

Developing Preference Band Model to Manage Collective Preferences

Wilfred Ng

Department of Computer Science and Engineering
The Hong Kong University of Science and Technology
Hong Kong
wilfred@cse.ust.hk

Abstract. Discovering user preference is an important task in various database applications, such as searching product information and rating goods and services. However, there still lacks of a unifying model that is able to capture both implicit and explicit user preference information and to support managing, querying and analysing the information obtained from different sources.

In this paper, we present a framework based on our newly proposed Preference Band Model (PBM), which aims to achieve several goals. First, the PBM can serve as a formal basis to unify both implicit and explicit user preferences. We develop the model using a matrix-theoretic approach. Second, the model provides means to manipulate different sources of preference information. We establish a set of algebraic operators on Preference-Order Matrices (POMs). Third, the model supports direct querying of collective user preference and the discovery of a preference band. Roughly, a preference band is a ranking on sets of equally preferred items discovered from a POM that presents collective user preference. We demonstrate the applicability of our framework by studying two real datasets.

1 Introduction

Discovering user preference is an important task in various database applications, such as searching product information and rating goods and services [2,7,8]. Many business activities also rely heavily on estimating overall user preference in order to import user preferred goods and to design appropriate selling tactics. Previous work that incorporates user preference into databases mainly falls into two main categories of approaches. It either assumes that users are able to formulate their preference explicitly in terms of formulas or constraints [2,7,8], which may not be easy in reality, or develops mining techniques to discover implicit users' preference from the log of item selection [9] and then generate an adaptive full ranking of items, which may be too restrictive. However, there still lacks of a unifying model that is able to capture both implicit and explicit user preference information and to support managing, querying and analysing the information obtained from different sources.

We identify three problems that arise from handling preferences obtained from different sources and times. First, individual user preference can be modelled as a partial or full ranking between items [7,8]; but how do we model a large amount of such ranking information, which represents very different, possibly noisy and conflicting user

preferences? Second, how do we compare and contrast expected preference and real user preference? This also gives challenges to focus on some target preference data for analysing. Third, how do we rank and classify collective user preference and express the result in some form of simple but useful knowledge?

In this paper, we present a holistic framework shown in Figure 1, which involves a formal model, a set of algebraic operators and various algorithmic techniques to tackle the above problems. The framework aims to manage preference information in a systematic way and to support better analyses of collective user preferences.

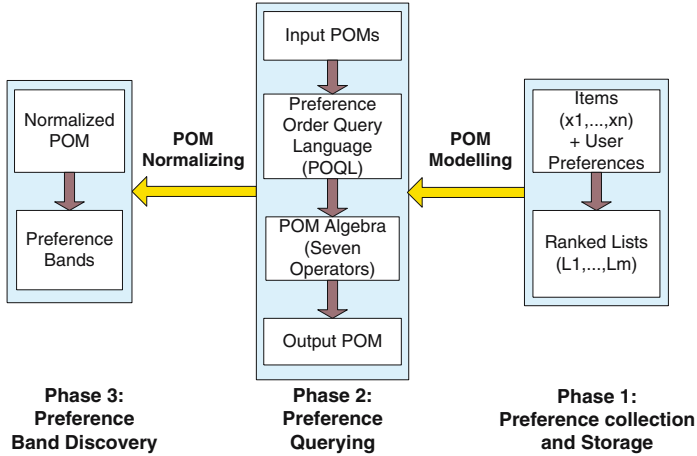


Fig. 1. The framework for managing collective user preference based on PBM

The basic idea of our approach is to convert a collection of user ranked items into a Preference-Order Matrix (POM). An entry a_{ij} in a POM is the frequency that one item x_i (more preferred) precedes another item x_j (less preferred) in the collected information, which captures the fact that users prefer x_i to x_j in a number of choice records.

As shown in the blueprint presented in Figure 1, we classify the process based on our Preference Band Model (PBM) into three phases. In the first phase, the PBM integrates implicit and explicit preference information represented by a set of user ranked lists, each of which is a user preference with possible choice order on target items. The information is modelled within a unifying matrix structure called POMs. In the second phase, we use a high level declarative query language called POQL to support manipulation of preference data and to generate a more effective POM for further analyses. The language is executed and can be optimized via a set of POM operators, which shares the similar spirit of execution of SQL expressions in relational databases. In the third phase, we need to normalize the POM in order to have it run in an adapted PIVOT algorithm [1], which classify and rank items into a preference band. The algorithm is known to be an efficient and effective method for discovering order knowledge from similar matrix structures, which is capable of tolerating noise of order in some reasonable limit [5,4].

Our main contribution in this paper is the development work of the PBM that supports managing user preference information and discovering preference bands.

- We formalise the novel concept of preference-order matrices (POMs) that express collected preference information. POMs support the discovery of preference bands that represent the maximal consistent collective preference obtained from a given set of user preference ranked lists.
- We establish a new set of seven algebraic operators such as the sum, the union, the intersection and the selection on POMs. We develop a new declarative language POQL based on the POM operators.
- We identify efficient algorithmic techniques for discovering preference bands, which shows that the discovery of preference bands from a POM is feasible in practice. Finally, we present the preliminary results of two interesting applications, which show the idea of preference bands is easily applicable.

The rest of the paper is organised as follows. Section 2 gives some preliminary concepts of user preference. Section 3 defines a set of formal POM operators. Section 4 discusses the issues arising from the three phases of the PBM framework. Section 5 presents our results of finding preference bands on two applications. Section 6 reviews some related work on user addressing preference and the algorithmic techniques for finding preference bands. In Section 7, we give our concluding remarks.

2 Preliminaries

User preferences can be expressed in an implicit and explicit order of items. For example, explicit user preference is formulated in a lexicographic order of preference attributes [8]. Preference SQL [7] is equipped with a “preferring” clause that allows user to specify soft constraints reflecting multiple preference terms. All these approaches would result in a preference ranking of items. Implicit preferences are some ranking order that is inferred from users’ track record of choosing items against some background of item order [9], which also depends on the interpretation of the chosen items.

We assume that a preference rank list can be obtained via some users’ implicit or explicit preference information and formalize the concept as follows:

Definition 1. (Preference Ranked List) Given a set of items $I = \{x_1, \dots, x_n\}$ with an imposed order $<_I$. A choice $C \subseteq I$ is a list of items ordered according to some user preference $<_C$, where $(C, <_C)$ may or may not equal to $(I, <_I)$. We generate a *preference ranked list* T by appending $(I - C)$ into C such that $T = \langle C, (I - C) \rangle$, in which the items of $(I - C)$ are ordered according to $<_I$. We denote by $L = \{T_1, \dots, T_m\}$ a collection of m preference ranked lists.

Note that in Definition 1, $<_C$ can represent the order of choosing items in C by the user, which is an implicit preference, or it can represent the order according to some explicit preference criteria. L is a collective user preference on I . The following example helps illustrate the ideas of individual and collective preferences.

Example 1. Let $I = \{a, b, c, d, e, f\}$ be the items stored in a database and two customers search some items on the web. Let the items be ranked according to the default order $I = \{a > b > c > d > e > f\}$, which can be interpreted as an

arbitrary ranking or a ranking derived from some commercial considerations. One customer explicitly states his or her preference criteria, which result in a preferred ranking $C_1 = \{c > f > e > d > b > a\}$. Another customer does not state any preference and just browses the items ranked by the default order $<_I$. However, s/he only checks a subset of items of I in some different order given by $C_2 = \{d > b > f > e\}$.

We now show in Figure 2 two sets of preference ranked lists L_1 and L_2 , each of which contains five preference ranked lists. L_1 and L_2 can represent collective customer preference information obtained from two sources. It can be checked that T_1 and T_8 are the preference ranked lists derived from C_1 and C_2 respectively. T_5 can also be viewed as a precise match between a user's preference and the shop's estimated user preference that is expressed in the default item order.

$L_1 =$	$T_1: c > f > e > d > b > a$	$L_2 =$	$T_6: c > f > e > d > b > a$
	$T_2: a > c > b > d > f > e$		$T_7: a > c > b > d > f > e$
	$T_3: a > c > d > b > f > e$		$T_8: d > b > f > e > a > c$
	$T_4: d > f > c > a > b > e$		$T_9: d > f > c > a > e > b$
	$T_5: a > b > c > d > e > f$		$T_{10}: b > d > f > e > c > a$

Fig. 2. An example of two sets of five preference ranked lists

A collection of preference ranked lists L from a given source can be modelled as a Preference Order Matrices (POMs) M_L , which is a matrix structure derived from L to support further querying and preference band discovery.

Definition 2. (Preference Order Matrix) Given L and let $|L| = n$. A *Preference-Order Matrix* (POM) M_L with respect to L (or simply M if L is understood) is an $n \times n$ matrix such that each entry $a_{ij} \in M$ is equal to the number of occurrences that item x_i precedes item x_j in L . We define $a_{ij} \in M$ to be $|L|$ whenever $i = j$. M is said to be *normalized* if each entry $norm(a_{ij}) = \frac{a_{ij}}{n}$ (or equivalently, $\frac{a_{ij}}{a_{ii}}$ or $\frac{a_{ij}}{a_{jj}}$) whenever $i \neq j$, and $norm(a_{ij}) = \frac{1}{2}$ whenever $i = j$.

Clearly, it follows from Definition 2 that $0 \leq a_{ij} \leq |L|$, $a_{ij} + a_{ji} = |L|$ for any distinct i, j . The normalization implies that each entry $a_{ij} \in norm(M_L)$ is equal to the probability that item x_i precedes item x_j in L . Given $norm(M_L)$. The entries satisfies the following conditions: (1) (normalization constraint) $0 \leq a_{ij} \leq 1$; (2) (linearity constraint) $a_{ij} + a_{ji} = 1$; and (3) (triangle constraint) $a_{ij} \leq a_{ik} + a_{kj}$.

Normalised POMs are important in this work. For simplicity of presentation we assume all T being a total order. However, we remark that if a partial order is resulted from more general user preferences [7] we still can use a corresponding set of linear extensions to represent the order information of equally preferred items. I.e. two equally preferred items x_i and x_j can still be captured by normalised POMs in the entry $a_{ij} = \frac{1}{2}$. Specifically, a partial order can be directly represented by a normalised POM as follows: $a_{ij} \in M$ is equal to 1 if $x_i \leq x_j$, $a_{ij} \in M$ is equal to 0 if $x_j \leq x_i$, and $a_{ij} \in M$ is equal to $\frac{1}{2}$ otherwise (i.e. x_i and x_j are incomparable). In other words, normalised POMs are flexible enough to capture more general user preference information.

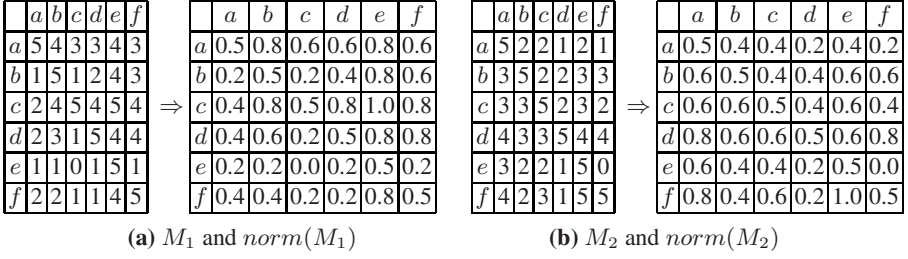


Fig. 3. The POMs M_1 and M_2 obtained from L_1 and L_2 , and their normalized counterparts

Example 2. We now continue Example 1 to derive the POMs and their normalized form from L_1 and L_2 . The POMs M_1 and M_2 and their normalized forms are given in Figures 3(a) and (b).

One objective of this work is to establish a set of operators that are able to manipulate POMs and support in-depth analyses of collective user preference modelled in POMs. Another objective is to discover preference bands obtained from a given normalised POM. A preference band is intuitively a linear order of sets of equally preferred items taking into collective user preference. The motivation of defining a preference band is that in many adaptive searching applications, modelling the preference with a total order might be too restrictive in the sense that all items in I should be comparable. However, it is also too complex to manage, analyse and present the preference information expressed in general partial order.

We now give the formal definition of a preference band as follows.

Definition 3. (Preference Band) A preference band B of I with granularity k is an ordered partition of I given by $\langle P_1, \dots, P_k \rangle$, where $P_i \subseteq I$ is a non-empty set called a band element in B . A band order of the items in I arising from B , denoted as $<_B$, can be inferred from B as follows: For any pair of items x and y in the same band element P , x and y are incomparable (i.e. $x \not<_B y$ and $y \not<_B x$). However, if x and y are from two distinct band elements P_i and P_j such that $i < j$, then $x <_B y$.

Essentially, a preference band adopts an appropriately low granularity to eliminate noises and ranking contradictions of items among preference ranked lists collected in L . However, there is an inherent cost of losing some information of the preference orders in individual T when putting items in a band element; as an extreme we can put all items in one single band element which is obviously no use. Thus, we need to establish a penalty cost function for measuring the quality of a preference band. A simple one is adopted as below but more sophisticated functions can be studied as a future work.

If M_1 and M_2 are two normalized POMs, we define the distance measure given by $D(M_1, M_2) = \sum_{i \neq j} |a_{ij}^{M_1} - a_{ij}^{M_2}|$.

Given a normalised POM M . We aim to find the best preference band describing M . We now formulate the Preference Band Discovery (PBD) problem.

Definition 4. (PBD problem) Let M_L be a normalised POM on I obtained from L and M_B be a normalised POM obtained from a preference band B of I . Find a preference

band such that $D(M_L, M_B)$ is minimized. (Note that B is a partial order and can be represented by POM M_B .)

In other words, the main task of PBD is to find a preference band B on the items of I such that the penalty cost function is minimized, which is equal to the sum of values in the following three cases: (i) $|(a_{ij})_M - 1|$ if $x_i <_B x_j$, (ii) $|(a_{ij})_M|$ if $x_j <_B x_i$, and (iii) $|(a_{ij})_M - \frac{1}{2}|$ otherwise.

Example 3. Consider our running example L_1 and L_2 given in Figure 2. Suppose there is a preference band $B = \langle \{a, b, c, d\} > \{e, f\} \rangle$ (which may not be the best one according to D). Then, it can be checked that $D(M_{L_1}, M_B) = 6.6$ and $D(M_{L_2}, M_B) = 10$. The difference is reasonable, since it can be checked that L_1 has a noise level of 40% with respect to the preference band B because two of the five preference ranked lists, namely T_1 and T_4 , are not the linear extension [5] of B , whereas L_2 has an even worse noise level of 80% with respect to B because only T_7 is the linear extension of B . This motivates us to develop an effective algorithm to generate B from a given M_L .

3 The POM Operators

In this section, we define seven POM operators, namely the sum, the union, the difference, the intersection, the complement, the projection and the selection, each of which takes one or two given POMs as input parameters, and returns a POM.

The POM operators are easy to understand and to compute and their output result serves as a basis for carrying out the discovery of preference band. The sum operator is additive to preference data, that is, the result is incremental with respect to preference information. This gives advantage to handle preference data in a progressive manner, say temporal change in preference can be detected. The overall, change, common, opposite preference can be formalised by the sum, the union, the difference, the intersection, and the complement operations. We can focus on the preference information of a particular set of items by using the projection and the selection operations. On the other hand, matrices are an elegant notion studied in a well-established branch of mathematics. It implies that many interesting results in matrix theory can be used to strengthen the foundation of PBMs as a development work.

We need two binary operators, denoted as *min* and *max*, to represent the usual minimum and maximum of two given integers in defining the POM operators. From now on, we use throughout the paper M_1 and M_2 to represent two POMs defined over the same set of items I . The sum operation is given in Definition 5.

Definition 5. (Sum) The *sum* of two POMs M_1 and M_2 , denoted as $M_1 + M_2$, is defined as a POM M_3 over I such that for all $i, j \in \{1, \dots, n\}$, $(a_{ij})_3 = (a_{ij})_1 + (a_{ij})_2$.

An interesting property of the sum operation is that the sum of two POMs M_1 and M_2 , which originates from the two respective preference ranked lists L_1 and L_2 , is equal to the matrix which originates from the preference ranked list L_3 containing the combination of the information in L_1 and L_2 (i.e. $L_3 = L_1 \cup L_2$). This implies that we are able to perform analyses on “overall” preference information by summing up

“individual” pieces of preference information, in this sense we say that the sum operator is *additive* with respect to the preference data.

Another possibility to combine preference information is to consider the maximum number of occurrences of precedence between items, taking into account the “overlap” effect. We now define this operation by the union operator as follows.

Definition 6. (Union) The *union* of two POMs M_1 and M_2 , denoted as $M_1 \cup M_2$, is defined as a POM M_3 over I such that $(a_{ii})_3 := \max((a_{ii})_1, (a_{ii})_2)$ for all $i \in \{1, \dots, n\}$ and for all pairs of distinct $i, j \in \{1, \dots, n\}$, we obtain $(a_{ij})_3$ as follows:

Let $x := \max((a_{ij})_1, (a_{ij})_2)$ and $y := \max((a_{ji})_1, (a_{ji})_2)$.

If $x > y$ then $(a_{ij})_3 := x$ and $(a_{ji})_3 := \max((a_{ii})_1, (a_{ii})_2) - x$;

otherwise (i.e. $x \leq y$) $(a_{ji})_3 := y$ and $(a_{ij})_3 := \max((a_{ii})_1, (a_{ii})_2) - y$.

Note that the comparison $x > y$ makes sure that the most dominant precedence in M_1 and M_2 is chosen as the output. In contrast, we can consider the less dominant precedence but it is common to both M_1 and M_2 . This becomes another operation called intersection defined as follows.

Definition 7. (Intersection) The *intersection* of two POMs M_1 and M_2 , denoted as $M_1 \cap M_2$, is defined as a POM M_3 over I such that $(a_{ii})_3 := \min((a_{ii})_1, (a_{ii})_2)$ for all $i \in \{1, \dots, n\}$ and for all pairs of distinct $i, j \in \{1, \dots, n\}$, we obtain $(a_{ij})_3$ as follows:

Let $x := \max((a_{ij})_1, (a_{ij})_2)$ and $y := \max((a_{ji})_1, (a_{ji})_2)$.

If $x < y$ then $(a_{ij})_3 := x$ and $(a_{ji})_3 := \min((a_{ii})_1, (a_{ii})_2) - x$;

otherwise (i.e. $x \geq y$) $(a_{ji})_3 := y$ and $(a_{ij})_3 := \min((a_{ii})_1, (a_{ii})_2) - y$.

We now introduce the difference operator, which is useful to contrast two preference ranked lists. For example, we can use the difference operator to find the temporal change in preference data obtained at two different time intervals or to compare the preference data obtained from two user groups having different preference profiles.

Definition 8. (Difference) Assume $|L_1| > |L_2|$. The *difference* of two POMs M_1 and M_2 , denoted as $M_1 - M_2$, is defined as a POM M_3 over I such that $(a_{ii})_3 := (a_{ii})_1 - (a_{ii})_2$ for all $i \in \{1, \dots, n\}$ and for all pairs of distinct $i, j \in \{1, \dots, n\}$, we obtain $(a_{ij})_3$ as follows:

Let $x := (a_{ij})_1 - (a_{ij})_2$ and $y := (a_{ji})_1 - (a_{ji})_2$.

Case $x > 0$ and $y > 0$: $(a_{ij})_3 := x$ and $(a_{ji})_3 := y$;

Case $x > 0$ and $y \leq 0$: $(a_{ij})_3 := \min(x, (a_{ii})_1 - (a_{ii})_2)$ and $(a_{ji})_3 := (a_{ii})_1 - (a_{ii})_2 - \min(x, (a_{ii})_1 - (a_{ii})_2)$;

Case $y > 0$ and $x \leq 0$: $(a_{ij})_3 := (a_{ii})_1 - (a_{ii})_2 - \min(y, (a_{ii})_1 - (a_{ii})_2)$ and $(a_{ji})_3 := \min(y, (a_{ii})_1 - (a_{ii})_2)$.

Notably, the operator is undefined for $|L_1| \leq |L_2|$. This prevents the happening of meaningless negative entries in a POM. Thus, there is no need to consider the possibility of $y \leq 0$ and $x \leq 0$ in the definition. It also implies that either x or y is strictly positive as shown in the three cases in Definition 8.

The following is an interesting operator that defines the “reverse” of the preference order for all items in I .

Definition 9. (Complement) The *complement* of a POM M_1 , denoted as $\neg M_1$, is defined as a POM M_2 over I such that for all $i, j \in \{1, \dots, n\}$, $(a_{ij})_2 = (a_{ji})_1$.

Similar to relational algebra, the POM operators are employed to extract a target set of items. The operators projection and selection aim to realise similar objectives. But unlike operating on relational tables, an POM operation needs to preserve the entries of the non-target items in the output, since POM is a matrix. To tackle this problem, we simply fill in the entries of these items according to the default order $<_I$. I.e. All the items in $(I - X)$ in the projection are ordered according to $<_I$ in the output POM.

Definition 10. (Projection) The *projection* of a POM M_1 , denoted as $\pi_X(M_1)$ where $X \subseteq I$ is a set of items, is defined as a POM M_2 as follows:

Case $x_i, x_j \in P$: $(a_{ij})_2 = (a_{ij})_1$.

Case $x_i \in P$ but $x_j \notin P$: $(a_{ij})_2 = (a_{ij})_1 + (a_{ji})_2$ and $(a_{ji})_2 = 0$.

Case $x_i, x_j \notin P$: $(a_{ij})_2 = (a_{ij})_1 + (a_{ji})_2$ if $x_j <_I x_i$; otherwise (i.e. $x_i <_I x_j$) $(a_{ji})_2 = 0$.

The tricky point in defining the POM selection is that we should not allow arbitrary entry comparison, since it is difficult to know or seldom need to concern the absolute occurrence of precedence in L . Instead, we choose the relative ratio of precedence between items as the parameter used in the selection predicate, which always falls in a unit interval. This also captures the intuition that how much x is more preferred to y .

Definition 11. (Selection) The *selection* of a POM M_1 , denoted as $\sigma_{\theta x}(M_1)$, where θ is a comparator such as $>$, $<$, \geq , \leq , or $=$, and $x \in [0, 1]$ is a positive number, is defined as a POM M_2 as follows:

For all $i \in \{1, \dots, n\}$ $(a_{ii})_2 = (a_{ii})_1$.

For all distinct $i, j \in \{1, \dots, n\}$, if $\left| \frac{(a_{ij})_1 - (a_{ji})_1}{(a_{ij})_1 + (a_{ji})_1} \mid \theta x \right.$, then $(a_{ij})_2 = (a_{ij})_1$;

otherwise (i.e. $\left| \frac{(a_{ij})_1 - (a_{ji})_1}{(a_{ij})_1 + (a_{ji})_1} \mid \bar{\theta} x \right.$ if $x_j <_I x_i$ $(a_{ij})_2 = (a_{ij})_1 + (a_{ji})_2$; or else (i.e. $x_i <_I x_j$) $(a_{ji})_2 = 0$,

where the notation $\bar{\theta}$ denotes the complement of θ (e.g. $\bar{\theta}$ is “ \leq ” when θ is “ $>$ ”).

It follows from Definition 11 $\sigma_{\geq 0}(M) = M$, where all entries are trivially selected. The result of “ $\sigma_{>1}(M)$ ” represents the preference being reduced to the default order $<_I$. This can be used as the *identity* POM ID for revealing some interesting properties of the union and intersection operators in the following discussion.

We let $\delta_1 \in \{+, \cup, \cap\}$ and $\delta_2 \in \{\neg, \pi_X\}$ and $\delta_3 \in \{\neg, \sigma_p, \pi_X\}$. The following properties of the POM operations can be verified by Definitions 5 to 11.

Associative Property. The POMs on the same I are associative under *sum*, *union* and *intersection* which is shown in the equations as follows:

$$(M_1 \delta_1 M_2) \delta_1 M_3 = M_1 \delta_1 (M_2 \delta_1 M_3).$$

Commutative Property. The POMs on the same I are also commutative under *sum*, *union* and *intersection*.

$$(1) M_1 \delta_1 M_2 = M_2 \delta_1 M_1. (2) \delta'_3(\delta_3(M)) = \delta_3(\delta'_3(M)).$$

Distributive Property. *Projection* and *negation* are distributive under *sum*, *union* and *intersection*.

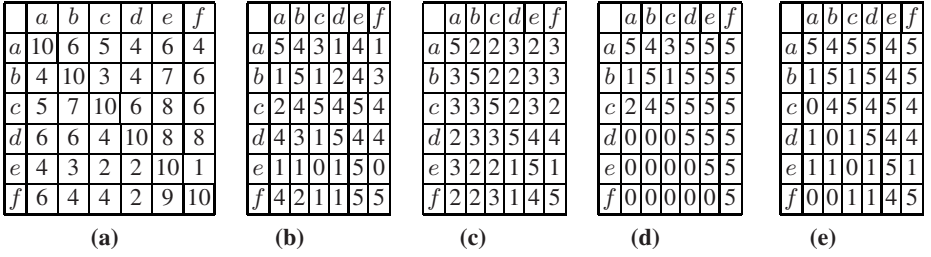


Fig. 4. The results of the some binary and unary POM operations (a) $M_1 + M_2$ (b) $M_1 \cup M_2$ (c) $M_1 \cap M_2$ (d) $\pi_{abc}(M_1)$ (e) $\sigma_{\theta > 0.2}(M_1)$

$$\delta_2(M_1 \delta_1 M_2) = (\delta_2(M_1)) \delta_1 (\delta_2(M_2)).$$

In general, $\sigma_p(M_1 \delta_1 M_2) \neq (\sigma_p(M_1)) \delta_1 (\sigma_p(M_2))$ and $\pi_X(M_1 - M_2) \neq (\pi_X(M_1)) - (\pi_X(M_2))$. However, $\neg(M_1 - M_2) = (\neg(M_1)) - (\neg(M_2))$ and $\neg(\neg(M)) = M$.

We now let $\sigma_{>1}(M) = ID$ and present the identity property.

Identity Property. Union and intersection satisfy the identity property.

- (1) $M \cup ID = ID$. (2) $M \cap ID = M$.

Example 4. We now make use the POMs M_1 and M_2 given in Figure 3 to show some of the results of $M_1 \delta_1 M_2$ and $\delta_3(M_2)$ in Figure 4.

4 Implementation Issues

In this section, we briefly discuss the techniques and methods that address the challenges arising from the three phases of the PBM framework depicted in Figure 1.

4.1 Phase 1: How to Obtain and Store POMs?

We can construct the POM M by using L in two ways. The first way is to adopt a brute force approach to bookkeep a_{ij} when scanning the total order in which x_i precedes x_j . As the algorithm for discovering preference bands needs to take a normalized POM as the input, we still need to compute the normalized values of the entries in M in the final phase. So we may use a more direct way to obtain $norm(M)$. We can sample possible orderings of the user preference $T \in I$ and in this case $norm(a_{ij})$ is set to be the fraction of L in which x_i precedes x_j , such as Markov Chain Monte Carlo method used in [10]. In this approach, we avoid the heavy computation to track all precedences. In fact it is inevitable to have noise in preference data in real applications. Sampling techniques seem to be more appropriate to generate normalized POMs.

The storage scheme of a matrix greatly affects the performance of POM operators. In practice, there are usually many items to be considered. As I is large and the number of preferred items is small, the POMs may be very sparse with respect to the entries that actually represent user preference data. As the entries of diagonal and lower triangular portion of a POM can be deduced from $|L|$, we may develop a concise version that targets on only those entries that contain preference information. Specifically, we fill in zeros in entry of (i) the lower triangular portion of the POM and (ii) the diagonal running

top left to bottom right. The first condition is valid, since we have $a_{ji} = |L| - a_{ij}$. The second condition is also valid since we have $a_{ii} = |L|$. There will be more zero entries in the upper triangular portion of M if the ranking order is skew.

4.2 Phase 2: How to Use the POM Operators to Formulate Queries?

We develop a declarative language on POMs in our framework and term the language the Preference Order Query Language (or simply the POQL). A POQL expression is executed via the seven POM operators defined in Section 3, which shares the same principle of translating a SQL expression into a sequence of relational algebra operations.

We now define the POQL syntax in Backus Naur Form (BNF):

```

<query> ::= <selectClause><fromClause>[<conditionClause>]
<selectClause> ::= SELECT <itemList>
<queryList> ::= query [, query. . .]
<fromClause> ::= FROM <matrixIdentifier> | FROM <operator> <matrixList>
<conditionClause> ::= WHERE PREFERENCE RATIO <compOp> a number in [0,1]
<itemList> ::= itemIdentifier [, itemIdentifier. . .] | *
<matrixList> ::= matrixIdentifier [, matrixIdentifier. . .]
<operator> ::= SUM|UNION|DIFFERENCE|INTERSECT|COMPLEMENT
<compOp> ::= |=|<=|>|=|<|=

```

There are three main clauses in a query expression: the *select*, *from*, and *condition*. Among them the *select* and the *from* clauses are compulsory, while the *condition* clause is optional. Similar to SQL, POQL is a simple declarative language but is expressive enough to formulate query on finding preference information stored as POMs.

We execute a POQL expression by translating it into a sequence of POM operations using Algorithm 1. Suppose $X \subseteq L$ is a set of items which appears in the select clause, \mathcal{M} is a set of m POMs which appears in the from clause, where $\mathcal{M} =$

Algorithm 1. Translate POQL Algorithm

Input: A POQL expression q

LET $q = \langle \text{selectClause} \rangle \langle \text{fromClause} \rangle [\langle \text{conditionClause} \rangle]$

$\langle \text{selectClause} \rangle := \text{"SELECT } X \text{ | *"}$

$\langle \text{fromClause} \rangle := \text{"FROM OPER } \mathcal{M}"$

$\langle \text{conditionClause} \rangle := \text{"WHERE PREFERENCE RATIO } \theta x"$

Procedure:

Step 1 : For the fromClause, **CASE OPER OF:**

$\epsilon : TEMP := "M_1"$

COMPLEMENT : $TEMP := "\neg M_1"$

SUM : $TEMP := "M_1 + M_2 + \dots + M_m"$

UNION : $TEMP := "M_1 \cup M_2 \cup \dots \cup M_m"$

DIFF : $TEMP := "M_1 - M_2 - \dots - M_m"$

INTERSECT : $TEMP := "M_1 \cap M_2 \cap \dots \cap M_m"$

Step 2 : For the selectClause:

IF "*" **THEN** $TEMP := TEMP$

ELSE $TEMP := "\pi_X(TEMP)"$

Step 3 : **IF** there is a whereClause **THEN** $TEMP := "\sigma_{\theta x}(TEMP)"$

Output: $TEMP$ expression

$\{M_1^*, M_2^*, \dots, M_m^*\}$ and M^* means either M or *COMPLEMENT* M , and $OPER \in \{\epsilon, SUM, UNION, DIFFERENCE, INTERSECT\}$. Note that the input POQL query expression is assumed to be syntactically valid. If $OPER = \epsilon$, then $m = 1$.

4.3 Phase 3: How to Generate Preference Bands?

The challenge of this final phase is to tackle the PBD problem in Definition 4. However, the problem is equivalent to the Bucket Order Discovery (BOD) problem, which aims to find a linear order of buckets from a given set of linear orders. The BOD problem has been proved to be NP-hard [1]. Thus, we have to resort to heuristic algorithms. Given a set of full rankings T_1 and T_2 over I . Gionis et al. [5] proposed a heuristic algorithm called Bucket Pivot, which is adapted from Ailon’s randomized algorithm Pivot [1] (or simply called the PIVOT algorithm).

Essentially, the PIVOT algorithm starts with a random selected element and then compares the elements with others and finally generates three classes of “left”, “same” and “right” classes. To apply the techniques in our context, the PIVOT algorithm can be employed to exploit precedence probabilities stored in the input normalized POM M and then to discover the output preference band recursively, in which a preference band can be regarded as a bucket order B . The three classes of buckets correspond to the lower band (LB), current band (CB) and upper band (UB) elements.

Specifically, the adapted PIVOT algorithm runs in a quick-sort-like manner. In each recursion, a random pivot item x_i is first selected and the following three band elements, $\langle UB \rangle \langle CB \rangle \langle LB \rangle$, are also created for x_i . Other items are then compared with the pivot x_i . An item x_j will be put into CB if the precedence probability a_{ij} satisfies the following inequality: $0.5 - \beta \leq a_{ij} < 0.5 + \beta$, where $\beta \in [0, 0.5]$ is a parameter that describes the degree of precision of precedence probability relative to 0.5. The drawback of using PIVOT is that it is not easy to set a suitable β value and the choice of β affects the final number of band elements in B . Clearly if $\beta = 0$ we cannot generate any approximated band element to provide insight of the collective preference in the input POM. On the other hand, if $\beta = 0.5$ we have a very imprecise (and is likely to be large) band element that contains all x_j having $a_{ij} \neq 1$ (which has a high probability).

Algorithm 2. The Preference Band Discovery Algorithm $PBD(X, M, \beta)$

Input: A normalized POM M of I , $X \subseteq I$ and $\beta = 0.25$ by default

Procedure:

Step 1 : IF $X = \emptyset$ RETURN \emptyset

Step 2 : Pick $x_i \in I$ randomly as a pivot DO

$UB \leftarrow \emptyset$

$CB \leftarrow \{x_i\}$

$LB \leftarrow \emptyset$

Step 3 : FOR ALL items $x_j \in (X - \{x_i\})$ DO

IF $0.5 + \beta \leq a_{ij}$ THEN $UB \leftarrow UB \cup \{x_j\}$

ELSE IF $0.5 - \beta \leq a_{ij} < 0.5 + \beta$ THEN $CB \leftarrow CB \cup \{x_j\}$

ELSE IF $a_{ij} < 0.5 - \beta$ THEN $LB \leftarrow LB \cup \{x_j\}$

Output: A preference band given by $\langle PBD(UB, M, \beta), CB, PBD(LB, M, \beta) \rangle$

Fortunately, according to [5], the default value $\beta = 0.25$ is shown to be able to make PIVOT achieve good approximation ratio with respect to the penalty cost function D .

We now present the adapted PIVOT algorithm for PBD problem in Algorithm 2, which adapts PIVOT in our contest. We simply assume $\beta = 0.25$ for running PBD.

5 Applications of Preference Bands

We now present two applications of the PBM to discover the user preference in (1) a real rating movie dataset and (2) a real clickstream dataset. A clickstream is the evidence of user preference in choosing a web page to browse.

(1) Movie Preference. We consider the ranking of movie obtained from MovieLens dataset¹. The dataset consists of 6040 users, 3900 movies, and 1 million movie votes. Each of the movie votes has a weight and a score, where the weight denotes whether or not the user actually saw the movie, and where the score denotes the user's movie rating. Using Algorithm 2, we ignore the weight and set it to 1 for all users. There are actually very different preference in ranking the movies. We target on the following set of ten movies to find the PB:

$I = \{\text{Strange Days (SD), Lawnmover Man 2 (LM), Flintstones (FS), Free Willy (FW), Godfather (GF), 3 Musketeers (MU), 101 Dalmatians (DM), Empire Strikes Back (ES), Barb Wire (BW), and Jungle Book (JB)}\}$.

Figure 5 shows the normalised POM for the rankings of the 10 movies. We obtain the best preference band of four band elements of movies as follows: $\{\{ES, GF\} > \{JB\} > \{SD, DM, MU\} > \{FW, FS, BW, \} > \{LM\}\}$. The band elements are not so trivial to be obtained from the huge amount of data or even from the POM given in Figure 5, which is obtained after the processing efforts in Phases 1 and 2.

(2) Web Page Preference. We consider the user page visits of msnbc.com on a single day. The clickstream dataset is obtained from MSNBC data². Each user has exactly one

	<i>ES</i>	<i>GF</i>	<i>JB</i>	<i>SD</i>	<i>MU</i>	<i>DM</i>	<i>FW</i>	<i>FS</i>	<i>LM</i>	<i>BW</i>
<i>ES</i>	0.50	0.32	0.78	0.92	0.90	0.92	0.95	0.96	1.00	0.99
<i>GF</i>	0.68	0.50	0.86	0.94	0.90	0.92	0.95	0.94	1.00	0.98
<i>JB</i>	0.22	0.14	0.50	0.69	0.77	0.88	0.92	0.94	0.96	0.95
<i>SD</i>	0.08	0.06	0.31	0.50	0.69	0.71	0.77	0.92	0.96	0.91
<i>MU</i>	0.10	0.10	0.23	0.31	0.5	0.65	0.81	0.88	0.96	0.85
<i>DM</i>	0.08	0.08	0.12	0.29	0.35	0.5	0.68	0.70	0.80	0.88
<i>FE</i>	0.05	0.05	0.08	0.23	0.19	0.32	0.5	0.64	0.77	0.67
<i>FS</i>	0.04	0.06	0.06	0.08	0.12	0.30	0.36	0.5	0.83	0.67
<i>LM</i>	0.00	0.00	0.04	0.04	0.04	0.20	0.23	0.17	0.50	0.19
<i>BW</i>	0.01	0.02	0.05	0.09	0.05	0.12	0.33	0.33	0.81	0.5

Fig. 5. The normalized POM of movie preference to discover preference bands

¹ <http://www.grouplens.org/node/73#attachments>

² <http://kdd.ics.uci.edu/databases/msnbc/msnbc.html>

User 1:	frontpage	→	news	→	sports	→	weather
User 2:	news	→	politics	→	weather		
⋮	⋮			⋮			⋮
User n :	frontpage	→	politics	→	weather		

Fig. 6. Collective users' browsing habits to form a preference band

sequence of page visits, where each page visit is recorded as a URL category. There are 17 categories in total, and each category is encoded by a distinct integer from 1 to 17. In other words, $I = \{1, \dots, 17\}$ with usual numerical order and L consists of a set of T sequences, each of which is a permutation of I .

First, we keep the first occurrence of a category in each sequence T and denote the resulting sequence as C . Usually, $C \subseteq I$ and thus $|C|$ is less than 17. We denote the $(17 - |C|)$ categories as S and sort the categories in S according to the numerical order of their integer code. Then, we append S to C to form a permutation of I as a result, which is a T sequence. We select a subset of the sequences whose $|C|$ contains at least 14 categories in order to increase the effect of the user preference in L . In total, we identify 160 such T sequences.

Examples of URL categories are “frontpage” and “msn-sports”. The page visit sequence is a total order if we do not consider duplicate page visits. The order of page visits reflects a user's browsing habit. By clustering page-visit orders, we generate a user preferred categorization. For example, in Figure 6, we may use a preference band to indicate the common browsing habits of the set of users, who may usually start from reading some news, either about politics or about sports, and end up with a look at the weather conditions.

The best preference band discovered in MSNBC data is found as follows: $\{\{frontpage, news\} > \{tech, local, on-air, misc, weather, msn-news, health, living, business, sports, summary, travel\} > \{opinion\} > \{msn-sports, bbs\}\}$. It shows that the item “news” is put in the first band element in addition to the “frontpage”. The result that both “frontpage” and “news” are in the top band element is not so trivial but seems to match real-life experience: while users can first visit “frontpage” and then go from the “frontpage” to “news”, users may have bookmarked “news” and visit “news” directly. In contrast, if we use a preference ranking instead of a preference band to represent such browsing preference, “frontpage” would most likely be put in the top of the preference rank, which wrongly assumes that users normally start with the “frontpage” and then go to the “new”. Thus, using a preference band is a more natural and meaningful representation of collective preference.

6 Related Work

Preferences are receiving much attention in querying, since DBMSs need to provide better information services in advanced applications [7]. In particular, preference SQL [7] is equipped with a “preferring” clause that allows user to specify soft constraints reflecting multiple preference terms. Implicit preference has been commonly studied in

the area of searching such as using clickthrough data to mine user preference in web browsing [9]. In reality, implicit and explicit user preferences are important but there lacks of a formal basis to accommodate and manipulate both of them.

Our previous work [8] models single user preference as a hierarchy of its underlying data values and formalises the notion of Prioritized Preferences (PPs). We then consider multiple user preferences in ranking tuples in a relational table. We examine the impact of a given set of PPs on possible choices in ranking a database relation and develop a new notion of Choice Constraints (CCs) in a relation.

The PBD problem can be translated into the bucket order discovery (BOD) problem that is studied in recent years. The problem of obtaining a single bucket order from a collection of input total orders has been considered by Fagin et al. [3] in their general framework of comparing and aggregating partial rankings. In our recent work [4], we develop a new algorithm, called GAP, to tackle the BOD problem. The GAP algorithm consists of a two phase ranking aggregation and involves the use of a novel rank gap heuristic for segmenting multiple quantile orders.

7 Concluding Remarks

In this paper, we propose a new PBM framework that consists of three phases of work for managing a collection of user preference data that are modelled as POM matrices. We also establish a set of operations which serve as the formal tool to analyse and manipulate preference data. We further develop a declarative query language based on the POM operators and a new concept of preference band to discover and classify collective user preference. We discuss various technical issues related to the three phases.

There are indeed many interesting issues that deserve further study. In the modelling aspect, we need to clarify if POQL is expressive enough to obtain important preference information from given POMs. In the deployment aspect, we still need to study the efficiency of the POM operators and devise effective optimization strategy for the execution of the operations. The algorithm of finding preference bands is also related to the research work of discovering bucket order in data mining area. In the application aspect, we are exploring more interesting possibilities to apply our framework, in addition to usual databases and web applications discussed in Section 5.

Acknowledgements. Thanks to Qiong Fang for updating the experimental results presented in Section 5 in this published version. This work is partially supported by HKUST grant under grant No. DAG04/05.EG10.

References

1. Ailon, N., Charikar, M., Newman, A.: Aggregating inconsistent information: Ranking and clustering. In: ACM STOC, pp. 684–693 (2005)
2. Chomicki, J.: Preference formulas in relational queries. *ACM Transaction Database System* 28(4), 427–466 (2003)
3. Fagin, R., Kumar, R., Mahdian, M.: Comparing and aggregating with ties. In: ACM PODS, pp. 47–58 (2004)

4. Feng, J., Fang, Q., Ng, W.: Discovering Bucket Orders from Full Rankings. In: ACM SIGMOD (2008)
5. Gionis, A., Mannila, H., Puolamaki, K., Ukkonen, A.: Algorithms for discovering bucket orders from data. In: ACM SIGKDD, pp. 561–566 (2006)
6. Goharian, N., Jain, A., Sun, Q.: Comparative analysis of sparse matrix algorithms for information retrieval. *Journal of Sys., Cyb. and Inf.* 1(1) (2003)
7. Kießling, W., Köstler, G.: Foundations of Preference in Database Systems. In: Proc. of VLDB (2002)
8. Ng, W.: Prioritized Preferences and Choice Constraints. In: Parent, C., Schewe, K.-D., Storey, V.C., Thalheim, B. (eds.) ER 2007. LNCS, vol. 4801, pp. 261–276. Springer, Heidelberg (2007)
9. Tan, Q., et al.: Applying Co-training to Clickthrough Data for Search Engine Adaptation. In: Lee, Y., Li, J., Whang, K.-Y., Lee, D. (eds.) DASFAA 2004. LNCS, vol. 2973, pp. 519–532. Springer, Heidelberg (2004)
10. Puolamäki, K., Fortelius, M., Mannila, H.: Seriation in paleontological data using markov chain monte carlo methods. *PLoS Computational Biology* 2(2) (2006)
11. van Zuylen, A., et al.: Deterministic pivoting algorithms for constrained ranking and clustering problems. In: ACM-SIAM SODA, pp. 405–414 (2007)