

Choosing the Best of Both Worlds: Diverse and Novel Recommendations through Multi-Objective Reinforcement Learning

Dušan Stamenković*

University of Novi Sad, Serbia

dusan.stamenkovic@dm.uns.ac.rs

Alexandros

Karatzoglou

Google Research, United

Kingdom

alexkz@google.com

Ioannis Arapakis

Telefonica Research, Spain

ioannis.arapakis@telefonica.com

Xin Xin

Shandong University, China

xinxin@sdu.edu.cn

Kleomenis Katevas

Telefonica Research, Spain

kleomenis.katevas@telefonica.com

摘要

*警告：该PDF由GPT-Academic开源项目调用大语言模型+Latex翻译插件一键生成，版权归原文作者所有。翻译内容可靠性无保障，请仔细鉴别并以原文为准。项目Github地址 https://github.com/binary-husky/gpt_academic/。当前大语言模型: gpt-4o-mini, 当前语言模型温度设定: 0。为了防止大语言模型的意外谬误产生扩散影响，禁止移除或修改此警告。

*Full affiliation is Faculty of Sciences, University of Novi Sad. This work is done when taking an internship in Telefonica Research, Spain, and it was supported by C4IOT.

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 ACM 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.

Woodstock '18, June 03–05, 2018, Woodstock, NY

© 2018 Association for Computing Machinery.

ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

<https://doi.org/10.1145/1122445.1122456>

自推荐系统（Recommender Systems, RS）诞生以来，推荐的相关性准确性一直是评估RS算法质量的黄金标准。然而，专注于物品相关性往往会在其他重要指标上付出重大代价：用户被困在“过滤气泡”中，选择范围大幅度减少，从而降低用户体验质量并导致用户流失。推荐，尤其是基于会话/顺序的推荐，是一个复杂的任务，具有多个且往往是冲突的目标，而现有的最先进方法未能解决这些问题。

在这项工作中，我们挑战上述问题，引入标量化多目标强化学习（*Scalarized Multi-Objective Reinforcement Learning, SMORL*）用于RS设置，这是一种新颖的强化学习（Reinforcement Learning, RL）框架，可以有效解决多目标推荐任务。所提出的SMORL智能体（agent）增强了标准推荐模型，采用额外的RL层，使其同时满足三个主要目标：准确性、多样性和新颖性。我们将该框架与四个最先进的基于会话的推荐模型进行整合，并与一个仅专注于准确性的单一目标RL智能体进行比较。我们在两个真实世界数据集上的实验结果显示，

聚合多样性有显著增加, 准确性适度提高, 推荐的重复性降低, 并证明了增强多样性和新颖性作为互补目标的重要性。

CCS CONCEPTS

• **Information systems** → **Recommender systems**; • **Retrieval models and ranking**; • **Diversity and novelty in information retrieval**;

KEYWORDS

Recommendation; Reinforcement Learning; Multi-Objective Reinforcement Learning; Diversity; Novelty

ACM Reference Format:

Dušan Stamenković, Alexandros Karatzoglou, Ioannis Arapakis, Xin Xin, and Kleomenis Katevas. 2018. Choosing the Best of Both Worlds: Diverse and Novel Recommendations through Multi-Objective Reinforcement Learning. In *Woodstock '18: ACM Symposium on Neural Gaze Detection, June 03–05, 2018, Woodstock, NY*. ACM, New York, NY, USA, 14 pages. <https://doi.org/10.1145/1122445.1122456>

1 INTRODUCTION

无论是在娱乐、社交网络还是电子商务的背景下, 现代网络用户面临的选择数量往往令人不知所措。与普遍认为的“更多的选项总是更好”的看法相反, 从大量选项中所做的选择可能导致选择过载 [17], 并削弱用户理性决策的能力。简单来说, 当用户面临大量选择的情况 (例如, 购买无尽的产品或消费媒体内容) 时, 他们更有可能感觉自己做出了错误的决定并经历后悔, 这会降低他们与在线服务或平台的体验质量。当用户倾向于考虑所有替代选项的成本和收益时, 这一问题变得更加严重。

推荐系统 (RS) 通过作为第二阶策略 [35], 缓解这一选择悖论 [32], 促进对相关信息的访问并改善浏览体验 [16, 43]。因此, 在选项丰富可能导致不满意选择或更糟糕的放弃的情况下, 用户体验最终取决于推荐系统过滤不相关内容并仅推荐被视为理想项的能力。迄今为止, 推荐系统领域的研究重点主要放在设计能够识别和推荐相关内容的算法上。然而, 在这样

做的过程中, 他们往往过度优化 (大部分情况下) 主流指标, 如准确性, 牺牲其他内容导出的定性方面。在这项工作中, 术语“准确性”表示推荐系统在离线测试集中的相关项排名表现, 不应与分类任务中的准确性混淆。

近年来, 推荐的多样性和新颖性被认为是促进用户参与的重要因素, 因为推荐多样的相关项更有可能满足用户的变动需求。例如, Hu and Pu [15] 报告了推荐多样性与可用性、感知有用性和使用意图之间的强正相关。因此, 向刚刚购买了浓缩咖啡机的用户推荐严格相关的物品, 最终很可能会推荐更多的咖啡机, 而理想的推荐集合应包括咖啡杯、清洁设备、咖啡豆, 等等。在前一种情况下, 用户只会与可用物品空间的一个小子空间进行互动 [28], 根据“边际效益递减法则”, 随着用户一次又一次地接触到类似内容, 推荐的效用最终将下降。

基于会话的推荐被引入作为推荐系统的另一种与行业相关的方法。在基于会话的推荐中, 一个顺序模型 (例如, RNN [14] 或 transformer [19, 37]) 以自我监督的方式进行训练, 以预测序列中下一个项, 而不是一些“外部”标签 [14, 19, 43]。该训练过程受到语言建模任务的启发, 其中, 给定一个词序列输入, 语言模型预测下一个最可能出现的词 [24]。然而, 这种训练方法也可能产生次优推荐, 因为损失函数完全由模型预测与序列中的实际项之间的不匹配定义。在这种损失函数下训练的模型仅关注匹配用户可能生成的点击序列, 而放弃其他期望的目标。例如, 服务提供商可能希望促进那些将汇聚到购买、提高用户满意度、多样化用户与物品互动并促进长期参与的推荐。然而, 为了使推荐系统优化以达到上述目标, 需要用可微分函数捕捉这些目标, 而这并不是一件简单的事情。因此, 在重要目标只能以非可微分函数/指标的形式展示的领域, 多目标优化 (MOO) 的使用受到很大限制。

多样性和推荐项目列表的新颖性与销售多样性的增加相关 [9], 并通过推荐所谓“长尾”中的不太流行项目来解决“赢家通吃”的问题。来自多样化推荐列表的一个项目更有可能是新颖的, 即用户通常不会与

之互动的项目。这得到了以往工作的支持，研究表明大多数用户欣赏新颖且不太流行的推荐 [21, 44]。使用简单的监督学习训练的推荐模型可能在满足上述推荐期望和许多在线任务的多目标性质方面遇到困难。

为了解决当前的挑战，我们扩展了在推荐系统 (RS) 环境中利用强化学习 (RL) 的想法，并引入了一种标量化多目标强化学习 (Scalarized Multi-Objective RL, SMORL) 方法。SMORL 使用单个 RL 智能体 (agent) 同时满足三个潜在冲突的目标：i) 促进点击，ii) 多样化推荐集，以及 iii) 引入新颖项目，同时优化相关性。该模型专注于所选择的奖励 (reward)，同时保持高的相关性排名性能。更具体地说，给定一个生成的序列或基于会话的推荐模型，模型的 (最终) 隐藏状态可以视为它的输出层，因为它与最后的 (密集 softmax) 层相乘以生成推荐 [14, 19, 43]。我们为这些模型增添了多个最终输出层。传统的自监督头通过交叉熵损失进行训练以执行排名，而 SMORL 部分则同时训练以修改自监督头的排名。RL 头可以视为正则化器，用于引入更多样化和新颖的推荐，而基于排名的监督头可以为参数更新提供更强的学习信号 (包括负信号)。在推荐系统的背景下，使用多目标强化学习 (MORL) 而不是多目标优化 (MOO) 的主要优势之一是能够使用不可微分的函数作为 RL 智能体用于正则化基础模型的奖励系统。

之前平衡准确性与多样性和新颖性的尝试包括对最终推荐集的重新排名，或训练多个模型并使用遗传算法聚合这些模型 [31]，而我们的方法依赖于训练单个模型，并使用 SMORL 框架来平衡主要推荐目标。我们认为，这个框架可以很容易地推广到其他领域，如音乐、视频和新闻推荐 (通过使用嵌入系统 [3, 23])，在这些领域中，多样性和新颖性是高价值的指标。总之，我们的工作做出了以下贡献：

- 我们设计了一种新颖的多样性奖励 (diversity reward)，该奖励利用了物品嵌入空间 (item embedding space)。
- 我们设计了一种用于评估推荐系统 (RS) 的新指标，该指标衡量推荐的重复性。

- 据我们所知，我们首次在推荐系统 (RS) 的背景下应用多目标强化学习 (MORL)，并探索这一方法所提供的众多可能性和未来研究方向。
- 我们引入了 SMORL，它驱动自监督的推荐系统 (RS) 产生更准确、多样化和新颖的推荐。我们将四个最先进的推荐模型 (model) 集成到所提出的框架中。
- 我们在两个真实的电子商务数据集上进行实验，展示了更少重复的推荐集，在聚合多样性指标上有显著改善 (高达 20%)，同时保持，甚至提升了所有四个状态 (state) -of-the-art (最先进) 模型的准确性。

2 RELATED WORK

多个基于深度学习的方法已被提出，用于有效地建模用户交互序列以适应推荐系统 (RS)。Hidasi et al. [14] 使用门控循环单元 (GRU) [8] 来建模用户会话，而 Tang and Wang [36] 和 Yuan et al. [43] 则使用卷积神经网络 (CNN) 来捕捉序列信号。Kang and McAuley [19] 在序列推荐领域利用了著名的 Transformer [37]，取得了良好的结果。所有这些模型都可以作为基础模型，其输入是用户-项目交互的序列，输出是描述相应用户状态的潜在表征 (representation)。

也有多个尝试将强化学习 (RL) 应用于推荐系统。在 off-policy 设置中，Chen et al. [5] 和 Zhao et al. [45] 提出了使用倾向分数进行 off-policy 校正，但由于高方差导致训练困难。有模型 (model-based) RL 方法 [6, 34, 47] 首先建立一个模型以模拟环境，从而避免任何与 off-policy 训练相关的问题。然而，这两阶段方法在很大程度上依赖于模拟器的准确性。Xin et al. [42] 引入了 SQN 和 SAC，这两个自我监督 (self-supervised) RL 框架为推荐系统增强了推荐模型，增加了两个输出层 (heads)。第一个头是基于交叉熵监督损失，而另一个 RL 头则基于双重 Q 学习 (Double Q-learning) [11]。尽管 SQN 和 SAC 提高了性能，但它们仅通过促进用户可能进行的点击和购买来增加准确性。然而，准确的推荐系统不一定是有效的：真正的价值在于提

示用户可能不会自行发现的项目，即推荐的新颖性和多样性 [12]。提高准确性通常会降低多样性和新颖性，这可能在 RL 被应用于常规化基于会话的推荐系统时出现（见第 5 节中的讨论）。聚合多样性的减少会影响用户体验和对推荐系统的满意度 [15]。Anderson et al. [1] 也报告称，当前的推荐会抑制多样化的用户-项目交互。

最近，多样化推荐和引入新颖推荐被认为是改善推荐系统的重要因素。早期的努力集中在后处理方法上，旨在平衡准确性和多样性 [2, 29, 33]。为了缓解在排名函数上显著累积损失的问题，提出了个性化排名方法 [7]。Chen et al. [4] 尝试解决只考虑对偶测量多样性而忽视项目间相关性的后处理方法的问题，提出了概率模型——行列式点过程（Determinantal Point Process）[22]，该模型使用核矩阵捕捉项目间的相关性。一旦学习了这个矩阵，许多抽样技术便可生成多样的项目集合 [4, 38, 40]。这些模型在最佳情况下实现了准确性和多样性之间的权衡。另一方面，SMORL 显著增加了多样性，并略微提高了准确性。

在 RL 设置中，Zheng et al. [46] 关注于探索-利用策略，以促进多样性，通过随机选择当前推荐项目邻域中的随机项目候选。Hansen et al. [10] 提出了一个基于 RL 的抽样排名器，生成多样项目的排名列表。该模型是一个简单的排名器，模型本身并没有学习生成多样项目集合，而学习过程利用了 REINFORCE 算法 [41]，该算法已知存在高方差问题。最后，以往在推荐系统中优化多个目标的尝试依赖于帕累托优化（Pareto-Optimization），使用网格搜索 [31] 或多梯度下降 [25]。然而，按定义，一个帕累托最优解并不一定比其他帕累托最优解在所有目标上更优。

3 MULTI-OBJECTIVE RL FOR RS

让 \mathcal{I} 表示整个物品集，那么用户-物品交互序列可以表示为 $\mathbf{x}_{1:t} = \{x_1, x_2, \dots, x_{t-1}, x_t\}$ ，其中 $x_i \in \mathcal{I} (0 < i \leq t)$ 是在时间戳 i 上交互的¹ 物品的索引。下一个物品推

荐的目标是向用户推荐最适合其当前兴趣的物品 \mathbf{x}_{t+1} ，给定之前交互序列 $\mathbf{x}_{1:t}$ 。

从多目标强化学习（MORL）的角度来看，下一个物品推荐任务可以被表述为一个多目标马尔可夫决策过程（MOMDP）[39]，其中推荐智能体与环境 \mathcal{E} （用户）通过顺序推荐物品来最大化折扣累计奖励。MOMDP 可以通过元组 $(\mathcal{S}, \mathcal{A}, \mathbf{P}, \mathbf{R}, \rho_0, \gamma)$ 定义，其中：

- \mathcal{S} : 一个连续状态空间，描述用户状态。用户在时间戳 t 的状态可以定义为 $\mathbf{s}_t = G(\mathbf{x}_{1:t}) \in \mathcal{S} (t > 0)$ ，其中 G 是一个将在 Section 4 中讨论的顺序模型（sequential model）。
- \mathcal{A} : 一个包含候选项的离散动作空间。智能体（agent）的动作 a 是推荐所选项。在离线强化学习（offline RL）设置中，我们要么从用户-项目交互中提取时间戳 t 的动作，即 $a_t = x_{t+1}$ ，要么通过将其设置为从自监督层获得的最佳预测来进行设置。状态-动作对 (\mathbf{s}_t, a_t) 的“好坏”由其多目标 Q 价值函数 $Q(\mathbf{s}_t, a_t)$ 来描述。
- $\mathbf{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ 是状态转移概率 $p(\mathbf{s}_{t+1} | \mathbf{s}_t, a_t)$ ，即当智能体（agent）采取动作（action） a_t 时，从 \mathbf{s}_t 到 \mathbf{s}_{t+1} 的状态（state）转移概率。
- $\mathbf{R} : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}^m$ 是向量值奖励函数²，其中 $\mathbf{r}(\mathbf{s}, a)$ 表示在状态 \mathbf{s} 下采取动作 a 所获得的即时奖励。
- ρ_0 是初始状态（state）分布，其中 $\mathbf{s}_0 \sim \rho_0$ 。
- $\gamma \in [0, 1]$ 是未来奖励的折扣因子。对于 $\gamma = 0$ ，智能体（agent）只考虑即时奖励，而对于 $\gamma = 1$ ，所有未来奖励都被完全考虑，除了当前动作（action）的奖励。

目标是寻找 MORL（多目标强化学习）智能体（agent）在 MOMDP（多目标马尔可夫决策过程）中以目标策略 $\pi_\theta(a|\mathbf{s})$ 的形式解决问题，使得根据 $\pi_\theta(a|\mathbf{s})$ 采样的轨迹能够实现最大期望累积奖励（reward）：

$$\max_{\pi_\theta} \mathbb{E}_{\tau \sim \pi_\theta} [f(\mathbf{R}(\tau))], \text{ where } \mathbf{R}(\tau) = \sum_{t=0}^{|\tau|} \gamma^t \mathbf{r}(\mathbf{s}_t, a_t)$$

²每个分量对应一个目标。

¹在现实世界场景中，可能有不同种类的交互。例如，在电子商务中，交互可以是点击、购买、购物车添加等。在音乐推荐中，交互可以通过歌曲的播放时间、歌曲的听歌次数等来表征。

其中 $f: \mathbb{R}^m \mapsto \mathbb{R}$ 是一个标量化函数，而 $\theta \in \mathbb{R}^d$ 表示策略 (policy) 参数。期望是在轨迹 $\tau = (s_0, a_0, s_1, a_1, \dots)$ 上取的，该轨迹是通过根据目标策略执行动作 (action) 获得的： $s_0 \sim \rho_0, a_t \sim \pi_\theta(\cdot | s_t), s_{t+1} \sim P(\cdot | s_t, a_t)$ 。

标量化函数 f 将多目标 Q 价值 (Q-values) $Q(s_t, a_t)$ 和奖励 (reward) 函数 $r(s_t, a_t)$ 映射到一个标量值，即用户效用。在本文中，我们专注于线性 f ；每个目标 i 被赋予一个重要性，即权重 (weight) $w_i, i = 1, \dots, m$ ，使得标量化函数变为 $f_w(x) = w^\top x$ ，其中 $w = [w_1, \dots, w_m]$ 。

4 MODEL AND TRAINING

我们将下一项推荐任务视为一个（自监督）多类分类问题，并构建一个顺序模型，该模型接收用户-项目交互序列 $\mathbf{x}_{1:t} = [x_1, x_2, \dots, x_{t-1}, x_t]$ 作为输入，并生成 n 个分类对数 $\mathbf{y}_{t+1} = [y_1, y_2, \dots, y_n] \in \mathbb{R}^n$ ，其中 n 是候选项目的数量。然后，我们可以从 \mathbf{y}_{t+1} 中选择 top- k 项作为我们在时间戳 $t + 1$ 的推荐列表。每个候选项目对应于一个类别。

通常可以使用生成序列模型 $G(\cdot)$ 将输入序列映射到隐藏状态 $\mathbf{s}_t = G(\mathbf{x}_{1:t})$ 。这作为一个通用编码器函数。基于获得的隐藏状态，我们可以利用一个简单的解码器将隐藏状态映射到分类对数 $\mathbf{y}_{t+1} = d(\mathbf{s}_t)$ 。可以将解码器函数 d 定义为一个简单的全连接层或与候选项目嵌入的内积 [14, 19, 43]。在本工作中，我们使用全连接层。最后，我们通过优化分类对数 \mathbf{y}_{t+1} 基于交叉熵损失 L_s 来训练我们的推荐模型。交叉熵损失的优化将促进正对数的高价值，而用户未互动的项目将被“惩罚”，从而产生强烈的负学习信号。这个负信号对于基础模型的学习至关重要，因为 SMORL 头仅对正动作提供强梯度，即 top-1 项。此外，由于顺序推荐模型 G 已经将输入序列编码为潜在表征 \mathbf{s}_t ，我们直接使用 \mathbf{s}_t 作为 RL 头的当前状态，而无需引入单独的 RL 模型。我们堆叠额外的全连接层来计算位于 G 之上的一维 Q -值：

$$Q_z(s_t, a_t) = \delta(s_t \mathbf{H}_z + b_z) = \delta(G(\mathbf{x}_{1:t}) \mathbf{H}_z + b_z)$$

其中 $z \in \{\text{acc}, \text{div}, \text{nov}\}$ ， δ 表示激活函数，而 \mathbf{H}_z 和 b_z 是 Q-learning 输出层的可学习参数。SMORL 部分随后将计算出的准确性 (accuracy)、多样性 (diversity) 和新颖性 (novelty) Q -值堆叠成一个向量值 Q -值函数：

$$Q(s_t, a_t) = [Q_{\text{acc}}(s_t, a_t), Q_{\text{div}}(s_t, a_t), Q_{\text{nov}}(s_t, a_t)] \quad (1)$$

为了学习向量值 Q 函数并解决 MORL 任务，标量化深度 Q 学习 (SDQL) [27] 扩展了流行的 DQN 算法 [26]，通过引入一个标量化函数 f 。在每个时间步 t ， Q 网络在从经验池 (buffer) D 中获得的经验元组 (s_t, a_t, r_t, s_{t+1}) 的小批量上优化损失 L_{SDQL} ：

$$\begin{aligned} L_{\text{SDQL}} &= (f(y_t^{\text{SDQL}}(s_t, a_t) - \gamma Q(s_{t+1}, a_t)))^2 \\ &= (w^\top (y_t^{\text{SDQL}}(s_t, a_t) - \gamma Q(s_{t+1}, a_t)))^2 \end{aligned} \quad (2)$$

在这里， $y_t^{\text{SDQL}}(s_t, a_t) = r_t + \gamma Q'(s_{t+1}, \arg\max_{a'} [w^\top Q'(s_{t+1}, a')])$ ，其中 Q' 是目标网络。向固定目标网络进行训练防止了状态间近似误差的快速传播，而采样经验进行训练（经验回放）则提高了样本效率，并减少了训练样本间的相关性。

在生成推荐时，我们仍然从监督头中返回前 k 个项目。SMORL 头作为基础推荐模型 G 的正则化器，通过评估推荐的顶部项目质量，根据预定义的奖励 (reward) 设置和标量化函数 f ，即目标的重要性，对其进行微调。

4.1 Reinforcing Accuracy

为了使基础模型 G 学会提供更相关的推荐，我们扩展了 [42] 并将准确性奖励定义为

$$r_{\text{acc}}(s_t, a_t) = r_{\text{acc}}(a_t) = 1, \quad a_t \text{ is a clicked item} \quad (3)$$

根据奖励的定义，当模型 (model) 匹配序列中的下一个点击项目时，将会获得奖励。Xin 等人 [42] 建议同时对点击和购买使用奖励。然而，在本工作中，我们提出了一种可以很容易从电子商务推广到其他相

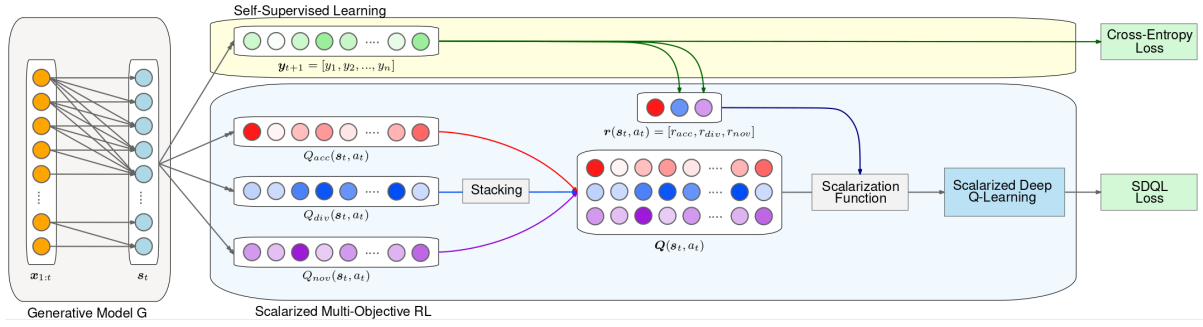


图 1: 推荐系统的 SMORL (SMORL4RS) 训练例程, 适用于序列或基于会话的推荐系统。生成模型 G 将用户-项目交互序列 $x_{1:t}$ 映射到潜在状态 s_t 。通过使用全连接层, s_t 被映射到 **logits** y_{t+1} 以及一维 Q -值: Q_{acc} 、 Q_{div} 和 Q_{nov} 。多样性和新颖性奖励是通过 **logits** 获得的最佳预测计算的。向量值 Q -值函数设定为: $Q = [Q_{acc}, Q_{div}, Q_{nov}]$ 。SDQL 损失通过标量化函数和 SDQL 算法获得, 并与交叉熵损失一起用于训练基模型。

关领域的推荐系统 (RS) 的方法。我们注意到, 通过强化推荐项目的相关性, 可以显著阻碍用户探索平台的能力, 因为推荐与用户最近的兴趣相似。我们在第 5 节中探讨了这一观点。因此, 对于模型来说, 学习如何推荐多样化的项目集合以及那些更可能从未被用户发现的项目是至关重要的。

4.2 Reinforcing Diversity

为了使 SMORL 头部促进多样化的推荐集, 我们首先训练一个 GRU4Rec 模型 [14], 并保存嵌入层 E_{div} 。然后, 我们冻结 E_{div} 的权重以停止参数的进一步更新。我们将奖励 (reward) r_{div} 定义为

$$r_{div} = r_{div}(s_t, p_t) = 1 - \cos(l_t, p_t) = 1 - \frac{e_{l_t}^\top e_{p_t}}{\|e_{l_t}\| \|e_{p_t}\|} \quad (4)$$

在这里, l_t 是会话中最后点击的项目, p_t 是从自监督层获得的最佳预测, 而 e_x 是项目 x 的嵌入, 来自于 E_{div} 。我们不使用目前正在训练的模型的嵌入来计算 r_{div} 。在训练过程的初始阶段, 这会不稳定, 从而产生不可靠的多样性奖励。该奖励加强了推荐的多样性, 而不仅仅是单个列表。我们考虑基于最佳预测 p_t 和前 k 个推荐的多样性奖励系统, 而不仅是最后点击的项目 l_t , 但我们并未观察到性能的改善。

4.3 Reinforcing Novelty

给定一个项目, 用户可能之前在另一组推荐中见过它, 但选择不点击, 或者在某个平台上已经遇到过它。因此, 在实际应用中, 无法跟踪用户可能已经看到的项目, 以及推荐那些肯定是新颖的项目。为了解决这个问题并将新颖性和偶然性引入推荐集, 我们采取了一种概率方法。不太流行的项目更有可能是新颖的, 并导致项目受欢迎程度的更平衡分布。我们使用二元化的项目频率作为我们 MORL (多目标强化学习) 头的奖励 (reward), 我们定义如下:

$$r_{nov} = r_{nov}(s_t, p_t) = \begin{cases} 0.0 & p_t \text{ in top } x\% \text{ of most popular items} \\ 1.0 & \text{otherwise} \end{cases}$$

其中 p_t 是从自监督层获得的顶部预测项 (top predicted item)。 x 的选择基于从训练集推断出的项目流行度的经验分布, 即我们将其设置为长尾开始的近似百分位数。此工作中使用的两个数据集具有相似的分布, 因此我们设置 $x := 10$ 。可以看出, 准确性奖励 (accuracy reward) 依赖于会话中的下一个项目, 而多样性 (diversity) 和新颖性 (novelty) 奖励则依赖于自监督层的顶部预测。

Algorithm 1: Training Procedure of SMORL

Input : user-item interaction sequence set \mathcal{X} ,
recommendation model G ,
SMORL head \mathbf{Q} , supervised head \mathbf{S} ,
predefined parameters α and \mathbf{w}

Output: all parameters in the learning space Θ

```

1 Initialize all trainable parameters
2 Create  $G'$ ,  $\mathbf{Q}'$ , as copies of  $G$  and  $\mathbf{Q}$ , respectively
3 repeat
4   Draw a mini-batch of  $(\mathbf{x}_{1:t}, a_t)$  from  $\mathcal{X}$ 
5    $\mathbf{s}_t = G(\mathbf{x}_{1:t})$ ,  $\mathbf{s}'_t = G'(\mathbf{x}_{1:t})$ 
6    $\mathbf{s}_{t+1} = G(\mathbf{x}_{2:t+1})$ ,  $\mathbf{s}'_{t+1} = G'(\mathbf{x}_{2:t+1})$ 
7   Generate random variable  $z \in (0, 1)$ 
   uniformly
8   if  $z < 0.5$  then
9      $a^* = \operatorname{argmax}_a [\mathbf{Q}(\mathbf{s}_{t+1}, a) \cdot \mathbf{w}]$ 
10     $\text{pred} = \operatorname{argmax} \mathbf{S}(\mathbf{s}_t)$ 
11    Set reward  $\mathbf{r}_t = \text{stack}(r_{\text{acc}}, r_{\text{div}}, r_{\text{nov}})$ 
12     $L_{\text{SDQL}} = \mathbf{w}^\top (\mathbf{r}_t + \gamma \mathbf{Q}'(\mathbf{s}'_{t+1}, a^*) - \mathbf{Q}(\mathbf{s}_t, a_t))^2$ 
13    Calculate  $L_s$ 
14     $L_{\text{SMORL}} = L_s + \alpha L_{\text{SDQL}}$ 
15    Perform updates by  $\nabla_{\Theta} L_{\text{SMORL}}$ 
16  else
17     $a^* = \operatorname{argmax}_a [\mathbf{Q}'(\mathbf{s}'_{t+1}, a) \cdot \mathbf{w}]$ 
18     $\text{pred} = \operatorname{argmax} \mathbf{S}(\mathbf{s}_t)$ 
19    Set reward  $\mathbf{r}_t = \text{stack}(r_{\text{acc}}, r_{\text{div}}, r_{\text{nov}})$ 
20     $L_{\text{SDQL}} =$ 
       $(\mathbf{w}^\top (\mathbf{r}_t + \gamma \mathbf{Q}(\mathbf{s}_{t+1}, a^*) - \mathbf{Q}'(\mathbf{s}'_t, a_t)))^2$ 
21    Calculate  $L_s$ 
22     $L_{\text{SMORL}} = L_s + \alpha L_{\text{SDQL}}$ 
23    Perform updates by  $\nabla_{\Theta} L_{\text{SMORL}}$ 
24  end
25 until converge;
26 Return all parameters in  $\Theta$ 

```

4.4 Scalarized Multi-Objective RL for RS

推荐本质上是一个多目标问题，因此，库存自我监督学习（self-supervised learning），甚至单目标强化学习（RL）方法都不能满足所有理想（或必要）目标。我们将提出的三个目标整合到一个单一的SMORL方法中，在每个时间戳下找到一个最佳动作（动作），该动作根据预定义的用户效用函数，或在此情况下，依据式(2)中 \mathbf{w} 的配置来考虑所有目标。SMORL高度可定制并适应特定提供者的目标——可以定义不同的奖励（奖励）系统，从而导致一个提供更相关、新颖、多样、意外或偶然推荐的RS。我们优化的最终损失是：

$$L_{\text{SMORL}} = L_s + \alpha L_{\text{SDQL}} \quad (5)$$

在这里， L_s 是交叉熵损失，而 α 是一个超参数，使我们能够控制 SMORL 部分的影响。为了增强学习的稳定性，我们交替训练两个可学习参数的副本。算法 1 描述了 SMORL 的训练过程。值得注意的是，在训练完成后，仅使用基础模型的自监督部分来生成推荐，同时通过 SMORL 部分观察与不同指标相关的效果。

这一训练框架可以集成到现有的推荐模型中，只要它们遵循之前讨论的通用架构。这适用于近年来引入的大多数基于会话或顺序的推荐模型。在这项工作中，我们使用交叉熵损失用于自监督部分，但其他模型可以结合不同的损失函数 [13, 30]。

此外，SMORL 是一个高度模块化的框架，用户可以重新加权和“停用”特定的强化学习 (RL) 目标，或在精心设计的奖励 (reward) 方案的帮助下添加更多的目标。最终，这一机制使推荐系统 (RS) 能够专注于提供者的特定短期和长期目标。然而，我们的实验结果表明，在大多数情况下，由三种 RL 目标正则化的模型在所有质量指标上表现最佳。

5 EXPERIMENTS

我们报告了我们实验的结果³ 在两个真实世界的连续电子商务数据集上。对于所有基础模型，我们使用自监督头生成推荐。我们解决以下研究问题：

RQ1: 当集成时，所提出的方法是否提高了基础模型的性能？

RQ2: 我们能否控制准确性、多样性和新颖性之间的平衡？

RQ3: 我们能否通过调整SMORL部分的梯度强度来增加其影响力？

5.1 Experimental Settings

5.1.1 *Datasets:* RC15⁴ 和 RetailRocket⁵，表 1。

RC15. 此数据集基于 2015 年 RecSys 挑战 (RecSys Challenge 2015)。该数据集是基于会话的，每个会话包含一系列的点击和购买⁶。我们丢弃长度小于 3 的会话，然后抽样 200K 个会话的子集。

RetailRocket. 此数据集是从一个真实的电子商务网站收集的。它包含查看和加入购物车的会话事件。为了与 RC15 数据集保持一致，我们将查看视为点击。我们删除互动次数少于三次 (3) 的项目，以及长度小于三 (3) 的序列。

5.1.2 *Quality of Recommendation Metrics.*

准确性度量。 推荐项目集的相关性通常通过两个度量来评估：命中率 (HR) 和归一化折扣累积增益 (NDCG)。HR@ k 是一种基于召回的度量，衡量真实项是否出现在推荐列表的前 k 个位置。我们将点击的 HR 定义为：

$$\text{HR}(\text{click}) = \frac{\text{\#hits among clicks}}{\text{\#clicks}}$$

另一方面，NDCG 是一种对排名敏感的度量，它对推荐列表中的顶级位置赋予更高的分数[18]

³实 现 可 以 在 <https://drive.google.com/file/d/1lVeKlajOkZ4n9Rl2VmJvYR9i1aXWkR2j/view?usp=sharing> 找到

⁴<https://recsys.acm.org/recsys15/challenge/>

⁵<https://www.kaggle.com/retailrocket/ecommerce-dataset>

⁶在本研究中,我们仅考虑点击。

表 1: Dataset statistics.

Dataset	RC15	RetailRocket
#sequences	200,000	195,523
#items	26,702	70,852
#clicks	1,110,965	1,176,680
#purchase	43,946	57,269

多样性与新颖性指标。 推荐系统中的多样性可以在个体或整体层面上进行查看。例如，如果推荐系统向所有用户提供相同的一组十个不同的项目，则每个用户的推荐列表会是多样的，即它具有较高的个体多样性。然而，系统只能从整个项目池中推荐十个项目，因此整体多样性可能微不足道。因此，在我们的实验中，我们使用 Coverage@ k (CV@ k) 来测量整体多样性， $k \in \{1, 5, 10, 20\}$ 。更具体地说，我们在两个集合上测量 CV@ k ：所有项目的集合和不太受欢迎项目的集合。覆盖率可以计算为所有顶级- k 推荐的验证或测试序列所覆盖的所有项目（不太受欢迎项目）的百分比。

推荐的重复性。 我们引入重复性 (Repetitiveness, R)，这是一种评估推荐有用性的新的度量。我们认为该指标是判断推荐系统如何易于创建过滤泡沫的良好替代，因为它测量了推荐列表中顶级- k 位置的每个会话平均重复次数。我们测量 R@ k ， $k \in \{5, 10, 20\}$ ，并将其定义为：

$$R@K = \frac{1}{N} \sum_{i=1}^N \text{\#repetitions in top-}k \text{ items of session } i \quad (6)$$

其中 N 是测试（或验证）集中的会话总数。

5.1.3 *Evaluation Protocols.* 我们使用 5 折交叉验证 (5-fold cross-validation) 进行性能评估，训练、验证和测试的比例为 8:1:1。我们报告所有折叠的平均性能。

5.1.4 *Baselines.* 我们在四个最先进的 (state-of-the-art) 生成 (generative) 序列推荐模型中集成了 SMORL：

- GRU4Rec [14]: 该方法使用 GRU（门控递归单元）对输入序列进行建模。GRU4Rec 的最终隐

藏状态被视为输入序列的潜在表征 (latent representation)。

- Caser [36]: 这种最近引入的基于卷积神经网络 (CNN, Convolutional Neural Network) 的方法通过对前一个项目的嵌入矩阵 (embedding matrix) 应用卷积操作来捕捉序列信号。
- NextItNet [43]: 该方法通过使用膨胀卷积神经网络 (dilated CNN) 来扩大感受野, 并使用残差连接 (residual connection) 来增加网络深度, 从而增强了Caser。
- SASRec [19]: 该基准模型受到自注意力 (self-attention) 的启发, 并使用transformer [37]架构对用户-物品交互序列进行编码。transformer编码器的输出被视为潜在表征 (latent representation)。

5.1.5 Parameter settings. 对于这两个数据集, 输入序列由目标时间戳之前的最后 10 个项目组成。如果序列长度小于 10, 我们用一个填充项进行补充。我们使用 Adam 优化器 (optimizer) 训练所有模型 [20]。小批量 (mini-batch) 大小设定为 256。学习率对于 RC15 设置为 0.01, 对于 RetailRocket 设置为 0.005。我们在 RC15 的每 5,000 批更新和 RetailRocket 的每 10,000 批更新上评估验证集。为了确保公平比较, 所有模型的项目嵌入 (embedding) 大小均设置为 64。对于 GRU4Rec 模型, 隐状态 (hidden state) 的大小设置为 64。对于 Caser, 我们使用一个垂直卷积 (vertical convolution) 滤波器和 16 个水平 (horizontal) 滤波器, 其高度设置为 {2, 3, 4}。丢弃率 (drop-out ratio) 设置为 0.1。对于 NextItNet, 我们使用作者报告的相同参数。对于 SASRec, 自注意力 (self-attention) 中的头数设置为 1, 根据其原始论文 [19]。我们将折扣因子 (discount factor) γ 设置为 0.5, 如 Xin et al. [42] 推荐的那样。

5.2 Performance Comparison (RQ1)

对于这两个数据集, SQN 方法 [42] 在向用户推荐相关项目方面优于基线模型。然而, 通过提高基线模型的准确性, 它导致模型在多样性 (diversity) 和新颖性 (novelty) 上出现“漂移”。这导致基线模型的覆盖率指标显著下降 (最高可达 20%), 无论是对所有项目还是对不太受欢迎的项目。结合这一事实, 推荐的重复性增加表明, 仅仅强化准确性可能明显影响感知的体验质量。此外, 显然应该同时优化模型以实现多样性和新颖性, 从而在相对指标之间取得平衡。在表 2 和表 3 中, 我们看到通过使用 SMORL 方法, 我们不仅在准确性、多样性和新颖性之间取得平衡, 而且在所有指标上持续优于相应的基线模型, 并在某种程度上提高了它们的准确性。与基线模型相比, 多样性和新颖性的提高高达 20%, 而与 SQN 模型相比, 提升高达 40%。基线模型准确性的提高可以归因于大多数用户的兴趣多样, 推荐系统 (RS) 产生的推荐无法满足这些兴趣 [1]。图 2 显示了在 RC15 测试集上获得的累计多样性和新颖性奖励的差异。当使用 SMORL 框架训练基础模型时, 我们注意到累计多样性和新颖性奖励显著增加。此外, 表 2 和表 3 中的结果表明, 强化多样性和新颖性在这些指标上引入了显著的改善, 这些指标与感知的体验质量和参与度高度相关。

5.3 Reinforcing a Subset of Objectives (RQ2)

使用SMORL的一个优点是其目标平衡能力, 它通过在公式(2)中使用不同配置的 \mathbf{w} 对目标进行重新加权。在我们的设置中, \mathbf{w} 的第一个条目对应于准确性目标的权重, 第二个对应于多样性, 第三个对应于新颖性目标。我们使用以下参数 \mathbf{w} 的配置进行实验:

$$\mathbf{w} \in \{(0, 1, 0), (0, 0, 1), (0, 1, 1), (1, 1, 0), (1, 0, 1)\} \quad (7)$$

在这里, 我们旨在展示在强化三项重要目标的子集时性能的差异。我们在此分析中不包括 $\mathbf{w} = (1, 0, 0)$,

表 2: 在RC15数据集上的推荐性能。NG是NDCG。CV是覆盖率。粗体字表示最高得分。

Models	accuracy				diversity				novelty				repetitiveness		
	HR@10	NG@10	HR@20	NG@20	CV@1	CV@5	CV@10	CV@20	CV@1	CV@5	CV@10	CV@20	R@5	R@10	R@20
GRU	0.3793	0.2279	0.4581	0.2478	0.2481	0.4330	0.5188	0.5942	0.1777	0.3707	0.4654	0.5492	12.11	25.63	53.67
GRU-SQN	0.3946	0.2394	0.4741	0.2587	0.2406	0.4025	0.4710	0.5364	0.1656	0.3363	0.4122	0.4849	12.20	25.81	53.67
GRU-SMORL	0.4007	0.2433	0.4793	0.2632	0.2825	0.4758	0.5577	0.6334	0.2086	0.4176	0.5086	0.5927	11.29	23.81	48.67
Caser	0.3593	0.2177	0.4371	0.2372	0.2631	0.4349	0.5019	0.5608	0.1912	0.3724	0.4466	0.5120	14.38	29.65	60.33
Caser-SQN	0.3668	0.2223	0.4448	0.2420	0.2154	0.3525	0.4057	0.4557	0.1411	0.2810	0.2154	0.3953	14.45	29.79	60.33
Caser-SMORL	0.3664	0.2224	0.4425	0.2417	0.3174	0.5157	0.5944	0.6685	0.2476	0.4621	0.5495	0.6316	13.77	28.56	58.67
NtItNet	0.3885	0.2332	0.4684	0.2535	0.2950	0.4914	0.5705	0.6427	0.2313	0.4354	0.5228	0.6030	10.03	22.02	46.67
NtItNet-SQN	0.4083	0.2492	0.4878	0.2693	0.2737	0.4572	0.5183	0.5715	0.2082	0.3975	0.4649	0.5239	10.19	22.32	47.67
NtItNet-SMORL	0.4116	0.2505	0.4898	0.2703	0.3385	0.5639	0.6518	0.7283	0.2720	0.5156	0.6131	0.6981	9.97	21.73	45.67
SASRec	0.4257	0.2599	0.5053	0.2801	0.2971	0.5208	0.6046	0.6792	0.2298	0.4679	0.5607	0.6436	10.62	23.24	49.67
SASRec-SQN	0.4288	0.2630	0.5073	0.2829	0.2701	0.4527	0.5194	0.5755	0.2018	0.3922	0.4660	0.5283	10.94	23.85	50.67
SASRec-SMORL	0.4315	0.2651	0.5104	0.2851	0.3380	0.5755	0.6508	0.7158	0.2698	0.5285	0.6120	0.6842	10.38	22.79	48.67

表 3: 在RetailRocket数据集上的推荐性能。NG是NDCG。CV是覆盖率（Coverage）。粗体字表示最高分。

Models	accuracy				diversity				novelty				repetitiveness		
	HR@10	NG@10	HR@20	NG@20	CV@1	CV@5	CV@10	CV@20	CV@1	CV@5	CV@10	CV@20	R@5	R@10	R@20
GRU	0.2673	0.1878	0.3082	0.1981	0.2439	0.4695	0.5699	0.6632	0.1837	0.4139	0.5238	0.6267	14.25	29.44	60.33
GRU-SQN	0.2967	0.2094	0.3406	0.2205	0.2180	0.4114	0.4975	0.5763	0.1526	0.3489	0.4430	0.5299	14.62	30.19	60.33
GRU-SMORL	0.3060	0.2103	0.3535	0.2224	0.2796	0.5369	0.6419	0.7353	0.2154	0.4871	0.6029	0.7064	13.53	28.02	57.67
Caser	0.2302	0.1675	0.2628	0.1758	0.2327	0.4379	0.5133	0.5718	0.1643	0.3773	0.4605	0.5252	16.16	33.24	60.33
Caser-SQN	0.2454	0.1778	0.2803	0.1867	0.2088	0.3880	0.4511	0.5021	0.1387	0.3219	0.3914	0.4479	16.88	34.50	70.33
Caser-SMORL	0.2657	0.1898	0.3052	0.1998	0.2855	0.5411	0.6324	0.7138	0.2224	0.4917	0.5925	0.6827	15.90	32.47	60.33
NtItNet	0.3007	0.2060	0.3506	0.2186	0.2867	0.5113	0.6033	0.6837	0.2305	0.4595	0.5605	0.6495	12.25	25.76	50.67
NtItNet-SQN	0.3129	0.2150	0.3586	0.2266	0.2802	0.5255	0.6077	0.6750	0.2184	0.4747	0.5651	0.6395	12.27	25.93	50.67
NtItNet-SMORL	0.3183	0.2222	0.3659	0.2342	0.3429	0.6335	0.7351	0.8129	0.2800	0.5938	0.7062	0.7924	10.92	22.89	47.67
SASRec	0.3085	0.2107	0.3572	0.2227	0.2767	0.5305	0.6300	0.7149	0.2171	0.4806	0.5899	0.6838	15.67	32.27	60.33
SASRec-SQN	0.3302	0.2279	0.3803	0.2406	0.2393	0.4617	0.5490	0.6254	0.1753	0.4040	0.5001	0.5847	15.60	32.20	60.33
SASRec-SMORL	0.3521	0.2477	0.4028	0.2605	0.3037	0.5724	0.6672	0.7476	0.2366	0.5261	0.6311	0.7202	12.58	26.69	50.67

因为 SMORL 与 [42] 中的 SQN 方法等价，我们的结果在所有模型中表现出完全相同的行为。

我们在本工作中涉及的目标之间存在复杂的关系。例如，训练过程开始时的相关性和多样性是相关的，

即，更多样化的推荐会产生更相关的推荐，而随着训练的进展，它们之间的相关性变为负数。多样性和新颖性是交织在一起的目标，例如，一组多样化的推荐项目更可能包含新颖的项目。另一方面，项目的受欢迎

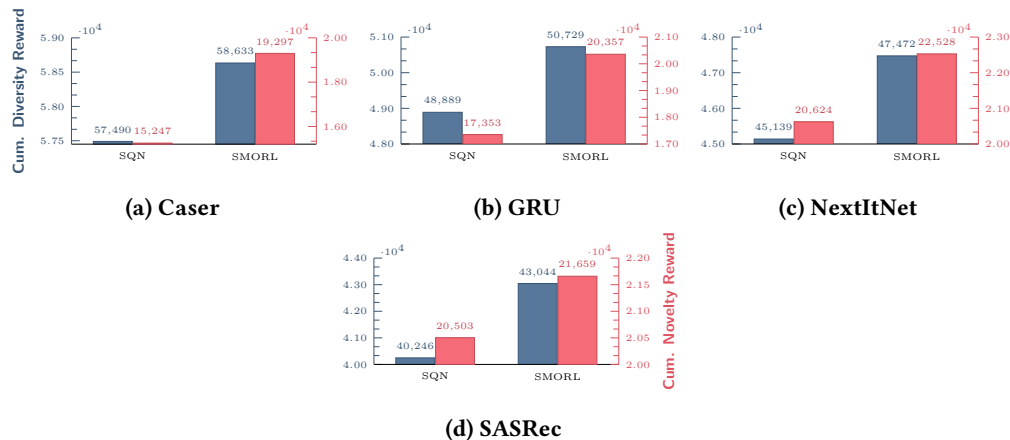


图 2: 在SMORL和SQN框架下, 对基模型在RC15数据集上的累计奖励 (cumulative rewards) 进行比较。

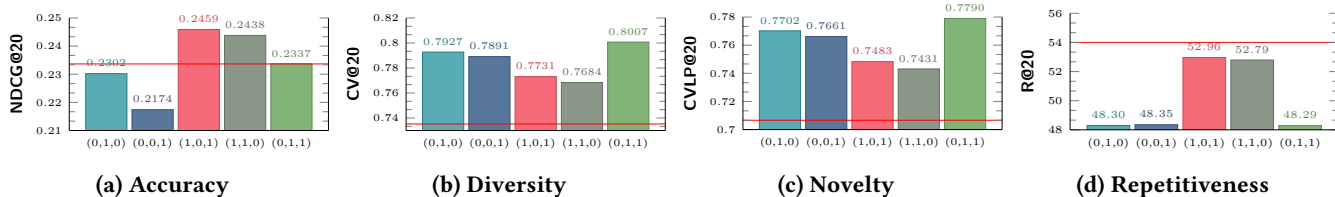


图 3: 当增强目标子集的性能比较 (通过使用来自 Eq.(7) 的不同配置的 w 实现)。红线表示股票 NextItNet 模型的相关指标 - 未使用 SMORL4RS 框架进行训练。

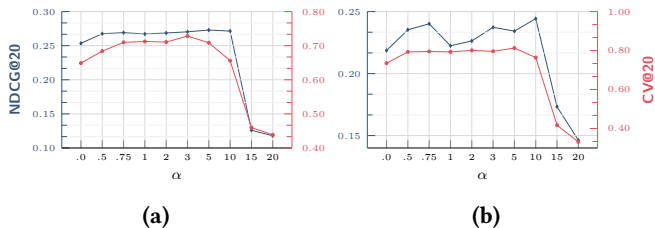


图 4: 在RC15 (4a) 和RetailRocket (4b) 数据集上具有不同SMORL梯度强度的NextItNet

受欢迎程度遵循一个幂律分布, 因此, 受欢迎程度较低的项目占数据集的 90%, 这意味着很可能是新颖的项目本质上是多样的。考虑到所提出的方法不是一个纯粹的 MORL 模型, 而是强制基模型捕捉不同 (并且通常是相互竞争的) 目标的正则化器, 优化和平衡多个目标的复杂性构成了一个重大的研究挑战。在本节中, 我们的目标是展示我们可以控制每个目标的影响程度, 而不是如何找到理想的平衡。通过这种控制的能力, 出现了许多工程可能性, 例如部署多个 SMORL4RS 智能体, 并在线决定用户是否应从一个旨在新颖性、多样性或准确性优化的智能体中接收推荐。

图 3 显示了由 SMORL 智能体 (agent) 使用上述权重配置 w 对 RetailRocket 数据集进行正规化的 NextItNet-SMORL 模型, 而对于 RC15 数据集和其他模型也可以观察到类似的行为。更具体地说, 图 3a 指出, 如果我们仅针对新颖性对模型进行正规化, 我们将牺牲其推荐相关项目的的能力。如果我们仅强化多样性, 这种现象也是存在的, 但 NDCG@20 指标的下降并不明显。另一方面, 如果我们共同优化多样性和新颖性, 我们并未观察到基础模型的准确性下降。此外, 如果我们将准确性目标包含在其他两个目标之中, 我们观察到相关指标的增加。从图 3b 和 3c 中, 我们注意到包含准确性目标会以牺牲多样性和新颖性为代价, 而共同优化多样性和新颖性则在这些指标方面产生最佳结果。类似地, 通过包含准确性目标, 我们与优化多样性和新颖性组合的 NextItNet-SMORL 模型相比, 增加了重复性。

5.4 Gradient Intensity Investigation (RQ3)

在所有基本模型和两个数据集上，SDQL损失由自监督损失主导，这表明从公式(5)中优化参数 α 可能会改善SMORL部分对基本模型的影响。图4显示了NextItNet-SMORL模型在两个数据集上关于NDCG@20和CV@20指标的表现，当我们改变SDQL梯度的强度时。如预期，按 $\alpha < 1$ 乘以SDQL时，效果降低，与基本模型相比我们没有显著改善。对于 $\alpha \in \{1, 2, 3, 5, 10\}$ ，可以看到两个指标的增加，当 $\alpha = 5$ 时获得了最佳平衡。对于更高的 α 值，由于从自监督损失获得的梯度信号丧失，我们观察到质量显著下降，这表明需要有自监督部分来学习基本的排名。对于RC15数据集，也可以进行类似的分析。

在大多数情况下， α 参数的最佳价值等于1 - 在RC15数据集上的SASRec，RetailRocket上的GRU4Rec，以及RetailRocket数据集上的Caser。然而，对于RC15上的GRU4Rec和Caser，最佳值等于0.75，对于RC15上的SASRec和NextItNet为3，而RetailRocket上的SASRec等于10。因此，对于真实世界的使用案例，数据集通常包含数百万项时，更高的 α 值可能是最佳的。更复杂的模型，如NextItNet和SASRec需要更高的 α 值。

此外，监督损失与SDQL损失的联合优化本身也是一个研究课题。最后，我们计划探索使用非线性和个性化的标量化函数。

6 CONCLUSIONS & FUTURE WORK

我们首先正式化了下一个项目推荐任务，并将其作为多目标MDP任务呈现。SMORL方法作为正则化器，用于将期望的特性引入推荐模型，特别是为了实现推荐的相关性、多样性和新颖性之间的平衡。我们将SMORL与四个最先进的推荐模型集成，并在两个真实的电子商务数据集上进行了实验。我们的实验发现，三项相互冲突目标的联合优化对于提高与用户满意度高度相关的指标至关重要，同时还保持内容的相关性。未来的工作为探索SMORL范式在推荐系统（RS）环境中的应用带来了广阔的可能性，并将包括与不同目标的进一步实验以及在不同领域（如音乐平台）中应用SMORL。

REFERENCES

- [1] Ashton Anderson, Lucas Maystre, Ian Anderson, Rishabh Mehrotra, and Mounia Lalmas. 2020. Algorithmic effects on the diversity of consumption on spotify. In *Proceedings of The Web Conference 2020*. 2155–2165.
- [2] Azin Ashkan, Branislav Kveton, Shlomo Berkovsky, and Zheng Wen. 2015. Optimal Greedy Diversity for Recommendation.. In *IJCAI*, Vol. 15. 1742–1748.
- [3] Aishwariya Budhrani, Akashkumar Patel, and Shivam Ribadiya. 2020. Music2Vec: Music Genre Classification and Recommendation System. In *2020 4th International Conference on Electronics, Communication and Aerospace Technology (ICECA)*. IEEE, 1406–1411.
- [4] Laming Chen, Guoxin Zhang, and Hanning Zhou. 2017. Fast greedy map inference for determinantal point process to improve recommendation diversity. *arXiv preprint arXiv:1709.05135* (2017).
- [5] Minmin Chen, Alex Beutel, Paul Covington, Sagar Jain, Francois Belletti, and Ed H Chi. 2019. Top-k off-policy correction for a REINFORCE recommender system. In *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining*. 456–464.
- [6] Xinshi Chen, Shuang Li, Hui Li, Shaohua Jiang, Yuan Qi, and Le Song. 2019. Generative adversarial user model for reinforcement learning based recommendation system. In *International Conference on Machine Learning*. PMLR, 1052–1061.
- [7] Peizhe Cheng, Shuaiqiang Wang, Jun Ma, Jiankai Sun, and Hui Xiong. 2017. Learning to recommend accurate and diverse items. In *Proceedings of the 26th international conference on World Wide Web*. 183–192.
- [8] Kyunghyun Cho, Bart Van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. 2014. On the properties of neural machine translation: Encoder-decoder approaches. *arXiv preprint arXiv:1409.1259* (2014).
- [9] Daniel M Fleder and Kartik Hosanagar. 2007. Recommender systems and their impact on sales diversity. In *Proceedings of the 8th ACM conference on Electronic commerce*. 192–199.
- [10] Christian Hansen, Rishabh Mehrotra, Casper Hansen, Brian Brost, Lucas Maystre, and Mounia Lalmas. 2021. Shifting Consumption towards Diverse Content on Music Streaming Platforms. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*. 238–246.
- [11] Hado V Hasselt. 2010. Double Q-learning. In *Advances in neural information processing systems*. 2613–2621.
- [12] Jonathan L Herlocker, Joseph A Konstan, Loren G Terveen, and John T Riedl. 2004. Evaluating collaborative filtering recommender systems. *ACM Transactions on Information Systems (TOIS)* 22, 1 (2004), 5–53.
- [13] Balázs Hidasi and Alexandros Karatzoglou. 2018. Recurrent neural networks with top-k gains for session-based recommendations. In *Proceedings of the 27th ACM international conference on information and knowledge management*. 843–852.
- [14] Balázs Hidasi, Alexandros Karatzoglou, Linas Baltrunas, and Domonkos Tikk. 2015. Session-based recommendations with recurrent neural networks. *arXiv preprint arXiv:1511.06939* (2015).
- [15] Rong Hu and Pearl Pu. 2011. Helping Users Perceive Recommendation Diversity.
- [16] Yujing Hu, Qing Da, Anxiang Zeng, Yang Yu, and Yinghui Xu. 2018. Reinforcement learning to rank in e-commerce search engine: Formalization, analysis, and application. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 368–377.
- [17] Sheena S Iyengar and Mark R Lepper. 1999. Rethinking the value of choice: a cultural perspective on intrinsic motivation. *Journal of personality and social psychology* 76, 3 (1999), 349.
- [18] Kalervo Järvelin and Jaana Kekäläinen. 2002. Cumulated gain-based evaluation of IR techniques. *ACM Transactions on Information Systems (TOIS)* 20, 4 (2002), 422–446.
- [19] Wang-Cheng Kang and Julian McAuley. 2018. Self-attentive sequential recommendation. In *2018 IEEE International Conference on Data Mining (ICDM)*. IEEE, 197–206.
- [20] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [21] Neal Lathia, Stephen Hailes, Licia Capra, and Xavier Amatriain. 2010. Temporal diversity in recommender systems. In *Proceedings of the 33rd international ACM SIGIR conference on Research and development in information retrieval*. 210–217.
- [22] Frédéric Lavancier, Jesper Møller, and Ege Rubak. 2015. Determinantal point process models and statistical inference. *Journal of the Royal Statistical Society: Series B: Statistical Methodology* (2015), 853–877.
- [23] Ye Ma, Lu Zong, Yikang Yang, and Jionglong Su. 2019. News2vec: News network embedding with subnode information. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. 4845–4854.
- [24] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781* (2013).
- [25] Nikola Milojkovic, Diego Antognini, Giancarlo Bergamin, Boi Faltings, and Claudiu Musat. 2019. Multi-gradient descent for multi-objective recommender systems. *arXiv preprint arXiv:2001.00846* (2019).
- [26] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. 2013. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602* (2013).
- [27] Hossam Mossalam, Yannis M Assael, Diederik M Roijers, and Shimon Whiteson. 2016. Multi-objective deep reinforcement learning. *arXiv preprint arXiv:1610.02707* (2016).
- [28] Eli Pariser. 2011. *The filter bubble: What the Internet is hiding from you*. Penguin UK.
- [29] Lijing Qin and Xiaoyan Zhu. 2013. Promoting diversity in recommendation by entropy regularizer. In *Twenty-Third International Joint Conference on Artificial Intelligence*. Citeseer.
- [30] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2012. BPR: Bayesian personalized ranking from implicit feedback. *arXiv preprint arXiv:1205.2618* (2012).
- [31] Marco Tulio Ribeiro, Anisio Lacerda, Adriano Veloso, and Nivio Ziviani. 2012. Pareto-efficient hybridization for multi-objective recommender systems. In *Proceedings of the sixth ACM conference on Recommender systems*. 19–26.
- [32] Barry Schwartz. 2004. The paradox of choice: Why less is more. *New York: Ecco* (2004).
- [33] Chaofeng Sha, Xiaowei Wu, and Junyu Niu. 2016. A Framework for Recommending Relevant and Diverse Items.. In *IJCAI*, Vol. 16. 3868–3874.

- [34] Wenjie Shang, Yang Yu, Qingyang Li, Zhiwei Qin, Yiping Meng, and Jieping Ye. 2019. Environment reconstruction with hidden confounders for reinforcement learning based recommendation. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 566–576.
- [35] Cass R Sunstein and Edna Ullmann-Margalit. 1999. Second-order decisions. *Ethics* 110, 1 (1999), 5–31.
- [36] Jiayi Tang and Ke Wang. 2018. Personalized top-n sequential recommendation via convolutional sequence embedding. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining*. 565–573.
- [37] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*. 5998–6008.
- [38] Romain Warlop, Jérémie Mary, and Mike Gartrell. 2019. Tensorized determinantal point processes for recommendation. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 1605–1615.
- [39] C Ch White, CC III WHITE, and KIM KW. 1980. Solution procedures for vector criterion Markov decision processes. (1980).
- [40] Mark Wilhelm, Ajith Ramanathan, Alexander Bonomo, Sagar Jain, Ed H Chi, and Jennifer Gillenwater. 2018. Practical diversified recommendations on youtube with determinantal point processes. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*. 2165–2173.
- [41] Ronald J Williams. 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning* 8, 3-4 (1992), 229–256.
- [42] Xin Xin, Alexandros Karatzoglou, Ioannis Arapakis, and Joemon M Jose. 2020. Self-Supervised Reinforcement Learning for Recommender Systems. *arXiv preprint arXiv:2006.05779* (2020).
- [43] Fajie Yuan, Alexandros Karatzoglou, Ioannis Arapakis, Joemon M Jose, and Xiangnan He. 2019. A simple convolutional generative network for next item recommendation. In *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining*. 582–590.
- [44] Yuan Cao Zhang, Diarmuid Ó Séaghdha, Daniele Quercia, and Tamas Jambor. 2012. Auralist: introducing serendipity into music recommendation. In *Proceedings of the fifth ACM international conference on Web search and data mining*. 13–22.
- [45] Xiangyu Zhao, Liang Zhang, Zhuoye Ding, Long Xia, Jiliang Tang, and Dawei Yin. 2018. Recommendations with negative feedback via pairwise deep reinforcement learning. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 1040–1048.
- [46] Guanjie Zheng, Fuzheng Zhang, Zihan Zheng, Yang Xiang, Nicholas Jing Yuan, Xing Xie, and Zhenhui Li. 2018. DRN: A deep reinforcement learning framework for news recommendation. In *Proceedings of the 2018 World Wide Web Conference*. 167–176.
- [47] Lixin Zou, Long Xia, Zhuoye Ding, Jiaxing Song, Weidong Liu, and Dawei Yin. 2019. Reinforcement learning to optimize long-term user engagement in recommender systems. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2810–2818.