

Recommending Packages to Groups

Shuyao Qi
University of Hong Kong
qisy@connect.hku.hk

Nikos Mamoulis
University of Ioannina
nikos@cs.uoi.gr

Evaggelia Pitoura
University of Ioannina
pitoura@cs.uoi.gr

Panayiotis Tsaparas
University of Ioannina
tsap@cs.uoi.gr

Abstract—The success of recommender systems has made them the focus of a massive research effort in both industry and academia. Recent work has generalized recommendations to suggest packages of items to single users, or single items to groups of users. However, to the best of our knowledge, the interesting problem of recommending a *package to a group of users* (P2G) has not been studied to date. This is a problem with several practical applications, such as recommending vacation packages to tourist groups, entertainment packages to groups of friends, or sets of courses to groups of students. In this paper, we formulate the P2G problem, and we propose probabilistic models that capture the preference of a group towards a package, incorporating factors such as user impact and package viability. We also investigate the issue of recommendation *fairness*. This is a novel consideration that arises in our setting, where we require that no user is consistently slighted by the item selection in the package. We present aggregation algorithms for finding the best packages and compare our suggested models with baseline approaches stemming from previous work. The results show that our models find packages of high quality which consider all special requirements of P2G recommendation.

I. INTRODUCTION

Consider a group of people who would like to dine at a restaurant and then have drinks at a nearby bar. Given the potentially numerous options, the group would favor a recommendation of a (restaurant, bar) pair, which is consistent with the preferences of its members and does not make a member unhappy with respect to the rest of the group.

Despite the vast amount of work on recommender systems, to the best of our knowledge, this *package-to-group* (P2G) recommendation problem has not been studied before, although there is work on recommending a package of items to a single user (e.g., [3], [20], [7]) and recommending a single item to a group of users (e.g., [2], [14]). In addition, there are studies on helping a group of users to select a bundle of items (e.g., [17], [18]). However, they assume that the users are given a set of items and together they decide the items to select, which is a different problem from P2G recommendation.

Specifically, given a group of users U , the goal of P2G recommendation is to suggest one or more packages of items to U , which are suitable for U 's members. This problem has several applications beyond the night-out scenario. For example: (i) A tour operator wants to create a package for a group of tourists, consisting of hotels, restaurants, attractions and activities; (ii) An academic institution that organizes a summer school wants to create a curriculum that meets the interests of a group of students; (iii) A movie channel wants

to package together a set of movies to offer to a group of movie-goers, or a large family.

In line with previous work on package recommendation [19], [12], we assume the existence of constraints limiting possible item combinations that can be included in a package. Constraints may be either *hard* or *soft*. Hard constraints should definitely be satisfied by a set of items in order for it to form a valid package. Soft constraints express desirable, but not necessary properties for an item set. In terms of hard constraints, without loss of generality, we focus on the special, but practical case, where the items are divided into categories, and a valid package is formed by selecting one item from each category. For instance, in the night-out example, the group may be interested in a package which includes a restaurant and a bar (i.e., one item from category “restaurants” and one item from category “bars”); the tourist group visit a city and they are interested in visiting a museum, dining at a restaurant, and finally resting at a good hotel; the summer school may consist of courses covering different areas (e.g., Theory, Systems); the movie-goers may want to watch a thriller and a comedy. The soft constraints we consider in the paper are defined based on the relationships between the items in a package. For example, in a venue-package recommendation problem, a set of places far from each other is less likely to be selected by a group, compared to a package of nearby places. In this case, we say that the package is less *viable*. The assumption of constraints is not compulsory and is independent of our proposals, as we show in Sections III-D and VI-E. For instance, the category constraint can be easily replaced with selecting a number of items regardless of their categories, by virtually assuming that all items belong to a single large category.

Based on the above, we present two probabilistic models for P2G recommendation: one that first computes the probabilities that the group of users likes individual items, before deriving the probability that the group would select a package of items, and one which first forms item packages that are favored by the individual group members before identifying those that have high likelihood to be selected by the group. Our experimental results show that the first model is superior because it seamlessly takes into consideration all special factors of P2G recommendation (e.g., user impact, package viability). In addition, we design and implement algorithms for the models on a database of individual user ratings on items. The algorithms efficiently combine items into candidate packages for recommendation, while avoiding the exploration of the entire search space with the help of pruning bounds.

A unique and novel characteristic of P2G recommendation is that it raises the issue of *fairness*. User groups may be heterogeneous, consisting of people with potentially dissimilar tastes. Thus, for a package I that is overall good for the group (i.e., the average group member preferences on its items are high), there could be one or more members that do not like any of the items in I ; these users would be frustrated if I is selected by the group. In this case, we consider the package to be *unfair*. On the other hand, if each group member finds at least one item in the package that she likes, we consider such a package to be *fair*. We formalize fairness for P2G recommendation, inspired by the corresponding concept in fair division of resources [16] and adapt our models accordingly. Note that our fairness definition is very different from that in group decision making [17], where it is assumed that the group would do item selection for multiple times and to be fair, the unsatisfied users will have higher priority in the next decision. On the contrary, the fairness problem in P2G recommendation is one time and defined on item basis.

Our contributions in this paper are the following: (i) This is the first work that formulates and studies P2G recommendation. (ii) We propose probabilistic models that incorporate factors such as user impact, package viability, and fairness. (iii) We design efficient P2G recommendation algorithms that scale for large data. In the rest of the paper, we first formally define P2G recommendation in Section II. Section III presents our two probabilistic models and introduces package viability. In Section IV, we define fairness and show how it can be integrated into the models. Section V presents algorithms that efficiently implement the models. Section VI presents our experimental evaluation. Related work is reviewed in Section VII and the paper concludes in Section VIII.

II. PROBLEM STATEMENT

We assume a collection \mathcal{I} of items and a collection \mathcal{U} of users, who express their preferences to items from \mathcal{I} through ratings. A rating $r(u, i)$ of user u for item i may be *explicit*, i.e., u has used and evaluated item i , or *implicit*, i.e., predicted by a classic recommender (e.g., collaborative filtering [15]).

Given a *group* (set) U of users in \mathcal{U} , we consider recommending to U a *package* (set) I of items in \mathcal{I} . Recommended packages must be *valid*, i.e., have specific properties. In this paper, we study the case where items belong to categories taken from a set \mathcal{C} (e.g., $\mathcal{C} = \{\text{restaurant, bar, theater, museum}\}$). Without loss of generality, we assume that each $i \in \mathcal{I}$ belongs to a single category $c_i \in \mathcal{C}$. The group U inputs a query specifying the set of categories $C \subseteq \mathcal{C}$ where the items of the package should be drawn from (e.g., $C = \{\text{bar, restaurant}\}$). For the ease of discussion, we assume that each item belongs to only one category and a feasible package must contain one item per category (e.g., the users want to visit one bar and one restaurant). More general problem instances will be elaborated in Section III-D.

Formally, a P2G recommendation task takes as input a group of users $U \subseteq \mathcal{U}$, a set of ratings, and a set of user-specified item categories $C \subseteq \mathcal{C}$, and recommends to U the k

most preferable among all feasible packages. We now present generic probabilistic models which define the preference of a group U over a package I .

III. PROBABILISTIC MODELS

Given a target group U and a query input by U specifying category set C , the objective is to derive the probability distribution $\Pr(I|U, C)$ of the group U to select the package I over C . The probability $\Pr(I|U, C)$ obviously depends on the preference of each user $u \in U$ for the individual items $i \in I$. Given a $u \in U$ and an item i from a specific category $c_i \in C$, the probability of u independently selecting i over other items in c_i can be defined as

$$\Pr(i|u, c_i) = \frac{r(u, i)}{\sum_{i' \in c_i} r(u, i')} \quad (1)$$

Here $r(u, i)$ is u 's (explicit or implicit) rating on i . Note that a $\Pr(i|u, c_i)$ is defined for every category c_i . Intuitively, u is more likely to accept a recommendation $i \in c_i$ with higher $r(u, i)$ compared to $r(u, i')$ for other items $i' \in c_i$. Next, we present two models for computing $\Pr(I|U, C)$ based on $\Pr(i|u, c_i)$ and other factors, such as the influence between users in the group, and the likelihood that a set of items are appealing together as a package.

A. Group Rating (GR) Model

In the *group rating* (GR) model, we first define the probability that the group U will select an item i . Then we combine the probabilities of individual items, to derive the likelihood of a package.

Item to Group (I2G) Probability. Given group U and a category c_i , the probability of U selecting $i \in c_i$ is $\Pr(i|U, c_i)$. Here $\Pr(i|U, C) = \Pr(i|U, c_i)$, that is, the probability of item i being selected depends only on its own category and not in the full set of categories C . The above probability can be computed based solely on the probabilities of the users in U selecting the item (e.g., see [5]). In our model we adopt the approach in [10], [21], where it is assumed that different group members may have different impact on the group's decision. In simple words, one or more members of the group, who could be considered as experts on a category, may influence the group in selecting an item in this category. For example, the preference of a group member who is a "foodie" will count more in selecting a specific restaurant.

Following this intuition, we model the group selection as a stochastic process where a user u is first selected as the representative of the group with probability $\Pr(u|U, c_i)$, and then the group selects an item according to u 's item distribution. Therefore, we have:

$$\Pr(i|U, c_i) = \sum_{u \in U} \Pr(u|U, c_i) \Pr(i|u, c_i) \quad (2)$$

In this work, we assume that the probability $\Pr(u|U, c_i)$ of a user $u \in U$ is proportional to the activity of the user in category c_i , relative to the other members in the group. This captures the relative *expertise* of the user in the group for

this category, which determines her influence in the group. Specifically, let η_{u,c_i} denote the number of explicit ratings user u has given for items in category c_i . We have that

$$\Pr(u|U, c_i) = \frac{\eta_{u,c_i}}{\sum_{u' \in U} \eta_{u',c_i}}$$

Note that Equation (2) is general enough to model different scenarios, depending on the definition of the probability $\Pr(u|U, c_i)$. For example, we can set $\Pr(u|U, c_i)$ to the uniform distribution, where all users influence equally the final selection. Or, we may assume that user influence is independent of the category as $\Pr(u|U)$. In fact, we also considered an approach similar to that in [21], where we used topic-modeling to extract the user-topic and item-topic distributions, and then defined the user influence probability based on the user-item distribution. Experimentally, this approach gave us similar recommendation results on our test data, because the categorization of items is correlated to their underlying topics.

Package to Group (P2G) Probability. To derive the probability $\Pr(I|U, C)$ of the package I to be selected by the group U , for the moment we assume that items are selected independently. Therefore, given $\Pr(i|U, c_i)$, we have:

$$\Pr(I|U, C) = \prod_{i \in I} \Pr(i|U, c_i) \quad (3)$$

B. User Package (UP) Model

The GR model assumes that items are selected independently, according to the preferences of a representative user, who is chosen according to her expertise and influence in the item category. The *user package* (UP) model reverses the above generative process. In UP, the group first chooses a representative user u with probability $\Pr(u|U, C)$. The representative user will decide for the *whole package*. We assume that the representative user selects each item independently for now, according to her own preferences. Formally:

$$\begin{aligned} \Pr(I|U, C) &= \sum_{u \in U} \Pr(u|U, C) \Pr(I|u, C) \\ &= \sum_{u \in U} \left\{ \Pr(u|U, C) \prod_{i \in I} \Pr(i|u, c_i) \right\} \quad (4) \end{aligned}$$

Accordingly, once the group has selected a representative u , the selection probability for the package depends only on u , i.e. $\Pr(I|u, U) = \Pr(I|u)$. Also, $\Pr(i|u, C) = \Pr(i|u, c_i)$ similar to GR. We can again adjust the probability $\Pr(u|U, C)$ to model different scenarios. Different from GR, however, UP considers the user impact on packages instead of items. Therefore, $\Pr(u|U, C)$ is defined based on the influence of u on all target categories C collectively.

$$\Pr(u|U, C) = \frac{\sum_{c \in C} \eta_{u,c}}{\sum_{u' \in U} \sum_{c \in C} \eta_{u',c}}$$

where as before $\eta_{u,c}$ is the user u 's overall influence on category c .

Essentially, the selection is a two-step process: The package to user (P2U) phase computes the probability $\Pr(I|u, C)$ that

user u selects package I , for all users in U ; the package to group (P2G) phase computes the overall preference probability $\Pr(I|U, C)$ of the group by taking the combination of the user preferences weighted by the user impact probabilities.

Note that the UP model gives more power to the user selected as representative, since the package selection disregards other group members. As a result, GR and UP may produce very different package selection probabilities, even in the case of uniform user impact probabilities. Consider the example in Table I(a), where a group of two users $U = \{u_1, u_2\}$ wants to select a package over two categories $C = \{X, Y\}$, each having two items, (x_1, x_2) and (y_1, y_2) respectively, with the preference probabilities shown in the table. Assume that $\Pr(u_1|U) = \Pr(u_2|U) = 1/2$, in all categories. Table I(b) shows the probabilities of GR and UP for each possible package.

		u_1	u_2	Packages		GR	UP
X	x_1	1	0	$I_1(x_1, y_1)$	1/4	1/2	
	x_2	0	1	$I_2(x_1, y_2)$	1/4	0	
Y	y_1	1	0	$I_3(x_2, y_1)$	1/4	0	
	y_2	0	1	$I_4(x_2, y_2)$	1/4	1/2	

(a) $\Pr(i|u, c_i)$

(b) $\Pr(I|U, C)$

TABLE I
COMPARISON OF GR AND UP

The example shows that, for the UP model, a package that no user likes as a whole (I_2 and I_3) will have very low (zero) probability, while the packages with high probability are actually favored by a single user (I_1 by u_1 and I_4 by u_2). On the other hand, in the GR model, a package becomes acceptable as long as there is at least one user that likes some item in the package (e.g., u_1 likes x_1 and u_2 likes y_2), balancing the preferences of the users better.

Overall, the UP model has the following drawbacks: (1) It assumes each user selects the package as a whole, so that the users' impact on different categories cannot be evaluated; (2) For the same reason, a user will never select a low rated item by her, that is, a user will never compromise for the sake of the group; (3) The top packages for different users may not overlap, especially for dissimilar users, leaving the group dissatisfied as a whole. We therefore expect UP to produce worse packages than GR in practice.

C. Package Viability

So far, we have considered only the preferences of the users over individual items (Equation (3) and section III-B), assuming independence between items. However, in real-life, some items are more likely to be selected together than others. For example, a restaurant and a movie theater have higher chances to form a preferable package if they are spatially close. Motivated by this fact, we define the probability $\Pr(V|I)$ that a package I is *viable as a whole*. One possible evaluation of $\Pr(V|I)$ is to consider the pairwise *relevance* between items in I .¹ Here, V denotes a random variable, which is 1 if the

¹In general, for a set of n items, the viability can be defined by aggregating their pairwise relevance or by defining an n -ary function. In this paper, for simplicity and due to the application domain of our case studies in the experiments, we follow the first approach.

package is viable and 0 otherwise. The relevance between two items can be derived by a function on their features (e.g., their spatial distance), or by recorded statistics (e.g., joint probability). Take the case of recommending a package of places as an example. If we regard the relevance between any pair of items to be inversely proportional to their distance (measured by a function $dist(\cdot)$), we can define $\Pr(V|I)$ as:

$$\Pr(V|I) \propto e^{-\max_{i,i' \in I} \{dist(i,i')\}} \quad (5)$$

Intuitively, if the maximum distance between any pair of places in a package is large, the package has low probability to be appealing. There can be other measurements of $\Pr(V|I)$ as well. For example, we can consider the visiting order of the items in I , or relate the viability to traveling time cost instead of distance between items; we can also define $\Pr(V|I)$ based on the historical probability where items in I are selected together [22]. Our models are independent of the specific definition of $\Pr(V|I)$ and for ease of discussion, we use Equation (5) as an exemplary viability definition in this paper.

Let us now formalize the probability that a group U will select a package I and the package is viable. Assuming that viability depends only on the package, we have:

$$\begin{aligned} \Pr(I, V|U, C) &= \Pr(V|I, U, C) \Pr(I|U, C) \\ &= \Pr(V|I) \Pr(I|U, C) \end{aligned} \quad (6)$$

In the rest of the paper, both the GR and UP models are augmented with package viability according to Section III-C.

D. Generality

So far we have assumed that (1) each item i belongs to a single category and (2) only one item is recommended from each category. In a real-life scenario, these assumptions may be too restrictive. Our models can be easily adapted to apply to more general cases.

Firstly, suppose that an item may belong to multiple categories, e.g., a place is regarded both a restaurant and a bar. If the group U accepts a duplicate item serving different purposes, then the models do not require any adaptation; an item may appear in the recommended package multiple times, from different categories. If, on the contrary, U would not accept any duplicate item in a package, the models can still work with a minor adaptation that filters out packages containing duplicate items. For example, given a package I , its viability probability $\Pr(V|I)$ is set to 0 if an item appears more than once in I .

Next, suppose that the group U is looking for multiple items in one category, or simply looking for items without any category constraint. Without loss of generality, we assume that U wants to find n items in category c_i (or C if U sets no category constraints). In this case, we virtually replicate c_i n times and apply the same models on categories set $C' = \{c'_{i1}, \dots, c'_{in}\}$. As a result, n items will be selected from category c_i . However, since it becomes possible to select an item from c_i multiple times, the aforementioned filtering method should be applied to avoid selecting duplicate

items. The above strategy also extends to the generic case of recommending arbitrary number (0 to $|c_i|$) of items from multiple categories. In Section VI-E, we show experimentally the performance of our models without category constraint.

IV. FAIRNESS IN RECOMMENDATIONS

Both GR and UP find the top packages without considering which users are the most or least happy with the items in the packages. For a selected package I , it is possible that a given user $u \in U$ does not like *any* of the items in I , or that u is the least satisfied user in U for *all* items in the package. Therefore, although U as a whole may like package I , the package selection is not *fair* to user u . In a real-life scenario, where the users in the group care for each other's preferences, we should recommend a package which is both attractive and fair to the majority of the group members.

For a user u and a package I , we say that I is fair to u , if there exists at least one item $i \in I$, such that u 's rating on i is ranked in the top- $\Delta\%$ of u 's ratings on all items. The rationale is that the existence of at least one item in the package for which u has high rating would make the user tolerant to the existence of other items that she may not prefer, considering that there are other members in the group who may like these items. Given the group U and a package I , we denote by $U_f \subseteq U$ the users to whom I is fair. A *fairness* measure $fair(U, I)$ is accordingly defined:

$$fair(U, I) = \frac{|U_f|}{|U|}, \quad (7)$$

meaning that the more users I is fair to, the better I is for U .

Lastly, we define the fairness-aware score of a package as

$$score_{fair}(U, I) = \Pr(I, V|U, C) \cdot fair(U, I), \quad (8)$$

i.e., we look for packages that are both highly preferable and fair. Note that the above equation is applicable to both GR and UP models. It scores a package I based on both its relevance to the group members U (according to GR or UP), and its fairness to U . In the rest of the paper, we denote the GR and UP package selection models augmented with fairness as GR-Fair and UP-Fair, respectively.

Fairness is inspired by the classic fair division problem in Economics [16]. Fair division splits one or more heterogeneous resources to a number of people who have different preferences to different parts of the resources, such that everybody believes that they have a fair share. Our P2G selection problem is reminiscent to fair division, because every user in the group has different preferences in the items. However, in P2G, the group members share the items in the suggested package, instead of the items being divided.

V. ALGORITHMS

Given a group U of users in \mathcal{U} , a set of categories C , a database of items \mathcal{I} and the user ratings over the items, our goal is to find the top- k packages that maximize $score(U, I)$ according to Equation (8), following either model GR or UP. An efficient implementation is critical because the number of

candidate packages is exponential to the number of categories. We now present efficient branch-and-bound algorithms for ranking packages based on GR and UP.

A. Algorithms for GR

Recall that GR includes two phases: the I2G phase which finds in each category the probability $\Pr(i|U, c_i)$ of each item being selected, and the P2G phase which combines items into packages. The final scoring function (Equation (8)) considers three factors in the P2G phase, (1) the group preference $\Pr(i|U, c_i)$ (Equation (2)), (2) the package viability $\Pr(V|I)$ (Section III-C), and (3) fairness (Equation (8)). As a result, combining the best items found in the I2G phase into packages does not necessarily lead to the best packages.

A baseline algorithm (GR-BA) would first calculate the I2G probability for each item relevant to U , then consider each possible package by calculating $score(U, I)$ and finally select the top- k packages accordingly. Once there are at least k package candidates, a lower bound θ of the current k th maximum probability can be used to prune any package I for which $score(U, I) \leq \theta$. In addition, we can prune *partial* packages, which contain items only from some categories, before they are expanded to complete packages, if even by including into them the best possible items for each remaining category, their scores cannot exceed θ . Even with these optimizations, GR-BA does not give priority to the most preferred items by the users and thus considers most of the possible packages.

As an alternative to GR-BA, we propose a 2-level incremental algorithm GR-INC, which prioritizes items and packages with respect to their potential probability of being selected and computes the I2G and P2G phases concurrently. In particular, the I2G phase is implemented as an (incremental) top- k selection query [4], which generates for each category a list of its items in decreasing probability order of being selected by the group U , according to Equation (2). The I2G phase takes as input $|U|$ sorted lists of item ratings, one per user in U ; each list includes only the items in one of the input categories c_i . The P2G phase is implemented as an (incremental) top- k join query [6] where viability is considered in the aggregation score of the joined item combinations. P2G takes as input the items output by the I2G phase on each category and combines them. Algorithm 1 is a pseudo code of GR-INC, using Procedure 2 as a module, which implements the I2G phase. Figure 1 illustrates an example where there is a group U with three users u_1 - u_3 looking for recommendations from categories bar $c_b(b_1, \dots, b_5)$ and restaurant $c_r(r_1, \dots, r_5)$, assuming that a package $I(b_i, r_j)$ is always viable (i.e., $\Pr(V|I) = 1$). For simplicity, in Algorithm 1, we assume that the packages are ranked and selected in decreasing order of $\Pr(I, V|U, C)$ (not $score(U, I)$). Its adaptation to a GR-Fair algorithm (i.e., find the top packages considering fairness) is straightforward. GR-Fair ranks the packages by $score_{fair}(U, I)$ based on Equation (8) (Line 13) and sets θ_I as the k th maximum score in R (Line 16). Based on Equation (8), we can use a tighter bound $T_{I, fair} = \prod_{c \in C} ub_c \cdot ub_{fair}$ (replacing Line 17) where ub_{fair} is the maximum fairness degree of unseen packages.

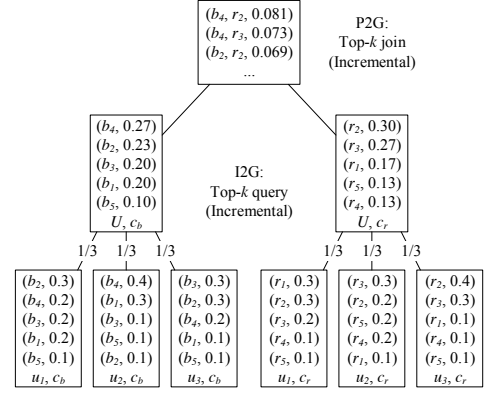


Fig. 1. GR-INC

ALGORITHM 1: Incremental Algorithm for GR (GR-INC)

Input : U, C, k
Output: R

- 1 min-heap $R \leftarrow \emptyset$
- 2 **for each** $c \in C$ **do**
- 3 initialize a max-heap $H_c \leftarrow \emptyset$
- 4 initialize a buffer $B_c \leftarrow \emptyset$
- 5 $ub_c = \infty$
- 6 $\theta_I = -\infty, T_I = \infty$
- 7 **while** $T_I > \theta_I$ **do**
- 8 $c =$ select the next category
- 9 $L_c =$ GR-INC-I2G(U, c, k, H_c)
- 10 **for each item** $i \in L_c$ **do**
- 11 insert i into B_c
- 12 **for each package** I with i and items from $B_{c'}, c' \neq c$ **do**
- 13 calculate $\Pr(I, V|U, C)$ // Section III-C
- 14 insert I into R and pop from R if $|R| > k$
- 15 $ub_c = \min_{i \in L_c} \Pr(i|U, c_i)$
- 16 $\theta_I = k$ th largest probability in R
- 17 $T_I = \prod_{c \in C} ub_c$
- 18 **return** R

ub_{fair} is initially 1, and is decreased by $1/|U|$ if a user u exhausts all her top- $\Delta\%$ items, as in this case none of the unseen items could be fair to u .

B. Algorithms for UP

We can design algorithms for the UP (and UP-Fair) model in a similar manner as for the GR model. The baseline UP-BA algorithm, in the first step, finds for each user $u \in U$ the relevant packages \mathcal{I}_u and their probabilities of being selected by u , i.e., $\Pr(I|u, C)$. Then, UP-BA ranks the packages in $\cup_{u \in U} \mathcal{I}_u$. UP-BA can be improved by ordering each \mathcal{I}_u by $\Pr(I|u, C)$, so that the best packages can be found in a top- k query fashion. Similar to GR-INC, we propose an efficient UP-INC algorithm, which follows a 2-level procedure prioritizing items and packages of high probability to be selected by each user. The P2U phase of UP is implemented as an (incremental) viability top- k join query that computes for each user, packages in decreasing degree of being liked

PROCEDURE 2: GR-INC-I2G

Input : U, c, k, H_c
Output: L_c, H_c

- 1 $\theta_i = -\infty, T_i = \infty$
- 2 $ub_u = \infty$ for each $u \in U$
- 3 **while** $T_i > \theta_i$ **do**
- 4 **for each** $u \in U$ **do**
- 5 access the next item $i \in c$ rated by u
- 6 calculate $\Pr(i|U, c)$ // Equation (2)
- 7 insert i into H_c
- 8 $ub_u = \Pr(i|u, c)$
- 9 $\theta_i = k$ th largest I2G probability in H_c
- 10 $T_i = \sum_{u \in U} \{\Pr(u|U, c) \cdot ub_u\}$ // Equation (2)
- 11 move the top- k items in H_c to L_c
- 12 **return** L_c

by her. On top of that, the P2G phase is implemented as an (incremental) top- k selection query where the packages being liked by the group as a whole are progressively selected. Figure 2 presents an example for UP with the same data and setup as in Figure 1.

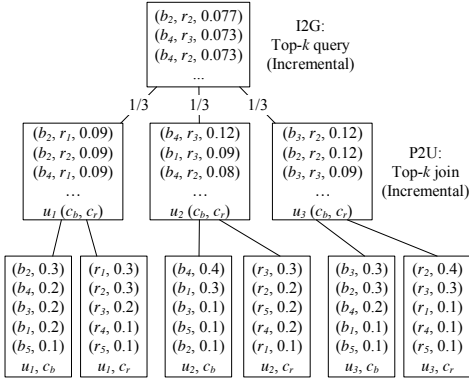


Fig. 2. UP-INC

VI. EXPERIMENTAL EVALUATION

This section evaluates our P2G models and algorithms. Section VI-A details the setup of our analysis. Section VI-B studies the effectiveness of the proposed models. Section VI-C evaluates the effect of considering fairness in the models and presents a user study. Finally, Section VI-D tests the efficiency of our algorithms.

A. Setup

We use two real datasets: Yelp² and MovieLens³ in our evaluation. For Yelp, we use as items venues from the city of Phoenix and consider five categories (restaurants, shopping, beauty & spa, health & medicine, nightlife) with the most venues. Yelp originally contains about 100K users, 17K places and 476K reviews with a numerical rating. Because the number of reviews is small, we employ collaborative filtering (CF) [15] to get additional review ratings for each user. In particular,

²http://www.yelp.com/dataset_challenge

³<http://grouplens.org/datasets/movielens/>

we use Mahout⁴ to build an item-based CF recommender and retrieve for each user u all item ratings that are not present in the dataset. For the items that are neither explicitly rated by u in the dataset nor recommended by CF, we set zero as u 's rating. Finally, we end up having 53M non-zero ratings in total. For MovieLens, we use movies as items from the five most popular genres (drama, comedy, thriller, romance, action), which contain about 138K users, 33K movies and 31M reviews. The same CF recommendation process results in 51M ratings in total. To prevent bias toward any user, in both datasets, we normalize the ratings of every user to $[0, 1]$. All algorithms were implemented in C++ and the tests ran on a machine with Intel Core i7-3770 3.40GHz and 16GB main memory, running Ubuntu Linux.

TABLE II
PARAMETERS IN EXPERIMENTS (DEFAULT VALUES IN BOLD)

description	parameter	values
Group size	$ U $	1, 2, 3 , 4, 5
Number of categories	$ C $	1, 2, 3 , 4, 5
Fairness threshold (%)	Δ	1, 5, 10 , 50, 100

Table II summarizes all parameters involved in our study. On each test, we vary one parameter, while keeping the others to their default values. Each test computes the top-10 recommended packages to a random group U of users. We consider two classes of user groups. Groups in the SIM class consist of users that have similar preferences to items. Each SIM group is generated by randomly selecting a user and then iteratively picking the next user as the one for which the item preference vector has the maximum cosine similarity to the selected users so far. DSIM user groups are generated in the same way, however, using minimum instead of maximum similarity when selecting the next user to add to the group.

B. Model Evaluation

We study the effectiveness of our proposed GR and UP models, by first focusing on the basic models where fairness is not considered. In the evaluation, we include a baseline approach (BASE) which is based on the state-of-the-art group recommendation technique (COM [21]). For each category $c \in C$, COM is used to select the best item for U . These items are then combined to form the top package. The 2nd-best item of each category is then combined with the best items from the other ones to form additional packages and so on until k packages are computed. BASE aims at maximizing the preferences of the group to the individual items in the suggested packages. Note that the original COM model is designed for the scenario where the topics (categories in our case) are not specified by the group of users and thus need to be inferred from group or user-topic distributions. BASE adapts the COM model for our problem by limiting to one topic for recommendation in each category. Still, BASE ignores the possible relationships between items (see Section III-C); thus, the top items per category selected by BASE do not necessarily form good packages.

⁴<http://mahout.apache.org>

We compare BASE, GR, and UP in terms of package quality using two metrics: the average group-item rating $R(U, I)$ and the average item distance $dist(I)$. $R(U, I)$, indicating how much the members of U like the individual items in I , is the average of group rating $\rho(U, i)$ to each item $i \in I$, weighted by the user impacts:

$$R(U, I) = avg \sum_{i \in I} \rho(U, i) = avg \sum_{i \in I} \sum_{u \in U} \Pr(u|U, c_i) r(u, i)$$

The average distance $dist(I) = avg \sum_{i, i' \in I} dist(i, i')$ between the items in the package I indicates how viable it is for them to be chosen together (i.e., items far from each other could be a bad choice). For Yelp, $dist(i, i')$ is the Euclidean distance between the items (venues). For MovieLens, we define $dist(i, i') = 1 - sim(i, i')$, where sim is the similarity between movies i and i' , calculated via the Movie-Topic matrix extracted using Latent Dirichlet Allocation (LDA). In this LDA model, we use movie items as *documents* and users who have rated a movie as its *words*. Note that $R(U, I)$ and $dist(I)$ are two indicators of package quality, in terms of group rating on items and package viability, respectively. BASE is expected to generate packages with the best $R(U, I)$, because it is designed to combine items most liked by the groups regardless of the relationship among them. A desirable model should have similar $R(U, I)$ to BASE and at the same time find packages with small $dist(I)$ (i.e., high viability).

Figure 3 shows the average $R(U, I)$ over the packages recommended by BASE, GR and UP, respectively, on Yelp and MovieLens. Since each model recommends a set of top-10 packages, we average $R(U, I)$ (and $dist(I)$) over all these packages. BASE performs the best because of its design goal, however GR finds packages of nearly the same group-item rating. UP, on the other hand, always performs the worst because it only considers user impact at the package level, failing to address cases where different users have different impact on the various categories in C . As expected, for SIM groups, the models perform similarly, as it is easy to find packages where all items satisfy all group members.

Figure 4 compares the models based on the average distance $dist$ between items. Since BASE ignores relationships between the items, the packages it selects may contain items that are far from each other and have high $dist(I)$ values. UP fails to find packages with items close to each other, which are liked by the group as a whole, but not that much by individual group members (i.e., representatives); hence, its performance w.r.t. $dist(I)$ is worse than that of GR. On MovieLens, $dist(I)$ tends to be larger than on Yelp, because it is harder for two movies to be very similar to each other, compared to finding venues in Yelp that are spatially close. In addition, we observe that the relative performance among the models is the same regardless of the similarity between group members (SIM/DSIM). Overall, GR performs the best considering $dist(I)$, while being only marginally inferior to BASE w.r.t. $R(U, I)$. In the rest of the experiments, we only show results for the more interesting case of DSIM groups.

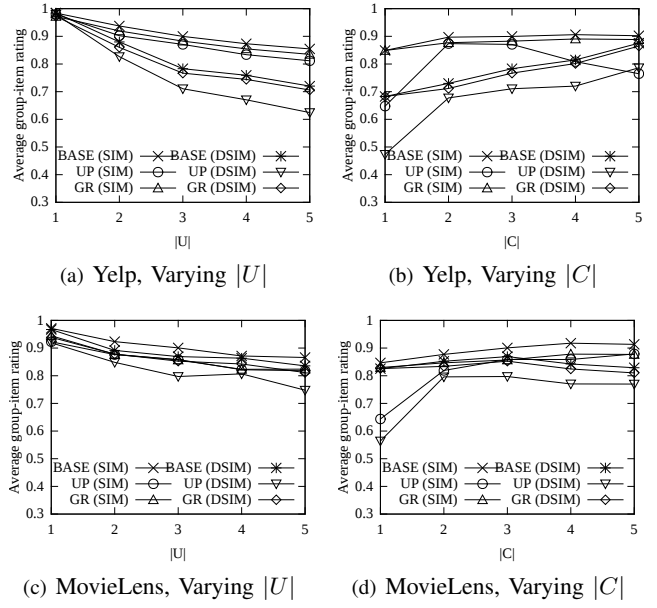


Fig. 3. Group Rating

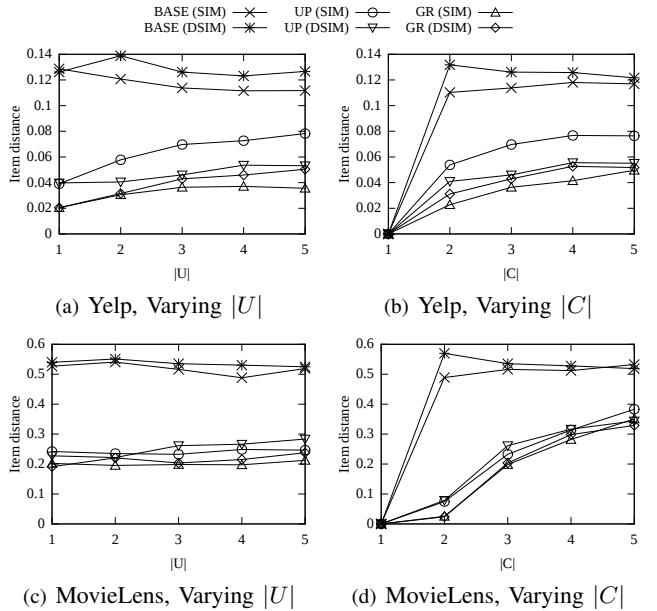


Fig. 4. Item Distance

C. Fairness Evaluation

In this section, we compare the basic GR and UP models presented in Section III with the variations GR-Fair and UP-Fair that consider fairness (see Equation (8)). Our goal is to understand the tradeoff between quality of recommendation and fairness. Figure 5 shows the package quality in terms of R , for all three versions of GR and UP. Figure 6 shows the average fairness degree of the packages; $fair(U, I)$, defined in Equation (7), with $\Delta = 10$. In order to consider a metric of fairness independent of the ones optimized in our algorithms, we also compute the mean highest rank of an item $i \in I$ for a user u . Formally, $hrank(U, I) = avg_{u \in U} \min_{i \in I} rank(u, i)$ where $rank(u, i)$ is defined as the rank of i among all items

rated by u in category c_i (normalized to $(0, 1]$, the lower the better). Intuitively, if a user u is happy in at least one category, at least one item will have high rank. Figure 7 shows $hrank$ for our algorithms.

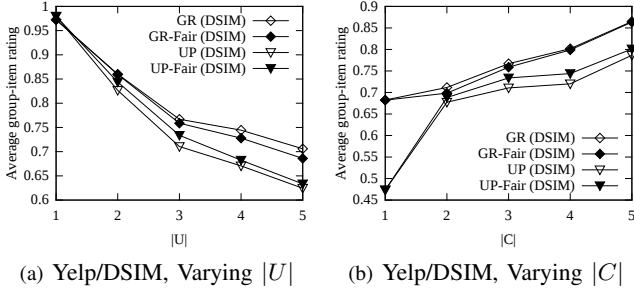


Fig. 5. Fair Models: Group Rating

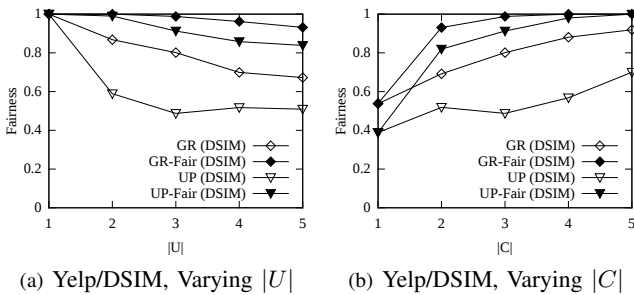


Fig. 6. Fair Models: Fairness Degree

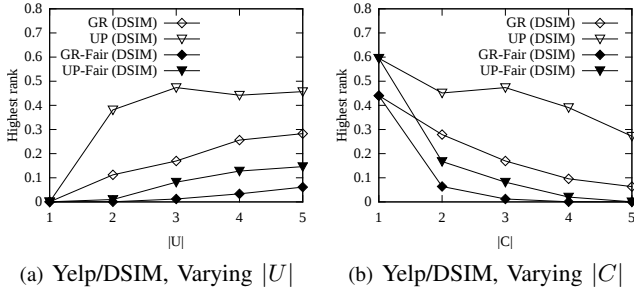


Fig. 7. Fair Models: Highest Rank

The first observation from these plots is that introducing fairness to GR reduces the quality, as it prevents the model from selecting packages of higher quality which are not fair to some users. Nevertheless, the loss in quality is relatively small. On the other hand, the gains in fairness are significant: GR-Fair improves both $hrank$ (the lower the better) and $fair$. Surprisingly, we observe that the addition of fairness *improves* the quality of UP (i.e., UP-Fair performs better than UP). Note that UP is inherently unfair (it has the worst performance in all fairness metrics – Figures 6 and 7), since it bases the selection on the preferences of a single user. The introduction of fairness counter-balances the drawbacks of UP, and forces the selection process to consider better packages.

Figure 8 evaluates the effect of the fairness threshold Δ , which controls the tradeoff between package quality and fairness. The figure shows the $hrank$ and quality values (based

on Section III-C) against Δ . For small Δ , an item must be ranked very high by u to make u happy; on the other hand, if Δ is large, fairness becomes looser. As expected, with Δ increasing, quality improves and fairness deteriorates. $\Delta = 10$ gives a good tradeoff between the two.

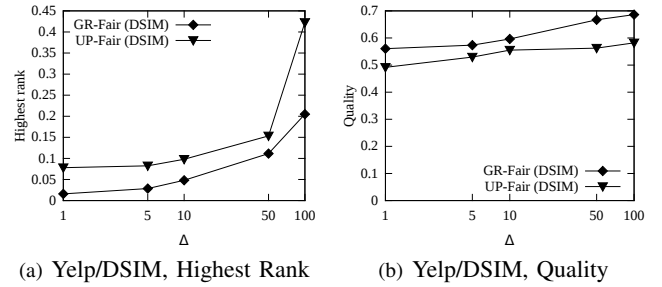


Fig. 8. GR-Fair Varying Δ

We repeated the above tests on MovieLens; the results are consistent with those on Yelp. In sum, GR-Fair finds packages such that users are more likely to be happy by at least one item, while not compromising quality compared to GR.

User Study. We also conducted a user study with 30 participants (students) to test the effectiveness of our models and the importance of fairness. First, we asked each participant to rate 70 popular movies belonging to 5 different genres (action, animation, comedy, romance, thriller). The participants were divided into 10 groups of 2-4 users each. For each group, movie packages with 3-4 genres were generated using BASE, GR-Fair, and UP-Fair. We also used a RAND model which selects movies randomly and a least-misery (LM) model that minimizes the maximum compromise a member makes for the group. We asked each group to assess the created packages by providing (1) an overall *rating* (PR) of the package and (2) a characterization of its *fairness* (PF). We did not provide any information on how the packages were generated and presented them to the groups in random order. Figure 9 depicts the average of the PR and PF values (0–1) given by the users. We also report the $R(U, I)$ and $dist(I)$ values of the packages as defined in Section VI-B. In terms of PR, the GR-Fair model outperforms all other models, i.e., it generated the packages that the groups liked the most. This is consistent with the fairness (PF), where GR-Fair also gives the best result, indicating that group satisfaction is correlated with fairness. Lastly, the relative values of $R(U, I)$ and $dist(I)$ are consistent with our experiments on Yelp and MovieLens.

D. Efficiency Evaluation

Finally, we evaluate the efficiency of algorithms GR-BA, GR-INC, UP-BA and UP-INC that implement GR and UP models (Section V). In terms of CPU cost, as Figures 10(a) and 10(b) show, GR-INC outperforms GR-BA by up to an order of magnitude, especially for large values of $|U|$ or $|C|$. As opposed to GR-BA, GR-INC accesses and calculates items/packages in an incremental fashion only when necessary, and stops once the bounding condition is satisfied. Similarly, UP-INC outperforms UP-BA. Note that the UP model is more expensive than GR to compute, because UP prioritizes

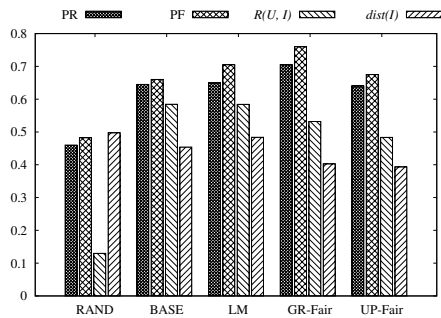


Fig. 9. User Study

packages favored by a single user, however, most of these items/packages are not favored by the other users and do not participate in the results. Package recommendations are more costly on MovieLens (Figure 10(c)) compared to Yelp. This is again due to the different item distance distribution between Yelp and MovieLens; on MovieLens, it is more likely that packages have larger distance and thus lower viability, rendering the termination condition during package formation harder to hold. Finally, GR-INC and UP-INC outperform GR-BA and UP-BA, respectively, in terms of accesses to item ratings (Figure 10(d)). Summing up, (1) GR-INC and UP-INC greatly outperform baseline implementations of GR and UP and (2) GR is not only better than UP in terms of quality of suggested packages, but also it is much faster to compute.

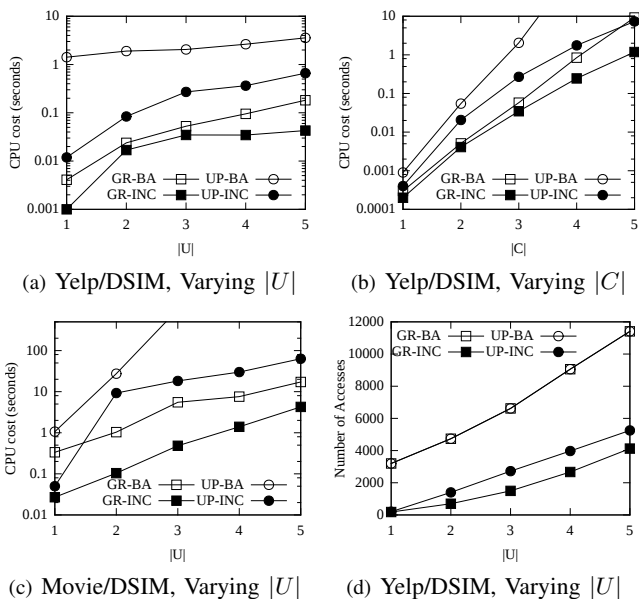


Fig. 10. Cost of Algorithms

E. The Case of No Category Constraints

Lastly, as discussed in Section III-D, our definitions and models are also applicable in the more general case, where there are no category constraints. Figures 11–13 show the performance of our models on the Yelp dataset, for the case where the users specify the desired number of items $|I|$ to be drawn from the general pool of items regardless of their

categories. Observe that the relative performance of the models is similar to that of the category-constrained case. Specifically, GR-Fair is best at finding packages that are more likely to satisfy each user by at least one item (i.e., being fair) without compromising quality compared to GR.

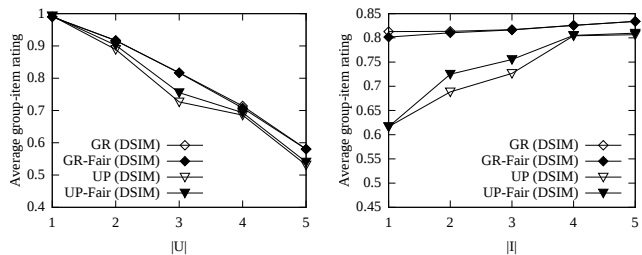


Fig. 11. Fair Models Without Category Constraints: Group Rating

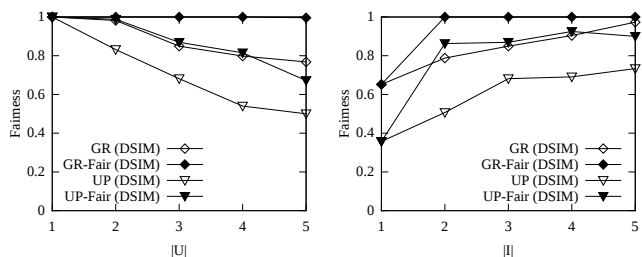


Fig. 12. Fair Models Without Category Constraints: Fairness Degree

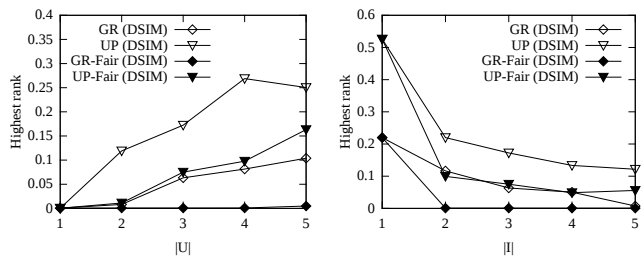


Fig. 13. Fair Models Without Category Constraints: Highest Rank

VII. RELATED WORK

Package to User Recommendation One category of previous work deals with recommending a package of items to a single user. The recommender by [19] finds packages of items that collectively maximize the user’s interest, but whose total cost does not exceed a given budget. Budget-based package selection, considering diversity and complementarity is studied in [1]. In [3], it is shown that selecting a package of items is a hard problem because of the larger search space; strict user-defined constraints can reduce this complexity (e.g., see the work of [12]). To avoid searching the whole space, Xie et al. [20] propose a learning process for predicting the interestingness of packages to users. Interdonato et al. [7] form packages for different models under item type compatibility

and given contextual constraints, and then rank them based on the user’s ratings and model/item property preferences. Zhu et al. [22] study the problem of recommending packages to a user by maximizing the expected reward of the packages. The reward expectation of a package depends on the probability of the user buying all its items together, which can be derived from the transactions history. Package viability (discussed in Section III-C) is a generalization of the reward defined in [22].

Item to Group Recommendation Another line of work deals with the recommendation of single items to a group of users. Some approaches [8] combine the ratings of all group members, in order to derive the ratings of a single artificial *representative* user for the group; then, a base recommender is used. Other methods compute recommendations for each group member separately and then aggregate them [11]. For the computation of the combined rating, [2] also considers the agreement between group members. Some recent works (e.g., [9]) use the social relationships between members to derive group recommendations. Using feedback from users to improve group recommendation has been studied in [13], [14]. Gorla et al. [5] define the probability of a group liking an item, based on the item’s relevance to each user as an individual. The I2G component of our GR model (see Section III-A) is an extension of [5] where we also consider the impact of each user on the different categories. Liu et al. [10] propose a personal impact weighted topic model, where each user has different impact on the group’s selection of topics and thus items; i.e., the group selection may be more biased to the preferences of the more influential user.

Yuan et al. [21] propose an improved consensus model (COM) which differentiates the preference of a user to a topic as an individual or a group member and defines topic-specific user impacts. Our P2G recommendation problem differs from those studied in [10] and [21]; in our case, the group requests recommendations of items from particular categories. Therefore, we use a probabilistic model with users’ item preference in each category, instead of a topic model with users’ topic and item distribution, to derive the group preferences. Finally, [21] considers content information (e.g. venue distance) to improve group recommendation. However, such information is derived from the user selection history and is used to infer the user’s historical preference. This is different from our definition of package viability, which models the potential of a set of items being selected as a package.

VIII. CONCLUSION

In this paper, we studied the problem of recommending one or more packages of items to a group of users. We proposed two probabilistic models (GR and UP), both of which incorporate individual ratings by users to items, user impacts, and package viability. In addition, we introduced fairness which is a unique but important feature of the P2G problem. Algorithms were proposed to efficiently implement the two models. Our experiments show that the GR-Fair model finds packages of superior quality in terms of user satisfaction, package viability, and fairness, compared to baseline approaches and UP models.

In addition, our algorithms GR-INC and UP-INC clearly outperform baseline implementations. We plan to study additional classes of P2G problems, e.g., when items are selected based on soft/hard budget constraints. We also plan to investigate more issues related to fair P2G recommendation, for example algorithms that find packages of maximum fairness.

ACKNOWLEDGEMENTS

This work was supported by grant 17205015 from Hong Kong RGC and by Marie Curie Reintegration Grant projects titled JMUGCS and LBSKQ which have received research funding from the European Union.

REFERENCES

- [1] S. Amer-Yahia, F. Bonchi, C. Castillo, E. Feuerstein, I. Méndez-Díaz, and P. Zabala. Composite retrieval of diverse and complementary bundles. *IEEE TKDE*, 26(11):2662–2675, 2014.
- [2] S. Amer-Yahia, S. B. Roy, A. Chawla, G. Das, and C. Yu. Group recommendation: Semantics and efficiency. *PVLDB*, 2(1):754–765, 2009.
- [3] T. Deng, W. Fan, and F. Geerts. On the complexity of package recommendation problems. *SIAM J. Comput.*, 42(5):1940–1986, 2013.
- [4] R. Fagin, A. Lotem, and M. Naor. Optimal aggregation algorithms for middleware. *J. Comput. Syst. Sci.*, 66(4):614–656, 2003.
- [5] J. Gorla, N. Lathia, S. Robertson, and J. Wang. Probabilistic group recommendation via information matching. In *WWW*, pages 495–504, 2013.
- [6] I. F. Ilyas, W. G. Aref, and A. K. Elmagarmid. Supporting top-k join queries in relational databases. In *VLDB*, pages 754–765, 2003.
- [7] R. Interdonato, S. Romeo, A. Tagarelli, and G. Karypis. A versatile graph-based approach to package recommendation. In *ICTAI*, pages 857–864, 2013.
- [8] A. Jameson and B. Smyth. Recommendation to groups. In *The Adaptive Web, Methods and Strategies of Web Personalization*, pages 596–627, 2007.
- [9] K. Li, W. Lu, S. Bhagat, L. V. S. Lakshmanan, and C. Yu. On social event organization. In *KDD*, pages 1206–1215, 2014.
- [10] X. Liu, Y. Tian, M. Ye, and W.-C. Lee. Exploring personal impact for group recommendation. In *CIKM*, pages 674–683, 2012.
- [11] M. O’Connor, D. Cosley, J. A. Konstan, and J. Riedl. PolyLens: A recommender system for groups of user. In *ECSCW*, pages 199–218, 2001.
- [12] A. G. Parameswaran, P. Venetis, and H. Garcia-Molina. Recommendation systems with complex constraints: A course recommendation perspective. *ACM TOIS*, 29(4):20, 2011.
- [13] J. A. Recio-García, G. Jiménez-Díaz, A. A. Sánchez-Ruiz-Granados, and B. Díaz-Agudo. Personality aware recommendations to groups. In *RecSys*, pages 325–328, 2009.
- [14] S. B. Roy, S. Thirumuruganathan, S. Amer-Yahia, G. Das, and C. Yu. Exploiting group recommendation functions for flexible preferences. In *ICDE*, pages 412–423, 2014.
- [15] B. M. Sarwar, G. Karypis, J. A. Konstan, and J. Riedl. Item-based collaborative filtering recommendation algorithms. In *WWW*, pages 285–295, 2001.
- [16] H. Steinhaus. The problem of fair division. *Econometrica*, 16(1):101–104, 1948.
- [17] M. Stettinger. Choicla: Towards domain-independent decision support for groups of users. In *RecSys*, pages 425–428, 2014.
- [18] M. Stettinger, A. Felfernig, G. Leitner, and S. Reiterer. Counteracting anchoring effects in group decision making. In *UMAP*, pages 118–130, 2015.
- [19] M. Xie, L. V. S. Lakshmanan, and P. T. Wood. Breaking out of the box of recommendations: from items to packages. In *RecSys*, pages 151–158, 2010.
- [20] M. Xie, L. V. S. Lakshmanan, and P. T. Wood. Generating top-k packages via preference elicitation. *PVLDB*, 7(14):1941–1952, 2014.
- [21] Q. Yuan, G. Cong, and C.-Y. Lin. Com: A generative model for group recommendation. In *KDD*, pages 163–172, 2014.
- [22] T. Zhu, P. Harrington, J. Li, and L. Tang. Bundle recommendation in ecommerce. In *SIGIR*, pages 657–666, 2014.