# How to Diversify any Personalized Recommender? A User-centric Pre-processing approach

Manel Slokom
Laura Hollink
manel.slokom@cwi.nl
laura.hollink@cwi.nl
Centrum Wiskunde & Informatica
The Netherlands

## ABSTRACT

In this paper, we introduce a novel approach to improve the diversity of Top-N recommendations while maintaining recommendation performance. Our approach employs a user-centric pre-processing strategy aimed at exposing users to a wide array of content categories and topics. We personalize this strategy by selectively adding and removing a percentage of interactions from user profiles. This personalization ensures we remain closely aligned with user preferences while gradually introducing distribution shifts.

Our pre-processing technique offers flexibility and can seamlessly integrate into any recommender architecture. To evaluate our approach, we run extensive experiments on two publicly available data sets for news and book recommendations. We test various standard and neural network-based recommender system algorithms. Our results show that our approach generates diverse recommendations, ensuring users are exposed to a wider range of items. Furthermore, leveraging pre-processed data for training leads to recommender systems achieving performance levels comparable to, and in some cases, better than those trained on original, unmodified data. Additionally, our approach promotes provider fairness by facilitating exposure to minority or niche categories.

## KEYWORDS

Recommendations, Diversity, Fairness, Pre-processing, User-centered

## 1 INTRODUCTION

In today's digital age, personalized recommender systems form a solution to the information overload problem. They filter the enormous amount of information available online to only the top-N items - e.g. news articles, movies or books - that are most relevant to a user [19, 22]. Personalized recommendations contribute to improving user experience by providing accurate recommendations that align with user preferences.

While accurate recommendations ensure relevance to users' interests, excessive personalization risks closing users within filter bubbles, limiting their exposure to diverse content [4, 24, 26]. Diversity in recommender systems has recently been an active research area. This trend aligns with a larger movement in the recommender systems community, where the goal is not only to improve accuracy but also other factors, which we refer to as 'beyond accuracy', such as diversity and fairness [15, 44, 47].

To achieve diversity, several factors play a role, both relating to the recommender systems itself and to contextual factors. Editors, for instance, can shape news recommendations by curating content [43]. The user interface design influences user interaction and exploration [18]. The multi-stakeholder nature of recommender systems involves various providers, platform operators, and users, which introduces further complexity [1]. Relating to the recommender system itself, various diversification approaches have emerged. Some studies use post-processing approaches such as re-ranking to diversify the recommendations [21, 46]. Others build diversity directly into the recommendation algorithm itself, e.g. through multi-objective optimization, in so-called in-processing approaches [28, 32, 48]. Pre-processing approaches involve modifying input data before the recommendation process. Several works have proposed pre-processing approaches to improve beyond-accuracy factors such as user fairness [10, 27], provider fairness [4], and privacy [30].

In this paper, we present a pre-processing approach to diversify the output of a recommender system. We alter user profiles that are input to the recommender system to achieve greater diversity in the recommender's output. We opt for pre-processing for two reasons. Firstly, pre-processing has shown good results for other beyond-accuracy metrics, such as fairness [4] and privacy [30, 38]. Secondly, a pre-processing approach can be combined with a wide range of algorithms. Any personalized recommender system architecture takes a user profile as input, and can thus be used in combination with our approach to diversification.

Our approach is user-centric, focusing on personalized addition of items to user profiles. This is because the input data for recommender systems is highly personal. A user profile is more than a list of historical interactions with items; it encodes parts of the user's preferences, abilities, and characteristics. Altering a user profile risks diluting this rich encoding of the user's taste. By adopting a personalized approach, we mitigate this risk, maintaining accuracy without compromising the integrity of the user profile. We opt

to alter the user profile directly instead of altering processed versions of it, e.g. a latent representation. This increases the potential for transparency and explainability, allowing users to review their profiles, including any additions or removals.

Our main research question[1] focuses on how to alter user profiles to improve the diversity of top-N recommendations without compromising recommendation performance. We present two variants of our approach: one involves adding interactions to user profiles, while the other includes both adding and removing interactions to remain close to the original profile size. We explore different levels of adjustment ranging from 1% to 10%. We perform extensive evaluations to test which variant achieves optimal accuracy and diversity. We evaluate the effectiveness of our approach on two public data sets: MIND for news recommendation, where we aim to diversify news categories, and GoodBook for book recommendation, where we aim to diversify genres. As mentioned above, our approach can be integrated with any personalized recommender system. In this study, we test combinations with seven algorithms ranging from standard collaborative filtering to neural network recommenders.

We report accuracy metrics (MRR, nDCG, and HR), normative diversity metrics (calibration via divergence metrics), descriptive diversity metrics (coverage and Gini index), and provider fairness metrics (fair-nDCG). Our findings indicate that using pre-processed data for training can lead to recommender systems achieving comparable or improved performance compared to those trained on original data. Regarding diversity, calibration metrics consistently improved for certain algorithms, while coverage and the Gini index showed mixed results across different levels of pre-processing. Additionally, pre-processed data consistently resulted in higher fair-nDCG scores, indicating enhanced exposure fairness and better representation of minority categories.

## 2 BACKGROUND AND RELATED WORK

The literature on recommender systems has moved on from evaluating only the accuracy of recommended items towards evaluating also other factors, which we call 'beyond accuracy' factors, such as diversity, fairness [12, 15, 17, 34]. This section, first, provides an overview of existing work on pre-, in-, and post-processing approaches to increase beyond-accuracy factors. Next, we discuss the interplay between diversity, filter bubbles, and fairness. Finally, we examine existing measures used to evaluate diversity in recommender systems.

### 2.1 (Pre, In, Post)-processing approaches

Prior work has shown that specific minority user or provider groups, distinguished by sensitive attributes such as gender or race, experience disproportionate effects of indirect or unintentional discrimination [4, 11, 16]. In a common categorization, we distinguish between three approaches to detecting and mitigating bias and discrimination in recommender systems [47]: (i) data pre-processing phase, (ii) in-processing during model learning and optimization, and (iii) post-processing phase.

*Pre-processing approaches* involve modifying input data before the recommendation process to reduce bias [47]. Inspired by data

poisoning attacks, in [27], the authors introduce antidote/fake data, into the training set of a matrix factorization algorithm. Their approach involves augmenting the training set with new users and simulating ratings on existing items. These ratings are selected in such a way as to improve a socially relevant aspect of the recommendations, i.e., individual vs. group fairness, given to the original users. Inspired by data obfuscation in privacy-preserving techniques, in [30], the authors propose personalized blurring (PerBlur), a gender obfuscation approach aimed at protecting user privacy while maintaining recommendation quality and recommending less gender stereotypical items. In [10], the authors use a re-sampling strategy to adjust the proportion of users in groups. They generate gender-balanced training data and retrain the recommendation algorithms accordingly. Their findings indicate that while resampling the data led to a slight reduction in recommender accuracy, it did not introduce new gender disparities in performance for LastFM data and appeared to mitigate such differences in the Movielens data set. In [4], the authors evaluate discrimination in provider fairness by examining disparities in relevance, visibility, and exposure among minority groups. The disparity consists of items of minority groups of providers receiving unfairly low relevance. To do so, they first propose a pre-processing strategy that up-samples interactions where the minority group is the predominant. Then, they add an in-processing component that aims to control the relevance given to the items of the minority group, which is proportional to the minority group's contribution to the catalog. Their results show that up-sampling brings benefits to disparate impacts and coverage while maintaining the recommendation accuracy.

Pre-processing approaches offer flexibility by modifying input data exclusively, allowing existing recommender system algorithms to operate on adjusted data [22, 47]. Our study focuses on a pre-processing approach designed to integrate into any personalized recommender system algorithm without extensive modifications.

*In-processing approaches* involve incorporating beyond-accuracy considerations directly into the recommendation algorithm itself. In [25], the authors introduce a multi-objective optimization solution for music recommendations, balancing diversity and similarity to user preferences. In [12], the authors focus on shifting user consumption towards less familiar content in music streaming. They consider two key factors, taste similarity, which refers to how similar a piece of music is to the type of music the user has historically streamed, and popularity, which measures how many users have recently streamed the piece of content. In [48], the authors propose diversifying the recommendations by optimizing the selection of neighbors in the graph, category-boosted negative sampling, and adversarial learning on top of Graph Convolutional Networks. They focus on category diversification by making items of minority categories more reachable. More recently, the authors in [32] introduce SMORL, a Scalarized Multi-Objective Reinforcement Learning for recommender system setting. SMORL addresses multi-objective recommendation tasks by augmenting the recommender system with an additional layer to satisfy the accuracy, diversity, and novelty of recommendations. In [28], the authors introduce diversity into a two-tower news recommender architecture. This architecture incorporates a category loss function during the training phase, which aligns the representation of news items across a spectrum of news categories. Their findings demonstrate that this approach can

---

provide accurate and diverse news recommendations. In [41], the authors introduce LeaDivRec, a news recommendation approach designed to generate diversity-aware recommendations in an end-to-end fashion. They apply a diversity-aware regularization method to encourage the model to produce controllable diversity recommendations. Even though the in-processing approaches allow more control over the accuracy-beyond accuracy trade-off during training, they are designed for specific models and cannot be generalized to other models [47].

*Post-processing approaches* use a re-rank or post-processing technique to alter the recommendations after the generation of candidate items [48]. To balance the accuracy and diversity of the recommendations, the order and position of the items are determined using heuristics. In [21], the authors propose FairMatch, a graph-based approach to improve aggregate diversity in recommendations. In [45, 46], the authors investigate the challenge of generating a fair ranking while preserving high utility, considering protected attribute(s). They propose to mitigate the systematic bias through a ranked group fairness criterion. Post-processing approaches typically incorporate a tuning parameter aimed at balancing the accuracy-beyond accuracy trade-off. However, achieving beyond accuracy factors such as fairness or diversity without compromising recommendation performance can be challenging, and selecting an appropriate tuning parameter is crucial as it substantially impacts the recommendations performance [19, 41].

## 2.2 On the interplay between diversity, filter bubbles, and provider fairness

The impact of diversifying the recommendation goes beyond maintaining the accuracy and the user satisfaction into other aspects including filter bubbles and fairness.

Considering the interplay between filter bubbles and diversity: according to [26], filter bubbles are an environment created by personalization algorithms, where users are exposed solely to familiar information or opinions. Filter bubbles can potentially harm the democratic process and exacerbate polarization [26]. In [23], the authors state that "a technological filter bubbles is a decrease in the *diversity* of a user's recommendations over *time*, in any dimension of *diversity*, resulting from the choices made by different recommendation *stakeholders*." In [24], The authors provide an in-depth analysis to measure the technological filter bubbles in a longitudinal study of online news websites. They formalize variety as a property of diversity and introduce a binomial mixed-effects regression model to measure decreases in the variety of the recommendations provided to users over time. In [20], the authors examine the formation of filter bubbles in the Brazilian political news domain. They introduce a new metric based on the homogenization of recommended items' to measure the filter bubbles. Their findings indicate that diversification can reduce the homogenization.

Considering the interplay between fairness and diversity, we believe diversity is related to subject or provider fairness. In [9], the authors emphasize that diversity and fairness are rooted in distinct normative considerations. However, they note that certain diversity approaches could potentially contribute to improving provider fairness. In [37], the authors propose an algorithm that optimizes user and item fairness as a convex optimization problem.

Then, a ranking policy is derived via a Birkhoff-von Neumann decomposition algorithm which optimizes diversity. In [16], the authors propose a heuristic optimization that aims to solve the problem of provider fairness in multi-stakeholder recommender systems. The provided solution generates fair recommendations across different item providers in terms of coverage and diversity of users to whom the items are recommended.

## 2.3 Measuring diversity

In its early age, recommender systems' diversity was considered the opposite of similarity (1 - similarity) [31]. The similarity measures the proximity of the recommendations with the known preferences of users. In [34], the authors define novelty and diversity based on three concepts: preferences choice, discovery, and relevance. In [49], the authors introduce the Intra-List Similarity (ILS) and its counterpart, the Intra-List Distance. These metrics aim to measure the average diversity within a recommendation set ensuring that rearranging recommendations' positions does not affect the measure. Other metrics include relative diversity which evaluates diversity concerning a set of items, expected intra-list diversity which is sensitive to the ranking, and aggregate diversity which quantifies item diversity across recommendation lists [33, 34]. Lastly, the inverse of the Gini index calculates distributional inequality [5]. In [33], the authors propose an extensive analysis of the MIND news data pointing to the role of social media and news recommender system platforms in fostering a robust democratic discourse. [33] highlights that traditional descriptive diversity metrics rely on single-number measures, providing only a partial solution to the broader challenges recommender systems face beyond accuracy. Recent work on diversity in news recommendations has extended the literature by focusing on normative diversity. In [13], the authors highlight a number of principles designed for exposure diversity in recommender systems. In [35, 36], the authors introduce RADio, a rank-aware divergence framework to evaluate recommendations diversity according to normative goals.

For the present study, we recognize that no single metric captures all aspects of diversity. We evaluate our approach using two descriptive diversity metrics: item coverage and the Gini index. Also, we use rank-aware divergence metrics, e.g., Jensen-Shannon (JS) and KL-divergence (KL) as normative diversity metrics to measure the concept of calibration across categories of items.

## 3 PERSONALIZED USER-CENTRIC DATA PRE-PROCESSING

This section presents the problem formulation and outlines our user-centric pre-processing approach for recommender system data.

## 3.1 Problem formulation

Our objective is to develop a personalized approach to improve the diversity of recommended items without impacting the recommendation performance. Our approach extends user profiles through strategic addition and removal of interactions with items. We distinguish a one-step variant of the approach, in which we only add items to user profiles, and a two-step variant, in which we add and remove items. Figure 1 visualizes our approach.
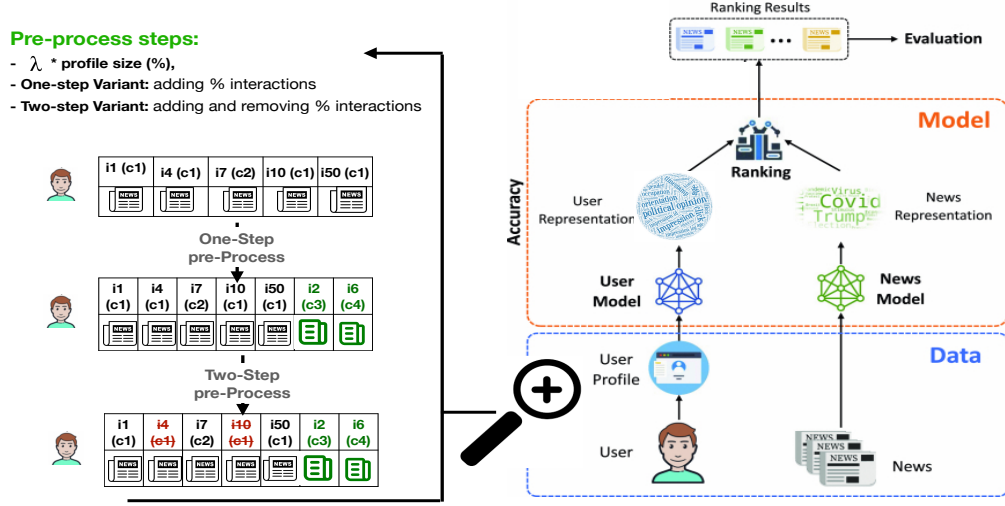
**Figure 1: Overall diagram for recommendation. To the left, we zoom in on our user-centered pre-processing. The one-step and two-step variants show our approach to altering the user profile.** $\{i1, i2, i4, i6, i7, i10, i50\}$ **represent items in target user** $u$ **profile.** $\{i1, i4, i7, i10, i50\}$ **are** $u$**'s past historical clicks.** $\{c1, c2, c3, c4\}$ **represent items** `categories`**. With one-step pre-processing,** $i2$**, and** $i6$ **are the personalized indicative items added to** $u$**. Following the two-step pre-processing,** $i4$ **and** $i10$ **can be removed from** $u$**.**

Let's assume that $U$ is a set of users, and $I$ is a set of items. $C$ represents the set of categories for each item $i$ in $I$. $R_{ui}$ represents the interaction that user $u$ has with item $i$ such that $R_{ui} = 1$ if user $u$ interacts with item $i$, 0 otherwise. $S_{ui}$ is a score representing the relevance of item $i$ to user $u$. $F_{ui}$ is a variable indicating whether item $i$ belongs to a category of interest to user $u$ ($F_{ui} = 1$ if item $i$ belongs to a relevant category, 0 otherwise). $W_{ic}$ is a weight representing the importance of category $c$ to item $i$. $S_{ui}$, $F_{ui}$, and $W_{ic}$ are key input variables for our personalized lists of indicative items for each user (section 3.2). $\lambda$ specifies the percentage of adjustment required for each user. The objective function can be formulated as:

$$\max \sum_{u \in U} \sum_{i \in I} \lambda \cdot S_{ui} \cdot R_{ui} \cdot F_{ui}$$

subject to the following constraints:
1. Personalized user interaction constraint:

$$\sum_{i \in I} R_{ui} = \text{desired interaction level}, \quad \forall u \in U$$

2. Category diversity constraint:

$$\sum_{i \in I} W_{ic} \cdot R_{ui} \cdot F_{ui} \geq \text{minimum desired diversity}, \quad \forall u \in U, \forall c \in C$$

3. Adjustment limit constraint (for two-step variant):

$$\sum_{i \in I} R_{ui} \leq \text{maximum adjustment limit}, \quad \forall u \in U$$

## 3.2 One-step pre-processing

We present the basic skeleton of our user-centric data pre-processing in Algorithm 1. The input to the algorithm includes the level of modification $\lambda$, expressed as a percentage by which the user profile is to be extended. Additionally, it involves lists of indicative users per category and lists of items per category. The lists of indicative

users per category are generated by training a logistic regression model on labeled training data. The coefficients $\beta = \{\beta_0, \beta_1, ..., \beta_M\}$ of the logistic regression capture the extent to which each user is correlated with the class attribute categories of items. These coefficients are instrumental in selecting users for the lists and ordering them based on the strength of the association. We extend user profiles by adding interactions until they are $\lambda$ percent longer than the original profile. We initiate our one-step pre-processing, detailed in Algorithm 1. The algorithm starts by generating users' personalized lists of indicative items. Our approach is founded on the principle that items added to the user profile should closely align with user preferences. This alignment increases the likelihood of maintaining recommendations accuracy when using pre-processed data for training. In the one-step variant, we add interactions to a user's profile consistent with their preferences. This involves employing a personalized list of indicative items for each user, denoted as $Personalized_L^u$. The creation of this list involves intersecting a personalized list of preferred items for each user with lists of users indicative for each category and a list of items per category. This personalized list is ordered based on the probability that the user would have interacted with the item. To generate the personalized list, we use the UserKNN algorithm since it provides a confidence score. It is selected due to its simplicity in providing the count of neighbors for prediction as a confidence score [30]. We sort the items in descending order based on counts of neighbors to obtain $List_{NCounts}^u$, representing our personalized list for each user.

---

**Algorithm 1:** One-step pre-processing algorithm.

---

**Input:**

- User preferences (confidence-score based user preferences), List of users indicative for each category of items (*user categories*), List of items per category (*item categories*)
- $\lambda$: level of pre-processing
- Original user-item matrix $\mathcal{D}$ ($\mathbb{N}$ users, $\mathbb{M}$ items)
- Initial count: user profile size at time $t = 0$

**Output:** One-step pre-processed data $\mathcal{D}'$ ($\mathbb{N}$ users, $\mathbb{M}$ items)

    `// 0.` **`Users' personalized lists of indicative items`**

The confidence score for recommendation based on UserKNN;

**for** *(user* u *in* $\mathbb{N}$*)* **do**
    **for** *(item* i *in* $\mathbb{M}$*)* **do**
        Similarity computation finds nearest neighbor candidates;
        Sort selected items based on the number of possessed neighbor candidates;

$List^u_{NCounts}$ contains a list of counts for each user $u$;

**for** *(user* u *in* $\mathbb{N}$*)* **do**
    Fix a cutoff on the List of users indicative for each category of items `// we set the cutoff to Top-50 for MIND News data and to top-100 for GoodBood Data, in the rest of the experiments`
    Create new personalized list of indicative items for $u$: $Personalized^u_L$;
    Find *non member categories* as categories not in *user categories*
    **if** *(u is a non member of categories)* **then**
        Choose a random category from non member categories.
        Find items of the chosen category from *item categories*.
        `//` $Personalized^u_L$ `= intersection between items in the chosen category and items in user preferences`
    **for** *item* $i \in List^u_{NCounts}$ **do**
        $Personalized^u_L = Personalized^u_L$. add $(i)$

               `// 1.` **`One-step: adding extra items`**

**for** *(user u in* $\mathbb{N}$*)* **do**
    count = initial count [u] $* \lambda$
    added = 0
    **while** *added < count* **do**
        i = picks the item in the first position in $Personalized^u_L$
        **if** $\mathcal{D}'[u, i] == 0$ **then**
            $\mathcal{D}'[u, i] == 1$
        added += 1
    Total added += added

---

**Algorithm 2:** Two-step pre-processing algorithm.

---

**Input:**

- One-step pre-processed data $\mathcal{D}'$ ($\mathbb{N}$ users, $\mathbb{M}$ items)
- Total added: total number of extra interactions added, Interaction count: user profile size after adding $p\%$ extra items.
- Removal threshold

**Output:** Two-step pre-processed data $\mathcal{D}''$ ($\mathbb{N}$ users, $\mathbb{M}$ items)

    `// 2.` **`Two-step pre-processing: Removing certain items`**

**for** *user u in* $\mathcal{N}$ **do**
    **if** *Interaction count $\geq$ Removal threshold* **then**
        `// Removal threshold is chosen by us to be 20.`
        remove count += 1

To be removed = Total added / remove count `// To be removed: contains the number of interactions that will be removed from individual user profiles.`

**for** *user u in* $\mathcal{N}$ **do**
    **if** *(Interaction count $\geq$ Removal threshold)* **then**
        removed = 0
        **while** *(removed < To be removed [u] )* **do**
            i = pick an item that is not recently added `// i depends on the removal mode: random or greedy`
            **if** $\mathcal{D}'[u, i]\, != 0$ **then**
                $\mathcal{D}''[u, i] == 0$
                removed += 1
            `//` $\mathcal{D}''$ `is` $\mathcal{D}'$ `after applying the removal`

---

indicators of preference. We randomly sample 80% of each user profile for the training set and we keep the remaining 20% for the test set. It is important to note that the choice of static splitting plays a key role in preventing the data pre-processing from adding items into the test set. Our pre-process is solely applied to the training set, with careful consideration to avoid adding items already present in the test set. As a result, the test set remains consistent and invariant across all experimental conditions. Further details about our pre-processed data using both one-step and two-step variants are presented in the Supplementary material (Table 1).

**Table 1: Statistics of the data sets.**

| Data Sets | | #Users | #Click history | Sparsity (%) | #Items | Features | |
|---|---|---|---|---|---|---|---|
| | | | | | | #Categories | #Impressions |
| **MIND News** | *Training Data* | 1000 | 9368 | 99.58 | 26740 | 17 | 7105 |
| **Subset** | *Test Data* | 5000 | 15557 | 99.82 | 18723 | 16 | 7538 |
| **GoodBook** | *Training Data* | 943 | 8477 | 98.77 | 729 | 31 | - |
| **Subset** | *Test Data* | 943 | 4715 | 99.27 | 688 | 31 | - |

We provide the frequency of news categories and book genres in the training and testing data in our supplementary material. In the MIND data, there are 17 categories, whereas in the GoodBook data, there are 31 genres. A news article is assigned to a single category, while a book is associated with multiple genres.

## 3.3 Two step pre-processing

We proceed to the second variant of our data pre-processing approach, denoted as "two-step pre-processing," detailed in Algorithm 2. This algorithm takes the data generated by one-step pre-processing as input and focuses on the removal of specific interactions. The removal process helps in maintaining the density of the modified data close to that of the original data. The removal is performed in such a way that the total number of user interactions for each item remains proximate to the total in the original dataset. We distribute the removal of interactions evenly across all users. To account for users with very short profiles, we set a threshold = 20 that exempts them from having interactions removed.

## 4 EXPERIMENTAL SETUP

In this section, we start by providing an overview of our datasets. Next, we describe recommender system algorithms used for both news and book recommendations in our experiments.

## 4.1 Data Sets

We selected two publicly available datasets: the MIND news data [42] for news recommendations and the GoodBook data for book recommendations. Statistics for our datasets can be found in Table 1. The published MIND data includes train, validation, and test splits. For the GoodBook dataset, we transform explicit ratings into implicit feedback, considering ratings with a cutoff of rating >= 3 as implicit

## 4.2 Recommender System Algorithms

In this section, we describe our recommender system's pipeline, including the selection of recommendation algorithms, hyperparameter tuning, and candidate item selection strategy. We select a number of recommendation algorithms, ranging from standard collaborative filtering techniques to advanced neural-based recommendation models. This selection allows us to demonstrate the adaptability of our pre-processing approach across various recommendation architectures. For news recommendations, we use the state-of-the-art neural network recommendation algorithms:

**Table 2: Recommendation performance for the MIND Data. We use one and two-step pre-processing with $\lambda$ values of {0%, 1%, 2%, 5%, 10%}. We measure the accuracy of the Top-K (K= 5 and 10) recommendation using nDCG and MRR. We measure diversity using: (1) a normative framework focusing on calibration by categories of news, including KL divergence (*KL*) and Jensen Shannon (*JS*), and (2) a descriptive measure using item coverage (*Cov*). Higher scores indicate better performance. We emphasize the recommendations with the highest scores across conditions (original vs pre-processed data).**

| | | One-step Pre-process | | | | | | Two-step Pre-process | | | | | |
| | | Accuracy | | | Calibration @10 | | Descriptive @10 | Accuracy | | | Calibration @10 | | Descriptive @10 |
| *Algorithms* | *λ (%)* | MRR | nDCG@5 | nDCG@10 | KL | JS | Cov | MRR | nDCG@5 | nDCG@10 | KL | JS | Cov |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **NRMS** | 0 | 0.2744 | 0.2970 | 0.3643 | 2.1947 | 0.3918 | 1342 | 0.2744 | 0.2970 | 0.3643 | 2.1947 | 0.3918 | 1342 |
| | 1 | 0.2773 | 0.2995 | *0.3639* | 2.2307 | 0.3949 | 1548 | 0.2769 | 0.2982 | 0.3630 | 2.2427 | 0.3959 | 1519 |
| | 2 | 0.2783 | 0.2999 | 0.3651 | 2.2490 | 0.3960 | 1557 | 0.2795 | 0.3019 | 0.3651 | 2.2672 | 0.3978 | 1502 |
| | 5 | 0.2759 | 0.2963 | 0.3620 | 2.2063 | 0.3939 | 1509 | 0.2767 | 0.2993 | 0.3643 | 2.2595 | 0.3985 | 1492 |
| | 10 | 0.2752 | 0.2984 | 0.3616 | 2.2579 | 0.3972 | 1554 | 0.2722 | 0.2955 | 0.3587 | 2.2221 | 0.3950 | 1521 |
| **NPA** | 0 | 0.2776 | 0.2997 | 0.3638 | 2.1994 | 0.3906 | 1314 | 0.2776 | 0.2997 | 0.3638 | 2.1994 | 0.3906 | 1314 |
| | 1 | 0.2677 | 0.2885 | 0.3513 | *2.1975* | 0.3898 | 1323 | 0.2712 | 0.2906 | 0.3553 | 2.1652 | 0.3893 | 1293 |
| | 2 | 0.2689 | 0.2890 | 0.3536 | *2.1702* | 0.3894 | 1269 | 0.2690 | 0.2880 | 0.3528 | *2.1841* | *0.3905* | 1295 |
| | 5 | *0.2726* | 0.2943 | 0.3573 | *2.1805* | *0.3905* | 1282 | 0.2697 | 0.2886 | 0.3533 | 2.1560 | *0.3889* | *1310* |
| | 10 | *0.2764* | 0.3009 | 0.3624 | *2.1878* | 0.3904 | *1304* | *0.2697* | 0.2807 | 0.3473 | 2.1556 | 0.3901 | 1287 |
| **LSTUR** | 0 | 0.2751 | 0.2960 | 0.3569 | 2.2243 | 0.3919 | 1586 | 0.2751 | 0.2960 | 0.3569 | 2.2243 | 0.3919 | 1586 |
| | 1 | *0.2730* | 0.2943 | 0.3574 | *2.2182* | *0.3912* | 1589 | 0.2727 | 0.2931 | 0.3524 | 2.1931 | *0.3914* | 1590 |
| | 2 | 0.2725 | 0.2931 | 0.3554 | 2.1940 | *0.3911* | 1624 | 0.2775 | 0.2992 | 0.3601 | 2.2159 | *0.3895* | 1646 |
| | 5 | 0.2810 | 0.3015 | 0.3635 | 2.2431 | 0.3948 | *1581* | 0.2798 | 0.3017 | 0.3638 | 2.2289 | 0.3926 | 1575 |
| | 10 | 0.2713 | 0.2925 | 0.3540 | 2.2354 | 0.3925 | 1600 | 0.2766 | 0.2989 | 0.3601 | 2.2663 | 0.3955 | 1552 |

**Table 3: Recommendation performance for the GoodBook Data. We use one and two-step pre-processing with $\lambda$ values of {0%, 1%, 2%, 5%, 10%}. We measure the accuracy of the Top-K (K= 5 and 10) recommendation using HR and nDCG. We measure diversity using: (1) a normative framework focusing on calibration by categories, including KL divergence (*KL*) and Jensen Shannon (*JS*), and (2) a descriptive measure using item coverage (*Cov*) and Gini index (*Gini*). Higher scores indicate better performance. We emphasize the recommendations with the highest scores across conditions (original vs pre-processed data).**

| | | One-step Pre-process | | | | | | | | | | | | Two-step Pre-process | | | | | | | | | | | |
| | | Top-K = 5 | | | | | | Top-K = 10 | | | | | | Top-K = 5 | | | | | | Top-K = 10 | | | | | |
| | | Accuracy | | Calibration | | Descriptive | | Accuracy | | Calibration | | Descriptive | | Accuracy | | Calibration | | Descriptive | | Accuracy | | Calibration | | Descriptive | |
| *Algorithms* | *λ (%)* | HR | nDCG | KL | JS | Cov | Gini | HR | nDCG | KL | JS | Cov | Gini | HR | nDCG | KL | JS | Cov | Gini | HR | nDCG | KL | JS | Cov | Gini |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **MostPop** | 0 | 0.1209 | 0.0347 | 0.0547 | 0.1652 | 20 | 0.2547 | 0.1909 | 0.0507 | 0.0626 | 0.1784 | 31 | 0.3493 | 0.1209 | 0.0347 | 0.0547 | 0.1652 | 20 | 0.2547 | 0.1909 | 0.0507 | 0.0626 | 0.1784 | 31 | 0.3493 |
| | 1 | 0.1209 | 0.0347 | 0.2335 | 0.3348 | 20 | 0.2547 | 0.1909 | 0.0507 | 0.1985 | 0.3091 | 31 | 0.3493 | 0.1251 | 0.0331 | 0.2771 | 0.3587 | 17 | 0.2999 | 0.1994 | 0.0485 | 0.1988 | 0.3090 | 29 | 0.3707 |
| | 2 | 0.1188 | 0.0343 | 0.2343 | 0.3353 | 20 | 0.2548 | 0.1994 | 0.0512 | 0.1994 | 0.3095 | 31 | 0.3494 | 0.1347 | 0.0345 | 0.2359 | 0.3361 | 16 | 0.3184 | 0.1951 | 0.0482 | 0.1999 | 0.3094 | 29 | 0.3699 |
| | 5 | 0.1188 | 0.0335 | 0.1972 | 0.3131 | 20 | 0.2562 | 0.2036 | 0.0514 | 0.2005 | 0.3102 | 31 | 0.3499 | 0.1241 | 0.0302 | 0.1959 | 0.3107 | 15 | 0.3341 | 0.2078 | 0.0468 | 0.1809 | 0.2976 | 25 | 0.4246 |
| | 10 | 0.1166 | 0.0328 | 0.1981 | 0.3150 | 18 | 0.2835 | 0.1941 | 0.0491 | 0.1968 | 0.3029 | 31 | 0.3501 | 0.1389 | 0.0309 | 0.1896 | 0.3038 | 12 | 0.4113 | 0.2153 | 0.0445 | 0.1766 | 0.2883 | 20 | 0.5227 |
| **ItemKNN** | 0 | 0.5483 | 0.1810 | 0.0646 | 0.1751 | 564 | 0.3577 | 0.6839 | 0.2200 | 0.0819 | 0.1944 | 683 | 0.3471 | 0.5483 | 0.1810 | 0.0646 | 0.1751 | 564 | 0.3577 | 0.6839 | 0.2200 | 0.0819 | 0.1944 | 683 | 0.3471 |
| | 1 | 0.5793 | 0.1816 | 0.1819 | 0.3021 | 574 | 0.3565 | 0.6839 | 0.2200 | 0.1711 | 0.2899 | 665 | 0.3503 | *0.5398* | *0.1798* | 0.1821 | 0.3031 | 576 | *0.3560* | 0.6978 | 0.2213 | 0.1711 | 0.2899 | 665 | 0.3503 |
| | 2 | 0.5356 | 0.1780 | 0.1835 | 0.3029 | 564 | 0.3435 | 0.6649 | 0.2176 | 0.1712 | 0.2892 | 663 | 0.3345 | *0.5451* | *0.1787* | 0.1837 | 0.3036 | 568 | 0.3429 | 0.6734 | 0.2178 | 0.1719 | 0.3094 | 663 | 0.3384 |
| | 5 | 0.5122 | 0.1704 | 0.1855 | 0.3038 | 548 | 0.3070 | 0.6320 | 0.2034 | 0.1741 | 0.2906 | 662 | 0.2886 | 0.5090 | 0.1667 | 0.1859 | 0.3044 | 553 | 0.3229 | 0.6532 | 0.2032 | 0.1747 | 0.2916 | 665 | 0.3503 |
| | 10 | 0.4507 | 0.1482 | 0.1886 | 0.3052 | 465 | 0.2498 | 0.5705 | 0.1783 | 0.1746 | 0.2907 | 602 | 0.2249 | 0.4846 | 0.1526 | 0.1826 | 0.3031 | 537 | 0.3108 | 0.6129 | 0.1840 | 0.1729 | 0.2905 | 663 | 0.2961 |
| **ImplicitMF** | 0 | 0.5588 | 0.1808 | 0.0667 | 0.1790 | 498 | 0.4755 | 0.6914 | 0.2208 | 0.0853 | 0.1998 | 608 | 0.4629 | 0.5588 | 0.1808 | 0.0667 | 0.1790 | 498 | 0.4755 | 0.6914 | 0.2208 | 0.0853 | 0.1998 | 608 | 0.4629 |
| | 1 | 0.5525 | 0.1784 | 0.1817 | 0.3028 | 495 | 0.4832 | *0.6893* | 0.2208 | 0.1693 | 0.2890 | 607 | 0.4661 | *0.5503* | *0.1785* | 0.1795 | 0.3021 | 498 | 0.4816 | 0.6850 | 0.2168 | 0.1701 | 0.2892 | 603 | 0.4695 |
| | 2 | 0.5471 | 0.1768 | 0.1792 | 0.3021 | 498 | 0.4752 | 0.6797 | 0.2152 | 0.1716 | 0.2894 | 605 | 0.4666 | *0.5525* | *0.1749* | 0.1793 | 0.3019 | 504 | 0.4724 | 0.6861 | 0.2169 | 0.1720 | 0.2911 | 605 | 0.4679 |
| | 5 | 0.5058 | 0.1486 | 0.1847 | 0.3046 | 505 | 0.4642 | 0.6489 | 0.1833 | 0.1732 | 0.2911 | 630 | 0.4533 | 0.4995 | 0.1452 | 0.1837 | 0.3041 | 518 | 0.4579 | 0.6447 | 0.1768 | 0.1724 | 0.2912 | 643 | 0.4554 |
| | 10 | 0.4454 | 0.1238 | 0.1870 | 0.3039 | 542 | 0.4592 | 0.5769 | 0.1556 | 0.1724 | 0.2895 | 649 | 0.4654 | 0.4051 | 0.1121 | 0.1823 | 0.3020 | 559 | 0.4789 | 0.5408 | 0.1372 | 0.1696 | 0.2875 | 657 | 0.4889 |
| **BPR** | 0 | 0.5355 | 0.1839 | 0.0673 | 0.1795 | 395 | 0.4494 | 0.6808 | 0.2228 | 0.0837 | 0.1979 | 503 | 0.4463 | 0.5355 | 0.1839 | 0.0673 | 0.1795 | 395 | 0.4494 | 0.6808 | 0.2228 | 0.0837 | 0.1979 | 503 | 0.4463 |
| | 1 | 0.5620 | 0.1882 | 0.1779 | 0.2997 | 387 | 0.4784 | 0.7052 | 0.2262 | 0.1686 | 0.2885 | 495 | 0.4579 | 0.5567 | 0.1849 | 0.1781 | 0.3001 | 401 | 0.4762 | *0.6744* | *0.2199* | 0.1700 | 0.2893 | 499 | 0.4481 |
| | 2 | 0.5429 | *0.1812* | 0.1770 | 0.3002 | 395 | 0.4519 | *0.6925* | 0.2245 | 0.1683 | 0.2883 | 495 | 0.4489 | 0.5578 | *0.1789* | 0.1774 | 0.3002 | 418 | 0.4434 | 0.6734 | *0.2130* | 0.1656 | 0.2884 | 509 | 0.4353 |
| | 5 | 0.5207 | 0.1670 | 0.1817 | 0.3027 | 424 | 0.4386 | 0.6532 | 0.2041 | 0.1666 | 0.2868 | 538 | 0.4358 | 0.5048 | 0.1606 | 0.1821 | 0.3017 | 439 | 0.4462 | 0.6405 | 0.1984 | 0.1689 | 0.2884 | 569 | 0.4264 |
| | 10 | 0.4909 | 0.1533 | 0.1829 | 0.3023 | 442 | 0.4222 | 0.6129 | 0.1872 | 0.1665 | 0.2876 | 576 | 0.3987 | 0.3701 | 0.1068 | 0.1829 | 0.3013 | 465 | 0.4142 | 0.5016 | 0.1318 | 0.1657 | 0.2859 | 615 | 0.4048 |

*NRMS* uses multi-head self-attention networks to learn news and user representations [40]. *NPA* uses personalized attention networks to learn news and user representations [39]. *LSTUR* is a news recommendation approach capturing users' both long-term and short-term preferences [2]. We repeat each experiment 3 times and show average performance. We use the recommenders package from Microsoft.[2] For hyper-parameter tuning, we set epochs to 50. We follow the parameters as suggested by recommenders.

For book recommendations, we use state-of-the-art collaborative filtering algorithms: *MostPop* is a non-personalized algorithm recommending the most popular items. *ItemKNN* is the item-item K-nearest neighbor algorithm for top-N recommendations [8]. *ImplicitMF* is implicit matrix factorization trained with alternating least squares (ALS) [14]. *BPR* is a matrix factorization algorithm using Bayesian personalized ranking for implicit data [29]. We use the Lenskit toolkit for our implementation of standard recommender

---

[2]https://github.com/recommenders-team/recommenders

system algorithms.[3] For ItemKNN, implicitMF, and BPR, we perform hyperparameter tuning on the training set to optimize their performance. We adjust parameters such as the number of nearest neighbors for ItemKNN, the features and iterations for implicitMF, and the epochs, batch size, and features for BPR.

In our experiments, we used the AllItems strategy for our candidate item selection, as proposed in [3]. AllItems generates a list of candidate items $L_u$ for each user $u$, excluding items that user $u$ has previously interacted with in the training set $Tr_u$. Formally, given a set of all items $I$ in the data set and the training set vector $Tr_u$ representing the set of items that user $u$ has interacted with, the list of candidate items $L_u$ for user $u$ can be formulated as $L_u = I \setminus Tr_u$.

*Accuracy of recommendations.* We follow existing works on MIND news recommendations, we use the Normalized Discounted Cumulative Gain (nDCG), and Mean Reciprocal Rank (MRR) to measure the accuracy of the recommendations [42]. As for book recommendations, we use the Hit rate (HR) and nDCG. nDCG measures the ranking quality of the top recommended items.

*Diversity and Fairness in the recommendation list.* beyond optimizing for the accuracy of the recommendations, we look at measuring the diversity and fairness in our recommendation lists. We use the discounted cumulative fairness [44]. Given that we have $C$ categories of items (either news articles or books), each item is assigned to one or many categories. We compute the fair-nDCG-score by counting the unique categories a user is exposed to in the top recommendation list. We vary $K$ from $\{1, .., 100\}$. Fair-nDCG-score accumulates the number of items belonging to the protected group $G+$ ($G+ \in C$) at discrete positions in the ranking (at $k = \{1, 5, ..100\}$) and discounts these numbers accordingly, to favor the representation of the protected group at higher positions [44]. We note that for MIND News data, we consider the niche categories:'$kids'$, '$weather'$, '$video'$, '$music'$, '$autos'$, '$movies'$, '$middleeast'$, '$northamerica'$ as our protected group that we aim to see in the top recommendation list. For GoodBook data, we consider '$music'$, '$poetry'$, '$horror'$, '$spirituality'$, '$sports'$, '$christian'$, '$comics'$, '$manga'$, '$cookbooks'$, '$psychology'$, '$art'$ as our protected group.

We utilize *calibration* as our target normative diversity, evaluating the extent to which the recommendations are tailored to the users' preferences. We focus on calibration as it measures the alignment of recommendations with users' preferences inferred from their click history. Calibration can have two aspects: the divergence of the recommended articles' categories and complexity [36]. We specifically address the latter aspect since item categories are provided in the metadata. We compare the differences in distributions between users' preferences and generated recommendations. We calculate the Kullback-Leibler (KL) and the Jensen-Shannon (JS) divergence metrics. KL and JS measure the divergence between two probability mass functions, but with JS offering a divergence score bounded between 0 and 1 and does not introduce a consequential skew in the distribution of users' past historical interactions vs. recommended items [36]. Additionally, we compute two descriptive metrics: item coverage and the Gini index. Item coverage measures the proportion of unique items recommended to users, and the Gini index quantifies the inequality in item distribution among users. We report $1 -$ the Gini index.

## 5 RECOMMENDATION PERFORMANCE

In Table 2 and Table 3, we present our recommendation results using one-step and two-step pre-processing techniques on the MIND and GoodBook data sets, respectively. We measure the accuracy of the recommendation using MRR and nDCG. The one-step pre-processing involves adding a percentage $\lambda = \{0\%, 1\%, 2\%, 5\%, 10\%\}$ of interactions to users' profiles. We see that recommender system algorithms react differently to different $\lambda$. For NRMS and LSTUR, we observe improvements in nDCG@5 and nDCG@10 with certain levels of $\lambda$ compared to the original data. NRMS experiences improvement at 1%, 2%, 5%, and 10%. LSTUR exhibits improved performance at $\lambda = 5\%$. For the NPA, the performance in accuracy metrics remains relatively stable across different lambda levels. On the GoodBook data, the recommendation performance of ItemKNN and ImplicitMF varies across different levels of $\lambda$. BPR consistently maintains higher accuracy across lambda levels, particularly at 1% and 2%. This emphasizes the need to carefully select $\lambda$ in order to maintain recommendation quality.

Extending the analysis to two-step pre-processing, we observe a substantial effect on recommendation performance. While NRMS and LSTUR show improved nDCG at certain $\lambda = 2\%, 5\%, 10\%$, the performance of NPA remains comparable to that of original data. Similarly, on the GoodBook data, we see that the accuracy trends across $\lambda$ levels exhibit similar patterns to the one-step variant. Especially the BPR algorithm continues to demonstrate superior accuracy for $\lambda = 1\%$. Our results demonstrate the importance of carefully considering thresholds and personalized lists of indicative items per user to strike a balance between user interaction augmentation and reduction. It is important to note that the recommender system platform needs to select only one $\lambda$. By analyzing various $\lambda$, we aim to illustrate the impact of our pre-processing technique on improving recommendation performance and diversity.

## 6 MEASURING DIVERSITY IN RECOMMENDATIONS

In the previous section, we explored the impact of one-step and two-step pre-processing on recommendation performance. To evaluate diversity, we use the normative framework focusing on calibration and descriptive metrics focusing on coverage and the Gini index.

Table 2 and Table 3 present our diversity analysis for news and book recommendations, respectively. We observe variations in the results of calibration (KL, JS) and the descriptive metrics (Cov, Gini) when comparing recommendation outputs generated from pre-processed data with those from the original data. Recommender system algorithms exhibit diverse behaviors across different $\lambda$ values. Focusing on news recommendations, NRMS consistently yields higher calibration and coverage scores for $\lambda = 1\%, 2\%, 5\%, 10\%$. LSTUR provides its most diverse recommendations at $\lambda = 5\%, 10\%$. Recommendations generated using the NPA algorithm at different $\lambda$ values did not outperform recommendations generated using original data but remained comparable. For book recommendations, we note positive reactions to diversity from recommender system algorithms to varying $\lambda$ values. For instance, we see a global agreement on generating more calibrated recommendation outputs with pre-processed data ($\lambda = 1, 2, 5, 10$) compared to those using the original data for MostPop, ItemKNN, ImplicitMF, and BPR. However, this

trend does not apply uniformly across all descriptive metrics. While ImplicitMF and BPR exhibit optimal coverage and Gini index results for $\lambda = 1\%, 2\%$, ItemKNN performs best at $\lambda = 1\%$. Conversely, for higher $\lambda$ values, particularly $\lambda = 10\%$, the performance of ItemKNN (ImplicitMF and BPR) gradually decreases. This phenomenon may be attributed to the algorithms' sensitivity to popularity bias and the change in the long-tail distribution [6, 7].

## 7 MEASURING FAIRNESS

To evaluate fairness, we use the fair-nDCG-score, a metric designed to measure the exposure of minority target categories in the recommendation lists. This metric discussed in section 4.2 considers the representation of a protected group, denoted as $G+$, at different positions in the ranking, discounting these numbers to favor a fair representation of the protected group at higher positions [44].

Figure 2 and Figure 3 illustrate our fair-nDCG results for both MIND News and GoodBook Data. Each pair of figures corresponds to a specific algorithm, with one figure depicting the results for recommendations using one-step pre-processing, and the other for two-step pre-processing. The fair-nDCG scores are evaluated across different levels of $\lambda$ (1%, 2%, 5%, 10%) for recommendations ranging from 1 to 100. We observe that recommendations generated using pre-processed data consistently exhibit higher fair-nDCG scores, indicating increased exposure fairness. For instance, in Figure 2 (NRMS), recommendations using 10% one-step pre-processing (same for two-step) depict the highest fair-nDCG score, closely followed by 1% and 2% pre-processing levels. This emphasizes that pre-processing not only improves the accuracy of the recommendations but also contributes to increased exposure to minority categories.

Similar observations hold for GoodBook Data in Figure 3. For instance, in the case of BPR (Figure 3 (`Bottom`), recommendations with 1% and 5% One-step pre-processing exhibit the highest fair-nDCG-scores, signifying fair recommendations. Additionally, recommendations with 2% and 5% two-step pre-processing demonstrate noteworthy exposure fairness.

## 8 CONCLUSION

In this paper, we proposed a novel user-centered pre-processing approach to diversify the output of a recommender system. The central component of our approach involves personalizing the process of adding and removing $\lambda$ % of interactions from a user profile. Our results demonstrated that by keeping data pre-processing closely aligned with user preferences, it has the potential to maintain or even improve upon the performance on the original profiles. Our experiments have also shown that our pre-processing approach achieves diverse recommendations and simultaneously promotes provider fairness. For diversity, through extensive analysis utilizing normative and descriptive diversity measures, we have demonstrated that our approach has the potential to increase item coverage and Gini index and amplify divergence. For fairness, we have demonstrated that our approach actively promotes the recommendation of niche categories.

Our paper opens an important new vista for future work. Our approach has the potential to bridge disciplinary boundaries by engaging with other fields, such as *conducting user studies* to gather qualitative insights into how users perceive the recommendations
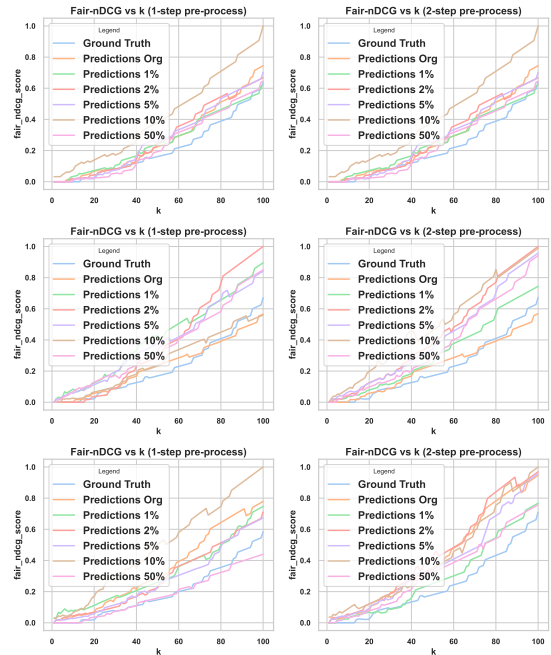


**Figure 2: Results of fair-nDCG using NRMS (`Top`), NPA (`Middle`), and LSTUR (`Bottom`) on MIND News data for different recommendation lists. The fair-nDCG-score is measured for the recommendation of a different *categories* of news in the top-k list.**

generated by our approach, or *ethical considerations* to examine the societal impact of diverse recommendations on users and society. Also, considering the foundational role of our processing of users' preferences, future work could explore the long-term dynamics of our recommendation evolution and whether diversity, fairness and accuracy are maintained.

## REFERENCES

[1] Himan Abdollahpouri, Gediminas Adomavicius, Robin Burke, Ido Guy, Dietmar Jannach, Toshihiro Kamishima, Jan Krasnodebski, and Luiz Pizzato. 2019. Beyond Personalization: Research Directions in Multistakeholder Recommendation. *arXiv preprint arXiv:1905.01986* (2019).

[2] Mingxiao An, Fangzhao Wu, Chuhan Wu, Kun Zhang, Zheng Liu, and Xing Xie. 2019. Neural News Recommendation with Long- and Short-term User Representations. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, Anna Korhonen, David Traum, and Lluís Màrquez (Eds.). 336–345.

[3] Alejandro Bellogin, Pablo Castells, and Ivan Cantador. 2011. Precision-Oriented Evaluation of Recommender Systems: An Algorithmic Comparison. In *Proceedings of the Fifth ACM Conference on Recommender Systems* (Chicago, Illinois, USA). 333–336.

[4] Ludovico Boratto, Gianni Fenu, and Mirko Marras. 2021. Interplay between upsampling and regularization for provider fairness in recommender systems. *User Modeling and User-Adapted Interaction* 31, 3 (2021), 421–455.

[5] Fleder Daniel and Hosanagar Kartik. 2009. Blockbuster Culture's Next Rise or Fall: The Impact of Recommender Systems on Sales Diversity. *Management Science* 55, 5 (2009), 697-712 pages.

[6] Savvina Daniil, Mirjam Cuper, Cynthia C. S. Liem, Jacco van Ossenbruggen, and Laura Hollink. 2023. Reproducing popularity bias in recommendation: the effect of evaluation strategies. *ACM Transaction on Recommender Systems* (2023).

[7] Yashar Deldjoo, Alejandro Bellogin, and Tommaso Di Noia. 2021. Explaining recommender systems fairness and accuracy through the lens of data characteristics. *Information Processing & Management* 58, 5 (2021), 102662.
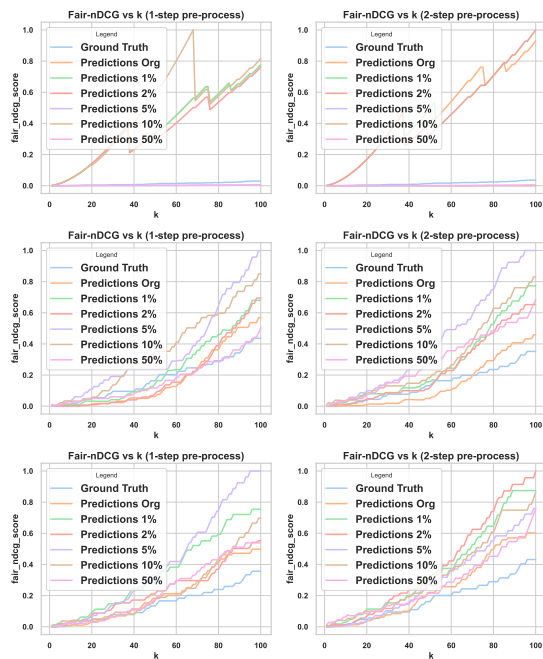
**Figure 3: Results of fair-nDCG-score using MostPop (`Top`), implicitMF (`Middle`), and BPR (`Bottom`) on GoodBook Data for different recommendation lists. The fair-nDCG is measured for the recommendation of different _genres_ of books in the topk list. The fairness is related to recommending "niche" or non-stereotypical categories.**

[8] Mukund Deshpande and George Karypis. 2004. Item-based top-N recommendation algorithms. *ACM Transaction on Information Systems* 22, 1 (2004), 143–177.

[9] Michael D. Ekstrand, Anubrata Das, Robin Burke, and Fernando Diaz. 2022. *Fairness in Recommender Systems.* Springer US, 679–707.

[10] Michael D. Ekstrand, Rezvan Joshaghani, and Hoda Mehrpouyan. 2018. Privacy for All: Ensuring Fair and Equitable Privacy Protections. In *Proceedings of the International Conference on Fairness, Accountability and Transparency (Proceedings of Machine Learning Research, Vol. 81)*, Sorelle A. Friedler and Christo Wilson (Eds.). 35–47.

[11] Michael D. Ekstrand, Mucun Tian, Ion Madrazo Azpiazu, Jennifer D. Ekstrand, Oghenemaro Anuyah, David McNeill, and Maria Soledad Pera. 2018. All The Cool Kids, How Do They Fit In?: Popularity and Demographic Biases in Recommender Evaluation and Effectiveness. In *Proceedings of the International Conference on Fairness, Accountability and Transparency (Proceedings of Machine Learning Research, Vol. 81)*, Sorelle A. Friedler and Christo Wilson (Eds.). PMLR, 172–186.

[12] 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 ACM International Conference on Web Search and Data Mining.* 238–246.

[13] Natali Helberger, Kari Karppinen, and Lucia D'Acunto. 2018. Exposure diversity as a design principle for recommender systems. *Information, Communication & Society* 21, 2 (2018), 191–207.

[14] Yifan Hu, Yehuda Koren, and Chris Volinsky. 2008. Collaborative Filtering for Implicit Feedback Datasets. In *IEEE International Conference on Data Mining.* 263–272.

[15] Marius Kaminskas and Derek Bridge. 2016. Diversity, Serendipity, Novelty, and Coverage: A Survey and Empirical Analysis of Beyond-Accuracy Objectives in Recommender Systems. *ACM Transactions on Interactive Intelligent Systems* 7, 1, Article 2 (2016), 42 pages.

[16] Evangelos Karakolis, Panagiotis Kokkinakos, and Dimitrios Askounis. 2022. Provider Fairness for Diversity and Coverage in Multi-Stakeholder Recommender Systems. *Applied Sciences* 12, 10 (2022), 4984.

[17] Matevž Kunaver and Tomaž Požrl. 2017. Diversity in recommender systems – A survey. *Knowledge-Based Systems* 123 (2017), 154–162.

[18] Yu Liang and Martijn C. Willemsen. 2022. Exploring the longitudinal effects of nudging on users' music genre exploration behavior and listening preferences. In *Proceedings of the 16th ACM Conference on Recommender Systems (RecSys '22).* 3–13.

[19] Qinghua Liu, Andrew Henry Reiner, Arnoldo Frigessi, and Ida Scheel. 2019. Diverse personalized recommendations with uncertainty from implicit preference data with the Bayesian Mallows model. *Knowledge-Based Systems* 186 (2019), 104960.

[20] Gabriel Machado Lunardi, Guilherme Medeiros Machado, Vinicius Maran, and José Palazzo M de Oliveira. 2020. A metric for filter bubble measurement in recommender algorithms considering the news domain. *Applied Soft Computing* 97 (2020), 106771.

[21] Masoud Mansoury, Himan Abdollahpouri, Mykola Pechenizkiy, Bamshad Mobasher, and Robin Burke. 2020. FairMatch: A Graph-Based Approach for Improving Aggregate Diversity in Recommender Systems. In *Proceedings of the ACM International Conference on User Modeling, Adaptation and Personalization.* 154–162.

[22] Xiangfu Meng, Hongjin Huo, Xiaoyan Zhang, Wanchun Wang, and Jinxia Zhu. 2023. A survey of personalized news recommendation. *Data Science and Engineering* 8, 4 (2023), 396–416.

[23] Lien Michiels, Jens Leysen, Annelien Smets, and Bart Goethals. 2022. What Are Filter Bubbles Really? A Review of the Conceptual and Empirical Work. In *Adjunct Proceedings of the 30th ACM Conference on User Modeling, Adaptation and Personalization.* 274–279.

[24] Lien Michiels, Jorre Vannieuwenhuyze, Jens Leysen, Robin Verachtert, Annelien Smets, and Bart Goethals. 2023. How Should We Measure Filter Bubbles? A Regression Model and Evidence for Online News. In *Proceedings of the 17th ACM Conference on Recommender Systems.* 640–651.

[25] Ricardo S Oliveira, Caio Nóbrega, Leandro Balby Marinho, and Nazareno Andrade. 2017. A Multiobjective Music Recommendation Approach for Aspect-Based Diversification.. In *Proceedings of the International Society for Music Information Retrieval Conference.* 414–420.

[26] Eli Pariser. 2011. *The filter bubble: What the Internet is hiding from you.* penguin UK.

[27] Bashir Rastegarpanah, Krishna P. Gummadi, and Mark Crovella. 2019. Fighting Fire with Fire: Using Antidote Data to Improve Polarization and Fairness of Recommender Systems. In *Proceedings of the 12th ACM International Conference on Web Search and Data Mining (WSDM '19).* 231–239. https://doi.org/10.1145/3289600.3291002

[28] Shaina Raza, Syed Raza Bashir, and Usman Naseem. 2022. Accuracy meets Diversity in a News Recommender System. In *Proceedings of the 29th International Conference on Computational Linguistics.* 3778–3787.

[29] Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. 2009. BPR: Bayesian Personalized Ranking from Implicit Feedback. In *Proceedings of the International Conference on Uncertainty in Artificial Intelligence.* AUAI Press, 452–461.

[30] Manel Slokom, Alan Hanjalic, and Martha Larson. 2021. Towards user-oriented privacy for recommender system data: A personalization-based approach to gender obfuscation for user profiles. *Information Processing & Management* 58, 6 (2021).

[31] Barry Smyth and Paul McClave. 2001. Similarity vs. Diversity. In *Proceedings of the 4th International Conference on Case-Based Reasoning: Case-Based Reasoning Research and Development.* Springer-Verlag, 347–361.

[32] Dusan Stamenkovic, Alexandros Karatzoglou, Ioannis Arapakis, Xin Xin, and Kleomenis Katevas. 2022. Choosing the Best of Both Worlds: Diverse and Novel Recommendations through Multi-Objective Reinforcement Learning. In *Proceedings of the 15th ACM International Conference on Web Search and Data Mining.* 957–965.

[33] Celina Treuillier, Sylvain Castagnos, Evan Dufraisse, and Armelle Brun. 2022. Being Diverse is Not Enough: Rethinking Diversity Evaluation to Meet Challenges of News Recommender Systems. In *Adjunct Proceedings of the 30th ACM Conference on User Modeling, Adaptation and Personalization* (Barcelona, Spain). 222–233.

[34] Saúl Vargas and Pablo Castells. 2011. Rank and Relevance in Novelty and Diversity Metrics for Recommender Systems. In *Proceedings of the ACM International Conference on Recommender Systems.* 109–116.

[35] Sanne Vrijenhoek, Gabriel Bénédict, Mateo Gutierrez Granada, and Daan Odijk. 2023. RADio* – An Introduction to Measuring Normative Diversity in News Recommendations. *ACM Transaction on Recommender Systems* (2023).

[36] Sanne Vrijenhoek, Gabriel Bénédict, Mateo Gutierrez Granada, Daan Odijk, and Maarten De Rijke. 2022. RADio – Rank-Aware Divergence Metrics to Measure Normative Diversity in News Recommendations. In *Proceedings of the ACM International Conference on Recommender Systems.* 208–219.

[37] Lequn Wang and Thorsten Joachims. 2021. User Fairness, Item Fairness, and Diversity for Rankings in Two-Sided Markets. In *Proceedings of the ACM International Conference on Theory of Information Retrieval.* 23–41.

[38] Udi Weinsberg, Smriti Bhagat, Stratis Ioannidis, and Nina Taft. 2012. BlurMe: Inferring and Obfuscating User Gender Based on Ratings. In *Proceedings of the*

*ACM International Conference on Recommender Systems.* 195–202.

[39] Chuhan Wu, Fangzhao Wu, Mingxiao An, Jianqiang Huang, Yongfeng Huang, and Xing Xie. 2019. NPA: Neural News Recommendation with Personalized Attention. In *Proceedings of the ACM International Conference on Knowledge Discovery & Data Mining.* 2576–2584.

[40] Chuhan Wu, Fangzhao Wu, Suyu Ge, Tao Qi, Yongfeng Huang, and Xing Xie. 2019. Neural News Recommendation with Multi-Head Self-Attention. In *Proceedings of the ACM International Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (Eds.). 6389–6394.

[41] Chuhan Wu, Fangzhao Wu, Tao Qi, and Yongfeng Huang. 2022. End-to-end learnable diversity-aware news recommendation. *arXiv preprint arXiv:2204.00539* (2022).

[42] Fangzhao Wu, Ying Qiao, Jiun-Hung Chen, Chuhan Wu, Tao Qi, Jianxun Lian, Danyang Liu, Xing Xie, Jianfeng Gao, Winnie Wu, and Ming Zhou. 2020. MIND: A Large-scale Dataset for News Recommendation. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (Eds.). Association for Computational Linguistics, 3597–3606.

[43] Shangyuan Wu, Edson C Tandoc Jr, and Charles T Salmon. 2019. Journalism reconfigured: Assessing human–machine relations and the autonomous power of automation in news production. *Journalism studies* 20, 10 (2019), 1440–1457.

[44] Yao Wu, Jian Cao, and Guandong Xu. 2023. Fairness in Recommender Systems: Evaluation Approaches and Assurance Strategies. *ACM Transaction on Knowledge Discovery Data* 18, 1, Article 10 (2023), 37 pages.

[45] Meike Zehlike, Francesco Bonchi, Carlos Castillo, Sara Hajian, Mohamed Megahed, and Ricardo Baeza-Yates. 2017. FA*IR: A Fair Top-k Ranking Algorithm. In *Proceedings of the ACM on Conference on Information and Knowledge Management.* 1569–1578.

[46] Meike Zehlike, Tom Sühr, Ricardo Baeza-Yates, Francesco Bonchi, Carlos Castillo, and Sara Hajian. 2022. Fair Top-k Ranking with multiple protected groups. *Information Processing & Management* 59, 1 (2022), 28 pages.

[47] Yuying Zhao, Yu Wang, Yunchao Liu, Xueqi Cheng, Charu Aggarwal, and Tyler Derr. 2023. Fairness and diversity in recommender systems: a survey. *arXiv preprint arXiv:2307.04644* (2023).

[48] Yu Zheng, Chen Gao, Liang Chen, Depeng Jin, and Yong Li. 2021. DGCN: Diversified Recommendation with Graph Convolutional Networks. In *Proceedings of the Web Conference 2021.* Association for Computing Machinery, 401–412.

[49] Cai-Nicolas Ziegler, Sean M. McNee, Joseph A. Konstan, and Georg Lausen. 2005. Improving recommendation lists through topic diversification. In *Proceedings of the 14th International Conference on World Wide Web.* 22–32.