# Inductive Conformal Recommender System

Venkateswara Rao Kagita[a], Arun K Pujari[b], Vineet Padmanabhan[c], Vikas Kumar[d,*]

[a]*National Institute of Technology, Warangal, India*
[b]*Mahindra University, Hyderbad, India*
[c]*University of Hyderabad, Hyderbad, India*
[d]*University of Delhi, Delhi, India*

## Abstract

Traditional recommendation algorithms develop techniques that can help people to choose desirable items. However, in many real-world applications, along with a set of recommendations, it is also essential to quantify each recommendation's (un)certainty. The conformal recommender system uses the experience of a user to output a set of recommendations, each associated with a precise confidence value. Given a significance level $\varepsilon$, it provides a bound $\varepsilon$ on the probability of making a wrong recommendation. The conformal framework uses a key concept called *nonconformity measure* that measures the strangeness of an item concerning other items. One of the significant design challenges of any conformal recommendation framework is integrating nonconformity measures with the recommendation algorithm. This paper introduces an inductive variant of a conformal recommender system. We propose and analyze different nonconformity measures in the inductive setting. We also provide theoretical proofs on the error-bound and the time complexity. Extensive empirical analysis on ten benchmark datasets demonstrates that the inductive variant substantially improves the performance in computation time while preserving the accuracy.

Keywords: Recommender System, Inductive Conformal Prediction, Conformal Recommender System

## 1. Introduction

Recommending quality services to improve customer satisfaction is of prime concern for the overall success of any online community. In this context, an automatic recommendation has become even more indispensable. Recommender systems are software tools that use past behaviour (usage information) of individuals to provide personalized recommendations for a large variety of available products such as movies, books, music, etc. There have been numerous research proposals on recommendation problem focusing on improving recommendation accuracy [1, 2, 3]. With the upcoming importance on accountability and

---

explainability of AI techniques, deployment of a plain recommendation whatsoever accurate it may be on a testing platform will not be satisfactory without value additions such as explanations, confidence, or sensitivity. Among the desirable features of the future of recommender systems, providing a confidence measure (or, equivalently, a probable error bound) on recommendation is essential. Most of the existing recommender systems do not offer any such measure to indicate the level of confidence till very recently when the present authors propose *Conformal Recommender Systems (CRS)* [4]. Though some of the earlier systems endeavour to provide confidence [5, 6, 7], the confidence values so provided are not related to or bound to the error values. On the other hand, *Conformal Recommender Systems (CRS)* satisfies a *validity* property that ensures that the error value does not exceed a predetermined significance level $\varepsilon$. In other words, the correctness-confidence of a recommendation is 1-$\varepsilon$. It is observed that though CRS is an advancement in research in the area of Recommender Systems, the underlying process is computationally intensive and expensive. Having established the point that a *valid* quantitative measure of confidence can be computed using the principles of conformal prediction, the need arises to provide a computationally efficient method of accomplishing this task. The objective of the present work is to investigate efficient alternative techniques retaining the strength of CRS. This paper proposes an *inductive* variant of a conformal recommender system that is computationally efficient and retains the validity property of CRS.

For a set of training examples $S = \{z_1, \ldots, z_n\}$, where $z_i$ is a pair $(x_i, y_i)$ with $x_i \in \mathbb{R}^d$ is a vector of $i^{th}$ example and $y_i$ is the corresponding class label, any common predictor predicts a class label $y_{n+i}$ for unclassified objects $x_{n+i}, i \geq 1$. In contrast, conformal predictors give a set of class labels as prediction regions and corresponding probability-bounds of error. A (1-$\varepsilon$) confidence prediction region is defined by the probability that the correct label is not in the prediction region does not exceed $\varepsilon$. To predict the class label of an unclassified object, say $x_{n+1}$ to one of the class labels, say $y_c$, based on the available information in terms of $S$, conformal predictors define suitable numeric measure to compute a nonconformity measure for each pair of training example and class-label. Intuitively, it is a measure of how well an unknown data $x_{n+1}$ conforms to any training example $x_i$ when $x_{n+1}$ is assigned class label $y_c$. This is done by measuring the change in predicting behaviour of $S$ when $z_i$ is replaced by $z_{n+1}$. The predicting behaviour is observed by applying any of the conventional predictors. The nonconformity measures for all such pairs are analysed to compute $p$-values and then to determine (1-$\varepsilon$) region subsequently.

Two important observations can be made from the foregoing discussion. First, the process is hinged on the definition of suitable nonconformity measure. We observe that depending on the context, it is sometimes easier to use a conformity measure instead of a nonconformity measure, though both processes are equivalent intuitively. For the sake of notational convenience, we use the term nonconformity measure to refer to both situations. Second, the measure is required to be computed for all pairs of $x_i$ and $y_c$ in order to determine the p-values. In order to show the relevance and applicability of the principle of conformal prediction, a nonconformity measure is introduced by Kagita et al. [4] in the context of recommender systems using precedence information. Based on the rating data of a set of users on a set of items, a nonconformity measure is calculated for all possible rec-

ommendations by examining how well the tentative recommendation conforms to all other known recommendations and earlier ratings for any user. The underlying algorithm uses precedence mining as proposed in [8]. A different nonconformity measure is defined in [9], wherein the matrix factorization is used as the underlying algorithm.

The main contributions of the present work are as follows. First, we analyze different possibilities of defining nonconformity measures in the context of inductive CRS by using the *precedence relations* among items. As stated earlier, defining a conformity measure is observed to be more relevant than a nonconformity measure in some situations. We adopt different probability measures using pairwise precedence statistics characterized by Parameswaran et al. [8] for defining suitable (non)conformity measures. Second, we introduce the concept of *inductive conformal recommendation*, which is a computationally efficient alternative to the CRS framework. Further, we theoretically and empirically establish the crucial properties of the conformal framework, i.e., validity and efficiency. To verify its efficacy, we conducted extensive experiments on seven bench-mark datasets using various standard evaluation metrics. We show that the proposed inductive conformal recommender system improves the computational time while retaining a similar level of accuracy.

The rest of the paper is structured as follows. In Section 2, we briefly discuss the related work. Section 3 describes the key concepts required to build the proposed system. Section 3.1 presents the background on conformal prediction framework. In Section 3.2, we discuss the underlying precedence mining based recommender systems. We discuss the existing conformal recommender system in Section 4. We introduce the proposed inductive conformal recommender system in Section 5. We report experimental results in Section 6. Finally, Section 7 concludes and indicates several issues for future work.

## 2. Related Work: Confidence Measure in Recommender System

Recommender systems are generally employed to provide tailor-made suggestions that can assist the user in decision making [10, 11]. These systems exploit the user's consumption experience collected via implicit or explicit feedback data to infer their preferences [12, 13, 14]. However, most of these systems are less transparent because of the unavailability of confidence with which an item is recommended [9, 15]. Despite the enormous application of recommender systems, a limited number of methods are available that associate confidence value with the recommendations. In this section, a brief review of the earlier work concerning confidence measures in the recommender system is presented. Readers' familiarity with recommender system is assumed here.

To measure the effect of confidence and uncertainty measures, McNee et al. [16] involved an elementary confidence computation in existing collaborative filtering algorithms and have shown that a confidence display increases user satisfaction. In [7], the authors have considered the previously collected user's rating as noisy evidence of the user's actual rating and proposed a Belief Distribution algorithm that explicitly outputs the uncertainty in each predicted rating along with the predicted rating value. Adomavicius et al. [17] proposed a rating variance-based confidence measure to refine the prediction generated by any traditional recommendation algorithm. Symeonidis et al. [18] constructed a feature profile of

each user, and then the prediction is justified by considering the correlation between users and features. Shani et al. [19] suggested measuring the significance level of recommendation by running a significance test between the results of different recommender algorithms. OrdRec [20] provides a richer expressive power by producing a full probability distribution of the expected item ratings. Mazurowski [5] compared the concept of confidence in collaborative filtering with similar concepts in other fields within machine learning. Bayesian confidence intervals-based evaluation method has been proposed to measure recommendation algorithms' performance. The author also proposed three different resampling-based methods to estimate the confidence of individual predictions [5]. In [6], for a target user, the confidence in prediction for an item is defined based on $k$-nearest neighbors of the user. In [15], a content-based fuzzy recommendation model is proposed that utilize *similarity* and *dissimilarity* score between user and item for the rating prediction task. For every unknown (user, item)-pair, the prediction confidence is computed based on the difference between the actual ratings given by that user and their corresponding predictions by the fuzzy model. A recommendation model is proposed in [21] to integrate the trust and certainty information for confidence modeling. Mesas et al. [22] explored the prediction confidence from the perspective of the system. The idea is to embed awareness into the recommendation models that help in deciding the more reliable suggestions rather than all potential recommendations. A Course Recommender system is proposed in [23], where a course-specific regression model is trained over the course contents and students' academic interests for the grade predictions. To complement the model predicted grades, the authors have employed an Inductive Confidence Machine (ICM) [24] to construct prediction intervals attune with each student. In [9], two variants of conformal framework, namely transductive and inductive, are proposed in the matrix factorization (MF) setting that associate a confidence score to each predicted rating. The method proposed in [9] can be seen as a two-stage procedure. At first, a MF-based model is applied over the partially filled rating matrix to get the rating prediction for each (user, item)-pair. These predictions are then used to calculate the confidence score for individual predicted ratings. A confidence-aware MF model is proposed in [25], which can be seen as a comprehensive framework that optimizes the accuracy of rating prediction and estimates the confidence over predicted rating simultaneously. Costa et al. [26] proposed an ensemble-based co-training approach for the rating prediction problem. In the co-training phase, two or more recommender algorithms are trained to predict the rating for all unobserved user-item pairs. The training set for the next iteration of the co-training is then augmented with the $M$ most confident predictions. The confidence is calculated based on the deviation from the baseline estimate and the rating predicted by the recommendation algorithm. *However, none of these works provide confidence to the recommendation set. They focus on providing confidence to the individual rating prediction, and it is non-trivial and cumbersome to obtain the confidence of recommendation from confidence regions of point predictions.*

In this work, we focus on providing confidence to the recommendation, not for rating prediction. The only work that focuses on providing confidence to the recommendation is our previous work on conformal recommender system [4], wherein a conformal framework is introduced for the recommender systems, and a new nonconformity measure is proposed for

the conformal recommender system. It is also shown that the proposed nonconformity satisfies the desirable properties of conformal prediction, such as *exchangeability*, *validity*, and *efficiency*. Nonetheless, the framework proposed in [4] suffers from similar shortcomings of traditional conformal predictions and requires high computation times. We briefly describe the approach in the following section.

## 3. Foundational Concepts

In this section, we first introduce the basic concepts related to *conformal prediction*, the main framework we use to build our proposed confidence-based recommender system. We then give a brief description of *precedence mining*, a collaborative filtering model, on which we apply our conformal prediction framework for producing confidence-based recommendations.

### 3.1. Conformal Prediction

In this section, a brief account of the principle of conformal prediction is reported in order to provide the relevant background. We start with a training example $z_i$ as a pair $(x_i, y_i)$ with $x_i \in \mathbb{R}^d$ is a feature vector of $i^{th}$ example and $y_i$ is the corresponding class label. Given the training set $S = \{z_1, \ldots, z_n\}$, a prediction or classification task is to predict a class label $y_{n+i}$ for unclassified objects $x_{n+i}, i \geq 1$. A conformal predictor provides a subset of class labels for each unclassified object $x_{n+i}$ and the error that the correct label is not in this set does not exceed $\varepsilon$. Let us consider one unclassified object $x_{n+1}$ and the task is to examine whether a class label $y_c$ is a member of $(1 - \varepsilon)-$prediction region. Let $z_{n+1}^c = (x_{n+1}, y_c)$, where $y_c$ is tentatively assigned to $x_{n+1}$. The nonconformity measure for an example $z_i \in \{S \cup z_{n+1}^c\}$ is a measure of how well $z_i$ conforms to $\{S \cup z_{n+1}^c\} \setminus z_i, \forall i \in [1, n+1]$. From another point of view, it can be seen as a measure of how well $z_{n+1}^c$ conforms to $z_i \in S$. This is done by measuring the change in predicting behaviour of $S$ when $z_i$ is replaced by $z_{n+1}^c$. The $p$-value is the proportion of $z_i \in S$ having nonconformity score worse than that of $z_{n+1}^c$ for all possible values of $y_c$ (all class labels). The set of labels whose $p$-value higher than $\varepsilon$ forms $(1 - \varepsilon)-$prediction region. Intuitively, the predicting behaviour is observed by applying any of the conventional predictors which uses S as the training set. The conformal prediction algorithm makes $(n+1) \times n_c \times C$ calls to the underlying prediction algorithm, where $n_c$ is the number of candidate items, and $C$ is the number of class labels. The conformal prediction framework has been well-studied from different perspectives in recent years [27, 28, 29, 30].

On the other hand, the *inductive conformal framework* avoids the computational overhead [27] of initial proposal of conformal prediction. In an inductive setting the training set $S = \{z_1, \ldots, z_n\}$ is divided into two sets, namely *proper* training set $S^t = \{z_1, z_2, \ldots, z_m\}$ and *calibration* set $S^c = \{z_{m+1}, \ldots, z_{m+l}\}$, $n = m + l$. The former is used to learn the prediction model and the latter is used for computation of $p$-values. The system uses an underlying conventional prediction algorithm to learn a model using proper training set $S^t$. The same model is then used to determine (non)conformity measures and $p$-value for every example in $S^c$ and $z_{n+1}$ with respect to $S^t$. As a result, the framework learns the underlying model only once, leading to a significant reduction in computation time and effort.

**Example 1.** *Consider a problem of classifying samples as cancerous (+ve) or noncancerous (−ve) based on the tumor size and other pathological features. Let $x_i$ be the feature vector describing an $i^{th}$ instance, and $y_i \in \{+ve, -ve\}$ is the corresponding label. Let $S = \{(x_1, y_1), (x_2, y_2), \ldots, (x_{10}, y_{10})\}$ be the training set containing observations of ten different individuals. Given the new instance, say $x_{11}$, the task is to classify it as either +ve or −ve. Assume, Support Vector Machine (SVM) is the underlying classifier that also assigns a nonconformity value as the deviation between the actual and the predicted class label. The conformal prediction framework labels the new instance with all possible classes and sees which conforms more to the existing ones. At first, it considers $z_{11}^{+ve} = (x_{11}, +ve)$ and adds it to the training set. After that the the SVM classifier is trained for each $i^{th}$ instance with the training set $\{S \cup z_{11}^{+ve}\} \setminus z_i, \forall i \in [1, 11]$ and measures the removed example's nonconformity ($\alpha_i$). Finally, it calculates the p-value concerning the label C (let's say $P^{+ve}$), $P^{+ve} = \frac{|i|1 \leq i \leq 11, \alpha_i \geq \alpha_{11}|}{11}$. The conformal predictor repeats the same procedure concerning another label $NC$, i.e., for $(x_{11}, -ve)$ and determines the corresponding p-value ($P^{-ve}$). The prediction region includes the labels with a p-value greater than the significance level $\varepsilon$. We can also observe that for the given example, the conformal predictor requires training of 22 $(11 \times 2 = (n + 1) \times C)$ SVMs for one candidate item and hence, for $n_c$ candidate items it would be $(n + 1) \times C \times n_c$ SVMs.*

*On the other hand, the Inductive Conformal Predictor divides the dataset into two sets, namely proper training set $S^t$ and calibration set $S^c$. It then trains the underlying model with $S^t$ and uses the same to evaluate the nonconformity of $S^c$ and new instance $x_{11}$. Hence, only one SVM classifier is learnt using the set $S^t$ and called for $|S^c| + C$ times. For $n_c$ candidates it is $n_c \times (|S^c| + C)$, which is a drastic improvement over conformal predictor in terms of computation time.*

### 3.2. Precedence Mining based Recommender Systems [8]

The precedence mining model [8] is a Collaborative Filtering (CF) based model that maintains precedence statistics, i.e., the temporal count of all the pairs of items. The precedence mining model estimates the probability of future consumption based on past behaviour. For example, a person who has seen *Godfather I* is more likely to watch *Godfather II* in the future. In most of the traditional CF techniques, the aim is to find users having similar profiles as the active user $u$, and then restrict its search to items consumed by this subset of users and not consumed by $u$. Thus, certain consumption patterns of items exhibited by the whole set of users are not captured as the search is restricted. The precedence mining model overcomes these shortcomings and attempts to capture pairwise precedence relations frequently occurring among all users. It calculates a recommendation score for each item based on the precedence statistics, and then the set of items having scores greater than the threshold are recommended.

**Example 2.** *Figure 1 shows the difference between traditional collaborative filtering and precedence mining. The leftmost table in the figure is a toy example in which we provide the profiles of different users. Let $u_a$ denote the active user $u_a$. Each row in the table can be interpreted as a sequence of movies that the user has watched. For instance, $u_a$ has watched*

6

movies $m_1$, $m_2$, and $m_3$ in the given order. The figure in the middle demonstrates the working of collaborative filtering. Here, we assume that the set of users who have at least two movies in common with the active user are in its neighbours. The most popular movie among the neighbours are then recommended to the active user. A careful observation of Figure 1 reveals that the movie $m_5$ is popular among the neighbours $u_1$ and $u_2$ and therefore collaborative filtering recommends movie $m_5$ to the user $u_a$. It is to be noted that the search of items in collaborative filtering is limited to the neighbours space. In contrast, precedence mining looks for patterns in which one item follows the other in the whole user space. The rightmost image in the Figure 1 demonstrates the idea of precedence mining. We highlight the patterns that occur at least thrice using different colors, for instance, $m_1$ and $m_7$.

| User | Order of Consumption | | | | |
|---|---|---|---|---|---|
| Ua | m1 | m2 | m3 | | |
| U1 | m1 | m7 | m2 | m5 | |
| U2 | m1 | m3 | m6 | m5 | |
| U3 | m3 | m4 | m6 | m9 | |
| U4 | m1 | m4 | m7 | m9 | |
| U5 | m8 | m1 | m7 | m4 | m9 |
| U6 | m3 | m7 | m6 | | |

| User | Order of Consumption | | | | |
|---|---|---|---|---|---|
| Ua | m1 | m2 | m3 | | |
| U1 | m1 | m7 | m2 | m5 | |
| U2 | m1 | m3 | m6 | m5 | |
| U3 | m3 | m4 | m6 | m9 | |
| U4 | m1 | m4 | m7 | m9 | |
| U5 | m8 | m1 | m7 | m4 | m9 |
| U6 | m3 | m7 | m6 | | |

| User | Order of Consumption | | | | |
|---|---|---|---|---|---|
| Ua | m1 | m2 | m3 | | |
| U1 | m1 | m7 | m2 | m5 | |
| U2 | m1 | m3 | m6 | m5 | |
| U3 | m3 | m4 | m6 | m9 | |
| U4 | m1 | m4 | m7 | m9 | |
| U5 | m8 | m1 | m7 | m4 | m9 |
| U6 | m3 | m7 | m6 | | |

Figure 1: Comparison of collaborative filtering (middle) and precedence mining (right) approaches for a toy example (left).

Recommender systems based on precedence relations is concerned with mining precedence relations among items consumed by users and thereafter recommends new items having high *relevance score* computed using precedence statistics. The nicety and novelty of this approach is the use of pairwise precedence relations between items. We describe the score computation formally as follows.

Let $O = \{o_1, o_2, \ldots, o_{mobject}\}$ be the set of items and $U = \{u_1, u_2, \ldots, u_{muser}\}$ be the set of users. $profile(u_j)$ is a sequence of items known to have been consumed by user $u_j$. Let $O_j$ be the set of items consumed by $u_j$. A recommender system is concerned with recommending items to a user based on profiles of different users. A recommender system aims at selecting items for recommendation such that these items are absent in $profile(u)$ and are expectedly preferred to other items by the user for whom it is recommended. Let $Support(o_i)$ be the number of users that have consumed item $o_i$ and $PrecedenceCount(o_i, o_h)$ be the number of users having consumed item $o_i$ preceding $o_h$. The precedence probability for item $o_i$ preceding $o_h$ is denoted as $PP(o_i|o_h)$. We define $PP(o_i|o_h)$, and $Score(o_i, u_j)$ as follows.

$$PP(o_i|o_h) = \frac{PrecedenceCount(o_i, o_h)}{Support(o_h)}, \tag{1}$$

$$Score(o_i, u_j) = \frac{Support(o_i)}{muser} \times \prod_{o_l \in O_j} PP(o_l|o_i). \tag{2}$$

7

The objects with high score are recommended. If the score for a given unutilized item is low, it is highly unlikely to be of interest to the user. We now consider an example which illustrates the working of the preceding precedence mining based recommender system.

**Example 3.** *Consider the following* PrecedenceCount *and* Support *statistics calculated based on the preferences given by thirty users* $U = \{u_1, u_2, \ldots, u_{30}\}$ *over ten items* $O = \{o_1, o_2, \ldots, o_{10}\}$.

$$PrecedenceCount = \begin{bmatrix} 0 & 9 & 8 & 11 & 7 & 8 & 6 & 7 & 7 & 3 \\ 8 & 0 & 10 & 11 & 9 & 7 & 7 & 6 & 8 & 4 \\ 8 & 8 & 0 & 5 & 7 & 6 & 5 & 6 & 4 & 3 \\ 5 & 11 & 12 & 0 & 6 & 8 & 6 & 6 & 3 & 2 \\ 7 & 9 & 7 & 13 & 0 & 8 & 7 & 6 & 9 & 4 \\ 5 & 9 & 6 & 8 & 6 & 0 & 5 & 6 & 4 & 1 \\ 4 & 6 & 5 & 7 & 5 & 5 & 0 & 4 & 4 & 1 \\ 4 & 8 & 6 & 10 & 8 & 6 & 5 & 0 & 7 & 2 \\ 7 & 11 & 10 & 16 & 7 & 7 & 6 & 5 & 0 & 3 \\ 2 & 1 & 0 & 2 & 1 & 2 & 1 & 3 & 1 & 0 \end{bmatrix}$$

$$Support = \begin{bmatrix} 20 & 25 & 21 & 25 & 22 & 18 & 15 & 18 & 20 & 6 \end{bmatrix}.$$

*Let $u_1$ be the target user and $O_1 = \{o_1, o_3, o_5, o_7, o_9\}$ be set of items consumed by $u_1$. For $u_1$, the candidate items for recommendation are $O \setminus O_1 = \{o_2, o_4, o_6, o_8, o_{10}\}$. The score of an item $o_2$ which not consumed by user $u_1$ is then calculated as*

$$Score(o_2, u_1) = \frac{Support(o_2)}{30} \times PP(o_1 \mid o_2) \times PP(o_3 \mid o_2) \times PP(o_5 \mid o_2) \times PP(o_7 \mid o_2) \times PP(o_9 \mid o_2)$$

$$= \frac{25}{30} \times \frac{9}{25} \times \frac{8}{25} \times \frac{9}{25} \times \frac{6}{25} \times \frac{11}{25} = 0.0036.$$

*Similarly, $Score(o_4, u_1) = 0.0068$, $Score(o_6, u_1) = 0.0043$, $Score(o_8, u_1) = 0.0016$, and $Score(o_{10}, u_1) = 0.0028$. Hence, it ranks the items in the order of $o_4$, $o_6$, $o_2$, $o_{10}$, and $o_8$.*

The problem with this approach is even if one of the precedence probabilities (PPs) is zero, the whole score becomes zero. To avoid this problem, Parameswaran et al. [8] proposed to consider only top-I precedence probabilities in the product term, where $I$ is a hyper-parameter to tune. In our experiments, we have fine tune the value $I$ to be 1.

## 4. Conformal Recommender System

The principle of conformal prediction is applied to recommender system in [4]. Here, we briefly report the proposal of CRS. The readers are requested to refer [4] for details. Let $O$ be the set of items, $n_i = |O|$ be the total number of items, $n_u$ be the number of users and $O_j = \{o_1, o_2, \ldots, o_n\}$ be the set of items consumed by a user $u_j$. Given $O$, $u_j$, $O_j$, and the significance level $\varepsilon$, the problem is to recommend a set of items $\Gamma^\varepsilon$ with $(1 - \varepsilon)$ confidence.

For a given user $u_j$, $O_j$ is split into two sets based on the precedence of usage of the items. The first set $O_j^{train} = \{o_1, o_2, \ldots, o_n\}$ is used as the training set. The set $O_j^{candidates} = \{c_1, c_2, \ldots, c_k\}$, of candidate items that are consumed by $u_j$ after use of items in $O_j^{train}$ and are not part of the training set. The conformal recommendation process is to determine the confidence measure of recommending a new object $o_{n+1}$ for user $u_j$. The first step of CRS is to see how well the object $o_{n+1}$ and the training set $O_j^{train}$ conform to each other. Let $O_j^{train+} = \{O_j^{train} \bigcup o_{n+1}\}$ be the appended set. Nonconformity measure is computed for each $o_h \in O_j^{train+}$ by ignoring $o_h$ in $O_j^{train+}$ and examining the recommendability of $c_i$ when the profile is $O_j^h$, where

$$O_j^h = O_j^{train+} \setminus \{o_h\} = \{o_1, o_2, \ldots, o_{h-1}, o_{h+1}, \ldots, o_{n+1}\}. \tag{3}$$

With precedence mining [31, 32, 33] as the *underlying algorithm*, the measure of recommendability of $c_i$ is a numerical value, $Score(c_i, u_j)$ and higher value of $Score$ implies greater chance of being recommended. The $Score$ is calculated with reference to each of $h$ and for each tentative profile $O_j^h$, $Score^h$ is defined as

$$\alpha_h = Score^h(c_i, u_j) = \frac{Sup(c_i)}{m_{user}} \times \prod_{o_l \in O_j^h} PP(o_l | c_i).$$

**Definition 1.** *(CRS nonconformity measure [4]). Given a subset $O_j^{train}$ of user $u_j$ profile; a set of objects $O_j^{candidates} = \{c_1, c_2, \ldots, c_k\}$, that are consumed by $u_j$ after use of items in $O_j^{train}$ and are not part of the training set; and a new object $o_{n+1} \in O_j$, the nonconformity measure $\mathcal{A}(o_1, o_2, \ldots, o_{n+1})$ w.r.t. $c_i \in O_j^{candidates}$ is $(\alpha_1, \alpha_2, \ldots, \alpha_{n+1})$, where $\alpha_h = Score^h(c_i, u_j)$.*

The computed nonconformity scores $\alpha_h, h \in [1, n+1]$ are used to compute the $p$-value as the proportion of examples with $\alpha_h \geq \alpha_{n+1}, h \in [1, n+1]$. A $p$-value is computed for each $c_i \in O_j^{candidates}$ and then we employ two different aggregation techniques to select the final $p$-value from several $p$-values. If the selected $p$-value is greater than $\varepsilon$, then $o_{n+1}$ is included in the $(1 - \varepsilon)$ confidence recommendation region. The procedure is repeated for every new item $o_{n+i}, i \geq 1$ to get $(1 - \varepsilon)$ confidence recommendation set.

**Example 4.** *We consider the precedence statistics given in Example 3 for this example also. Let $O_1^{train} = \{o_1, o_3, o_5\} \subset O_1$ and $O_1^{candidates} = \{o_7, o_9\}$. Let $o_2$ be the candidate item for recommendation. We append $o_2$ with $O_1^{train}$ as $O_1^{train+} = \{o_1, o_3, o_5, o_2\}$. Nonconformity of an item $o_h \in O_1^{train+}$ is measured by the recommendability of a candidate item $c \in O_1^{candidates}$ using the profile $O_1^{train+} \setminus \{o_h\}$. For example, nonconformity of an item $o_1$ concerning the recommendability of $o_7$ is computed as*

$$\alpha_1 = Score^1(o_7, u_1) = \frac{Support(o_7)}{30} \times PP(o_3 \mid o_7) \times PP(o_5 \mid o_7) \times PP(o_2 \mid o_7)$$
$$= \frac{15}{30} \times \frac{5}{15} \times \frac{7}{15} \times \frac{7}{15} = 0.036.$$

*Nonconformity score of $o_3$ is*

$$\alpha_3 = Score^3(o_7, u_1) = \frac{Support(o_7)}{30} \times PP(o_1 \mid o_7) \times PP(o_5 \mid o_7) \times PP(o_2 \mid o_7) == 0.043.$$

*Similarly, Nonconformity score of $o_5$ is $\alpha_5 = Score^5(o_7, u_1) = 0.031$ and nonconformity score of $o_2$ is $\alpha_2 = Score^2(o_7, u_1) = 0.031$. The p-value of $o_2$ concerning the recommendability of $o_7$ is computed as follows.*

$$p(o_2, o_7) = \frac{\left| \left\{ o_h \middle| o_h \in O_1^{train+}, Score^h(o_7, u_1) \geq Score^2(o_7, u_1) \right\} \right|}{|O_1^{train+}|} = \frac{4}{4} = 1,$$

*Similarly, we compute the p-value of $o_2$ concerning the recommendability of $o_9$ that is $p(o_2, o_9) = 0.75$. In order to get the final p-value from $p(o_2, o_7)$ and $p(o_2, o_9)$, CRS-max [4] employs a maximum strategy and CRS-med [4] employs a median strategy. Therefore, the final p-value according to CRS-max and CRS-med are 1 and 0.875 respectively. Similarly, we compute the p-value for all the candidate items for recommendation and recommend the items whose p-value is greater than $\varepsilon$ with the confidence of $(1 - \varepsilon)$.*

## 5. Inductive Conformal Recommender System

This section presents the proposed *inductive conformal recommender system (ICRS)* to gauge the confidence of recommendations. The proposed conformal approach determines a recommendation set $\Gamma^\varepsilon$ with $(1 - \varepsilon)$ confidence for a given significance level $\varepsilon$. A pivotal component of the conformal framework is the nonconformity measure quantifying the reliability in prediction. We use precedence relations among the items to determine the nonconformity score. Precedence relations capture the temporal patterns in user transactions. Besides, precedence relations based recommender systems do not require rating information, which is indeed challenging to obtain in a real-time scenario. Furthermore, these systems are ranking systems and thus allow us to define confidence for recommendation instead of a rating prediction. These are the various reasons for choosing precedence relations to represent nonconformity measures.

The brief idea of the proposed approach is as follows. We split $O_j$ into *proper training set* $O_j^t = \{o_1, o_2, \ldots, o_m\}$ and *calibration set* $O_j^c = \{o_{m+1}, o_{m+2}, \ldots, o_{m+l}\}$, wherein $O_j^c$ is the set of items known to be consumed after $O_j^t$ and $n = m + l$. The idea is to compute the (non)conformity measure for every item in the calibration set along with a new item $o_{n+1}$ and determine $o_{n+1}$'s *p-value*: the proportion of items having (non)conformity score better than or equal to that of a new item. Subsequently, we include item $o_{n+1}$ in the $\Gamma^\varepsilon$ recommendation region if the p-value of $o_{n+1}$ exceeds $\varepsilon$.

The following subsections elaborate on the notions of (non)conformity measures and the p-value and describe the complete procedure. Subsection 5.1 defines the various (non)conformity measures to determine the conformity or strangeness of an object concerning the training set. Subsection 5.2 defines *p*-value, which quantifies the conformity

score of a new item concerning the training set of items and defines the recommendation set $\Gamma^\varepsilon$ with $(1 - \varepsilon)$ confidence. Subsection 5.3 gives the flowchart of the proposed system and describes the proposed algorithm. In Subsection 5.4, we describe the two important measures of any conformal prediction framework, *validity* and *efficiency*, in the recommender systems setting. Finally, we proffer theoretical time complexity analysis of the proposed approach against the existing methods in Subsection 5.5.

## 5.1. Nonconformity Measures

Nonconformity measure is a measurable function $\mathcal{A}$ that determines a new object's relation with the proper training set in terms of a scalar value. There are several ways traditional algorithms can construct nonconformity measures; each of these measures defines a unique ICRS. It is worth mentioning that a particular (non)conformity measure only affects the ICRS model's efficiency, and the validity of the results remains unaffected. We propose different conformity/nonconformity measures in this section and analyze the efficiency. We use *precedence count* $PC(o_i, o_h)$ and *precedence probability* $PP(o_i|o_h)$ that determines the precedence relation among items to define various (non)conformity measures. When we compute these quantities for each item in the user profile, we get multiple values. We use different aggregation techniques as a design choice to calculate the (non)conformity value using multiple precedence statistics. For the simplicity of notations, we refer to conformity measure as *CM* and nonconformity measure as *NCM* in the subsequent discussion.

We adapt the score function proposed by Parameswaran et al. [8] that estimates the relevance of an item to the user profile to establish the first conformity measure. We define the conformity score of an item $o_h$ for a user $u_j$ profile as follows.

$$CM1(o_h) = \frac{Sup(o_h)}{n_u} \times \prod_{o_l \in O_j^t}^{(\text{I})} PP(o_l|o_h),$$

where $\prod^{(I)}$ denotes the multiplication of top-I quantities in the product term. We validate the algorithm for different I values and take $I$ as 1 in the experiments. The score is high when it conforms more to the training set. Note that every measure that we define here is with respect to a target user $u_j$. Furthermore, we determine an object's conformity in terms of the precedence count of an item with the set of items consumed by the user. The precedence count $(PC(o_i, o_h))$ defined previously represents the number of times an item $o_h$ appeared after $o_i$ in user profiles. The higher the number, the more likely it is that $o_h$ appears after $o_i$. Hence, we use precedence count to determine a conformity measure. We compute the precedence count of an item $o_h$ to every item $o_i$ in the proper training set $O_j^t$ of user $u_j$ and then aggregate them to get a numerical score. Using the different aggregation strategies such as *minimum*, *median*, *mean* and *maximum*, we arrive at the following conformity measures: *CM2, CM3, CM4,* and *CM5*, respectively. The detailed formulation of these measures is given in Annexure 1. We also use the precedence probability of an object with respect to the user profile to determine the conformity score of an object. Precedence probability $PP(o_h \mid o_i)$ of an item $o_h$ with respect to an item $o_i$ indicates how likely an item $o_h$ follows an item $o_i$. Hence, we use precedence probabilities of an item $o_h$ with respect to individual items in the

user profile to define the conformity measures. We again use different aggregation strategies to summarize the precedence probability scores of $o_h$ with respect to multiple items in the user profile. The process resulted in four different conformity scores, *CM6 (minimum)*, *CM7 (median)*, *CM8(mean)*, and *CM9(maximum)* with the corresponding aggregation operator mentioned in the brackets. The detailed formulation is given in Annexure 1. We then employ probability of $o_h$ given that $o_i$ is present in the target user profile without preceding $o_h$ to determine the conformity score of an item $o_h$ concerning the training data. We compute this score with respect to each and every item in the training data and employ different aggregation strategies resulting in four different conformity scores, *CM10 (minimum)*, *CM11 (median)*, *CM12(mean)*, and *CM13(maximum)*. Finally, we consider the probability that an item $o_i$ appears in the profile without succeeding an item $o_h$ ($\frac{Sup(o_i)-PC(o_i,o_h)}{n_u}$) as the potential nonconformity measure for an item $o_h$. Since there are multiple $o_i$'s in the user profile/training set, we use different aggregation strategies and define the nonconformity measures *NCM14*, *NCM15*, *NCM16*, and *NCM17* as given in Annexure 1.

**Lemma 1.** *(Non)conformity of items $\{o_{m+1}, \ldots, o_{n+1}\}$ is invariant of permutation, i.e., for any permutation $\pi$ of $\{m+1, \ldots, n+1\}$ i.e., $\mathcal{A}(o_{m+1}, o_{m+2}, \ldots, o_{n+1}) = (\alpha_{m+1}, \alpha_{m+2}, \ldots, \alpha_{n+1})$ $\Rightarrow \mathcal{A}(o_{\pi(m+1)}, o_{\pi(m+2)}, \ldots, o_{\pi(n+1)}) = (\alpha_{\pi(m+1)}, \ldots, \alpha_{\pi(n+1)})$.*

*Proof.* It is easy to see that the nonconformity scores are invariant of permutation $\pi$ of $\{o_{m+1}, \ldots, o_{n+1}\}$. All the proposed conformity/nonconformity measures are independent of the calibration set $\{o_{m+1}, \ldots, o_{n+1}\}$ and only makes use of the proper training set. Hence changing the permutation of a calibration set does not effect the nonconformity scores and remains the same. Therefore the proposed (non)conformity scores are invariant of permutation of $\{o_{m+1}, \ldots, o_{n+1}\}$. $\square$

---

**Algorithm 1:** Inductive Conformal Recommender Systems.

---

   **Input:** $O$, *target user* $u_j$, $O_j$, $\varepsilon$
   **Output:** Recommendation set ($\Gamma^\varepsilon$)
   split $O_j$ into two sets $O_j^t$ and $O_j^c$;
   $\Gamma^\varepsilon \leftarrow \emptyset$ ;
   **for** *each $o_h$ in $O_j^c$* **do**
      | Compute $\alpha_h$ using any of the (non)conformity measures;
   **end**
   **for** *each $o \in O \setminus O_j$* **do**
      | Compute (non)conformity score of an item $o$;
      | Compute $p(o)$ using Equation 4 or 5;
      | **if** $p(o) > \varepsilon$ **then** $\Gamma^\varepsilon \leftarrow \Gamma^\varepsilon \cup \{o\}$ ;
   **end**

---

## 5.2. p-value and Recommendation Set

Let $\alpha_h$ be the conformity or nonconformity value of an item $o_h$. For nonconformity measures, the proportion of examples having a nonconformity value greater than the new example defines the $p$-value,

$$p(o_{n+1}) = \frac{\left|\{h|m+1 \leq h \leq n+1, \alpha_h \geq \alpha_{n+1}\}\right|}{l+1}. \tag{4}$$

In the case of conformity measure, we define it as the proportion of examples having conformity value less than the new example,

$$p(o_{n+1}) = \frac{\left|\{h|m+1 \leq h \leq n+1, \alpha_h \leq \alpha_{n+1}\}\right|}{l+1}. \tag{5}$$

For a target user $u_j$, the recommendation set is then constructed by computing the $p$-value for every unused item. All the items whose $p$-value is greater than the predetermined significance level $\varepsilon$ will form a recommendation region $\Gamma^\varepsilon$.

$$\Gamma^\varepsilon = \{o \mid p(o) > \varepsilon\}.$$

## 5.3. Algorithm

In this section, we describe the algorithm by using the concepts defined in the previous sections. Algorithm 1 outlines the main flow of the proposed method. At first, we divide the dataset into a proper training set and calibration set. Next, we compute every item's nonconformity value in the calibration set and for every candidate item. We then compute the $p$-value for every candidate item and determine the recommendation set. The flowchart of the proposed algorithm is shown in Figure 2.
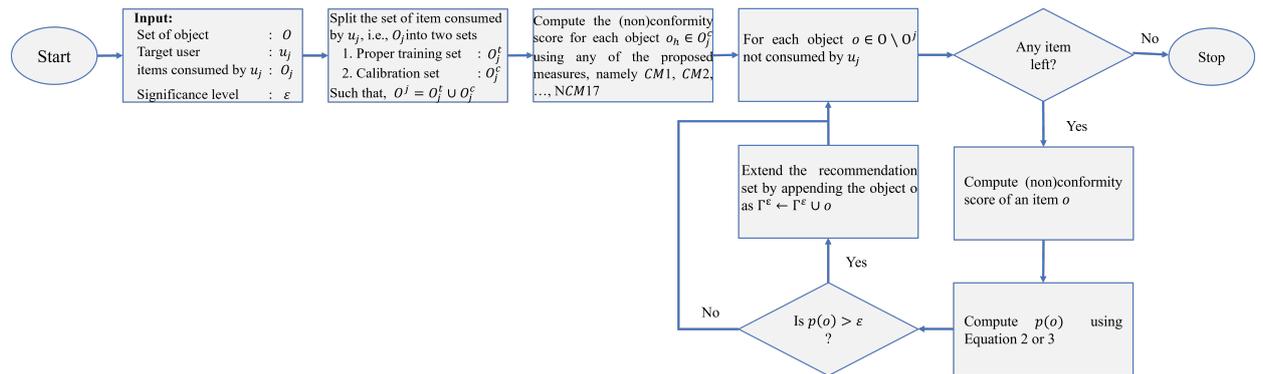


Figure 2: Inductive Conformal Recommender System.

**Example 5.** *We consider the precedence statistics given in Example 3. We divide the target user $u_1$ profile into $O_1^t = \{o_1, o_3, o_5\}$ and $O_1^c = \{o_7, o_9\}$. Let us compute the nonconformity*

13

*values, and p-value with respect to $o_2$ and the conformity measure **CM1**. The procedure is similar for other (non)conformity measures also.*

$$CM1(o_7) = \frac{Support(o_7)}{30} \times top1\big(PP(o_1 \mid o_7), PP(o_3 \mid o_7), PP(o_5 \mid o_7)\big) = \frac{15}{30} \times \frac{7}{15} = 0.23.$$

*Similarly, $CM1(o_9) = 0.3$ and $CM1(o_2) = 0.3$. Hence,*

$$p(o_2) = \frac{\big|\{o_h \mid o_h \in \{O_1^c \cup \{o_2\}\} \wedge CM1(o_h) \leq CM1(o_2)\}\big|}{|O_1^c \cup \{o_2\}|} = \frac{3}{3} = 1.$$

*We can compute the p-value for all the other candidate items and include those items whose p-value is greater than $\varepsilon$ in the recommendation set.*

### 5.4. Validity and Efficiency

We have already shown that the (non)conformity measures defined above satisfy the invariant property in Lemma 1. Hence, following the line of argument given by Vovk et al. [34], it is easy to see that our proposed method ICRS satisfies the validity property.

**Lemma 2.** *If objects $o_{m+1}, o_{m+2}, \ldots, o_{n+1}$ are independently and identically distributed (i.i.d.) in terms of their precedence relations with individual items in the history, then the probability of error that $o_{n+1} \notin \Gamma^\varepsilon(o_1, o_2, \ldots, o_m)$ will not exceed $\varepsilon \in [0, 1]$ i.e., $P(P(o_{n+1}) \leq \varepsilon) \leq \varepsilon$.*

*Proof.* An error occurs when $P(o_{n+1}) \leq \varepsilon$. That is, when $\alpha_{n+1}$ is among the $\lfloor \varepsilon(l+1) \rfloor$ largest elements of the set $\{\alpha_{m+1}, \alpha_{m+2}, \ldots, \alpha_{n+1}\}$. When all the objects are in i.i.d in terms of precedence relations with the set of items consumed by an user, all permutations of the set $\{\alpha_{m+1}, \ldots, \alpha_{n+1}\}$ are equiprobable. Thus, the probability that $\alpha_{n+1}$ is among the $\lfloor \varepsilon(l+1) \rfloor$ largest elements does not exceed $\varepsilon$, which is therefore the probability of error. $\qquad \square$

In addition to satisfying the validity property, it is desirable to have an efficient recommendation set. In the conformal framework setting, a narrow set with higher confidence is more efficient. We empirically analyze the validity and efficiency properties in Section 6.

### 5.5. Time Complexity Analysis

In this section, we analyze the time complexity of the proposed method against transductive conformal recommender systems [4] and the underlying precedence mining based algorithm [8]. For simplicity, we assume that the calibration set size is the same as that of candidate-set ($|O_j^{candidates}|$) in the Conformal Recommender System [4]. We know that $m$ is the size of the proper training set, and $l$ is the calibration set size. Let $n_c$ be the number of candidate items i.e., $n_c = n_i - n$. Since the complexity of measuring nonconformity scores varies from measure to measure, we assume it to be $O(t)$. With $O(t)$ as the complexity of nonconformity measure, the inductive conformal predictor takes $O((l + n_c)t)$ time complexity to determine all the required $p$-values and make recommendations. On the other hand, transductive conformal recommender systems take $O(n_c lmt)$ complexity with

$O(t)$ as the nonconformity measure's complexity. Kagita et al. [4] reduce it to $O(n_c l m)$ using the relation between the score and precedence probability, but it is higher than the inductive conformal recommender system. The complexity of the precedence mining based recommender system is $O(n_c n)$.

Table 1: Summary of experimental datasets.

| Dataset | Users | Items | Records |
|---|---|---|---|
| Personality-2018 | 1820 | 35196 | 1,028,752 |
| Flixsters | 20,618 | 28,331 | 1,048,575 |
| MovieLens 10M | 71,567 | 10,681 | 10,000,054 |
| MovieLens 20M | 138,494 | 26,745 | 20,000,262 |
| MovieLens 25M | 162,000 | 62,000 | 25,000,096 |
| MovieLens-Latest-V1 | 229,061 | 26,780 | 21,063,128 |
| MovieLens-Latest-V2 | 280,000 | 58,000 | 27,753, 445 |

## 6. Empirical Study

In this section, we empirically evaluate the efficacy of proposed *Inductive Conformal Recommender System (ICRS)*. We provide an in-depth quantitative evaluation with regard to the prediction accuracy and running time on seven real-world datasets of varying size. The characteristics of these datasets are reported in Table 1. In all our experiments, we converted the multi-class (different ratings) datasets into one class by setting a threshold to 0. The prediction accuracy of the comparing algorithms are evaluated based on the ranking-based performance metrics that is *Average Precision (AP), Area Under Curve (AUC), Normalized Discounted Cumulative Gain (NDCG)* and *Reverse Reciprocal (RR)*. We also evaluated the performance based on top-K recommendation metrics that is *Precion@K, Recall@K* and *F1@K* [35]. We compared our proposed method ICRS with the underlying Precedence Mining Model [8] and the Conformal Recommender Systems (CRS-max and CRS-med) [4]. In ICRS, to fine-tune the values of parameters $n$ and $k$, we experimented with different combinations and selected $n$ to be 30% of the profile and $k$ to be 30% and the remaining 40% is the test data. All the results reported here are the average of 500 randomly selected instances. We use a notation $ICRS < x >$ to denote an inductive conformal recommender system that uses (non)conformity measure $x$. For example, $ICRS1$ uses conformity measure 1 (CM1).

The remainder of the section is structured as follows. In Section 6.1, we report the experimental evaluation of the validity and efficiency of the proposed methods. Section 6.2 report comparative experimental results in terms of *ranking-based metrics, top-k recommendation metrics* and *execution time*.

15

## 6.1. Validity and Efficiency

This subsection empirically evaluates the validity and efficiency of the proposed approach. We adapt the definitions of validity and efficiency given by Kagita et al. [4]. Figure 3 and Figure 4 shows the validity and efficiency of the proposed approach respectively, over seven different datasets. We report the validity and efficiency related to *ICRS1, ICRS3, ICRS7, ICRS11* and *ICRS15*[1]. It can be seen from Figure 3 that the error is proportional to $\varepsilon$ and in the relative bound of $\varepsilon$. Figure 4 reports the error related to efficiency. It can be seen from the figures that even for smaller values of $\varepsilon$, most of the irrelevant items are filtered out hence, resulting in a small error. We also observed that, for higher values of $\varepsilon$, the recommendation set is more informative for all the strategies.
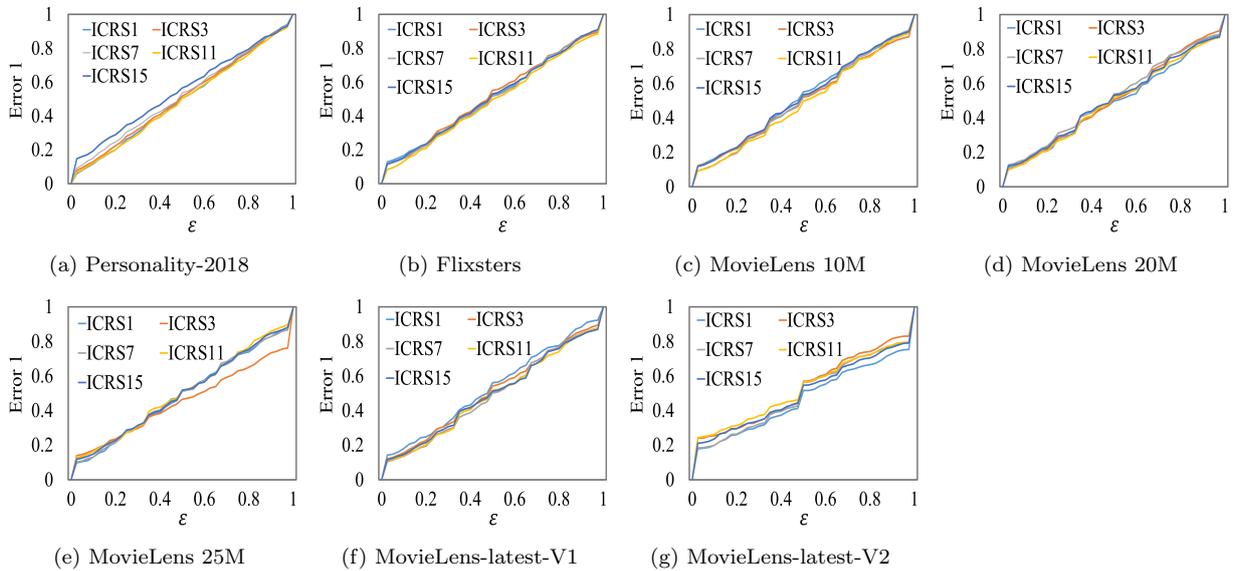


Figure 3: Evaluation of recommendation validity for different datasets

## 6.2. Comparative Analysis

In this section, we carried out experiments to demonstrate that the proposed methods achieve comparable results with significantly reduced execution times. Table 3 gives the findings related to ranking-based evaluation measures over seven datasets. Each result is composed of *mean* and *rank*. The rank reflects relative performance of an algorithm over a dataset for a given evaluation measure. In the case of ties, we have assigned the average rank. Furthermore, the entries in boldface highlight best results among all the algorithms being compared.

To carry out comparative analysis in more well-founded ways, we employed *Friedman test* which is widely-accepted as the favorable statistical test for comparing more than two

---

[1]ICRS3, ICRS7, ICRS11, and ICRS15 use the median strategy. Similar results have been observed for other strategies also.

|  | (a) Personality-2018 | (b) Flixsters | (c) MovieLens 10M | (d) MovieLens 20M |

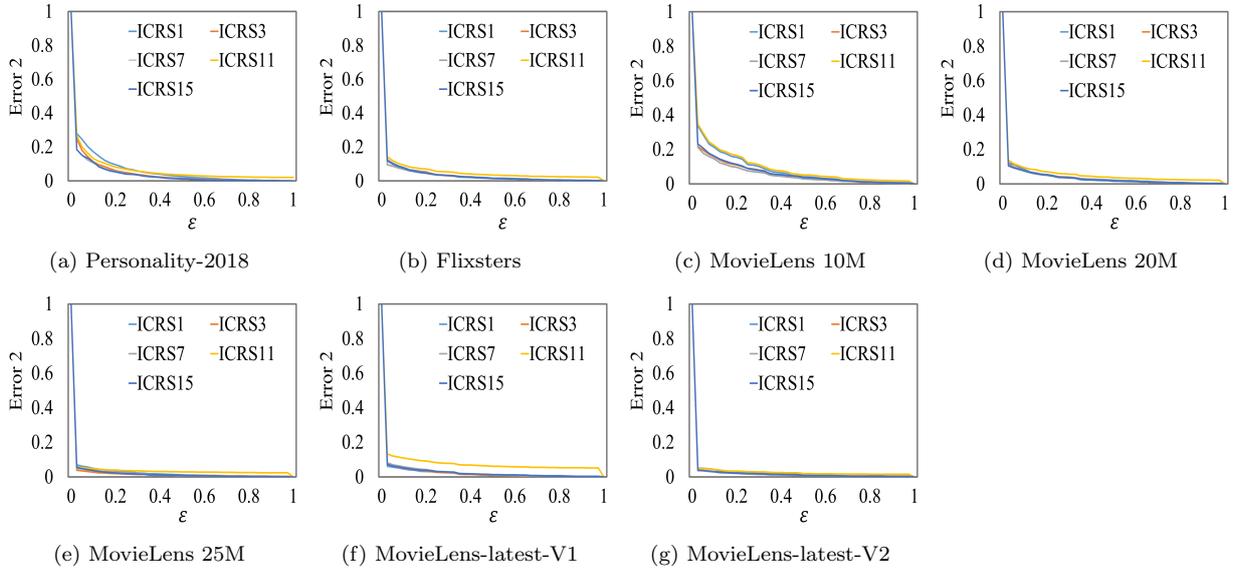|  | (e) MovieLens 25M | (f) MovieLens-latest-V1 | (g) MovieLens-latest-V2 |

Figure 4: Efficiency of recommendation for different datasets

algorithms over multiple data sets [36]. For each evaluation criterion, *Friedman statistics* $F_F$ and the corresponding critical value are reported in Table 2. It can be observed that at significance level $\alpha = 0.05$, Friedman test rejects the null hypothesis of "equal" performance for each evaluation metric. This leads to the use of post-hoc tests to assess the pairwise differences between two algorithms within a multiple comparison test. We use the Nemenyi test to check whether the proposed methods achieves a competitive performance against the algorithms being compared [36]. The performance of two algorithms is significantly different if the corresponding average ranks differ by at least the critical difference $CD = q_\alpha \sqrt{\frac{\mathcal{K}(\mathcal{K}+1)}{6\mathcal{N}}}$, where the value $q_\alpha$ is based on the Studentized range statistic divided by $\sqrt{2}$. For Nemenyi test with $\mathcal{K} = 20$, we have $q_\alpha = 3.5438$ at significance level $\alpha = 0.05$ and thus $CD = 11.2065$ [36].

Table 2: Summary of the Friedman Statistics $F_F(\mathcal{K} = 20, \mathcal{N} = 7)$ and the Critical Value in Terms of Each Evaluation Metric ($\mathcal{K}$: # Comparing Algorithms; $\mathcal{N}$: # Data Sets).

| Metric | $F_F$ | Critical Value ($\alpha = 0.05$) |
|--------|-------|----------------------------------|
| AP | 10.1378 | |
| AUC | 11.1020 | |
| NDCG | 11.3810 | 1.6785 |
| RR | 15.2027 | |

Figure 5 gives the CD diagrams [36] for each evaluation criterion, where the average rank of each comparing algorithm is marked along the axis (lower ranks to the left). It can be seen from the Figure 5 that the proposed methods achieve better performance than *CRS-Med*
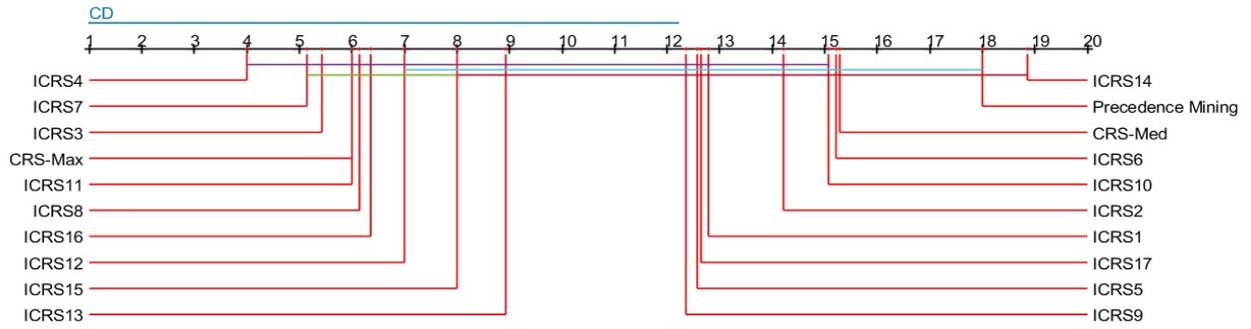
17

Table 3: Experimental results of each comparing algorithm in terms of AP, AUC, NDCG, and RR.

**AP**

| Comparing algorithm | Personality-2018 | | Flixsters | | MovLens 10M | | MovieLens 20M | | MovieLens 25M | | MovieLens-Latest | | MovieLens-Latest-V2 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Precedence Mining | 0.19 | 12 | 0.15 | 14 | 0.06 | 20 | 0.03 | 20 | 0.10 | 20 | 0.01 | 20 | 0.08 | 20 |
| CRS-Med | 0.13 | 19 | 0.16 | 9 | 0.15 | 16 | 0.15 | 16 | 0.14 | 13 | 0.15 | 15 | 0.17 | 19 |
| CRS-Max | 0.15 | 14 | **0.20** | **1** | 0.18 | 7.5 | 0.18 | 5 | **0.16** | 4.5 | **0.18** | 4 | 0.21 | 6 |
| ICRS1 | 0.23 | 8 | 0.19 | 2.5 | 0.15 | 16 | 0.15 | 16 | 0.13 | 16 | 0.15 | 15 | 0.19 | 16 |
| ICRS2 | 0.14 | 16.5 | 0.08 | 18.5 | 0.16 | 12.5 | 0.16 | 12 | 0.14 | 13 | 0.15 | 15 | 0.20 | 12 |
| ICRS3 | **0.26** | **1.5** | 0.18 | 4.5 | 0.18 | 7.5 | 0.18 | 5 | **0.16** | 4.5 | 0.17 | 9 | 0.21 | 6 |
| ICRS4 | **0.26** | **1.5** | 0.18 | 4.5 | **0.19** | **2.5** | 0.18 | 5 | **0.16** | 4.5 | **0.18** | 4 | 0.21 | 6 |
| ICRS5 | 0.22 | 10.5 | 0.19 | 2.5 | 0.15 | 16 | 0.15 | 16 | 0.13 | 16 | 0.15 | 15 | 0.20 | 12 |
| ICRS6 | 0.16 | 13 | 0.08 | 18.5 | 0.15 | 16 | 0.16 | 12 | 0.13 | 16 | 0.15 | 15 | 0.19 | 16 |
| ICRS7 | 0.24 | 5 | 0.16 | 9 | **0.19** | **2.5** | 0.18 | 5 | **0.16** | 4.5 | **0.18** | 4 | 0.21 | 6 |
| ICRS8 | 0.24 | 5 | 0.15 | 14 | **0.19** | **2.5** | **0.19** | **1** | **0.16** | 4.5 | **0.18** | 4 | 0.20 | 12 |
| ICRS9 | 0.14 | 16.5 | 0.15 | 14 | 0.17 | 11 | 0.16 | 12 | 0.15 | 10 | 0.16 | 11 | 0.20 | 12 |
| ICRS10 | 0.14 | 16.5 | 0.09 | 17 | 0.15 | 16 | 0.15 | 16 | 0.14 | 13 | 0.15 | 15 | 0.20 | 12 |
| ICRS11 | 0.22 | 10.5 | 0.16 | 9 | 0.18 | 7.5 | 0.18 | 5 | **0.16** | 4.5 | **0.18** | 4 | **0.22** | **1.5** |
| ICRS12 | 0.23 | 8 | 0.15 | 14 | 0.18 | 7.5 | 0.18 | 5 | **0.16** | 4.5 | **0.18** | 4 | 0.21 | 6 |
| ICRS13 | 0.14 | 16.5 | 0.15 | 14 | 0.18 | 7.5 | 0.17 | 9.5 | **0.16** | 4.5 | 0.17 | 9 | **0.22** | **1.5** |
| ICRS14 | 0.07 | 20 | 0.05 | 20 | 0.14 | 19 | 0.14 | 19 | 0.11 | 19 | 0.13 | 19 | 0.19 | 16 |
| ICRS15 | 0.24 | 5 | 0.16 | 9 | 0.18 | 7.5 | 0.17 | 9.5 | 0.15 | 10 | 0.17 | 9 | 0.21 | 6 |
| ICRS16 | 0.23 | 8 | 0.16 | 9 | **0.19** | **2.5** | 0.18 | 5 | 0.15 | 10 | **0.18** | 4 | 0.21 | 6 |
| ICRS17 | 0.25 | 3 | 0.17 | 6 | 0.16 | 12.5 | 0.15 | 16 | 0.12 | 18 | 0.15 | 15 | 0.18 | 18 |

**AUC**

| Comparing algorithm | Personality-2018 | | Flixsters | | MovLens 10M | | MovieLens 20M | | MovieLens 25M | | MovieLens-Latest | | MovieLens-Latest-V2 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Precedence Mining | 0.92 | 2 | **0.96** | **1** | 0.64 | 20 | 0.90 | 14 | **0.98** | **1** | 0.65 | 20 | **0.92** | **1** |
| CRS-Med | 0.62 | 16 | 0.87 | 7.5 | 0.72 | 19 | 0.75 | 20 | 0.77 | 20 | 0.75 | 19 | 0.67 | 20 |
| CRS-Max | 0.57 | 18 | 0.83 | 13.5 | 0.84 | 14 | 0.86 | 16 | 0.88 | 16 | 0.87 | 15 | 0.80 | 18 |
| ICRS1 | **0.93** | **1** | 0.91 | 2.5 | **0.90** | **3** | **0.92** | 5.5 | 0.93 | 6 | **0.93** | 4.5 | 0.86 | 6 |
| ICRS2 | 0.67 | 15 | 0.60 | 17 | 0.83 | 15 | 0.87 | 15 | 0.89 | 15 | 0.88 | 14 | 0.84 | 14.5 |
| ICRS3 | 0.91 | 4 | 0.88 | 5.5 | **0.90** | **3** | **0.92** | 5.5 | 0.94 | 2 | 0.92 | 11 | 0.86 | 6 |
| ICRS4 | 0.91 | 4 | 0.89 | 4 | **0.90** | **3** | **0.92** | 5.5 | 0.92 | 11.5 | **0.93** | 4.5 | 0.85 | 11 |
| ICRS5 | 0.91 | 4 | 0.91 | 2.5 | 0.89 | 9.5 | **0.92** | 5.5 | 0.93 | 6 | 0.92 | 11 | 0.85 | 11 |
| ICRS6 | 0.59 | 17 | 0.52 | 18 | 0.81 | 16 | 0.84 | 17 | 0.86 | 17 | 0.86 | 16 | 0.81 | 16 |
| ICRS7 | 0.90 | 6 | 0.86 | 10 | **0.90** | **3** | **0.92** | 5.5 | 0.92 | 11.5 | **0.93** | 4.5 | 0.86 | 6 |
| ICRS8 | 0.88 | 8.5 | 0.79 | 15 | 0.89 | 9.5 | **0.92** | 5.5 | 0.92 | 11.5 | 0.92 | 11 | 0.84 | 14.5 |
| ICRS9 | 0.87 | 11 | 0.87 | 7.5 | 0.89 | 9.5 | **0.92** | 5.5 | 0.93 | 6 | **0.93** | 4.5 | 0.85 | 11 |
| ICRS10 | 0.53 | 19 | 0.47 | 19 | 0.78 | 17 | 0.82 | 18.5 | 0.84 | 18 | 0.84 | 17 | 0.80 | 18 |
| ICRS11 | 0.87 | 11 | 0.84 | 12 | 0.89 | 9.5 | 0.91 | 12 | 0.93 | 6 | **0.93** | 4.5 | 0.87 | 2.5 |
| ICRS12 | 0.85 | 14 | 0.75 | 16 | 0.89 | 9.5 | 0.91 | 12 | 0.91 | 14 | 0.92 | 11 | 0.85 | 11 |
| ICRS13 | 0.86 | 13 | 0.86 | 10 | **0.90** | **3** | **0.92** | 5.5 | 0.93 | 6 | **0.93** | 4.5 | 0.87 | 2.5 |
| ICRS14 | 0.51 | 20 | 0.45 | 20 | 0.74 | 18 | 0.82 | 18.5 | 0.83 | 19 | 0.83 | 18 | 0.80 | 18 |
| ICRS15 | 0.88 | 8.5 | 0.83 | 13.5 | 0.89 | 9.5 | **0.92** | 5.5 | 0.92 | 11.5 | 0.92 | 11 | 0.85 | 11 |
| ICRS16 | 0.87 | 11 | 0.86 | 10 | 0.89 | 9.5 | **0.92** | 5.5 | 0.93 | 6 | **0.93** | 4.5 | 0.86 | 6 |
| ICRS17 | 0.89 | 7 | 0.88 | 5.5 | 0.89 | 9.5 | 0.91 | 12 | 0.93 | 6 | **0.93** | 4.5 | 0.86 | 6 |

**NDCG**

| Comparing algorithm | Personality-2018 | | Flixsters | | MovLens 10M | | MovieLens 20M | | MovieLens 25M | | MovieLens-Latest | | MovieLens-Latest-V2 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Precedence Mining | 0.66 | 12 | 0.53 | 13.5 | 0.16 | 20 | 0.37 | 20 | 0.48 | 19.5 | 0.27 | 20 | 0.41 | 20 |
| CRS-Med | 0.58 | 18 | 0.55 | 7 | 0.54 | 17.5 | 0.53 | 18 | 0.52 | 15.5 | 0.53 | 17.5 | 0.48 | 19 |
| CRS-Max | 0.59 | 16.5 | **0.61** | **1** | **0.59** | **4** | **0.58** | 3.5 | **0.56** | 2.5 | **0.58** | **1** | 0.54 | 7 |
| ICRS1 | 0.70 | 5.5 | 0.57 | 3 | 0.56 | 13.5 | 0.55 | 13.5 | 0.53 | 12.5 | 0.55 | 12 | 0.52 | 14.5 |
| ICRS2 | 0.59 | 16.5 | 0.44 | 17 | 0.55 | 15.5 | 0.55 | 13.5 | 0.52 | 15.5 | 0.53 | 17.5 | 0.53 | 12 |
| ICRS3 | **0.72** | **1.5** | 0.56 | 5 | **0.59** | **4** | 0.57 | 8 | 0.55 | 7.5 | 0.57 | 5.5 | 0.54 | 7 |
| ICRS4 | **0.72** | **1.5** | 0.57 | 3 | **0.59** | **4** | **0.58** | 3.5 | 0.55 | 7.5 | 0.57 | 5.5 | 0.54 | 7 |
| ICRS5 | 0.69 | 9.5 | 0.57 | 3 | 0.56 | 13.5 | 0.55 | 13.5 | 0.53 | 12.5 | 0.54 | 14.5 | 0.53 | 12 |
| ICRS6 | 0.60 | 15 | 0.43 | 18.5 | 0.55 | 15.5 | 0.55 | 13.5 | 0.51 | 18 | 0.54 | 14.5 | 0.51 | 17 |
| ICRS7 | 0.70 | 5.5 | 0.55 | 7 | **0.59** | **4** | **0.58** | 3.5 | 0.55 | 7.5 | 0.57 | 5.5 | 0.54 | 7 |
| ICRS8 | 0.70 | 5.5 | 0.52 | 15.5 | **0.59** | **4** | **0.58** | 3.5 | 0.55 | 7.5 | 0.57 | 5.5 | 0.53 | 12 |
| ICRS9 | 0.62 | 13.5 | 0.54 | 10.5 | 0.58 | 9.5 | 0.55 | 13.5 | 0.55 | 7.5 | 0.56 | 10.5 | 0.54 | 7 |
| ICRS10 | 0.57 | 19 | 0.43 | 18.5 | 0.54 | 17.5 | 0.54 | 17 | 0.52 | 15.5 | 0.54 | 14.5 | 0.52 | 14.5 |
| ICRS11 | 0.69 | 9.5 | 0.54 | 10.5 | **0.59** | **4** | **0.58** | 3.5 | **0.56** | 2.5 | 0.57 | 5.5 | **0.56** | **1** |
| ICRS12 | 0.68 | 11 | 0.52 | 15.5 | **0.59** | **4** | **0.58** | 3.5 | **0.56** | 2.5 | 0.57 | 5.5 | 0.54 | 7 |
| ICRS13 | 0.62 | 13.5 | 0.53 | 13.5 | 0.58 | 9.5 | 0.57 | 8 | **0.56** | 2.5 | 0.57 | 5.5 | 0.55 | 2.5 |
| ICRS14 | 0.50 | 20 | 0.38 | 20 | 0.51 | 19 | 0.51 | 19 | 0.48 | 19.5 | 0.50 | 19 | 0.51 | 17 |
| ICRS15 | 0.70 | 5.5 | 0.54 | 10.5 | 0.58 | 9.5 | 0.56 | 10 | 0.54 | 11 | 0.56 | 10.5 | 0.54 | 7 |
| ICRS16 | 0.70 | 5.5 | 0.54 | 10.5 | 0.58 | 9.5 | 0.57 | 8 | 0.55 | 7.5 | 0.57 | 5.5 | 0.55 | 2.5 |
| ICRS17 | 0.70 | 5.5 | 0.55 | 7 | 0.57 | 12 | 0.55 | 13.5 | 0.52 | 15.5 | 0.54 | 14.5 | 0.51 | 17 |

**RR**

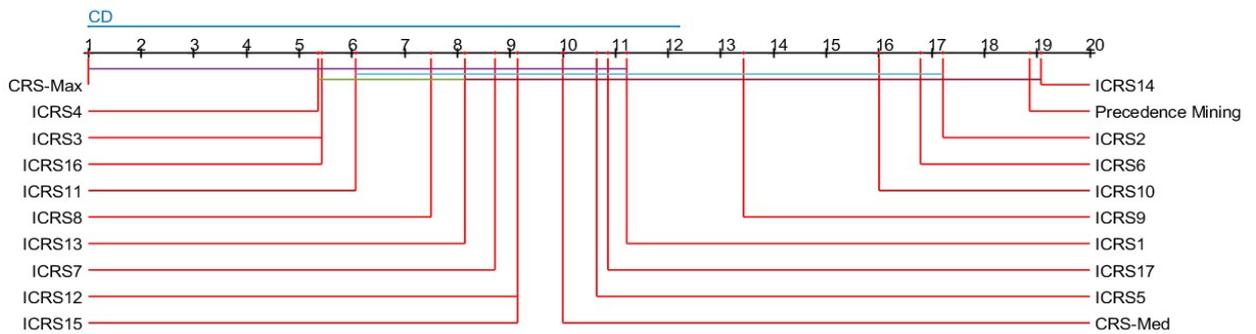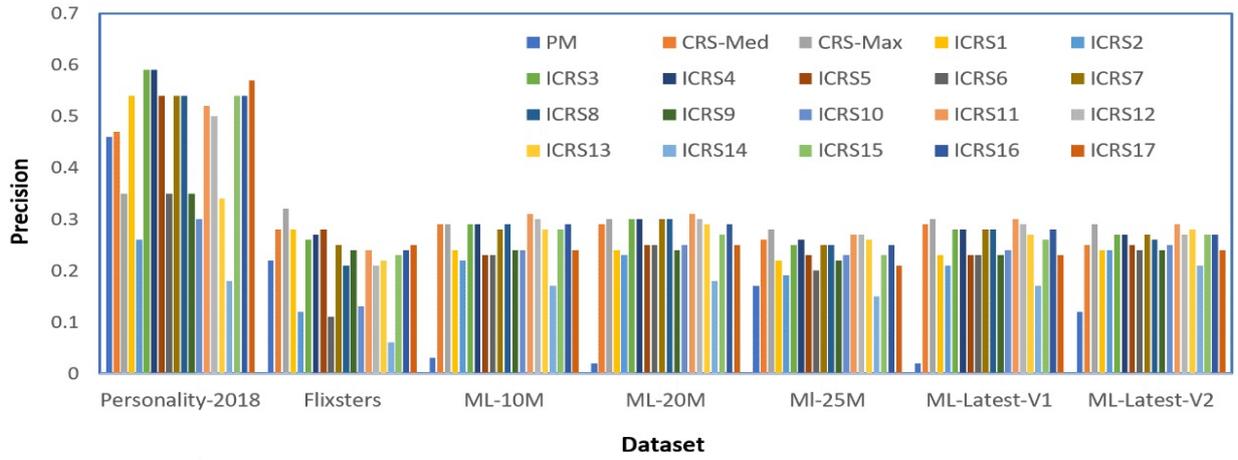| Comparing algorithm | Personality-2018 | | Flixsters | | MovLens 10M | | MovieLens 20M | | MovieLens 25M | | MovieLens-Latest | | MovieLens-Latest-V2 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Precedence Mining | 0.82 | 12 | 0.19 | 20 | 0.31 | 20 | 0.24 | 20 | 0.34 | 20 | 0.22 | 20 | 0.25 | 20 |
| CRS-Med | 0.78 | 14 | 0.60 | 2 | 0.59 | 14.5 | 0.64 | 9 | 0.59 | 5.5 | 0.62 | 6.5 | 0.55 | 18.5 |
| CRS-Max | **0.97** | **1** | **0.92** | **1** | **0.76** | **1** | **0.76** | **1** | **0.71** | **1** | **0.75** | **1** | **0.72** | **1** |
| ICRS1 | 0.85 | 10.5 | 0.51 | 4.5 | 0.60 | 12.5 | 0.60 | 14 | 0.56 | 12 | 0.59 | 12 | 0.60 | 13 |
| ICRS2 | 0.51 | 19 | 0.37 | 16 | 0.53 | 18 | 0.57 | 16.5 | 0.49 | 17.5 | 0.51 | 18 | 0.58 | 15.5 |
| ICRS3 | 0.91 | 2.5 | 0.47 | 7.5 | 0.65 | 4 | 0.65 | 6.5 | 0.59 | 5.5 | 0.62 | 6.5 | 0.63 | 5.5 |
| ICRS4 | 0.91 | 2.5 | 0.47 | 7.5 | 0.64 | 6.5 | 0.66 | 4 | 0.61 | 3 | 0.62 | 6.5 | 0.62 | 7.5 |
| ICRS5 | 0.85 | 10.5 | 0.52 | 3 | 0.59 | 14.5 | 0.60 | 14 | 0.57 | 10 | 0.59 | 12 | 0.61 | 10.5 |
| ICRS6 | 0.62 | 17 | 0.32 | 18 | 0.54 | 17 | 0.60 | 14 | 0.49 | 17.5 | 0.53 | 17 | 0.56 | 17 |
| ICRS7 | 0.86 | 7.5 | 0.45 | 10.5 | 0.63 | 8.5 | 0.67 | 2.5 | 0.56 | 12 | 0.61 | 9.5 | 0.61 | 10.5 |
| ICRS8 | 0.86 | 7.5 | 0.42 | 15 | 0.65 | 4 | 0.65 | 6.5 | 0.59 | 5.5 | 0.62 | 6.5 | 0.62 | 7.5 |
| ICRS9 | 0.69 | 15 | 0.48 | 6 | 0.60 | 12.5 | 0.56 | 18 | 0.54 | 16 | 0.56 | 16 | 0.61 | 10.5 |
| ICRS10 | 0.57 | 18 | 0.36 | 17 | 0.56 | 16 | 0.57 | 16.5 | 0.55 | 14.5 | 0.57 | 14.5 | 0.58 | 15.5 |
| ICRS11 | 0.86 | 7.5 | 0.44 | 12.5 | 0.64 | 6.5 | 0.67 | 2.5 | 0.58 | 8.5 | 0.63 | 3 | 0.67 | 2 |
| ICRS12 | 0.79 | 13 | 0.43 | 14 | 0.63 | 8.5 | 0.65 | 6.5 | 0.63 | 2 | 0.61 | 9.5 | 0.61 | 10.5 |
| ICRS13 | 0.65 | 16 | 0.44 | 12.5 | 0.65 | 4 | 0.62 | 10.5 | 0.59 | 5.5 | 0.63 | 3 | 0.63 | 5.5 |
| ICRS14 | 0.40 | 20 | 0.29 | 19 | 0.47 | 19 | 0.50 | 19 | 0.43 | 19 | 0.46 | 19 | 0.55 | 18.5 |
| ICRS15 | 0.88 | 5 | 0.46 | 9 | 0.61 | 10.5 | 0.61 | 12 | 0.56 | 12 | 0.59 | 12 | 0.64 | 3.5 |
| ICRS16 | 0.89 | 4 | 0.45 | 10.5 | 0.67 | 2 | 0.65 | 6.5 | 0.58 | 8.5 | 0.63 | 3 | 0.64 | 3.5 |
| ICRS17 | 0.86 | 7.5 | 0.51 | 4.5 | 0.61 | 10.5 | 0.62 | 10.5 | 0.55 | 14.5 | 0.57 | 14.5 | 0.59 | 14 |

(a) AP



(b) AUC



(c) NDCG



(d) RR

Figure 5: CD diagrams of the comparing algorithms under each evaluation criterion.

and *Precedence Mining* models over most of the evaluation metrics. We can also observe that the proposed approaches achieve similar performance to *CRS-Max* or even derive a better rank in most cases, especially the median-based inductive approaches (*ICRS3,ICRS7, ICRS11*, and *ICRS15*) and mean-based inductive approaches (*ICRS4,ICRS8, ICRS12*, and *ICRS16*). We have also observed similar results with Maximum strategy based methods (*ICRS1, ICRS5, ICRS9, ICRS13*, and *ICRS17*), whereas Minimum strategy based approaches (*ICRS2, ICRS6, ICRS10*, and *ICRS14*) performing poorly among the seventeen proposed approaches. This comprehensive analysis reveals that Median and Mean-based approaches capture the true precedence relations of a new item with respect to a user profile compared to the Minimum strategy. The reason could be that sometimes there is a higher chance of a user consuming items that are not of his regular interest but due to other users' influence (like family, friends, etc.) or situational context . These items do not follow good precedence relations with users' actual interests. Therefore, it is evident that there is a higher chance of minimum-based strategies capturing such precedence relations and that do not represent the allure of a new item concerning the user profile. In other words, there may be some noisy/outlier points in the user profile, and minimum-based strategies are more attractive to these points and therefore not suitable to measure (non)conformity. Though the same could be valid with the *Maximum* based strategies, it is more likely that even a single item in the profile can influence to consume another item that follows higher precedence relations. For example, a Deep Learning course may have a higher precedence relation with a Machine Learning course, and that could be an influencing factor for a student to opt for a Deep Learning course irrespective of other courses in the student profile. Experiment results corroborate our claims.
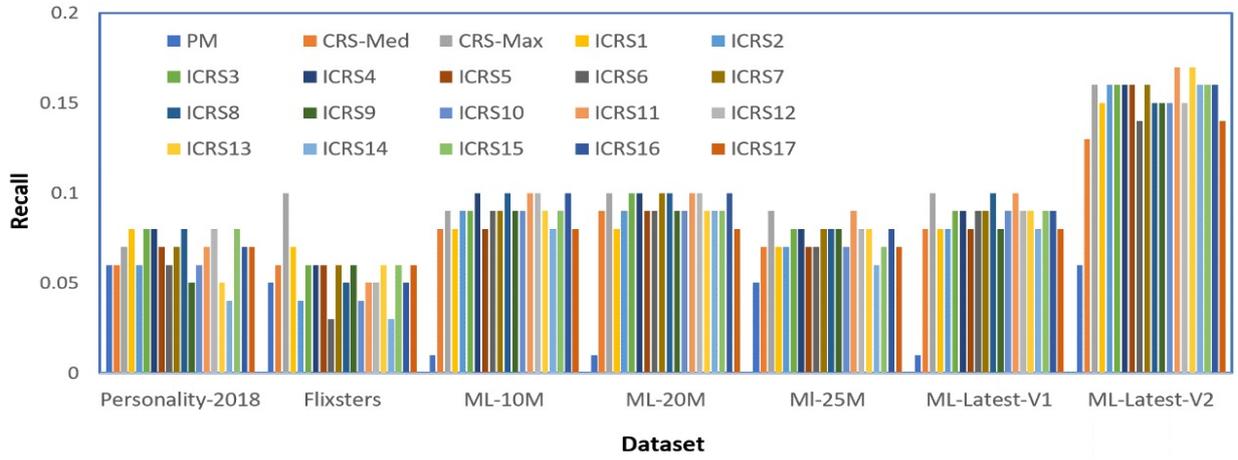
In the second set of experiments, we compare the performance in terms of top-$k$ recommendation measures, namely *precision@10*, *recall@10* and *F1@10* (Figure 6). We observe the similar results with varying the number of recommendations. It can be observed from the figures that conformal approaches methods outperform the underlying precedence mining model. Furthermore, findings reveal that inductive variants are comparable with the *CRS-max* and *CRS-med*. Finally, we compared the execution time (in milliseconds) of the different approaches. It can be seen from the Figure 7 that the inductive conformal recommender systems are much faster than traditional conformal recommender systems and better than the precedence mining model. Altogether, the results corroborate our claim that the inductive variant achieves a similar level of accuracy compared to its counterparts but significantly reduces the execution time.

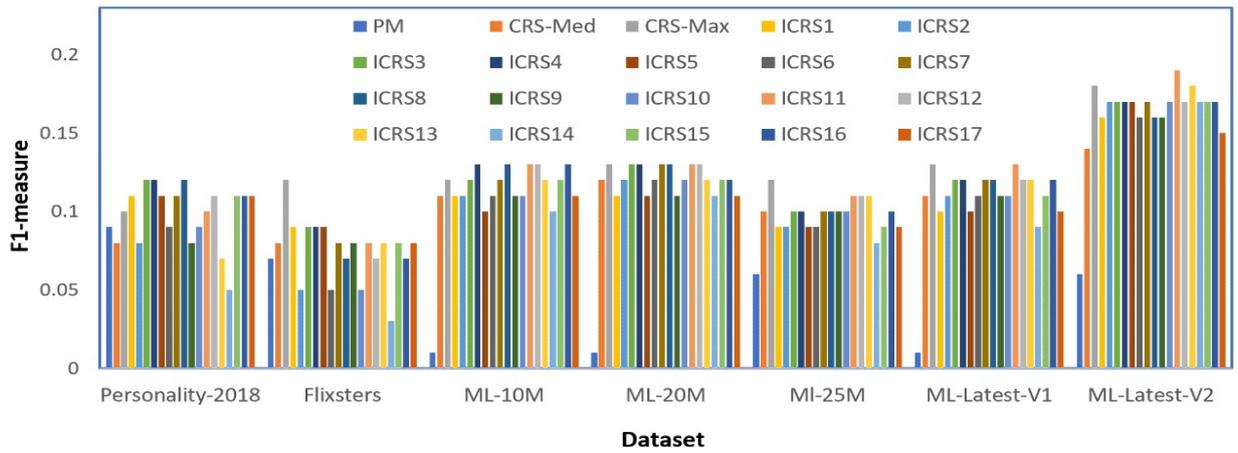## 7. Conclusions and Future Work

In this paper, we propose an inductive variants of the conformal recommender system that complements the recommendation by quantifying the (un)reliability in predictions. One natural limitation with the existing transductive variants is the computation time that prevents their applicability in the time constraint domains. We address this limitation and propose an inductive variant that maintains the same moderate level of predictive accuracy but reduces the computation time to a large extent. Our conformal approach exemplifies

(a) precision@10



(b) recall@10



(c) F$_1$@10

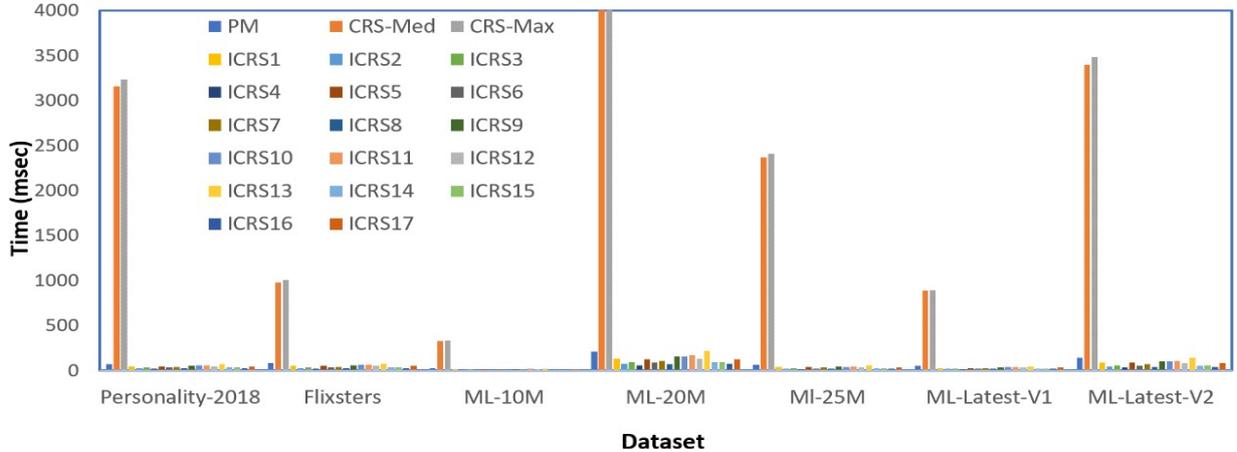Figure 6: Performance of each comparing algorithm in terms of top-k recommendation metrics.

Figure 7: Execution time comparison for different datasets

confidence in terms of the bounds on the error. Conformity/nonconformity measures are key component of any conformal recommendation framework, and the prediction accuracy largely depends on how well these measures are defined. In this work, we examined sevneteen different (non)conformity measures using the precedence relations among objects. We theoretically proved that the proposed (non)conformity measures adhere to the principle of validity under certain assumptions. Further, we emphasized our theoretical results with an empirical demonstration. Rigorous experiments on several real-world datasets demonstrated that the inductive conformal recommendation algorithms outperform the precedence mining based recommender system and non-inductive methods in terms of execution time. We observed that a few of the inductive variants outperforming the other approaches in terms of other crucial measures of the recommender system when the basic assumptions of the model are satisfied.

The current proposal sets a lot of scope for future research. Attaining the notion of confidence in different recommendation models by determining suitable (non)conformity measures is one of the exacting directions for enthusiastic researchers. Investigating the conformal prediction for group recommender systems is a direction worth studying. Exploring the conformal approach for different matrix factorization-based methods is another exciting direction to pursue.

## Acknowledgements

# References

[1] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishi Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, et al. Wide & deep learning for recommender systems. In *Proceedings of the 1st workshop on deep learning for recommender systems*, pages 7–10. ACM, 2016.

[2] Alexandros Karatzoglou and Balázs Hidasi. Deep learning for recommender systems. In *Proceedings of the eleventh ACM conference on recommender systems*, pages 396–397. ACM, 2017.

[3] Vikas Kumar, Arun K Pujari, Sandeep Kumar Sahu, Venkateswara Rao Kagita, and Vineet Padmanabhan. Collaborative filtering using multiple binary maximum margin matrix factorizations. *Information Sciences*, 380:1–11, 2017.

[4] Venkateswara Rao Kagita, Arun K Pujari, Vineet Padmanabhan, Sandeep Kumar Sahu, and Vikas Kumar. Conformal recommender system. *Information Sciences*, 405:157–174, 2017.

[5] Maciej A Mazurowski. Estimating confidence of individual rating predictions in collaborative filtering recommender systems. *Expert Systems with Applications*, 40(10):3847–3857, 2013.

[6] Antonio Hernando, JesúS Bobadilla, Fernando Ortega, and Jorge Tejedor. Incorporating reliability measurements into the predictions of a recommender system. *Information Sciences*, 218:1–16, 2013.

[7] Matthew R McLaughlin and Jonathan L Herlocker. A collaborative filtering algorithm and evaluation metric that accurately model the user experience. In *Proceedings of the 27th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 329–336. ACM, 2004.

[8] A. G. Parameswaran, G. Koutrika, B. Bercovitz, and H. Garcia-Molina. Recsplorer: recommendation algorithms based on precedence mining. In *SIGMOD*, pages 87–98, 2010.

[9] Tadiparthi VR Himabindu, Vineet Padmanabhan, and Arun K Pujari. Conformal matrix factorization based recommender system. *Information Sciences*, 467:685–707, 2018.

[10] Paul Resnick and Hal R Varian. Recommender systems. *Communications of the ACM*, 40(3):56–58, 1997.

[11] Jie Lu, Dianshuang Wu, Mingsong Mao, Wei Wang, and Guangquan Zhang. Recommender system application developments: a survey. *Decision Support Systems*, 74:12–32, 2015.

[12] Douglas W Oard, Jinmook Kim, et al. Implicit feedback for recommender systems. In *Proceedings of the AAAI workshop on recommender systems*, volume 83, pages 81–83. WoUongong, 1998.

[13] Vikas Kumar, Arun K Pujari, Sandeep Kumar Sahu, Venkateswara Rao Kagita, and Vineet Padmanabhan. Proximal maximum margin matrix factorization for collaborative filtering. *Pattern Recognition Letters*, 86:62–67, 2017.

[14] Dionisis Margaris, Dionysios Vasilopoulos, Costas Vassilakis, and Dimitris Spiliotopoulos. Improving collaborative filtering's rating prediction accuracy by introducing the common item rating past criterion. In *2019 10th International Conference on Information, Intelligence, Systems and Applications (IISA)*, pages 1–8. IEEE, 2019.

[15] Sundus Ayyaz, Usman Qamar, and Raheel Nawaz. Hcf-crs: A hybrid content based fuzzy conformal recommender system for providing recommendations with confidence. *PloS one*, 13(10):e0204849, 2018.

[16] Sean M McNee, Shyong K Lam, Catherine Guetzlaff, Joseph A Konstan, and John Riedl. Confidence displays and training in recommender systems. In *Proc. INTERACT*, volume 3, pages 176–183, 2003.

[17] Gediminas Adomavicius, Sreeharsha Kamireddy, and YoungOk Kwon. Towards more confident recommendations: Improving recommender systems using filtering approach based on rating variance. In *Proc. of the 17th Workshop on Information Technology and Systems*, pages 152–157, 2007.

[18] Panagiotis Symeonidis, Alexandros Nanopoulos, and Yannis Manolopoulos. Providing justifications in recommender systems. *IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems and Humans*, 38(6):1262–1272, 2008.

[19] Guy Shani and Asela Gunawardana. Evaluating recommendation systems. In *Recommender systems handbook*, pages 257–297. Springer, 2011.

[20] Yehuda Koren and Joe Sill. Ordrec: an ordinal model for predicting personalized item rating distributions. In *Proceedings of the fifth ACM conference on Recommender systems*, pages 117–124. ACM, 2011.

[21] Faezeh Sadat Gohari, Fereidoon Shams Aliee, and Hassan Haghighi. A new confidence-based recommendation approach: Combining trust and certainty. *Information Sciences*, 422:21–50, 2018.

[22] Rus M Mesas and Alejandro Bellogín. Exploiting recommendation confidence in decision-aware recommender systems. *Journal of Intelligent Information Systems*, 54(1):45–78, 2020.

[23] Raphaël Morsomme and Evgueni Smirnov. Conformal prediction for students' grades in a course recommender system. In *Conformal and Probabilistic Prediction and Applications*, pages 196–213, 2019.

[24] Harris Papadopoulos, Kostas Proedrou, Volodya Vovk, and Alex Gammerman. Inductive confidence machines for regression. In *European Conference on Machine Learning*, pages 345–356. Springer, 2002.

[25] Chao Wang, Qi Liu, Runze Wu, Enhong Chen, Chuanren Liu, Xunpeng Huang, and Zhenya Huang. Confidence-aware matrix factorization for recommender systems. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.

[26] Arthur F Da Costa, Marcelo G Manzato, and Ricardo JGB Campello. Boosting collaborative filtering with an ensemble of co-trained recommenders. *Expert Systems with Applications*, 115:427–441, 2019.

[27] Harris Papadopoulos. Inductive conformal prediction: Theory and application to neural networks. In *Tools in artificial intelligence*. IntechOpen, 2008.

[28] Azene Zenebe, Ant Ozok, and Anthony F Norcio. Personalized recommender systems in e-commerce and m-commerce: a comparative study. In *Conference on Human-Computer Interaction (HCI International)*, 2005.

[29] Glenn Shafer and Vladimir Vovk. A tutorial on conformal prediction. *Journal of Machine Learning Research*, 9(Mar):371–421, 2008.

[30] Vladimir Vovk, Alex Gammerman, and Glenn Shafer. *Algorithmic learning in a random world.* Springer Science & Business Media, 2005.

[31] Venkateswara Rao Kagita, Vineet Padmanabhan, and Arun K. Pujari. Precedence mining in group recommender systems. In Pradipta Maji, Ashish Ghosh, M. Narasimha Murty, Kuntal Ghosh, and Sankar K. Pal, editors, *Pattern Recognition and Machine Intelligence*, pages 701–707, Berlin, Heidelberg, 2013. Springer Berlin Heidelberg.

[32] Venkateswara Rao Kagita, Arun K. Pujari, and Vineet Padmanabhan. Group recommender systems: A virtual user approach based on precedence mining. In Stephen Cranefield and Abhaya Nayak, editors, *AI 2013: Advances in Artificial Intelligence*, pages 434–440, Cham, 2013. Springer International Publishing.

[33] Venkateswara Rao Kagita, Arun K. Pujari, and Vineet Padmanabhan. Virtual user approach for group recommender systems using precedence relations. *Information Sciences*, 294:15 – 30, 2015. Innovative Applications of Artificial Neural Networks in Engineering.

[34] V. Vovk, A. Gammerman, and G Shaffer. *Algorithmic Learning in a Random World.* Springer, 2005.

[35] Huayu Li, Richang Hong, Defu Lian, Zhiang Wu, Meng Wang, and Yong Ge. A relaxed ranking-based factor model for recommender system from implicit feedback. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence*, IJCAI'16, page 1683–1689. AAAI Press, 2016.

[36] Janez Demšar. Statistical comparisons of classifiers over multiple data sets. *The Journal of Machine Learning Research*, 7:1–30, 2006.

## Appendix A. Formal defintions of (non)conformity

We give the formal definitions of the (non)conformity measures described in Section 3 as follows.

$$CM2(o_h) = \underset{o_i \in O_j^t}{minimum} \ PC(o_i, o_h).$$

$$CM3(o_h) = \underset{o_i \in O_j^t}{median} \ PC(o_i, o_h).$$

$$CM4(o_h) = \underset{o_i \in O_j^t}{mean} \ PC(o_i, o_h).$$

$$CM5(o_h) = \underset{o_i \in O_j^t}{maximum} \ PC(o_i, o_h).$$

$$CM6(o_h) = \underset{o_i \in O_j^t}{minimum} \ PP(o_h \mid o_i).$$

$$CM7(o_h) = \underset{o_i \in O_j^t}{median} \ PP(o_h \mid o_i).$$

$$CM8(o_h) = \underset{o_i \in O_j^t}{mean} \ PP(o_h \mid o_i).$$

$$CM9(o_h) = \underset{o_i \in O_j^t}{maximum} \ PP(o_h \mid o_i).$$

$$CM10(o_h) = \underset{o_i \in O_j^t}{minimum} \ \frac{PC(o_i, o_h)}{Sup(o_i) - PC(o_h, o_i)}.$$

$$CM11(o_h) = \underset{o_i \in O_j^t}{median} \ \frac{PC(o_i, o_h)}{Sup(o_i) - PC(o_h, o_i)}.$$

$$CM12(o_h) = \underset{o_i \in O_j^t}{mean} \ \frac{PC(o_i, o_h)}{Sup(o_i) - PC(o_h, o_i)}.$$

$$CM13(o_h) = \underset{o_i \in O_j^t}{maximum} \ \frac{PC(o_i, o_h)}{Sup(o_i) - PC(o_h, o_i)}.$$

$$NCM14(o_h) = \underset{o_i \in O_j^t}{minimum} \ \frac{Sup(o_i) - PC(o_i, o_h)}{n_u}.$$

$$NCM15(o_h) = \underset{o_i \in O_j^t}{median} \ \frac{Sup(o_i) - PC(o_i, o_h)}{n_u}.$$

$$NCM16(o_h) = \underset{o_i \in O_j^t}{mean} \ \frac{Sup(o_i) - PC(o_i, o_h)}{n_u}.$$

$$NCM17(o_h) = \underset{o_i \in O_j^t}{maximum} \ \frac{Sup(o_i) - PC(o_i, o_h)}{n_u}.$$