

Efficient Privacy-Preserving Matrix Factorization via Fully Homomorphic Encryption

[Extended Abstract] *

Sungwook Kim¹ Jinsu Kim¹ Dongyoung Koo²
sw14.kim@samsung.com jinsu86.kim@samsung.com dykoo@nslab.kaist.ac.kr

Yuna Kim¹ Hyunsoo Yoon² Junbum Shin¹
yuna1.kim@samsung.com hyoon@kaist.ac.kr junbum.shin@samsung.com

¹Software R&D Center, Samsung Electronics, Seoul, South Korea

²School of Computing, Korea Advanced Institute of Science and Technology (KAIST), Daejeon, South Korea

ABSTRACT

Recommendation systems become popular in our daily life. It is well known that the more the release of users' personal data, the better the quality of recommendation. However, such services raise serious privacy concerns for users. In this paper, focusing on matrix factorization-based recommendation systems, we propose the first privacy-preserving matrix factorization using fully homomorphic encryption. On inputs of encrypted users' ratings, our protocol performs matrix factorization over the encrypted data and returns encrypted outputs so that the recommendation system knows nothing on rating values and resulting user/item profiles. It provides a way to obfuscate the number and list of items a user rated without harming the accuracy of recommendation, and additionally protects recommender's tuning parameters for business benefit and allows the recommender to optimize the parameters for quality of service. To overcome performance degradation caused by the use of fully homomorphic encryption, we introduce a novel data structure to perform computations over encrypted vectors, which are essential operations for matrix factorization, through secure 2-party computation in part. With the data structure, the proposed protocol requires dozens of times less computation cost over those of previous works. Our experiments on a personal computer with 3.4 GHz 6-cores 64 GB RAM show that the proposed protocol runs in 1.5 minutes per iteration. It is more efficient than Nikolaenko et al.'s work proposed in CCS 2013, in which it took about 170 minutes on two servers with 1.9 GHz 16-cores 128 GB RAM.

*Full version available at
<http://eprint.iacr.org/2016/232>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ASIA CCS '16, May 30-June 03, 2016, Xi'an, China

© 2016 Copyright held by the owner/author(s). Publication rights licensed to ACM. ISBN 978-1-4503-4233-9/16/05...\$15.00

DOI: <http://dx.doi.org/10.1145/2897845.2897875>

CCS Concepts

•Security and privacy → Privacy-preserving protocols; •Theory of computation → Cryptographic protocols;

Keywords

homomorphic encryption; matrix factorization; gradient descent; privacy

1. INTRODUCTION

A lot of recommendation services, including mobile ads and movie recommendations, have been used in our daily life, and a collaborative filtering (CF) is one of commonly used algorithms for such services [2, 36]. The principle of CF algorithm is to infer user's preference from gathered other users' history and to give a recommendation for new items similar users prefer.

To address privacy issues caused by CF [26, 31, 40], there have been numerous works to propose privacy-preserving CF algorithms [6, 21, 24, 28, 32]. It has been discovered that anonymization of users' data is not sufficient to protect the privacy [6, 26, 37]. Differential privacy [21, 24, 32] has been proposed for minimizing the chances of identifying any single user's data. However it takes trade-off between privacy and accuracy of estimated prediction and does not offer confidentiality of user's rating values by definition.

Recently, a remarkable approach based on cryptographic techniques, that is garbled circuits, was proposed by Nikolaenko et al. [28]. They targeted a *matrix factorization*, mainly used among several CF algorithms. It addresses how to protect both user's rating values and which items are rated by the user. However it seriously requires improvement of efficiency and utility for practical use. In addition, it only supports a fixed number of iteration of matrix factorization, which does not guarantee the quality of matrix factorization.

In this paper, we focus our interest on how to efficiently perform matrix factorization to compute user and item profiles from the ratings given by users in a privacy-preserving way with full functionality. We develop a fully homomorphic encryption-based matrix factorization protocol among users, a recommendation system (*RecSys*), and crypto-service provider (*CSP*).

We revisit the following three security requirements for secure matrix factorization for practical use: (i) *Privacy*, (ii) *Accuracy*, and (iii) *Completeness*. Privacy consists of considerations for both users and the RecSys. On a user’s side, it typically comprises a user’s rating values, which items a user has rated, the number of items a user has rated, and user/item profiles. On RecSys’s side, tuning parameters for matrix factorization can be considered as sensitive information, because the efficacy of the outputs depends on the parameters that the RecSys chooses [38]. Accuracy points out that a privacy protection mechanism should not weaken the accuracy of the recommendation as far as possible. Completeness means a mechanism needs to support a full learning process for recommendation, repeatedly training and validation until it converges. During the process, the tuning parameters can be optimized for faster convergence [8, 43].

To the best of our knowledge there is no known mechanism for privacy-preserving matrix factorization which meets all of previously mentioned requirements. In this work we initiate a study in the direction of efficient privacy-preserving matrix factorization to take all requirements into account.

1.1 Contributions

Main contributions and features of the proposed protocol are as follows:

New construction from fully homomorphic encryption. We first construct a privacy-preserving matrix factorization via *fully homomorphic encryption* (FHE). Our construction performs matrix factorization to compute user/item profiles using approximation techniques (e.g. gradient descent) as described in Section 3.3. Our protocol runs a secure two-party computation protocol between the RecSys and CSP, when it is too burdensome to perform the above operations by FHE solely. The efficient use of FHE is quite non-trivial; authors in [28] introduced garbled circuit instead of FHE, pointing out that FHE schemes for simpler algebraic computations are not as efficient as garbled circuit approaches.

The main obstacles for gradient descent computation using FHE are (i) the inner product operation of vectors according to the index, which is induced from the user-item matrix and (ii) the fixed point operation due to vectors over real numbers. For inner product operation, a naive use of Single Instruction Multiple Data (SIMD) operation supported by FHE [34] does not help for reducing computation cost since it gets to perform one homomorphic multiplication for one inner product operation only, not multiple inner products, after all (refer to Section 4.1 for details). For the fixed point arithmetic, it requires the division-like operation over encrypted data and as far as we know there is no efficient method to evaluate it homomorphically. Thus naive approach using FHE does not work efficiently in the sense that the computation and communication costs linearly depend on the number of ratings.

We resolve the above limitations of FHE in gradient descent by introducing a novel data structure. It enables to enjoy the full use of slots supported by the SIMD operation in FHE ciphertext while allowing one homomorphic operation for multiple operations over vectors during gradient descent. As a result, both the running time and the communication overhead of the proposed protocol are linear in M/L , where M and L are the number of ratings and slots which FHE supports, respectively. Since L is typically a

number in thousands in our implementation, the data structure reduces computation cost by dozens of times. For the fixed point arithmetic on encrypted data, we use secure 2-party computation between the RecSys and CSP. That is, the CSP performs fixed point arithmetic over masked plaintexts without learning any information about actual plaintexts and then returns the encrypted results to the RecSys. We believe that our technique can be generally applied to various privacy-preserving data analytics protocols using gradient descent.

Additionally, we provide an optimization for the ciphertext size of FHE itself. In most FHE constructions [10, 11, 16], the size of ciphertext grows quadratically in the size of message space, and it causes large computational and communication overhead. To obtain small ciphertexts for the case of large message space, we represent the message in several smaller spaces using Chinese remainder theorem. This technique yields linear growth in the ciphertext size.

Privacy and security. Our protocol guarantees the secrecy of rating values from users and all intermediate computation results during joint computation between the RecSys and CSP. The final results, user/item profiles, are also returned under FHE encryption so that the RecSys and CSP cannot learn both user and item profiles in the clear. On RecSys’s side, the optimized parameters for quality recommendation can be kept secret, because the RecSys itself can set tuning parameters for matrix factorization without disclosure from the use of FHE.

Finally the extension of the proposed protocol delivers the ability of launching various techniques to obfuscate the number of items and the list of items a user rated. Among the obfuscation techniques, we describe how to apply an injection of fake ratings considered in differential privacy approach [9, 32], and the description can be helpful for applying another techniques easily.

Accuracy and completeness. As mentioned previously, the use of differential privacy or the injection of fake ratings suffer from distortion of results of data analysis. We ensure the accuracy of the recommendation by introducing the encrypted indicator of rated items. The indicator sets 0 or 1 by users, depending on whether the item is rated or not. Through the protocol, the indicator vector works correctly without revealing the value at all to both the RecSys and CSP. If the rating is real, it gets to involve to the computation for analysis. Otherwise its influence on the computation is removed.

Typically, matrix factorization iteratively computes user/item profiles using gradient descent. The iteration ends up with convergence of the profiles. In our protocol, the RecSys can check the convergence by verifying a stopping criteria without knowing the raw values of users’ ratings or user/item profiles. Thus it provides the ability to complete a learning process by optimizing the gradient descent parameters.

Efficiency. We implement our protocol and evaluate the performance. We use the MovieLens dataset for our experiment. Our matrix factorization protocol without obfuscation of rating information takes about 1.5 minutes with 580 MB communication between the RecSys and CSP per iteration for 14K ratings dataset. This improves execution time and communication overhead drastically over the work in [28], which takes 2.9 hours with 128 GB. The extended pro-

toCol takes 1.9 minutes with 623 MB for 10% fake ratings, and 2.5 minutes with 910 MB for 50%.

We further estimate performance with taking into account parallelism and compare to recent progress on parallelization of Nikolaenko et. al’s protocol using GraphSC [27] and optimization technique of circuits for the garbled circuit protocol called TinyGarble [35]. The estimate shows that parallelized version of our construction is expected to be about 33 times faster than that of Nikolaenko et al.’s protocol together with TinyGarble.

2. RELATED WORK

Our work is to present an efficient and practical matrix factorization for collaborative filtering based-recommendation while preserving user privacy about item ratings and the analyzed preference throughout the two-party protocol between the RecSys and CSP.

Privacy with cryptography. A lot of works have been done to solve the privacy issues for various machine learning algorithms, and some of their basic operations can be protected using cryptography efficiently: classification [4], regression [3, 29].

The first work for privacy-preserving matrix factorization is proposed by Nikolaenko et. al. [28], which used a mixed protocol of additive homomorphic encryption [7, 30] and garbled circuits [42] for secure two-party computation, and applied oblivious sorting network for supporting sparsity of rating data. They tried to improve efficiency of the implementation using multi-threaded FastGC [22], and parallelization of the sorting network. However, the excessive computation and communication cause impracticality, for example, it takes 2.9 hours per iteration for 14K ratings. The performance comparisons between [28] and ours will be given in Section 8.

Nayak et al. [27] proposed more practical approach, named GraphSC, by applying both a GraphLab [9] for a parallelism and oblivious approach [18] for the security. They have shown that the previous work proposed by Nikolaenko et al. [28] can be extended using GraphSC and showed that it can make overall parts of secure matrix factorization parallel on large dataset of 1M ratings, which is a large improvement from the original one [28]. Such a parallelism for our work is analyzed in Section 8.3.

Besides matrix factorization, crypto-based privacy protection techniques have been applied to other collaborative filtering methods [13, 39]. Erkin et al. [13] applied homomorphic encryption to user-based collaborative filtering using cosine similarity between users’ ratings. For reducing computational and communication overhead, multiple numerical values are packed in a single ciphertext. It differs from our data packing in that their packing can be applied only to vector-scalar multiplications, not batch inner products needed for matrix factorization.

Veugen et al. [39] proposed a generic framework for secure multi-party computations based on secret sharing, especially focusing on the recommendation systems. It provides security in the malicious model, where one of servers may not follow the protocol specifications and may control a couple of users to deduce more personal data. It also focuses on efficient online recommendation when a user asks for recommendation with updated ratings. Differently from their work, we focus on efficient computation of matrices for user

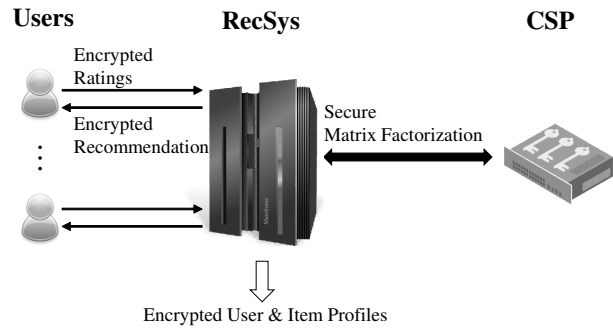


Figure 1: Three Parties in Our Protocol Design.

and item profiles without leakage of private data before recommendation.

Privacy with non-cryptography. Differential privacy is mainly considered for non crypto-based privacy protection [12, 25], including the recommendation privacy [21, 24, 32, 41]. It prevents attacker from the inference of users’ privacy by adding fake users’ ratings as long as it does not affect significantly to the distribution of the system’s output. However, differential privacy approach does not guarantee confidentiality of users’ data to the RecSys, and reduce accuracy of recommendation as much as fake ratings added. In our work, users can add fake ratings along with indicator vectors, 1 for real rating and 0 for fake to obfuscate which items are rated. This does not affect the accuracy of recommendation, since the fake ratings are removed when matrix factorization runs.

3. PRELIMINARIES

3.1 Setting

The proposed protocol and system consists of three ends: Users, the Recommendation System (RecSys), and Crypt-Service Provider (CSP), as shown in Figure 1. User’s rating information is uploaded to the RecSys, which provides a recommendation service based on users’ uploaded information. In our model, we will use a matrix factorization algorithm to infer the value of un-rated item: the details will be explained in Section 3.3.

Users’ rating information will be encrypted using a public key homomorphic encryption, so it is very important to protect such a key. Therefore, we will use the CSP, which is secure server and the key will be stored and used only inside the CSP to minimize the risk of key exposure.

Regarding the security, we assume the standard adversary model used in [28]: Following the definition of honest-but-curious adversary model [17], either the RecSys or CSP can be compromised by an adversary, but they will follow the prescribed protocol correctly. In addition, the adversary might look at the data exchanged with the other parties, and try to uncover users’ data including the ratings and to infer recommendation and matrix factorization results including both user and item profiles. However, the RecSys does not misbehave with users’ ratings such as altering, dropping, or copying them to control the output. The CSP can also always keep the secrets. We also assume the RecSys and CSP do not collude, and such an assumption makes sense in

business environment because they value the reputation. In addition, the CSP usually could be a governmental organization guarding the privacy protection of users [39].

3.2 Tools and Notations

In the protocol we use somewhat (possibly fully) homomorphic encryption and additively homomorphic encryption schemes, denoted by HE and AHE, respectively. With public keys hpk for HE, we call $\text{HE}(\text{hpk}, \text{msg})$ a HE-ciphertext or a HE-encryption of a plaintext msg . Given two HE-ciphertexts $\text{HE}(\text{hpk}, \text{msg}_1)$ and $\text{HE}(\text{hpk}, \text{msg}_2)$, we call $C_1 + (\times) C_2$ *homomorphic addition* (multiplication) which is a HE-encryption of $\text{msg}_1 + (\times) \text{msg}_2$. When it is clear in the context, we just say addition or multiplication and omit the public key part in encryption algorithm, i.e., $\text{HE}(\text{msg})$. Similar notations are used for AHE with a public key apk .

In [34], Smart et al. suggested homomorphic encryption supporting Single Instruction Multiple Data (SIMD) operation, which enables to pack many plaintexts into each ciphertext. We say HE has $\#$ slot L if L plaintexts can be packed.

A vector is written in bold, i.e., \mathbf{a} . Given two vectors $\mathbf{a} = (a_1, \dots, a_k)$ and $\mathbf{b} = (b_1, \dots, b_k)$, we define

$$\mathbf{a} \times_c \mathbf{b} = (a_1 b_1, \dots, a_k b_k)$$

and $\text{sum}(\mathbf{a}) = \sum_{i=1}^k a_i$. We then have $\langle \mathbf{a}, \mathbf{b} \rangle = \text{sum}(\mathbf{a} \times_c \mathbf{b})$.

For a vector \mathbf{a} we denote $\text{HE}(\mathbf{a})$ by a HE-encryption of \mathbf{a} , where each entry of \mathbf{a} is encrypted to each slot of a ciphertext. Then $\text{HE}(\mathbf{a})$ consists of $k/\#$ slot HE-ciphertexts. Note that $\text{HE}(\mathbf{a}) + \text{HE}(\mathbf{b})$ is a HE-encryption of $\mathbf{a} + \mathbf{b}$ and $\text{HE}(\mathbf{a}) \times \text{HE}(\mathbf{b})$ is a HE-encryption of $\mathbf{a} \times_c \mathbf{b}$.

A fixed point representation over a real number f is denoted by \bar{f} . For a vector $\mathbf{a} = (a_1, \dots, a_k)$, we denote $\bar{\mathbf{a}}$ by component-wise fixed point representation. In this work, we use the binary fixed point representation and denote α by the bit length of scaling factor.

3.3 Matrix Factorization

We briefly introduce a matrix factorization method to be used in our system, which is also considered in [28]. Suppose n users rate a subset of m possible items. We use the following notations:

- $[n] = \{1, \dots, n\}$, $[m] = \{1, \dots, m\}$: the sets of users and items, respectively,
- $r_{ij} \in \mathbb{R}$: the user i 's rating for the item j ,
- $\mathcal{M} \subset [n] \times [m]$: a set of the user/item pairs for which a rating has been generated,
- M : the total number of ratings, i.e., $M = |\mathcal{M}|$,
- $\mathbf{u}_i, \mathbf{v}_j \in \mathbb{R}^d$: profiles for the user i and item j , respectively, where d is the dimension of profiles
- $U \in \mathbb{R}^{n \times d}$: the user-profile matrix (user profiles) whose i -th row is \mathbf{u}_i ,
- $V \in \mathbb{R}^{m \times d}$: the item-profile matrix (item profiles) whose j -th row is \mathbf{v}_j

Given the ratings $\{r_{ij} : (i, j) \in \mathcal{M}\}$, matrix factorization computes the user profiles U and item profiles V . The resulting profiles are used to predict the user i 's rating for item j , that is, $\langle \mathbf{u}_i, \mathbf{v}_j \rangle$, for $(i, j) \notin \mathcal{M}$. This can be done

by fitting $\langle \mathbf{u}_i, \mathbf{v}_j \rangle$ on the existing ratings r_{ij} for $(i, j) \in \mathcal{M}$, that is, solving the regularized least squares minimization:

$$\min_{U, V} \frac{1}{M} \sum_{(i, j) \in \mathcal{M}} (r_{ij} - \langle \mathbf{u}_i, \mathbf{v}_j \rangle)^2 + \lambda \sum_{i \in [n]} \|\mathbf{u}_i\|_2^2 + \mu \sum_{j \in [m]} \|\mathbf{v}_j\|_2^2$$

for some $\lambda, \mu > 0$.

Among methods to solve the above, we use gradient descent [23], a popular method in practice. It iteratively adopts the profiles U and V through the following adaptation rule:

$$\begin{aligned} \mathbf{u}_i(t) &= \mathbf{u}_i(t-1) - \gamma \nabla_{\mathbf{u}_i} F(U(t-1), V(t-1)), \\ \mathbf{v}_j(t) &= \mathbf{v}_j(t-1) - \gamma \nabla_{\mathbf{v}_j} F(U(t-1), V(t-1)), \end{aligned} \quad (1)$$

where $\gamma > 0$ a small gain factor and

$$\begin{aligned} \nabla_{\mathbf{u}_i} F(U, V) &= \nabla_{\mathbf{u}_i} \\ &= \sum_{j: (i, j) \in \mathcal{M}} \mathbf{v}_j (\langle \mathbf{u}_i, \mathbf{v}_j \rangle - r_{ij}) + \lambda \mathbf{u}_i, \\ \nabla_{\mathbf{v}_j} F(U, V) &= \nabla_{\mathbf{v}_j} \\ &= \sum_{i: (i, j) \in \mathcal{M}} \mathbf{u}_i (\langle \mathbf{u}_i, \mathbf{v}_j \rangle - r_{ij}) + \mu \mathbf{v}_j, \end{aligned} \quad (2)$$

where $U(0)$ and $V(0)$ consist of uniformly random norm 1 rows.

The number of iterations depends on stopping criteria. One of typical criteria is to set a (small) threshold for gradient norms, which is applied to the proposed protocol in Section 4. The protocol stops when the number of iteration exceeds some pre-defined number or when both

$$\sum_{i: (i, j) \in \mathcal{M}} \|\nabla_{\mathbf{u}_i} F(U, V)\|_2^2 \quad \text{and} \quad \sum_{j: (i, j) \in \mathcal{M}} \|\nabla_{\mathbf{v}_j} F(U, V)\|_2^2$$

become smaller than some threshold values, respectively.

Note that the values of d, γ, λ, μ , and stopping criteria can be a know-how for a recommendation service provider [38].

4. DESIGN COMPONENTS

4.1 Naive Approach and Challenges

Fully homomorphic encryption. We can consider a naive way to perform gradient descent using fully homomorphic encryption. The simplest way is to encrypt each component in vectors $\mathbf{u}_i, \mathbf{v}_j$ independently. In this case, we need to deal dM HE-ciphertexts to compute gradient descent, which is computational burden.

Instead, the SIMD operation proposed in [34] can be used to improve the inner product algorithm, that is, if we pack elements in two vectors of length d into two ciphertexts, respectively, the computation overhead is reduced by a factor of $(2 \log d/d)$ ¹. In our case, on the other hand, encrypting one vector in one ciphertext is quiet inefficient, since the dimension of $\mathbf{u}_i, \mathbf{v}_j$ is quite small compared to the number of slots underlying HE and the number of ciphertexts grows in M .

To make use of full message slots, one may consider packing several user and item profile vectors into a single ciphertext, respectively. Let $\mathbf{c}_u = \text{HE}(\mathbf{u}), \mathbf{c}_v = \text{HE}(\mathbf{v})$ for $\mathbf{u} = (\mathbf{u}_1 || \dots || \mathbf{u}_n)$ and $\mathbf{v} = (\mathbf{v}_1 || \dots || \mathbf{v}_m)$ where $||$ means vector concatenation. To compute $\langle \mathbf{u}_i, \mathbf{v}_j \rangle$ with $i < j$, we

¹The inner product on packed two vectors of length d can be carried out through $\log d$ shift in plaintext slots and $\log d$ multiplications instead of d .

operate $d \cdot (j - i)$ left-shift of plaintext slots on c_v , say c'_v , and then evaluate homomorphically multiplication of c_v and c'_v . Even though we have $c_u \times c'_v$ whose slots from $d \cdot (i - 1)$ to $d \cdot i$ are $\mathbf{u}_i \times_c \mathbf{v}_j$, the rest of slots is hard to be exploited in computing gradient descent. This means that we need to evaluate one homomorphic multiplication for only one inner product, which does not enjoy the benefit of SIMD at all.

Another concern in adapting HE is to perform fixed point arithmetic in encrypted form. Since there is no efficient homomorphic encryption supports rational number arithmetic, we need to transform data into integer form. In this case, it essentially needs additional post-process like integer division to maintain the same scaling factor. Under HE, this post-process seems heavy computations as in [14].

Stopping criteria. In general, the matrix factorization for recommendation runs several iterations while tuning some parameters such as γ, λ and μ until one gets meaningful user/item profiles. Furthermore, these tuning parameters might be know-how of the recommendation service provider. Therefore, it is important to construct recommendation system that can adjust tuning parameters during matrix factorization without leakage of the private parameters.

4.2 Data Structure for Gradient Descent

To overcome challenges in the previous section, we come up with an elaborate data structure to exploit slots fully so that the protocol requires only M/L HE-ciphertexts, instead of M , for #slot L . The idea behind our approach is to set a long vector whose segments are profile vectors preserving index information induced from the user-item matrix and perform joint computation with the CSP. If inner products and fixed point operations are required, the RecSys calls the CSP as encryption/decryption oracles and computes a part of operations with the CSP.

We define a total order on a set \mathcal{M} as follows: for (i_1, j_1) and $(i_2, j_2) \in \mathcal{M}$,

$$\begin{cases} (i_1, j_1) > (i_2, j_2) & \text{if } i_1 > i_2 \text{ or } j_1 > j_2 \text{ when } i_1 = i_2, \\ (i_1, j_1) = (i_2, j_2) & \text{if } i_1 = i_2 \text{ and } j_1 = j_2. \end{cases}$$

From now on we assume elements in \mathcal{M} are sorted in ascending order. Then for a set \mathcal{M} , we define 4 vectors as follows:

- $\mathbf{U} = \parallel_{(i,j) \in \mathcal{M}} \mathbf{u}_i$ and $\mathbf{V} = \parallel_{(i,j) \in \mathcal{M}} \mathbf{v}_j$, where \parallel means vector concatenation. Note that, depending on its order, every d -dimensional vector \mathbf{u}_i (\mathbf{v}_j) in \mathbf{U} (\mathbf{V}) corresponds to the unique pair (i, j) in \mathcal{M} .
- $\hat{\mathbf{U}} = \parallel_{(i,j) \in \mathcal{M}} \hat{\mathbf{u}}_i$, where $\hat{\mathbf{u}}_i = \mathbf{u}_i$ if i of the corresponding index (i, j) appears first and the d -dimensional zero vector otherwise. We define $\hat{\mathbf{V}}$ similarly, according to the index j from $(i, j) \in \mathcal{M}$.

Suppose $\mathcal{M}_I = \{i_1, \dots, i_{\ell_u}\}$ and $\mathcal{M}_J = \{j_1, \dots, j_{\ell_v}\}$ are the sets of distinct user index i 's and item index j 's of $(i, j) \in \mathcal{M}$ in ascending order, respectively. Consider a dM -dimensional vector $\mathbf{A} = \mathbf{a}_1 \parallel \dots \parallel \mathbf{a}_M$, where \mathbf{a}_k is a d -dimensional vector. Then we can match \mathbf{a}_k to $(i, j) \in \mathcal{M}$ according to order, which is a 1-1 correspondence. With this correspondence, we define additional 4 operations on \mathbf{A} as follows:

- $\mathbf{agg}_u, \mathbf{agg}_v$: We define aggregations of the user and

item profiles in \mathbf{A} according to \mathcal{M} in \mathbf{agg}_u and \mathbf{agg}_v :

$$\begin{aligned} \mathbf{agg}_u(\mathbf{A}, \mathcal{M}) &= \sum_{j:(i_1, j) \in \mathcal{M}} \mathbf{a}_j \parallel \dots \parallel \sum_{j:(i_{\ell_u}, j) \in \mathcal{M}} \mathbf{a}_j, \\ \mathbf{agg}_v(\mathbf{A}, \mathcal{M}) &= \sum_{i:(i, j_1) \in \mathcal{M}} \mathbf{a}_i \parallel \dots \parallel \sum_{i:(i, j_{\ell_v}) \in \mathcal{M}} \mathbf{a}_i. \end{aligned}$$

- $\mathbf{rec}_u, \mathbf{rec}_v$: We define reconstitutions of the user/item profiles by

$$\begin{aligned} \mathbf{rec}_u(\mathbf{agg}_u, \mathcal{M}) &= \parallel_{(i,j) \in \mathcal{M}} \mathbf{A}'_i, \\ \mathbf{rec}_v(\mathbf{agg}_v, \mathcal{M}) &= \parallel_{(i,j) \in \mathcal{M}} \mathbf{A}''_j \end{aligned}$$

where

$$\begin{aligned} \mathbf{agg}_u(\mathbf{A}, \mathcal{M}) &= \mathbf{A}'_1 \parallel \dots \parallel \mathbf{A}'_{\ell_u}, \\ \mathbf{agg}_v(\mathbf{A}, \mathcal{M}) &= \mathbf{A}''_1 \parallel \dots \parallel \mathbf{A}''_{\ell_v}. \end{aligned}$$

When it is clear in the context we omit \mathcal{M} of \mathbf{agg} and \mathbf{rec} .

4.3 Fixed Point Arithmetic on Encrypted Data

We use fixed point representation (FPR) of real number instead of floating point representation at the cost of small errors as in [28]. We use the integral version of FPR and denote the FPR with fixed α -bit precision by \bar{a} for a real number $a \in \mathbb{R}$. Here, we just say the FPR for integral version of that. The operations on FPR are as follows:

- Addition/Subtraction: $\overline{\bar{a} \pm \bar{b}} = \overline{\bar{a} \pm b}$
- Multiplication: $\overline{\bar{a} \cdot \bar{b}} = \lfloor \bar{a} \cdot \bar{b} / 2^\alpha \rfloor$

where $\lfloor \cdot \rfloor$ is a round down function. In our construction, we encrypt the FPR of user/item profiles and perform homomorphic additions/subtractions and multiplications on them. When encrypting this type of data, the main obstacle is multiplication because of integer division by 2^α and rounding, which can be considered as the α -bit right shift. In [15], the authors proposed homomorphic computation of right shift on integers using bit-extraction. On the other hand, α multiplicative depth of underlying FHE is consumed and it needs α^2 homomorphic multiplications, which results in huge computational overhead.

We design the RecSys to work jointly with the CSP to enhance the efficiency when performing multiplications on FPR values. We need one encryption and one decryption with one interaction in handling fixed point multiplication as follows: (i) the RecSys sends masked encryption of multiplication result to the CSP to protect the privacy of data. (ii) the CSP decrypts the masked ciphertext and computes right shift on the plaintext. (iii) the CSP encrypts resulting data and send it to the RecSys. (iv) the RecSys removes the mask and obtain the result. If the right shift is applied to masked integer, the overflow in the LSB of integer may occur. However it does not affect to the quality of matrix factorization significantly because of its small magnitude. The detail analysis and experimental results will be presented in Section 8.1.

5. OUR PROTOCOL

We present our privacy-preserving matrix factorization protocol. Our protocol consists of the setup, rating-upload, and matrix factorization phases. The recommendation phase after finishing matrix factorization is described in the full version of the paper.

5.1 Setup

In the setup phase, the CSP generates its HE-public/secret key pairs and AHE-public/secret key pairs. For any public key encryption PE, the users also generate their PE-public/secret key pairs.

The RecSys specifies the following public parameters: (i) the dimension of profiles d , (ii) the number of items m , the number of users n , (iii) the number of bits used to represent the integer and fractional parts of a real number in ratings as well as the fractional part of real numbers λ , μ , and γ in gradient descent computation.

We denote by α and β the number of bits used to represent the fractional part of real numbers in ratings and γ , respectively. We do not set such numbers for λ, μ assuming that they are smaller than α (see Step 6 of Figure 4). Note that λ, μ , and γ are privately selected by the RecSys and kept secret from users and the CSP.

5.2 Rating-Upload

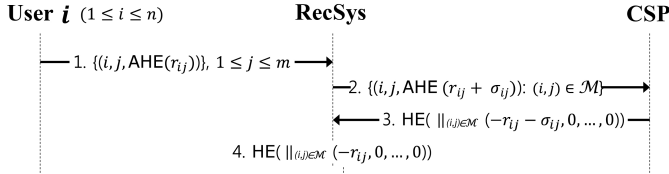


Figure 2: The Rating-Upload Phase

In the rating-upload phase (Figure 2), the users upload their ratings so that the RecSys performs matrix factorization to obtain the user and item profiles. Since users' ratings are secret information, the RecSys are given encrypted ratings via the CSP. During the phases no one learns r_{ij} except the user i . In the below protocol, all public key encryption schemes use only public/secret keys of the CSP. So we omit key argument (the 1st argument) in encryption. The detailed description is as follows:

- For $i \in [n]$ and $j \in [m]$, each user i encrypts ratings r_{ij} 's with AHE under the public key of the CSP. The users send $(i, j, \text{AHE}(r_{ij}))$'s to the RecSys.

- Receiving

$$\{(i, j, \text{AHE}(r_{ij})) : (i, j) \in \mathcal{M}\},$$

the RecSys generates random masks σ_{ij} 's. It then computes and sends

$$\{(i, j, \text{AHE}(r_{ij} + \sigma_{ij})) : (i, j) \in \mathcal{M}\}$$

to the CSP.

- The CSP decrypts AHE-ciphertexts to obtain

$$\{(i, j, r_{ij} + \sigma_{ij}) : (i, j) \in \mathcal{M}\}.$$

It then sets the set of d -dimensional vectors

$$\{(-r_{ij} - \sigma_{ij}, 0, \dots, 0) : (i, j) \in \mathcal{M}\}.$$

Finally it computes and sends

$$\text{HE}(\|_{(i,j) \in \mathcal{M}} (-r_{ij} - \sigma_{ij}, 0, \dots, 0))$$

to the RecSys.

- From the knowledge of σ_{ij} 's the RecSys can build the set of d -dimensional vectors

$$\{(\sigma_{ij}, 0, \dots, 0) : (i, j) \in \mathcal{M}\}.$$

It then adds

$$\text{HE}(\|_{(i,j) \in \mathcal{M}} (\sigma_{ij}, 0, \dots, 0))$$

to a given

$$\text{HE}(\|_{(i,j) \in \mathcal{M}} (-r_{ij} - \sigma_{ij}, 0, \dots, 0)).$$

The result is a HE-encryption of

$$\|_{(i,j) \in \mathcal{M}} (-r_{ij}, 0, \dots, 0).$$

Let $\mathbf{r} = \|_{(i,j) \in \mathcal{M}} (-r_{ij}, 0, \dots, 0)$. At last the RecSys selects uniformly random norm 1 vectors $\mathbf{u}_i(0)$ and $\mathbf{v}_j(0)$ for $i \in [n]$ and $j \in [m]$ and sets $\mathbf{U}(0)$, $\hat{\mathbf{U}}(0)$, $\mathbf{V}(0)$, and $\hat{\mathbf{V}}(0)$.

5.3 Matrix Factorization

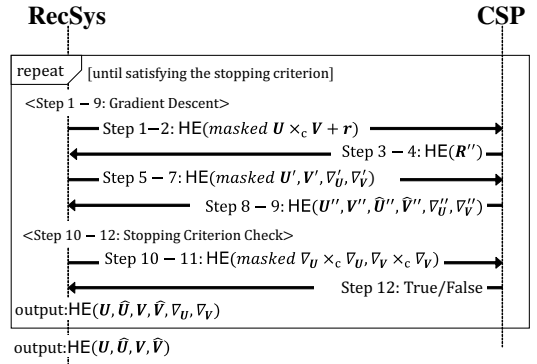


Figure 3: The Matrix Factorization Phase

Given \mathcal{M} , let

$$\mathbf{U}(t-1) = \|_{(i,j) \in \mathcal{M}} \mathbf{u}_i(t-1), \quad \mathbf{V}(t) = \|_{(i,j) \in \mathcal{M}} \mathbf{v}_j(t-1)$$

from $\mathbf{u}_i(t-1)$'s and $\mathbf{v}_j(t-1)$'s in (1). $\hat{\mathbf{U}}(t-1)$ and $\hat{\mathbf{V}}(t-1)$ are induced vectors from $\mathbf{U}(t-1)$ and $\mathbf{V}(t-1)$, respectively as defined in the previous section. For items which nobody rates, item profiles are directly computed by the RecSys, since no rating r_{ij} is involved in the gradient descent computation (1) and (2). For items which are rated, we consider the vector representation of ∇_{u_i} and ∇_{v_j} for $i \in \mathcal{M}_I$ and $j \in \mathcal{M}_J$ in (2) as

$$\nabla_{\mathbf{U}}(t-1) = \nabla_{u_{i_1}}(t-1) \parallel \dots \parallel \nabla_{u_{i_{\ell_u}}}(t-1) \quad (i_k \in \mathcal{M}_I),$$

$$\nabla_{\mathbf{V}}(t-1) = \nabla_{v_{j_1}}(t-1) \parallel \dots \parallel \nabla_{v_{j_{\ell_v}}}(t-1) \quad (j_k \in \mathcal{M}_J).$$

Now we are in a position to present how to compute HE-encryption of $\{\mathbf{U}, \mathbf{V}, \hat{\mathbf{U}}, \hat{\mathbf{V}}, \nabla_{\mathbf{U}}, \nabla_{\mathbf{V}}\}$ at t from HE-encryption at $t-1$ for a positive integer t . Note that the RecSys can initialize each vector since it does not require rating information at $t=0$.

In the protocol, random masking vectors are required when the RecSys sends HE-ciphertexts to the CSP. Given \mathcal{M} and the message space \mathcal{B} of HE, we define the distribution $\mathcal{D}(\mathcal{M}, \mathcal{B})$, which works as follows:

1. chooses M making vectors $\mathbf{ev}^{(i,j)} = (\mathbf{ev}_1^{(i,j)}, \dots, \mathbf{ev}_d^{(i,j)})$ for $(i, j) \in \mathcal{M}$, where $\mathbf{ev}_k^{(i,j)}$'s are random elements from \mathcal{B} .
2. outputs $\mathbf{ev} = \parallel_{(i,j) \in \mathcal{M}} \mathbf{ev}^{(i,j)}$.

The protocol consists of the secure 2-party computation between the RecSys and CSP, as shown in Figure 3. To keep intermediate values secret from the CSP, the RecSys mostly sends the HE-encryption of masked values. Then the CSP decrypts them. The CSP performs fixed point arithmetic and gets back them to the RecSys in HE-encrypted form. Given HE-ciphertexts the RecSys can remove masks since it creates them. It is important to note that this mask-removing approach causes small errors to the accuracy of the resulting profiles. However it is easy to see that errors are limited and negligible. For details refer to Section 7.

The detailed procedure for the matrix factorization phase is described in Figure 4. For an insightful example see the full version of the paper.

5.4 Enhancing Privacy of Rated Items

Besides users' ratings on items, other information can be used to infer users' partial identities. In [40], it has been shown that one can infer user's gender from which items the user rated. To address the issue, the use of differential privacy or injection of fake ratings has been suggested. However these methods distort results of data analysis, i.e., U and V of matrix factorization.

The protocol presented in the previous section permits all the above approaches, that is, in the rating-upload phase users can blur their rating information. Note that noising users' ratings is unnecessary in the protocol since all ratings are encrypted, hence, kept secret from the RecSys and CSP. Thus we extend the previous protocol by taking the injection of fake ratings in this section. To prevent the distortion of resulting profiles, the protocol removes the effects of fake ratings in the matrix factorization phase from the use of input indicators produced by the users in the encrypted form².

The protocol then is able to additionally hide the following information from the above approach: (i) which items the users rated exactly and (ii) the number of ratings given by the users. Note that this approach increases the number of ratings in encrypted form. However it increases the running time of the protocol only linearly with the number of fake items.

There are slight differences in both the rating-upload and matrix factorization phases. We give the description of the extended protocol focusing on the difference.

The rating-upload phase. In the rating-upload phase the users generate their fake and real ratings attached with indicators $\tau_{ij} \in \{0, 1\}$. 0 indicates the fake and 1 the real. Users send the AHE-encryption of the ratings under CSP's public key. After the CSP decrypts them, it figures out the HE-encryption of the masked rating vector as it does in the original protocol. In addition, it reconstructs the masked indicator vector τ according to \mathcal{M} so that every u_i in $\mathbf{U}(t)$ corresponds to $(1, \dots, 1)$ if it is real or $(0, \dots, 0)$ otherwise in the clear forms. The detailed procedures are as follows (all encryptions are made under CSP's public key):

²This approach has been mentioned in [28], however, was not included in the design of the protocol for efficiency issue.

1. User i sends $(i, j, \text{AHE}(r_{ij}), \text{AHE}(\tau_{ij}))$ to the RecSys.
2. Adding errors σ_{ij} and σ'_{ij} to AHE-encryptions, the RecSys sends $\{(i, j, \text{AHE}(r_{ij} + \sigma_{ij}), \text{AHE}(\tau_{ij} + \sigma'_{ij}))\}$ to the CSP.
3. The CSP sets two dM -dimensional vectors

$$\parallel_{(i,j) \in \mathcal{M}} (-r_{ij} - \sigma_{ij}, 0, \dots, 0),$$

$$\parallel_{(i,j) \in \mathcal{M}} (\tau_{ij} + \sigma'_{ij}, \dots, \tau_{ij} + \sigma'_{ij}),$$

where all vectors in concatenation are d -dimensional. The CSP encrypts them with HE and sends the ciphertexts to the RecSys.

4. The RecSys removes errors in the encrypted vectors as it does in the original protocol. The rest of the phase is identical.

The matrix factorization phase. Let τ be a HE-encryption of the error-free indicator vectors obtained in the previous phase, i.e, $\tau = \text{HE}(\parallel_{(i,j) \in \mathcal{M}} (\tau_{ij}, \dots, \tau_{ij}))$. The only difference happens when computing $\nabla'_{\mathbf{U}}(t)$ and $\nabla_{\mathbf{U}}(t)$, that is, the multiplication of \mathbf{R} vector to \mathbf{V} and \mathbf{U} vectors at Step 6 in Figure 4. That is, the RecSys computes

$$\text{HE}(\tau \times_c \mathbf{V}(t-1) \times_c \mathbf{R}(t-1) + \alpha \lambda \hat{\mathbf{U}}(t-1)),$$

$$\text{HE}(\tau \times_c \mathbf{U}(t-1) \times_c \mathbf{R}(t-1) + \alpha \mu \hat{\mathbf{V}}(t-1)).$$

Note that multiplying $\text{HE}(\tau)$ requires the appropriate set-up of multiplication depth for HE, but does not the enlarged message space.

6. ANALYSIS OF THE PROTOCOL

In this section we give an analysis of our protocol in Section 5.

6.1 Complexity

We analyze computation and communication complexities of the proposed protocol in Table 1. We only counted the number of cryptographic operations performed by each party. Given fixed d , the table shows that the proposed protocol requires $O(M/L)$ computation and communication complexities.³

RU and MF are abbreviations for rating-upload, matrix factorization, respectively. `enc` and `dec` denote encryption and decryption algorithm, respectively. `add`, `smul`, and `mul` point out addition, scalar multiplication, and multiplication over encrypted data, respectively. `|AHE|` and `|HE|` mean the sizes of AHE and HE-ciphertexts, respectively. We remark that masks themselves are encryption, so we do not count adding and removing masks by the RecSys as a cryptographic operation.

In MF, workloads for both the RecSys and CSP are comparable. In fact we observe that `HE.dec` is more expensive than any other homomorphic operations (see Section 8.3) because it requires an evaluation of polynomials to recover plaintext in each slot. However our evaluation in Section 8.2 shows that the workload for the CSP in our construction is quite small compared to that from garbled circuit approaches.

³ M includes the number of fake ratings in the case of the extended protocol in Section 5.4.

On inputs \mathbf{r} , HE-encryption of $\{\mathbf{U}(t-1), \hat{\mathbf{U}}(t-1), \mathbf{V}(t-1), \hat{\mathbf{V}}(t-1)\}$, the RecSys and CSP jointly and securely compute HE-encryption of $\{\mathbf{U}(t), \hat{\mathbf{U}}(t), \mathbf{V}(t), \hat{\mathbf{V}}(t)\}$. At the end of the below procedure, the RecSys confirms the stopping criteria and decides to continue the next iteration.

1. The RecSys determines parameters λ , μ , and γ within the ranges defined in the setup phase. It computes

$$\text{HE}(\mathbf{U}(t-1)) \times \text{HE}(\mathbf{V}(t-1)) - 2^\alpha \cdot \text{HE}(\mathbf{r}),$$

where a positive integer α is used to match fixed point representation of entries in \mathbf{r} to those in $\mathbf{U}(t-1) \times_c \mathbf{V}(t-1)$ (See Section 4.3 for details). Let $\mathbf{R}'(t-1) = \mathbf{U}(t-1) \times_c \mathbf{V}(t-1) - 2^\alpha \mathbf{r}$.

2. The RecSys samples $\epsilon(t-1) \stackrel{\$}{\leftarrow} \mathcal{D}(\mathcal{M}, \mathcal{B})$ and stores it. The RecSys then computes

$$\text{HE}(\mathbf{U}(t-1)) \times \text{HE}(\mathbf{V}(t-1)) - 2^\alpha \cdot \text{HE}(\mathbf{r}) + \text{HE}(\epsilon(t-1)),$$

and sends HE-ciphertexts to the CSP.

3. The CSP decrypts given HE-ciphertexts from the RecSys. The CSP then performs significant arithmetic over it, which yields

$$\|_{(i,j) \in \mathcal{M}} \{\overline{\mathbf{u}_i(t-1) \times_c \mathbf{v}_j(t-1) - (r_{ij}, 0, \dots, 0) + \epsilon^{(i,j)}(t-1)}\}.$$

4. The CSP computes $R_{ij} = \text{sum}(\overline{\mathbf{u}_i(t-1) \times_c \mathbf{v}_j(t-1) - (r_{ij}, 0, \dots, 0) + \epsilon^{(i,j)}(t-1)})$'s for $(i, j) \in \mathcal{M}$, and sets $\mathbf{R}''(t-1) = \|_{(i,j) \in \mathcal{M}} (R_{ij}, \dots, R_{ij})$, where each vector in concatenation is d -dimensional. At last the CSP computes $\text{HE}(\mathbf{R}''(t-1))$, and sends them to the RecSys.

5. The RecSys computes $\epsilon''^{(i,j)}(t-1) = \text{sum}(\overline{\epsilon^{(i,j)}(t-1)})$, sets the concatenation of dimension d vectors $\epsilon''(t-1) = \|_{(i,j) \in \mathcal{M}} (\epsilon''^{(i,j)}(t-1), \dots, \epsilon''^{(i,j)}(t-1))$. It then computes $\text{HE}(\mathbf{R}''(t-1)) - \text{HE}(\epsilon''(t-1))$, which yields a HE-encryption of

$$\mathbf{R}(t-1) = \|_{(i,j) \in \mathcal{M}} (\overline{\langle \mathbf{u}_i(t-1), \mathbf{v}_j(t-1) \rangle - r_{ij}}, \dots, \overline{\langle \mathbf{u}_i(t-1), \mathbf{v}_j(t-1) \rangle - r_{ij}}),$$

where each vector in concatenation is d -dimensional.

6. The RecSys computes HE-encryptions of

$$\nabla'_{\mathbf{U}}(t) = \mathbf{V}(t-1) \times_c \mathbf{R}(t-1) + 2^\alpha \lambda \hat{\mathbf{U}}(t-1), \quad \nabla'_{\mathbf{V}}(t) = \mathbf{U}(t-1) \times_c \mathbf{R}(t-1) + 2^\alpha \mu \hat{\mathbf{V}}(t-1).$$

It then computes HE-encryptions of

$$\mathbf{U}'(t) = 2^{\alpha+\beta} \cdot \hat{\mathbf{U}}(t-1) - \gamma 2^\beta \cdot \nabla'_{\mathbf{U}}(t), \quad \mathbf{V}'(t) = 2^{\alpha+\beta} \cdot \hat{\mathbf{V}}(t-1) - \gamma 2^\beta \cdot \nabla'_{\mathbf{V}}(t).$$

7. The RecSys samples $\delta_{\mathbf{U}}(t), \delta_{\mathbf{V}}(t), \theta_{\mathbf{U}}(t), \theta_{\mathbf{V}}(t) \stackrel{\$}{\leftarrow} \mathcal{D}(\mathcal{M}, \mathcal{B})$. It then computes and sends the CSP

$$\{\text{HE}(\mathbf{U}'(t) + \delta_{\mathbf{U}}(t)), \text{HE}(\mathbf{V}'(t) + \delta_{\mathbf{V}}(t))\}, \quad \{\text{HE}(\nabla'_{\mathbf{U}}(t) + \theta_{\mathbf{U}}(t)), \text{HE}(\nabla'_{\mathbf{V}}(t) + \theta_{\mathbf{V}}(t))\}.$$

8. After decrypting HE-ciphertexts the CSP performs fixed point arithmetic on them to get

$$\overline{\mathbf{U}'(t) + \delta_{\mathbf{U}}(t)}, \quad \overline{\mathbf{V}'(t) + \delta_{\mathbf{V}}(t)}, \quad \overline{\nabla'_{\mathbf{U}}(t) + \theta_{\mathbf{U}}(t)}, \quad \overline{\nabla'_{\mathbf{V}}(t) + \theta_{\mathbf{V}}(t)}.$$

It then computes $\mathbf{U}''(t) = \text{rec}(\text{agg}_u(\overline{\mathbf{U}'(t) + \delta_{\mathbf{U}}(t)}))$ and $\hat{\mathbf{U}}''(t)$ in turn. Similarly it computes $\mathbf{V}''(t) = \text{rec}(\text{agg}_v(\overline{\mathbf{V}'(t) + \delta_{\mathbf{V}}(t)}))$ and $\hat{\mathbf{V}}''(t)$.

9. The CSP computes

$$\nabla''_{\mathbf{U}}(t) = \text{agg}_u(\overline{\nabla'_{\mathbf{U}}(t) + \theta_{\mathbf{U}}(t)}), \quad \nabla''_{\mathbf{V}}(t) = \text{agg}_v(\overline{\nabla'_{\mathbf{V}}(t) + \theta_{\mathbf{V}}(t)}).$$

Note that when mask vectors in $\mathbf{U}''(t)$, $\hat{\mathbf{U}}''(t)$, $\mathbf{V}''(t)$, $\hat{\mathbf{V}}''(t)$, $\nabla''_{\mathbf{U}}(t)$, and $\nabla''_{\mathbf{V}}(t)$ are removed, those vectors equal to $\mathbf{U}(t)$, $\hat{\mathbf{U}}(t)$, $\mathbf{V}(t)$, $\hat{\mathbf{V}}(t)$, $\nabla_{\mathbf{U}}(t)$, and $\nabla_{\mathbf{V}}(t)$, respectively. At last it sends HE-encryption of $\mathbf{U}''(t)$, $\hat{\mathbf{U}}''(t)$, $\mathbf{V}''(t)$, $\hat{\mathbf{V}}''(t)$, $\nabla''_{\mathbf{U}}(t)$, and $\nabla''_{\mathbf{V}}(t)$.

10. From the knowledge of $\delta_{\mathbf{U}}(t)$, $\delta_{\mathbf{V}}(t)$, $\theta_{\mathbf{U}}(t)$, and $\theta_{\mathbf{V}}(t)$, the RecSys removes aggregated mask vectors in given HE-ciphertexts to get HE-encryptions of desired vectors $\mathbf{U}(t)$, $\hat{\mathbf{U}}(t)$, $\mathbf{V}(t)$, $\hat{\mathbf{V}}(t)$, $\nabla_{\mathbf{U}}(t)$, and $\nabla_{\mathbf{V}}(t)$.

11. The RecSys sets threshold values for the user and item profiles, say ω_u and ω_v . It generates mask vectors \mathbf{w}_u and \mathbf{w}_v , whose dimensions are equal to those of $\nabla_{\mathbf{U}}(t)$ and $\nabla_{\mathbf{V}}(t)$, respectively. It computes and sends the CSP

$$\text{HE}(\nabla_{\mathbf{U}}(t) \times_c \nabla_{\mathbf{U}}(t) + \mathbf{w}_u), \quad \text{HE}(\nabla_{\mathbf{V}}(t) \times_c \nabla_{\mathbf{V}}(t) + \mathbf{w}_v)$$

together with $s_u = \omega_u + \text{sum}(\mathbf{w}_u)$ and $s_v = \omega_v + \text{sum}(\mathbf{w}_v)$.

12. The CSP computes

$$s'_u = \text{sum}(\nabla_{\mathbf{U}}(t) \times_c \nabla_{\mathbf{U}}(t) + \mathbf{w}_u), \quad s'_v = \text{sum}(\nabla_{\mathbf{V}}(t) \times_c \nabla_{\mathbf{V}}(t) + \mathbf{w}_v).$$

Finally, it returns a boolean vector $(s_u - s'_u \geq 0?, s_v - s'_v \geq 0?)$ to the RecSys.

Figure 4: Matrix Factorization

		RecSys	CSP	Communication Cost
RU	AHE.dec	-	M	$M AHE + \frac{dM}{L} HE $
	HE.enc	-	dM/L	
MF	HE.enc	-	$7dM/L$	$\frac{14dM}{L} HE $
	HE.dec	-	$7dM/L$	
	HE.add	$5dM/L$	-	
	HE.smul	$7dM/L$	-	
	HE.mul	$5dM/L$	-	

Table 1: Computation and Communication Complexities

RU: rating-upload, MF: matrix factorization, enc and dec: encryption and decryption, resp., add, smul, mul: addition, scalar multiplication, and multiplication over encrypted data, resp., $|AHE|$ and $|HE|$: the sizes of AHE and HE-ciphertexts, resp., M : # of ratings, d : the dimension of profile vectors, L : # of slots

6.2 Security

As mentioned in Section 1, the purpose of privacy preserving matrix factorization is to protect the followings: user’s rating values, the list and number of items a user has rated, tuning parameters of the RecSys side, intermediate values and user/item profiles.

Our protocol protects all user’s rating values, intermediate values, and resulting user/item profiles using homomorphic encryptions (AHE and HE) and random mask all in the three phases. In a rating uploading phase, users’ ratings are AHE-encrypted under CSP’s public key, so the RecSys cannot reveal the data; and the CSP can obtain no information about users’ ratings even if it decrypts AHE-ciphertexts, because the rating values are masked by the RecSys. Note that mask works if there is no colluding attack of the RecSys and CSP, which is already assumed in Section 3. The similar argument works in the matrix factorization and recommendation phases. All intermediate values are kept secret, since the RecSys computes gradient descent on HE-encrypted data homomorphically and the values are masked when it needs help of the CSP. Furthermore, user/item profiles, output of matrix factorization, are obtained in HE-encrypted form and used in recommendation phase with masking for an individual user.

Regarding protection for the rated item list and number of ratings each user generates, our construction adopts an injecting method of fake ratings along with indicators. The construction without hiding the item profile in [28] has inferior security level than that of ours, since the adversary can infer whether a target user rate on a target item or not as described in [28, Section 4.1]. They extended their construction to protect the rated item list using random mask provided by the CSP at the cost of additional computation and communication overhead [28, Section 4.2]. The extended construction seems to have higher security level with respect to the rated item list than that of ours.

However there are some applications where only a specific item set, possibly small, is sensitive to leak whether rated or not. In such applications, we only obfuscate the small portion of the whole item set in our construction, hence achieve the same security with the extended version of [28]. For

example, we can apply the collaborative filtering to medical check-up to predict patients’ medical condition. The fact whether patients have checked typical medical condition such as blood pressure and sugar does not leak any private information on the patients. On the other hand, the fact that they have checked HIV infection or mental status is significantly related to privacy. In this case, we obfuscate only these sensitive checkup items to keep the user privacy.

When computing gradient descent in matrix factorization phase, the RecSys itself computes constant HE-multiplications with tuning parameters (γ, μ, λ) without help of the CSP. Therefore the optimized tuning parameters for quality recommendation can be kept secret from the CSP.

7. IMPLEMENTATION

In this section, we describe an implementation of our construction. In our protocol, we use two encryption scheme, an additively homomorphic encryption (AHE) and a fully (or somewhat) homomorphic encryption (FHE). For AHE in our protocol, a *partially* additive homomorphic property, which allows a constant addition to encrypted data without knowing any secret key, is enough, since the additive property is only used for large random masking. Therefore, we use hash-ElGamal encryption scheme [7]. We use a RLWE-based homomorphic encryption for FHE, which provides several optimizations such as efficient key switching, message packing and message slot movement [5, 15, 16, 34]. We implement our protocol using HELib [20], a homomorphic encryption library based on RLWE. The library is written in C++ and uses the Gnu MP library [19] and the NTL mathematical library [33].

The main bottleneck in deploying FHE is huge computation and communication overhead due to its large ciphertext. There are several factors that determines the size of ciphertext: multiplicative depth, security level and message space (\mathbb{Z}_p) . If the bit length of message space becomes k times larger, then the size of ciphertext is increased k^2 times.⁴ Since the computational overhead heavily depends on the size of ciphertext, the reduction of the size is the most important to enhance the efficiency. We consider chinese remainder theorem (CRT) to represent a large message. When applying CRT to message space, we obtain only linear overhead increase in ciphertext size. For a large message msg , we compute $\text{HE.enc}(\text{msg}_1), \dots, \text{HE.enc}(\text{msg}_k)$ where $\text{msg}_i = \text{msg} \bmod p_i$ for relatively primes p_1, \dots, p_k of the same size. Since the Chinese remainder algorithm is a ring isomorphism, this optimization does not affect the correctness of our construction.

8. EVALUATION

In this section we discuss accuracy of computation results including user and item profiles in the proposed protocol and provide experiment results.

8.1 Error Estimation

Compared to the work in [28], our protocol introduces more error in profiles depending on d , the dimension of profiles, due to fixed point arithmetic between the RecSys and

⁴This can be easily verified from the parameter selection of [16]

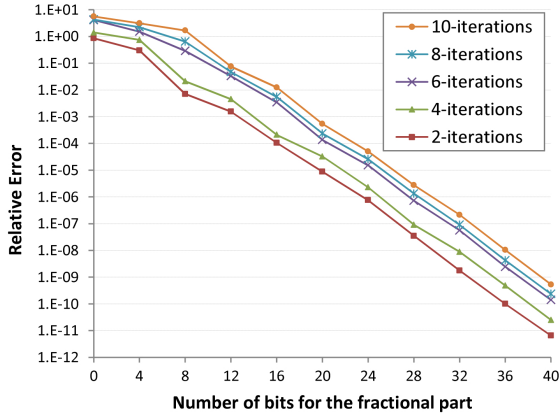


Figure 5: Relative Errors due to the Fixed Point Representation

CSP described in 4.3. The Recsys computes

$$\sum_{k=1}^d \overline{(f_{ijk} + \epsilon_k - \bar{\epsilon}_k)}$$

not $\sum_{k=1}^d \overline{f_{ijk}}$ for $\mathbf{f}_{ij} = (f_{ijk})_{k=1,\dots,d}$ and a mask vector $(\epsilon_k)_{k=1,\dots,d}$ in Step 3–4 of matrix factorization phase. This computation increases the error size in profiles. For $a, \epsilon \in \mathbb{R}$, let us consider the difference of \bar{a} and $\overline{(a + \epsilon - \bar{\epsilon})}$. We have

$$\left| \bar{a} - \overline{(a + \epsilon - \bar{\epsilon})} \right| \leq 1,$$

which means the overflow is bounded by $2^{-\alpha}$. Therefore the error in our protocol becomes at most $d \cdot 2^{-\alpha}$ by triangle inequality and this does not affect significantly on the accuracy of matrix factorization.

We evaluated the relative error of our matrix factorization to the original matrix factorization that operates in the clear with the float point representation. We use the relative error, which is used in [28]. Given user profiles U and item profiles V , the squared error is

$$E(U, V) = \sum_{(i,j) \in \mathcal{M}} (r_{ij} - \langle \mathbf{u}_i, \mathbf{v}_j \rangle)^2.$$

The relative error is then defined as

$$\frac{|E(U^*, V^*) - E(U, V)|}{|E(U, V)|},$$

where U^* and V^* are computed using our protocol with fixed-point representation and U and V are computed using gradient descent executed in the clear over floating point arithmetic. Figure 5 shows the relative error of our matrix factorization for the 100K MovieLens dataset. It shows that the relative errors occurred during 10 iterations is small (less than 10^{-4}) for more than 20-bit fractional part.

8.2 Experiment

Parameters and encoding of message. Let us consider message space of the i -th FHE, $R_{p_i}^{r_i} = \mathbb{Z}_{p_i}^{r_i}[x]/\Phi_{k_i}(x)$, where $\Phi_{k_i}(x)$ is the k_i -th cyclotomic polynomial. Since the bit length of message space for the dataset is 36, we use

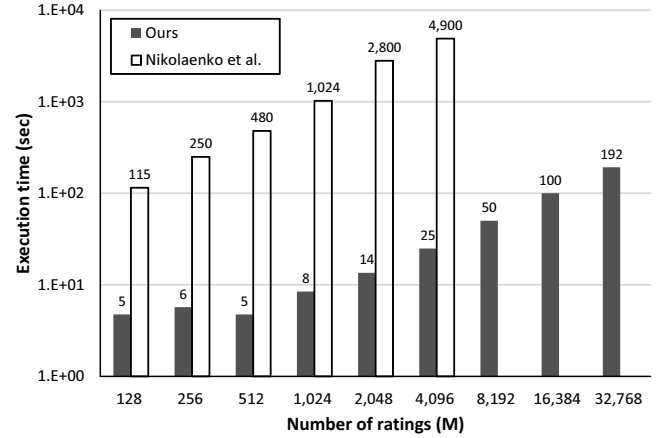


Figure 6: Comparison of Execution Time to [28]

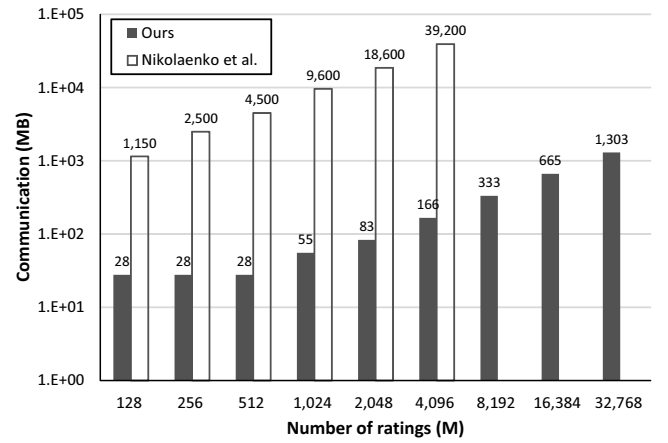


Figure 7: Comparison of Communication Cost to [28]

three FHE whose message space moduli are relatively coprime p_i 's of size 16-bit to represent 48-bit message space. In this case, we fix $r_i = 1$ for all i and find p_i and k_i to have large enough message slots in R_{p_i} . We choose p_i and k_i such that the number of message slots of underlying FHE is at least 6000 for all i . To evaluate our protocol we need FHE that allows only 3 multiplicative depth. Putting them together, the FHE-ciphertext size is 700 KB. For AHE, we use hash-ElGamal from OpenSSL [1] with slight customization⁵. Every encryption scheme supports 80-bit security.

Experiment results. We implement the proposed protocol and test the performance. We use MovieLens 100K dataset for experiment, which is popular real dataset in the research area. It contains 100K ratings on 1682 movies made by 943 users. For comparison with [28], we evaluate the time and communication performance for the above dataset restricted to 40 most popular movies. The restricted dataset contains 14683 ratings generated by 940 users. The rating-upload phase finishes in 4.3 seconds. For the matrix factorization, we measure the run time for one iteration to compare the result of the previous protocol [28]. Our matrix factorization

⁵The description can be found in the full version of the paper.

protocol takes 1.5 minutes with 580 MB for communication. Note that our matrix factorization protocol includes the procedure to verify the stopping criterion.

The test is performed on a single machine with 3.4GHz 6-cores 64GB RAM, where the operation system is Ubuntu and the program language is C++. Our platform equips about 2 times faster CPU than that of [28]. After normalizing the factor, we can estimate our matrix factorization protocol runs at least 50 times faster and consumes much less communication cost⁶, compared to [28]. Figure 6 and 7 compare the execution time and the communication cost of our matrix factorization to those of [28] for the same sets, respectively.

As describe in Section 5.4, each user can add fake ratings to enhance his own privacy. Our extended protocol takes 1.8 minutes with 623 MB for 10% fake ratings and 2.3 minutes with 910 MB for 50%, which is almost linear growth to the number of ratings.

8.3 Further Analysis and Comparison

We consider further potential improvement of the proposed protocol toward parallelism. We also give comparison to recent results on parallelized privacy-preserving matrix factorization [27] and optimization of circuits for the garbled circuit protocol [35].

Parallelism. In the matrix factorization phase of our protocol, all operations over ciphertexts such as HE.enc, HE.dec, HE.add, HE.smul, and HE.mul can be parallelized since ciphertexts are treated independently. For example, consider Step 1 of the matrix factorization in Figure 4 to compute

$$\text{HE}(\mathbf{U}(t-1)) \times \text{HE}(\mathbf{V}(t-1)) - 2^\alpha \cdot \text{HE}(\mathbf{r}).$$

Suppose there are N' machines acting as the RecSys and each $\text{HE}(\mathbf{U}(t-1))$, $\text{HE}(\mathbf{V}(t-1))$, and $\text{HE}(\mathbf{r})$ consists of N HE-ciphertexts, say $\text{HE}(\mathbf{U}(t-1))_i$, $\text{HE}(\mathbf{V}(t-1))_i$, and $\text{HE}(\mathbf{r})_i$ for $i = 1, \dots, N$. Parallelism then can be achieved by sending

$$\{\text{HE}(\mathbf{U}(t-1))_i, \text{HE}(\mathbf{V}(t-1))_i, \text{HE}(\mathbf{r})_i\}_{i=j+1, \dots, j+N/N'}$$

to the j -th machine and gathering the result

$$\{\text{HE}(\mathbf{U}(t-1))_i \times \text{HE}(\mathbf{V}(t-1))_i - 2^\alpha \cdot \text{HE}(\mathbf{r})_i\}_{i=j+1, \dots, j+N/N'}$$

computed by the j -th machine for $j = 1, \dots, N'$. The CSP side can be parallelized in a similar way.

On the contrary to the case of ciphertexts, operations over plaintexts should be performed by a single machine since they include reconstitution of data according to the order defined on \mathcal{M} . However run time for operations over plaintexts is negligible compared to that over ciphertexts; From the experiment for 14K dataset, we observe that the RecSys and CSP totally spend 15.5 and 73.8 seconds for matrix factorization (1.5 minutes overall), respectively. Run times over plaintexts are 1.57 seconds for the RecSys and 0.16 for the CSP. Run time over plaintexts has only 1.9% of total run time, which implies that time complexity of our matrix factorization is almost reduced by a factor of the number of machines involved in parallel operations.

Further comparison from the estimate. In 2015 Songhori et al. presented techniques, called TinyGarble, for optimizing circuits for protocols with garbled circuits [35]. It minimizes the number of non-XOR gates in a circuit. Their

⁶When $M = 4096$, [28] already consumes 40 GB.

proof-of-concept implementation showed that the technique decreases the number of non-XOR gates by 80%, which results in 5 times speed-up in generating and evaluation of a garbled circuit. The run time of Nikolaenko et al.'s protocol together with TinyGarble is estimated to about 34.8 minutes. It is still 10 times slower than our protocol (3 minutes after normalization due to CPU clock).

At the same time Nayak et al. proposed a parallel secure computation framework called GraphSC. It focuses on secure 2-party computation protocol based on garbled circuit [27]. They evaluated matrix factorization over 1M ratings on GraphSC, which takes 13 hours using a cluster, which consists of 7 machines of 128 processors (1.9–2.6 GHz CPUs) with 2048 cores overall. From the experiment result, estimated run time for 14K ratings using GraphSC is about 10 minutes. Since our experiment takes 1.5 minutes (3 minutes after normalization from CPU clock), our protocol even on a single machine with 6 cores is still about 3 times faster.

GraphSC together with TinyGarble may reduce the run time of matrix factorization further, that is, 13 hours for 1M rating to 2.6 hours. We estimate our protocol under the same environment of GraphSC, that is, 7 machines with 2048 cores. Our experiment on 14K ratings over a single machine with 6 cores (3.4 GHz) shows that it takes about 175.3 seconds for operations over ciphertexts and 3.4 seconds over plaintexts with CPU clock normalization. Thus, for 1M dataset, parallelization of our protocol is expected to take 36.7 seconds for operation over ciphertexts and 244.2 seconds over plaintexts since operations over plaintexts is not parallelizable. Thus, by adapting parallelization, the run time of our protocol is estimated to 4.7 minutes, which is about 33 times faster than the expected run time of GraphSC together with TinyGarble for 1M dataset.

9. CONCLUSION

A new efficient and practical privacy preserving matrix factorization protocol is proposed in this paper. The protocol protects users' rating values, user/item profiles and service's tuning parameters. Using various obfuscation techniques, we reasonably protect the rated item list and number users rated. We improved the efficiency by optimizing the use of fully homomorphic encryption, so the performance of ours is at least 50 times faster than that of the previous work [28]. Furthermore, ours satisfies both accuracy and completeness requirements.

10. REFERENCES

- [1] *OpenSSL: The Open Source toolkit for SSL/TLS*.
- [2] G. Adomavicius and A. Tuzhilin. Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions. *IEEE Trans. on Knowl. and Data Eng.*, 17(6):734–749, 2005.
- [3] J. W. Bos, K. Lauter, and M. Naehrig. Private predictive analysis on encrypted medical data. *Cryptology ePrint Archive*, Report 2014/336, 2014.
- [4] R. Bost, R. A. Popa, S. Tu, and S. Goldwasser. Machine learning classification over encrypted data. In *NDSS 2015*, 2015.
- [5] Z. Brakerski, C. Gentry, and V. Vaikuntanathan. (leveled) fully homomorphic encryption without bootstrapping. In *ITCS 2012*, pages 309–325, 2012.

- [6] J. Canny. Collaborative filtering with privacy. In *IEEE Security Privacy 2002*, pages 45–57, 2002.
- [7] B. Chevallier-Mames, P. Paillier, and D. Pointcheval. Encoding-free elgamal encryption without random oracles. In *PKC 2006*, pages 91–104, 2006.
- [8] W.-S. Chin, Y. Zhuang, Y.-C. Juan, and C.-J. Lin. A learning-rate schedule for stochastic gradient methods to matrix factorization. In *PAKDD 2015*, pages 442–455, 2015.
- [9] R. Chow, M. A. Pathak, and C. Wang. A practical system for privacy-preserving collaborative filtering. In *ICDM 2012*, pages 547–554, 2012.
- [10] J. Coron, T. Lepoint, and M. Tibouchi. Scale-invariant fully homomorphic encryption over the integers. In *PKC 2014*, pages 311–328, 2014.
- [11] J. Coron, D. Naccache, and M. Tibouchi. Public key compression and modulus switching for fully homomorphic encryption over the integers. In *EUROCRYPT 2012*, pages 446–464, 2012.
- [12] C. Dwork, F. McSherry, K. Nissim, and A. Smith. Calibrating noise to sensitivity in private data analysis. In *TCC 2006*, pages 265–284, 2006.
- [13] Z. Erkin, T. Veugen, T. Toft, and R. Legendijk. Generating private recommendations efficiently using homomorphic encryption and data packing. *IEEE Trans. Inf. Forensics Security*, 7(3):1053–1066, 2012.
- [14] C. Gentry, S. Halevi, and N. P. Smart. Better bootstrapping in fully homomorphic encryption. In *PKC 2012*, pages 1–16, 2012.
- [15] C. Gentry, S. Halevi, and N. P. Smart. Fully homomorphic encryption with polylog overhead. In *EUROCRYPT 2012*, pages 465–482, 2012.
- [16] C. Gentry, S. Halevi, and N. P. Smart. Homomorphic evaluation of the AES circuit. In *CRYPTO 2012*, pages 850–867, 2012.
- [17] O. Goldreich. *Foundations of Cryptography: Volume 2, Basic Applications*. Cambridge University Press, New York, NY, USA, 2004.
- [18] O. Goldreich and R. Ostrovsky. Software protection and simulation on oblivious rams. *J. ACM*, 43(3):431–473, 1996.
- [19] T. Granlund and the GMP development team. *GMP – The GNU Multiple Precision Arithmetic Library*.
- [20] S. Halevi and V. Shoup. Algorithms in helib. In *CRYPTO 2014*, pages 554–571, 2014.
- [21] J. Hua, C. Xia, and S. Zhong. Differentially private matrix factorization. In *IJCAI 2015*, pages 1763–1770, 2015.
- [22] Y. Huang, D. Evans, J. Katz, and L. Malka. Faster secure two-party computation using garbled circuits. In *USENIX Security 2011*, 2011.
- [23] Y. Koren, R. Bell, and C. Volinsky. Matrix factorization techniques for recommender systems. *IEEE Computer*, 42(8):30–37, 2009.
- [24] F. McSherry and I. Mironov. Differentially private recommender systems: Building privacy into the netflix prize contenders. In *KDD 2009*, pages 627–636, 2009.
- [25] N. Mohammed, R. Chen, B. Fung, and P. S. Yu. Differentially private data release for data mining. In *KDD 2011*, pages 493–501, 2011.
- [26] A. Narayanan and V. Shmatikov. Robust de-anonymization of large sparse datasets. In *IEEE Security Privacy 2008*, pages 111–125, 2008.
- [27] K. Nayak, X. S. Wang, S. Ioannidis, U. Weinsberg, N. Taft, and E. Shi. GraphSC: Parallel secure computation made easy. In *IEEE Security Privacy 2015*, pages 377–394, 2015.
- [28] V. Nikolaenko, S. Ioannidis, U. Weinsberg, M. Joye, N. Taft, and D. Boneh. Privacy-preserving matrix factorization. In *ACM CCS 2013*, pages 81–8121–812.
- [29] V. Nikolaenko, U. Weinsberg, S. Ioannidis, M. Joye, D. Boneh, and N. Taft. Privacy-preserving ridge regression on hundreds of millions of records. In *IEEE Security Privacy 2013*, pages 334–348, 2013.
- [30] P. Paillier. Public-key cryptosystems based on composite degree residuosity classes. In *EUROCRYPT 1999*, pages 223–238, 1999.
- [31] N. Ramakrishnan, B. Keller, B. Mirza, A. Y. Grama, and G. Karypis. Privacy risks in recommender systems. *IEEE Internet Comput.*, 5(6):54–63, 2001.
- [32] Y. Shen and H. Jin. Privacy-preserving personalized recommendation: An instance-based approach via differential privacy. In *ICDM 2014*, pages 540–549, 2014.
- [33] V. Shoup. *NTL: A Library for doing Number Theory*.
- [34] N. P. Smart and F. Vercauteren. Fully homomorphic simd operations. *Des. Codes Cryptography*, 71(1):57–81, 2014.
- [35] E. M. Songhori, S. U. Hussain, A. Sadeghi, T. Schneider, and F. Koushanfar. Tinygarble: Highly compressed and scalable sequential garbled circuits. In *IEEE Security Privacy 2015*, pages 411–428, 2015.
- [36] X. Su and T. M. Khoshgoftaar. A survey of collaborative filtering techniques. *Adv. in Artif. Intell.*, 2009:1–19, 2009.
- [37] L. Sweeney. K-anonymity: A model for protecting privacy. *Int. J. Uncertain. Fuzziness Knowl.-Based Syst.*, 10(5):557–570, 2002.
- [38] G. Takacs, I. Pillaszy, B. Nemeth, and D. Tikk. Investigation of various matrix factorization methods for large recommender systems. In *ICDMW 2008.*, pages 553–562, 2008.
- [39] T. Veugen, R. de Haan, R. Cramer, and F. Muller. A framework for secure computations with two non-colluding servers and multiple clients, applied to recommendations. *IEEE Trans. Inf. Forensics Security*, 10(3):445–457, 2015.
- [40] U. Weinsberg, S. Bhagat, S. Ioannidis, and N. Taft. Blurme: inferring and obfuscating user gender based on ratings. In *RecSys 2012*, pages 195–202, 2012.
- [41] Y. Xin and T. Jaakkola. Controlling privacy in recommender systems. In *NIPS 2014*, pages 2618–2626, 2014.
- [42] A. C. Yao. Protocols for secure computations. In *FOCS 1982*, pages 160–164, 1982.
- [43] H. Yun, H.-F. Yu, C.-J. Hsieh, S. V. N. Vishwanathan, and I. Dhillon. Nomad: Non-locking, stochastic multi-machine algorithm for asynchronous and decentralized matrix completion. *VLDB 2014*, 7(11):975–986, 2014.