

# Improving Collaborative Filtering's Rating Prediction Accuracy by Introducing the Common Item Rating Past Criterion

Dionisis Margaris  
Department of Informatics and Telecommunications  
University of Athens  
Athens, Greece  
[margaris@di.uoa.gr](mailto:margaris@di.uoa.gr)

Dionysios Vasilopoulos, Costas Vassilakis and Dimitris  
Spiliotopoulos  
Department of Informatics and Telecommunications  
University of the Peloponnese  
Tripoli, Greece  
[dvasilop@uop.gr](mailto:dvasilop@uop.gr), [costas@uop.gr](mailto:costas@uop.gr), [dspiliot@uop.gr](mailto:dspiliot@uop.gr)

**Abstract**—Collaborative filtering computes personalized recommendations by taking into account ratings expressed by users. Collaborative filtering algorithms firstly identify people having similar tastes, by examining the likeness of already entered ratings. Users with highly similar tastes are termed “near neighbors” and recommendations for a user are based on his near neighbors’ ratings. In order to measure similarity between users, so as to determine a user’s NNs, a similarity metric is used. Insofar, similarity metrics proposed in the literature either consider all user ratings equally or take into account temporal variations within the users’ or items’ ratings history. However users’ ratings are co-shaped according to the experiences that they had in the past; therefore if two users enter similar (or dissimilar) ratings for an item while having experienced –to a large extent– the same items in the past, this constitutes stronger evidence about user similarity (or dissimilarity). Insofar however, no similarity metric takes into account this aspect. In this work, we (1) propose an algorithm that considers the common item rating past in the rating prediction computation process, aiming to improve rating prediction quality, and (2) evaluate the proposed algorithm against seven widely used datasets, both dense and sparse, and considering two widely used user similarity metrics.

**Keywords**— *Collaborative Filtering, Items’ Rating Sequence, Common Item Rating Past Criterion, Pearson Correlation Coefficient, Cosine Similarity, Evaluation.*

## I. INTRODUCTION

Collaborative filtering (CF) computes personalized recommendations, by taking into account users’ past tastes and likings, which are recorded in the form of ratings entered in a ratings database. CF has been proven to be the most successful approach for building recommender systems (RSs) [1]. CF algorithms firstly identify people having similar tastes, by examining the resemblance of already entered ratings; for each user  $u$ , other users having highly similar tastes with  $u$  are designated as  $u$ ’s nearest neighbors (NNs). Afterwards, in order to predict the rating that user  $u$  would give to an item  $i$  that  $u$  has not reviewed yet, the ratings entered for  $i$  by  $u$ ’s NNs are combined [1], under the assumption that users who have exhibited similar tastes in the past, are likely to do so in the future as well [2,3].

In order to measure similarity between users, so as to compute a user’s NNs, a correlation coefficient is used. A correlation

coefficient maps pairs of entities (users or items) to a similarity metric, typically falling in the range of  $[0, 1]$  or in the range  $[-1, 1]$ ; in both cases, the highest value in the range denotes highly similar entities, while the lowest value denotes entirely dissimilar ones.

The most commonly used similarity metrics in CF are the Pearson Correlation Coefficient (PCC) and the Cosine Similarity (CS), where the PCC adjusts the ratings of a user  $u$  by the mean value of all ratings entered by  $u$ , so as to tackle the issue that some users may rate items higher than others, while the Cosine Similarity (CS) does not [1,2,4].

However, a user’s rating behavior is co-shaped by the items he has interacted with and rated (e.g. a movie viewer’s rating behavior is defined by the movies that he has watched) at any specific time and therefore the future ratings are effectively biased by the formerly watched content; however this information is not taken into account by either of these metrics, while this also holds for all similarity metrics that have been proposed in the literature.

In this work, we introduce the *Common Item Rating Past Criterion (CIRPaC)*, which focuses on the sequence that each user’s ratings have been entered in the database. More specifically, in the process of formulating a prediction for the rating that user  $u$  would assign to item  $i$ , for each of  $u$ ’s NNs  $NN_{u,k}$  the algorithm checks the degree of similarity between (a) the set of items that  $u$  has rated and (b) the set of items that  $NN_{u,k}$  has rated *up to the point that*  $NN_{u,k}$  had entered his rating for item  $i$ . The higher the similarity of these two sets, the more the similarity between the two aforementioned users will be “rewarded” (increased), under the rationale that their ratings for item  $i$  (factual for  $NN_{u,k}$ , predicted for  $u$ ) will have been influenced by the same set of experiences.

To illustrate the concept of the *CIRPaC criterion*, let us consider the case where we want to predict the rating that user  $u_1$  would give to the item  $i_1$ , which corresponds to the historic period drama series “Tudors” (<http://www.imdb.com/title/tt0758790/>), in order to decide whether it should be recommended to him or not, and only two users, namely  $u_2$  and  $u_3$ , have watched and rated this series (in order to be used as NNs). All of our three users, have watched and rated the also historic period drama series “Game of Thrones”

(<http://www.imdb.com/title/tt0944947/>), however  $u_2$  has watched and rated “Tudors” *before* “Game of Thrones”, while  $u_3$  has followed the inverse order, i.e. he has watched and rated “Tudors” *after* “Game of Thrones” (resembling the case of our active user,  $u_1$ , who is asking a recommendation for it). Taking into account that the “Game of Thrones” series is considered by many people as the best historic period drama series of all times, achieving an IMDB score of 9.5/10, we expect that a user that has already seen “Game of Thrones” will rate “Tudors” by different standards than another user that has not. This example indicates that the rating that a user sets to an item is not rating history-independent; on the contrary it clearly depends on his experience on the items that the user has interacted with and rated up to that time. Therefore, in this paper we (1) introduce the concept of *CIRPaC*, which aims precisely at quantifying and exploiting the degree of similarity between the sets of items that two users had experienced up to the time of rating registration and prediction and (2) investigate how we can incorporate the *CIRPaC* into the rating prediction computation process, so as to leverage the prediction accuracy of CF recommender systems.

To validate our approach, we present an extensive comparative evaluation among:

1. the proposed algorithm,
2. the rating abstention interval-based algorithm presented in [26], which (i) is a state-of-the-art algorithm exploiting temporal, within user history information, to achieve prediction error reduction in the context of CF-based rating predictions and (ii) has been shown to surpass the performance of other state-of-the-art algorithms. It is noted here that the algorithm presented in [26] necessitates the existence of data regarding the interaction among users within social networks and also exhibits a drop in coverage (i.e. loses some potential to compute personalized recommendations for users),
3. the dynamic average-based algorithm presented in [5], which (i) is a state-of-the-art algorithm targeting improvement of prediction accuracy in the context of CF, (ii) does not need extra information, regarding users or items (e.g. item categories or user social relationships) and (iii) does not deteriorate the prediction coverage,
4. the plain CF algorithm,

considering both the PCC and CS metrics. Rating prediction accuracy is measured using both the Mean Absolute Error (MAE) and Root-Mean-Square Error (RMSE) accuracy metrics, since both of them provide valuable insight on the performance of the prediction algorithm: the MAE measures the average absolute deviation between a predicted rating and the user’s true rating, while the RMSE squares the errors before summing them, thus putting more emphasis on large errors [6].

Finally, it is worth noting that the proposed approach can be combined with other algorithms that have been proposed for improving performance, rating prediction accuracy and recommendation quality in CF-based systems, including clustering techniques [7,8], exploitation of social network (SN) data [9-11] or pruning of old user ratings [12-14].

The rest of the paper is structured as follows: section 2 overviews related work, while section 3 presents the proposed algorithm. Section 4 presents the algorithm tuning and evaluation

procedure and the results obtained and, finally, section 5 concludes the paper and outlines future work.

## II. RELATED WORK

The accuracy of CF-based systems is a topic that has attracted considerable research efforts and various features of the ratings database or external linked data sources have insofar been exploited for improving prediction accuracy [3,15,16]. Whitby et al. [17] propose a filtering technique that applies to both unfairly positive and unfairly negative ratings in Bayesian reputation system. It is based on a reputation system and integrates a reputation system’s filtering method, under the assumption that ratings provided by different raters on a given agent will follow more or less the same probability distribution.

Koren [18] proposes a new neighbourhood-based model, which is based on formally optimizing a global cost function and leads to improved prediction accuracy, while maintaining merits of the neighbourhood approach such as explainability of predictions and ability to handle new ratings (or new users) without retraining the model. In addition, he suggests a factorized version of the neighbourhood model, which improves its computational complexity while retaining prediction accuracy. Liu et al. [19] present a new user similarity model to improve the recommendation performance when only few ratings are available to calculate the similarities for each user. The model considers the local context information of user ratings, as well as the global preference of user behaviour.

Research has shown that exploiting time in the rating prediction computation can improve prediction accuracy, due to concept drift; concept drift is the phenomenon when the relation between the input data and the target variable changes over time [20]. Change of interests [20,21] is a typical example of concept drift. Koenigstein et al. [22] consider the temporal dimension in the context of RSs by capturing different temporal dynamics of music ratings, along with information from the taxonomy of music-related items; both these dimensions are exploited by a rich bias model. Minku et al. [23] present a new categorization for concept drift, separating drifts according to different criteria into mutually exclusive and non-heterogeneous categories. Moreover, they present a diversity analysis in the presence of different types of drifts and it shows that, before the drift, ensembles with less diversity obtain lower test errors.

Pruning techniques have also been proposed [7,8] to eliminate old-aged ratings from the database, that do not reflect users’ current likings, hence they contribute towards formulating predictions with high absolute errors.

Knowledge-based RSs constitute a different approach to rating prediction accuracy improvement. Margaritis et al [24] present a novel algorithm for making accurate knowledge-based leisure time recommendations to social media users. The proposed algorithm considers qualitative attributes of the places (e.g., price, service, atmosphere), the profile and habits of the user for whom the recommendation is generated, place similarity, the physical distance of locations within which places are located, and the opinions of the user’s influencers. Rodríguez et al. [25] present AKNOBAS, a Knowledge-based Segmentation Recommender System, which follows trends using Intelligent Clustering Techniques for Information Systems.

With the advent of SN, SN recommendation has received considerable research attention. Konstas et al. [27] investigate the role of SN relationships in developing a track recommendation system using CF and taking into account both the social annotations and friendships inherent in the social graph established among users, items and tags. Arazy et al. [28] outline a conceptual RS design within which the structure and dynamics of a SN contribute to the dimensions of trust propagation, source's reputation and tie strength between users, which are then taken into account by the system's prediction component to generate recommendations. Quijano-Sanchez et al. [29] enhance a content-based RS by including in the recommendation algorithm the trust between individuals, users' interaction and aspects of each user's personality. Margaritis et al. [10] show that more reliable and successful recommendations can be produced when utilizing distinct sets of influencers per interest category, instead of using a single set of influencers for every recommendation to be made. Another use of SN recommendation can be found in [30], where a novel recommendation algorithm is presented, which exploits data sourced from web services [37,38] provided by the Internet of Things in order to produce more accurate venue recommendations. Recently, users' ratings variability was included in the rating prediction computation process, through which rating prediction quality is improved [5]. Additionally, the work in [26] introduces the concept of rating abstention intervals, i.e. periods of rating inactivity on behalf of the users, which indicate a shift of interest. The computation of rating abstention intervals entails the exploitation of temporal, within-user history information. The algorithm presented in [26] additionally computes and exploits metrics regarding influence levels among users, on the basis of data regarding interaction between users in social networks. Combining these two features, the work in [26] achieves considerable rating prediction accuracy improvements, surpassing other state-of-the-art algorithms.

However, none of the works mentioned above considers the aspect of shared experiences prior to the rating of each item. The present paper fills this gap by presenting an algorithm that leverages the similarity score of the users who have this content in common and evaluates its performance using different user similarity metrics and datasets.

### III. THE PROPOSED ALGORITHM

In CF, predictions for a user  $U$  are computed based on a set of users who have rated items similarly with  $U$ ; this set of users is termed "near neighbors of  $U$ " (NNs). The similarity metric most widely used in CF-based systems is the Pearson correlation coefficient [5], where the similarity between two users  $U$  and  $V$  is expressed as:

$$\text{simP}(U, V) = \frac{\sum_k (r_{U,k} - \bar{r}_U) * (r_{V,k} - \bar{r}_V)}{\sqrt{\sum_k (r_{U,k} - \bar{r}_U)^2 * \sum_k (r_{V,k} - \bar{r}_V)^2}} \quad (1)$$

where  $k$  ranges over items that have been rated by both  $U$  and  $V$ , while  $\bar{r}_U$  and  $\bar{r}_V$  are the mean value or ratings entered by users  $U$  and  $V$ , respectively. Then, for user  $U$ , his NN users  $NN_U$  are chosen, selecting the users having the highest similarity values with  $U$ . Typically, only users having a positive similarity with  $U$  are considered for inclusion in  $NN_U$ . Similarly, the Cosine Similarity metric [5] is expressed as:

$$\text{simC}(U, V) = \frac{\sum_k (r_{U,k} * r_{V,k})}{\sqrt{\sum_k (r_{U,k})^2 * \sum_k (r_{V,k})^2}} \quad (2)$$

As can be seen in formulas (1) and (2), the PCC adjusts the ratings of each user  $u$  by the mean value of all  $U$ 's ratings, in order to tackle the issue that some users may follow more strict practices when entering ratings, while other users may be more lenient; on the other hand, CS does not make such an adjustment.

Afterwards, in order to compute a rating prediction  $p_{U,i}$  for the rating of user  $U$  on item  $i$ , formula (3) is employed:

$$p_{U,i} = \bar{r}_U + \frac{\sum_{V \in NN_U} \text{sim}(U, V) * (r_{V,i} - \bar{r}_V)}{\sum_{V \in NN_U} \text{sim}(U, V)} \quad (3)$$

The proposed algorithm modifies the prediction computation phase, by including provisions for taking into account the items' rating sequence within the users' rating sets. More specifically, formula (3) is modified as follows:

$$p_{U,i} = \bar{r}_U + \frac{\sum_{V \in NN_U} \text{sim}(U, V) * \text{CIRPaC\_bonus}(V, U, i) * (r_{V,i} - \bar{r}_V)}{\sum_{V \in NN_U} \text{sim}(U, V) * \text{CIRPaC\_bonus}(V)} \quad (4)$$

where the  $\text{CIRPaC\_bonus}(V, U, i)$  parameter is the weight-bonus assigned to the each NN user  $V$  of user  $U$ , depending on the similarity of the between (a) the set of items that  $U$  has rated and (b) the set of items that  $V$  has rated up to the point that  $V$  had entered  $r_{V,i}$  in the ratings database.

In more detail, the computation of the  $\text{CIRPaC\_bonus}(V, U, i)$  is performed as follows: let  $\text{hist}(U)$  and  $\text{hist}(V)$  be, respectively, the sets of items that users  $U$  and  $V$  have rated, and  $\text{comRat}(U, V) = \text{hist}(U) \cap \text{hist}(V) = \{x_1, x_2, \dots, x_l\}$  be the set of items that have been commonly rated by these users.  $\text{comRat}(U, V)$  can be partitioned into two disjoint subsets with respect to item  $i$  for which the rating prediction is computed: the first subset  $\text{pre}(V, U, i)$  contains all items that had been rated by user  $V$  before he rated item  $i$ , while correspondingly  $\text{post}(V, U, i)$  contains all items that had been rated by user  $V$  after he rated item  $i$ . Formally, sets  $\text{pre}(V, U, i)$  and  $\text{post}(V, U, i)$  are defined as follows:

$$\begin{aligned} \text{pre}(V, U, i) &= \{x \in \text{comRat}(U, V) : \text{timestamp}(r_{v,x}) < \text{timestamp}(r_{v,i})\} \\ \text{post}(V, U, i) &= \{x \in \text{comRat}(U, V) : \text{timestamp}(r_{v,x}) > \text{timestamp}(r_{v,i})\} \end{aligned} \quad (5)$$

Effectively,  $\text{pre}(V, U, i)$  corresponds to the shared experiences between users  $U$  and  $V$  that user  $V$  had perceived before his (factual) rating for item  $i$ , while  $\text{comRat}(U, V)$  reflects the shared experiences between users  $U$  and  $V$  that user  $U$  has perceived before his (predicted) rating of item  $i$ . We can measure the similarity of these sets using the Jaccard index [35], where the similarity of two sets  $A$  and  $B$  is calculated as follows:

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} \quad (6)$$

Considering that  $\text{comRat}(U, V) \supseteq \text{pre}(V, U, i)$ , equation (6) can be rewritten as:

$$\text{histSim}(V, U, i) = J(\text{comRat}(U, V), \text{pre}(V, U, i)) = \frac{|\text{pre}(V, U, i)|}{|\text{comRat}(U, V)|} \quad (7)$$

We note here that we have also experimented with other set similarity measures, such as the Sorensen Similarity Index [11],

and the results were in close agreement with those obtained while using the Jaccard index (differences were capped by 1.6% in all cases); hence, for conciseness purposes, we only report on the results obtained while using the Jaccard index.

Finally,  $CIRPaC\_bonus(V, U, i)$  is computed by multiplying the  $histSim(V, U, i)$  metric by a constant termed  $CIRPaC\_base$ ; the optimal value for this constant is determined experimentally, and the relevant experiments are reported in section IV. Since however  $CIRPaC\_bonus(V, U, i)$  is designed to be a weight amplification parameter for cases when users  $V$  and  $U$  have considerable common experiences before the rating of item  $i$ , the  $CIRPaC\_bonus(V, U, i)$  is bounded from below by the value of 1.0. Thus, formally,  $CIRPaC\_bonus(V, U, i)$  is computed as follows:

$$CIRPaC\_bonus(V, U, i) = \max(histSim(V, U, i) * CIRPaC\_base, 1.0) \quad (8)$$

In the next section, we investigate candidate  $CIRPaC\_base$  values, in order to identify the optimal setting for this parameter and assess the performance of the proposed algorithm.

#### IV. ALGORITHM TUNING AND PERFORMANCE EVALUATION

In this section, we report on the experiments that were designed to:

1. determine the optimal  $CIRPaC\_base$  parameter value, in order to tune the algorithm and
2. compute the prediction improvement achieved due to the consideration of the  $CIRPaC$  criterion.

In order to determine the optimal parameter value, we experimentally explored the parameter value solution space, by iteratively selecting parameter value assignments and examining the effect that the particular parameter value assignments have on rating prediction quality. To quantify rating prediction quality, we employed two widely used error metrics, namely the Mean Absolute Error (MAE), and the Root Mean Squared Error (RMSE). The use of two different metrics allows us to gain more extensive insight on the prediction accuracy achieved by each parameter setting, since the MAE metric handles all error scales in a uniform fashion, whereas the RMSE metric penalizes more severely larger errors.

To compute the algorithm's prediction error, in terms of MAE and RMSE, we exercised the standard "hide one" technique [2,5,6]: each user's last rating in the database was hidden and then its value was predicted on the basis of the values of other, non-hidden ratings. We also performed a second experiment where, for each user  $u$ , a random rating  $r_{u,x}$  was selected and hidden, while in parallel all ratings of user  $u$  that had been

entered after  $r_{u,x}$  were dropped; subsequently, the hidden rating was again predicted its value on the basis of the values of other, non-hidden ratings. The results obtained from the two experiments were in close agreement (the differences observed were less than 2% in all cases), therefore for conciseness purposes we report only on the results of the first experiment. All our experiments were run on seven datasets. Five of these datasets are obtained from Amazon [31,32] and two from MovieLens [33,34]; the Amazon datasets are relatively sparse, while the MovieLens datasets are relatively dense. We choose to test both sparse and dense datasets (a dataset  $DS$  is deemed to be very sparse if  $d(DS) \ll 1\%$ , where  $d(DS)$  is the density of the dataset, defined as  $d(DS) = \frac{\#ratings}{\#users * \#items}$  [4]), in order to establish that the proposed algorithm can be used in every dataset.

The seven datasets used in our experiments are summarized in Table I and have the following characteristics:

- they are up to date (published between 1996 and 2016),
- they are widely used for benchmarking in CF research,
- they contain each rating's timestamp, which is necessary for the operation of the proposed algorithm and
- they differ in regards to the type of item domain of the dataset (videogames, movies, music and books) and size (ranging from 2 MB to 486 MB in plain text format).

Each dataset was initially preprocessed, and users found to have less than 10 ratings were dropped, since predictions formulated for users with few ratings are known to demonstrate high error levels [2,3]. This procedure did not have any effect on the MovieLens dataset, since it includes only users that have submitted at least 20 ratings. It is worth noting that we repeated the same experiments with datasets where users having between 5 and 10 ratings were retained, in order to gain insight on the proposed algorithm's behavior under contexts more akin to a "cold start" situation. In these contexts, the absolute average prediction error expectedly increased, however the error reduction levels were found to be in close agreement with those reported in the following subsections ( $\pm 0.5\%$  of the gains reported for the respective cases where users had at least 10 ratings each). These findings indicate that the proposed algorithm can be useful in cold start contexts, however further investigation on this aspect is needed; in our future work we will elaborate on this aspect.

For our experiments we used a machine equipped with six Intel Xeon E7 - 4830 @ 2.13 GHz CPUs, 256 GB of RAM and one 900 GB HDD with a transfer rate of 200 MBps, which hosted the datasets and ran the rating prediction algorithms.

TABLE I. DATASETS SUMMARY

Dataset name	#Users	#Items	#Ratings	Avg. #Ratings / User	Density	DB size (in text format)
Amazon "Videogames" [31,32]	8.1K	50K	157K	19.4	0.039%	4 MB
Amazon "CDs and Vinyl" [31,32]	41.2K	486K	1.3M	31.6	0.006%	32 MB
Amazon "Movies and TV" [31,32]	46.4K	134K	1.3M	28.0	0.021%	31 MB
Amazon "Books" [31,32]	295K	2.33M	8.7M	29.5	0.001%	227 MB
Amazon "Digital Music" [31,32]	6.2K	35K	86K	13.9	0.040%	2 MB
MovieLens "Latest 100K – Recommended for education and development" [33,34]	700	9K	100K	142.8	1.587%	2 MB
MovieLens "Latest 20M – Recommended for new research" dataset [33,34]	138K	27K	20M	144.9	0.537%	486 MB

In the remainder of this section, we present and discuss the results obtained from applying the algorithm presented above on these seven datasets.

#### A. The Amazon “Videogames” dataset

Fig. 1 illustrates the effect of the value of the *CIRPaC\_base* parameter on the quality of rating predictions produced by the proposed algorithm, under both similarity metrics, as this is reflected by the MAE and the RMSE error metrics, using the plain CF algorithm’s performance as a yardstick. Considering the quality of the rating predictions, the best performance for both similarity metrics is achieved when the *CIRPaC\_base* parameter is set to 130% (giving thus a maximum of 30% bonus). Under this setting, when using the PCC the MAE drops by 4.35% and the RMSE by 4.6%. As far as the CS metric is concerned, the respective reductions are 3.99% and 3.64%.

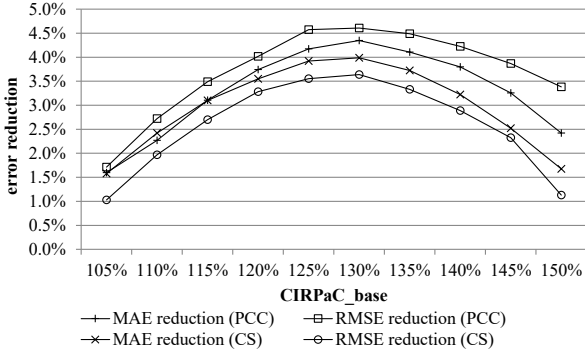


Fig. 1. MAE and RMSE reduction achieved by the proposed algorithm, under different *CIRPaC\_base* values and under both similarity metrics for the Amazon “Videogames” dataset.

#### B. The Amazon “CDs and Vinyl” dataset

Fig. 2 depicts the effect of the *CIRPaC\_base* parameter value on the quality of rating predictions produced by the proposed algorithm, under both similarity metrics, as this is reflected by the MAE and the RMSE error metrics. Again, the plain CF algorithm’s performance is used as a baseline. Regarding rating prediction quality, the best performance for both similarity metrics is achieved when the *CIRPaC\_base* parameter is set to 125%; under this setting, when using the PCC the MAE drops by 5.54% and the RMSE by 4.84%. The respective reductions are for the CS metric are 4.16% and 3.76%.

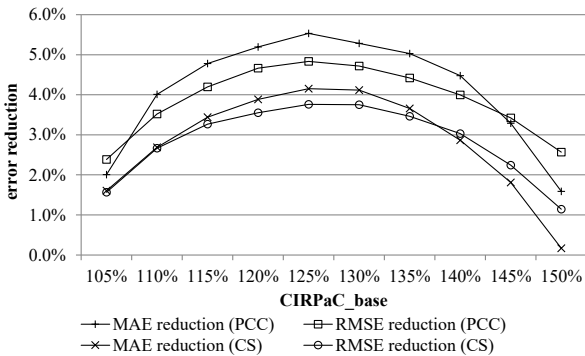


Fig. 2. MAE and RMSE reduction achieved by the proposed algorithm, under different *CIRPaC\_base* values and under both similarity metrics for the Amazon “CDs and Vinyl” dataset.

#### C. The Amazon “Movies and TV” dataset

Fig. 3 displays the effect of the value of the *CIRPaC\_base* parameter on the quality of rating predictions computed by the proposed algorithm, for both similarity metrics, as this is reflected by the MAE and the RMSE metrics. The plain CF algorithm’s performance is used as a baseline. Considering the quality of the rating predictions, the best performance for both metrics is achieved when the *CIRPaC\_base* parameter is set to 125%. Under this configuration, when using the PCC the MAE drops by 2.38% and the RMSE by 2.23%. As far as the CS metric is concerned, the respective reductions are 2% and 2.04%.

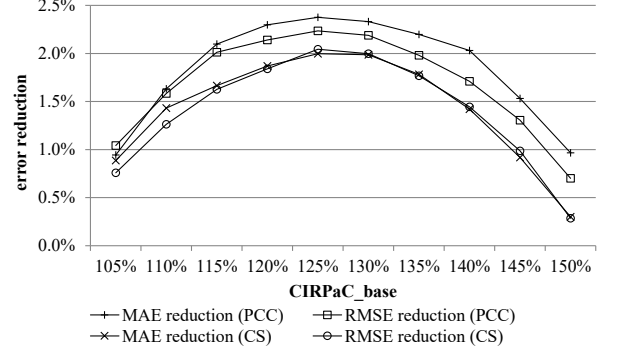


Fig. 3. MAE and RMSE reduction achieved by the proposed algorithm, under different *CIRPaC\_base* values and under both similarity metrics for the Amazon “Movies and TV” dataset.

#### D. The Amazon “Books” dataset

Fig. 4 illustrates the effect of the *CIRPaC\_base* parameter value on the quality of rating predictions formulated by the proposed algorithm, for both similarity metrics, as this is reflected by the MAE and the RMSE error metrics. The plain CF algorithm’s performance is used as a yardstick. Regarding the quality of the rating predictions, the best performance for both similarity metrics is achieved when the *CIRPaC\_base* parameter is set to 125%; under this setting, when using the PCC the MAE drops by 3.88% and the RMSE by 3.01%. As far as the CS metric is concerned, the reductions are 2.88% and 2.0%, respectively.

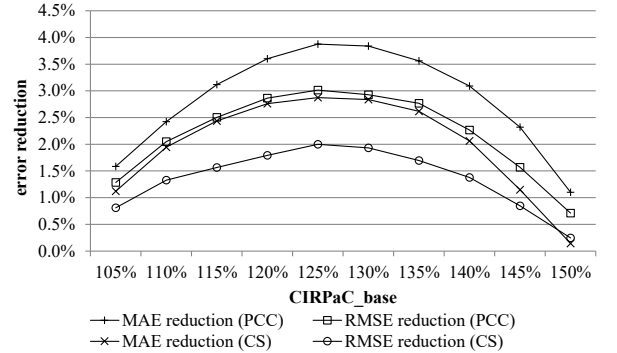


Fig. 4. MAE and RMSE reduction achieved by the proposed algorithm, under different *CIRPaC\_base* values and under both similarity metrics for the Amazon “Books” dataset.

#### E. The Amazon “Digital Music” dataset

Fig. 5 depicts the effect of the *CIRPaC\_base* parameter value on the quality of rating predictions produced by the proposed algorithm, for both similarity metrics, as this is manifested by

the MAE and the RMSE error metrics. The performance of the plain CF algorithm is used as a yardstick.

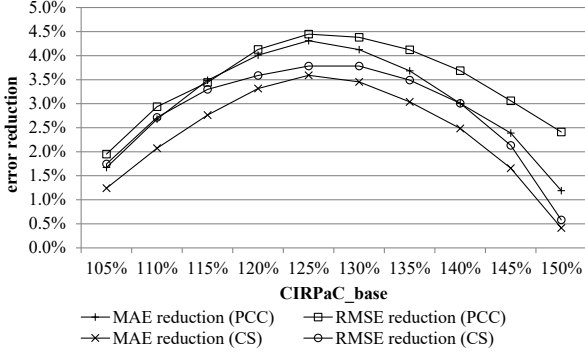


Fig. 5. MAE and RMSE reduction achieved by the proposed algorithm, under different *CIRPaC\_base* values and under both similarity metrics for the Amazon “Digital Music” dataset.

Regarding rating prediction quality, the best performance for both similarity metrics is achieved when the *CIRPaC\_base* parameter is set to 125%; when this setting is applied, under the PCC the MAE drops by 4.3% and the RMSE by 4.44%; under the CS metric, the respective reductions are 3.59% and 3.78 %.

#### F. The MovieLens “Latest 100K – Recommended for education and development” dataset

Fig. 6 shows the effect of the value of the *CIRPaC\_base* parameter on the quality of rating predictions computed by the proposed algorithm, for both similarity metrics, as this is reflected by the MAE and the RMSE error metrics. Again, the plain CF algorithm’s performance is used as a baseline.

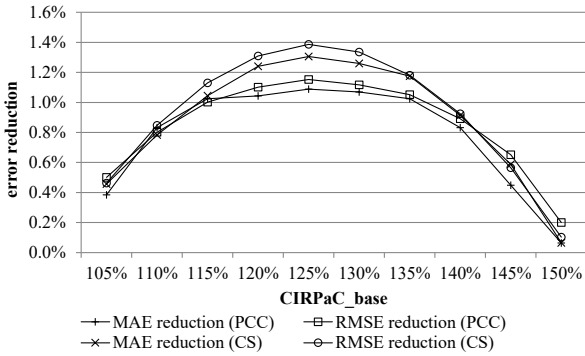


Fig. 6. MAE and RMSE reduction achieved by the proposed algorithm, under different *CIRPaC\_base* values and under both similarity metrics for the MovieLens “Latest 100K – Recommended for education and development” dataset.

Considering the quality of the rating predictions, the best performance for both similarity metrics is achieved when the *CIRPaC\_base* parameter is set to 125%. Under this configuration, when using the PCC the MAE drops by 1.09% and the RMSE by 1.15%. The respective reductions for the CS metric are 1.31% and 1.39%.

#### G. The MovieLens “Latest 20M – Recommended for new research” dataset

Fig. 7 demonstrates the effect of the value of the *CIRPaC\_base* parameter on the quality of rating predictions computed by the proposed algorithm, for both similarity metrics, as

this is expressed by the MAE and the RMSE metrics. The performance of the plain CF algorithm is used as a baseline.

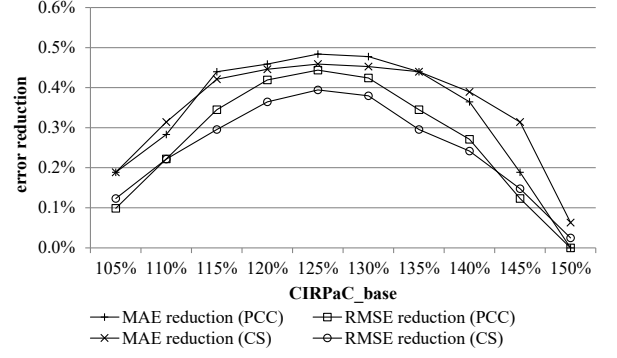


Fig. 7. MAE and RMSE reduction achieved by the proposed algorithm, under different *CIRPaC\_base* values and under both similarity metrics for the MovieLens “Latest 20M – Recommended for new research” dataset.

Regarding the quality of the rating predictions, the best performance for both similarity metrics is achieved when the *CIRPaC\_base* parameter is set to 125%; under this setting, when using the PCC the MAE drops by 0.48% and the RMSE by 0.44%. The corresponding reduction under the CS similarity metric are 0.46% and 0.39%.

#### H. Results overview and comparison with previous work

In this section, we overview the results presented in the previous paragraphs and we compare these results with the ones produced by the CF variability algorithm, proposed in [5]. This algorithm was chosen for the comparison since it (i) is a state-of-the-art algorithm targeting improvement of prediction accuracy in the context of CF, (ii) does not need extra information, regarding the users or the items (e.g. item categories or user social relationships) and (iii) does not deteriorate the prediction coverage.

Considering the optimal value of the *CIRPaC\_base* parameter, based on the results presented in the previous subsections, we can clearly see that it lays around 125%-130%. More specifically the setting 125% proved to be the optimal one in 6 out of the 7 datasets tested, while in the remaining one (the Amazon “Videogames” dataset), the 130% setting proved to be the best. On a global average, the setting of 125% has a performance edge over the setting of 130%, ranging from 1.5% (MAE reduction under the CS similarity metric) to 2.3% (RMSE reduction under the PCC metric); hence in the next experiments the *CIRPaC\_base* parameter will be set to 125%.

From the result analysis in subsections IV.A to IV.G, we can observe that the proposed algorithm achieves higher reductions in sparse datasets rather than in dense ones. An initial analysis revealed that in the two dense datasets considered, there is a smaller probability that experience sequences among users coincide; however deeper analysis on this aspect is required; this aspect will be investigated in the context of our future work.

Fig. 8 depicts the improvement in the MAE achieved by the proposed algorithm, when compared to the CF variability algorithm, proposed in [5], taking the performance of the plain CF algorithm as a baseline and using the PCC as the similarity metric, since this is the one tested in [5]. Clearly, the proposed algorithm achieves the best results, in all the datasets tested, with its

MAE reduction being 39.4% higher than that achieved by the CF variability algorithm (3.15% against 2.26% in absolute figures). At individual dataset level, the performance edge of the proposed algorithm against the CF variability algorithm ranges from 13% to 128%.

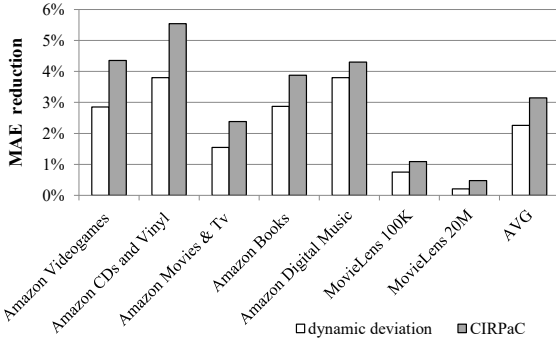


Fig. 8. MAE reduction achieved by the proposed algorithm, in comparison to the Dynamic Deviation Algorithm, proposed in [5].

Fig. 9 depicts the respective improvement in the RMSE achieved by the proposed algorithm, when compared to the CF variability algorithm, proposed in [5], again taking the performance of the plain CF algorithm as a baseline, again using the PCC as the similarity metric. Again, the proposed algorithm clearly achieves the best results, in all the datasets tested, with its MAE reduction being 100% higher than that achieved by the CF variability algorithm (2.96% against 1.48% in absolute figures). At individual dataset level, the performance edge of the proposed algorithm against the CF variability algorithm ranges from 39% to 175%.

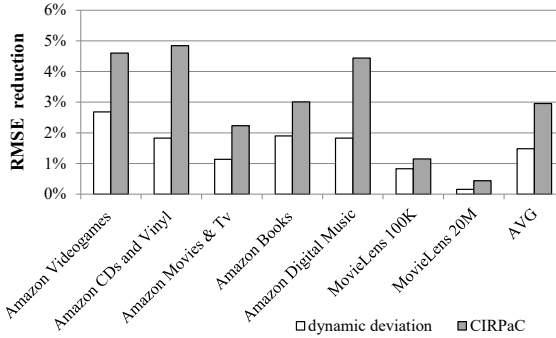


Fig. 9. RMSE reduction achieved by the proposed algorithm, in comparison to the Dynamic Deviation Algorithm, proposed in [5].

Finally, we compare the performance of the proposed algorithm against the algorithm presented in [26], which is a state-of-the-art algorithm exploiting temporal, within user history information, to achieve prediction error reduction in the context of CF-based rating predictions, and has also been shown to surpass the performance of other state-of-the-art algorithms. The proposed algorithm achieves an average MAE improvement of 3.15% over all tested datasets, while the respective gains of the algorithm presented in [26] are 2.99%. While the relative difference is limited to 5.4%, it is stressed here that the algorithm presented in [26] requires and exploits additionally information from social networks regarding the influence levels among users, which are not always available. Additionally, the algorithm presented in [26] exhibits a coverage drop which is considerable in the context of sparse datasets; on the other hand, the proposed

algorithm fully maintains the coverage levels. It is worth also noting that the algorithm presented in [26] has been shown to surpass the performance of other state-of-the-art algorithms, such as the ones in [36] and [13].

## V. CONCLUSION AND FUTURE WORK

In this paper we introduced the *Common Item Rating Past Criterion (CIRPaC)*, which considers the effect that items already experienced by a user have on the ratings that he assigns to other items. We have also proposed an algorithm which includes this criterion in the rating prediction computation process, in order to improve prediction accuracy.

The proposed algorithm has been validated through a set of experiments, using two user similarity metrics and seven datasets, both sparse and dense. These experiments showed that the inclusion of items' rating sequence introduces considerable prediction accuracy gains. More specifically, the experiment results have shown that the proposed algorithm delivers a significant MAE reduction, ranging from 0.48% to 5.54%, with an average of 3.15%, and a RMSE reduction, ranging from 0.44% to 4.84% with an average of 2.96%, as far as the PCC metric is concerned. The respective average error reductions, when the CS metric is used, are 2.62% and 2.43% (in all the above percentages, the plain CF algorithm is used as a baseline).

We have also compared the performance of the proposed algorithm against two other state-of-the-art algorithms targeting prediction error reduction, and the proposed algorithm has exhibited superior performance against both of them. More specifically in the comparison against the user rating variability algorithm [5], the proposed algorithm has proved to consistently outperform the user rating variability algorithm across all datasets tested, by margins ranging from 13% to 175%. In the comparison against the algorithm presented in [26], which exploits temporal, within-user history information, the proposed algorithm again proved to achieve better results (by 5.4% on average), even though the algorithm presented in [26] exploits additional information from social networks regarding the influence levels among users, which are not always available.

The proposed algorithm can be directly incorporated in a CF-based RS, since (1) it needs no extra information about the users or the items, (2) it needs minimal additional dataset pre-processing time, computing only the items' rating sequence within each user rating set, (3) it needs no extra storage space (4) it is easy to implement, through the modification of existed CF-based systems and (5) it can be combined with other algorithms that have been proposed for improving rating prediction accuracy and/or coverage.

Our future work will focus on exploring alternative algorithms for reducing prediction error in CF datasets. Furthermore, we are planning to evaluate the algorithm's performance under more user similarity metrics, such as Spearman coefficient and Euclidean distance [6]. Adaptation of the proposed approach for use with matrix factorization techniques [18] is also considered. Both can be utilized in broader applications of prediction methods [39-41]. Finally, the combination of the proposed technique with other algorithms that have been proposed for improving rating prediction accuracy, recommendation quality or prediction coverage in CF-based RSs will be examined.



## REFERENCES

- [1] M. Balabanovic and Y. Shoham, "Fab: content-based, collaborative recommendation," *Communications of the ACM*, vol. 40(3), pp. 66-72, 1997.
- [2] M. Ekstrand, R. Riedl and J. Konstan, "Collaborative Filtering Recommender Systems," *Foundations and Trends in Human-Computer Interaction*, vol. 4(2), pp. 81-173, 2011.
- [3] K. Yu, A. Schwaighofer, V. Tresp, X. Xu and H.P. Kriegel, "Probabilistic Memory-Based Collaborative Filtering," *IEEE Transactions on Knowledge Data Engineering*, vol. 16(1), 56-69, 2004.
- [4] J.B. Schafer, D. Frankowski, J. Herlocker and S. Sen, "Collaborative Filtering Recommender Systems," *The Adaptive Web, Lecture Notes in Computer Science*, vol. 4321, pp. 291-324, 2007.
- [5] D. Margaritis and C. Vassilakis, "Improving Collaborative Filtering's Rating Prediction Accuracy by Considering Users' Rating Variability," *Proceedings of the 4th IEEE International Conference on Big Data Intelligence and Computing*, pp. 1022-1027, 2018.
- [6] J. Herlocker, J. Konstan, L. Terveen and J. Riedl, "Evaluating collaborative filtering recommender systems," *ACM Transactions in Information Systems*, vol. 22(1), pp. 5-53, 2004.
- [7] S. Gong, "A Collaborative Filtering Recommendation Algorithm Based on User Clustering and Item Clustering," *Journal of Software*, vol. 5(7), pp. 745-752, 2010.
- [8] D. Margaritis, P. Georgiadis and C. Vassilakis, "A Collaborative Filtering Algorithm with Clustering for Personalized Web Service Selection in Business Processes," *Proceedings of the 9th IEEE International Conference on Research Challenges in Information Science*, pp. 169-180, 2015.
- [9] E. Bakshy, D. Eckles, R. Yan and I. Rosenn, "Social Influence in Social Advertising: Evidence from Field Experiments," *Proceedings of the 13th ACM Conference on Electronic Commerce*, pp. 146-161, 2012.
- [10] D. Margaritis, C. Vassilakis and P. Georgiadis, "Recommendation information diffusion in social networks considering user influence and semantics," *Social Network Analysis and Mining*, vol. 6(1), 108, pp. 1-22, 2016.
- [11] T.A. Sorensen, "A method of establishing groups of equal amplitude in plant sociology based on similarity of species content, and its application to analyses of the vegetation on Danish commons," *K dan Vidensk Selsk Biol Skr* 5, pp. 1-34, 1948.
- [12] D. Margaritis and C. Vassilakis, "Pruning and aging for user histories in collaborative filtering," *Proceedings of the 2016 IEEE Symposium Series on Computational Intelligence*, pp. 1-8, 2016.
- [13] D. Margaritis and C. Vassilakis, "Enhancing User Rating Database Consistency through Pruning," *Transactions on Large-Scale Data and Knowledge-Centered Systems*, vol. XXXIV, pp. 33-64, 2017.
- [14] D. Margaritis and C. Vassilakis, "Improving Collaborative Filtering's Rating Prediction Quality in Dense Datasets, by Pruning Old Ratings," *Proceedings of the 22nd IEEE Symposium on Computers and Communications*, pp. 1168-1174, 2017.
- [15] R. Dias and M. Fonseca, "Improving Music Recommendation in Session-Based Collaborative Filtering by Using Temporal Context," *Proceedings of the 25th IEEE International Conference on Tools with Artificial Intelligence*, pp. 783-788, 2013.
- [16] D. Margaritis, C. Vassilakis and P. Georgiadis, "Query personalization using social network information and collaborative filtering techniques," *Future Generation Computer Systems*, vol. 78(1), pp. 440-450, 2018.
- [17] A. Whitby, A. Jøsang and J. Indulska, "Filtering Out Unfair Ratings in Bayesian Reputation Systems," *Proceedings of the Workshop on Trust in Agent Societies, at the Autonomous Agents and Multi Agent Systems Conference (AAMAS2004)*, vol. 4, 2004.
- [18] Y. Koren, R. Bell and C. Volinsky, "Matrix factorization techniques for recommender systems," *Computer*, vol. 42(8), pp. 30-37, 2009.
- [19] H. Liu, Z. Hu, A. Mian, H. Tian and X. and Zhu, "A new user similarity model to improve the accuracy of collaborative filtering," *Knowledge-Based Systems*, vol. 56, pp. 156-166, 2014.
- [20] J. Gama, I. Zliobaite, A. Bifet, M. Pechenizki and A. Bouchachia, "A Survey on Concept Drift Adaptation," *ACM Computing Surveys*, vol. 1(1), Article 1, 2013.
- [21] L. Li, L. Zheng, F. Yang and T. Li, "Modeling and broadening temporal user interest in personalized news recommendation," *Expert Systems with Applications*, vol. 41 (7), pp. 3168-3177, 2014.
- [22] N. Koenigstein, G. Dror and Y. Koren, "Yahoo! Music recommendations: modeling music ratings with temporal dynamics and item taxonomy," *Proceedings of the fifth ACM conference on Recommender systems (RecSys '11)*, pp. 165-172, 2011.
- [23] L. Minku, A. White and X. Yao, "The Impact of Diversity on Online Ensemble Learning in the Presence of Concept Drift," *IEEE Transactions on Knowledge and Data Engineering*, vol. 22(5), pp. 730-742, 2010.
- [24] D. Margaritis, C. Vassilakis and P. Georgiadis, "Knowledge-Based Leisure Time Recommendations in Social Networks," *Current Trends on Knowledge-Based Systems: Theory and Applications*, pp. 23-48, 2017.
- [25] A. Rodríguez, J. Torres, E. Jimenez, M. Gomez and G. Alor-Hernandez, "AKNOBAS: A knowledge-based segmentation recommender system based on intelligent data mining techniques," *Computer Science and Information Systems*, vol. 9(2), pp. 713-740, 2012.
- [26] D. Margaritis and C. Vassilakis, "Exploiting Rating Abstention Intervals for Addressing Concept Drift in Social Network Recommender Systems," *Informatics*, vol. 5(2), 21, 2018.
- [27] I. Konstas, V. Stathopoulos and J.M. Jose, "On social networks and collaborative recommendation," *Proceedings of the 32nd international ACM SIGIR conference on Research and development in information retrieval*, pp. 195-202, 2009.
- [28] O. Arazy, N. Kumar and B. Shapira, "Improving Social Recommender Systems," *IT professional*, vol. 11(4), 2009.
- [29] L. Quijano-Sanchez, J.A. Recio-Garcia and B. Diaz-Agudo, "Group recommendation methods for social network environments," *3rd Workshop on Recommender Systems and the Social Web at the 5th ACM International Conference on Recommender Systems*, pp. 1-24, 2011.
- [30] D. Margaritis and C. Vassilakis, "Exploiting Internet of Things Information to Enhance Venues' Recommendation Accuracy," *Service Oriented Computing & Applications*, vol. 11(4), pp. 393-409, 2017.
- [31] Amazon product data. Available online: <http://jmcauley.ucsd.edu/data/amazon/links.html> (accessed on April 4, 2019).
- [32] J.J. McAuley, C. Targett, Q. Shi and A. Van den Hengel, "Image-Based Recommendations on Styles and Substitutes," *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 43-52, 2015.
- [33] MovieLens datasets. Available online: <http://grouplens.org/datasets/movielens/> (accessed on April 4, 2019).
- [34] F. Harper and J. Konstan, "The MovieLens Datasets: History and Context," *ACM Transactions on Interactive Intelligent Systems*, vol. 5(4), Article no. 19, 2016.
- [35] P.N. Tan, M. Steinbach and V. Kumar, "Introduction to Data Mining," Addison-Wesley, 2005.
- [36] Y. Koren, "Collaborative Filtering with Temporal Dynamics," *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 447-456, 2009.
- [37] D. Margaritis, C. Vassilakis and P. Georgiadis, "Adapting WS-BPEL scenario execution using collaborative filtering techniques," *Proceedings of the 7th IEEE International Conference on Research Challenges in Information Science*, pp. 174-184, 2013.
- [38] D. Margaritis, C. Vassilakis and P. Georgiadis, "An integrated framework for QoS-based adaptation and exception resolution in WS-BPEL scenarios," *Proceedings of the 28th ACM Symposium on Applied Computing*, pp. 1900-1906, 2013.
- [39] D. Antonakaki, D. Spiliotopoulos, C.V. Samaras, S. Ioannidis and P. Fragopoulou, "Investigating the Complete Corpus of Referendum and Elections Tweets," *Proceedings of the IEEE/ACM Conference on Advances in Social Networks Analysis and Mining*, pp. 100-105, 2016.
- [40] G. Scheffbeck, D. Spiliotopoulos and T. Risse, "The Recent Challenge in Web Archiving: Archiving the Social Web," *Proceedings of the International Council on Archives Congress*, pp. 20-24, 2012.
- [41] D. Antonakaki, D. Spiliotopoulos, C.V. Samaras, P. Pratikakis, S. Ioannidis and P. Fragopoulou, "Social media analysis during political turbulence," *PloS one*, vol. 12(10), pp. 1-23, 2017.