq R^m$  其中： $m \ll n$  举例来说，对于一个具有100个特征的问题，我们可以通过深度网络将其映射到一个只有10个特征的空间中。这样，我们可以大大降低问题的维度，同时保留其主要信息。

空间转换在ColdNAS中起着关键作用，表 有助于我们更好地理解命题在多大程度上帮助减少搜索空间，我们取层数  $L = 4$ ，比率计算为  $\frac{\text{original space}}{\text{transformed space}} = \frac{6^{C \times 4}}{2^{4 \times 4}}$ 。

注意，当  $C = 1$  时，转换不会减少空间。然而，这样的情况并无意义，因为由于调整函数的灵活性不足，其性能较差。转换后的空间通过将架构参数的数量从  $6 \times C \times L$  减少到  $4 \times L$ ，实现了高效和稳健的可微搜索。同时，空间大小从  $6^{C \times L}$  转换为  $2^{4 \times L}$ ，对于任何  $C > 1$ ，这都是一种减少。

### Construction of the Supernet

对于每个  $M^l$ ，由于最多有 4 个操作且它们是对称的，我们只需按任意顺序确定是否进行操作。我们通过引入可微参数来优化权重进行搜索。对于预测器的第  $l$  层，我们有：

$$\hat{h}^{l,k+1} = \alpha^{l,k+1} (\hat{h}^{l,k} \circ_{op^{k+1}} \phi_i^{k+1}) + (1 - \alpha^{l,k+1}) \hat{h}^{l,k},$$

$$\text{衡量 } M^l \text{ 中操作 } o_{op^{k+1}} \text{ 的权重, } k \in \{0, 1, 2, 3\}, \text{ 并且 } \{o_{op^{k+1}}\}_{k=0}^3 = \{\max, \min, \odot, +\}.$$

$$\text{为简化记号, 我们设 } \hat{h}^{l,0} = h^l, h_i^l = \hat{h}^{l,4}.$$

## 知乎

Algorithm 总结了完整的算法。首先，我们优化超网并做出选择以确定结构，然后使用确定的结构重建模型并重新训练以进行推理。

---

**Algorithm 1** Training procedure of ColdNAS.
 

---

**Input:** Learning rate  $\beta$ , number of operations to keep  $K$ .

- 1: Construct the supernet by (8) and randomly initialize all parameters  $\Theta = \{\{\alpha^{l,k}\}_{k=1,l=0}^{4,L-1}, \theta_E, \theta_A, \theta_P\}$ .
- 2: **while** Not converge **do**
- 3:   **for** Every  $T_i \in \mathcal{T}^{\text{train}}$  **do**
- 4:     Calculate  $\Phi_i$  by (1).
- 5:     Calculate  $\hat{y}_{i,j}$  for every  $v_j$  in  $Q_i$  by (2).
- 6:     Calculate loss  $\mathcal{L}_i$  by (10).
- 7:   **end for**
- 8:    $\mathcal{L}^{\text{train}} = \frac{1}{|\mathcal{T}^{\text{train}}|} \sum_{i=1}^{|\mathcal{T}^{\text{train}}|} \mathcal{L}_i$
- 9:   Update all parameters  $\Theta \leftarrow \Theta - \beta \nabla_{\Theta} \mathcal{L}^{\text{train}}$ .
- 10: **end while**
- 11: Determine the modulation structure by keeping operations corresponding to Top- $K$   $\alpha^{l,k}$  and remove the others.
- 12: Construct the model with determined modulation structure and randomly initialize all parameters  $\Theta = \{\theta_E, \theta_A, \theta_P\}$ .
- 13: Train the model in the same way as Step 2 ~ 10.
- 14: **Return:** The trained model.

知乎 @SmartMindAI

学习率  $\beta$ ，保留的操作数  $K$ 。通过构建超网并随机初始化所有参数

$$\Theta = \{\{\alpha^{l,k}\}_{k=1,l=0}^{4,L-1}, \theta_E, \theta_A, \theta_P\}.$$

计算损失  $\mathcal{L}_i$ 。  $\mathcal{L}^{\text{train}} = \frac{1}{|\mathcal{T}^{\text{train}}|} \sum_{i=1}^{|\mathcal{T}^{\text{train}}|} \mathcal{L}_i$  更新所有参数  $\Theta \leftarrow \Theta - \beta \nabla_{\Theta} \mathcal{L}^{\text{train}}$ 。

通过保留对应于Top- $K$   $\alpha^{l,k}$ 的操作并移除其他操作来确定调整结构。构建具有确定调整结构的模型并随机初始化所有参数  $\Theta = \{\theta_E, \theta_A, \theta_P\}$ 。以与步骤 2 ~ 10 相同的方式训练模型。

空间转换带来了巨大的优化，使得双层目标函数\*得以显著降低。传统的可微分架构搜索通过上层变量 $\alpha$ 和下层变量 $\theta$ 来优化双层目标函数。

$$\min_{\alpha} \mathcal{L}^{\text{val}}(\theta^*(\alpha), \alpha), \text{ s.t. } \theta^*(\alpha) = \underset{\theta}{\operatorname{argmin}} \mathcal{L}^{\text{train}}(\theta, \alpha),$$

我们只需要训练超网络来优化训练集上的损失函数\*  $\mathcal{L}^{\text{train}}$ ，并通过端到端的离散训练实现目标。对于训练集中的每个任务  $T_i \in \mathcal{T}^{\text{train}}$ ，我们首先将  $\mathcal{S}_i$  输入到自适应网络  $A$ ，通过等式 (eq:adapara) 生成  $\Phi_i$ 。然后对于每个项目  $v_j \in Q_i$ ，我们通过等式 (eq:forward) 将  $(u_i, v_j)$  和  $\Phi_i$  输入到预测器  $P$  进行预测。最后，我们使用预测值  $\hat{y}_{i,j}$  和真实标签  $y_{i,j}$  之间的均方误差\* (MSE) 作为损失函数。

$$\mathcal{L}_i = \frac{1}{M} \sum_{j=1}^M (y_{i,j} - \hat{y}_{i,j})^2,$$

我们通过梯度下降法\*更新所有参数。当超网收敛后，我们选取  $\{\alpha^{l,k+1}\}_{k=1,l=0}^{4,L-1}$  中的前  $K$  个最大值所对应的操作，联合确定所有的  $M^l$ 。随后，我们重新训练模型以构建最终的用户冷启动模型。在推理阶段，对于一组新的任务  $\mathcal{T}^{\text{test}}$ ，我们将其中的  $T_i$  以及整个  $\mathcal{S}_i$  和  $(u_i, v_j)$  作为输入，获取每个物品  $v_j \in Q_i$  的预测值  $\hat{y}_{i,j}$ 。

## Discussion

架构等价性，如ColdNAS。如果找到一组等价架构，评估其中任何一个架构就足够了。ColdNAS的搜索空间专为冷启动问题中的调整结构设计，我们已证明了原始空间和转换空间的等价性，从而大大减小了空间大小。

Experiments

我们在三个基准数据集上进行了实验，以回答以下问题：

- RQ1: ColdNAS选择了何种调整结构？与当前最先进的冷启动模型相比，ColdNAS的表现如何？
- RQ2: ColdNAS的搜索空间和算法为我们提供了怎样的理解？
- RQ3: 超参数<sup>+</sup>对ColdNAS的影响如何？

Datasets

- MovieLens：收集了包含1百万条用户电影评分的数据库，特征包括性别、年龄、职业、邮政编码<sup>+</sup>、发布年份、评分、类型、导演和演员。
- BookCrossing：BookCrossing社区用户对书籍评分的集合，特征包括年龄、地点、出版年份、作者和出版商。
- Last.fm：Last.fm在线系统中用户对艺术家听力的集合，特征仅包括用户和项目ID。按照@lin2021task的说法，我们在Last.fm的查询集中生成了负样本。

Table 4: Summary of datasets used in this paper.

Dataset	# User (Cold)	# Item	# Rating	# User Feat.	# Item Feat.
MovieLens	6040 (52.3%)	3706	1000209	4	5
BookCrossing	278858 (18.6%)	271379	1149780	2	3
Last.fm	1872 (15.3%)	3846	42346	1	1

数据分割 按照文献1的设置，我们将数据集分为训练集、验证集<sup>+</sup>、测试集<sup>+</sup>，比例为7：1：2。验证集用于判断超级网的收敛。训练集、验证集、测试集没有重叠的用户。

对于MovieLens和Last.fm，我们选择了交互历史长度在40, 200之间的用户。对于BookCrossing，我们将交互历史长度在[50,1000)之间的用户放入训练集，将交互历史长度在[2,50)之间的用户分为70%、10%和20%，分别放入训练集、验证集和测试集。评估指标 我们使用平均绝对误差<sup>+</sup>(MAE)、均方误差(MSE)、归一化折扣累积增益nDCG3和nDCG5评估性能。

Performance Comparison (RQ1)

我们将TheName与其他用户冷启动方法进行比较，包括：

- (i)传统的深度冷启动模型DropoutNet
- (ii) 基于FSL的方法如MeLU, MetaCS, MetaHIN, MAMO和TaNP。

我们还与ColdNAS的一个变体ColdNAS-Fixed进行了比较，该变体在每层都使用公式（2）中的固定FiLM函数，而不是我们搜索到的调整函数。



MovieLens	MSE	100.90 <sub>(0.70)</sub>	95.02 <sub>(0.03)</sub>	95.05 <sub>(0.04)</sub>	91.89 <sub>(0.06)</sub>	90.20 <sub>(0.22)</sub>	89.11 <sub>(0.18)</sub>	91.05 <sub>(0.13)</sub>	87.96 <sub>(0.12)</sub>
	MAE	85.71 <sub>(0.48)</sub>	77.38 <sub>(0.25)</sub>	77.42 <sub>(0.26)</sub>	75.79 <sub>(0.27)</sub>	75.34 <sub>(0.26)</sub>	74.78 <sub>(0.14)</sub>	75.65 <sub>(0.30)</sub>	74.29 <sub>(0.20)</sub>
	nDCG <sub>3</sub>	69.21 <sub>(0.76)</sub>	74.43 <sub>(0.59)</sub>	74.46 <sub>(0.78)</sub>	74.69 <sub>(0.32)</sub>	74.95 <sub>(0.13)</sub>	75.60 <sub>(0.07)</sub>	75.11 <sub>(0.09)</sub>	76.16 <sub>(0.03)</sub>
	nDCG <sub>5</sub>	68.43 <sub>(0.48)</sub>	73.52 <sub>(0.41)</sub>	73.45 <sub>(0.56)</sub>	73.63 <sub>(0.22)</sub>	73.84 <sub>(0.16)</sub>	74.29 <sub>(0.12)</sub>	73.89 <sub>(0.12)</sub>	74.74 <sub>(0.09)</sub>
BookCrossing	MSE	15.38 <sub>(0.23)</sub>	15.15 <sub>(0.02)</sub>	15.20 <sub>(0.08)</sub>	14.76 <sub>(0.07)</sub>	14.82 <sub>(0.05)</sub>	14.75 <sub>(0.05)</sub>	14.44 <sub>(0.16)</sub>	14.15 <sub>(0.08)</sub>
	MAE	3.75 <sub>(0.01)</sub>	3.68 <sub>(0.01)</sub>	3.66 <sub>(0.01)</sub>	3.50 <sub>(0.01)</sub>	3.51 <sub>(0.02)</sub>	3.48 <sub>(0.01)</sub>	3.49 <sub>(0.02)</sub>	3.40 <sub>(0.01)</sub>
	nDCG <sub>3</sub>	77.66 <sub>(0.18)</sub>	77.69 <sub>(0.15)</sub>	77.68 <sub>(0.12)</sub>	77.66 <sub>(0.19)</sub>	77.68 <sub>(0.09)</sub>	77.48 <sub>(0.06)</sub>	77.65 <sub>(0.09)</sub>	77.83 <sub>(0.01)</sub>
	nDCG <sub>5</sub>	80.87 <sub>(0.15)</sub>	81.10 <sub>(0.15)</sub>	80.97 <sub>(0.09)</sub>	80.95 <sub>(0.04)</sub>	81.01 <sub>(0.05)</sub>	81.16 <sub>(0.21)</sub>	81.12 <sub>(0.06)</sub>	81.32 <sub>(0.10)</sub>
Last.fm	MSE	21.91 <sub>(0.38)</sub>	21.69 <sub>(0.34)</sub>	21.68 <sub>(0.12)</sub>	21.43 <sub>(0.23)</sub>	21.64 <sub>(0.10)</sub>	21.58 <sub>(0.20)</sub>	21.62 <sub>(0.16)</sub>	20.91 <sub>(0.05)</sub>
	MAE	43.02 <sub>(0.52)</sub>	42.28 <sub>(1.21)</sub>	42.28 <sub>(0.76)</sub>	42.07 <sub>(0.49)</sub>	42.30 <sub>(0.28)</sub>	42.15 <sub>(0.56)</sub>	42.32 <sub>(0.34)</sub>	41.78 <sub>(0.24)</sub>
	nDCG <sub>3</sub>	75.13 <sub>(0.48)</sub>	80.15 <sub>(2.09)</sub>	80.81 <sub>(0.97)</sub>	82.01 <sub>(0.56)</sub>	80.73 <sub>(0.80)</sub>	81.03 <sub>(0.33)</sub>	80.77 <sub>(0.32)</sub>	82.90 <sub>(0.09)</sub>
	nDCG <sub>5</sub>	69.03 <sub>(0.31)</sub>	75.03 <sub>(0.68)</sub>	75.01 <sub>(0.64)</sub>	75.98 <sub>(0.33)</sub>	75.45 <sub>(0.29)</sub>	75.98 <sub>(0.41)</sub>	75.48 <sub>(0.21)</sub>	76.77 <sub>(0.10)</sub>

结果显示，ColdNAS在所有数据集和指标上都显著优于其他方法。此外，ColdNAS相对于ColdNAS-Fixed的一致性能提升验证了搜索调整结构以适应数据集而不是使用固定结构的必要性。

Table 3: Modulation structure with Top-4 operations searched on the three benchmark datasets respectively. We also show the modulation structure of ColdNAS-Fixed, which is the same regardless of the dataset used.

	$M^0$	$M^1$	$M^2$	$M^3$
MovieLens	$\min(\max(h^0, \phi_i^{0,1}), \phi_i^{0,2}) + \phi_i^{0,3}$	$h^1 + \phi_i^{1,1}$	$h^2$	$h^3$
BookCrossing	$\min(h^0, \phi_i^{0,1})$	$h^1 + \phi_i^{1,1}$	$h^2 \odot \phi_i^{2,1} + \phi_i^{2,2}$	$h^3$
Last.fm	$h^0 \odot \phi_i^{0,1}$	$h^1 + \phi_i^{1,1}$	$\max(h^2, \phi_i^{2,1}) + \phi_i^{2,2}$	$h^3$
ColdNAS-Fixed	$h^0 \odot \phi_i^{0,1} + \phi_i^{0,2}$	$h^1 \odot \phi_i^{1,1} + \phi_i^{1,2}$	$h^2 \odot \phi_i^{2,1} + \phi_i^{2,2}$	$h^3 \odot \phi_i^{3,1} + \phi_i^{3,2}$

Searching in ColdNAS (RQ2)

Choice of Search Strategy

我们在《TheName》中优化了方程（3）并与其他搜索策略进行了比较。这些策略包括在原始空间中进行随机搜索（C=4），在转换空间中进行随机搜索，以及使用双层目标优化超网的"ColdNAS-Bilevel"。我们发现，在转换空间中进行搜索比在C=4时进行随机搜索更有效，而且ColdNAS和ColdNAS-Bilevel的性能都优于随机搜索。ColdNAS和ColdNAS-Bilevel的测试MSE收敛到相似值，但ColdNAS更快，这验证了在ColdNAS中直接对所有参数使用梯度下降的有效性。

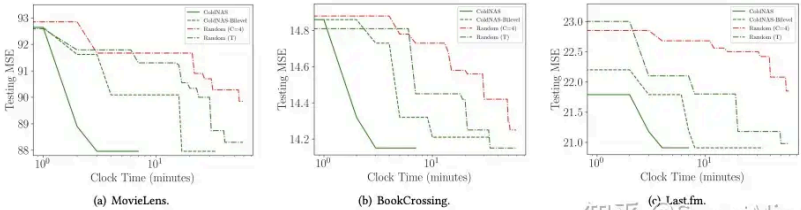


Figure 3: Testing MSE vs clock time of different search strategies. ColdNAS, ColdNAS-Bilevel and Random-t operate on the transformed space, while Random operates on the original space with C = 4 in every  $M^i$ .

Necessity of Search Space Transformation

我们已经选定了搜索算法，接下来我们需要特别考察搜索空间变换在时间和性能方面的必要性。

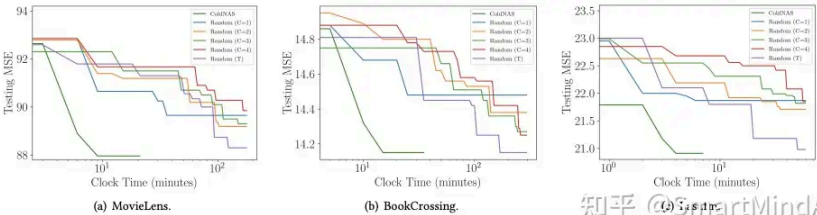


Figure 4: Testing MSE vs clock time on various search spaces: transformed space, and original spaces with different C.

Understanding Proposition

我们首先证明了适应网络A的表现力足够强大。如图所示，尽管A中的层数不同，但搜索到的调整操作相同，且性能差异小。因此，我们选择了2层，既具有足够的表现力，又具有较小的参数大

Sensitivity Analysis (RQ3)

我们还探讨了预测器深度的影响。图绘制了使用具有不同L层数的预测器的效果。可以看出，在一定范围内选择不同的L对性能的影响不大，而选择L=4已经足够好，可以获得如表中所示的最佳结果。

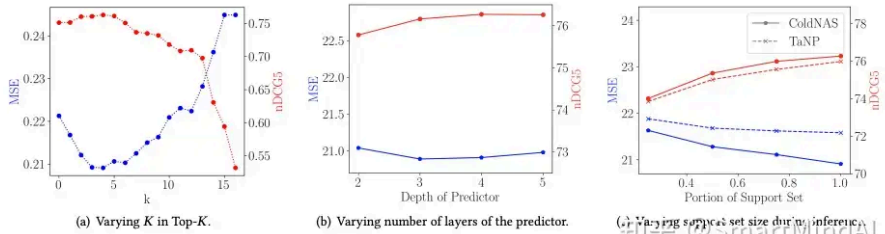


Figure 5: Model sensitivity analysis of ColdNAS on Last.fm.

Conclusion

我们提出了ColdNAS，一个解决用户冷启动推荐问题的调整框架。我们使用超网络将用户历史交互映射到用户特定参数，用于调整预测器。我们设计了一个搜索空间，能够找到更有效的调整结构。理论证明了我们可将搜索空间转换为一个更小的空间，以高效搜索调整结构。实验结果表明，ColdNAS在基准数据集上表现最佳，且易于部署在推荐系统中。

原文《ColdNAS: Search to Modulate for User Cold-Start Recommendation》



冷启动

冷启动关注的是产品早期获取早期核心用户，以及如何运营的问题。

浏览量 367 万  
精华 3,721  
讨论量 3721

关注话题

主后私信

喜欢 收藏 申请转载 ...



还没有评论，发表第一个评论吧

推荐阅读