



# FedeRank: User Controlled Feedback with Federated Recommender Systems

Vito Walter Anelli, Yashar Deldjoo, Tommaso Di Noia, Antonio Ferrara<sup>(✉)</sup>,  
and Fedelucio Narducci

Politecnico di Bari, Bari, Italy

{vitowalter.anelli,yashar.deldjoo,tommaso.dinoia,antonio.ferrara,  
fedelucio.narducci}@poliba.it

**Abstract.** Recommender systems have shown to be a successful representative of how data availability can ease our everyday digital life. However, data privacy is one of the most prominent concerns in the digital era. After several data breaches and privacy scandals, the users are now worried about sharing their data. In the last decade, Federated Learning has emerged as a new privacy-preserving distributed machine learning paradigm. It works by processing data on the user device without collecting data in a central repository. We present FedeRank (<https://split.to/federank>), a federated recommendation algorithm. The system learns a personal factorization model onto every device. The training of the model is a synchronous process between the central server and the federated clients. FedeRank takes care of computing recommendations in a distributed fashion and allows users to control the portion of data they want to share. By comparing with state-of-the-art algorithms, extensive experiments show the effectiveness of FedeRank in terms of recommendation accuracy, even with a small portion of shared user data. Further analysis of the recommendation lists' diversity and novelty guarantees the suitability of the algorithm in real production environments.

**Keywords:** Recommender systems · Collaborative filtering · Federated learning · Learning to rank

## 1 Introduction

Recommender Systems (RSs) are well-known information-filtering systems widely adopted for mitigating the information-overload problem. Indeed, the broad amount of items and services has caused a cognitive impairment that takes the name of over-choice, or choice overload. RSs have proved to be very useful in making possible personalized access to these catalogs of items. These systems are generally hosted on centralized servers and train their models by exploiting massive proprietary and sensitive data. However, public awareness related to data collection was spurred and increased. In recent years, an increasing number of countries have introduced regulations to protect user privacy and data security.

Representative examples are the GDPR in the European Union [15], the CCPA in California [8], and the Cybersecurity Law in China [42]. Such policies prohibit free data circulation and force personal data to remain isolated and fragmented.

In this context, Google has recently proposed Federated Learning (FL) as a privacy-by-design technique which tackles data isolation while meeting the need for big data [23, 33]. FL trains a global machine-learning model by leveraging both users' data and personal devices' computing capabilities. Unlike previous approaches, it keeps data on the devices (e.g., smartphones, tablets, etc.) without sharing it with a central server. Today, FL is considered the best candidate to face the data privacy, control and property challenges posed by the data regulations.

Among the recommendation paradigms proposed in the literature, Collaborative Filtering (CF) demonstrated a very high accuracy [32, 47]. The strength of CF recommendation algorithms is that users who expressed similar tastes in the past tend to agree in the future as well. One of the most prominent CF approaches is the Latent Factor Model (LFM) [26]. LFMs uncover users and items latent representation, whose linear interaction can explain observed feedback.

In this paper, we introduce FedeRank, a novel factorization model that embraces the FL paradigm. A disruptive effect of employing FedeRank is that users participating in the federation process can decide if and how they are willing to disclose their private sensitive preferences. Indeed, FedeRank mainly leverages non-sensitive information (e.g., items the user has not experienced). Here, we show that even only a small amount of sensitive information (i.e., items the user has experienced) lets FedeRank reach a competitive accuracy. How incomplete data impacts the performance of the system is an entirely unexplored field. Analogously, it is still to establish the minimum amount of data necessary to build an accurate recommendation system [46]. At the same time, preserving privacy at the cost of a worse tailored recommendation may frustrate users and reduce the "acceptance of the recommender system" [35]. In this work, instead of focusing on how to protect personal information from privacy breaches (that is explored in other active research fields), we investigate how to guarantee the users the control and property of their data as determined by regulations. The work's contributions are manifold due to the number of open challenges that still exist with the FL paradigm. To summarize, our contributions in this paper include:

- the development of the first, to the best of our knowledge, federated pair-wise recommendation system;
- an analysis of the impact of client-side computation amount;
- an investigation on the existing relation between incomplete data and recommendation accuracy, and an analysis of the algorithmic bias on the final recommendation lists, based on the data deprivation amount.

To this extent, we have carried out extensive experiments on three real-world datasets (*Amazon Digital Music*, *LibraryThing*, and *MovieLens 1M*) by considering two evaluation criteria: (a) the accuracy of recommendations measured by

exploiting precision and recall, (b) beyond-accuracy measures to evaluate the novelty, and the diversity of recommendation lists. The experimental evaluation shows that FedeRank provides high-quality recommendations, even though it leaves users in control of their data.

## 2 Related Work

In the last decades, academia and industry have proposed several competitive recommendation algorithms. Among the Collaborative Filtering algorithms, the most representative examples are undoubtedly Nearest Neighbors systems, Latent Factor Models, and Neural Network-based recommendation systems. The Nearest Neighbors scheme has shown its competitiveness for decades. After them, factorization-based recommendation emerged thanks to the disruptive idea of Matrix Factorization (MF). Nevertheless, several generalized/specialized variants have been proposed, such as FM [37], SVD++ [24], PITF [40], FPMC [39]. Unfortunately, rating-prediction-oriented optimization, like SVD, has shown its limits in the recommendation research [34]. Consequently, a new class of *Learning to Rank* algorithms has been developed in the last decade, mainly ranging from point-wise [28] to pair-wise [38], through list-wise [41] approaches. Among pair-wise methods, BPR [38] is one of the most adopted, thanks to its outstanding capabilities to correctly rank with an acceptable computational complexity. Finally, in the last years, methods exploiting the architectures of deep neural networks have established themselves in search and recommendation research.

To make RSs work properly (easing the user decision-making process and boosting the business), they need to collect user information related to attributes, demands, and preferences [20], jeopardizing the user’s privacy. In this scenario—and, more generally, in any scenario with a system learning from sensitive data—FL was introduced for meeting privacy shortcomings by horizontally distributing the model’s training over user devices [33]. Beyond privacy, FL has posed several other challenges and opened new research directions [21]. In the last years, it has extended to a more comprehensive idea of privacy-preserving decentralized collaborative ML approaches [45], ranging from horizontal federations, where different devices (and local datasets) share the same feature space, to vertical federations, where devices share training samples that differ in feature space.

Some researchers focused the attention on the decentralized and distributed matrix-factorization approaches [12, 16]. However, in this work, we focus on federated learning principles theoretically and practically different from classical distributed approaches. Indeed, FL assumes the presence of a coordinating server and the use of private and self-produced data on each node. In general, distributed approaches do not guarantee these assumptions. Ammad-uddin *et al.* [3] propose a federated implementation of collaborative filtering, whose security limits are analyzed in [11], which uses the SVD-MF method for implicit feedback [19]. Here, the training is a mixture of Alternating Least Squares (ALS) and Stochastic Gradient Descent (SGD) for preserving users’ privacy. Nevertheless, incomprehensibly, almost no work addressed top-N recommendation exploiting the “Learning to rank” paradigm. In this sense, one rare example is the work

by Kharitonov *et al.* [22], who recently proposed to combine evolution strategy optimization with a privatization procedure based on differential privacy. The previous work focuses neither on search or recommendation. Perhaps, like ours, it can be classified as a federated learning-to-rank algorithm. Finally, Yang *et al.* [46] identified some recent FL challenges and open research directions.

### 3 Approach

In this section, we introduce the fundamental concepts regarding the Collaborative Filtering recommendation using a Federated Learning scheme. Along with the problem definition, the notation we adopt is presented.

The recommendation problem over a set of users  $\mathcal{U}$  and a set of items  $\mathcal{I}$  is defined as the activity of finding for each user  $u \in \mathcal{U}$  an item  $i \in \mathcal{I}$  that maximizes a utility function  $g : \mathcal{U} \times \mathcal{I} \rightarrow \mathbb{R}$  [36]. Let  $\mathbf{X} \in \mathbb{R}^{|\mathcal{U}| \times |\mathcal{I}|}$  be the user-item matrix containing for each element  $x_{ui}$  an implicit feedback (e.g., purchases, visits, clicks, views, check-ins) of user  $u \in \mathcal{U}$  for item  $i \in \mathcal{I}$ . Therefore,  $\mathbf{X}$  only contains binary values,  $x_{ui} = 1$  and  $x_{ui} = 0$  denoting whether user  $u$  has consumed or not item  $i$ , respectively.

The recommendation model is based on Factorization approach, originally introduced by Matrix Factorization [27], that became popular in the last decade thanks to its state-of-the-art recommendation accuracy [29]. This technique aims to build a model  $\Theta$  in which each user  $u$  and each item  $i$  is represented by the embedding vectors  $\mathbf{p}_u$  and  $\mathbf{q}_i$ , respectively, in the shared latent space  $\mathbb{R}^F$ . Let assume  $\mathbf{X}$  can be factorized such that the dot product between  $\mathbf{p}_u$  and  $\mathbf{q}_i$  can explain any observed user-item interaction  $x_{ui}$ , and any non-observed interaction can be estimated as  $\hat{x}_{ui}(\Theta) = b_i(\Theta) + \mathbf{p}_u^T(\Theta) \cdot \mathbf{q}_i(\Theta)$  where  $b_i$  is a term denoting the bias of the item  $i$ .

Among pair-wise approaches for learning-to-rank the items of a catalog, Bayesian Personalized Ranking [38] is the most broadly adopted, thanks to its capabilities to correctly rank with *acceptable* computational complexity. Given a training set defined by  $\mathcal{K} = \{(u, i, j) \mid x_{ui} = 1 \wedge x_{uj} = 0\}$ , BPR minimizes the ranking loss by exploiting the criterion  $\max_{\Theta} G(\Theta)$ , with  $G(\Theta) = \sum_{(u,i,j) \in \mathcal{K}} \ln \sigma(\hat{x}_{uij}(\Theta)) - \lambda \|\Theta\|^2$ , where  $\hat{x}_{uij}(\Theta) = \hat{x}_{ui}(\Theta) - \hat{x}_{uj}(\Theta)$  is a real value modeling the relation between user  $u$ , item  $i$  and item  $j$ ,  $\sigma(\cdot)$  is the sigmoid function, and  $\lambda$  is a model-specific regularization parameter to prevent overfitting. Pair-wise optimization applies to a wide range of recommendation models, including factorization. Hereafter, we denote the model  $\Theta = \langle \mathbf{P}, \mathbf{Q}, \mathbf{b} \rangle$ , where  $\mathbf{P} \in \mathbb{R}^{|\mathcal{U}| \times F}$  is a matrix whose  $u$ -th row corresponds to the vector  $\mathbf{p}_u$ , and  $\mathbf{Q} \in \mathbb{R}^{|\mathcal{I}| \times F}$  is a matrix in which the  $i$ -th row corresponds to the vector  $\mathbf{q}_i$ . Finally,  $\mathbf{b} \in \mathbb{R}^{|\mathcal{I}|}$  is a vector whose  $i$ -th element corresponds to the value  $b_i$ .

#### 3.1 FedeRank

FedeRank redesigns the original factorization approach for a federated setting. Indeed, the initial factorization model and its variants use a single, centralized

model, which does not guarantee users to control their data. FedeRank splits the pair-wise learning model  $\Theta$  among a central server  $S$  and a federation of users  $\mathcal{U}$ . Federated learning aims to optimize a global loss function by using data distributed among a federation of users’ devices. The rationale is that the server no longer collects private users’ data. Rather, it aggregates the results of some steps of local optimizations performed by clients, preserving privacy, ownership, and locality of users’ data [6]. Formally, let  $\Theta$  be the machine learning model parameters, and  $G(\Theta)$  be a loss function to minimize. In Federated learning, the users  $\mathcal{U}$  of a federation collaborate to minimize  $G$  (under the coordination of a central server  $S$ ) without sharing or exchanging their raw data. From an algorithmic point of view,  $S$  shares  $\Theta$  with the federation of devices. Then, the optimization problem of minimizing  $G$  is locally solved. Since each user participates to the federation with her personal data and with her personal client device, we will interchangeably use the terms “client”, “user”, and “device”.

To set up the framework, we consider the central server  $S$  holding a model  $\Theta_S = \langle \mathbf{Q}, \mathbf{b} \rangle$ , where  $\mathbf{Q} \in \mathbb{R}^{|\mathcal{I}| \times F}$  is a matrix in which  $i$ -th row represents the embedding  $\mathbf{q}_i$  for item  $i$  in the catalog, while the element  $b_i$  of  $\mathbf{b} \in \mathbb{R}^{|\mathcal{I}|}$  is the bias of item  $i$ . That is, the information on  $S$  only characterizes the items of the catalog. On the other hand, each user  $u \in \mathcal{U}$  holds a local model  $\Theta_u = \langle \mathbf{p}_u \rangle$ , where  $\mathbf{p}_u \in \mathbb{R}^F$  corresponds to the representation of user  $u$  in the latent space of dimensionality  $F$ . Each user holds a private interaction dataset  $\mathbf{x}_u \in \mathbb{R}^{|\mathcal{I}|}$ , which—compared to a centralized recommender system—corresponds to the  $\mathbf{X}$ ’s  $u$ -th row. The user  $u$  leverages her private dataset  $\mathbf{x}_u$  to build the local training set  $\mathcal{K}_u = \{(u, i, j) \mid x_{ui} = 1 \wedge x_{uj} = 0\}$ . Finally, the overall number of interactions in the system can be obtained by exploiting the local datasets. Let us define it as  $X^+ = \sum_{u \in \mathcal{U}} |\{x_{ui} \mid x_{ui} = 1\}|$ .

The training procedure iterates for  $E$  epochs, in each of which  $rpe$  rounds of communication between the server and the devices are performed. A round of communication envisages a **Distribution to Devices**  $\rightarrow$  **Federated Optimization**  $\rightarrow$  **Transmission to Server**  $\rightarrow$  **Global Aggregation** sequence. The notation  $\{\cdot\}_S^t$  denotes an object computed by the server  $S$  at round  $t$ , while  $\{\cdot\}_u^t$  indicates an object computed by a specific client  $u$  at round  $t$ .

(1) **Distribution to Devices.** Let  $\{\mathcal{U}^-\}_S^t$  be a subset of  $\mathcal{U}$  with cardinality  $m$ , containing  $m$  clients  $u \in \mathcal{U}$ . The set  $\{\mathcal{U}^-\}_S^t$  can be either defined by  $S$ , or the result of a request for availability sent by  $S$  to clients in  $\mathcal{U}$ . Each client  $u \in \{\mathcal{U}^-\}_S^t$  receives from  $S$  the latest snapshot of  $\{\Theta_S\}_S^{t-1}$ .

(2) **Federated Optimization.** Each user  $u \in \{\mathcal{U}^-\}_S^t$  generates the set  $\{\mathcal{K}_u^-\}_u^t$  containing  $T$  random triples  $(u, i, j)$  from  $\mathcal{K}_u$ . It is worth underlining that Rendle [38] suggests, for a centralized scenario, to train the recommendation model by randomly choosing the training triples from  $\mathcal{K}$ , to avoid data is traversed item-wise or user-wise, since this may lead to slow convergence. Conversely, in a federated approach, we require to train the model user-wise. Indeed, the learning is separately performed on each device ( $u$ ), that only knows the data in  $\mathcal{K}_u$ . Thanks to the user-wise traversing, FedeRank can decide who controls (the designer or the user) the number of triples  $T$  in the training set  $\{\mathcal{K}_u^-\}_u^t$ , to tune

the degree of local computation. With the local training set, the user  $u$  can compute her contribution to the overall optimization of  $\Theta_S$  with the following update:

$$\{\Delta\Theta_S\}_u^t = \{\Delta\langle\mathbf{Q}, \mathbf{b}\rangle\}_u^t := \sum_{(u,i,j) \in \{\mathcal{K}_u^-\}_u^t} \frac{\partial}{\partial\Theta_S} \ln \sigma(\hat{x}_{uij}(\{\Theta_S\}_S^{t-1}; \{\Theta_u\}_u^{t-1})), \quad (1)$$

plus a regularization term. At the same time, the client  $u$  updates its local model  $\Theta_u$ , and incorporates it in the current model by using:

$$\{\Delta\Theta_u\}_u^t = \{\Delta\langle\mathbf{p}_u\rangle\}_u^t := \sum_{(u,i,j) \in \{\mathcal{K}_u^-\}_u^t} \frac{\partial}{\partial\Theta_u} \ln \sigma(\hat{x}_{uij}(\{\Theta_S\}_S^{t-1}; \{\Theta_u\}_u^{t-1})), \quad (2)$$

plus a regularization term. The partial derivatives in Eq. 1 and 2 are straightforward, and can be easily computed by following the scheme proposed by Rendle *et al.* [38]. At the end of the federated computation, given a shared learning rate  $\alpha$ , each client can update its local model  $\Theta_u$ —containing the user profile—by aggregating the computed update:

$$\{\Theta_u\}_u^t := \{\Theta_u\}_u^{t-1} + \alpha\{\Delta\Theta_u\}_u^t. \quad (3)$$

(3) **Transmission to Server.** In a purely distributed architecture, each user in  $\mathcal{U}^-$  returns to  $S$  the computed update. Here, instead of sending  $\{\Delta\Theta_S\}_u^t$ , each user transmits a modified version  $\{\Delta\Theta_S^\Phi\}_u^t$ . To introduce this aspect of FedeRank, let us define  $\mathcal{F} = \{i, \forall(u, i, j) \in \{\mathcal{K}_u^-\}_u^t\}$ , and a randomized object  $\Phi = \langle\mathbf{Q}^\Phi, \mathbf{b}^\Phi\rangle$ , with  $\mathbf{Q}^\Phi \in \mathbb{R}^{|\mathcal{I}| \times F}$ , and  $\mathbf{b}^\Phi \in \mathbb{R}^{|\mathcal{I}|}$ . Each row  $\mathbf{q}_i^\Phi$  of  $\mathbf{Q}^\Phi$  and each element  $b_i^\Phi$  of  $\mathbf{b}^\Phi$  assume their value according to the probabilities:

$$\begin{aligned} P(\mathbf{q}_i^\Phi = \mathbf{1}, b_i^\Phi = 1 \mid i \in \mathcal{F}) &= \pi, & P(\mathbf{q}_i^\Phi = \mathbf{0}, b_i^\Phi = 0 \mid i \in \mathcal{F}) &= 1 - \pi, \\ P(\mathbf{q}_i^\Phi = \mathbf{1}, b_i^\Phi = 1 \mid i \notin \mathcal{F}) &= 1 \end{aligned} \quad (4)$$

Based on  $\{\mathbf{Q}^\Phi\}_u^t$  and  $\{\mathbf{b}^\Phi\}_u^t$ ,  $\Delta\Theta_S^\Phi$  can be computed as it follows:

$$\{\Delta\Theta_S^\Phi\}_u^t = \{\Delta\Theta_S\}_u^t \odot \{\Phi\}_u^t := \langle\{\Delta\mathbf{Q}\}_u^t \odot \{\mathbf{Q}^\Phi\}_u^t, \{\Delta\mathbf{b}\}_u^t \odot \{\mathbf{b}^\Phi\}_u^t\rangle, \quad (5)$$

where the operator  $\odot$  denotes the Hadamard product. This transformation is motivated by a possible privacy issue. The update  $\Delta\mathbf{Q}$  computed in Eq. 1 by user  $u$  is a matrix whose rows are non-zero in correspondence of the items  $i$  and  $j$  of all the triples  $(u, i, j) \in \mathcal{K}_u^-$  [38]. An analogous behavior can be observed for the elements of  $\Delta\mathbf{b}$ . Focusing on the non-zero elements, we observe that, for each triple  $(u, i, j) \in \mathcal{K}_u^-$ , the updates  $\{\Delta\mathbf{q}_i\}_u^t$  and  $\{\Delta\mathbf{q}_j\}_u^t$ , as well as  $\{\Delta b_i\}_u^t$  and  $\{\Delta b_j\}_u^t$ , show the same absolute value with opposite sign [38]. In fact, this makes completely distinguishable for the server the consumed and the non-consumed items of user  $u$ , allowing  $S$  to reconstruct  $\mathcal{K}_u^-$ , thus raising a privacy issue. Since our primary goal is to put users in control of their data, we leave users the possibility to choose a fraction  $\pi$  of positive item updates to send.

The remaining positive item updates (a fraction  $1 - \pi$ ) are masked by setting them to zero, by means of the transformation in Eq. 5. Conversely, the negative updates are always sent to  $S$ , since their corresponding rows are always multiplied by a  $\mathbf{1}$  vector. Indeed, these updates are related to non-consumed items, which are indistinguishably negative or missing values, assumed to be *non-sensitive* data.

(4) **Global Aggregation.** Once  $S$  has received  $\{\Delta\theta_S^\Phi\}_u^t$  from all clients  $u \in \mathcal{U}^-$ , it aggregates the received updates in  $\mathbf{Q}$  and  $\mathbf{b}$  to build the new global model, with  $\alpha$  being the learning rate:

$$\{\theta_S\}_S^t := \{\theta_S\}_S^{t-1} + \alpha \sum_{u \in \mathcal{U}^-} \{\Delta\theta_S^\Phi\}_u^t. \quad (6)$$

Although Federated Learning was conceived as a privacy-by-design paradigm for distributed machine learning, it still does not provide formal privacy guarantees. Malicious actors might acquire different information. They may be able to analyze remote devices or communication flows in the network or infer users' private data by inspecting updates received on the server [21]. Possible solutions include the encryption of data on local devices, the network traffic, or the adoption of secure multi-party computation [7]. Moreover, local differential privacy can guarantee that even if an adversary can inspect the communication between a user and the central server, she can learn only limited information [13, 14, 43]. To date, FedeRank explicitly provides the needed APIs to work, out of the box, with encryption communication libraries, thus providing state-of-the-art privacy guarantees. We have chosen this solution since discussing privacy and security implications in FL is beyond our scope. In this way, the system designer can choose the privacy solution while disregarding the underlying machine learning model. Moreover, FedeRank can also be easily adapted to guarantee local differential privacy. Indeed, it is not due to chance the choice of putting the user in control of  $\Phi$ . Suppose the  $\Phi$  object also considers sending fake information. In that case, FedeRank becomes utterly compliant with the randomized response technique, which guarantees differential privacy [44].

## 4 Experiments

**Datasets.** We have investigated the performance of FedeRank considering three well-known datasets: *Amazon Digital Music* [31], *LibraryThing* [48], and *MovieLens 1M* [18]. The former includes the users' satisfaction feedback for a catalog of music tracks available with Amazon Digital Music service. It contains 1,835 users and 41,488 tracks, with 75,932 ratings ranging from 1 to 5. *LibraryThing* collects the users' ratings on a book catalog. It captures the interactions of 7,279 users on 37,232 books. It provides more than two million ratings with 749,401 unique ratings in a range from 1 to 10. The latter is *MovieLens 1M* dataset, which collects users' ratings in the movie domain: it contains 1,000,209 ratings ranging from 1 to 5, 6,040 users, and 3,706 items. We have filtered out users

**Table 1.** Characteristics of the datasets used in the offline experiments:  $|\mathcal{U}|$  is the number of users,  $|\mathcal{I}|$  the number of items,  $X^+$  the amount of positive feedback.

Dataset	$ \mathcal{U} $	$ \mathcal{I} $	$X^+$	$\frac{X^+}{ \mathcal{U} }$	$\frac{X^+}{ \mathcal{I} }$	$\frac{X^+}{ \mathcal{I} \cdot  \mathcal{U} }$ %
Amazon DM	1,835	41,488	75,932	41.38	1.83	0.000997%
LibraryThing	7,279	37,232	749,401	102.95	20.13	0.002765%
MovieLens 1M	6,040	3,706	1,000,209	165.60	269.89	0.044684%

with less than 20 ratings (considering them as cold-users). Table 1 shows the characteristics of the resulting datasets adopted in the experiments.

**Baseline Algorithms.** We compared FedeRank with representative centralized algorithms to position its performance with respect to the state-of-the-art techniques: **VAE** [30], a non-linear probabilistic model taking advantage of Bayesian inference to estimate the model parameters; **User-kNN** and **Item-kNN** [25], two neighbor-based CF algorithms, that exploit cosine similarity to compute similarity between users or items; **BPR-MF** [38], the centralized vanilla BPR-MF implementation; and **FCF** [3], the only federated recommendation approach, to date, based on MF<sup>1</sup>. We have evaluated FedeRank considering  $|\mathcal{U}^-| = 1$ . That is, in each round of communication we involve only a single client to avoid noisy results. We thereby guarantee the sequential training, needed to compare against centralized per-wise techniques. We have investigated with two different FedeRank settings. In the **first setting**, we have set  $T = 1$ , i.e., each client extracts solely one triple  $(u, i, j)$  from its dataset when asked for training the model; with this special condition, we test whether FedeRank is effectively comparable to BPR. Moreover, to make the comparison completely fair, we extract triples as proposed by Rendle *et al.* [38]. The **second setting** follows a real Federated scenario where the client local computation is not limited to a single triple. Specifically, the number  $T$  of triples extracted by each client is set to  $\frac{X^+}{|\mathcal{U}|}$ .

**Reproducibility and Evaluation Metrics.** To train FedeRank, we have adopted a realistic temporal hold-out 80–20 splitting for training set and test set [17]. We have further split the training set adopting a temporal hold-out strategy on a user basis to pick the last 20% of interactions as a validation set. Hence, we have explored a grid in the range  $\{0.005, \dots, 0.5\}$ . Then, to ensure a fair comparison, we have used the same learning rate to train FedeRank. We have set up the remaining parameters as follows: the user- and positive item-regularization parameter is set to 1/20 of the learning rate; conversely, the negative item-regularization parameter is set to 1/200 of the learning rate as suggested by Anelli *et al.* [4]. Moreover, for each setting, we have selected the best model in the first 20 epochs. Finally, the number of latent factors is equal to 20. This value reflects a trade-off between latent factors’ expressiveness and memory space limits (given by a realistic Federated Learning environment). We have measured

<sup>1</sup> Since no source code is available, we reimplemented it in the reader’s interest.



the recommendation accuracy by exploiting: Precision ( $P@N$ ) (the proportion of relevant items in the recommendation list), and Recall ( $R@N$ ), that measures the relevant suggested items. Regarding diversity, we have adopted Item Coverage ( $IC$ ) and Gini Index ( $G$ ). The former provides the overall number of diverse recommended items, and it highlights the degree of personalization expressed by the model [1]. The latter measures how unequally an RS provides users with different items [10], being higher values corresponding to more tailored lists.

#### 4.1 Performance of Federated Learning to Rank

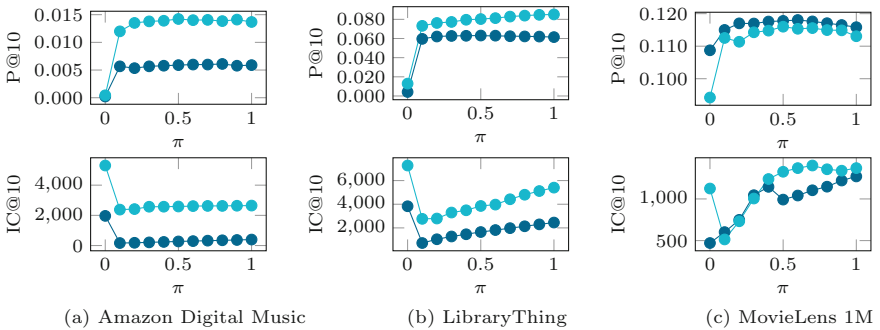
We begin our experimental evaluation by investigating the efficacy of FedeRank, and we assess whether its performance is comparable to baseline algorithms. Table 2 depicts the results in terms of accuracy and diversity. The Table is visually split into two parts. The algorithms in the bottom part (BPR-MF, FCF, and the two settings of FedeRank) are the factorization-based models. The upper part provides the positioning of FedeRank to the other state-of-the-art approaches. Focusing on the factorization-based methods, we can note that BPR-MF outperforms FedeRank for  $T = 1$ , but it remains at about 67% and 88% of the centralized algorithm for *Amazon Digital Music* and *LibraryThing*, respectively. However, the realistic Federated setting is with  $T = X^+ / |\mathcal{U}|$ . Here, FedeRank consistently improves the recommendation performance with respect to BPR-MF and FCF, over the three datasets. Actually, for *Amazon Digital Music* and *LibraryThing* FedeRank improves accuracy metrics of about 50% and 25% with respect to BPR-MF. The achievement can be explained as an advantage brought by the increased local computation. It is worth noticing that these results partially contradict Rendle *et al.* [38] since they hypothesize that traversing user-wise the training triples would worsen the recommendation performance. The same accuracy improvements are not visible in *MovieLens 1M*, where we witness results comparable or worse than BPR-MF, probably due to the overfitting caused by the very high ratio between ratings and items. FedeRank with increased computation still results robust with respect to the  $IC$  metric, since, in general, it outperforms or remains comparable to FCF and BPR-MF.

**Table 2.** Recommendation performance for baselines and FedeRank on the three datasets. For each value of  $T$ , the experiment with the best  $\pi$  is shown.

	Amazon Digital Music				LibraryThing				MovieLens 1M			
	P@10	R@10	IC@10	G@10	P@10	R@10	IC@10	G@10	P@10	R@10	IC@10	G@10
Random	0.00005	0.00005	14186	0.28069	0.00054	0.00028	31918	0.60964	0.00871	0.00283	3666	0.85426
Most popular	0.00469	0.00603	24	0.00023	0.05013	0.03044	36	0.00031	0.10224	0.03924	118	0.00569
User-kNN	0.01940	0.02757	4809	0.04115	0.14193	0.10115	3833	0.01485	0.12613	0.06701	737	0.04636
Item-kNN	0.02147	0.03171	4516	0.03801	0.20214	0.14778	12737	0.09979	0.08873	0.05475	2134	0.19292
VAE	0.01580	0.02289	3919	0.04179	0.10834	0.07711	7800	0.04638	0.11735	0.06192	1476	0.09259
BPR-MF	0.00921	0.01298	739	0.00415	0.07009	0.04303	3082	0.01359	<b>0.11911</b>	0.05817	<b>1444</b>	<b>0.08508</b>
FCF	0.00839	0.01222	<b>2655</b>	0.01861	<b>0.10760</b>	0.04392	829	0.01305	0.10760	0.04392	829	0.01305
FedeRank												
$T = 1$	0.00610	0.00889	349	0.00136	0.06309	0.03738	1650	0.00512	0.11805	<b>0.05902</b>	1041	0.06608
$T = X^+ /  \mathcal{U} $	<b>0.01422</b>	<b>0.02060</b>	2586	<b>0.02153</b>	0.08512	<b>0.05627</b>	<b>5404</b>	<b>0.02784</b>	0.11599	0.05571	1326	0.02513

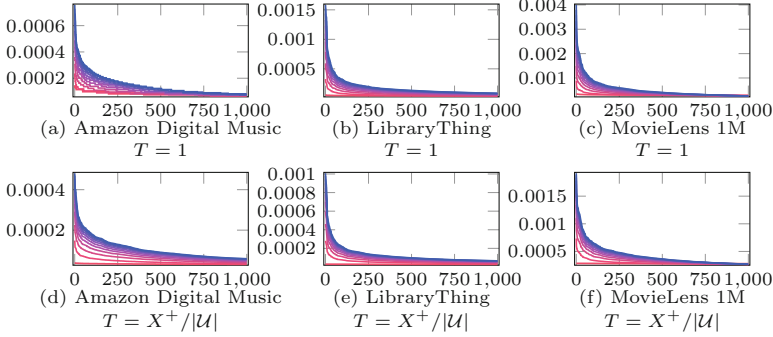
## 4.2 Analysis of Positive Feedback Transmission Ratio

We have extensively analyzed the behavior of FedeRank when tuning  $\pi$  for sending progressive fractions of positive feedback in  $[0.0, \dots, 1.0]$  with step 0.1. We believe that the most important dimensions for this analysis are accuracy (Precision), and aggregate diversity (Item Coverage). Figure 1 reports the results for the two experimented settings. Even here, *Amazon Digital Music* and *LibraryThing* show similar trends. The accuracy of the recommendation progressively increases reaching the maximum with fractions 0.8 and 0.5, respectively, for  $T = 1$ , and with fractions 0.9 and 1.0 for  $T = X^+ / |\mathcal{U}|$ . First, this suggests that, at the beginning of the training, some positive feedback is needed for establishing the value of an item. Notwithstanding, even with  $\pi = 0.1$  (i.e., sharing just 10% of private information), we witness a jump in recommendation accuracy (one order of magnitude), reaching up to 92% of the best accuracy. We should also observe another significant behavior. With a fraction of 0.0, we observe a high value of  $IC$ , with poor recommendation accuracy. It suggests that the system could not capture population preferences, and it behaves similarly to Random. However, even with a small fraction of positive feedback like 0.1, we observe a significant decrease in diversity metrics. The system learns which items are popular and starts suggesting them. Moreover, if we observe large fractions, we may notice that diversity increases as we feed the system with more information. For *MovieLens 1M*, it is worth noticing that FedeRank shows accuracy performance extremely close to the best value by sharing only 10% of positive interactions. This behavior may be due to several reasons. Firstly, *MovieLens 1M* is a relatively dense dataset in the recommendation scenario (it has a sparsity of 0.955). Secondly, it shows a very high user-item ratio [2] (i.e., 1.63) compared to *Amazon Digital Music* (0.04), and *LibraryThing* (0.20), and it shows high values for the average number of ratings per user (132.87), and ratings per item (216, 56). All these clues suggest that the system learns how to rank items even without the need for the totality of ratings. If we focus on diversity metrics,  $IC$  and Gini, we may notice that diversity is progressively increasing from fraction 0.1 to 1.0.

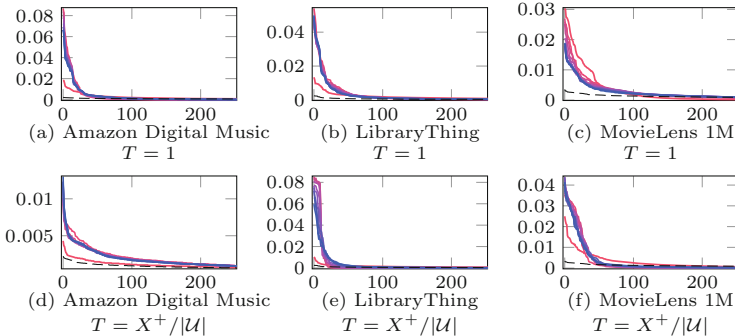


**Fig. 1.** F1 performance at different values of  $\pi$  in the range  $[0.1, 1]$ . Dark blue is  $T = 1$ , light blue is  $T = X^+ / |\mathcal{U}|$ . (Color figure online)

It suggests that the system recommends a small number of popular items with a fraction of 0.1, while it provides more diversified recommendation lists considering larger portions of positive user feedback. At this stage of the analysis, we can draw an interesting consideration: in general, the highest accuracy values do not correspond to the fraction of 1.0. Specifically, the experiments show that, initially, the recommender struggles to suggest relevant items without positive feedback (fraction 0.0). However, with a small injection of feedback, the system starts to work well. Nonetheless, in *Amazon Digital Music* and *LibraryThing*, if we increase the fraction, we witness an increase concerning accuracy only until a certain value of  $\pi$ . Although this consideration, we observe an increase in diversity metrics when we continue to increase the value of  $\pi$ . Since it has a small or even detrimental impact on accuracy, those items might be unpopular items erroneously suggested to users.



**Fig. 2.** Normalized number of item updates during the training: the 1,000 most updated items for different values of  $\pi$  (from  $\pi = 0.0$  in red to  $\pi = 1.0$  in blue). (Color figure online)



**Fig. 3.** Normalized number of recommendations for each item (colored curves from  $\pi = 0.0$  in red to  $\pi = 1.0$  in blue) vs. normalized amount of positive feedback per item (black dashed curve). The 250 most popular items are shown. (Color figure online)

### 4.3 Study on FedeRank Algorithmic Bias

In this section, we study how incomplete transmission of user feedback affects the item popularity in the recommendations and during the learning process. It is essential to discover whether the exploitation of a FL approach influences the algorithmic bias, determining popular items to be over-represented [5,9]. We have re-trained FedeRank with all the previously considered  $\pi$ . For each experiment, we analyzed the data flow between the clients and the server. Afterward, we have extracted the number of updates for each item. Figure 2 illustrates the occurrences for the 1,000 most updated items. In the Figure, the curve colors denote the different  $\pi$ , while the values represent the update frequency during the training process for each item on the horizontal axis. Analogously, we considered the final top-10 recommendation list of each user. Following the same strategy, we analyzed the occurrences of the items in the recommendation. Then, we ordered items from the most to the least recommended, and we plotted the occurrences of the first 250 in Fig. 3. To compare the different datasets, we have normalized the values considering the overall dataset occurrences. Figure 2 shows that data disclosure, i.e., the value of  $\pi$ , highly influences the information exchanged during the training process. Additionally, the update frequency curve exhibits a constant behavior for all the datasets, when  $\pi = 0.0$ . This trend suggests that items are randomly updated without taking into account any information about item popularity. This behavior explains the high *IC* entirely observed in Fig. 1 for  $\pi = 0.0$ . The curve for  $\pi = 0.1$  shows that the exchanged data is enough to provide the system with information about item popularity. The curves suggest that the information on item popularity is being injected into the system. By increasing the value of  $\pi$ , the trend becomes more evident. Due to the original rating distribution, the system initially exchanges more information about the very popular items. To analyze the algorithmic bias, we can observe Fig. 3, where the colored curves represent the frequency of item recommendation on the horizontal axis, and the black dashed curve the amount of positive feedback for that item in the dataset. Remarkably, item popularity in recommendation lists does not vary as we may expect based on the previous analysis. The setting  $\pi = 0.0$  is an exception, as extensively explained before. Since in *Amazon Digital Music* and *LibraryThing* the updates sent by the clients are randomly selected among the negative items, FedeRank acts like a Random recommender. The system cannot catch popularity information and it struggles to make the right items popular. Finally, we can focus on the curves for  $\pi > 0$ . It is noteworthy that the  $\pi$  curves behave similarly, and they propose the same proportion of popular items. The curves show the model absorbs the initial variation in exchanged item distribution, unveiling an unknown aspect of factorization models.

## 5 Conclusion and Future Work

In this paper, we have tackled the problem of putting users in control of their private data for a recommendation scenario. Witnessing the growing concern about privacy, users might want to exploit their sensitive data and share only a

small fraction of it. In such a context, classic CF approaches are no more feasible. To overcome these problems, we have proposed FedeRank, a novel recommendation framework that respects the FL paradigm. With FedeRank, private user feedback remains on user devices unless they decide to share it. Nevertheless, FedeRank ensures high-quality recommendations despite the constrained setting. We have extensively studied the performance of FedeRank by comparing it with other state-of-the-art methods. We have then analyzed the impact of progressive reduction of user feedback and studied the effects on the diversity of the recommendation results. Finally, we have investigated whether the federated algorithm imposes an algorithmic bias to the recommended lists. The study paves the way for further research directions. On the one hand, the results' analysis suggests that centralized recommender systems are not performing at their best. Feeding recommender systems with all the available feedback, without any filtering, may lead to a performance worsening. On the other hand, the competitive results of FedeRank suggest that the FL-based algorithms show a recommendation quality that makes them suitable to be adopted on a massive scale.

## References

1. Adomavicius, G., Kwon, Y.: Improving aggregate recommendation diversity using ranking-based techniques. *IEEE Trans. Knowl. Data Eng.* **24**(5), 896–911 (2012)
2. Adomavicius, G., Zhang, J.: Impact of data characteristics on recommender systems performance. *ACM Trans. Manage. Inf. Syst.* **3**(1), 31–317 (2012)
3. Ammad-ud-din, M., et al.: Federated collaborative filtering for privacy-preserving personalized recommendation system, CoRR, abs/1901.09888 (2019)
4. Anelli, V.W., Noia, T.D., Sciascio, E.D., Pomo, C., Ragone, A.: On the discriminative power of hyper-parameters in cross-validation and how to choose them. In: *Proceedings of the 13th ACM Conference on Recommender Systems*, pp. 447–451. ACM (2019)
5. Baeza-Yates, R.: Bias in search and recommender systems. In: *14th ACM Conference on Recommender Systems, RecSys 2020, Virtual Event, Brazil, 22–26 September 2020*, p. 2 (2020)
6. Bonawitz, K., et al.: Towards federated learning at scale: system design, CoRR, abs/1902.01046 (2019)
7. Bonawitz, K., et al.: Practical secure aggregation for privacy-preserving machine learning. In: *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security, CCS 2017, Dallas, TX, USA, 30 October–03 November 2017*, pp. 1175–1191 (2017)
8. California State Legislature: The California Consumer Privacy Act of 2018 (2018)
9. Cañamares, R., Castells, P.: Should I follow the crowd?: a probabilistic analysis of the effectiveness of popularity in recommender systems. In: *2018 41st International ACM SIGIR Conference on Research & Development in Information Retrieval, SIGIR 2018, Ann Arbor, MI, USA, July 08–12*, pp. 415–424 (2018)
10. Castells, P., Hurley, N.J., Vargas, S.: Novelty and diversity in recommender systems. In: Ricci, F., Rokach, L., Shapira, B. (eds.) *Recommender Systems Handbook*, pp. 881–918. Springer, Boston (2015)
11. Chai, D., Wang, L., Chen, K., Yang, Q.: Secure federated matrix factorization. *IEEE Intell. Syst.*, 1 (2020)

12. Duriakova, E., et al.: PDMFRec: a decentralised matrix factorisation with tunable user-centric privacy. In: Proceedings of the 13th ACM Conference on Recommender Systems, RecSys 2019, Copenhagen, Denmark, 16–20 September 2019, pp. 457–461 (2019)
13. Dwork, C., Roth, A.: The algorithmic foundations of differential privacy. *Found. Trends Theor. Comput. Sci.* **9**(3–4), 211–407 (2014)
14. Erlingsson, Ú., Pihur, V., Korolova, A.: RAPPOR: randomized aggregatable privacy-preserving ordinal response. In: Proceedings of the 2014 ACM SIGSAC Conference on Computer and Communications Security, Scottsdale, AZ, USA, 3–7 November 2014, pp. 1054–1067 (2014)
15. European Commission. 2018 reform of EU data protection rules (2018)
16. Fierimonte, R., Scardapane, S., Uncini, A., Panella, M.: Fully decentralized semi-supervised learning via privacy-preserving matrix completion. *IEEE Trans. Neural Netw. Learn. Syst.* **28**(11), 2699–2711 (2017)
17. Gunawardana, A., Shani, G.: Evaluating recommender systems. In: Ricci, F., Rokach, L., Shapira, B. (eds.) *Recommender Systems Handbook*, pp. 265–308. Springer, Boston (2015)
18. Harper, F.M., Konstan, J.A.: The MovieLens datasets: history and context. *ACM Trans. Interact. Intell. Syst. (TIIS)* **5**(4), 1–19 (2015)
19. Hu, Y., Koren, Y., Volinsky, C.: Collaborative filtering for implicit feedback datasets. In: Proceedings of the 8th IEEE International Conference on Data Mining, ICDM 2008, 15–19 December 2008, Pisa, Italy, pp. 263–272. IEEE Computer Society (2008)
20. Jeckmans, A.J.P., Beye, M., Erkin, Z., Hartel, P.H., Lagendijk, R.L., Tang, Q.: Privacy in recommender systems. In: Ramzan, N., van Zwol, R., Lee, J., Clüver, K., Hua, X. (eds.) *Social Media Retrieval, Computer Communications and Networks*, pp. 263–281. Springer, London (2013)
21. Kairouz, P.: *Advances and Open Problems in Federated Learning* (2019)
22. Kharitonov, E.: Federated online learning to rank with evolution strategies. In: Proceedings of the 12th ACM International Conference on Web Search and Data Mining, pp. 249–257 (2019)
23. Konečný, J., McMahan, H.B., Ramage, D., Richtárik, P.: Federated optimization: distributed machine learning for on-device intelligence, CoRR, abs/1610.02527 (2016)
24. Koren, Y.: Factorization meets the neighborhood: a multifaceted collaborative filtering model. In: Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Las Vegas, Nevada, USA, 24–27 August 2008, pp. 426–434 (2008)
25. Koren, Y.: Factor in the neighbors: scalable and accurate collaborative filtering. *ACM Trans. Knowl. Discov. Data (TKDD)* **4**(1), 1–24 (2010)
26. Koren, Y., Bell, R.M.: Advances in collaborative filtering. In: Ricci, F., Rokach, L., Shapira, B. (eds.) *Recommender Systems Handbook*, pp. 77–118. Springer, Heidelberg (2015)
27. Koren, Y., Bell, R.M., Volinsky, C.: Matrix factorization techniques for recommender systems. *IEEE Comput.* **42**(8), 30–37 (2009)
28. Koren, Y., Sill, J.: OrdRec: an ordinal model for predicting personalized item rating distributions. In: Mobasher, B., Burke, R.D., Jannach, D., Adomavicius, G. (eds.) Proceedings of the 2011 ACM Conference on Recommender Systems, RecSys 2011, Chicago, IL, USA, 23–27 October 2011, pp. 117–124. ACM (2011)
29. Kumar Bokde, D., Girase, S., Mukhopadhyay, D.: Role of matrix factorization model in collaborative filtering algorithm: a survey, CoRR, abs/1503.07475 (2015)

30. Liang, D., Krishnan, R.G., Hoffman, M.D., Jebara, T.: Variational autoencoders for collaborative filtering. In: Proceedings of the 2018 World Wide Web Conference, pp. 689–698 (2018)
31. McAuley, J., Targett, C., Shi, Q., Van Den Hengel, A.: Image-based recommendations on styles and substitutes. In: Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 43–52 (2015)
32. McFee, B., Barrington, L., Lanckriet, G.R.G.: Learning content similarity for music recommendation. *IEEE Trans. Audio Speech Lang. Process.* **20**(8), 2207–2218 (2012)
33. McMahan, H.B., Moore, E., Ramage, D., Hampson, S., et al.: Communication-efficient learning of deep networks from decentralized data. arXiv preprint [arXiv:1602.05629](https://arxiv.org/abs/1602.05629) (2016)
34. McNee, S.M., Riedl, J., Konstan, J.A.: Being accurate is not enough: how accuracy metrics have hurt recommender systems. In: CHI’06 Extended Abstracts on Human Factors in Computing Systems, pp. 1097–1101 (2006)
35. Muhammad, K., et al.: FedFast: going beyond average for faster training of federated recommender systems. In: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD 2020, Virtual Event, CA, USA, 23–27 August 2020, pp. 1234–1242 (2020)
36. Ning, X., Desrosiers, C., Karypis, G.: A comprehensive survey of neighborhood-based recommendation methods. In: Ricci, F., Rokach, L., Shapira, B. (eds.) *Recommender Systems Handbook*, pp. 37–76. Springer, Boston (2015). [https://doi.org/10.1007/978-1-4899-7637-6\\_2](https://doi.org/10.1007/978-1-4899-7637-6_2)
37. Rendle, S.: Factorization machines. In: The 10th IEEE International Conference on Data Mining, ICDM 2010, Sydney, Australia, 14–17 December 2010, pp. 995–1000 (2010)
38. Rendle, S., Freudenthaler, C., Gantner, Z., Schmidt-Thieme, L.: BPR: Bayesian personalized ranking from implicit feedback. In: Bilmes, J.A., Ng, A.Y. (eds.) *Proceedings of the 25th Conference on Uncertainty in Artificial Intelligence, UAI 2009*, Montreal, QC, Canada, 18–21 June 2009, pp. 452–461. AUAI Press (2009)
39. Rendle, S., Freudenthaler, C., Schmidt-Thieme, L.: Factorizing personalized markov chains for next-basket recommendation. In: Proceedings of the 19th International Conference on World Wide Web, WWW 2010, Raleigh, North Carolina, USA, 26–30 April 2010, pp. 811–820 (2010)
40. Rendle, S., Schmidt-Thieme, L.: Pairwise interaction tensor factorization for personalized tag recommendation. In: Proceedings of the 3rd International Conference on Web Search and Web Data Mining, WSDM 2010, New York, NY, USA, 4–6 February 2010, pp. 81–90 (2010)
41. Shi, Y., Larson, M., Hanjalic, A.: List-wise learning to rank with matrix factorization for collaborative filtering. In: Proceedings of the 4th ACM Conference on Recommender Systems, pp. 269–272 (2010)
42. Standing Committee of the National People’s Congress of Popular Republic of China. China internet security law (2017)
43. Differential Privacy Team: Learning with Privacy at Scale (2017). <https://machinelearning.apple.com/research/learning-with-privacy-at-scale>. Accessed Jan 2021

44. Wang, Y., Wu, X., Hu, D.: Using randomized response for differential privacy preserving data collection. In: Palpanas, T., Stefanidis, K. (eds.) Proceedings of the Workshops of the EDBT/ICDT 2016 Joint Conference, EDBT/ICDT Workshops 2016, Bordeaux, France, 15 March 2016, vol. 1558 of CEUR Workshop Proceedings. CEUR-WS.org (2016)
45. Yang, Q., Liu, Y., Chen, T., Tong, Y.: Federated machine learning: concept and applications. *ACM TIST* **10**(2), 12:1–12:19 (2019)
46. Yang, Q., Liu, Y., Cheng, Y., Kang, Y., Chen, T., Yu, H.: *Federated Learning*. Morgan & Claypool Publishers, San Rafael (2019)
47. Yuan, J., Shalaby, W., Korayem, M., Lin, D., AlJadda, K., Luo, J.: Solving cold-start problem in large-scale recommendation engines: a deep learning approach. In: 2016 IEEE International Conference on Big Data, BigData 2016, Washington DC, USA, 5–8 December 2016, pp. 1901–1910. IEEE Computer Society (2016)
48. Zhao, T., McAuley, J., King, I.: Improving latent factor models via personalized feature projection for one class recommendation. In: Proceedings of the 24th ACM International Conference on Information and Knowledge Management, pp. 821–830 (2015)