

Privacy Preserving Collaborative Filtering with k -Anonymity through Microaggregation

Fran Casino*, Josep Domingo-Ferrer*, Constantinos Patsakis[‡], Domènec Puig* and Agusti Solanas*[†]

*UNESCO Chair in Data Privacy

Dept. of Computer Engineering and Mathematics
Universitat Rovira i Virgili

Av.Països Catalans 26, 43007 Tarragona, Catalonia

[‡] Distributed Systems Group

School of Computer Engineering and Statistics
Trinity College, Dublin, Ireland

[†]Corresponding author e-mail: agusti.solanas@urv.cat

Abstract—Collaborative Filtering (CF) is a recommender system which is becoming increasingly relevant for the industry. Current research focuses on Privacy Preserving Collaborative Filtering (PPCF), whose aim is to solve the privacy issues raised by the systematic collection of private information. In this paper, we propose a new microaggregation-based PPCF method that distorts data to provide k -anonymity, whilst simultaneously making accurate recommendations. Experimental results demonstrate that the proposed method perturbs data more efficiently than the well-known and widely used distortion method based on Gaussian noise addition.

Keywords—Privacy Preserving Collaborative Filtering, Microaggregation, Electronic Commerce, Recommender Systems, Statistical Disclosure Control.

I. INTRODUCTION

Most people tend to take into serious consideration other people's recommendations, when they want to buy a product. The evaluation process prior to the purchase of a product is difficult when the opinions we get are diverse and/or from multiple sources. This is of particular relevance when the source of the recommendations is the Internet. On the one hand, the Internet provides a wealth of information on a huge variety of products and services that may be useful to potential buyers. On the other hand, this wealth of information may become a problem rather than a solution because it can hinder the decision making.

Recommender systems [1] derive from knowledge discovery in databases (KDD) [2], [3]. KDD systems are processes used by companies to analyse and discover understandable patterns within large collections of data that may help, for instance, to save money or to sell more products to customers. Collaborative Filtering (CF) [4], [5] is a recommender system which involves a large family of recommendation methods. The aim of CF is to make suggestions on items (*e.g.* books, music, or movies), based on the preferences of users that have already acquired and/or rated these items. CF appeared with the aim to provide automatic recommendations in a digital environment.

Currently CF is actively present throughout the Internet:

- It has a great impact on e-commerce (*e.g.* Amazon, Barnes & Noble).
- It is used to analyse information preferences. For instance, the websites visited by a user can be monitored and utilised to recommend similar sites to other users with similar interests.
- It is widely used in multimedia contexts, with sites like last.fm, MyStrands, Netflix and Moviefinder that leverage CF to benefit users in form of advice, and also provide a clear market overview to the companies.
- It is part of the *Web 2.0* concept, which is defined as the new way to use the Internet. This new trend gives high importance to active user participation in the infrastructure (*e.g.* blogs, social networks and information & service portals).

In order to make predictions, CF methods use large databases that store information regarding the relationships between sets of users and items. These data are modelled as matrices composed by n users and m items, and each cell (i, j) stores the evaluation of user i on item j . Therefore, a value is assigned, which can be within a range of values (*e.g.* between 0 and 10) or simply with binary votes (positive/negative, or bought/not bought) as in market basket databases. Table I shows a small example of a typical CF matrix. There are many examples of CF referenced databases [6] in the literature, like Eachmovie, MovieLens, Jester, and Netflix prize data. These databases are frequently used as benchmarks to evaluate the efficiency, quality and robustness of CF methods [7].

TABLE I. EXAMPLE OF A DATA MATRIX WHERE EACH ROW CORRESPONDS TO A USER U_i AND EACH COLUMN CORRESPONDS TO AN ITEM I_j . EACH MATRIX CELL STORES A VOTE/EVALUATION WITHIN A RANGE OF VALUES BETWEEN 1 AND 5. BLANKS REPRESENT ITEMS NOT YET RATED BY A PARTICULAR USER. ITEM 3, MARKED IN GRAY, HAS NOT BEEN EVALUATED YET.

	I_1	I_2	I_3	I_4	I_5	I_6
U_a	2	4			1	2
U_b		3		2		
U_c	3	2			2	2
U_d		5		1		1

CF methods are based on the assumption that similar users,

in the sense of similar interests or behaviours, will be interested in the same products. Therefore, items purchased by a user u_a can be recommended to another user u_b , if u_a and u_b have similar interest or similar behavior. Other approaches, invert this process, by making recommendations based on similarity between items. In this context, a group of similar users or items form a neighbourhood.

During the past years, the use of trust systems has witnessed an important increase on the Internet. A trust statement is defined as the explicit opinion expressed by a user on another user, regarding the perceived quality of certain features of that user [8]. The concept of trust is widely used, for instance, in search engines such as Google, which uses global trust metrics [9], and in e-commerce (*e.g.* ebay) in which users express their level of satisfaction after purchasing a product. By aggregating the trust statements expressed by every user, it is possible to produce a community or a neighborhood [10] as seen, for instance, in social networks. This way, information provided by networks of trust can be combined with CF matrices to create trust-aware recommender systems [8] that handle problems like data sparsity, cold start users and artificial identities, more efficiently.

CF methods can be classified into three main categories depending on the data they use to compute recommendations: memory-based methods, which use the full matrix with all ratings; model-based methods, which utilise statistical models and functions of the data matrix but not the complete data matrix; and hybrid methods, which combine the previous methods with content-based recommendation methods [6].

In memory-based CF, recommendations are made in two steps: (i) neighbourhood search and (ii) recommendation prediction. Given a user u_a , correlation and distance functions are used to compute her neighborhood. The most common correlation and distance functions used are the Pearson Correlation [5], the Cosine similarity [11] and the Euclidean distance. The similarity between users can also be computed in a much more efficient way, according to their behaviour when they vote. Examples of these are shown in [12], where the tendencies of users are calculated, or in [13], where the concordances between users are computed with privacy. Once the neighborhood of the user u_a is determined, recommendations can be computed using, for instance, the methods described in [5], [14]. These methods can be utilised to predict a vote or recommend top-N items for u_a .

Model-based CF methods [2], [6], [12], [15] build a model from the full matrix on which to make predictions. The emergence of these methods is justified by the constraints that memory-based CF methods have regarding scalability, computational costs and sparseness.

Hybrid CF methods combine memory-based and model-based methods, preserving the advantages of the algorithms involved, whilst at the same time minimizing their drawbacks and deficiencies. Examples of these methods are the Personality Diagnosis [16] and the Probabilistic memory-based model [17]. Hybrid methods can also be obtained by the combination of memory-based and model-based methods with content-based recommender systems. Some well-known examples [6] are: Filterbots, Content-boosted, Fab and Ripper.

A. Contribution and Plan of the Article

Regardless of the CF method, there are several limitations inherent to this kind of recommender systems. Some of the most important limitations [6], [8], [18] are *sparseness*, *scalability*, *cold start*, *shilling*, *synonymy*, *bribing*, *copy-profile attacks*, and the lack of privacy.

Amongst all the aforementioned open problems, in this paper we concentrate on the protection of the privacy of users involved in CF processes. Note that privacy is gaining much interest in many fields since the increasing collection of personal data (*e.g.* in smart cities [19], or location-based services (LBS) [20]).

The main contribution of this paper is a new microaggregation-based Privacy Preserving Collaborative Filtering method that guarantees k -anonymity to the users. The proposed method has proven to be more efficient in terms of privacy protection and information loss than a widely used PPCF method such as Gaussian noise addition. In addition to improve over Gaussian noise addition, our proposal remains simple, hence, easing its wide adoption.

The rest of the article is organised as follows. § II introduces basic concepts on PPCF and classifies the most relevant methods. Also, some background on statistical disclosure control and microaggregation is given. Next, in § III, we present our new PPCF method based on microaggregation. In § IV, we show the results of our approach and in § V we discuss the benefits of our approach by comparing it against the Gaussian noise addition method. Finally, § VI concludes the paper and provides directions for future research.

II. BACKGROUND

A. Privacy Preserving Collaborative Filtering

The widespread use of CF on the Internet provides great opportunities and benefits to both companies and users. However, a major drawback is introduced, the lack of user's privacy. The relevance of privacy in CF systems is emphasised by the growing pace at which information of each user is collected and stored. The careless management of personal information, might be illegal in many countries, but furthermore has severe consequences for the users whose information is disclosed, as well as companies. One of the main problems in CF is that customers believing that their preferences/profiles may be exposed, decide either not to give their assessment on a particular item or to do it incorrectly or inaccurately [21]. Therefore, the feeling of lack of privacy results in a reduction of both the number of evaluations as well as their quality. Privacy Preserving Collaborative Filtering (PPCF) methods aim at solving the privacy issues raised by the systematic collection of private information, which are required for the proper use of CF methods.

In dynamic markets, companies may be interested in cooperating to obtain better recommendations for their customers. Due to privacy and business concerns, data should not be disclosed between companies. In this context, data might be partitioned between various parties in different ways:

Vertical partitioning (VP) in which companies own disjoint sets of items but with the same users.

Horizontal partitioning (HP) in which different parties hold disjoint sets of users with opinions of the same items.

Arbitrary partitioning (AP) in which there is no pattern of how data are distributed. If the entire set is defined by an $m \times n$ user-item matrix, one party A holds a subset of users $m_a \leq m$ whilst another party B holds the rest $m_b = m - m_a$ and the same is applied for items.

There are several ways in the literature to protect privacy in databases. To anonymise profiles in large databases we can use methods to provide k -anonymity [22]. Due to their interesting properties, public key homomorphic cryptosystems [23], secure multi-party computations [24] and cryptographic protocols and voting systems [25] are often used on the Internet. Other methods add noise to the data so as to distort their values, in a way affecting as little as possible the statistical properties of the matrix, such as average ratings of the users and the items.

According to how information is stored and how recommendations are computed, we classify PPCF approaches into centralised or decentralised. We consider that a PPCF method is centralised if it utilises a third party to perform intermediate calculations between users or entities. A method is also considered to be centralised if ratings are stored in a single server where recommendations and predictions are computed. Cases in which data are partitioned are not considered as centralised because the data are distributed between different parties. Typically, centralised PPCF methods offer higher efficiency than their decentralised counterparts since similarity and prediction computations avoid several communication overheads. However, in centralised methods, data are managed by only one party that has total control over them, with the implied privacy issues if the data have a low protection level. Most centralised PPCF methods add noise in several ways to perturb the data. In [26] and [27] the authors propose the perturbation of the data by following an item-invariant Gaussian/uniform noise distribution. In the item-based PPCF scheme proposed by Zhang *et al.* in [28], perturbations are added depending on the importance of the ratings in the recommendation process.

On the other side, we have decentralised PPCF methods. These methods use the members of distributed networks, considered in most cases as users, to perform intermediate calculations and predictions. Decentralised schemes generally involve less information disclosure than their centralised counterparts, but they entail the use of expensive protocols and more complex calculations. Typically, in decentralised PPCF methods, users store their own ratings. This results in a series of shortcomings like the requirement of the active participation of users, which is needed to share their data and to perform intermediate calculations. If users are not active, the amount of data over which CF is applied decreases, drastically changing the accuracy of the recommendations. Well-known PPCF with partitioned data schemes, which involve different parties sharing their data to perform CF with more referrals, are also considered to be decentralised methods. Several approaches with partitioned market basket databases have been proposed in the literature. This kind of databases is suitable to make top-N recommendations with high accuracy and low computational cost due to its binary ratings contents. In the numerical rating context, we have several approaches with partitioned data

schemes like the ones proposed in [29], [30] and [31]. Finally, schemes in which users store their own ratings can be also found (e.g. [32], [33] and [34]).

B. Statistical Disclosure Control and Microaggregation

Statistical disclosure control [35], a.k.a. data anonymisation, seeks to transform microdata sets (*i.e.* datasets consisting of records corresponding to individual respondents) before publication in such a way that it is not possible to re-identify the respondent corresponding to any particular record in the anonymised published microdata set —identity disclosure— nor is it possible to discover the value of a confidential attribute (*e.g.* salary) for a *specific* respondent —attribute disclosure—. Prior to any anonymisation process, direct identifiers (name, passport no., etc.) need of course be suppressed from the dataset. However, some of the attributes that remain in the anonymised dataset may be *quasi-identifiers*, that is, attributes which may facilitate indirect re-identification of respondents through external data sources (available as intruders’ background knowledge) that combine those attributes with direct identifiers. Microaggregation is a family of anonymisation algorithms for datasets that works in two stages:

- 1) The set of records in a dataset is clustered in such a way that: i) each cluster contains at least k records; ii) records within a cluster are as similar as possible.
- 2) Records within each cluster are replaced by a representative of the cluster, typically the centroid record (*i.e.* the average of the cluster).

When microaggregation is applied to the projection of records on their quasi-identifier attributes, the resulting dataset is k -anonymous, that is, to an intruder each record in the dataset is indistinguishable within a group of k records in terms of the quasi-identifiers. In [36] a simple microaggregation heuristic called MDAV is described, in which all clusters have exactly k records, except the last one, which might have between k and $2k - 1$ records. More sophisticated microaggregation approaches can be found in [37], [38], [39], [40] and [41].

III. PROPOSED METHOD

In this section, we describe our proposal in detail. Our scheme could be classified as a centralised PPCF method. Its architecture is shown in Figure 1.

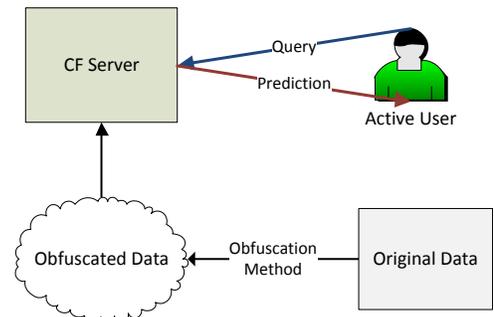


Fig. 1. Centralised PPCF scheme. The user makes a request on an item to the server, which responds with a personalised prediction.

Our approach is based on the aforementioned MDAV microaggregation algorithm and depicted in Figure 2. We slightly modify MDAV in the way the leftover records are managed. The original MDAV algorithm specifies that if at the end of the clustering process there are p records, with $k \leq p < 2k$, that do not belong to any group, they should form a final cluster C_f themselves. In our approach, in order to reduce the information loss, we first compute the mean of C_f , denoted by M_f , and we compute the distance between every C_f record with M_f . After that, we compare the distance between each member of C_f with all the already formed groups. If more than half of the records in C_f are closer to M_f than to any other group, we form a final cluster with the C_f elements. Otherwise, each record is clustered with the closest group amongst those already formed.

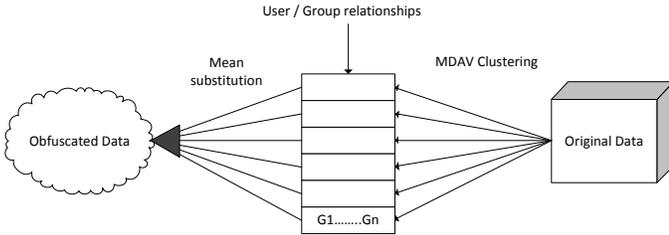


Fig. 2. MDAV clustering scheme. The steps flow from right (original dataset) to left (obfuscated dataset).

Our scheme, illustrated in Figure 3, works as follows:

- 1) Ensure that the dataset contains no missing values for any attribute in any record. This is necessary in order to compute the Euclidean distance between records. Imputation methods [42], [43] or non-personalised values can be utilised to fill the empty fields of the dataset matrix.
- 2) Once the matrix is completely filled, we compute the z-scores of each column (item) of the dataset, using the following expression

$$\text{z-score} = \frac{x_i - \mu}{\sigma}$$

where x_i is the i -th value of item x and μ and σ are the mean and the standard deviation of the item x , respectively. In this way, the mean and the standard deviation of the transformed item are 0 and 1, respectively.

- 3) Once the standardised matrix has been obtained, we are able to apply MDAV clustering. Users will be grouped into n clusters, with each cluster C_i consisting of the k most similar users, according to Euclidean distance. The value of k denotes the group cardinality. By selecting the most similar users, we maximise the group homogeneity and we therefore reduce the information loss. Once the group relationships are established, the mean values of each C_i , denoted as M_i , are computed. Afterwards, each value of C_i is replaced by the corresponding M_i .
- 4) The MDAV clustering process will result in a new dataset in which members of the same cluster C_i will have the same profiles. Therefore, after applying MDAV, this dataset will satisfy k -anonymity.

- 5) Finally, in order to make predictions, the results are de-standardised to obtain the final obfuscated dataset.

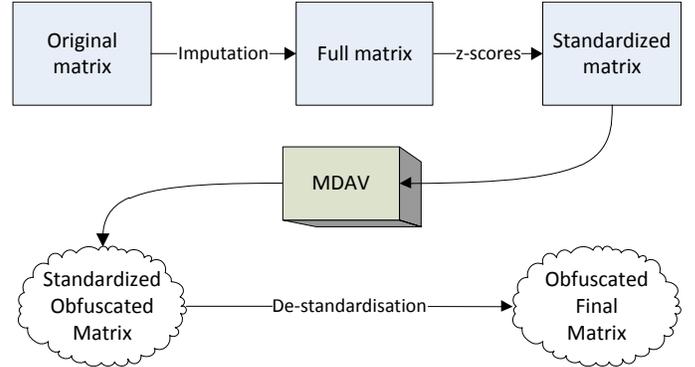


Fig. 3. Proposed method step by step.

IV. EXPERIMENTAL RESULTS

In this section, we show the experimental results of our method and we compare them with those obtained with the widely used Gaussian noise addition method (GNA), which uses a Gaussian distribution with zero mean and standard deviation σ (i.e. $\mathcal{N}(0, \sigma)$) to perturb the dataset.

Firstly, we show the results related to the privacy and the utility provided by the analysed methods in Section IV-A. Then, we assess the quality of the predictions in Section IV-B.

Experiments with GNA were repeated 50 times with each evaluated σ . As we already did in our proposal (*cf.* Figure 3), the dataset values are standardised before the Gaussian noise is added. In order to test the quality of our method, we will use the well-known Movielens dataset. Movielens was developed by Grouplens [5] and it is one of the reference sets in CF. Here, we will focus on *Movielens 100k*, which contains 100,000 ratings of 943 users on 1,682 films. The *Movielens 100k* range values are comprised between 1 and 5. This database is highly sparse, since more than 90% of the fields are empty. We establish the median of the range values (3) to fill the empty fields of the matrix. Once completely filled, the matrix contains a total of 1,586,126 values.

A. Privacy protection

In order to measure the quality of the privacy provided by a perturbation method we consider two factors, namely the information loss and the disclosure risk. The information loss is generally associated to the sum of squared errors (SSE). The SSE is commonly used as a measure of the distortion introduced on the original data. In the special case of microaggregation, the SSE is computed in vectorial notation as follows:

$$SSE = \sum_{i=1}^g \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)' (x_{ij} - \bar{x}_i) \quad (1)$$

where g is the number of subsets/clusters generated by the algorithm, n_i the number of elements in each cluster, x_{ij} the vector of the j -th user of the i -th cluster and \bar{x}_i is the average vector of the i -th cluster. More generally, given an original

dataset O represented by a matrix of $n \times m$ elements o_{ij} and a distorted/protected dataset P represented by a matrix of $n \times m$ elements p_{ij} , the SSE is computed as follows:

$$SSE = \sum_{i=1}^n \sum_{j=1}^m (o_{ij} - p_{ij})^2 \quad (2)$$

The disclosure risk (DR) measures the probability of rightly relating a record of the obfuscated/protected data matrix with a record of the original matrix. It is also known as the probability of re-identification, or the re-identification risk. For an attacker, the re-identification procedure consists in computing the distances (*e.g.* the Euclidean distance) between a given protected record p_i (corresponding to user i), and the target records o_j , that could be obtained from third party sources like census and the like. In our case we assume the best scenario for an attacker (the oracle scenario) in which he has the original dataset O and the distorted dataset P and he tries to link each record p_i in P with the records o_j in O .

For each record p_i in P the attacker determines the closest record o_j in O . If the determined closest record o_j is actually the original record belonging to p_i , the attacker succeeds and we say that p_i has been re-identified. To compute the disclosure risk, we try to re-identify all records and we compute the percentage of correct re-identifications. In terms of privacy and utility of the data, both the SSE and the DR should be low.

In the following tables and figures we show the results of the SSE and the DR for the analysed methods: our microaggregation-based method and the GNA. Table II shows the results obtained with our proposed method for different values of k , which represent the cardinality of the clusters, whilst Table III shows the results obtained using GNA with different values of σ . It can be clearly seen that the relation between SSE and DR is much better in the MDAV-based approach.

Figure 4 and Figure 5 respectively show the SSE and the DR of the MDAV-based PPCF for different values of k . It can be observed that their behavior is pretty antagonistic. When the SSE is increased, the DR is reduced accordingly.

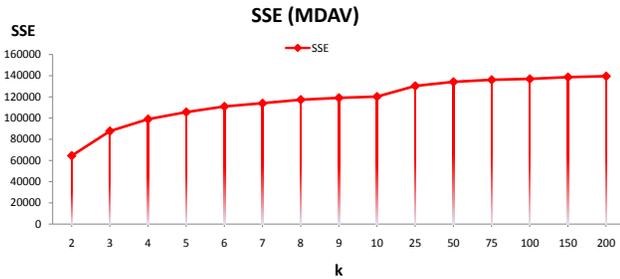


Fig. 4. SSE values of our method on *Movielens 100k* database.

Figure 6 and Figure 7 respectively show the SSE and DR for the GNA approach. Similarly to the MDAV-based approach, when the SSE grows the DR is decreased. However, as we discuss in Section V, the GNA method has to add much more distortion into the data (*i.e.* more SSE) than the MDAV-based approach to reach the same DR.

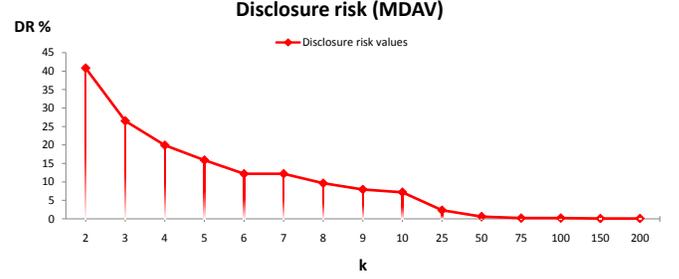


Fig. 5. DR values of our method on *Movielens 100k* database.

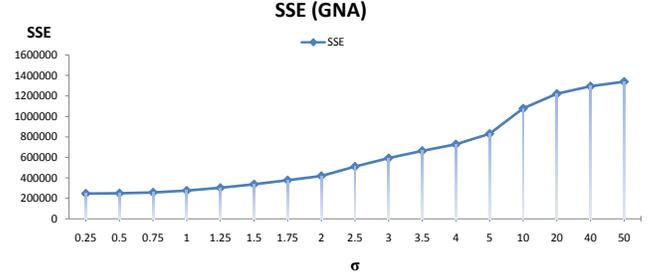


Fig. 6. SSE values of the GNA method on *Movielens 100k* database.

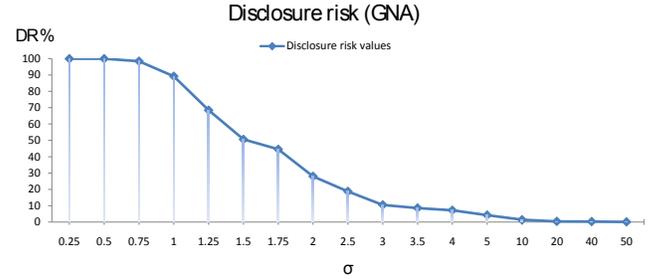


Fig. 7. DR values of the GNA method on *Movielens 100k* database.

B. Predictions accuracy

In the previous section, we have analysed the quality of the privacy and utility provided by the MDAV-based and the GNA approach. However, privacy and utility are just a dimension of the problem. Note that the protected data will be used by recommender systems to make predictions about which items a user would be more interested in. Thus, it is important not only to protect the privacy, but to provide accurate predictions too. In order to compare our proposal with the GNA method, we have selected obfuscated datasets with the same DR (*i.e.* in this case, for the sake of simplicity, we have selected DR = 7.21% because it corresponds to the value achieved with $k = 10$ for MDAV, and $\sigma = 4$ using GNA). Note that any other DR value could have been chosen as far as it is the same for both methods.

We have defined a training set with 80% of the item values and a prediction set with the remaining 20%. The predictions are computed only for the original values of each user. The prediction of the values is performed in two steps:

TABLE II. RESULTS OF MDAV BASED PPCF. THE SSE RESULTS ARE DISPLAYED IN A 10^3 SCALE.

ML 100k	MDAV														
k	2	3	4	5	6	7	8	9	10	25	50	75	100	150	200
SSE	64	87	99	105	110	114	117	119	120	130	134	136	136	138	139
DR %	40.82	26.51	19.93	15.9	12.19	12.19	9.65	7.95	7.21	2.33	0.63	0.21	0.21	0.1	0.1

TABLE III. RESULTS OF GNA BASED PPCF. THE SSE RESULTS ARE DISPLAYED IN A 10^3 SCALE.

ML 100k	GNA																
σ	0.25	0.5	0.75	1	1.25	1.5	1.75	2	2.5	3	3.5	4	5	10	20	40	50
SSE	246	248	257	275	302	336	376	418	509	592	663	727	830	1078	1221	1294	1339
DR%	100	100	98.51	89.28	68.5	50.58	44.53	27.99	18.76	10.49	8.58	7.21	4.24	1.4	0.42	0.31	0.1

- Find closest neighbor: Given a user u_i for which we want to predict some of its values, we consider the training set and we find its closest neighbour, say u_j .
- Assign his/her values: The predicted values for user u_i are those that correspond to u_j in the prediction set.

Once the prediction for all users is done, we compute the error between the values of the original dataset (*i.e.* those in the test set 20%) and the predicted values. To compute this error we apply the widely used mean absolute error (MAE), defined as follows:

$$MAE = \frac{\sum_{i=1}^n |p_i - r_i|}{n} \quad (3)$$

where n is the number of predicted elements, p_i is the predicted value over the element i , and r_i is the real value of i . The results are shown in Table IV.

TABLE IV. MAE OBTAINED VALUES, COMPARING THE PREDICTIONS MATRICES WITH THE ORIGINAL DATASET. THE % MAE HAS BEEN COMPUTED WITH RESPECT TO THE MOVIELENS RANGE VALUES (I.E 1-5).

Method	MAE	% MAE
MDAV, $k = 10$	0.89	22.25
GNA, $\sigma = 4$	1.08	27

V. COMPARISON AND DISCUSSION

In the previous section we presented the results obtained by our microaggregation-based approach and the classical GNA approach. In this section we briefly compare those results and show that the MDAV-approach is superior, both in terms of privacy and prediction accuracy.

In Figure 8, we can see a comparison between the SSE and the DR values of both methods. In the X-axis we represent the DR and in the Y-axis we show the SEE. This figure can be used to read the amount of noise, in terms of SSE, that is required by each method to achieve a given DR. For example, it can be seen that for a fixed value of $DR = 30\%$ the MDAV-approach roughly introduces an error in the order of $100K$ whilst the GNA approach requires $400K$.

The smallest possible DR value for the analysed dataset is $\frac{1}{943} \simeq 0.1\%$. In order to obtain that value, the MDAV-approach needs to form groups of $k = 150$ elements, which leads to

an SSE of 138,650. In contrast, the GNA obtains such DR value with an SSE of 1,339,008, which is almost one order of magnitude larger.

These results clearly show that the proposed approach perturbs data in a much more efficient way. Moreover, as already mentioned, our method ensures the users' privacy providing k -anonymity.

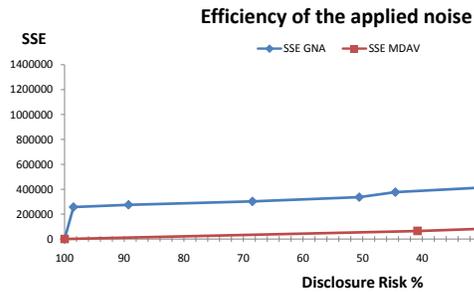


Fig. 8. Relation between SEE and DR for the analysed methods on the *Movielens 100k* database.

Regarding the quality of the prediction, Table IV shows the accuracy of the predicted values. It can be seen that when the predictions are done using the data protected with MDAV, the MAE is 22.25% with a DR value of 7.21%, which is a considerable privacy level. On the contrary, the values predicted by using the data protected with the GNA lead to an error of 27%, which is almost a 5% higher. Therefore, we may conclude that both the quality of the predictions and the quality of the privacy is better in the proposed method based on MDAV.

VI. CONCLUSIONS AND FUTURE WORK

Collaborative Filtering is a recommender system used to perform automatic recommendations to users in multiple contexts. Despite the great advantages of using CF, we have highlighted the important impact that it might have on users' privacy. Although a large amount of CF methods have been proposed, their study is still needed and there are many challenges to overcome. Probably, the most significant amongst them is proper protection of users' privacy.

Definitely, protecting users' privacy by hiding as much as possible their preferences, has an inherent trade-off, the quality of the recommendations decreases. Therefore, in this

article, we have proposed a new PPCF method based on microaggregation. The results obtained over the evaluated database demonstrate that the proposed method perturbs data in a much more efficient way than other well-known methods such as GNA. Moreover, our proposal achieves k -anonymity, which increases the users' privacy in such a way that GNA cannot guarantee.

Future work will focus on two different directions. The first one is to improve the efficiency of our method in order to be implemented in a decentralised scheme. The second direction is to analyse the influence of imputation methods on privacy and recommendations quality, then study networks of trust and efficient imputation policies.

ACKNOWLEDGMENTS AND DISCLAIMER

This work was partly supported by the Government of Catalonia under grant 2009 SGR 1135, by the Spanish Government through projects TIN2011-27076-C03-01 "CO-PRIVACY" and CONSOLIDER INGENIO 2010 CSD2007-00004 "ARES", by the European Commission under FP7 projects "DwB" and "Inter-Trust", and by La Caixa under program RECERCAIXA 2012. The authors are with the UNESCO Chair in Data Privacy, but they are solely responsible for the views expressed in this paper, which do not necessarily reflect the position of UNESCO nor commit that organisation. The second author is partially supported as an ICREA Acadèmia researcher by the Government of Catalonia.

REFERENCES

- [1] P. Resnick and H. Varian, "Recommender systems," *Communications of the ACM*, vol. 40, no. 3, pp. 56–58, 1997.
- [2] B. Sarwar, G. Karypis, J. Konstan, and J. Riedl, "Application of dimensionality reduction in recommender system—a case study," *ACM WebKDD 2000 Web Mining for ECommerce Workshop*, vol. 1625, no. 1, pp. 264–8, 2000.
- [3] U. Fayyad and G. Piatetsky-Shapiro, "Advances in knowledge discovery and data mining," *AAAI/MIT Press*, 1996.
- [4] D. Goldberg, D. Nichols, B. M. Oki, and D. Terry, "Using collaborative filtering to weave an information tapestry," *Communications of the ACM*, vol. 35, no. 12, pp. 61–70, 1992.
- [5] P. Resnick, N. Iacovou, and M. Suchak, "GroupLens: an open architecture for collaborative filtering of netnews," *Proceedings of the ACM conference on Computer supported cooperative work CSCW*, vol. pp. no. 3, pp. 175–186, 1994.
- [6] X. Su and T. M. Khoshgoftaar, "A Survey of Collaborative Filtering Techniques," *Advances in Artificial Intelligence*, vol. 2009, no. Section 3, pp. 1–19, 2009.
- [7] J. L. Herlocker, J. a. Konstan, L. G. Terveen, and J. T. Riedl, "Evaluating collaborative filtering recommender systems," *ACM Transactions on Information Systems*, vol. 22, no. 1, pp. 5–53, Jan. 2004.
- [8] P. Massa and P. Avesani, "Trust-aware collaborative filtering for recommender systems," *On the Move to Meaningful Internet Systems 2004: CoopIS, DOA, and ODBASE*, 492-508., vol. 3290, no. 8, pp. 492–508, 2004.
- [9] —, "Trust metrics on controversial users: balancing between tyranny of the majority," *International Journal on Semantic Web and Information Systems*, pp. 1–21, 2007.
- [10] J. Golbeck, "FilmTrust: Movie Recommendations from Semantic Web-based Social Networks," in *ISWC2005 Posters Demonstrations*, 2006, pp. 1314–1315.
- [11] J. Breese, D. Heckerman, and C. Kadie, "Empirical Analysis of Predictive Algorithms for Collaborative Filtering," *UAI 98*, pp. 43–52, 1998.
- [12] F. Ccheda, V. Carneiro, D. Fernández, and V. Formoso, "Comparison of collaborative filtering algorithms," *ACM Transactions on the Web*, vol. 5, no. 1, pp. 1–33, Feb. 2011.
- [13] N. Lathia, S. Hailes, and L. Capra, "Private Distributed Collaborative Filtering Using Estimated Concordance Measures," *Proceedings of the 2007 ACM conference on Recommender systems RecSys 07*, p. 1, 2007.
- [14] B. Sarwar, G. Karypis, J. Konstan, and J. Riedl, "Item-based collaborative filtering recommendation algorithms," *Proceedings of the 10th international conference on World Wide Web*, vol. 1581133480, no. 15, pp. 285–295, 2001.
- [15] Z. Xia, Y. Dong, and G. Xing, "Support vector machines for collaborative filtering," in *Proceedings of the 44th annual Southeast regional conference*, ser. ACM-SE 44. New York, NY, USA: ACM, 2006, pp. 169–174.
- [16] D. Y. Pavlov and D. M. Pennock, "A Maximum Entropy Approach to Collaborative Filtering in Dynamic, Sparse, High-Dimensional Domains," *Advances in Neural Information Processing Systems 15*, vol. 15, no. 2/3, pp. 1441–1448, 2003.
- [17] A. Schwaighofer, V. Tresp, and H. Kriegel, "Probabilistic memory-based collaborative filtering," *IEEE Transactions on Knowledge and Data Engineering*, vol. 16, no. 1, pp. 56–69, Jan. 2004.
- [18] C. Kaleli and H. Polat, "Providing naïve Bayesian classifier-based private recommendations on partitioned data," *Knowledge Discovery in Databases: PKDD 2007*, pp. 515–522, 2007.
- [19] A. Martínez-Ballesté, P. A. Pérez-Martínez, and A. Solanas, "The pursuit of citizens' privacy: a privacy-aware smart city is possible," *IEEE Communications Magazine*, vol. 51, no. 6, 2013.
- [20] A. Solanas and A. Martínez-Ballesté, "A ttp-free protocol for location privacy in location-based services," *Computer Communications*, vol. 31, no. 6, pp. 1181–1191, 2008.
- [21] L. Cranor, J. Reagle, and M. Ackerman, "Beyond concern: Understanding net users' attitudes about online privacy," *The Internet Upheaval: Raising Questions, Seeking Answers in Communications Policy*, Tech. Rep., 2000.
- [22] L. Sweeney, "k-anonymity: A model for protecting privacy," *International Journal of Uncertainty, Fuzziness and Knowledge-based Systems*, vol. 10, no. 5, pp. 557–570, 2002.
- [23] P. Paillier, "Public-Key Cryptosystems Based on Composite Degree Residuosity Classes," *Advances in Cryptology EUROCRYPT 99*, vol. 1592, pp. 223–238, 1999.
- [24] A. C. Yao, "Protocols for secure computations," *23rd Annual Symposium on Foundations of Computer Science (sfcs 1982)*, pp. 160–164, Nov. 1982.
- [25] B. Adida and D. Wikström, "How to shuffle in public," in *Artificial Intelligence*. Springer-Verlag, 2007, pp. 555–574.
- [26] H. Polat, "Privacy-preserving collaborative filtering using randomized perturbation techniques," *Third IEEE International Conference on Data Mining*, pp. 625–628, 2003.
- [27] H. Polat and L. Hall, "SVD-based Collaborative Filtering with Privacy," *Proceedings of the 2005 ACM symposium on Applied computing SAC 05*, pp. 791–795, 2005.
- [28] S. Zhang, J. Ford, and F. Makedon, "A privacy-preserving collaborative filtering scheme with two-way communication," *Proceedings of the 7th ACM conference on Electronic commerce - EC '06*, pp. 316–323, 2006.
- [29] I. Yakut and H. Polat, "Arbitrarily distributed data-based recommendations with privacy," *Data & Knowledge Engineering*, vol. 72, pp. 239–256, Feb. 2012.
- [30] —, "Privacy-Preserving Svd-Based Collaborative Filtering on Partitioned Data," *International Journal of Information Technology & Decision Making*, vol. 09, no. 03, pp. 473–502, May 2010.
- [31] J. Zhan, I. Wang, and C. Hsieh, "Towards efficient privacy-preserving collaborative recommender systems," *GrC 2008*, pp. 778–783, 2008.
- [32] J. Canny, "Collaborative filtering with privacy," *Security and Privacy, 2002. Proceedings. 2002 IEEE*, pp. 45–57, 2002.
- [33] S. Berkovsky, F. Ricci, Y. Eytani, and T. Kuflik, "Enhancing Privacy and Preserving Accuracy of a Distributed Collaborative Filtering," *Proceedings of the 2007 ACM conference on Recommender systems RecSys 07*, pp. 9–16, 2007.

- [34] C. Kaleli and H. Polat, "P2P collaborative filtering with privacy," *Turkish Journal of Electric Electrical Engineering and Compute Science*, vol. 18, no. 1, pp. 101–116, 2010.
- [35] A. Hundepool, J. Domingo-Ferrer, L. Franconi, S. Giessing, E. Schulte-Nordholt, K. Spicer, and P.-P. de Wolf, *Statistical Disclosure Control*. Wiley, 2012.
- [36] J. Domingo-Ferrer and V. Torra, "Ordinal, continuous and heterogeneous k -anonymity through microaggregation," *Data Mining and Knowledge Discovery*, vol. 11, no. 2, pp. 195–212, 2005.
- [37] A. Solanas and A. Martínez-Ballesté, "V-MDAV: Variable group size multivariate microaggregation," in *COMPSTAT 2006*, 2006, pp. 917–925.
- [38] A. Solanas and R. D. Pietro, "A linear-time multivariate microaggregation for privacy protection in uniform very large data sets," in *MDAI*, 2008, pp. 203–214.
- [39] A. Solanas, A. Martínez-Ballesté, and U. González-Nicalás, "A variable-MDAV-based partitioning strategy to continuous multivariate microaggregation with genetic algorithms," in *International Joint Conference on Neural Networks(IJCNN)*, 2010, pp. 1–7.
- [40] A. Solanas, A. Gavaldà, and R. Rallo, "Micro-som: A linear-time multivariate microaggregation algorithm based on self-organizing maps," in *ICANN (1)*, 2009, pp. 525–535.
- [41] A. Solanas, A. Martínez-Ballesté, and J. M. Mateo-Sanz, "Distributed architecture with double-phase microaggregation for the private sharing of biomedical data in mobile health," *IEEE Transactions on Information Forensics and Security*, vol. 8, no. 6, pp. 901–910, 2013.
- [42] A. Dempster, N. Laird, and D. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *Journal of the Royal Statistical Society Series BMethodological*, vol. 39, no. 1, pp. 1–38, 1977.
- [43] J. Montaquila and C. Ponikowski, "An evaluation of alternative imputation methods," in *Proceedings of the Survey Research Methods*, 1995.