# PDSR: A Privacy-Preserving Diversified Service Recommendation Method on Distributed Data

Lina Wang, Huan Yang, Yiran Shen, Chao Liu, Lianyong Qi, Xiuzhen Cheng, and Feng Li

**Abstract**—The last decade has witnessed a tremendous growth of service computing, while efficient service recommendation methods are desired to recommend high-quality services to users. It is well known that collaborative filtering is one of the most popular methods for service recommendation based on QoS, and many existing proposals focus on improving recommendation accuracy, i.e., recommending high-quality redundant services. Nevertheless, users may have different requirements on QoS, and hence diversified recommendation has been attracting increasing attention in recent years to fulfill users' diverse demands and to explore potential services. Unfortunately, the recommendation performances relies on a large volume of data (e.g., QoS data), whereas the data may be distributed across multiple platforms. Therefore, to enable data sharing across the different platforms for diversified service recommendation, we propose a *Privacy-preserving Diversified Service Recommendation* (PDSR) method. Specifically, we innovate in leveraging the Locality-Sensitive Hashing (LSH) mechanism such that privacy-preserved data sharing across different platforms is enabled to construct a service similarity graph. Based on the similarity graph, we propose a novel accuracy-diversity metric and design a $2$-approximation algorithm to select $K$ services to recommend by maximizing the accuracy-diversity measure. Extensive experiments on real datasets are conducted to verify the efficacy of our PDSR method.

**Index Terms**—Collaborative filtering, recommendation diversity, privacy preservation.

✦

## 1 INTRODUCTION

Recent years have witnessed a considerable development of services, thanks to the emerging computing paradigms and architectures, e.g., *Software Oriented Architecture* (SOA), Internet of Services (IoS) and Cloud Computing. As a growing variety of services are provided on the Web, it is of considerable importance to recommend appropriate services for a wide range of users [1].

It is well known that *Collaborative Filtering* (CF) has been widely used to build service recommender systems based on *Quality of Service* (QoS) [2], [3], [4]. The basic idea is to exploit the similarity among services to predict QoS values for particular users, while characterizing the service similarity relies on mining a big volume of historical QoS data. Traditional recommendation methods usually aim at improving recommendation accuracy such that they recommend only similar services with the highest predicted QoS to target users, whereas the users may have different requirements on the QoS [5], [6]. Furthermore, the recommendation highly relies on the QoS prediction; nevertheless, the prediction

usually has inevitable bias. Therefore, recommendation diversity is extensively studied in many recent proposals [7], [8], [9], [10], [11]. On one hand, it is able to fulfill the users' different QoS demands; on the other hand, it helps to explore potential services which is underestimated by the QoS prediction.

Unfortunately, diversifying service recommendation usually entails a large volume of data (e.g., QoS data), whereas the data may be distributed across different platforms [12]. Despite the benefits of enabling data sharing across the different platforms, the platforms owing the data are usually reluctant to share their data with each other, since private information could be inferred from the data by untrusted third-party platforms [13], [14]. Therefore, to enabling data sharing across different platforms, privacy preservation is one of the main concerns. Many efforts have been made to integrate privacy preservation techniques with CF. For example, [15], [16] apply homomorphic encryption techniques, and differential privacy mechanisms are adopted in [17], [18] to obfuscate the original data before sharing them with others. Nevertheless, the cryptography-based method entails high computation overhead, while the differential privacy mechanism usually impairs data utility. Another choice is to utilize hashing mechanisms. For example, *Locality-Sensitive Hashing* (LSH), have been utilized in [14], [19], [20] to ensure privacy preservation based on distributed data sources. It is usually applied to efficiently exploit the similarity among the services or users through light-weight computations, while the privacy preservation can be ensured thanks to its irreversibility. Unfortunately, how to effectively diversify the similarity relationship obtained by LSH is still an open problem.

In this paper, we propose a *Privacy-preserving Diversified Service Recommendation* (PDSR) method. The proposed PDSR

- *L. Wang, X. Cheng, and F. Li are with School of Computer Science and Technology, Shandong University, Qingdao, 266237, China.*
  *E-mail: linawang425@mail.sdu.edu.cn, {xzcheng, fli}@sdu.edu.cn.*
- *H. Yang is with College of Computer Science and Technology, Qingdao University, Qingdao, 266071, China.*
  *Email: cathy_huanyang@hotmail.com.*
- *Y. Shen is with School of Software, Shandong University, Jinan, 250101, China.*
  *E-mail: yiran.shen@sdu.edu.cn.*
- *C. Liu with Department of Computer Science and Technology, Ocean University of China, Qingdao, 266100, China.*
  *E-mail: liuchao@ouc.edu.cn*
- *L. Qi is with College of Computer Science and Technology, China University of Petroleum (East China), Qingdao, 266580, China.*
  *E-mail: 20220115@upc.edu.cn.*

method not only leverages a trade-off between recommendation accuracy and diversity, but also ensures privacy preservation for data sharing across different platforms. In particular, LSH mechanism is adopted to construct a graph to characterize the similarity among the services based on distributed data from different sources. Due to the irreversibility and the computational efficiency of the hash functions, the construction of the similarity graph entails light-weight computational overhead and enables privacy preservation across the different data sources. We innovate in designing a new diversity metric for the similarity graph, based on which, we formulate our $K$-*Diversified Service Recommendation* ($K$-DSR) problem where $K$ services are selected to maximize the accuracy-diversity evaluation of the recommendation. We prove the NP-hardness of the problem, and design a 2-approximation algorithm thanks to its monotone submodularity.

In summary, we make the following main contributions in our PDSR method as follows:

- We leverage the notion of LSH to efficiently construct a weighted service similarity graph across different data sources with their privacy preserved.
- We innovate in designing a new diversity metric for the similarity graph, based on which, a 2-approximation algorithm is proposed for diversified recommendation.
- We finally perform extensive experiments on real data to verify the efficacy of our proposed method in leveraging the accuracy-diversity trade-off.

The remaining of our paper is organized as follows. We first survey related literature and present an example to motivate the design of our service recommendation method in Sec. 2. We then give some preliminaries in Sec. 3. Moreover, we report the details of our proposed PDSR method in Sec. 4. Experiment results are then given in Sec. 5. We finally conclude this paper in Sec. 6.

## 2 RELATED WORK AND MOTIVATION

### 2.1 Diversified Service Recommendation

Some of state-of-the-art proposals mainly focus on developing innovative methods to enhance the result accuracy [4], [21], [22], [23], [24]. For example, in [22], a covering-based service recommendation method is proposed through neighbor-aware matrix factorization. [23] designs a deep learning-aided collaborative filtering framework, where neural networks is adopted to learn binary representations for both users and items. In [24], time correlation coefficient and $K$-means clustering with cuckoo search are employed to improve the performance of collaborative filtering by taking into account users' time-varying interests.

As mentioned in Sec. 1, an accurate recommendation is not necessarily a satisfactory one. Therefore, diversified recommendation have been attracting increasing attention in recent years. In [8], functional relevance, QoS utility, and diversity features of Web services are incorporated for recommending well diversified services to users. In [25], an acceleration algorithm is proposed to speed up the *Maximum A Posteriori* (MAP) inference for determinantal point process; the algorithm can be used for recommendation

diversification. The accelerated MAP inference algorithm is then adopted in [10]. In particular, [10], the fast MAP inference algorithm is utilized in a combinatorial bandit for diversified recommendation. [5] exploits the quality correlations among different dimensions of users' quality preferences, to diversify recommendation results. In [26], user exposure diversity and item concordance are exploited to endow the collaborative filtering with recommendation diversity. In [11], multiple recommendation lists are selected first, from which the most diversified one can be found according to a list diversity measure. In [7], rebalanced neighbor discovering, category boosted negative sampling and adversarial learning are performed on top of graph convolutional networks, so as to realize diversity in recommendation. In [9], an end-to-end dynamic diversified graph framework is proposed to construct the user-item graph dynamically based on the users and items embeddings, and a quantile progressive candidate selection operator is designed to efficiently select diverse items to recommend.

### 2.2 Privacy-Preserved Service Recommendation

The above studies rely on a huge volume of data, by exploiting which recommendation results are diversified. However, data are usually stored on different platforms, while data sharing across the different sources may considerably increases the risk of privacy leakage. Hence, one of another main concerns in service recommendation is the privacy issue. In [27], only a small fraction of data are shared publicly to reduce the disclosure of sensitive data. By this way, although privacy is preserved, but it is with a significant sacrifice in recommendation performance. Different from [27], [13] applies data obfuscation for the purpose of privacy preservation. Specifically, data are obfuscated through randomization techniques before being shared for recommendation. Likewise, the notion of differential privacy is leveraged in [28]; the data privacy is preserved through insert "noise" into the data before sharing them. However, in the above two methods, data are obfuscated or noised such that their utility is impaired. Additionally, when the data are stored in different platforms, sharing the obfuscated (or noised) data across different platforms still results in considerable communication overhead.

Another choice for privacy preservation is to utilize *Locality-Sensitive Hashing* (LSH). In [14], the similarity retention is LSH is employed such that different platforms can share the similarity among the services (which can be indicated by the hash values) with each other without disclosing the original data. Furthermore, the LSH-based method is developed by [29] by exploring mutiple dimensions of QoS data. Although these proposals have made considerable effort in combining LSH and CL for privacy-preserved recommendation, they concentrate in pursuing high recommendation accuracy and do not take into account recommendation diversity. Although the LSH mechanism is used to calculate both similarity and diversity in [30], this proposal does not exploit the application of the LSH mechanism in privacy-preservation across different data sources. Moreover, [30] calculates the diversity directly according to the LSH values, while we exploit the service similarity graph to define the recommendation diversity and gain better recommendation results, as will be shown in Sec. 4.

## 2.3 Motivation

We hereby give an example to explain the application scenario of our service recommendation method. As shown in Fig. 1, we assume that there are two platforms (or data sources), i.e., Netflix and IBM. We also suppose that there are three users, i.e., Tom, Bob and Jack, on Netflix platform and another two users John and Alice on IBM platform. Let $s_1, s_2, \cdots, s_M$ denote $M$ different services invoked by these users, respectively. Note that, a service may be invoked by multiple users on different platforms. For example, $s_1$ is invoked by Tom and Bob on Netflix platform and by John and Alice on IBM platform, and $s_M$ is invoked by Bob and Jack on Netflix platform and by John and Alice on IBM platform. The corresponding QoS data are then stored on different platforms. We suppose Netflix intends to recommend new services to Tom. Through the CF method, Netflix first computes the similarity (e.g., by constructing a similarity graph) among the services according to its local QoS data and then seeks for the services that are similar to the historical records of Tom. Specifically, for any new service for Tom, its QoS value can be predicted according to the ones of the other similar services. Usually, $K$ services with the highest predicted QoS values are finally selected and recommended to Tom.

As mentioned above, there are two challenging issues in the above recommendation process. On one hand, similar services are recommended to Tom in order to ensure recommendation accuracy, while the redundancy of the services may impair the experience of the recommendation for Tom. On the other hand, although data sharing between Netflix and IBM will be helpful for data-driven recommendation diversification, enabling the data sharing across different platforms efficiently is highly non-trivial especially due to privacy concerns. Therefore, to tackle the above challenges, we first leverage the notion of LSH to enable efficient data sharing between Netflix and IBM, based on which, we can characterize the similarity among the services more comprehensively. Furthermore, we propose a new diversity metric such that Netflix can mine diversified recommendation for Tom through the cumulative data more effectively.
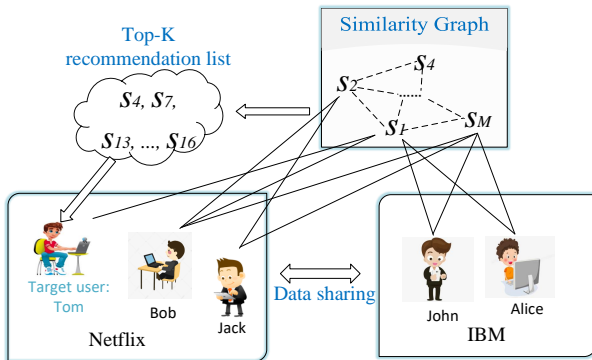


Fig. 1. Service recommendation across different data sources.

## 3 PRELIMINARIES

Before diving into the design and analysis of our proposed algorithm, we first introduce some preliminaries in this section, which will be very helpful later.

### 3.1 Locality-Sensitive Hashing

LSH is an efficient approximate nearest neighbor search method [31]. Given a set of points (e.g., services which can be represented as a high-dimensional vectors or points in our case), the goal of traditional hashing methods is to map the points to a set of values such that the points can be "spread" randomly, while LSH aims at mapping the points with the so-called "locality" guaranteed. Specifically, by LSH, the points which are "closed" to each other in the original data space can be hashed to the same value with a high probability, whereas there is a high probability for the ones "far away" from each other to be hashed to different values. In another word, it is more likely for "similar" points to have the same hash value.

Let $\mathsf{Dis}(s, s')$ denote the "distance" between the points $s$ and $s'$. The distance between $s$ and $s'$ is smaller when $s$ and $s'$ are more similar to each other. We formally define LSH in the following **Definition** 1.

**Definition 1** (Locality Sensitive Hashing). *Given a set of points $\mathcal{S}$, let $\mathcal{H} = \{h \mid \mathcal{S} \to \mathbb{Z}^+\}$ be a family of hash functions which maps the points $\mathcal{S}$ to a set of hash values $\mathbb{Z}^+$. $\mathcal{H}$ is said to be a $(\varepsilon_1, \varepsilon_2, p_1, p_2)$-LSH if*

$$\begin{cases} \mathsf{Pr}_{\mathcal{H}}(h(s) = h(s')) \geq p_1, & \text{if } \mathsf{Dis}(s, s') \leq \varepsilon_1 \\ \mathsf{Pr}_{\mathcal{H}}(h(s) = h(s')) \leq p_2, & \text{if } \mathsf{Dis}(s, s') \geq \varepsilon_2 \end{cases} \quad (1)$$

*holds for any $s, s' \in \mathcal{S}$, where $p_1 > p_2$. The size of the resulting hash table is the number of integers in $\mathbb{Z}^+$, and we also suppose $H$ is the number of functions in $\mathcal{H}$. Note that the probability is over the random choices of $h \in \mathcal{H}$.*

It is shown by the above definition that, in the LSH mechanism, similar points are mapped to the same value with high probability; therefore, we can figure out the similarity among any points (i.e., the services we intend to recommend) according to their hash values instead of their coordinates in the original data space (i.e., the private information of the services). In another word, for any services which are mapped to the same value through LSH, they are similar to each other with high probability. Furthermore, the hash functions are irreversible such that we can observe the similarity among the points without disclosing their "private" coordinates. As will be shown in Sec. 4, we leverage the notion of LSH to construct a graph to characterize the similarity among the services by enabling privacy-preserved data sharing across different platforms.

### 3.2 Expanded Set and Expansion Ratio

Let $\mathcal{G} = (\mathcal{M}, \mathcal{E})$ denote a graph, where $\mathcal{M}$ and $\mathcal{E}$ denote the set of vertices and the set of edges, respectively. We define expanded set and expansion ratio for any subset $\mathcal{M}' \subseteq \mathcal{M}$ as follows.

**Definition 2** (Expanded set). *Given a graph $\mathcal{G} = (\mathcal{M}, \mathcal{E})$, for any subset of vertices $\mathcal{M}' \subseteq \mathcal{M}$, the expanded set of $\mathcal{M}'$ can be defined as*

$$\mathsf{Exp}(\mathcal{M}') = \mathcal{M}' \bigcup \{i \in \mathcal{M}/\mathcal{M}' \mid \exists i' \in \mathcal{M}', (i, i') \in \mathcal{E}\} \quad (2)$$

**Definition 3** (Expansion ratio). *Let $|\mathsf{Exp}(\mathcal{M}')|$ denote the number of vertices in the expanded set of $\mathcal{M}' \subseteq \mathcal{M}$. We also suppose $M = |\mathcal{M}|$ denotes the number of vertices in graph $\mathcal{G}$. The expansion ratio of $\mathcal{M}'$ can be defined by*

$$\alpha(\mathcal{M}') = \frac{|\mathsf{Exp}(\mathcal{M}')|}{M} \tag{3}$$

According to the above definition, the expansion radio of a subset $\mathcal{M}' \subseteq \mathcal{M}$ actually measures its trend of expanding outwards. If we let the size of $\mathcal{M}'$, i.e., $M' = |\mathcal{M}'|$, be fixed, a larger expansion ratio of $\mathcal{M}'$ implies that the vertices in $\mathcal{M}'$ have more neighbors outside (i.e., in $\mathcal{M}/\mathcal{M}'$), and hence it is less likely for the vertices in $\mathcal{M}'$ to share the same neighbors in $\mathcal{M}/\mathcal{M}'$. For example, as illustrated in Fig. 2, we use red vertices and blue vertices to represent the selected subset and its expanded neighbors, respectively. The subset $\mathcal{M}'_1$ (consisting of the three red vertices) in Fig. 2 (b) has an expansion ratio $\alpha(\mathcal{M}'_1) = 8/13$ and the three vertices share the same neighbor, while its counterpart $\mathcal{M}'_2$ in Fig. 2 (c) has an expansion ratio $\alpha(\mathcal{M}'_2) = 10/13$ and none of the three red vertices in $\mathcal{M}'_2$ share the same neighbor, even $\mathcal{M}'_1$ and $\mathcal{M}'_2$ have the same number of outgoing edges.
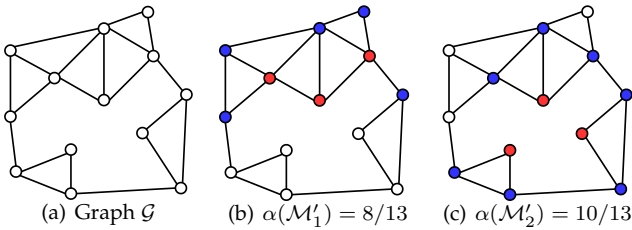


Fig. 2. An illustration of expanded set and expansion ratio.

### 3.3 Submodularity

The concept of submodular set function is very useful in approximation algorithms. In the following, we present a definition of monotone submodular set function, based on which we can design and analyze our graph-based recommendation algorithm.

**Definition 4** (Monotone submodular function). *Given a finite set $\mathcal{M}$, let $\mathsf{F} : 2^{\mathcal{M}} \to \mathbb{R}$ be a real valued set function. $\mathsf{F}(\cdot)$ is said to be a monotone submodular set function, if the following conditions meet.*

- *Monotonicity. For any two subsets $\mathcal{M}'_1, \mathcal{M}'_2 \subseteq \mathcal{M}$ such that $\mathcal{M}'_1 \subseteq \mathcal{M}'_2 \subseteq \mathcal{M}$, we have $\mathsf{F}(\mathcal{M}'_1) \leq \mathsf{F}(\mathcal{M}'_2)$.*
- *Submodularity. For any two subsets $\mathcal{M}'_1, \mathcal{M}'_2 \subseteq \mathcal{M}$ such that $\mathcal{M}'_1 \subseteq \mathcal{M}'_2 \subseteq \mathcal{M}$ and any $i \in \mathcal{M}/\mathcal{M}_2$, we have $\mathsf{F}(\mathcal{M}'_1 \bigcup \{i\}) - \mathsf{F}(\mathcal{M}'_1) \geq \mathsf{F}(\mathcal{M}'_2 \bigcup \{i\}) - \mathsf{F}(\mathcal{M}'_2)$.*

**Lemma 1.** *Given any two disjoint subset $\mathcal{M}_1 \subseteq \mathcal{M}$ and $\mathcal{M}_2 \subseteq \mathcal{M}$ such that $\mathcal{M}_1 \bigcap \mathcal{M}_2 = \emptyset$, we have*

$$\sum_{i \in \mathcal{M}_1} \left( \mathsf{F}(\{i\} \bigcup \mathcal{M}_1) - \mathsf{F}(\mathcal{M}_1) \right) \geq \mathsf{F}(\mathcal{M}_1 \bigcup \mathcal{M}_2) - \mathsf{F}(\mathcal{M}_2) \tag{4}$$

*if $\mathsf{F}(\cdot)$ is a monotone submodular set function.*

*Proof.* Supposing $\mathcal{M}_1 = \{i_1, i_2, \cdots, i_{M_1}\}$ and $\mathcal{M}_{1,\ell} = \{i_1, \cdots, i_\ell\}$, we have

$$\mathsf{F}(\mathcal{M}_1 \bigcup \mathcal{M}_2)$$
$$= \mathsf{F}(\mathcal{M}_2) + \sum_{\ell=1}^{M_1} (\mathsf{F}(\mathcal{M}_2 \bigcup \mathcal{M}_{1,\ell}) - \mathsf{F}(\mathcal{M}_2 \bigcup \mathcal{M}_{1,\ell-1})) \tag{5}$$

When $\mathsf{F}(\cdot)$ is a monotone submodular set fucntion, we have $\mathsf{F}(\mathcal{M}_2 \bigcup \{i_\ell\}) \geq \mathsf{F}(\mathcal{M}_2 \bigcup \mathcal{M}_{1,\ell})$ and $\mathsf{F}(\mathcal{M}_2 \bigcup \mathcal{M}_{1,\ell}) \geq \mathsf{F}(\mathcal{M}_2)$ for any $\ell = 1, \cdots, M_1$. The proof can be completed by substituting the above inequalities into (5). $\square$

## 4 OUR PROPOSED METHOD

In this section, we introduce the details of our PDSR method. We consider $M$ services and $N$ users across $R$ different platforms. Let $\mathcal{M} = \{1, 2, \cdots, M\}$, $\mathcal{N} = \{1, 2, \cdots, N\}$ and $\mathcal{R} = \{1, 2, \cdots, R\}$ denote the set of services, the set of users and the set of platforms, respectively. Suppose $\mathcal{N}_r \subseteq \mathcal{N}$ denotes the set of users on platform $r$, and $N_r$ is the cardinality of $\mathcal{N}_r$, i.e., the number of users on platform $r$. For each service $i \in \mathcal{M}$, it can be invoked by user $j \in \mathcal{N}_r$ through platform $r \in \mathcal{R}$. The resulting QoS data $d_{i,j,r} \geq 0$ is then observed by platform $r$ [1]. Therefore, for each service $i$, its QoS value over the $N_r$ users on platform $r$ can be represented by a $N_r$-dimensional vector $\vec{d}_{i,r} = (d_{i,j,r})_{j \in \mathcal{N}_r}$. Without loss of generality, we assume $d_{i,j,r} = 0$ if service $i$ is never invoked by user $j$ on platform $r$; otherwise, $d_{i,j,r} > 0$.

Before diving into the details of our proposed method, we first introduce its outline. As demonstrated in Sec. 4.1, each platform, e.g., $r$, first calculates a LSH vector for each of the services $i \in \mathcal{M}$ and then shares the vectors with each other. By concatenating the received vectors of any service $i \in \mathcal{M}$, platform $r$ can calculate an index $v_i$ for $\forall i$. The above procedure can be called repeatedly such that each platform $r$ can construct a similarity graph $\mathcal{G}$ where any two similar services are connected by an edge, as illustrated in Sec. 4.2. When some particular platform $r^\dagger$ intends to recommend services to user $j^\dagger$, it seeks for a subset of $K$ services by maximizing our accuracy-diversity measure over the similarity graph $\mathcal{G}$.

### 4.1 Step 1: Indexing Services Based on LSH

In this paper, we use *cosine similarity* to construct our LSH mechanism [32]. Specifically, given two points (or vectors) $v$ and $v'$ in a $N$-dimensional space, their cosine similarity can be defined by the cosine of the angle between them, i.e.,

$$\cos \theta_{v,v'} = \frac{v \cdot v'}{\|v\|\|v'\|} \tag{6}$$

where $\theta_{v,v'} \in [-\pi, \pi]$ denotes the angle between $v$ and $v'$. Cosine similarity $\cos \theta_{v,v'} \in [-1, 1]$ is a signed metric to measure the similarity between $v$ and $v'$. The larger is $\cos \theta_{v,v'}$, the smaller is the angle between $v$ and $v'$ and thus the more similar are $v$ and $v'$.

---

1. We hereby deliberately assume that QoS values are represented by a non-negative real numbers. It could be user rating, response time, price, reliability of service, or even their combinations. As will be shown later, our method can be readily extended to handle multi-dimensional QoS data, e.g., by considering the similarity between matrices.

TABLE 1
Frequently Used Notations

| | |
|---|---|
| $\mathcal{M} = \{1, 2, \cdots, M\}$ | A set of $M$ services |
| $\mathcal{N} = \{1, 2, \cdots, N\}$ | A set of $N$ users |
| $\mathcal{R} = \{1, 2, \cdots, R\}$ | A set of $R$ platforms |
| $\mathcal{N}_r \subset \mathcal{N}$ | A set of users on platform $r$ |
| $N_r = |\mathcal{N}_r|$ | The number of users on platform $r$ |
| $d_{i,j,r} \geq 0$ | QoS value of service $i$ for user $j$ on platform $r$ |
| $\vec{d}_{i,r} = \left( \vec{d}_{i,j,r} \right)_{j \in \mathcal{N}}$ | QoS vector of service $i$ on platform $r$ |
| $\mathcal{D} = \left\{ \vec{d}_{i,r} \right\}_{i \in \mathcal{N}, r \in \mathcal{R}}$ | QoS dataset |
| $\mathcal{H}_r = \{h_{r,1}, \cdots, h_{r,H_r}\}$ | A family of $H_r$ LSH functions on platform $r$ |
| $\mathcal{H}_r \left( \vec{d}_{i,r} \right) = \left( h_{r,1} \left( \vec{d}_{i,r} \right), \cdots, h_{r,H_r} \left( \vec{d}_{i,r} \right) \right)$ | LSH vector of service $i$ on platform $r$ |
| $v_i = \left( \mathcal{H}_1 \left( \vec{d}_{i,1} \right), \cdots, \mathcal{H}_R \left( \vec{d}_{i,R} \right) \right)$ | LSH index of service $i$ |
| $\mathcal{G} = (\mathcal{M}, \mathcal{E})$ | Similarity graph $\mathcal{G}$ consisting of vertices $\mathcal{M}$ and edges $\mathcal{E}$ |
| $j^\dagger, r^\dagger$ | Target user $j^\dagger$ and target platform $r^\dagger$ |
| $\mathcal{M}_{j^\dagger} = \{i \in \mathcal{M} \mid d_{i,j^\dagger,r^\dagger} = 0\}$ | Candidate services which we can recommend to user $j^\dagger$ on platform $r^\dagger$ |
| $\mathcal{M}_i = \{i' \in \mathcal{M} \mid (i, i') \in \mathcal{E}, d_{i',j^\dagger,r^\dagger} \neq 0\}$ | Services adjacent to $i$ in graph $\mathcal{G}$ which have been rated by user $j^\dagger$ on platform $r^\dagger$ |

Considering a hyperplane in the $N$-dimensional space which passes the origin and has its normal vector defined by $v^*$, it divides the space into two half-space, i.e., positive half-space and negative half-space . We denote by $h_{v^*}(v) : \mathbb{R}^N \to \{1, 0\}$ an indicator function parameterized by normal vector $v^*$, to specify the half-space containing point $v$. Specifically,

$$h_{v^*}(v) = \begin{cases} 1, & \cos \theta_{v,v^*} \geq 0 \\ 0, & \text{otherwise} \end{cases} \tag{7}$$

In another word, $v$ is in the positive half-space such that $h_{v^*}(v) = 1$, if $v$ is "similar" to $v^*$ with $\cos \theta_{v,v^*} \geq 0$ (and thus $-\frac{\pi}{2} \leq \theta_{v,v^*} \leq \frac{\pi}{2}$); otherwise, $h_{v^*}(v) = 0$. Given a random hyperplane defined by a random normal vector $v^*$ and two arbitrary points $v$ and $v'$, we have

$$\mathbb{P}(h_{v^*}(v) = h_{v^*}(v')) = 1 - \frac{|\theta_{v,v'}|}{\pi} \tag{8}$$

Therefore, all functions $\{h_{v^*} \mid v^* \in [-\pi, \pi]\}$ composes a LSH function family.

Recall that, for each service $i$, its quality on platform $r$ can be represented by a $N_r$-dimensional vector $\vec{d}_{i,r}$, where $N_r$ is the number of users on platform $r$. As shown in **Algorithm** 1, each platform $r$ randomly choose $H_r$ $N_r$-dimensional (normal) vectors to construct its $H_r$ LSH functions, i.e., $\mathcal{H}_r = \{h_{r,1}, h_{r,2}, \cdots, h_{r,H_r}\}$ (see Lin 2). Based on the family of LSH functions $\mathcal{H}_r$, the hash values of each service $i$ are calculated under each of the hash functions and the hash values compose a $H_r$-dimensional binary vector

$$\mathcal{H}_r \left( \vec{d}_{i,r} \right) = \left( h_{r,1} \left( \vec{d}_{i,r} \right), h_{r,2} \left( \vec{d}_{i,r} \right), \cdots, h_{r,H_r} \left( \vec{d}_{i,r} \right) \right)$$

as shown in Lines 3∼5. For each service $i \in \mathcal{M}$, it is indexed by concatenating the above vectors across different platforms. Specifically, as demonstrated in Line 7, the index of service $i$ can be represented by

$$v_i = \left( \mathcal{H}_1 \left( \vec{d}_{i,1} \right), \mathcal{H}_2 \left( \vec{d}_{i,2} \right), \cdots, \mathcal{H}_R \left( \vec{d}_{i,R} \right) \right) \tag{9}$$

Note that $v_i \in \{0, 1\}^{\sum_{r=1}^R H_r}$ for any service $i$.

---

**Algorithm 1:** A LSH-based algorithm for service indexing.

**Input:** Platforms $\mathcal{R}$, services $\mathcal{M}$, users $\mathcal{N}$, QoS data $\mathcal{D} = \left\{ \vec{d}_{i,r} \right\}_{i \in \mathcal{M}, r \in \mathcal{R}}$, the number of hash functions $H_r$ for any platform $r \in \mathcal{R}$.
**Output:** Similarity index $v_i$ for any $i \in \mathcal{M}$.

1 **for** $r = 1, 2, \cdots, R$ **do**
2     Choose $H_r$ random normal vectors and construct a family of LSH functions $\mathcal{H}_r$;
3     **for** $i = 1, 2, \cdots, M$ **do**
4        Calculate a LSH vector $\mathcal{H}_r \left( \vec{d}_{i,r} \right)$;
5     **end**
6 **end**
7 Calculate the index $v_i$ for service $i$, by concatenating the LSH vectors;

---

## 4.2 Step 2: Constructing Service Similarity Graph

In the above step, we construct a cosine similarity-based LSH mechanism, by which each service is mapped to a binary hash value such that it is highly likely for similar services to have the same hash value. Therefore, we can measure the similarity among the services according to their hash values, so as to serve the goal of privacy preservation as mentioned in Sec. 3.1.

In this step, we construct a graph $\mathcal{G} = (\mathcal{M}, \mathcal{E})$ to characterize the similarity among the services based on their LSH vectors. The graph contains $\mathcal{M}$ as vertices and we propose **Algorithm** 2 to add edges $\mathcal{E}$ to the graph according to the similarity among the vertices (i.e., the services). As shown in Line 1, $\mathcal{E}$ is first initialized by $\mathcal{E} \leftarrow \emptyset$. In each iteration $t = 1, \cdots, T$ of **Algorithm** 2, **Algorithm** 1 is called to create a LSH table $\text{Table}_t = \{(i, v_i)\}_{i \in \mathcal{M}}$ (see Line 3). The table is composed of $M$ entries and each entry is a tuple of $(i, v_i)$ . Then, for each pair of services $i$ and $i'$, if $(i, i') \notin \mathcal{E}$ and $v_i = v_{i'}$ (see Lines 4-8). It is demonstrated that, any two services $i$ and $i'$ is connected in the similarity graph if they have the same LSH indices in any one of the $T$ iterations.

As shown in **Algorithm** 1, for each service, its index is

---

**Algorithm 2:** Our method of constructing a weighted service similarity graph.

---
**Input:** Platforms $\mathcal{R}$, services $\mathcal{M}$, users $\mathcal{N}$, the number of iterations $T$.
**Output:** Service graph $\mathcal{G} = (\mathcal{M}, \mathcal{E})$.

1 $\mathcal{E} = \emptyset$;
2 **for** $t = 1, \cdots, T$ **do**
3      Call **Algorithm** 1 to calculate the LSH table $\text{Table}_t = \{(i, v_i)\}_{i \in \mathcal{M}}$;
4      **for** $(i, i') \in \mathcal{M} \times \mathcal{M}$ **do**
5          **if** $(v_i = v_{i'}) \&\& ((i, i') \notin \mathcal{E})$ **then**
6              $\mathcal{E} \leftarrow \mathcal{E} \bigcup (i, i')$;
7          **end**
8      **end**
9 **end**

---

calculated by concatenating its LSH vectors across the different platform. The service indexing algorithm is then called iteratively to construct a similarity graph in **Algorithm** 2. According to the definition of LSH mechanism (see **Definition** 1), similar services have the same LSH indices with high probability and will be connected in the similarity graph. In another word, each platform can construct a service similarity graph without sharing the original data with each other, while the privacy preservation can be ensured thanks to the irreversibility of the hash functions. Furthermore, as shown in **Algorithm** 1 and **Algorithm** 2, the numbers of hash functions $\{H_r\}_{r \in \mathcal{R}}$ and the number of hash tables $T$ both impact the constructions of the similarity graph. Their effect is described in **Theorem** 1.

**Theorem 1.** *Suppose $\theta_{i,i',r}$ represents the angle between $\vec{d}_{i,r}$ and $\vec{d}_{i',r}$ for any service $i, i' \in \mathcal{M}$. The probability for $i$ and $i'$ to have an edge in $\mathcal{G}$ can be written as*

$$\mathbb{P}((i, i') \in \mathcal{E}) = 1 - \left( 1 - \prod_{r=1}^{R} \left( 1 - \frac{|\theta_{i,i',r}|}{\pi} \right)^{H_r} \right)^{T} \quad (10)$$

*Proof.* For any services $i$ and $i'$ on platform $r$, suppose $\theta_{i,i',r} \in \left[ -\frac{\pi}{2}, \frac{\pi}{2} \right]$ denotes the angle between $\vec{d}_{i,r}$ and $\vec{d}_{i',r}$. Assume $X_{i,i',\ell} \in \{0, 1\}$ be a Bernoulli random variable indicating if $i$ and $i'$ have the same hash value under the $\ell$-th hash function; we then have $\mathbb{P}(X_{i,i',\ell} = 1) = 1 - \frac{|\theta_{i,i',r}|}{\pi}$ and $\mathbb{P}(X_{i,i',\ell} = 0) = \frac{|\theta_{i,i',r}|}{\pi}$. Since each platform chooses their hash functions independently, when there are $H_r$ hash functions on platform $r$, the probability of $v_i = v_{i'}$ can be written as

$$\mathbb{P}(v_i = v_{i'}) = \prod_{r=1}^{R} \left( 1 - \frac{|\theta_{i,i',r}|}{\pi} \right)^{H_r} \quad (11)$$

According to **Algorithm** 2, we perform the above similarity calculation $T$ times, and add an edge between $i$ and $i'$ if we judge they are similar to each other once. Hence, we have $(i, i') \in \mathcal{E}$ with probability (10)  $\square$

It is implied by the above theorem that $H_r$ and $T$ actually determine the "resolution" at which we make our similarity judgement and thus impact our QoS prediction. If we take very small values for $H_r$ and $T$, some dissimilar services (with large angles between their QoS vector) may be thought of as similar ones with the same index; while when the values of $H_r$ and $T$ are large, we may assign different index values to those similar services. We will report some empirical results in Sec. 5 to show the effect of $H_r$ and $T$.

## 4.3 Step 3: Recommending with Diversity

Now we have graph $\mathcal{G}$ where similar vertices (i.e., services) are connected by edges. We assume user $j^\dagger \in \mathcal{N}$ is the one to whom we intend to recommend new services (i.e., the services user $j^\dagger$ never invoked) through platform $r^\dagger$. Let $\overline{\mathcal{M}}_{j^\dagger} = \{i \in \mathcal{M} \mid d_{i,j^\dagger,r^\dagger} = 0\}$ be the set of target services. We first need to estimate the quality of these services. In particular, for any service $i \in \overline{\mathcal{M}}_{j^\dagger}$, we define $\mathcal{M}_i = \{i' \in \mathcal{M} \mid (i, i') \in \mathcal{E}, d_{i',j^\dagger,r^\dagger} \neq 0\}$ as the set of services which have been invoked by user $j^\dagger$ on platform $r^\dagger$ and are similar to service $i$. By exploiting the similarity graph, we then predict the quality of any service $i \in \overline{\mathcal{M}}_{j^\dagger}$ as follows

$$d_{i,j^\dagger,r^\dagger} = \frac{1}{M_i} \sum_{i' \in \mathcal{M}_i} d_{i',j^\dagger,r^\dagger} \quad (12)$$

where $M_i = |\mathcal{M}_i|$ denotes the cardinality of $\mathcal{M}_i$, i.e., the number of services in $\mathcal{M}_i$.

According to the predicted QoS values, one may want to recommend the services with the highest predicted QoS values to the target user; however, such a trivial choice may results in a homogeneous recommendation list which may not match the user's interests, as mentioned in Sec. 1. Therefore, we are now interest in seeking for a trade-off between accuracy and diversity. In the following, we first present our design of accuracy-diversity metric and formulate our $K$-Diversified Service Recommendation ($K$-DSR) problem in Sec. 4.3.1. Since the $K$-DSR problem is NP-hard (see **Theorem** 2), we propose a greedy approximation algorithm in Sec. 4.3.2 and prove its approximation ratio in Sec. 4.3.3.

### 4.3.1 Problem Formulation

Given any subset $\mathcal{K} \subseteq \mathcal{M}$, let $\text{Acc}(\mathcal{K})$ and $\text{Div}(\mathcal{K})$ be two real valued set functions to measure the accuracy and diversity of $\mathcal{K}$. Our $K$-Diversified Service Recommendation (K-DSR) problem then can be formulated as

$$\max_{\mathcal{K} \subseteq \overline{\mathcal{M}}_{j^\dagger}} \quad \mathsf{F}(\mathcal{K}) = \text{Acc}(\mathcal{K}) + \lambda \text{Div}(\mathcal{K}) \quad (13)$$

$$\text{s.t.} \quad |\mathcal{K}| = K \quad (14)$$

where $K$ is the number of services we intend to recommend and $\lambda$ denotes a factor to make a trade-off between recommendation accuracy and diversity [2]. Therein, as shown in Eq. (15), we define set function $\text{Acc}(\mathcal{K})$ to recommend services with high estimated QoS values to the target users, so as to ensure the recommendation accuracy.

$$\text{Acc}(\mathcal{K}) = \sum_{i \in \mathcal{K}} d_{i,j^\dagger,r^\dagger} \quad (15)$$

---

2. It is worthy to note that we hereby formulate the optimization for non-repetitive recommendation problems where we recommend the services which are never invoked by the target users, while our proposed method actually can be applied to the case where a service can be recommended repetitively.

As for the recommendation diversity $\mathsf{Div} : 2^{\overline{\mathcal{M}}_{j\dagger}} \to \mathbb{R}$, we take into account two types of diversity, i.e., *direct* diversity and *indirect* diversity. For any subset $\mathcal{K} \subseteq \overline{\mathcal{M}}_{j\dagger}$, we design a metric to measure its direct diversity based on *Jaccard dissimilarity*. Specifically, given two services $i$ and $i'$ as well as their QoS data $\vec{d}_{i,r\dagger} = (d_{i,1,r\dagger}, \cdots, d_{i,N,r\dagger})$ and $\vec{d}_{i',r\dagger} = (d_{i',1,r\dagger}, \cdots, d_{i',N,r\dagger})$ which are predicted upon the similar graph, their *Jaccard dissimilarity* can be calculated by

$$
\begin{aligned}
&\mathsf{J}\left(\vec{d}_{i,r\dagger}, \vec{d}_{i',r\dagger}\right) \\
=&1 - \frac{\sum_{j=1}^{N} \mathbb{I}((d_{i,j,r\dagger} \neq 0) \wedge (d_{i',j,r\dagger} \neq 0))}{N - \sum_{j=1}^{N} \mathbb{I}((d_{i,j,r\dagger} = 0) \wedge (d_{i',j,r\dagger} = 0))}
\end{aligned} \tag{16}
$$

where $\mathbb{I} : \{True, False\} \to \{1, 0\}$ is an indicator function. The direct diversity of $\mathcal{K}$ then is defined as

$$
\beta(\mathcal{K}) = \sum_{i,i' \in \mathcal{K}} \mathsf{J}\left(\vec{d}_{i,r\dagger}, \vec{d}_{i',r\dagger}\right) \tag{17}
$$

The subset $\mathcal{K}$ is with higher direct diversity if the services in the subset are less similar to each other. Nevertheless, considering the direct diversity only is not sufficient, since the direct diversity is calculated based on the estimates on the qualities of the services, which inevitably induces bias. Therefore, in addition to the direct diversity, we define a notion of indirect diversity as a complement to correct the bias. According to what we have mentioned in Sec. 3.2, it is more likely for the services in $\mathcal{K}$ to share the same expanded neighbors in graph $\mathcal{G}$ when $\mathcal{K}$ has smaller expansion ratio, while the service may be similar to each other if they are connected to and thus similar to the same expanded neighbors in graph $\mathcal{G}$. Hence, for any subset $\mathcal{K} \subseteq \mathcal{M}$, it is with higher (indirect) diversity if it has larger expansion ratio $\alpha(\mathcal{K})$. By combining both the direct and indirect diversities, we define

$$
\mathsf{Div}(\mathcal{K}) = \alpha(\mathcal{K}) + \xi\beta(\mathcal{K}) \tag{18}
$$

where $\xi > 0$ is a weight factor such that we can fine tune the trade-off between the direct and indirect diversity.

**Theorem 2.** *Our K-DSR problem defined by (13) and (14) is NP-hard.*

*Proof.* The basic idea of the proof is to show the NP-harness of the following two sub-problems

$$
\mathfrak{P}_1 : \max_{\mathcal{K} \subseteq \overline{\mathcal{M}}_{j\dagger}} \quad \alpha(\mathcal{K})
$$
$$
\text{s.t.} \quad |\mathcal{K}| = K
$$

and

$$
\mathfrak{P}_2 : \max_{\mathcal{K} \subseteq \overline{\mathcal{M}}_{j\dagger}} \quad \beta(\mathcal{K})
$$
$$
\text{s.t.} \quad |\mathcal{K}| = K
$$

The first problem $\mathfrak{P}_1$ is to find a subset $\mathcal{K} \subseteq \overline{\mathcal{M}}_{j\dagger}$ of size $K$ with the maximum expansion ratio. We prove its NP-hardness by demonstrating a reduction from the maximum coverage problem (which is a well-known NP-hard problem). In the maximum coverage problem, we assume $\mathcal{M}'$ denotes a set of $M'$ elements and $\mathcal{C} = \{\mathcal{C}_1, \mathcal{C}_2, \cdots, \mathcal{C}_C\}$ represent a collection of $C$ subsets of $\mathcal{M}'$ such that $\bigcup_\ell \mathcal{C}_\ell = \mathcal{M}'$. Given a positive integer $K < C$, the goal of the maximum coverage problem is to find $K$ out of the $C$ subsets such that their union has the maximum cardinality. To construct an instance of the maximum coverage problem, for each $\mathcal{C}_\ell \in \mathcal{C}$ and $i \in \mathcal{M}'$, we create node $a_\ell$ and $b_i$, respectively. If $i \in \mathcal{C}_\ell$, we add a directed edge from $a_\ell$ to $b_i$ and it is said that $b_i$ (resp. element $i \in \mathcal{M}'$) is "covered" by $a_\ell$ (resp. subset $\mathcal{C}_\ell$). Let $\mathcal{A} = \{a_\ell\}_\ell$ and $\mathcal{B} = \{b_i\}_i$. The edges are only from the nodes in $\mathcal{A}$ to the ones in $\mathcal{B}$; then, for any feasible solution to the maximum coverage problem $\widetilde{\mathcal{C}} \subseteq \mathcal{C}$, its expansion ratio is calculated by $|\mathcal{B}_{\widetilde{\mathcal{C}}}| / K$, where $\mathcal{B}_{\widetilde{\mathcal{C}}} \subseteq \mathcal{B}$ denotes the subset of the nodes in $\mathcal{B}$ covered by $\widetilde{\mathcal{C}}$, i.e., the union of the subsets in $\widetilde{\mathcal{C}}$. Therefore, finding a subset of $K$ nodes in $\mathcal{A}$ with the highest expansion ratio is equivalent to finding $K$ subsets in $\mathcal{C}$ with the highest coverage, and vice versa.

As for the second optimization problem $\mathfrak{P}_2$, it is actually a maximum dispersion problem, since the Jaccard distance (or dissimilarity) satisfies triangle inequality [33]. Therefore, the NP hardness of our problem can be proved. $\square$

### 4.3.2 Algorithm Design

Due to the NP-hardness of the optimization problem, we have to be content with approximation algorithms. We propose a greedy algorithm to address the problem in this section, as shown in **Algorithm** 3. Let $\mathcal{K}$ denote the recommended list which is initialized as an empty set (see Line 1). The algorithm proceeds iteratively. In each iteration, as shown in Line 3, we first find service $i^\dagger \in \overline{\mathcal{M}}_{j\dagger}$ such that the set function $\mathsf{F}'(\mathcal{K} \cup \{i^\dagger\})$ can be maximized, where $\mathsf{F}' : 2^{\overline{\mathcal{M}}_{j\dagger}} \to \mathbb{R}$ is defined as

$$
\mathsf{F}'(\mathcal{K}) = \frac{1}{2}\left(\mathsf{Acc}(\mathcal{K}) + \lambda\alpha(\mathcal{K})\right) + \lambda\xi\beta(\mathcal{K}) \tag{19}
$$

for $\forall \mathcal{K} \subseteq \overline{\mathcal{M}}_{j\dagger}$. We then add $i^\dagger$ into $\mathcal{K}$ and remove $i^\dagger$ from $\overline{\mathcal{M}}_{j\dagger}$, as demonstrated in Line 4 and Line 5, respectively.

---

**Algorithm 3:** Top-$K$ service recommendation algorithm

---

**Input:** Service similarity graph $\mathcal{G} = (\mathcal{M}, \mathcal{E})$, target user $j^\dagger$ and target platform $r^\dagger$, QoS data $\mathcal{D}$, the number of recommended entries $K$.

**Output:** Recommended list $\mathcal{K} \subseteq \mathcal{M}$ for user $j^\dagger$ on platform $r^\dagger$.

1 $\mathcal{K} = \emptyset$;
2 **while** $|\mathcal{K}| < K$ **do**
3     $i^\dagger = \arg\max_{i \in \overline{\mathcal{M}}_{j\dagger}} \mathsf{F}'(\mathcal{K} \bigcup \{i\})$;
4     $\mathcal{K} \leftarrow \mathcal{K} \cup \{i^\dagger\}$;
5     $\overline{\mathcal{M}}_{j\dagger} \leftarrow \overline{\mathcal{M}}_{j\dagger} / \{i^\dagger\}$;
6 **end**

---

### 4.3.3 Analysis

In this section, we will analyze the performance of **Algorithm** 3. Before showing its approximation ratio in **Theorem** 3, we first present two lemmas as follows, which will be very helpful in the proof of **Theorem** 3.

**Lemma 2.** *For $\forall \mathcal{K} \subseteq \mathcal{M}$, the following set function defined by*

$$
\Phi'(\mathcal{K}) = \mathsf{Acc}(\mathcal{K}) + \lambda\alpha(\mathcal{K}) = \sum_{i \in \mathcal{K}} d_{i,j\dagger,r\dagger} + \frac{\lambda|\mathsf{Exp}(\mathcal{K})|}{M} \tag{20}
$$

*is a monotone submodular set function.*

*Proof.* For $\forall \mathcal{K}_1 \subseteq \mathcal{K}_2 \subseteq \mathcal{M}$, we have $\mathsf{Exp}(\mathcal{K}_1) \subseteq \mathsf{Exp}(\mathcal{K}_2)$ and $\sum_{i \in \mathcal{K}_1} d_{i,j^\dagger,r^\dagger} \leq \sum_{i \in \mathcal{K}_2} d_{i,j^\dagger,r^\dagger}$. Furthermore, $\forall i \in \mathcal{M}/\mathcal{K}_2$, we have

$$\Phi'(\mathcal{K}_1 \cup \{i\}) - \Phi'(\mathcal{K}_1) - (\Phi'(\mathcal{K}_2 \cup \{i\}) - \Phi'(\mathcal{K}_2))$$
$$= \lambda \frac{|\mathsf{Exp}(\mathcal{K}_1 \cup \{i\})| - |\mathsf{Exp}(\mathcal{K}_1)|}{M}$$
$$- \lambda \frac{|\mathsf{Exp}(\mathcal{K}_2 \cup \{i\})| - |\mathsf{Exp}(\mathcal{K}_2)|}{M} \quad (21)$$

Therein, $|\mathsf{Exp}(\mathcal{K}_1 \bigcup\{i\})| - |\mathsf{Exp}(\mathcal{K}_1)| = |\mathsf{Exp}(\{i\}) - \mathsf{Exp}(\mathcal{K}_1)|$ and $|\mathsf{Exp}(\mathcal{K}_2 \bigcup\{i\})| - |\mathsf{Exp}(\mathcal{K}_2)| = |\mathsf{Exp}(\{i\}) - \mathsf{Exp}(\mathcal{K}_2)|$; we then obtain

$$\Phi'(\mathcal{K}_1 \cup \{i\}) - \Phi'(\mathcal{K}_1) - (\Phi'(\mathcal{K}_2 \cup \{i\}) - \Phi'(\mathcal{K}_2))$$
$$= \frac{\lambda(|\mathsf{Exp}(\{i\}) - \mathsf{Exp}(\mathcal{K}_1)| - |\mathsf{Exp}(\{i\}) - \mathsf{Exp}(\mathcal{K}_2)|)}{M} \geq 0 \quad (22)$$

which completes the proof according to **Definition** 4. $\square$

**Lemma 3.** *Given the Jaccard dissimilarity* $\mathsf{J}(\cdot, \cdot)$ *defined in Eq. (16), for any two disjoint nonempty sets* $\mathcal{K}_1, \mathcal{K}_2 \subseteq \mathcal{M}$*, we have the following inequality:*

$$(|\mathcal{K}_1| - 1)\Gamma(\mathcal{K}_1, \mathcal{K}_2) \geq |\mathcal{K}_2|\Gamma(\mathcal{K}_1) \quad (23)$$

*where* $\Gamma(\mathcal{K}_1) = \sum_{i,i' \in \mathcal{K}_1} \mathsf{J}(d_{i,j^\dagger,r^\dagger}, d_{i',j^\dagger,r^\dagger})$ *and* $\Gamma(\mathcal{K}_1, \mathcal{K}_2) = \sum_{i \in \mathcal{K}_1} \sum_{i' \in \mathcal{K}_2} \mathsf{J}(d_{i,j^\dagger,r^\dagger}, d_{i',j^\dagger,r^\dagger})$,

*Proof.* For any $i_2 \in \mathcal{K}_2$ and any $i_1, i'_1 \in \mathcal{K}_1$ such that $i_1 \neq i'_1$, we have

$$\mathsf{J}(d_{i_2,j^\dagger,r^\dagger}, d_{i_1,j^\dagger,r^\dagger}) + \mathsf{J}(d_{i_2,j^\dagger,r^\dagger}, d_{i'_1,j^\dagger,r^\dagger})$$
$$\geq \mathsf{J}(d_{i_1,j^\dagger,r^\dagger}, d_{i'_1,j^\dagger,r^\dagger}) \quad (24)$$

according to the triangle inequality; we thus obtain

$$\Gamma(\mathcal{K}_1) \leq \mathsf{J}(d_{i_1,j^\dagger,r^\dagger}, d_{i'_1,j^\dagger,r^\dagger})$$
$$= (|\mathcal{K}_1| - 1) \sum_{i_1 \in \mathcal{K}_1} \mathsf{J}(d_{i_1,j^\dagger,r^\dagger}, d_{i_2,j^\dagger,r^\dagger}) \quad (25)$$

for any $i_2 \in \mathcal{K}_2$. The lemma can be completed by summing the above inequality over $i_2 \in \mathcal{K}_2$. $\square$

Based on the **Definition** 4, the problem defined in Eq.(8) and Eq.(9) can be regarded as the max-sum diversification problem with monotone submodular set functions satisfying a cardinality constraint. Thus, **Theorem** 3 is as follows.

**Theorem 3.** *Our top-$K$ recommendation algorithm (see **Algorithm** 3) is a 2-approximation algorithm.*

*Proof.* Let $\mathcal{K}^*$ be the optimal solution to the problem K-DSR. For any service $i \in \mathcal{K}^*$, it is said to be "correct"; otherwise, it is "incorrect". Recall that $\mathcal{K}^\dagger$ denotes the recommendation list yielded by our greedy algorithm, and let $\mathcal{K}^\dagger_{[t]}$ be the (intermediate) output in the $t$-th iteration. Assume $\mathcal{K}_{[t]} = \mathcal{K}^* \bigcap \mathcal{K}^\dagger_{[t]}$, $\widetilde{\mathcal{K}}^\dagger_{[t]} = \mathcal{K}^\dagger_{[t]}/\mathcal{K}_{[t]}$ and $\mathcal{K}^*_{[t]} = \mathcal{K}^*/\mathcal{K}_{[t]}$. In another word, $\mathcal{K}_{[t]}$ and $\widetilde{\mathcal{K}}^\dagger_{[t]}$ denote the set of correct services and the set of incorrect ones we select till the $t$-th iteration, respectively, while $\mathcal{K}^*_{[t]}$ is the set of correct services we do not select till the $t$-th iteration. Since $\mathcal{K}_{[t]} \bigcup \mathcal{K}^*_{[t]} = \mathcal{K}^*$, it follows that

$$\Gamma\left(\mathcal{K}_{[t]}, \mathcal{K}^*_{[t]}\right) + \Gamma\left(\mathcal{K}_{[t]}\right) + \Gamma\left(\mathcal{K}^*_{[t]}\right) - \Gamma(\mathcal{K}^*) = 0 \quad (26)$$

Furthermore, according to **Lemma** 3, we have the following inequalities

$$\left(|\mathcal{K}^*_{[t]}| - 1\right) \cdot \Gamma\left(\widetilde{\mathcal{K}}^\dagger_{[t]}, \mathcal{K}^*_{[t]}\right) - |\widetilde{\mathcal{K}}^\dagger_{[t]}| \cdot \Gamma\left(\mathcal{K}^*_{[t]}\right) \geq 0 \quad (27a)$$

$$\left(|\mathcal{K}^*_{[t]}| - 1\right) \cdot \Gamma\left(\mathcal{K}_{[t]}, \mathcal{K}^*_{[t]}\right) - |\mathcal{K}_{[t]}| \cdot \Gamma\left(\mathcal{K}^*_{[t]}\right) \geq 0 \quad (27b)$$

$$(|\mathcal{K}_{[t]}| - 1) \cdot \Gamma\left(\mathcal{K}_{[t]}, \mathcal{K}^*_{[t]}\right) - |\mathcal{K}^*_{[t]}| \cdot \Gamma\left(\mathcal{K}_{[t]}\right) \geq 0 \quad (27c)$$

Before diving into the details of the proof, we first introduce some symbols and notations. For any $\mathcal{K} \subseteq \overline{\mathcal{M}}_j$ and $i \in \overline{\mathcal{M}}_j/\mathcal{K}$, suppose

$$\begin{cases} \Delta\Phi'_i(\mathcal{K}) = \Phi'(\{i\}\bigcup\mathcal{K}) - \Phi'(\mathcal{K}) \\ \Delta\beta_i(\mathcal{K}) = \beta(\{i\}\bigcup\mathcal{K}) - \beta(\mathcal{K}) \\ \Delta\mathsf{F}'_i(\mathcal{K}) = \mathsf{F}'(\{i\}\bigcup\mathcal{K}) - \mathsf{F}'(\mathcal{K}) = \frac{1}{2}\Delta\Phi'_i(\mathcal{K}) + \lambda\xi\Delta\beta_i(\mathcal{K}) \end{cases}$$

**Case** 1: We first take into account the first case where $|\mathcal{K}^*_{[t]}| = 1$ and thus $|\mathcal{K}^\dagger_{[t]}| = K - 1$ and $\mathcal{K}^\dagger_{[t]} \subseteq \mathcal{K}^*$ in the $t$-th iteration. In another word, $K - 1$ out of the $K$ correct services are selected by our greedy algorithm till the $t$-th iteration. We denote by $i^*$ the only element in $\mathcal{K}^*_{[t]}$. We assume $i_{[t+1]}$ is the service selected by our greedy algorithm in the following $(t + 1)$-th iteration; we then have

$$\Delta\mathsf{F}'_{i_{[t+1]}}\left(\mathcal{K}^\dagger_{[t]}\right) \geq \Delta\mathsf{F}'_{i^*}\left(\mathcal{K}^\dagger_{[t]}\right)$$
$$= \frac{1}{2}\Delta\Phi'_{i^*}\left(\mathcal{K}^\dagger_{[t]}\right) + \lambda\xi\Delta\beta_{i^*}\left(\mathcal{K}^\dagger_{[t]}\right)$$
$$\geq \frac{1}{2}\left(\Delta\Phi'_{i^*}\left(\mathcal{K}^\dagger_{[t]}\right) + \lambda\xi\Delta\beta_{i^*}\left(\mathcal{K}^\dagger_{[t]}\right)\right)$$
$$= \frac{1}{2}\Delta\mathsf{F}'_{i^*}\left(\mathcal{K}^\dagger_{[t]}\right) \quad (28)$$

and thus

$$\Delta\mathsf{F}_{i_{[t+1]}}(\mathcal{K}^\dagger_{[t]}) = \mathsf{F}\left(\{i_{[t+1]}\}\bigcup\mathcal{K}^\dagger_{[t]}\right) - \mathsf{F}\left(\mathcal{K}^\dagger_{[t]}\right)$$
$$= \Delta\Phi'_{i_{[t+1]}}(\mathcal{K}^\dagger_{[t]}) + \lambda\xi\Delta\beta_{i_{[t+1]}}(\mathcal{K}^\dagger_{[t]})$$
$$\geq \frac{1}{2}\Delta\Phi'_{i_{[t+1]}}(\mathcal{K}^\dagger_{[t]}) + \lambda\xi\Delta\beta_{i_{[t+1]}}(\mathcal{K}^\dagger_{[t]})$$
$$= \Delta\mathsf{F}'_{i_{[t+1]}}(\mathcal{K}^\dagger_{[t]})$$
$$\geq \frac{1}{2}\Delta\mathsf{F}_{i^*}(\mathcal{K}^\dagger_{[t]}) \quad (29)$$

according to which, we obtain

$$\mathsf{F}\left(\mathcal{K}^\dagger\right) = \mathsf{F}\left(\mathcal{K}^\dagger_{[t]}\right) + \Delta\mathsf{F}_{i_{[t+1]}}\left(\mathcal{K}^\dagger_{[t]}\right)$$
$$\geq \mathsf{F}\left(\mathcal{K}^\dagger_{[t]}\right) + \frac{1}{2}\Delta\mathsf{F}_{i^*}(\mathcal{K}^\dagger_{[t]})$$
$$\geq \frac{1}{2}\left(\mathsf{F}\left(\mathcal{K}^\dagger_{[t]}\right) + \Delta\mathsf{F}_{i^*}(\mathcal{K}^\dagger_{[t]})\right)$$
$$= \frac{1}{2}\left(\mathsf{F}\left(\mathcal{K}^\dagger_{[t]}\right) + \mathsf{F}(\mathcal{K}^*) - \mathsf{F}\left(\mathcal{K}^\dagger_{[t]}\right)\right)$$
$$= \frac{1}{2}\mathsf{F}(\mathcal{K}^*) \quad (30)$$

**Case 2**: We then consider the case with $|\mathcal{K}^*_{[t]}| > 1$. Since $|\mathcal{K}^*_{[t]}| > 1$, $K \geq 1$ and $|\mathcal{K}^*_{[t]}| - |\widetilde{\mathcal{K}}^\dagger_{[t]}| = K - t \geq 0$, we take them as multipliers and calculate

$$\frac{1}{|\mathcal{K}^*_{[t]}| - 1} \cdot (27a) + \frac{|\mathcal{K}^*_{[t]}| - |\widetilde{\mathcal{K}}^\dagger_{[t]}|}{K\left(|\mathcal{K}^*_{[t]}| - 1\right)} \cdot (27b)$$
$$+ \frac{t}{K(K-1)} \cdot (27c) + \frac{t|\mathcal{K}^*_{[t]}|}{K(K-1)} \cdot (26) \quad (31)$$

as

$$\gamma_0 \cdot \Gamma\left(\mathcal{K}_{[t]}, \mathcal{K}^*_{[t]}\right) + \Gamma\left(\widetilde{\mathcal{K}}^\dagger_{[t]}, \mathcal{K}^*_{[t]}\right)$$
$$+ \gamma_1 \cdot \Gamma(\mathcal{K}^*_{[t]}) - \frac{k|\mathcal{K}^*_{[t]}|}{K(K-1)} \cdot \Gamma(\mathcal{K}^*) \geq 0 \quad (32)$$

where

$$\gamma_0 = \frac{|\mathcal{K}^*_{[t]}| - |\widetilde{\mathcal{K}}^\dagger_{[t]}|}{K} + \frac{k(|\mathcal{K}_{[t]}| + |\mathcal{K}^*_{[t]}| - 1)}{K(K-1)} = 1 \quad (33)$$

and

$$\gamma_1 = \frac{k|\mathcal{K}^*_{[t]}|}{K(K-1)} - \frac{|\mathcal{K}_{[t]}|(|\mathcal{K}^*_{[t]}| - |\widetilde{\mathcal{K}}^\dagger_{[t]}|)}{K(|\mathcal{K}^*_{[t]}| - 1)} - \frac{|\widetilde{\mathcal{K}}^\dagger_{[t]}|}{|\mathcal{K}^*_{[t]}| - 1}$$
$$= \frac{k|\mathcal{K}^*_{[t]}|(K - |\mathcal{K}^*_{[t]}|)}{K(K-1)(|\mathcal{K}^*_{[t]}| - 1)} \geq 0 \quad (34)$$

as $|\mathcal{K}^*_{[t]}| = K - |\mathcal{K}_{[t]}|$, $|\mathcal{K}^*_{[t]}| = t = |\mathcal{K}_{[t]}|$ and $|\mathcal{K}^*_{[t]}| - |\widetilde{\mathcal{K}}^\dagger_{[t]}| + t = K$ naturally hold according to the definitions of $\mathcal{K}_{[t]}$, $\widetilde{\mathcal{K}}^\dagger_{[t]}$ and $\mathcal{K}^*_{[t]}$. Hence, we now have

$$\Gamma(\mathcal{K}^*_{[t]}, \mathcal{K}^\dagger_{[t]}) = \Gamma(\mathcal{K}_{[t]}, \mathcal{K}^*_{[t]}) + \Gamma(\widetilde{\mathcal{K}}^\dagger_{[t]}, \mathcal{K}^*_{[t]}) \geq \frac{t|\mathcal{K}^*_{[t]}|}{K(K-1)}\Gamma(\mathcal{K}^*) \quad (35)$$

Considering $\Phi'(\cdot)$ is a monotone submodular function, we have

$$\Phi'\left(\mathcal{K}^*_{[t]} \bigcup \mathcal{K}^\dagger_{[t]}\right) - \Phi'\left(\mathcal{K}^\dagger_{[t]}\right) \geq \Phi'(\mathcal{K}^*) - \Phi'\left(\mathcal{K}^\dagger\right) \quad (36)$$

since $\mathcal{K}^* \subseteq \mathcal{K}^*_{[t]} \bigcup \mathcal{K}_{[t]}$ and $\mathcal{K}^\dagger_{[t]} \subseteq \mathcal{K}^\dagger$. Moreover, according to **Lemma** 1, we also have

$$\sum_{i \in \mathcal{K}^*_{[t]}} \Delta\Phi'_i\left(\mathcal{K}^\dagger_{[t]}\right) = \sum_{i \in \mathcal{K}^*_{[t]}} \left(\Phi'\left(\{i\} \bigcup \mathcal{K}^\dagger_{[t]}\right) - \Phi'\left(\mathcal{K}^\dagger_{[t]}\right)\right)$$
$$\geq \Phi'\left(\mathcal{K}^*_{[t]} \bigcup \mathcal{K}^\dagger_{[t]}\right) - \Phi'\left(\mathcal{K}^\dagger_{[t]}\right)$$
$$\geq \Phi'(\mathcal{K}^*) - \Phi'\left(\mathcal{K}^\dagger\right) \quad (37)$$

Then, we obtain

$$\sum_{i \in \mathcal{K}^*_{[t]}} \Delta\mathsf{F}'_i\left(\mathcal{K}^\dagger_{[t]}\right)$$
$$= \sum_{i \in \mathcal{K}^*_{[t]}} \left(\mathsf{F}'\left(\{i\} \bigcup \mathcal{K}^\dagger_{[t]}\right) - \mathsf{F}'\left(\mathcal{K}^\dagger_{[t]}\right)\right)$$
$$= \frac{1}{2}\sum_{i \in \mathcal{K}^*_{[t]}} \Delta\Phi'_i\left(\mathcal{K}^\dagger_{[t]}\right) + \lambda\xi\Gamma\left(\mathcal{K}^*_{[t]}, \mathcal{K}^\dagger_{[t]}\right)$$
$$\geq \frac{1}{2}\left(\Phi'(\mathcal{K}^*) - \Phi'\left(\mathcal{K}^\dagger\right)\right) + \lambda\xi\frac{t|\mathcal{K}^*_{[t]}|}{K(K-1)}\Gamma(\mathcal{K}^*) \quad (38)$$

Let $i_{[t]}$ denote the service selected by our greedy algorithm in the $(t+1)$-th iteration such that $\Delta\mathsf{F}'_{i_{[t]}}(\mathcal{K}^\dagger_{[t]}) \geq \Delta\mathsf{F}'_i(\mathcal{K}^\dagger_{[t]})$ for any $i \in \mathcal{K}^*_{[t]}$. Therefore, we have

$$\Delta\mathsf{F}'_{i_{[t+1]}}(\mathcal{K}^\dagger_{[t]})$$
$$\geq \frac{1}{|\mathcal{K}^*_{[t]}|} \sum_{i \in \mathcal{K}^*_{[t]}} \Delta\mathsf{F}'_i(\mathcal{K}^\dagger_{[t]})$$
$$\geq \frac{1}{2|\mathcal{K}^*_{[t]}|}\left(\Phi'(\mathcal{K}^*) - \Phi'(\mathcal{K}^\dagger)\right) + \lambda\xi\frac{t}{K(K-1)}\Gamma(\mathcal{K}^*)$$
$$\geq \frac{1}{2K}\left(\Phi'(\mathcal{K}^*) - \Phi'(\mathcal{K}^\dagger)\right) + \lambda\xi\frac{t}{2K(K-1)}\Gamma(\mathcal{K}^*) \quad (39)$$

according to which, we can re-write $\mathsf{F}'(\mathcal{K}\dagger)$

$$\mathsf{F}'(\mathcal{K}^\dagger) = \sum_{t=0}^{K-1} \Delta\mathsf{F}'_{i_{[t+1]}}(\mathcal{K}^\dagger_{[t]})$$
$$\geq \frac{1}{2}\left(\Phi'(\mathcal{K}^*) - \Phi'(\mathcal{K}^\dagger)\right) + \frac{\lambda\xi}{2}\Gamma(\mathcal{K}^*) \quad (40)$$

Since $\mathsf{F}'(\mathcal{K}^\dagger) = \frac{1}{2}\Phi'(\mathcal{K}^\dagger) + \lambda\xi\Gamma(\mathcal{K}^\dagger)$, we obtain

$$\Phi'(\mathcal{K}^\dagger) + \lambda\xi\Gamma(\mathcal{K}^\dagger) \geq \frac{1}{2}\Phi'(\mathcal{K}^*) + \frac{\lambda\xi}{2}\Gamma(\mathcal{K}^*) = \frac{1}{2}\mathsf{F}(\mathcal{K}^*) \quad (41)$$

Finally,

$$\mathsf{F}(\mathcal{K}^\dagger) = \Phi'(\mathcal{K}^\dagger) + \lambda\xi\Gamma(\mathcal{K}^\dagger) \geq \frac{1}{2}\mathsf{F}(\mathcal{K}^*) \quad (42)$$

$\square$

We have proved the approximation ratio of **Algorithm** 3, and we then analyze its time complexity in **Theorem** 4.

**Theorem 4.** *The time complexity of our top-$K$ service recommendation algorithm (see **Algorithm** 3) is $\mathcal{O}(M^2K^3 + MNK^4)$.*

*Proof.* As demonstrated in **Algorithm** 3, in each iteration, it seeks for a service $i^\dagger \in \overline{\mathcal{M}}_{j^\dagger}$ such that $\mathsf{F}'(\mathcal{K}\bigcup\{i^\dagger\})$ is maximized by $\mathcal{K}\bigcup\{i^\dagger\}$, where $\mathcal{K}$ is the set of services which are already selected. The above procedure is repeated until $K$ services is selected in total.

Let $\mathcal{K}_{[t]} \subseteq \overline{\mathcal{M}}_{j^\dagger}$ denote the set of selected services in the $t$-th iteration and $K_{[t]} = |\mathcal{K}_{[t]}| = t$. In the $(t+1)$-th iteration, the time complexity of calculating $\mathsf{F}'(\mathcal{K}_{[t]}\bigcup\{i\})$ for $i \in \overline{\mathcal{M}}_{j^\dagger}/\mathcal{K}_{[t]}$ relies on the ones of computing $\mathsf{Acc}(\mathcal{K}_{[t]}\bigcup\{i\})$, $\alpha(\mathcal{K}_{[t]}\bigcup\{i\})$ and $\beta(\mathcal{K}_{[t]}\bigcup\{i\})$, as illustrated by (19). According to the definition of $\mathsf{Acc}$ shown in (15), the time complexity of calculating $\mathsf{Acc}$ is $\mathcal{O}(K_{[t]})$. Additionally, according to our definitions of $\alpha$ and $\beta$ as shown in (3) and (17), respectively, the time complexities are $\mathcal{O}(MK_{[t]})$ and $\mathcal{O}(NK_{[t]}^2)$, respectively. Therefore, the time complexity of calculating $\mathsf{F}'(\mathcal{K}_{[t]}\bigcup\{i\})$ in the the $(t+1)$-th iteration is $\mathcal{O}(MK_{[t]} + NK_{[t]}^2)$ and the one of calculating $i^\dagger = \arg\max_{i \in \overline{\mathcal{M}}_{j^\dagger}/\mathcal{K}_{[t]}} \mathsf{F}'(\mathcal{K}_{[t]}\bigcup\{i\})$ is thus $\mathcal{O}\left(\left(MK_{[t]} + NK_{[t]}^2\right)(\overline{M}_{j^\dagger} - t)\right)$. Since $\mathcal{K}_{[t+1]} = \mathcal{K}_{[t]}\bigcup\{i^\dagger\}$, the total time complexity across the $K$ iterations is $\mathcal{O}(M^2K^3 + MNK^4)$. $\square$

Note that $K$ is usually considered as a constant for any target user, which is much smaller than $M$ and $N$. The actual time complexity of our algorithm can be re-written as $\mathcal{O}(M^2 + MN)$.

# 5 EXPERIMENT

In this section, we first present our experiment settings and the metrics for evaluation in Sec. 5.1. We then introduce the state-of-the-art methods in 5.2 and compare them with our PDSR method in Sec. 5.3.

## 5.1 Experiment Settings and Metrics

In this paper, we conduct our experiments on both the WS-DREAM dataset [34] and the MovieLens dataset [35] to evaluate our PDSR method. The WS-DREAM dataset contains real-world QoS data, e.g., response time values, collected from $339$ users interacting with $5,825$ web services. To facilitate our comparison, each QoS value is normalized into a range of $[0, 1]$. The MovieLens dataset is another commonly used dataset for studying service recommendation. It involves the ratings given by $6,040$ users for $3,952$ movies. We thus treat these ratings as the QoS values. For each user $j$ in WS-DREAM dataset and the one in MovieLens dataset, we randomly choose $\mathcal{S}_j \subset \mathcal{M}_j$ such that $S = |\mathcal{S}_j| = 15$, where $\mathcal{M}_j$ denotes the set of the services invoked by user $j$, and let $\mathcal{S}_j$ be the test dataset. Specifically, since MovieLens dataset is rather sparse, we take into account only the users who have no less than $25$ service usage records such that we have sufficient historical QoS data to make recommendation decisions for each target user by our method. We assume there are two platforms and divide the WS-DREAM dataset into two subsets such that the platforms have $135$ users and $204$ users, respectively. Likewise, the MovieLens dataset is divided such that the two platforms have $2,249$ users and $3,375$ users, respectively. We randomly pick up $15$ users on each platform as target users, to each of which, we deliver a recommendation list containing $K = 5$ services.

We use the following four metrics to evaluate our PDSR method as well as the reference ones.

- *Mean Absolute Error* (MAE) and *Root Mean Squared Error* (RMSE). Since service recommendation highly rely on the predictions of the QoS values, we hereby adopt MAE and RMSE to measure the difference between our predictions and the actual values. Specifically, given a set of target users $\mathcal{N}_r^\dagger \subseteq \mathcal{N}$ on platform $r$, the MAE can be defined by

$$\text{MAE} = \frac{1}{N_r^\dagger S} \sum_{j \in \mathcal{N}_r^\dagger} \sum_{i \in \mathcal{S}_j} \left| d_{i,j,r} - d_{i,j,r}^* \right| \quad (43)$$

where $N_r^\dagger = |\mathcal{N}_r^\dagger|$ denotes the number of the target users on platform $r$, while $d_{i,j,r}$ and $d_{i,j,r}^*$ denote the predicted QoS value of service $i$ for user $j$ on platform $r$ and its ground truth, respectively. Similarly, the RMSE on platform $r$ can be calculated as

$$\text{RMSE} = \sqrt{\frac{1}{N_r^\dagger S} \sum_{j \in \mathcal{N}_r^\dagger} \sum_{i \in \mathcal{S}_j} \left( d_{i,j,r} - d_{i,j,r}^* \right)^2} \quad (44)$$

According to the above definitions, smaller MAE and RMSE indicate better prediction accuracy.

- *Average Quality of Service* (AQoS). We employ AQoS to reflect users' actual interests in the recommended services. In particular, we define AQoS as

$$\text{AQoS} = \frac{1}{K N_r^\dagger} \sum_{j \in \mathcal{N}_r^\dagger} \sum_{i \in \mathcal{K}_j^\dagger} d_{i,j,r}^* \quad (45)$$

where $\mathcal{K}_j^\dagger$ denotes the recommendation list for user $j \in \mathcal{N}_r^\dagger$ on platform $r$. Higher AQoS implies that the diverse services we recommend to target users have better actual quality.

- *Inter-List Diversity* (ILD). ILD is a commonly used metric to measure the diversity of a recommendation list [5]. Specifically,

$$\text{ILD} = \frac{1}{N_r^\dagger} \sum_{j \in \mathcal{N}_r^\dagger} \left( \frac{\sum_{i,i' \in \mathcal{K}_j^\dagger, i \neq i'} \mathsf{J}_{i,i',r}^*}{K(K-1)} \right) \quad (46)$$

where $\mathsf{J}_{i,i',r}^* = \mathsf{J}\left( \vec{d}_{i,r}^*, \vec{d}_{i',r}^* \right)$, as defined in Eq. 16, $\mathsf{J}\left( \vec{d}_{i,r}^*, \vec{d}_{i',r}^* \right)$ represents the *Jaccard dissimilarity* between services $i, i'$. The higher the ILD is, the higher the recommendation diversity is.

## 5.2 Reference Methods

In this section, we compare our method with the following six state-of-the-art ones.

- **BPR method** [36]: The *Bayesian Personalized Ranking* (BPR) method applies a Bayesian analysis to the personalized ranking problem. In particular, a maximum posterior estimator is derived to predict a personalized ranking.
- **MPR method** [37]: The *Multiple Pairwise Ranking* (MPR) method deeply exploits the unobserved item feedback to perform multiple pairwise ranking. Specifically, it divides the unobserved items into different parts to exploit the preference difference among multiple pairs of items, which is thus thought of as a multiple pairwise model.
- **EF method** [26]: In the *Entropy Fusion* (EF) method based on CF, accuracy-diversity trade-off are leveraged by taking into account both user exposure diversity and item concordance.
- **MF-based method** [38]: The *Matrix Factoring* (MF)-based method considers both the explicit and implicit information during the matrix factoring process;popularity bias with a weighting mechanism and neighborhood information are used to diversifying the recommendation results.
- **DPP-based method** [25] : Since DPP is a powerful tool for modeling diversity, [25] proposes an algorithm to accelerate the greedy MAP inference for DPP, and the algorithm can be applied to the task of diversified recommendation.
- **LSH-based method** [30]: In this method, the LSH mechanism is used to calculate both similarity and diversity. Additionally, thanks to the LSH mechanism, this method can be adapted to utilize across-platform data with privacy preserved.

Note that the BPR method and the MPR method only aim at improving recommendation accuracy, while the other four methods investigate the accuracy-diversity trade-off. Moreover, since the first five of the above methods do not take into account the privacy issue, we assume that they utilize only the original data on single platforms.

### 5.3 Experiment Results

#### 5.3.1 Comparison with Different Methods

We first compare our PDSR method with the six reference ones under the metrics including MAE, RMSE, AQoS, and ILD. We let $H_1 = H_2 = 3$, $T = 9$, $\lambda = 0.1$ and $\xi = 0.3$ for Platform 1 and $H_1 = H_2 = 3$, $T = 9$ and $\lambda = \xi = 0.3$ for Platform 2, when performing our experiments on the WS-DREAM dataset. As for the MovieLens dataset, we let $H_1 = H_2 = 4$, $T = 9$, $\lambda = 0.1$ and $\xi = 0.3$ for Platform 1 and $H_1 = H_2 = 5$, $T = 9$, $\lambda = 0.1$ and $\xi = 0.2$ for Platform 2. For each experiment data we report, we repeat the experiments fifty times and take an average over the experiment results.

As shown in Table 2, on the WS-DREAM dataset, our PDSR method has smaller MAE and RMSE than the reference ones (especially, the BPR, MRP, EF, and MF-based methods) on both the two platforms. This implies that our PDSR method makes QoS predictions more accurately, thanks to the privacy-preserved data sharing across the different platforms. Furthermore, compared with the BRP and MPR methods, our PDSR method takes into account recommendation diversity, and can explore potential services; therefore, it has higher AQoS and higher ILD. Additionally, the AQoS and ILD of our PDSR method are higher than the EF, MF-based and DPP methods where the recommendation diversity is also considered, since the PDSR method not only predicts the missed QoS values more accurately but also enables higher recommendation diversity through data sharing. Although the LSH-based method is extended to enable the data sharing with privacy preserved, it adopts the LSH mechanism directly to leverage the accuracy-diversity trade-off, while our method elaborately defines the accuracy-diversity metric upon the similarity graph as shown in Sec. 4.3. Hence, our method has higher AQoS and ILD than the LSH-based one.

The advantage of our privacy-preserved data sharing in prediction accuracy also can be revealed through the experiment results on the MovieLens dataset. Specifically, the MAE and RMSE of our PDSR method are smaller than the other reference ones on both of the platforms. Another lesson we learn from the experiment results is that, excessively high diversity (i.e., ILD) may result in decreased AQoS. For example, our PDSR method has higher recommendation diversity than the MF-based methods on Platform 1, but the AQoS of our method is slightly smaller than the one of the MF-based method. We will discuss about the trade-off between AQoS and ILD later in Sec. 5.3.3.

We also report the experiment results in terms of time cost. Since the reference methods (e.g., the first five ones) do not consider the data sharing across the two platforms, we omit the communication overhead induced by the data sharing for the purpose of comparability. We conduct our experiments on a desktop equipped with 2.40 GHz CPU and 16 GB RAM. As illustrated in Table 3, our PDSR method has much smaller time cost than the first five reference methods on both the two datasets. As for the LSH-based method, since it directly adopts LSH to measure both similarity and diversity, it results in slightly less time cost. Nevertheless, considering the advantage of our method under another metrics (e.g., MAE, RMSE, AQoS, and ILD as shown in Table 2), this is the price we have to pay.

#### 5.3.2 Fine Tuning Locality-Sensitive Hashing

Recall that, in our method, we index the services using $H_r$ hash functions on each platform $r$ (see **Algorithm** 1) and measure the similarity among the services by calling the subroutine of service indexing $T$ times across the different platforms (see **Algorithm** 2). We hereby reveal how $H_r$ (where $r = 1, 2$ in our case) impacts the performance of our method in terms of QoS prediction. Specifically, let $H_r \in \{3, 4, 5, 6\}$ for $r = 1, 2$ and $T \in \{6, 7, 8, 9, 10\}$. The results are shown in Fig. 3. Since we get similar results for the two datasets on both platforms, we report only the ones obtained with the WS-DREAM dataset on Platform 1.

According to **Algorithm** 1, the number of the hash functions and the number of the hash tables determine the "resolution" at which we characterize the similarity. It is shown in Fig. 3 that, fixing the number of the hash tables, i.e., $T$, when the number of hash functions used in **Algorithm** 1, i.e., $H_r$, is increased, the MAE and RMSE of our method is first decreased but then increased. When $H_r$ is smaller, introducing more hash functions helps us to characterize the service similarity more accurately. However, we may underestimate the service similarity when there are too many hash functions. Specifically, it is more likely for the services to have different hash values in this case, even they are similar to each other. Furthermore, $T$ actually has similar impact on the MAE and RMSE, as shown in **Theorem** 1. By and large, the actual impact of $T$ is slighter than the one of $H_r$, especially when $H_r$ is set such that a larger prediction error is induced (e.g., when $H_r = 3, 4, 6$). Nevertheless, when $H_r = 5$, we are still able to subtly reduce the MAE and RMSE of our prediction by increasing $T$. However, recklessly increasing $T$ finally results in a growth of the prediction error.

#### 5.3.3 An Accuracy-Diversity Trade-off

In our PDSR method, we leverage $\lambda$ to make a trade-off between accuracy and diversity. We hereby report only the AQoS and ILD values obtained on Platform 1 using the WS-DREAM dataset in Fig. 4, since the results with the MovieLens dataset and the ones on Platform 2 are rather similar. We let $\lambda = 0.1, 0.2, 0.3, 0.4$ and vary $\xi = 0.1, 0.2, 0.3, 0.4, 0.5$ to illustrate how the composition of the diversity impact the accuracy-diversity trade-off. It is demonstrated in Fig. 4 (a) that we obtain higher ILD by increasing $\lambda$. As mentioned in Sec. 1, taking into account diversity help us to explore the potential services with higher AQoS. Therefore, as shown in Fig. 4 (b), when $\xi = 0.2$, we obtain higher AQoS by increasing $\lambda$ from 0.1 to 0.2. Nevertheless, pursuing the recommendation diversity excessively may impair the quality of our service recommendation. For example, when $\xi = 0.3, 0.4, 0.5$, increasing $\lambda$ leads to recommendation results with lower AQoS.

TABLE 2
Comparison results in terms of MAE, RMSE, AQoS and ILD.

| Dataset | Methods | Platform 1 | | | | Platform 2 | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | MAE | RMSE | AQoS | ILD | MAE | RMSE | AQoS | ILD |
| WS-DREAM | BPR | 0.9354 | 1.0565 | 0.8945 | 0.4454 | 0.9308 | 1.0028 | 0.9123 | 0.4454 |
| | MRP | 1.486 | 1.8189 | 0.8952 | 0.4455 | 1.1678 | 1.3974 | 0.9142 | 0.4453 |
| | EF | 0.0885 | 0.1301 | 0.9025 | 0.4454 | 0.0748 | 0.1097 | 0.9186 | 0.4454 |
| | MF-based | 0.509 | 0.5235 | 0.9035 | 0.442 | 0.4893 | 0.506 | 0.9184 | 0.4435 |
| | DPP | 0.0783 | 0.1158 | 0.9025 | 0.4454 | 0.0703 | 0.1038 | 0.9089 | 0.4453 |
| | LSH-based | 0.0782 | 0.1158 | 0.8989 | 0.4471 | 0.0703 | 0.1039 | 0.8955 | 0.4477 |
| | PDSR | 0.0782 | 0.1158 | 0.9042 | 0.462 | 0.0703 | 0.1038 | 0.9191 | 0.4603 |
| MovieLens | BPR | 1.7207 | 2.1401 | 3.944 | 0.3843 | 1.8111 | 2.2051 | 3.9676 | 0.3818 |
| | MPR | 1.7553 | 2.2217 | 3.8969 | 0.3872 | 1.7543 | 2.2133 | 3.9298 | 0.3851 |
| | EF | 1.1173 | 1.311 | 3.6382 | 0.4165 | 1.2378 | 1.3332 | 3.75 | 0.4263 |
| | MF-based | 1.9795 | 2.1632 | 3.938 | 0.3883 | 2.1496 | 2.3298 | 3.9188 | 0.3971 |
| | DPP | 0.8748 | 1.044 | 3.856 | 0.4256 | 0.8713 | 1.0379 | 3.8293 | 0.4215 |
| | LSH-based | 0.8638 | 1.0565 | 3.8021 | 0.4369 | 0.8405 | 1.0422 | 3.7989 | 0.4279 |
| | PDSR | 0.8443 | 1.0248 | 3.9 | 0.4324 | 0.8326 | 1.0189 | 3.9138 | 0.4274 |

TABLE 3
Time cost.

| Dataset | Platform 1 | | | | | | | Platform 2 | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | BPR | MPR | EF | MF-based | DPP | LSH-based | PDSR | BPR | MPR | EF | MF-based | DPP | LSH-based | PDSR |
| WS-DREAM | 18.482 | 78.919 | 176.731 | 79.658 | 89.760 | 11.210 | 16.803 | 18.524 | 143.78 4 | 185.252 | 98.274 | 76.344 | 11.268 | 16.825 |
| MovieLens | 84.061 | 308.796 | 183.592 | 637.354 | 18.244 | 6.366 | 10.773 | 126.608 | 599.978 | 344.063 | 1007.899 | 25.178 | 6.2857 | 10.963 |

## 6 CONCLUSION

In this paper, we have proposed a *Privacy-preserving Diversified Service Recommendation* (PDSR) method. The method leverages LSH mechanism to enable data sharing among different data sources (or platforms) with privacy preserved. Specifically, according to the LSH values, we can construct a graph to characterize the similarity among the services, based on which we predict missing QoS values accurately. Furthermore, we have designed a novel diversity metric and proposed a 2-approximation algorithm to select $K$ services by maximizing the objective of accuracy-diversity trade-off. We have performed extensive experiments on real datasets to verify the efficacy of our PDSR method.

Whereas our PDSR method currently exploit diverse services with diverse qualities to fulfill users' different demands, we are on the way of extending our definition of diversity by taking into account more attribute of the services. For example, we can consider the functional features of the services such that the services with different functional features have higher diversity measurements. Moreover, the users may be interested in the services which provide a combination of specific functions. In this case, to fulfill the users' requirement on service functions, we need to recommend multiple sets of services such that each of the sets can provide the specific functions but has diversified QoS measurement.

## REFERENCES

[1] S. Ghafouri, S. Hashemi, and P. Hung. A Survey on Web Service QoS Prediction Methods. *IEEE Trans. on Service Computing*, 15(4):2439–2454, 2022.

[2] G. Linden, B. Smith, and J. Jeremy. Amazon. com Recommendations: Item-to-Item Collaborative Filtering. *IEEE Internet Computing*, 7(1):76–80, 2003.

[3] J. El Hadad, M. Manouvrier, and M. Rukoz. TQoS: Transactional and QoS-aware Selection Algorithm for Automatic Web Service Composition. *IEEE Trans. on Services Computing*, 3(1):73–85, 2010.

[4] Y. Ma, S. Wang, P. Hung, C. Hsu, Q. Sun, and F. Yang. A Highly Accurate Prediction Algorithm for Unknown Web Service QoS Values. *IEEE Trans. on Services Computing*, 9(4):511–523, 2015.

[5] Y. Zhang, L. Wu, Q. He, F. Chen, S. Deng, and Y. Yang. Diversified Quality Centric Service Recommendation. In *Proc. of IEEE ICWS*, pages 126–133, 2019.

[6] G. Kang, J. Liu, B. Cao, and Y. Xiao. Diversified QoS-Centric Service Recommendation for Uncertain QoS Preferences. In *Proc. of IEEE SCC*, pages 288–295, 2020.

[7] Y. Zheng, C. Gao, L. Chen, D. Jin, and Y. Li. DGCN: Diversified Recommendation with Graph Convolutional Networks. In *Proc. of the 30th WWW*, pages 401–412, 2021.

[8] G. Kang, M. Tang, J. Liu, X. Liu, and B. Cao. Diversifying Web Service Recommendation Results via Exploring Service Usage History. *IEEE Trans. on Service Computing*, 9(4):566–579, 2016.

[9] R. Ye, Y. Hou, T. Lei, Y. Zhang, Q. Zhang, J. Guo, H. Wu, and H. Luo. Dynamic Graph Construction for Improving Diversity of Recommendation. In *Proc. of the 15th ACM RecSys*, pages 651–655, 2021.

[10] Y. Liu, Y. Xiao, Q. Wu, C. Miao, J. Zhang, B. Zhao, and H. Tang. Diversified Interactive Recommendation with Implicit Feedback. In *Proc. of the 34th AAAI*, pages 4932–4939, 2020.

[11] S. Wu, H. Kou, C. Lv, W. Huang, L. Qi, and H. Wang. Service Recommendation with High Accuracy and Diversity. *Wireless Communications and Mobile Computing*, 2020:1–10, 2020.

[12] L. Qi, X. Zhang, W. Dou, and Q. Ni. A Distributed Locality-Sensitive Hashing-Based Approach for Cloud Service Recommendation From Multi-Source Data. *IEEE Journal on Selected Areas in Communications*, 35(11):2616–2624, 2017.

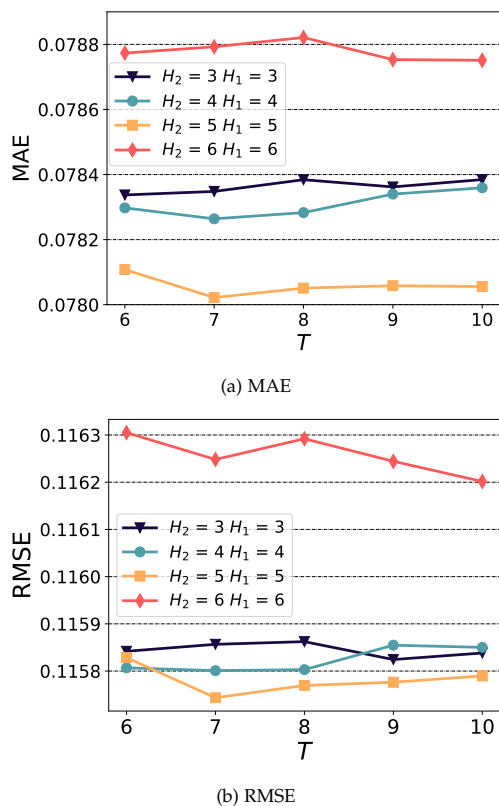[13] J. Zhu, P. He, Z. Zheng, and M. Lyu. A Privacy-Preserving QoS

(a) MAE



(b) RMSE

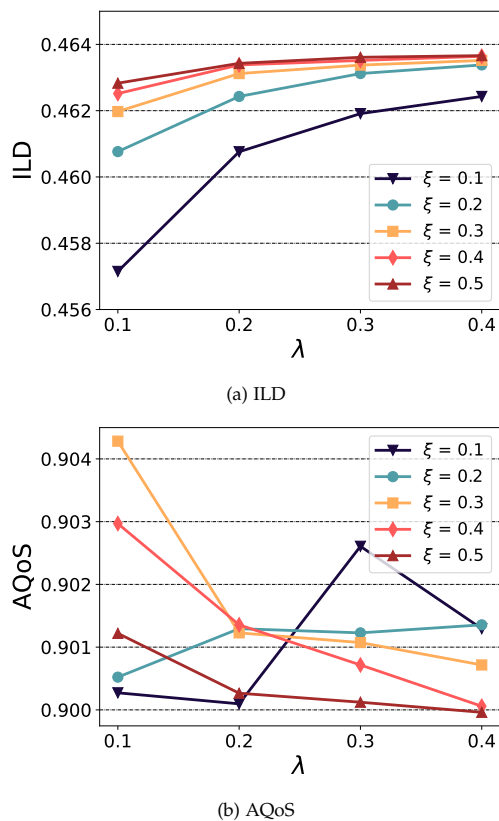Fig. 3. Prediction accuracy with different values of $H$ and $T$.



(a) ILD



(b) AQoS

Fig. 4. Trade-off between accuracy and diversity.

Prediction Framework for Web Service Recommendation. In *Proc. of IEEE ICWS*, pages 241–248, 2015.

[14] L. Qi, H. Xiang, W. Dou, C. Yang, Y. Qin, and X. Zhang. Privacy-Preserving Distributed Service Recommendation Based on Locality-Sensitive Hashing. In *Proc. of IEEE ICWS*, pages 49–56, 2017.

[15] J. Canny. Collaborative Filtering with Privacy. In *Proc of the 23rd IEEE S&P*, pages 45–57, 2002.

[16] S. Jumonji, K. Sakai, M. Sun, and W. Ku. Privacy-Preserving Collaborative Filtering Using Fully Homomorphic Encryption. In *Proc of the 38th IEEE ICDE*, pages 1551–1552, 2022.

[17] T. Zhu, G. Li, Y. Ren, W. Zhou, and P. Xiong. Differential Privacy for Neighborhood-based Collaborative Filtering. In *Proc. of ASONAM*, pages 752–759, 2013.

[18] C. Gao, C. Huang, D. Lin, D. Jin, and Y. Li. DPLCF: Differentially Private Local Collaborative Filtering. In *Proc of the 43rd ACM SIGIR*, pages 961–970, 2020.

[19] R. Chow, M. Pathak, and C. Wang. A Practical System for Privacy-Preserving Collaborative Filtering. In *Proc. of the 12th IEEE ICDM Workshops*, pages 547–554, 2012.

[20] A. Ahgasaryan, M. Bouzid, D. Kostadinov, M. Kothari, and A. Nandi. On the Use of LSH for Privacy Preserving Personalization. In *Proc. of the 12th IEEE TrustCom*, pages 362–371, 2013.

[21] Y. Zhong, Y. Fan, K. Huang, W. Tian, and J. Zhang. Time-Aware Service Recommendation for Mashup Creation in an Evolving Service Ecosystem. In *Proc. of the 21st IEEE ICWS*, pages 25–32, 2014.

[22] Y. Zhang, K. Wang, Q. He, F. Chen, S. Deng, Z. Zheng, and Y. Yang. Covering-Based Web Service Quality Prediction via Neighborhood-Aware Matrix Factorization. *IEEE Trans. on Services Computing*, 14(5):1333–1344, 2019.

[23] Y. Li, S. Wang, Q. Pan, H. Peng, T. Yang, and E. Cambria. Learning Binary Codes with Neural Collaborative Filtering for Efficient Recommendation Systems. *Knowledge-Based Systems*, 172:64–75, 2019.

[24] Z. Cui, X. Xu, F. Xue, X. Cai, Y. Cao, W. Zhang, and J. Chen. Personalized Recommendation System Based on Collaborative Filtering for IoT Scenarios. *IEEE Trans. on Services Computing*, 13(4):685–695, 2020.

[25] L. Chen, G. Zhang, and E. Zhou. Fast Greedy MAP Inference for Determinantal Point Process to Improve Recommendation Diversity. In *Proc. of the 32nd NeurIPS*, pages 5627–5638, 2018.

[26] R. Latha and R. Nadarajan. Analysing Exposure Diversity in Collaborative Recommender Systems—Entropy Fusion Approach. *Physica A: Statistical Mechanics and Its Applications*, 533:122052, 2019.

[27] W. Dou, X. Zhang, J. Liu, and J. Chen. HireSome-II: Towards Privacy-Aware Cross-Cloud Service Composition for Big Data Applications. *IEEE Trans. on Parallel and Distributed Systems*, 26(2):455–466, 2013.

[28] C. Li, B. Palanisamy, and J. Joshi. Differentially Private Trajectory Analysis for Points-of-Interest Recommendation. In *Proc. of IEEE BigData Congress*, pages 49–56, 2017.

[29] W. Zhong, X. Yin, X. Zhang, S. Li, W. Dou, R. Wang, and L. Qi. Multi-Dimensional Quality-Driven Service Recommendation with Privacy-Preservation in Mobile Edge Environment. *Computer Communications*, 157:116–123, 2020.

[30] L. Wang, X. Zhang, R. Wang, C. Yan, H. Kou, and L. Qi. Diversified Service Recommendation with High Accuracy and Efficiency. *Knowledge-Based Systems*, 204:106196, 2020.

[31] A. Gionis, P. Indyk, and R. Motwani. Similarity Search in High Dimensions via Hashing. In *Proc. of the 25th VLDB*, pages 518–529, 1999.

[32] M. Charikar. Similarity Estimation Techniques from Rounding Algorithms. In *Proc. of the 34th ACM STOC*, pages 380–388, 2002.

[33] R. Hassin, S. Rubinstein, and A. Tamir. Approximation Algorithms for Maximum Dispersion. *Operations Research Letters*, 21(3):133–137, 1997.

[34] WS-DREAM: Towards Open Datasets and Source Code for Web Service Research. https://wsdream.github.io/, 2017. [Online].

[35] MovieLens 25M Dataset. https://grouplens.org/datasets/movielens/25m/, 1997. [Online].

[36] S. Rendle, C. Freudenthaler, Z. Gantner, and L. Schmidt-Thieme. BPR: Bayesian Personalized Ranking from Implicit Feedback. In *Proc. of the 25th UAI*, page 452–461, 2009.

[37] R. Yu, Y. Zhang, Y. Ye, L. Wu, C. Wang, Q. Liu, and E. Chen. Multiple Pairwise Ranking with Implicit Feedback. In *Proceedings of the 27th ACM CIKM*, pages 1727–1730, 2018.

[38] F. Wang, L. Wang, G. Li, Y. Wang, C. Lv, and L. Qi. Edge-Cloud-Enabled Matrix Factorization for Diversified APIs Recommendation in Mashup Creation. *World Wide Web*, page 1809–1829, 2022.