

Secure Federated Matrix Factorization

Di Chai , Hong Kong University of Science and Technology

Leye Wang , Peking University

Kai Chen  and Qiang Yang, Hong Kong University of Science and Technology

To protect user privacy and meet law regulations, federated (machine) learning is obtaining vast interests in recent years. The key principle of federated learning is training a machine learning model without needing to know each user's personal raw private data. In this article, we propose a secure matrix factorization framework under the federated learning setting, called FedMF. First, we design a user-level distributed matrix factorization framework where the model can be learned when each user only uploads the gradient information (instead of the raw preference data) to the server. While gradient information seems secure, we prove that it could still leak users' raw data. To this end, we enhance the distributed matrix factorization framework with homomorphic encryption. We implement the prototype of FedMF and test it with a real movie rating dataset. Results verify the feasibility of FedMF. We also discuss the challenges for applying FedMF in practice for future research.

With the prevalence of government regulations and laws on privacy protection in the big data era (e.g., General Data Protection Regulation ([https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX:32016R0679R\(02\)](https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX:32016R0679R(02)))), privacy-preserving machine learning has obtained rapidly growing interests in both academia and industry. Among various techniques to achieve privacy-preserving machine learning, federated (machine) learning (FL) recently receives high attention. The original idea of FL was proposed by Google, (<https://ai.googleblog.com/2017/04/federated-learning-collaborative.html>) which targets at learning a centered model based on the personal information distributed at each user's mobile phone. While Google's solution targets individual consumers, WeBank formulated a to-business version of the solution known as Vertical federated learning, aimed at enterprise-level collaboration under privacy-preserving and security constraints.⁹ During the federated model training, no user's raw data is transferred to the central server or the other companies (i.e., other parties), thus ensuring privacy protection. Also, the learned privacy-preserving model can be

proved to hold almost similar predictive power compared to the traditional model learned on users' raw data. This highlights the practicality of FL as little predictive accuracy is sacrificed, especially compared to other accuracy-lossy privacy-preserving mechanisms such as differential privacy.

Starting from the original Google paper, many researchers have been devoted to this promising and critical area. Recently, a nice survey paper on FL has been published.⁸ While many research works have been done on FL, a popular machine learning technique, *matrix factorization* (MF),⁵ is still under-investigated in FL. Since MF is one of the prominent techniques widely employed in various applications such as item recommendation⁵ and environment monitoring,⁷ we highly believe that studying MF under FL is urgently required. This work is one of the pioneering research efforts toward this direction.

Taking recommendation systems as an example, two types of users' private information are leaked in traditional MF⁵: 1) *users' raw preference data*, and 2) *users' learned latent feature vectors*. As revealed by previous studies, either raw preference data or latent features can leak users' sensitive attributes, e.g., age, relationship status, political views, and sex orientations. This highlights the importance of protecting users' private information during MF. Prior studies have studied privacy-preserving MF in two main types:

- 1) *Obfuscation-based methods* obfuscate users' preference raw data before releasing it to the central server so as to ensure a certain level of privacy protection (e.g., differential privacy).³ The pitfall is that obfuscation inevitably leads to the loss of predictive power of the learned latent feature vectors. Hence, these methods usually need to make a tradeoff between privacy protection and model performance.
- 2) *Encryption-based methods* use advanced encryption schemes such as homomorphic encryption (HE) for implementing privacy-preserving MF.⁴ While they usually do not need to sacrifice predictive power for privacy protection, they commonly require a third-party crypto-service provider. This makes the system implementation complicated as such a provider is not easy to find in practice. (Some studies suggest that governments can perform this role, but this is still not the case in reality now.) Moreover, if the crypto-service provider colludes with the recommendation server, then no user privacy protection can be preserved.⁴

This research proposes a novel FL-based MF framework, called *FedMF* (federated matrix factorization). *FedMF* employs distributed machine learning and HE schemes. In brief, *FedMF* lets each user compute the gradients of his/her own rating information locally and then upload the gradients (instead of the raw data) to the server for training. To further enhance security, each user can encrypt the gradients with HE. With *FedMF*, two shortcomings of the traditional obfuscation- or encryption-based methods can be addressed: 1) no predictive accuracy is lost as we do not obfuscate data; 2) no third-party crypto-service provider is required as each user's device can handle the secure gradient computing task.

In summary, we have the following contributions.

- 1) To the best of our knowledge, we design the first FL-based secure MF framework, called *FedMF*, which can overcome the shortcomings of traditional obfuscation- or encryption-based mechanisms as aforementioned.
- 2) To implement *FedMF*, first, we design a user-level distributed MF framework where the model can be learned when each user only uploads the gradient information (instead of the raw preference data) to the server. Though gradient information seems secure, we prove that it could still leak users' raw data to some extent. (Similar work that proves gradients leak information also attracts many other researchers' interests, such

Algorithm 1. User-Level Distributed MF

Init: Server initializes item profile matrix V

Init: User initializes user profile matrix U

Output: Converged U and V

Server keeps latest item-profile for all users' download

User local update:

Download V from server, perform local updates:

$$u_i^t = u_i^{t-1} - \gamma \nabla_{u_i} F(U^{t-1}, V^{t-1})$$

$$\text{Gradient}_i = \gamma \nabla_{v_i} F(U^{t-1}, V^{t-1})$$

Send Gradient_i to server

Server update:

Receive Gradient_i from user- i

$$\text{Perform update : } v_i^t = v_i^{t-1} - \text{Gradient}_i$$

as Aono *et al.*,² proves gradients leak information in image deep learning domain. Here, we prove the privacy leakage from gradients in MF process.) Then, we enhance the distributed MF framework with HE to increase security.

- 3) We implement a prototype of *FedMF* (the code will be published in <https://github.com/Di-Chai/FedMF>). We test the prototype on a real movie rating dataset. Results verify the feasibility of *FedMF*.

It is worth noting that, similar to *FedMF*, Ammad-Ud-Din *et al.*¹ tried to develop a federated collaborative filtering system. However, Ammad-Ud-Din *et al.*¹ directly let users upload their gradient information to the server. As we will prove in this article, gradients can still reveal users' original preference data. Hence, a more secure system like *FedMF* is required to rigorously protect users' privacy.

PRELIMINARIES

In this section, we briefly introduce two techniques closely related to our research, *horizontal FL* and *additively HE*.

Horizontal FL

FL is a method that enables a group of data owners to train a machine learning model on their joint data and nobody can learn the data of the participants. *FL* can be categorized based on the distribution characteristics of the data.⁸ One category of *FL* is *horizontal FL*. *Horizontal FL* is introduced in the scenarios when the data from different contributors share the same feature space but vary in samples. In our secure MF recommendation system, the rating information distributed on each user's device has exactly the same feature space,

while different users are exactly different samples. So our MF in FL can be categorized as horizontal FL. Following the typical assumption of horizontal FL,⁸ we assume *all the users are honest and the system is designed to be secure against an honest-but-curious server*. (An honest-but-curious server means that the server is a legitimate participant in the system who will not deviate from the defined protocol but will attempt to learn all possible information from legitimately received messages.)

Additively HE

HE is often used in FL to protect user's privacy by providing encrypted parameter exchange. HE is a special kind of encryption scheme that allows any third party to operate on the encrypted data without decrypting it in advance. An encryption scheme is called homomorphic over an operation "*" if the following equation holds:

$$E(m_1) * E(m_2) = E(m_1^* m_2), \quad \forall m_1, m_2 \in M \quad (1)$$

where E is an encryption algorithm and M is the set of all possible messages.

Additively homomorphic encryption is homomorphic over addition. Typically, it consists of the following functions:

- $KeyGen \rightarrow (pk, sk)$: key generation process, where pk is the public key and sk is the secret key.
- $Enc(m, pk) \rightarrow c$: encryption process, where m is the message to be encrypted and c is the ciphertext.
- $Dec(c, sk) \rightarrow m$: decryption process
- $Add(c_1, c_2) \rightarrow c_a < (i.e. Enc(c_1 + c_2))$: add operation on ciphertext, c_a is the ciphertext of plaintexts' addition.
- $DecAdd(c_a, sk) \rightarrow m_a$: decrypt c_a , getting the addition of plaintexts.

USER-LEVEL DISTRIBUTED MF

We first introduce the MF optimization method used in our article, stochastic gradient descent.⁵ Based on it, we design a user-level distributed MF framework.

Stochastic Gradient Descent

Suppose we have n users, m items, and each user rated a subset of m items. For $[n] := \{1, 2, \dots, n\}$ as the set of users and $[m] := \{1, 2, \dots, m\}$ as the set of items, we denote $\mathcal{M} \in [n] \times [m]$ as user-item rating pairs which a rating has been generated, $M = |\mathcal{M}|$ as the total number of ratings and $r_{i,j}$ represents the rating generated by user i for item j .

Given the rating information $r_{ij} : (i, j) \in \mathcal{M}$, the recommendation systems are expected to predict the rating values of all the items for all the users. MF formulates this problem as fitting a bilinear model on the existing ratings. In particular, user profile matrix $U \in \mathbb{R}^{n \times d}$ and item profile matrix $V \in \mathbb{R}^{m \times d}$ are computed, where d is the dimension of the profile vector, the resulting profile matrices are used to predict user i 's rating on item j , which is $\langle u_i, v_j \rangle$. The computing process of U and V can be done by solving the following regularized least squares minimization:

$$\min_{U, V} \frac{1}{M} (r_{i,j} - \langle u_i, v_j \rangle)^2 + \lambda \|U\|_2^2 + \mu \|V\|_2^2 \quad (2)$$

where λ and μ are small positive values to rescale the penalizer. Stochastic gradient descent iteratively updates U and V with the following equations⁵:

$$u_i^t = u_i^{t-1} - \gamma \nabla_{u_i} F(U^{t-1}, V^{t-1}) \quad (3)$$

$$v_j^t = v_j^{t-1} - \gamma \nabla_{v_j} F(U^{t-1}, V^{t-1}) \quad (4)$$

where

$$\nabla_{u_i} F(U, V) = -2 \sum_{j:(i,j)} v_j (r_{ij} - \langle u_i, v_j \rangle) + 2\lambda u_i \quad (5)$$

$$\nabla_{v_j} F(U, V) = -2 \sum_{i:(i,j)} u_i (r_{ij} - \langle u_i, v_j \rangle) + 2\lambda v_j. \quad (6)$$

The number of iterations relies on the stopping criteria. A typical criteria is to set a small threshold ε , such that the training stops when the gradient $\nabla_{u_i} F$ and $\nabla_{v_j} F$ (or one of them) are smaller than ε .

Distributed MF

In the distributed MF scenario, users hold their rating information locally and the model is trained on their joint data. To achieve this goal, we leverage a distributed MF method, which decomposes the iterative updating process into two parts that are performed on the user side and the server side, respectively. In particular, (3) is executed on user i 's device, namely *user update*, and (4) is performed on the server, called *server update*. This decomposition prevents the server from directly knowing users' raw preference data or learned profiles.

Algorithm 1 shows our user-level distributed MF method. The server keeps providing the latest item profile matrix V for all the users to download. Having the latest downloaded V and his/her own rating information, each user i performs local updates and computes $Gradient_{u_i}$, which will be sent back to the server to update item profiles.

GRADIENTS LEAK INFORMATION

With the enacting of many privacy-preserving regularizations (e.g., GDPR), dealing with personal data is sensitive and critical. Because GDPR has strict restrictions on collecting and using personal data, e.g., companies cannot give users' data to other companies without users' explicit permission. Currently, few companies can meet the requirements, and the violation may face a hefty fine. According to the definition in GDPR, users' rating data belongs to personal data (Personal data means any information relating to an identified or identifiable natural person; an identifiable natural person is one who can be identified, directly or indirectly (more information can be found at <https://gdpr.eu/eu-gdpr-personal-data>)). because we can link the ratings to natural person even without identity information (e.g., using linkage attack), thus it needs to be carefully treated. In FL, users keep the rating data locally and only upload certain gradients. Thus, GDPR can be satisfied as long as the uploaded information leaks no private data. Algorithm 1 shows a framework that allows the server to build an MF recommendation system on a distributed dataset, i.e., users keep rating information locally. In this framework, users iteratively send *Gradient* information to the server in plain text. Next, we are going to prove that such a distributed MF system cannot protect users' rating information against the server, i.e., the server can deduce users' rating data using the *Gradient*. Zhu *et al.*¹⁰ have shown that raw data could be recovered from the gradients of deep neural networks (e.g., CNN) by training an attack model, and FL with DP cannot defend such attacks. Thus a privacy-preserving solution (i.e., HE) is required to avoid private data leakage from gradients and guarantee that the FedMF can satisfy the GDPR regulation.

For user-vector $u_{i\cdot}$, suppose the user rated item-set is M_i . We will have the following equation at time t

$$u_i^t(r_{ij} - \langle u_i^t, v_j^t \rangle) = G_j^t, \quad j \in M_i \quad (7)$$

where G is the gradient to be uploaded from the user i to the server for updating item-profiles. Similarly, at time $t + 1$

$$u_i^{t+1}(r_{ij} - \langle u_i^{t+1}, v_j^{t+1} \rangle) = G_j^{t+1}, \quad j \in M_i. \quad (8)$$

From (3) and (5), the correlation between U_t and U_{t+1} is

$$\begin{aligned} (1 - 2\lambda\gamma)u_i^t - u_i^{t+1} \\ = -2\gamma \sum_{j \in M_i} v_j^t * (r_{ij} - \langle u_i^t, v_j^t \rangle). \end{aligned} \quad (9)$$

Take a close look at the element-wise calculation of (7)–(9), where we denote the latent user, and item vectors (u_i and v_j) are D dimension

$$\begin{cases} u_{i1}^t(r_{ij} - \sum_{m=1}^D u_{im}^t v_{jm}^t) = G_{j1}^t \\ \vdots \\ u_{ik}^t(r_{ij} - \sum_{m=1}^D u_{im}^t v_{jm}^t) = G_{jk}^t \\ \vdots \\ u_{iD}^t(r_{ij} - \sum_{m=1}^D u_{im}^t v_{jm}^t) = G_{jD}^t \end{cases} \quad (10)$$

$$\begin{cases} u_{i1}^{t+1}(r_{ij} - \sum_{m=1}^D u_{im}^{t+1} v_{jm}^{t+1}) = G_{j1}^{t+1} \\ \vdots \\ u_{ik}^{t+1}(r_{ij} - \sum_{m=1}^D u_{im}^{t+1} v_{jm}^{t+1}) = G_{jk}^{t+1} \\ \vdots \\ u_{iD}^{t+1}(r_{ij} - \sum_{m=1}^D u_{im}^{t+1} v_{jm}^{t+1}) = G_{jD}^{t+1} \end{cases} \quad (11)$$

$$\begin{cases} (1 - 2\lambda\gamma)u_{i1}^t - u_{i1}^{t+1} \\ = -2\gamma \sum_{n=1}^N v_{n1}^t (r_{in} - \sum_{m=1}^D u_{im}^t v_{nm}^t) \\ \vdots \\ (1 - 2\lambda\gamma)u_{ik}^t - u_{ik}^{t+1} \\ = -2\gamma \sum_{n=1}^N v_{nk}^t (r_{in} - \sum_{m=1}^D u_{im}^t v_{nm}^t) \\ \vdots \\ (1 - 2\lambda\gamma)u_{iD}^t - u_{iD}^{t+1} \\ = -2\gamma \sum_{n=1}^N v_{nD}^t (r_{in} - \sum_{m=1}^D u_{im}^t v_{nm}^t). \end{cases} \quad (12)$$

Now we turn to analyze the k th entry of u_i , u_{ik} . From (10), we have

$$\frac{u_{ik}^t}{u_{i(k+1)}^t} = \frac{G_{ik}^t}{G_{i(k+1)}^t} \quad (13)$$

$$r_{ij} - \sum_{m=1}^D u_{im}^t v_{jm}^t = \frac{G_{jk}^t}{u_{ik}^t}. \quad (14)$$

Plug (13) into (12), we will have

$$(1 - 2\lambda\gamma)u_{ik}^t - u_{ik}^{t+1} = -2\gamma \frac{1}{u_{ik}^t} \sum_{n=1}^N v_{nk}^t G_{nk}^t. \quad (15)$$

Thus, we can represent u_{ik}^{t+1} using u_{ik}^t as

$$u_{ik}^{t+1} = (1 - 2\lambda\gamma)u_{ik}^t + 2\gamma \frac{1}{u_{ik}^t} \sum_{n=1}^N v_{nk}^t G_{nk}^t. \quad (16)$$

From (10) and (11), we have

$$\frac{G_{jk}^t}{u_{ik}^t} + \sum_{m=1}^D u_{im}^t v_{jm}^t = \frac{G_{jk}^{t+1}}{u_{ik}^{t+1}} + \sum_{m=1}^D u_{im}^{t+1} v_{jm}^{t+1}. \quad (17)$$

Plug (16) into (17)

$$\begin{aligned} & \frac{G_{jk}^t}{u_{ik}^t} - \frac{G_{jk}^{t+1}}{(1-2\lambda\gamma)u_{ik}^t + 2\gamma \frac{1}{u_{ik}^t} \sum_{n=1}^N v_{nk}^t G_{nk}^t} \\ &= \sum_{m=1}^D \left[((1-2\lambda\gamma)u_{im}^t + 2\gamma \frac{1}{u_{im}^t} \sum_{n=1}^N v_{nm}^t G_{nm}^t) v_{jm}^{t+1} - u_{im}^t v_{jm}^t \right]. \end{aligned} \quad (18)$$

Let $\alpha_k = 2\gamma \sum_{n=0}^{N-1} v_{nk}^t G_{nk}^t$

$$\begin{aligned} & \frac{G_{jk}^t}{u_{ik}^t} - \frac{G_{jk}^{t+1}}{u_{ik}^t + \frac{\alpha_k}{u_{ik}^t}} \\ &= \sum_{m=1}^D \left[((1-2\lambda\gamma)v_{jm}^{t+1} - v_{jm}^t)u_{im}^t + \frac{\alpha_m v_{jm}^{t+1}}{u_{im}^t} \right]. \end{aligned} \quad (19)$$

From (13), we can have

$$u_{im}^t = \frac{G_{jm}^t}{G_{jk}^t} u_{ik}^t. \quad (20)$$

Plug (20) into (19)

$$\begin{aligned} & \frac{G_{jk}^t}{u_{ik}^t} - \frac{G_{jk}^{t+1}}{u_{ik}^t + \frac{\alpha_k}{u_{ik}^t}} \\ &= \frac{u_{ik}^t}{G_{jk}^t} \sum_{m=1}^D \left[((1-2\lambda\gamma)v_{jm}^{t+1} - v_{jm}^t) G_{jm}^t \right] \\ & \quad + \frac{G_{jk}^t}{u_{ik}^t} \sum_{m=1}^D \left[\frac{\alpha_m v_{jm}^{t+1}}{G_{jm}^t} \right]. \end{aligned} \quad (21)$$

Denote β_j and γ_j as follows:

$$\begin{cases} \beta_j = \sum_{m=1}^D \left[((1-2\lambda\gamma)v_{jm}^{t+1} - v_{jm}^t) G_{jm}^t \right] \\ \gamma_j = \sum_{m=1}^D \left[\frac{\alpha_m v_{jm}^{t+1}}{G_{jm}^t} \right]. \end{cases} \quad (22)$$

We will have

$$\frac{G_{jk}^t}{u_{ik}^t} - \frac{G_{jk}^{t+1}}{u_{ik}^t + \frac{\alpha_k}{u_{ik}^t}} = \frac{u_{ik}^t}{G_{jk}^t} \beta_j + \frac{G_{jk}^t}{u_{ik}^t} \gamma_j. \quad (23)$$

Since we know there must be one real scalar of u_{ik}^t that satisfies (23). We can use some iterative methods to compute a numeric solution of (23), e.g., Newton's method.

After getting u_{ik}^t , we can use (10) to compute r_i , which can be written as

$$r_{ij} = \frac{G_{jk}^t}{u_{ik}^t} + \sum_{m=1}^D u_{im}^t v_{jm}^t. \quad (24)$$

In summary, knowing the gradients of a user uploaded in two continuous steps, we can infer this user's rating information. Thus, we propose a secure MF framework based on HE, which will be elaborated in the next section.

FEDMF: FEDERATED MATRIX FACTORIZATION

To overcome this information leakage problem, we propose to encode the gradients such that the server cannot inverse the encoding process. Then, the encoded data leaks no information. Meanwhile, the server should still be able to perform updates using the encoded gradients. One way to achieve such a goal is by using HE.

Figure 1 shows a framework of our method, called *FedMF*. Two types of participants are involved in this framework, the server and the users. As previously illustrated in the "Horizontal FL" section, we assume that the server is honest-but-curious, the users are honest, and the privacy of the users is protected against the server.

Key Generation: As the typical functions involved in HE (see "Additively HE" section), we first generate the *public key* and *secret key*. The key generation process is carried out on one of the users. The *public key* is known to all the participants including the server. And the *secret key* is only shared between users and needs to be protected against the server. After the keys are generated, different TLS/SSL secure channels will be established for sending the *public key* and *secret key* to the corresponding participants.

Parameter Initialization: Before starting the MF process, some parameters need to be initialized. The item profile matrix is initialized at the server side while the user profile matrix is initialized by each user locally.

Matrix Factorization: Major steps include the following.

- 1) The server encrypts item profile V using *public key*, getting the ciphertext C_V . From now on, the latest C_V is prepared for all users' download.
- 2) Each user downloads the latest C_V from the server, and decrypts it using *secret key*, getting the plaintext of V . V is used to perform local update and compute the gradient G . Then, G is encrypted using *public key*, getting ciphertext C_V . Then, a TLS/SSL secure channel is built, C_V is sent back to the server via this secure channel.
- 3) After receiving a user's encrypted gradient, the server updates the item profile using this

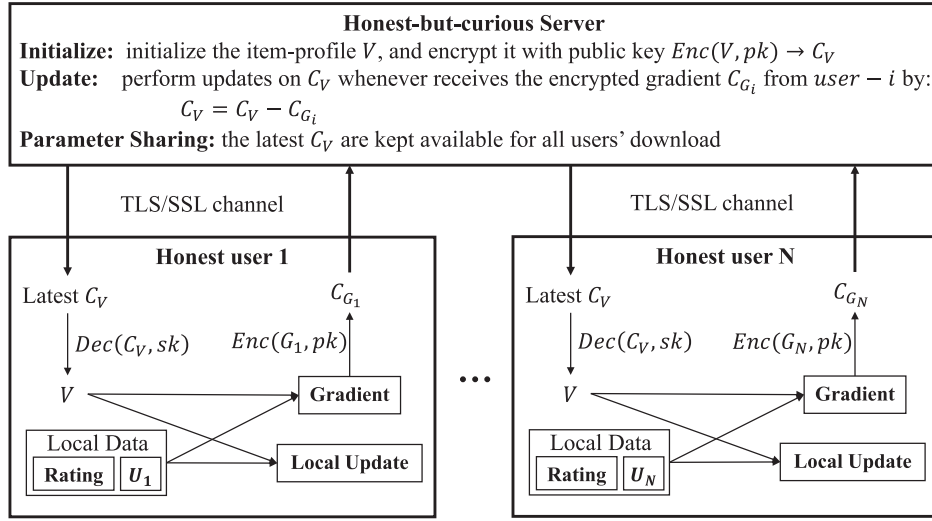


FIGURE 1. Overview of FedMF.

ciphertext : $C_V^{t+1} = C_V^t - C_G$. Afterward, the latest C_V is prepared for users' downloading.

- 4) Step 2 and 3 are iteratively executed until convergence.

Security Against Server: As shown in Figure 1, only ciphertext is sent to the server in FedMF. So no bit of information will be leaked to the server as long as our HE system ensures ciphertext indistinguishability against chosen plaintext attacks.²

No Accuracy Decline: We also claim that FedMF is accuracy equivalent to the user-level distributed MF. This is because the parameter updating process is the same as the distributed MF (see "User-Level Distributed Matrix Factorization" section) if the HE part is removed.

PROTOTYPE AND EVALUATION

In this section, we choose Paillier encryption method⁶ to instantiate a prototype of our system and use a real movie rating dataset to evaluate the prototype.

Prototype Implementation

We use Paillier encryption⁶ to build a prototype of FedMF. Paillier encryption is a probabilistic encryption schema based on composite residuosity problem. Given *public* key and encrypted plaintext, paillier encryption has the following homomorphic property operations:

- op1. $E(m_1) \cdot E(m_2) \pmod{n^2} = E(m_1 + m_2 \pmod{n})$
- op2. $E(m_1) \cdot g^{m_2} \pmod{n^2} = E(m_1 + m_2 \pmod{n})$
- op3. $E(m_1)^{m_2} \pmod{n^2} = E(m_1 m_2 \pmod{n})$

Typically, paillier encryption requires the plaintext to be a positive integer. But in our system, the data are all in the form of floating point numbers and some of them might be negative. Thus, we need to extend the encryption method to support our system.

Float. In brief, a base exponent was multiplied to the decimal, the integer part of the multiplication result I and the base exponent e was treated as the integer representation of the floating point number, i.e., (I, e) . In the encryption process, only I is encrypted, the ciphertext will be (C_I, e) . Then, the ciphertext with the same e can directly conduct operation op1 to get the encrypted summation of plaintext, ciphertext with different e needed to recalibrate such that e is the same. In practice, we use the same base exponent such that no information will leak from e .

Negative number. A *max number* parameter is set to handle the negative numbers. The *max number* can be set to half of n in *pk*, which means we assume all of our data is smaller than *max number*, such an assumption is easy to satisfy since n is usually set to a very large prime number. Then, we perform *mode n* on all the plaintext, all the positive number have no changes and all the negative numbers become positive numbers greater than *max number*. In the decryption process, if the decrypted plaintext is greater than *max number*, we minus n to get the correct negative plaintext.

FullText or PartText. Usually the rating or feedback comprises a sparse matrix⁵ which means the amount of feedback from a user could be very limited. Therefore, two different settings are implemented in our system. Both of them follow the overall steps of

TABLE 1. Time consumption (per iteration) on MovieLens-small dataset (second).

#Item	#Users	#Rating	Regular MF	Distributed MF	FedMF PartText	FedMF FullText
32	610	356	0.00429	0.00321	2.963	15.253
64	610	738	0.0082	0.00604	5.119	36.938
128	610	1420	0.0161	0.0104	9.518	95.021
256	610	2941	0.0308	0.0195	18.286	224.644
512	610	5227	0.0624	0.0378	35.759	487.755
1024	610	10615	0.131	0.0763	69.622	1004.051
2048	610	20758	0.251	0.154	133.113	2035.324
4096	610	42951	0.495	0.312	278.168	4080.073
8192	610	84994	0.997	0.618	560.156	8200.046

FedMF but are slightly different at the user uploading process. In one setting called *FullText*, users upload gradients for all the items; the gradient is set to 0 if a user does not rate an item. In the other setting called *PartText*, users only upload the gradients of the rated items. They both have advantages and disadvantages, *PartText* leaks information about which items the user has rated but has higher computation efficiency, *FullText* leaks no information but needs more computation time.

We utilize an open source python package, *python-paillier* (<https://github.com/n1analytics/python-paillier>) to accomplish the encryption part in our prototype system.

Evaluation

Dataset: We use two datasets in our experiments. The first one is MovieLens small dataset which is frequently used in other homomorphic-encrypted MF works such as by Kim *et al.*,⁴ and it contains 100 K ratings made by 610 users on 9724 movies. We also test our method on a large dataset: the MovieLens full dataset which contains 27 M ratings made by 280 K users on 58 K items.

Parameters: In Paillier encryption, we set the length of *public key* to 1024. The bandwidth of communication is set to 1 Gb/s. In the MF process, we set the dimension of the user and item profile to 100.

Environment: All the test experiments are performed on a server with 3.7 GHz 6-core CPU and 32 GB RAM, where the operating system is Windows and the programming language is Python. We used a module called *gmpy* (*gmpy* is a c-coded Python extension module that supports multiple-precision arithmetic, <https://github.com/aleaxit/gmpy>) to accelerate the HE part in Python such that it is as fast as C++ implementation.

Performance: Since neither distributed computing or HE mechanisms will affect the computation values, FedMF will output the same user and item profiles as

TABLE 2. Time consumption (per iteration) on MovieLens full dataset.

#Item	#Users	#Rating	Regular MF (second)	Distributed MF (second)	FedMF PartText (hour)	FedMF FullText (hour)*
58 K	5 K	493 K	5.833	3.520	0.4	~110
58 K	10 K	994 K	11.557	5.612	0.6	~220
58 K	20 K	1.9 M	23.646	9.125	1.0	~470
58 K	40 K	3.8 M	47.198	17.202	1.9	~970
58 K	80 K	7.7 M	93.813	31.350	2.8	~1500
58 K	160 K	15 M	184.706	63.340	8.1	~4000
58 K	283 K	27 M	324.601	114.592	15.3	~8400

* We note that running FedMF FullText is very time consuming (e.g., one iteration for all the users/items can be up to one year). So, currently we estimate the running time rather than actual running. Particularly, the user's local running time is linearly correlated with the number of ratings, thus, we sample 100 users and calculate the averaged time consumption on one rating and estimate the running time for all the users.

the original MF algorithm. Hence, the major objective of the experiments is testing the computation time of FedMF.

Fixing the number of users to 610, Table 1 shows the time consumption of each iteration of RegularMF, DistributedMF, FedMF-PartText, and FedMF-FullText (one iteration means all of the 610 users' uploaded gradients are used to update the item profiles once). We assume that all the users update in parallel and the server starts the aggregation only after receives all the users' gradients. Thus, the time consumption bottleneck of users' update is the one who has the most rating data, because it requires the most significant running time. For both *PartText* and *FullText*, the time consumption is quite good when there are not too many items, and the time efficiency decreases when more items are given. Roughly, the time consumed for each iteration linearly increases with the number of items. Compared with *Fulltext*, *PartText* is more efficient but it leaks some information. Particularly, *PartText* is nearly 15 times faster than the *Fulltext* solution. The results also show that FedMF-PartText is about four orders of magnitude slower compared with DistributedMF because of the significant time consumption overhead caused by HE. However, if not using HE, the security of users' rating data cannot be guaranteed due to the possible gradient leakage. Hence, this is a tradeoff between security and efficiency. Table 2 shows the time consumption on MovieLens full dataset. The time consumption of both PartText and FullText setting linearly increases with the number of users. When using the whole dataset (i.e., 280 K users), it requires more than 15 h to train one iteration in FedMF-PartText. This reveals that for applying FedMF or other HE-enhanced FL mechanisms to larger

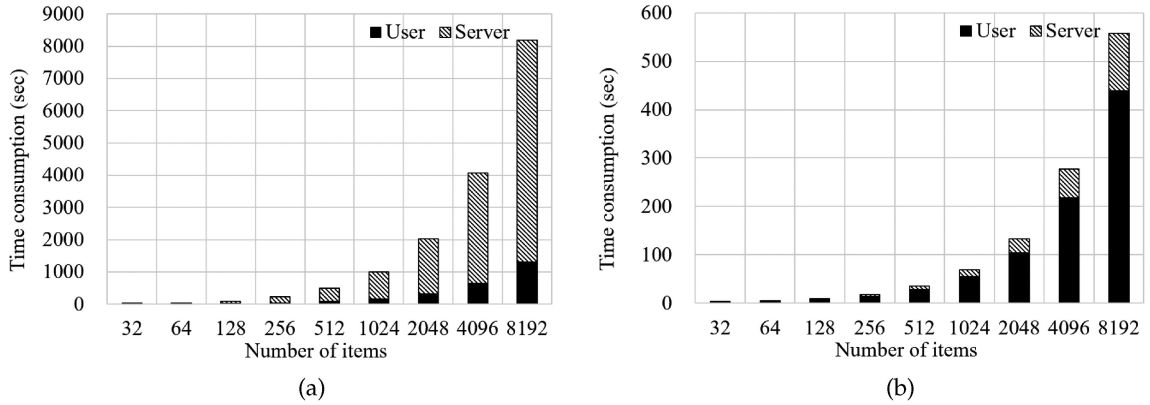


FIGURE 2. User-server time consumption ratio. (a) FullText. (b) PartText.

datasets, improving computation efficiency is an urgently required future research direction.

Figure 2(a) and (b) shows the ratio of the user and server updating time when the number of items changes. The communication time is dismissed from the figures because it is too small compared with the user and server updating time. For example, 256 MB of gradient data needs to be sent to the server when the item number is 8192, and it will cost about 4 s using a bandwidth of 1 Gb/s. In the FullText setting, the server update costs most of the time. In contrast, the user update consumes most of the time in the PartText setting. The reason is that the server's updating time significantly increases in the FullText setting, while

the users' updating time is barely influenced because of the bottleneck mentioned above.

To verify that our FedMF system has no accuracy decline, we performed experiments on a small-scale dataset (Such a small dataset can help to reduce model convergence time consumption.) which contains 40 movies, 7 users, and 130 ratings. For each user, we select three rating records for testing and the other ratings compose the training dataset. Root mean square error (RMSE) is used as metric to evaluate the model's performance. As illustrated in Figure 3 and Table 3, the RMSE of RegularMF/DistributedMF and FedMF are almost the same, showing that FedMF has no accuracy decline. It is worth nothing that the float numbers are transferred to integers using a fixed precision in pailliar encryption, thus the results of pailliar-encryption based methods (e.g., FedMF) will be slightly different from plaintext methods (e.g., DistributedMF).

Real world implementation: In real-world implementation, different from experiment settings, users' devices may not online simultaneously. Such circumstance will not affect FedMF system since the server can update item vector immediately after receiving gradient from one user instead of waiting all users' gradients received. The training process can be carried out when the users' devices are charging and remains extra computing resources.

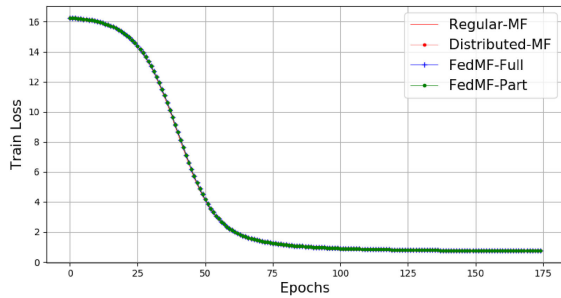


FIGURE 3. Comparison of training loss.

TABLE 3. Comparison of converged results.

	RegularMF	DistributedMF	FedMF PartText	FedMF FullText
RMSE	1.39692*	1.39761*	1.39654	1.39654

*The results of RegularMF and DistributedMF are slightly different because in DistributedMF we make the server updates item-vector only when all the users have their gradient uploaded, and in RegularMF the server updates item-vector immediately after receiving gradient from an arbitrary user.

CONCLUSION AND FUTURE WORK

In this article, we propose a novel secure MF framework in federated machine learning, called FedMF. More specifically, we first prove that a distributed MF system where users send gradients to the server in forms of plaintext will leak users' rating information. Then, we design an HE-based secure MF framework. We have proved that our system is secure against an

honest-but-curious server, and the accuracy is the same as the MF on users' raw data.

Experiments on real-world data show that FedMF's time efficiency is acceptable when the number of items is small. Also note that our system's time consumption linearly increases with the number of items. To make FedMF more practical in reality, we still face several challenges:

More efficient HE. FedMF shows large time consumption overhead because of the low efficiency of the HE algorithm. If we can improve the HE's efficiency when conducting operations on ciphertext, our system's performance will increase.

Between FullText and PartText. Our experiments have shown that *PartText* is much more efficient than *FullText*, but *PartText* reveals the set of items rated by a user. This information, without the exact rating scores, may still leak users' sensitive information. Perhaps we can ask users to upload more gradients than only the rated items, but not all the items, so as to increase efficiency compared to *FullText*, while not leaking the exactly rated item set.

More Secure Definitions. We currently use a typical horizontal FL secure definition, which assumes honest participants and an honest-but-curious server. Such a secure definition is relatively weak because there might be malicious users in real-world applications. Malicious users may collude with the server (i.e., share the private key) to attack other users, causing privacy issues. They also may perform backdoor attacks to the system, which will bring security issues. In the future, we will explore more challenging secure definitions, such as building a secure system where the server is honest-but-curious, and some participants are malicious, and the malicious participants may collude with the server.

ACKNOWLEDGMENTS

The work of L. Wang was supported by the NSFC under Grant 61972008. The work of Q. Yang was supported by the National Key Research and Development Program of China under Grant 2018AAA0101100. This research was partially supported by the collaboration project at HKUST-DiDi Joint Lab. Di Chai and Leye Wang contributed equally to this work.

REFERENCES

1. M. Ammad-ud-din *et al.*, "Federated collaborative filtering for privacy-preserving personalized recommendation system," 2019, *arXiv:1901.09888*.
2. L. T. Phong, Y. Aono, T. Hayashi, L. Wang, and S. Moriai, "Privacy-preserving deep learning via additively homomorphic encryption," *IEEE Trans. Inf. Forensics Secur.*, vol. 13, no. 5, pp. 1333–1345, May 2018.
3. A. Berlioz, A. Friedman, M. A. Kaafar, R. Boreli, and S. Berkovsky, "Applying differential privacy to matrix factorization," in *Proc. 9th ACM Conf. Recommender Syst.*, 2015, pp. 107–114.
4. S. Kim, J. Kim, D. Koo, Y. Kim, H. Yoon, and J. Shin, "Efficient privacy-preserving matrix factorization via fully homomorphic encryption," in *Proc. 11th ACM Asia Conf. Comput. Commun. Secur.*, 2016, pp. 617–628.
5. Y. Koren, R. Bell, and C. Volinsky, "Matrix factorization techniques for recommender systems," *Computer*, vol. 8, pp. 30–37, 2009.
6. P. Paillier, "Public-key cryptosystems based on composite degree residuosity classes," in *Proc. Int. Conf. Theory Appl. Cryptographic Techn.*, 1999, pp. 223–238.
7. L. Wang, D. Zhang, Y. Wang, C. Chen, X. Han, and A. M'hamed, "Sparse mobile crowdsensing: Challenges and opportunities," *IEEE Commun. Mag.*, vol. 54, no. 7, pp. 161–167, Jul. 2016.
8. Q. Yang, Y. Liu, T. Chen, and Y. Tong, "Federated machine learning: Concept and applications," *ACM Trans. Intell. Syst. Technol.*, vol. 10, no. 2, 2019, Art. no. 12.
9. Qiang Yang, Y. Liu, Y. Cheng, Yan Kang, T. Chen, and H. Yu, "Federated learning," *Synthesis Lectures Artif. Intell. Mach. Learn.*, vol. 13, no. 3, pp. 1–207, 2019.
10. L. Zhu, Z. Liu, and S. Han, "Deep leakage from gradients," *Adv. Neural Inf. Process. Syst.*, pp. 14774–14784, 2019, Art. no. 12.

DI CHAI is currently working toward a Ph.D. degree in computer science and engineering with Hong Kong University of Science and Technology, Hong Kong. His research interests include federated learning and privacy-preserving machine learning. He received the master's degree in science from Hong Kong University of Science and Technology, in 2018. Contact him at dchai@connect.ust.hk.

LEYE WANG is currently an Assistant Professor with Key Lab of High Confidence Software Technologies, Peking University, Beijing, China. His research interests include ubiquitous computing, mobile crowdsensing, and urban computing. He received a Ph.D. degree in computer science from the Institut Telecom SudParis and University Paris 6, France, in 2016, and was a Postdoc Researcher with Hong Kong University of Science and Technology, Hong Kong. Contact him at leye-wang@pku.edu.cn.

KAI CHEN is currently an Associate Professor with the Department of Computer Science and Engineering, Hong Kong University of Science and Technology, Hong Kong. He received a Ph.D. degree in computer science from Northwestern University, Evanston, IL, USA, in 2012. His research interest includes data center networking, machine learning systems, and privacy-preserving AI. Contact him at kai-chen@cse.ust.hk.

QIANG YANG is currently the Chief AI Officer with WeBank, Shenzhen, China, and a Chair Professor of Engineering with

Hong Kong University of Science and Technology, Hong Kong. His research interests include data mining and artificial intelligence including machine learning, planning, and case-based reasoning. He received a Ph.D. degree in computer science from the University of Maryland, College Park, MD, USA, in 1989. He was Founding Editor-in-Chief of *ACM Transactions on Intelligent Systems and Technology* and *IEEE TRANSACTIONS ON BIG DATA*. He is a Fellow of the IEEE, AAAI, and ACM. He is the corresponding author of this article. Contact him at qyang@cse.ust.hk.



IT Professional
TECHNOLOGY SOLUTIONS FOR THE ENTERPRISE

CALL FOR ARTICLES

IT Professional seeks original submissions on technology solutions for the enterprise. Topics include

- emerging technologies,
- cloud computing,
- Web 2.0 and services,
- cybersecurity,
- mobile computing,
- green IT,
- RFID,
- social software,
- data management and mining,
- systems integration,
- communication networks,
- datacenter operations,
- IT asset management, and
- health information technology.

We welcome articles accompanied by web-based demos. For more information, see our author guidelines at www.computer.org/itpro/author.htm.

WWW.COMPUTER.ORG/ITPRO

75 YEARS
IEEE COMPUTER SOCIETY

IEEE