

# Evaluating Aggregate Operations over Imprecise Data

ARBEE L.P. CHEN, JUI-SHANG CHIU AND FRANK S.C. TSENG

Department of Computer Science

National Tsing Hua University

Hsinchu, Taiwan 300

Republic of China

*Email: alpchen@cs.nthu.edu.tw*

## ABSTRACT

Imprecise data in databases were originally denoted as *null values*, which represent the meaning of “values unknown at present.” More generally, a partial value corresponds to a finite set of possible values for an attribute in which exactly one of the values is the “true” value. In this paper, we define a set of extended aggregate operations, namely *sum*, *average*, *count*, *maximum*, and *minimum*, which can be applied to an attribute containing partial values. Two types of aggregate operators are considered: *scalar aggregates* and *aggregate functions*. We study the properties of the aggregate operations and develop efficient algorithms for *count*, *maximum* and *minimum*. However, for *sum* and *average*, we point out that in general it takes exponential time complexity to do the computations.

**Index Terms:** relational databases, null values, partial values, scalar aggregates, aggregate functions, graph theory.

# 1 Introduction

Incomplete information can be classified into two aspects, i.e., imprecise information and uncertain information. Dubois and Prade [12] distinguish between imprecise and uncertain information by stating that the concept of *imprecision* is relevant to the content of an attribute value, while the concept of *uncertainty* is relevant to the degree of truth of its attribute value. In this paper, we focus our attention on the imprecision aspect. To the uncertainty aspect, two major approaches are the *possibility* approach based on the fuzzy set theory [33] and the *probability* approach based on the probability theory.

Imprecise data exist in the real world. *Null values* were originally adopted to represent the meaning of “values unknown at present” in database systems. Codd [6, 7] pioneers the work on extended relational algebra to manipulate null values. For a concise review of handling null values by algebraic approaches, refer to [21].

The concept of null values has been generalized to the concept of *partial values* by Grant [14]. Instead of being treated as an atomic value, an attribute value in a table is considered as a nonempty subset of the corresponding domain. A partial value is represented as an interval such that exactly one of the values in the interval is the “true” value. In our work, a partial value is considered as a finite set of *possible* values in which exactly one of the values is the “true” value [11]. Therefore, an *applicable null value* [8] is a partial value which corresponds to the whole domain. Also, a definite value  $a$  is a partial value which corresponds to a set containing  $a$  only.

Partial values have been considered by DeMichiel [11] for resolving domain mismatch problems in multidatabase systems, in which an algebraic approach for operating on partial values is proposed. Grant and Minker [15] consider query answering for indefinite (disjunctive) databases that contain disjunctive formulas in first-order logic in which each corresponds to a finite set of possible tuples under the relational model. Ola [23] presents an approach to processing relations containing exclusive disjunctive data (such as partial values). In [29], we discuss the elimination of redundant partial values. Moreover, in [32] we study the problem of incrementally refining partial values into more informative ones or into definite values by utilizing the knowledge of integrity constraints. Lipski [20] presents a general discussion on manipulating imprecise information including partial values.

To the uncertainty aspect, various kinds of fuzzy relational databases based on the possibility approach were proposed, such as Bosc *et al.* [3], Buckles and Petry [4], Prade and Testemale [25], and Zemankova and Kandel [34]. Moreover, the work with the probability approach includes Barbará *et al.* [1], Cavallo and Pittarelli [5], Tsai and Chen [28], and Tseng *et al.* [30]. In [28, 30], we generalize partial values into *probabilistic partial values* to provide a probabilistic approach to query processing in heterogeneous database systems.

In practice, relational algebra or calculus are inadequate for many important applications involving statistical information or aggregations. Therefore, modern query languages, like SQL [9], are “more than” *relational complete* and equipped with some useful aggregate operations. However, most of the previous works on the manipulation of incomplete information usually discuss the extended relational algebra and ignore the extended aggregate operators.

Aggregate operations can be applied to all the tuples in a table, to a subset of the table (specified by a WHERE clause), or to one or more groups of tuples in the table (specified by a GROUP BY clause). For each set of tuples to which an aggregate operation is applied, a single value is generated. Aggregate operations are often used with the GROUP BY clause for a more powerful functionality. Optionally, the keyword DISTINCT can be used with *count* to eliminate duplicate values before the aggregate operation is applied.

In this paper, we define a set of extended aggregate operations, namely *sum*, *average*, *count*, *maximum*, and *minimum*, which can be applied to attributes containing partial values. These aggregate operators are divided into the following two types (depending on whether the GROUP BY clause is used).

1. *scalar aggregates*— They take a set of partial values as input and produce a single partial value as a result.
2. *aggregate functions*— They compute an aggregation on attribute  $A_y$  of a relation  $T[\dots, A_x, A_y, \dots]$  containing partial values based on the distinct scalar values  $a_x$  of attribute  $A_x$ , and return  $\langle a_x, \eta_y \rangle$ , where  $\eta_y$  is the aggregation of  $A_y$  based on  $a_x$ . For databases containing partial values, the extended aggregate functions produce  $\eta_y$  as a partial value.

In [10], DeMichiel provides an *approximation* of aggregate operators over partial values. However, from the definitions in [10], only the *boundary possible values* (i.e., the maximum and minimum possible values) are returned for some aggregate operators. In comparison, our work generates *all possible values* for each aggregate operator instead of *boundary possible values*.

Özsoyoğlu *et al.* [24] study some aggregate operations over *set-valued* attributes. All values in the set are the true values. The *grouping attributes* if specified can be set-valued attributes. However, it is required that the *aggregate attribute* (to which an aggregate function is applied) should be a single-valued attribute. In comparison, in our work, both grouping attributes and aggregate attributes are partial values.

Rundensteiner and Bic [26, 27] study aggregate operations in possibilistic databases. The possibilistic relational model is characterized by a representation which allows for data values to be modeled by discrete possibility distributions. Their aggregate operators generate all possible values each with a possibility in exponential time. In comparison, based on partial values, our work generate all possible values without possibilities possibly in polynomial time.

Database practitioners are concerned primarily with the *performance* of a database system and many of the query evaluation algorithms for processing imprecise queries are very complex and inefficient. This inefficiency, as pointed out by Motro [22], is one of the reason for the slow acceptance of imprecision and uncertainty handling capabilities into commercial database systems. Therefore, this work is devoted to the accommodation of imprecise data in database systems with elaboration on speeding up its processing efficiency.

Since the brute-force computation of the extended aggregates is very time-consuming in general, we find some important properties and develop efficient algorithms (in polynomial time complexity) for *count*, *maximum*, and *minimum*. However, for *sum* and *average*, we point out that in general it takes exponential time complexity to do the computations.

The following section presents some basic concepts and definitions of partial values. Section 3 is devoted to the definition of our extended scalar aggregates and studies their properties. In Section 4, the definition of our extended aggregate functions and their evaluation algorithms are presented. Section 5 concludes our work.

## 2 Basic Concepts and Definitions

### 2.1 Basic Concept of Partial Values

Partial values model data imprecision in databases in the sense that, for an imprecise datum, its “true” value can be restricted to be in a specific set of possible values [11] or an interval of values [14]. In our work, we follow the definition of a partial value given in [11], which is formally stated as follows.

**DEFINITION 2.1** A *partial value*, denoted  $\eta = [a_1, a_2, \dots, a_p]$ , corresponds to a finite set of possible values,  $\{a_1, a_2, \dots, a_p\}$ ,  $p \geq 1$ , of the same domain for an attribute, in which exactly one of the values in  $\eta$  is the “true” value of  $\eta$ .

For a partial value  $\eta = [a_1, a_2, \dots, a_p]$ , a function  $\nu$  is defined on it [11], where  $\nu$  maps the partial value to its corresponding finite set of *possible* values; that is,  $\nu(\eta) = \{a_1, a_2, \dots, a_p\}$ . Recall that an *applicable null value* [8],  $\aleph$ , can be considered as a partial value with  $\nu(\aleph) = D$ , where  $D$  is the whole domain. A relation containing tuples with partial values is called a *partial relation* [11]. Otherwise, it will be called a *definite relation*. We assume that each partial relation contains primary key attributes with definite values.

The *cardinality* of a partial value  $\eta$  is defined as  $|\nu(\eta)|$  in [11]. When the cardinality of a partial value equals 1, i.e., there exists only one *possible* value, say  $a$ , in the partial value, then the partial value  $[a]$  actually corresponds to the definite “true” value  $a$ . A definite value  $a$  can be also represented as a partial value  $[a]$ . Moreover, a partial value with cardinality greater than 1 is referred to as a *proper partial value* in [11]. For any two proper partial values, say  $\eta_1$  and  $\eta_2$ , the “true” value of  $\eta_1$  may not be the same as the “true” value of  $\eta_2$  even if  $\nu(\eta_1) = \nu(\eta_2)$ .

### 2.2 Alternate Worlds of a Partial Relation

Let  $T$  be a partial relation on a set of attributes  $\{A_1, A_2, \dots, A_m\}$  containing a set of tuples  $\{t_1, t_2, \dots, t_n\}$  where  $t_i = \langle \eta_{i1}, \eta_{i2}, \dots, \eta_{im} \rangle$ ,  $1 \leq i \leq n$ . Let  $T[A_j]$  denote the projection of a partial relation  $T$  on an attribute  $A_j$ . We may enumerate all the possible cases that  $T[A_j]$  represents by the following definition.

**DEFINITION 2.2** Let  $T[A_j] = \{\eta_{1j}, \eta_{2j}, \dots, \eta_{nj}\}$  where  $\eta_{ij} = t_i.A_j$ ,  $t_i \in T$ . An *interpretation*  $\alpha$  of  $T[A_j]$  is an assignment of values from  $T[A_j]$  denoted by an  $n$ -ary vector  $\alpha = (a_{1j}, a_{2j}, \dots, a_{nj})$  where  $a_{ij} \in \nu(\eta_{ij})$ ,  $1 \leq i \leq n$ .

By Definition 2.2, for  $T[A_j] = \{\eta_{1j}, \eta_{2j}, \dots, \eta_{nj}\}$ ,  $\nu(\eta_{1j}) \times \nu(\eta_{2j}) \times \dots \times \nu(\eta_{nj})$  is the set of all interpretations of  $T[A_j]$ . Notice that  $\eta_{ij}$  and  $\eta_{lj}$ ,  $i \neq l$ , can both appear in  $T[A_j]$  even when  $\nu(\eta_{ij}) = \nu(\eta_{lj})$ .

**EXAMPLE 2.1** Suppose  $T[A_j] = \{\overbrace{[a, b]}^{\eta_{1j}}, \overbrace{[c, d]}^{\eta_{2j}}\}$ , there are four interpretations,  $\alpha_1 = (a, c)$ ,  $\alpha_2 = (a, d)$ ,  $\alpha_3 = (b, c)$ , and  $\alpha_4 = (b, d)$ .

For  $t_i \in T$ , we may enumerate all the possible tuples that  $t_i$  represents by the following definitions.

**DEFINITION 2.3** Let  $t_i = \langle \eta_{i1}, \eta_{i2}, \dots, \eta_{im} \rangle \in T$ . An interpretation,  $\alpha = (a_{i1}, a_{i2}, \dots, a_{im})$ , of  $t_i$  is an assignment of values from partial values in  $t_i$  such that  $a_{ij} \in \nu(\eta_{ij})$ ,  $1 \leq j \leq m$ .

**DEFINITION 2.4** Let  $\alpha_k$  be an interpretation of  $t_i$ ,  $t_i \in T$ .  $\alpha_k(t_i)$  is defined as a tuple  $\langle a_{i1k}, a_{i2k}, \dots, a_{imk} \rangle$  where  $a_{ij_k}$  denotes the value of attribute  $A_j$  for tuple  $t_i$  under interpretation  $\alpha_k$ . The set of all possible tuples represented by  $t_i$ ,  $\text{rep}(t_i)$ , is defined as  $\{\alpha_k(t_i) \mid \text{for each interpretation } \alpha_k \text{ of } t_i\}$ .

For a partial relation  $T$ , we may enumerate all the possible relations that  $T$  represents by the following definitions.

**DEFINITION 2.5** Let  $T$  be defined as the above. An interpretation  $\alpha$  of  $T$  is an assignment

of values from partial values in  $T$  denoted by an  $n \times m$  matrix  $\alpha = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1m} \\ a_{21} & a_{22} & \cdots & a_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nm} \end{bmatrix}$ ,

where  $a_{ij} \in \nu(\eta_{ij})$ ,  $1 \leq i \leq n$ ,  $1 \leq j \leq m$ .

By Definition 2.5, for partial relation  $T$ ,  $\left\{ \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1m} \\ a_{21} & a_{22} & \cdots & a_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nm} \end{bmatrix} \mid a_{ij} \in \nu(\eta_{ij}), 1 \leq i \leq n, 1 \leq j \leq m \right\}$  is the set of all interpretations of  $T$ .

**DEFINITION 2.6** Let  $\alpha_k$  be an interpretation of  $T$ .  $\alpha_k(T)$  is defined as a relation  $\{\alpha_k(t_i) \mid \alpha_k(t_i) = \langle a_{i1_k}, a_{i2_k}, \dots, a_{im_k} \rangle, t_i \in T\}$ . The set of all possible relations represented by partial relation  $T$ ,  $\text{rep}(T)$ , is defined as  $\{R_k \mid R_k = \alpha_k(T), \text{ for each interpretation } \alpha_k \text{ of } T\}$ .

**EXAMPLE 2.2** Consider the following partial relation  $T$ .

$$T$$

$A_1$	$A_2$
$a$	$c$
$b$	$[c, d]$

There are two interpretations,  $\alpha_1$  and  $\alpha_2$ , of  $T$ :

$$\alpha_1 = \begin{bmatrix} a & c \\ b & c \end{bmatrix} \text{ and } \alpha_2 = \begin{bmatrix} a & c \\ b & d \end{bmatrix}.$$

The possible relations represented by  $T$ ,  $\alpha_1(T) = R_1$  and  $\alpha_2(T) = R_2$ , are shown below:

$$R_1$$

$A_1$	$A_2$
$a$	$c$
$b$	$c$

$$R_2$$

$A_1$	$A_2$
$a$	$c$
$b$	$d$

□

### 3 Extended Scalar Aggregates over Partial Values

An extended scalar aggregate can be applied to an attribute containing partial values. It in turn produces a partial value.

#### 3.1 Definitions of the Extended Scalar Aggregates

Each extended scalar aggregate is defined as follows through the function  $\nu$ , which maps the result of an extended scalar aggregate (which is also a partial value) to the finite set of values to which it corresponds.

In the following,  $T$  denotes a partial relation and  $A_j$  denotes an attribute of  $T$ .

**DEFINITION 3.1** The *extended count* (for the number of distinct values) on  $A_j$  of  $T$ , denoted  $\text{count}(T.A_j)$ , is defined as:

$$\nu(\text{count}(T.A_j)) = \{|R_k[A_j]| \mid R_k \in \text{rep}(T)\}.$$

**DEFINITION 3.2** The *extended sum* on  $A_j$  of  $T$ , denoted  $sum(T.A_j)$ , is defined as:

$$\nu(sum(T.A_j)) = \left\{ \sum_{t_i \in R_k} a_{ij_k} \mid R_k \in rep(T) \right\}.$$

**DEFINITION 3.3** The *extended average* on  $A_j$  of  $T$ , denoted  $avg(T.A_j)$ , is defined as:

$$\nu(avg(T.A_j)) = \left\{ \frac{\sum_{t_i \in R_k} a_{ij_k}}{|R_k|} \mid R_k \in rep(T) \right\}.$$

**DEFINITION 3.4** The *extended maximum* on  $A_j$  of  $T$ , denoted  $max(T.A_j)$ , is defined as:

$$\nu(max(T.A_j)) = \left\{ \max_{t_i \in R_k} a_{ij_k} \mid R_k \in rep(T) \right\}.$$

**DEFINITION 3.5** The *extended minimum* on  $A_j$  of  $T$ , denoted  $min(T.A_j)$ , is defined as:

$$\nu(min(T.A_j)) = \left\{ \min_{t_i \in R_k} a_{ij_k} \mid R_k \in rep(T) \right\}.$$

The following example illustrates these extended aggregates.

**EXAMPLE 3.1** Consider the following partial relation.

<i>Employees</i>	
<i>name</i>	<i>salary</i>
<i>a</i>	[20k, 80k]
<i>b</i>	[60k, 75k]
<i>c</i>	[20k, 90k]

We compute the set of possible relations  $rep(Employee)$ , as follows.

$rep(Employee)$															
a	20k	a	20k	a	20k	a	20k	a	80k	a	80k	a	80k	a	80k
b	60k	b	60k	b	75k	b	75k	b	60k	b	60k	b	75k	b	75k
c	20k	c	90k	c	20k	c	90k	c	20k	c	90k	c	20k	c	90k

Then the extended scalar aggregates can be computed as follows. Let  $T = Employees$  and  $A_j = salary$ .

$$\begin{aligned} \nu(sum(Employee.salary)) &= \left\{ \sum_{t_i \in R_k} a_{ij_k} \mid R_k \in rep(T) \right\} \\ &= \{100k, 170k, 115k, 185k, 160k, 230k, 175k, 245k\}. \end{aligned}$$

$$\nu(count(Employee.salary)) = \{|R_k[A_j]| \mid R_k \in rep(T)\} = \{2, 3\}.$$

$$\begin{aligned} \nu(avg(Employee.salary)) &= \left\{ \frac{\sum_{t_i \in R_k} a_{ij_k}}{|R_k|} \mid R_k \in rep(T) \right\} \\ &= \left\{ \frac{100k}{3}, \frac{170k}{3}, \frac{115k}{3}, \frac{185k}{3}, \frac{160k}{3}, \frac{230k}{3}, \frac{175k}{3}, \frac{245k}{3} \right\}. \end{aligned}$$

$$\nu(max(Employee.salary)) = \left\{ \max_{t_i \in R_k} a_{ij_k} \mid R_k \in rep(T) \right\} = \{60k, 90k, 75k, 80k\}.$$

$$\nu(min(Employee.salary)) = \left\{ \min_{t_i \in R_k} a_{ij_k} \mid R_k \in rep(T) \right\} = \{20k, 60k, 75k\}.$$



That is,

$$\begin{aligned}
sum(Employees.salary) &= [100k, 170k, 115k, 185k, 160k, 230k, 175k, 245k], \\
count(Employees.salary) &= [2, 3], \\
avg(Employees.salary) &= \left[ \frac{100k}{3}, \frac{170k}{3}, \frac{115k}{3}, \frac{185k}{3}, \frac{160k}{3}, \frac{230k}{3}, \frac{175k}{3}, \frac{245k}{3} \right], \\
max(Employees.salary) &= [60k, 90k, 75k, 80k], \text{ and} \\
min(Employees.salary) &= [20k, 60k, 75k]. \quad \square
\end{aligned}$$

## 3.2 Properties and Algorithms for Evaluating Scalar Aggregates over Partial Values

When  $|T| = n$  and  $|\eta_{ij}| = p$  for all  $t_i \in T$ , there are  $p^n$  interpretations of  $T[A_j]$ . Therefore, brute force methods for computing the extended scalar aggregates are exponential. In the following, we develop polynomial time algorithms to compute *count*, *maximum* and *minimum* based on their properties. Since the cardinalities of  $\nu(sum(T.A_j))$  and  $\nu(avg(T.A_j))$  are equal to  $p^n$  in the worst case (see Example 3.1), no polynomial time algorithms can be found for *sum* and *average*.

### 3.2.1 Properties of the Scalar Aggregate—*Count*

In this section, we use techniques in graph theory [2] to develop an efficient algorithm for *count*. Before discussing the details, we define some terminologies as follows.

**DEFINITION 3.6** For a set  $S$  of vertices in a graph  $G(V, E)$ ,  $S \subseteq V$ , the *neighbor set* of  $S$ , denoted  $N(S)$ , is defined to be the set of all vertices adjacent to the vertices in  $S$ .

**DEFINITION 3.7** For a bipartite graph  $G = (X \cup Y, E)$ , we say a set  $S$ ,  $S \subseteq X$ , *covers*  $Y$  if  $N(S) = Y$ . If deleting any vertex in  $S$  would make  $S$  cease to *cover*  $Y$  then  $S$  is a *minimal* set that covers  $Y$ .  $S$  is a *minimum* set if it has the minimum cardinality among these cover sets. Note that  $Y$  can be an empty set.

**DEFINITION 3.8** Let  $Y = \{\eta_1, \eta_2, \dots, \eta_n\}$  and  $X = \bigcup_{1 \leq i \leq n} \nu(\eta_i) = \{a_1, a_2, \dots, a_q\}$ . The *membership graph of  $Y$  over  $X$*  is a bipartite graph  $G = (X \cup Y, E)$ , where

$$E = \{(a_i, \eta_j) \mid a_i \in \nu(\eta_j), 1 \leq i \leq q, 1 \leq j \leq n\}.$$

**DEFINITION 3.9** For a bipartite graph  $G = (X \cup Y, E)$ ,  $|X| \leq |Y|$ , we say there is a *complete matching*  $M$  from  $X$  to  $Y$  if there is a matching of cardinality  $|X|$ , that is, each vertex in  $X$  is adjacent to a distinct vertex in  $Y$ .

Hall [16] has given a necessary and sufficient condition under which there exists a complete matching  $M$  from  $X$  to  $Y$  for a bipartite graph  $G = (X \cup Y, E)$ .

**THEOREM 3.1 (Hall [16])** *Let  $G = (X \cup Y, E)$  be a bipartite graph. There exists a complete matching from  $X$  to  $Y$  if and only if  $|N(S)| \geq |S|$ ,  $\forall S \subseteq X$ .  $\square$*

Now we state the properties for *count* as follows. We discuss the lower bound of  $\nu(\text{count}(T.A_j))$  first. Then the upper bound of  $\nu(\text{count}(T.A_j))$  is addressed. Finally, we prove that each integer number between the lower bound and the upper bound is in  $\nu(\text{count}(T.A_j))$ .

**LEMMA 3.1** *Let  $Y = T[A_j] = \{\eta_{1j}, \eta_{2j}, \dots, \eta_{nj}\}$ ,  $X = \bigcup_{1 \leq i \leq n} \nu(\eta_{ij})$ , and  $G = (X \cup Y, E)$  be the membership graph of  $Y$  over  $X$ . Let  $l$  denote the cardinality of the minimum set  $S$ ,  $S \subseteq X$ , that covers  $Y$ , then for all  $c \in \nu(\text{count}(T.A_j))$ ,  $c \geq l$ .*

*Proof:* Suppose there is an element  $c$  of  $\nu(\text{count}(T.A_j))$  such that  $c < l$ . That is, there is an interpretation  $\alpha_k = (a_{1jk}, a_{2jk}, \dots, a_{nj_k})$  of  $T[A_j]$  with  $|R_k[A_j]| = c < l$ , where  $R_k[A_j] = \{a_{ijk} \mid 1 \leq i \leq n\}$ . Since  $a_{ijk} \in \nu(\eta_{ij})$ , for  $1 \leq i \leq n$ ,  $R_k[A_j]$  covers  $Y$ . But  $|R_k[A_j]| < l$ , which contradicts the assumption that the cardinality of the minimum set that covers  $Y$  is  $l$ . Hence, for all  $c \in \nu(\text{count}(T.A_j))$ ,  $c \geq l$ .  $\square$

The following lemma states the upper bound of  $\nu(\text{count}(T.A_j))$ .

**LEMMA 3.2** *Let  $Y = T[A_j] = \{\eta_{1j}, \eta_{2j}, \dots, \eta_{nj}\}$ ,  $X = \bigcup_{1 \leq i \leq n} \nu(\eta_{ij})$ , and  $G = (X \cup Y, E)$  be the membership graph of  $Y$  over  $X$ . Let  $u$  denote the cardinality of the maximum matching  $M$  in  $G$ , then for all  $c \in \nu(\text{count}(T.A_j))$ ,  $c \leq u$ .*

*Proof:* Suppose there is an element  $c$  of  $\nu(\text{count}(T.A_j))$  such that  $c > u$ . That is, there is an interpretation  $\alpha_k = (a_{1jk}, a_{2jk}, \dots, a_{nj_k})$  of  $T[A_j]$  with  $|R_k[A_j]| = c > u$ , where  $R_k[A_j] = \{a_{ijk} \mid 1 \leq i \leq n\}$ . In other words, for each  $a_{ljk} \in R_k[A_j]$ , there exists at least a distinct  $\eta_{ij} \in T[A_j]$  such that  $a_{ljk} \in \nu(\eta_{ij})$ . That implies, for all  $S \subseteq R_k[A_j]$ ,  $|S| \leq |N(S)|$  in the membership graph of  $T[A_j]$  over  $R_k[A_j]$ ,  $G' = (R_k[A_j] \cup T[A_j], E')$ .

Therefore, by Theorem 3.1, there exists a complete matching  $M^*$  from  $R_k[A_j]$  to  $T[A_j]$  in  $G'$ . By  $R_k[A_j] = \{a_{ij_k} \mid 1 \leq i \leq n\} \subseteq X$ ,  $M^*$  is also a matching in  $G$ . But  $|M^*| = |R_k[A_j]| > u$ , which contradicts the assumption that the maximum matching in  $G$  has cardinality  $u$ . Hence, for all  $c \in \nu(\text{count}(T.A_j))$ ,  $c \leq u$ .  $\square$

The following lemma gives a necessary and sufficient condition for an element to be in  $\nu(\text{count}(T.A_j))$ .

**LEMMA 3.3** *Let  $Y = T[A_j] = \{\eta_{1j}, \eta_{2j}, \dots, \eta_{nj}\}$  and  $X = \bigcup_{1 \leq i \leq n} \nu(\eta_{ij})$ . There is a complete matching from  $S$ ,  $S \subseteq X$ , to  $Y$  in the membership graph of  $Y$  over  $S$ ,  $G' = (S \cup Y, E')$ , and  $S$  covers  $Y$  in  $G = (X \cup Y, E)$  if and only if  $|S| \in \nu(\text{count}(T.A_j))$ .*

*Proof:* ( $\Leftarrow$ ) Suppose  $|S| = c \in \nu(\text{count}(T.A_j))$ . That is, there is an interpretation  $\alpha_k = (a_{1j_k}, a_{2j_k}, \dots, a_{nj_k})$  of  $T[A_j]$  with  $R_k[A_j] = S$ . In other words, for each  $a_{lj_k} \in S$ , there exists at least a distinct  $\eta_{ij} \in T[A_j]$  such that  $a_{lj_k} \in \nu(\eta_{ij})$ . That implies, for all  $S' \subseteq S$ ,  $|S'| \leq |N(S')|$ . Therefore, by Theorem 3.1, we can construct a complete matching  $M^*$  from  $R_k[A_j] = S$  to  $Y$  in  $G'$ . Besides, since  $S = R_k[A_j] = \{a_{ij_k} \mid 1 \leq i \leq n\}$  and  $a_{ij_k} \in \nu(\eta_{ij})$ , for  $1 \leq i \leq n$ ,  $S$  covers  $Y$  in  $G$ .

( $\Rightarrow$ ) Assume there is a complete matching  $M$  from  $S$  to  $Y$  in the membership graph of  $Y$  over  $S$  and  $S$  covers  $Y$  in  $G$ . By  $S$  covers  $Y$  in  $G$ , we obtain  $S \cap \eta_{ij} \neq \emptyset$ , for all  $\eta_{ij} \in Y$ . To show that  $|S| \in \nu(\text{count}(T.A_j))$ , we can generate an interpretation  $\alpha_k = (a_{1j_k}, a_{2j_k}, \dots, a_{nj_k})$  of  $T[A_j]$  such that

$$a_{ij_k} = \begin{cases} a & \text{if } (a, \eta_{ij}) \in M, \\ \text{any } a \in S \cap \eta_{ij} & \text{otherwise.} \end{cases}$$

Then, we obtain  $|R_k[A_j]| = |S| \in \nu(\text{count}(T.A_j))$ .  $\square$

The following theorem states that all integers between the lower bound and the upper bound are all in  $\nu(\text{count}(T.A_j))$ .

**THEOREM 3.2** *Let  $Y = T[A_j] = \{\eta_{1j}, \eta_{2j}, \dots, \eta_{nj}\}$ ,  $X = \bigcup_{1 \leq i \leq n} \nu(\eta_{ij})$ , and  $G = (X \cup Y, E)$  be the membership graph of  $Y$  over  $X$ . Let  $l$  denote the cardinality of the minimum set  $S$ ,  $S \subseteq X$ , that covers  $Y$  in  $G$ , and  $u$  denote the cardinality of the maximum matching  $M$  in  $G$ , then  $\nu(\text{count}(T.A_j)) = \{c \mid l \leq c \leq u\}$ .*

*Proof:* By Lemmas 3.1 and 3.2, we obtain for all  $c \in \nu(\text{count}(T.A_j))$ ,  $l \leq c \leq u$ . Now we need to show every integer between  $l$  and  $u$  (including  $l$  and  $u$ ) is an element of

$\nu(\text{count}(T.A_j))$ . To show that  $l \in \nu(\text{count}(T.A_j))$ , since  $S$  covers  $Y$  in  $G$ , we can generate an interpretation  $\alpha_k = (a_{1j_k}, a_{2j_k}, \dots, a_{nj_k})$  of  $T[A_j]$  such that  $a_{ij_k} \in S$ , for  $1 \leq i \leq n$ . That is,  $|R_k[A_j]| \leq |S| = l$ , where  $R_k[A_j] = \{a_{ij_k} \mid 1 \leq i \leq n\}$ . By the result of Lemma 3.1,  $|R_k[A_j]| \geq l$ . Therefore, we conclude  $|R_k[A_j]| = l \in \nu(\text{count}(T.A_j))$ .

To show that  $u \in \nu(\text{count}(T.A_j))$ , we can generate an interpretation  $\alpha_k = (a_{1j_k}, a_{2j_k}, \dots, a_{nj_k})$  of  $T[A_j]$  such that

$$a_{ij_k} = \begin{cases} a & \text{if } (a, \eta_{ij}) \in M, \\ \text{any } a \in \nu(\eta_{ij}) & \text{otherwise.} \end{cases}$$

Then, we obtain  $|R_k[A_j]| \geq |M| = u$ . By the result of Lemma 3.2,  $|R_k[A_j]| \leq u$ . Therefore, we conclude  $|R_k[A_j]| = u \in \nu(\text{count}(T.A_j))$ .

Finally, we need to show that all integers between  $l$  and  $u$  are elements of  $\nu(\text{count}(T.A_j))$ . Let  $V(M)$  be the set of vertices in  $M$ . Define  $U = V(M) \cap X$ , then  $|U| = u$ . For each  $c$ ,  $l < c < u$ , there exists a set  $S'$ ,  $S \subset S' \subset U \subseteq X$ , with  $|S'| = c$ . Since  $u \in \nu(\text{count}(T.A_j))$ , by Lemma 3.3, there is a complete matching from  $U$  to  $Y$  in the membership graph of  $Y$  over  $U$ . By Theorem 3.1, there is also a complete matching from  $S'$  to  $Y$  in the membership graph of  $Y$  over  $S'$ . Moreover, since  $S$  covers  $Y$  in  $G$ ,  $S'$  also covers  $Y$  in  $G$ . By Lemma 3.3, we conclude that  $|S'| \in \nu(\text{count}(T.A_j))$ .  $\square$

By Theorem 3.2, for a set of partial values  $T[A_j] = Y = \{\eta_{1j}, \eta_{2j}, \dots, \eta_{nj}\}$ , let  $X = \bigcup_{1 \leq i \leq n} \nu(\eta_{ij})$ . We can first construct the membership graph of  $Y$  over  $X$ ,  $G = (X \cup Y, E)$ , then compute the minimum set  $S$ ,  $S \subseteq X$ , that covers  $Y$  in  $G$  and compute the maximum matching  $M$  in  $G$ , then  $\nu(\text{count}(T.A_j)) = \{c \mid l \leq c \leq u\}$ , where  $l$  and  $u$  denote the cardinality of  $A$  and  $M$ , respectively. Although the minimum cover problem is NP-complete [13], the computation of a minimal set  $S$  can be done as follows.

**ALGORITHM 3.1** An Algorithm for Computing the Minimal Set  $S \subseteq X$  That Covers  $Y$  in a Membership Graph  $G = (X \cup Y, E)$ .

**Input:** A Membership Graph  $G = (X \cup Y, E)$ .

**Output:** The Minimal Subset  $S$  of  $X$  That Covers  $Y$  in  $G$ .

1.  $S = \emptyset$ ;
2. while  $(Y \neq \emptyset)$  do {

3. Choose a vertex  $a$  from  $X$  such that  $|N(\{a\})| = \max_{a_i \in X} |N(\{a_i\})|$ ;
4.  $S = S \cup \{a\}$ ;
5.  $E = E - \{(x, y) \mid x \in X \wedge y \in N(\{a\})\}$ ;
6.  $Y = Y - N(\{a\})$ ;
7.  $X = X - \{a\}$ ;
8. }
9. Output( $S$ );

In Example 3.1, to find  $\nu(\text{count}(\text{Employees.salary}))$ , we can construct the membership graph  $G = (X \cup Y, E)$  as depicted in Figure 1, where  $X = \{20k, 60k, 75k, 80k, 90k\}$  and  $Y = \{[20k, 80k], [60k, 75k], [20k, 90k]\}$ . Thus, one of the minimal subsets of  $X$  that covers  $Y$  is  $\{20k, 60k\}$  (which is also a minimum cover set in this example, shown as the two shaded nodes) and one of the maximum matching in  $G$  is  $M = \{(20k, [20k, 80k]), (60k, [60k, 75k]), (90k, [20k, 90k])\}$  (the three pairs associated with the lines of double-ended arrows). Therefore,  $\nu(\text{count}(\text{Employees.salary})) = \{2, 3\}$ .

An  $O(n^{5/2})$  algorithm, where  $n$  is the number of vertices, for the computation of maximum matching in a bipartite graph has been developed by Hopcroft and Karp [17]. The time complexity for computing maximum matching dominates that for computing a minimal subset of  $X$  that covers  $Y$  in  $G$ . Therefore, the time complexity of our algorithm for *count* is  $O(n^{5/2})$ .

### 3.2.2 Properties of the Scalar Aggregates—*Max* and *Min*

In this subsection, we develop two dual properties for the scalar aggregates—*Max* and *Min*, respectively. These properties state that their corresponding scalar aggregates can be solved in linear time.

**THEOREM 3.3** *Let  $T$  be a partial relation on attributes  $\{A_1, A_2, \dots, A_m\}$  containing tuples  $\{t_1, t_2, \dots, t_n\}$ .*

$$\nu(\text{max}(T.A_j)) = \{a \mid a \in \bigcup_{t_i \in T} \nu(\eta_{ij}), a \geq \max_{t_i \in T} \min \nu(\eta_{ij})\}.$$

*Proof:* By Definition 3.4,  $\nu(\max(T.A_j)) = \{\max_{t_i \in R_k} a_{ijk} \mid R_k \in \text{rep}(T)\}$ . Let  $\mathcal{A} = \{\max_{t_i \in R_k} a_{ijk} \mid R_k \in \text{rep}(T)\}$  and

$$\mathcal{B} = \{a \mid a \in \bigcup_{t_i \in T} \nu(\eta_{ij}), a \geq \max_{t_i \in T} \min \nu(\eta_{ij})\}.$$

We want to show  $\mathcal{A} = \mathcal{B}$ .

(“ $\subseteq$ ”) For any  $\max_{t_i \in R_k} a_{ijk} = a \in \mathcal{A}$ , we obtain

$$\begin{aligned} a &= \max \{a_{ijk} \mid a_{ijk} \in \nu(\eta_{ij}), t_i \in T, \text{ for interpretation } \alpha_k\} \\ &= \max \{a_{ijk} \mid a_{ijk} \in \nu(\eta_{ij}), a_{ijk} \geq \min \nu(\eta_{ij}), t_i \in T, \text{ for interpretation } \alpha_k\} \\ &= \max \{a_{ijk} \mid a_{ijk} \in \nu(\eta_{ij}), a_{ijk} \geq \max_{t_i \in T} \min \nu(\eta_{ij}), t_i \in T, \text{ for interpretation } \alpha_k\} \\ &\in \{a \mid a \in \bigcup_{t_i \in T} \nu(\eta_{ij}), a \geq \max_{t_i \in T} \min \nu(\eta_{ij})\} = \mathcal{B}. \end{aligned}$$

(“ $\supseteq$ ”) Assume  $a \in \mathcal{B}$ , then we obtain  $a \geq \max_{t_i \in T} \min \nu(\eta_{ij})$  and  $a \in \bigcup_{t_i \in T} \nu(\eta_{ij})$ . Hence, there exists an  $\eta_{lj} \in T[A_j]$  such that  $a \in \nu(\eta_{lj})$ . Therefore, there exists  $\alpha_k$  subject to

$$a_{ijk} = \begin{cases} a & \text{if } i = l, \\ \min \nu(\eta_{ij}) & \text{otherwise.} \end{cases}$$

By

$$a \geq \max_{t_i \in T} \min \nu(\eta_{ij})$$

we obtain  $a = \max_{t_i \in R_k} a_{ijk}$ . By  $\max_{t_i \in R_k} a_{ijk} \in \mathcal{A}$ , we obtain  $a \in \mathcal{A}$ .  $\square$

In Example 3.1,  $\max\{\min\{20k, 80k\}, \min\{60k, 75k\}, \min\{20k, 90k\}\} = 60k$ . Therefore,  $\nu(\max(\text{Employees.salary})) = \{60k, 75k, 80k, 90k\}$ .

Similarly, we can obtain the dual property for  $\nu(\min(T.A_j))$  as follows.

**THEOREM 3.4** *Let  $T$  be a partial relation on attributes  $\{A_1, A_2, \dots, A_m\}$  containing tuples  $\{t_1, t_2, \dots, t_n\}$ .*

$$\nu(\min(T.A_j)) = \{a \mid a \in \bigcup_{t_i \in T} \nu(\eta_{ij}), a \leq \min_{t_i \in T} \max \nu(\eta_{ij})\}.$$

*Proof:* The proof can be obtained in an analogous way as that of Theorem 3.3 by replacing “max”, “min”, and “ $\geq$ ” with “min”, “max”, and “ $\leq$ ”, respectively.  $\square$

In Example 3.1,  $\min\{\max\{20k, 80k\}, \max\{60k, 75k\}, \max\{20k, 90k\}\} = 75k$ . Therefore,  $\nu(\min(\text{Employees.salary})) = \{20k, 60k, 75k\}$ .

Let  $n_\Sigma = \sum_{t_i \in T} |\nu(\eta_{ij})|$ . By Theorems 3.3 and 3.4, we can compute *maximum* and *minimum* in  $O(n_\Sigma)$ . In comparison, the time complexity of the corresponding brute-force algorithms is  $O(n_\Pi)$ , where  $n_\Pi = \prod_{t_i \in T} |\nu(\eta_{ij})|$ .

## 4 Extended Aggregate Functions over Partial Values

### 4.1 Definitions of the Extended Aggregate Functions

Extended aggregate functions can be applied to a partial relation  $T[\dots, A_x, A_y, \dots]$ , where  $A_x$  and  $A_y$  are attributes containing partial values. As a GROUP BY clause is specified, tuples in  $T$  are partitioned into groups by attribute  $A_x$  and the aggregate operation is then applied to attribute  $A_y$ . That is, they compute aggregation on attribute  $A_y$  of partial relation  $T[\dots, A_x, A_y, \dots]$  based on the distinct scalar values  $a_x$  of attribute  $A_x$ , and return  $\langle a_x, \eta_y \rangle$ , where  $\eta_y$  is the aggregation of  $A_y$  based on  $a_x$ . For databases containing partial values, the extended aggregate functions produce  $\eta_y$  as a partial value. We use the same operator names as in Section 3.1, but each extended with a *by* clause.

**DEFINITION 4.1** For a partial relation  $T[\dots, A_x, A_y, \dots]$ , the *extended count on  $T.A_y$  by  $T.A_x$* , denoted  $count(T.A_y \text{ by } T.A_x)$ , is defined as:

$$count(T.A_y \text{ by } T.A_x) = \{ \langle a_x, \eta_y \rangle \mid a_x \in \bigcup_{t \in T} \nu(t.A_x) \wedge \nu(\eta_y) = \{ | S_k[A_y] \mid R_k \in \text{rep}(T) \} \} \text{ where } S_k = \{ t \mid t \in R_k \wedge t.A_x = a_x \}.$$

**DEFINITION 4.2** For a partial relation  $T[\dots, A_x, A_y, \dots]$ , the *extended sum on  $T.A_y$  by  $T.A_x$* , denoted  $sum(T.A_y \text{ by } T.A_x)$ , is defined as:

$$sum(T.A_y \text{ by } T.A_x) = \{ \langle a_x, \eta_y \rangle \mid a_x \in \bigcup_{t \in T} \nu(t.A_x) \wedge \nu(\eta_y) = \{ \sum_{t \in S_k} t.A_y \mid R_k \in \text{rep}(T) \wedge S_k \neq \emptyset \} \} \text{ where } S_k = \{ t \mid t \in R_k \wedge t.A_x = a_x \}.$$

**DEFINITION 4.3** For a partial relation  $T[\dots, A_x, A_y, \dots]$ , the *extended average on  $T.A_y$  by  $T.A_x$* , denoted  $avg(T.A_y \text{ by } T.A_x)$ , is defined as:

$$avg(T.A_y \text{ by } T.A_x) = \{ \langle a_x, \eta_y \rangle \mid a_x \in \bigcup_{t \in T} \nu(t.A_x) \wedge \nu(\eta_y) = \{ \frac{\sum_{t \in S_k} t.A_y}{|S_k|} \mid R_k \in \text{rep}(T) \wedge S_k \neq \emptyset \} \} \text{ where } S_k = \{ t \mid t \in R_k \wedge t.A_x = a_x \}.$$

$$\nu(\eta_y) = \left\{ \frac{\sum_{t \in S_k} t.A_y}{|S_k|} \mid R_k \in \text{rep}(T) \wedge S_k \neq \emptyset \right\} \text{ where } S_k = \{t \mid t \in R_k \wedge t.A_x = a_x\}.$$

**DEFINITION 4.4** For a partial relation  $T[\dots, A_x, A_y, \dots]$ , the *extended maximum on  $T.A_y$  by  $T.A_x$* , denoted  $\text{max}(T.A_y \text{ by } T.A_x)$ , is defined as:

$$\text{max}(T.A_y \text{ by } T.A_x) = \{ \langle a_x, \eta_y \rangle \mid a_x \in \bigcup_{t \in T} \nu(t.A_x) \wedge$$

$$\nu(\eta_y) = \left\{ \max_{t \in S_k} t.A_y \mid R_k \in \text{rep}(T) \wedge S_k \neq \emptyset \right\} \text{ where } S_k = \{t \mid t \in R_k \wedge t.A_x = a_x\}.$$

**DEFINITION 4.5** For a partial relation  $T[\dots, A_x, A_y, \dots]$ , the *extended minimum on  $T.A_y$  by  $T.A_x$* , denoted  $\text{min}(T.A_y \text{ by } T.A_x)$ , is defined as:

$$\text{min}(T.A_y \text{ by } T.A_x) = \{ \langle a_x, \eta_y \rangle \mid a_x \in \bigcup_{t \in T} \nu(t.A_x) \wedge$$

$$\nu(\eta_y) = \left\{ \min_{t \in S_k} t.A_y \mid R_k \in \text{rep}(T) \wedge S_k \neq \emptyset \right\} \text{ where } S_k = \{t \mid t \in R_k \wedge t.A_x = a_x\}.$$

The following example illustrates these extended aggregate functions.

**EXAMPLE 4.1** Consider the following partial relation.

<i>Employees</i>		
<i>name</i>	<i>pos</i>	<i>salary</i>
<i>a</i>	Mgr	[40k, 45k]
<i>b</i>	[Mgr, Engr]	32k
<i>c</i>	Engr	[20k, 32k]

$\text{rep}(\text{Employees})$  contains the following elements:

<i>a</i>   Mgr   40k	<i>a</i>   Mgr   40k	<i>a</i>   Mgr   45k	<i>a</i>   Mgr   45k
<i>b</i>   Mgr   32k	<i>b</i>   Mgr   32k	<i>b</i>   Mgr   32k	<i>b</i>   Mgr   32k
<i>c</i>   Engr   20k	<i>c</i>   Engr   32k	<i>c</i>   Engr   20k	<i>c</i>   Engr   32k
<i>a</i>   Mgr   40k	<i>a</i>   Mgr   40k	<i>a</i>   Mgr   45k	<i>a</i>   Mgr   45k
<i>b</i>   Engr   32k	<i>b</i>   Engr   32k	<i>b</i>   Engr   32k	<i>b</i>   Engr   32k
<i>c</i>   Engr   20k	<i>c</i>   Engr   32k	<i>c</i>   Engr   20k	<i>c</i>   Engr   32k

The aggregate functions over  $\text{Employees.salary}$  by  $\text{Employees.pos}$  are as follows.

<i>pos</i>	<i>salary</i>
Mgr	[72k, 77k, 40k, 45k]
Engr	[20k, 32k, 52k, 64k]



$count(Employees.salary \text{ by } Employees.pos) =$	<i>pos</i>	<i>salary</i>
	Mgr	[1, 2]
	Engr	[1, 2]

$avg(Employees.salary \text{ by } Employees.pos) =$	<i>pos</i>	<i>salary</i>
	Mgr	[ $\frac{72k}{2}, \frac{77k}{2}, \frac{40k}{1}, \frac{45k}{1}$ ]
	Engr	[ $\frac{20k}{1}, \frac{32k}{1}, \frac{52k}{2}, \frac{64k}{2}$ ]

$max(Employees.salary \text{ by } Employees.pos) =$	<i>pos</i>	<i>salary</i>
	Mgr	[40k, 45k]
	Engr	[20k, 32k]

$min(Employees.salary \text{ by } Employees.pos) =$	<i>pos</i>	<i>salary</i>
	Mgr	[32k, 40k, 45k]
	Engr	[20k, 32k]

Note that in this example the result of *avg* is  $[\frac{20k}{1}, \frac{32k}{1}, \frac{52k}{2}, \frac{64k}{2}]$ ,  $\frac{64k}{2}$  can be removed since  $\frac{64k}{2} = \frac{32k}{1}$ .

## 4.2 Properties and Algorithms for Evaluating Aggregate Functions over Partial Values

In the following, we develop polynomial time algorithms to compute *count*, *maximum* and *minimum* based on their properties. These properties are generalized from those presented in the previous section.

### 4.2.1 Properties of the Aggregate Function—*Count*

Before going through the details, we need the following definitions.

**DEFINITION 4.6** Let  $T[\dots, A_x, A_y, \dots]$  be a partial relation and  $R = \bigcup_{R_k \in \text{rep}(T)} R_k$ . The *marginal membership graph of  $T[A_x, A_y]$  over  $R[A_x, A_y]$  with respect to (w.r.t.)  $a_x$* , where  $a_x \in \bigcup_{t \in T} \nu(t.A_x)$ , is a bipartite graph  $H_{a_x} = (U \cup V, E)$ , where

$$\begin{aligned}
 U &= \{\tau \mid \tau \in R[A_x, A_y], \tau.A_x = a_x\}, \\
 V &= \{t \mid t \in T[A_x, A_y], a_x \in \nu(t.A_x)\}, \text{ and} \\
 E &= \{(\tau, t) \mid \tau \in \text{rep}(t), \tau \in U, t \in V\}.
 \end{aligned}$$

**DEFINITION 4.7** Let  $T[\dots, A_x, A_y, \dots]$  be a partial relation and  $R = \bigcup_{R_k \in \text{rep}(T)} R_k$ . The *strict marginal membership graph of  $T[A_x, A_y]$  over  $R[A_x, A_y]$  w.r.t.  $a_x$* , where  $a_x \in$

$\bigcup_{t \in T} \nu(t.A_x)$ , is a bipartite graph  $\mathcal{H}_{a_x} = (\mathcal{U} \cup \mathcal{V}, \mathcal{E})$ , where

$$\begin{aligned}\mathcal{U} &= \{\tau \mid \tau \in R[A_x, A_y], \tau.A_x = a_x\}, \\ \mathcal{V} &= \{t \mid t \in T[A_x, A_y], t.A_x = a_x\}, \text{ and} \\ \mathcal{E} &= \{(\tau, t) \mid \tau \in \text{rep}(t), \tau \in \mathcal{U}, t \in \mathcal{V}\}.\end{aligned}$$

According to our definitions, for an  $H_{a_x}$  and  $\mathcal{H}_{a_x}$ , it is always true that  $\mathcal{U} = U$  and  $\mathcal{V} \subseteq V$ . Therefore, for each  $a_x$ ,  $\mathcal{H}_{a_x}$  is always a subgraph of  $H_{a_x}$ .

We now state the properties for *count* as follows.

**LEMMA 4.1** *For a partial relation  $T[\dots, A_x, A_y, \dots]$  and each  $a_x \in \bigcup_{t \in T} \nu(t.A_x)$ , let  $R = \bigcup_{R_k \in \text{rep}(T)} R_k$ ,  $\mathcal{H}_{a_x} = (\mathcal{U} \cup \mathcal{V}, \mathcal{E})$  be the strict marginal membership graph of  $T[A_x, A_y]$  over  $R[A_x, A_y]$  w.r.t.  $a_x$ , and  $l_{a_x}$  denote the cardinality of the minimum  $\mathcal{S}$ ,  $\mathcal{S} \subseteq \mathcal{U}$ , that covers  $\mathcal{V}$  in  $\mathcal{H}_{a_x}$ , then  $\langle a_x, \eta_y \rangle \in \text{count}(T.A_y \text{ by } T.A_x)$  where, for all  $c \in \nu(\eta_y)$ ,  $c \geq l_{a_x}$ .*

*Proof:* Suppose, for a  $\langle a_x, \eta_y \rangle \in \text{count}(T.A_y \text{ by } T.A_x)$ , there is an element  $c$  of  $\nu(\eta_y)$  such that  $c < l_{a_x}$ . That is, there is an  $R_k$  with  $R_k[A_x, A_y] = \mathcal{M}_k \cup \mathcal{N}_k$ , where

$$\mathcal{M}_k = \{\langle a_x, a_{y_1} \rangle, \langle a_x, a_{y_2} \rangle, \dots, \langle a_x, a_{y_c} \rangle\}, \text{ and}$$

$$\mathcal{N}_k = \{\langle \bar{a}_x, a_{y_{c+1}} \rangle, \langle \bar{a}_x, a_{y_{c+2}} \rangle, \dots\}, \text{ where } \bar{a}_x \text{ denotes those elements which are not } a_x.$$

Since, for each  $t \in \mathcal{V}$ ,  $t.A_x = a_x$ , we obtain there is a  $\langle a_x, a_{y_j} \rangle \in \mathcal{M}_k$  such that  $\langle a_x, a_{y_j} \rangle \in \text{rep}(t)$ , which implies  $\mathcal{M}_k$  covers  $\mathcal{V}$  in  $\mathcal{H}_{a_x}$ . But  $|\mathcal{M}_k| = c < l_{a_x}$ , which contradicts that the cardinality of the minimum set that covers  $\mathcal{V}$  in  $\mathcal{H}_{a_x}$  is  $l_{a_x}$ . Hence, this lemma holds.  $\square$

**LEMMA 4.2** *For a partial relation  $T[\dots, A_x, A_y, \dots]$  and each  $a_x \in \bigcup_{t \in T} \nu(t.A_x)$ , let  $R = \bigcup_{R_k \in \text{rep}(T)} R_k$ ,  $H_{a_x} = (U \cup V, E)$  be the marginal membership graph of  $T[A_x, A_y]$  over  $R$  w.r.t.  $a_x$ , and  $u_{a_x}$  denote the cardinality of the maximum matching  $M$  in  $H_{a_x}$ , then  $\langle a_x, \eta_y \rangle \in \text{count}(T.A_y \text{ by } T.A_x)$  where, for all  $c \in \nu(\eta_y)$ ,  $c \leq u_{a_x}$ .*

*Proof:* Suppose, for a  $\langle a_x, \eta_y \rangle \in \text{count}(T.A_y \text{ by } T.A_x)$ , there is an element  $c \in \nu(\eta_y)$  such that  $c > u_{a_x}$ . That is, there is an  $R_k \in \text{rep}(T)$  with  $R_k[A_x, A_y] = M_k \cup N_k$ , where

$$M_k = \{\langle a_x, a_{y_1} \rangle, \langle a_x, a_{y_2} \rangle, \dots, \langle a_x, a_{y_c} \rangle\}, \text{ and}$$

$$N_k = \{\langle \bar{a}_x, a_{y_{c+1}} \rangle, \langle \bar{a}_x, a_{y_{c+2}} \rangle, \dots\}, \text{ where } \bar{a}_x \text{ denotes those elements which are not } a_x.$$

In other words, for each  $\langle a_x, a_{y_j} \rangle \in M_k$ , there exists at least a distinct  $t \in T[A_x, A_y]$

such that  $\langle a_x, a_{y_j} \rangle \in \text{rep}(t)$  (i.e.,  $a_x \in \nu(t.A_x)$  and  $a_{y_j} \in \nu(t.A_y)$ ). That implies, for all  $S \subseteq M_k$ ,  $|S| \leq |N(S)|$  in the membership graph of  $V$  over  $M_k$ ,  $G' = (M_k \cup V, E')$ . Therefore, by Theorem 3.1, there exists a complete matching  $M^*$  from  $M_k$  to  $V$  in  $G'$ . Besides, since  $M_k \subseteq U$ ,  $M^*$  is also a matching in  $H_{a_x}$ . But  $|M^*| = |M_k| = c > u_{a_x}$ , which contradicts that the maximum matching in  $H_{a_x}$  has cardinality  $u_{a_x}$ . Therefore, this lemma holds.  $\square$

**LEMMA 4.3** *For a partial relation  $T[\dots, A_x, A_y, \dots]$  and each  $a_x \in \bigcup_{t \in T} \nu(t.A_x)$ , let*

1.  $R = \bigcup_{R_k \in \text{rep}(T)} R_k$ ,
2.  $H_{a_x} = (U \cup V, E)$  be the marginal membership graph of  $T[A_x, A_y]$  over  $R[A_x, A_y]$  w.r.t.  $a_x$ , and
3.  $\mathcal{H}_{a_x} = (\mathcal{U} \cup \mathcal{V}, \mathcal{E})$  be the strict marginal membership graph of  $T[A_x, A_y]$  over  $R[A_x, A_y]$  w.r.t.  $a_x$  (Recall that, according to our definitions,  $\mathcal{U} = U$  and  $\mathcal{V} \subseteq V$ ).

Then, there is a complete matching from  $S$ ,  $S \subseteq U$ , to  $V$  in the membership graph of  $V$  over  $S$ ,  $G' = (S \cup V, E')$  (a subgraph of  $H_{a_x}$ ), and  $S$  covers  $\mathcal{V}$  in  $\mathcal{H}_{a_x}$  if and only if  $\langle a_x, \eta_y \rangle \in \text{count}(T.A_y \text{ by } T.A_x)$  and  $|S| \in \nu(\eta_y)$ .

*Proof:* ( $\Leftarrow$ ) For a tuple  $\langle a_x, \eta_y \rangle \in \text{count}(T.A_y \text{ by } T.A_x)$ , suppose  $|S| = c \in \nu(\eta_y)$ . That is, there is an  $R_k \in \text{rep}(T)$  with  $R_k[A_x, A_y] = M_k \cup N_k$ , where

$$M_k = S = \{\langle a_x, a_{y_1} \rangle, \langle a_x, a_{y_2} \rangle, \dots, \langle a_x, a_{y_c} \rangle\}, \text{ and}$$

$$N_k = \{\langle \bar{a}_x, a_{y_{c+1}} \rangle, \langle \bar{a}_x, a_{y_{c+2}} \rangle, \dots\}, \text{ where } \bar{a}_x \text{ denotes those elements which are not } a_x.$$

In other words, for each  $\langle a_x, a_{y_j} \rangle \in S$ , there exists at least a distinct  $t \in T[A_x, A_y]$  such that  $\langle a_x, a_{y_j} \rangle \in \text{rep}(t)$  (i.e.,  $a_x \in \nu(t.A_x)$  and  $a_{y_j} \in \nu(t.A_y)$ ). That implies, for all  $S' \subseteq S$ ,  $|S'| \leq |N(S')|$  in  $G' = (S \cup V, E')$ . Therefore, by Theorem 3.1, there exists a complete matching  $M^*$  from  $S$  to  $V$  in  $G'$ . Since, for each  $t \in \mathcal{V}$ ,  $t.A_x = a_x$ , we obtain there is a  $\langle x, a_{y_j} \rangle \in S$ , such that  $\langle a_x, a_{y_j} \rangle \in \text{rep}(t)$ , which implies  $S$  covers  $\mathcal{V}$  in  $\mathcal{H}_{a_x}$ .

( $\Rightarrow$ ) Assume there is a complete matching  $M$  from  $S$  to  $V$  in  $G'$  and  $S$  covers  $\mathcal{V}$  in  $\mathcal{H}_{a_x}$  where  $S = \{\langle a_x, a_{y_1} \rangle, \langle a_x, a_{y_2} \rangle, \dots\} \subseteq U$ . To show that  $|S| \in \nu(\eta_y)$ , we can generate an interpretation  $\alpha_k(T[A_x, A_y]) = \{\langle a_{1x_k}, a_{1y_k} \rangle, \langle a_{2x_k}, a_{2y_k} \rangle, \dots, \langle a_{nx_k}, a_{ny_k} \rangle\}$

such that

$$\langle a_{ix_k}, a_{iy_k} \rangle = \begin{cases} \langle a_x, a_{y_l} \rangle & \text{if } (\langle a_x, a_{y_l} \rangle, t_i) \in M, \\ \text{any } \langle a_x, a_{y_l} \rangle \in (S \cap \text{rep}(t_i)), & \text{if } (\langle a_x, a_{y_l} \rangle, t_i) \notin M \text{ and } t_i \in \mathcal{V}, \\ \text{any } \langle a_{ix_l}, a_{iy_l} \rangle \in \text{rep}(t_i) \text{ with } a_{ix_l} \neq a_x & \text{if } t_i \in T[A_x, A_y] - \mathcal{V}. \end{cases}$$

Then, we obtain  $|\{t_i.A_y \mid t_i \in R_k, t_i.A_x = a_x\}| = |S| \in \nu(\eta_y)$ .  $\square$

**THEOREM 4.1** For a partial relation  $T[\dots, A_x, A_y, \dots]$  and each  $a_x \in \bigcup_{t \in T} \nu(t.A_x)$ , let

1.  $R = \bigcup_{R_k \in \text{rep}(T)} R_k$ ,
2.  $\mathcal{H}_{a_x} = (\mathcal{U} \cup \mathcal{V}, \mathcal{E})$  be the strict marginal membership graph of  $T[A_x, A_y]$  over  $R[A_x, A_y]$  w.r.t.  $a_x$ ,
3.  $l_{a_x}$  denote the cardinality of the minimum set  $\mathcal{S}$ ,  $\mathcal{S} \subseteq \mathcal{U}$ , that covers  $\mathcal{V}$  in  $\mathcal{H}_{a_x}$ ,
4.  $H_{a_x} = (U \cup V, E)$  be the marginal membership graph of  $T[A_x, A_y]$  over  $R[A_x, A_y]$  w.r.t.  $a_x$ , and
5.  $u_{a_x}$  denote the cardinality of the maximum matching  $M$  in  $H_{a_x}$ ,

then  $\langle a_x, \eta_y \rangle \in \text{count}(T.A_y \text{ by } T.A_x)$  and  $\nu(\eta_y) = \{c \mid l_{a_x} \leq c \leq u_{a_x}\}$ .

*Proof:* By Lemmas 4.1 and 4.2, we obtain, for each  $\langle a_x, \eta_y \rangle \in \text{count}(T.A_y \text{ by } T.A_x)$ ,  $l_{a_x} \leq c \leq u_{a_x}$ , for all  $c \in \nu(\eta_y)$ . Now we need to show every integer between  $l_{a_x}$  and  $u_{a_x}$  (including  $l_{a_x}$  and  $u_{a_x}$ ) is an element of  $\nu(\eta_y)$ .

We first show that  $l_{a_x} \in \nu(\eta_y)$ . Since  $\mathcal{S}$  covers  $\mathcal{V}$  in  $\mathcal{H}_{a_x}$ , we can generate an interpretation  $\alpha_k(T[A_x, A_y]) = \{\langle a_{1x_k}, a_{1y_k} \rangle, \langle a_{2x_k}, a_{2y_k} \rangle, \dots, \langle a_{nx_k}, a_{ny_k} \rangle\}$  such that

$$\langle a_{ix_k}, a_{iy_k} \rangle = \begin{cases} \langle a_x, a_{iy_l} \rangle \in (S \cap \text{rep}(t_i)), & \text{if } t_i \in \mathcal{V}, \\ \text{any } \langle a_{ix_l}, a_{iy_l} \rangle \in \text{rep}(t_i) \text{ with } a_{ix_l} \neq a_x & \text{if } t_i \in T - \mathcal{V}. \end{cases}$$

That is,  $|\{t_i.A_y \mid t_i \in R_k, t_i.A_x = a_x\}| = c \leq |\mathcal{S}| = l_{a_x}$ . By the result of Lemma 4.1,  $c \geq l_{a_x}$ , we conclude  $c = l_{a_x} \in \nu(\eta_y)$ .

To show that  $u_{a_x} \in \nu(\eta_y)$  by employing the maximum matching  $M$  in  $H_{a_x} = (U \cup V, E)$ , we can generate an interpretation  $\alpha_k(T[A_x, A_y]) = \{\langle a_{1x_k}, a_{1y_k} \rangle, \langle a_{2x_k}, a_{2y_k} \rangle, \dots, \langle a_{nx_k}, a_{ny_k} \rangle\}$  such that

$$\langle a_{ix_k}, a_{iy_k} \rangle = \begin{cases} \langle a_x, a_{iy_l} \rangle & \text{if } (\langle a_x, a_{y_l} \rangle, t_i) \in M, \\ \text{any } \langle a_{ix_l}, a_{iy_l} \rangle \in \text{rep}(t_i), t_i \in T[A_x, A_y] & \text{otherwise.} \end{cases}$$

Then, we obtain  $|\{t_i.A_y \mid t_i \in R_k, t_i.A_x = a_x\}| = c \geq |M| = u_{a_x}$ . By the result of Lemma 4.2,  $c \leq u_{a_x}$ . Therefore, we conclude  $c = u_{a_x} \in \nu(\eta_y)$ .

Finally, we need to show that all integers between  $l_{a_x}$  and  $u_{a_x}$  are elements of  $\nu(\eta_y)$ . Let  $V(M)$  be the set of vertices in  $M$ . Define  $W = V(M) \cap U$ , then  $|W| = u_{a_x}$ . For each  $c$ ,  $l_{a_x} < c < u_{a_x}$ , there exists a set  $\mathcal{S}'$ ,  $\mathcal{S} \subset \mathcal{S}' \subset W \subseteq U$ , with  $|\mathcal{S}'| = c$ . Since  $u_{a_x} \in \nu(\eta_y)$ , by Lemma 4.3, there is a complete matching from  $W$  to  $V$  in the membership graph of  $V$  over  $W$ . By Theorem 3.1, there is also a complete matching from  $\mathcal{S}'$  to  $V$  in the membership graph of  $V$  over  $\mathcal{S}'$ . Moreover, since  $\mathcal{S}$  covers  $\mathcal{V}$  in  $\mathcal{H}_{a_x}$ ,  $\mathcal{S}'$  also covers  $\mathcal{V}$  in  $\mathcal{H}_{a_x}$ . By Lemma 4.3, we conclude that  $|\mathcal{S}'| = c \in \nu(\eta_y)$ .  $\square$

Theorem 4.1 presents an algorithm for computing the answer of the aggregate function *count* over partial values. That is, for a partial relation  $T[\dots, A_x, A_y, \dots]$ , and each  $a_x \in \bigcup_{t \in T} \nu(t.A_x)$ , let  $R = \bigcup_{R_k \in \text{rep}(T)} R_k$ , we can

1. construct

- $\mathcal{H}_{a_x} = (\mathcal{U} \cup \mathcal{V}, \mathcal{E})$  as the *strict marginal membership graph* of  $T[A_x, A_y]$  over  $R[A_x, A_y]$  w.r.t.  $a_x$  and
- $H_{a_x} = (U \cup V, E)$  as the *marginal membership graph* of  $T[A_x, A_y]$  over  $R[A_x, A_y]$  w.r.t.  $a_x$ , then

2. compute

- $l_{a_x}$ , which denotes the cardinality of the minimum set  $\mathcal{S}$ ,  $\mathcal{S} \subseteq \mathcal{U}$ , that covers  $\mathcal{V}$  in  $\mathcal{H}_{a_x}$  and
- $u_{a_x}$ , which denotes the cardinality of the maximum matching  $M$  in  $H_{a_x}$ .

Then  $\langle a_x, \eta_y \rangle \in \text{count}(T.A_y \text{ by } T.A_x)$  and  $\nu(\eta_y) = \{c \mid l_{a_x} \leq c \leq u_{a_x}\}$ .

Although the minimum cover problem is NP-complete [13], the computation of a minimal set  $\mathcal{S}$  can be obtained by Algorithm 3.1 with input parameter  $\mathcal{H}_{a_x}$ .

In Example 4.1, to find  $\text{count}(\text{Employees.salary by Employees.pos})$ , we can construct  $H_{\text{Mgr}} = (U \cup V, E)$  as depicted in Figure 2(a), where  $U = \{\langle \text{Mgr}, 40\text{k} \rangle, \langle \text{Mgr}, 45\text{k} \rangle, \langle \text{Mgr}, 32\text{k} \rangle\}$  and  $V = \{\langle \text{Mgr}, [40\text{k}, 45\text{k}] \rangle, \langle [\text{Mgr}, \text{Engr}], 32\text{k} \rangle\}$ . Besides,  $\mathcal{H}_{\text{Mgr}} = (\mathcal{U} \cup \mathcal{V}, \mathcal{E})$  is as depicted in Figure 2(b), where  $\mathcal{U} = \{\langle \text{Mgr}, 40\text{k} \rangle, \langle \text{Mgr}, 45\text{k} \rangle, \langle \text{Mgr}, 32\text{k} \rangle\}$  and  $\mathcal{V} = \{\langle \text{Mgr}, [40\text{k}, 45\text{k}] \rangle\}$ . Thus, one of the minimal subset of  $\mathcal{U}$  that

covers  $\mathcal{V}$  is  $\{\langle \text{Mgr}, 40\text{k} \rangle\}$  (which is also a minimum cover set in this example, shown as the shaded node in Figure 2(b)) and one of the maximum matching in  $H_{\text{Mgr}}$  is  $M = \{(\langle \text{Mgr}, 40\text{k} \rangle, \langle \text{Mgr}, [40\text{k}, 45\text{k}] \rangle), (\langle \text{Mgr}, 32\text{k} \rangle, \langle [\text{Mgr}, \text{Engr}], 32\text{k} \rangle)\}$  (the two pairs associated with the lines of double-ended arrows in Figure 2(a)). Therefore,  $\langle \text{Mgr}, [1, 2] \rangle \in \text{count}(\text{Employees.salary by Employees.pos})$ . The other answer tuple for ‘Engr’ can be computed accordingly.

Since our algorithm for *count* is also dominated by the maximum matching algorithm for each element in  $\bigcup_{t \in T} \nu(t.A_x)$ , the time complexity for computing  $\text{count}(T.A_y \text{ by } T.A_x)$  is  $O(m \cdot n^{5/2})$ , where  $m = |\bigcup_{t \in T} \nu(t.A_x)|$  and  $n$  is the number of vertices.

#### 4.2.2 Properties of the Aggregate Functions—*Max* and *Min*

In this subsection, we develop two dual properties for the aggregate functions—*Max* and *Min*, respectively. Let  $\max \emptyset = \infty$  and  $\min \emptyset = -\infty$ .

**THEOREM 4.2** *For a partial relation  $T[\dots, A_x, A_y, \dots]$ ,*

$$\begin{aligned} \max(T.A_y \text{ by } T.A_x) &= \{\langle a_x, \eta_y \rangle \mid a_x \in \bigcup_{t \in T} \nu(t.A_x) \wedge \\ &\nu(\eta_y) = \{a \mid a \in \bigcup_{t \in T \wedge a_x \in \nu(t.A_x)} \nu(t.A_y) \wedge \\ &a \geq \max_{t \in T \wedge t.A_x = a_x} \min \nu(t.A_y)\}. \end{aligned}$$

*Proof:* For a given  $a_x \in \bigcup_{t \in T} \nu(t.A_x)$ , let

$$\begin{aligned} \mathcal{A} &= \left\{ \max_{t \in R_k \wedge t.A_x = a_x} t.A_y \mid R_k \in \text{rep}(T) \right\} \text{ and} \\ \mathcal{B} &= \left\{ a \mid a \in \bigcup_{t \in T \wedge a_x \in \nu(t.A_x)} \nu(t.A_y), a \geq \max_{t \in T \wedge t.A_x = a_x} \min \nu(t.A_y) \right\}. \end{aligned}$$

We want to show  $\mathcal{A} = \mathcal{B}$ .

(“ $\subseteq$ ”) For any  $\max_{t \in R_k \wedge t.A_x = a_x} t.A_y \in \mathcal{A}$ ,  $R_k \in \text{rep}(T)$ , we obtain

$$\begin{aligned} &\max_{t \in R_k \wedge t.A_x = a_x} t.A_y \\ &= \max\{t_i.A_y \mid t_i \in R_k, t_i.A_x = a_x\} \\ &= \max\{a_{iy_k} \mid a_{iy_k} \in \nu(t_i.A_y), t_i \in T, a_{ix_k} = a_x\} \\ &= \max\{a_{iy_k} \mid a_{iy_k} \in \nu(t_i.A_y), t_i \in T, a_{ix_k} = a_x, a_{iy_k} \geq \min \nu(t_i.A_y)\} \end{aligned}$$

$$\begin{aligned}
&= \max\{a_{iy_k} \mid a_{iy_k} \in \nu(t_i.A_y), t_i \in T, a_{ix_k} = a_x, a_{iy_k} \geq \max_{t \in T \wedge t.A_x = a_x} \min \nu(t.A_y)\} \\
&\in \max\{a \mid a \in \bigcup_{t \in T \wedge a_x \in \nu(t.A_x)} \nu(t.A_y), a \geq \max_{t \in T \wedge t.A_x = a_x} \min \nu(t.A_y)\} \\
&= \mathcal{B}.
\end{aligned}$$

(“ $\supseteq$ ”) Assume  $a \in \mathcal{B}$ , then we obtain

$$a \in \bigcup_{t \in T \wedge a_x \in \nu(t.A_x)} \nu(t.A_y) \text{ and } a \geq \max_{t \in T \wedge t.A_x = a_x} \min \nu(t.A_y).$$

Hence, there exists a  $t_l \in T$  such that  $a_x \in \nu(t_l.A_x)$  and  $a \in \nu(t_l.A_y)$ . Therefore, there exists an  $R_k \in \text{rep}(T)$  such that, for all  $a_{ix_k} = a_x$ ,

$$a_{iy_k} = \begin{cases} a & \text{if } i = l, \\ \min \nu(t_i.A_y) & \text{otherwise.} \end{cases}$$

By

$$a \geq \max_{t \in T \wedge t.A_x = a_x} \min \nu(t.A_y) \geq \min_{t \in T \wedge a_x \in \nu(t.A_x)} \nu(t.A_y),$$

we obtain  $\max_{t_i \in R_k \wedge a_{ix_k} = a_x} a_{iy_k} = a$ . Therefore,  $a \in \mathcal{A}$ .  $\square$

In Example 4.1, for ‘Mgr’,

$$\max_{t \in \text{Employees} \wedge t.pos = \text{‘Mgr’}} \min \nu(t.salary) = \max\{\min\{40k, 45k\}\} = 40k.$$

Therefore,  $\langle \text{Mgr}, [40k, 45k] \rangle \in \max(\text{Employees.salary by Employees.pos})$ . For ‘Engr’,

$$\max_{t \in \text{Employees} \wedge t.pos = \text{‘Engr’}} \min \nu(t.salary) = \max\{\min\{20k, 32k\}\} = 20k.$$

Therefore,  $\langle \text{Engr}, [20k, 32k, 32k] \rangle \in \max(\text{Employees.salary by Employees.pos})$ .

Similarly, for a partial relation  $T[\dots, A_x, A_y, \dots]$ , we can obtain the dual property for  $\min(T.A_y \text{ by } T.A_x)$  as follows.

**THEOREM 4.3** For a partial relation  $T[\dots, A_x, A_y, \dots]$ ,

$$\min(T.A_y \text{ by } T.A_x) = \{\langle a_x, \eta_y \rangle \mid a_x \in \bigcup_{t \in T} \nu(t.A_x) \wedge$$

$$\nu(\eta_y) = \{a \mid a \in \bigcup_{t \in T \wedge a_x \in \nu(t.A_x)} \nu(t.A_y) \wedge$$

$$a \leq \min_{t \in T \wedge t.A_x = a_x} \max \nu(t.A_y)\}.$$

*Proof:* The proof can be obtained in an analogous way as that of Theorem 4.2 by replacing “max”, “min”, and “ $\geq$ ” with “min”, “max”, and “ $\leq$ ”, respectively.  $\square$

In Example 4.1, for ‘Mgr’,

$$\min_{t \in \text{Employees} \wedge t.pos = \text{'Mgr'}} \max \nu(t.salary) = \min\{\max\{40k, 45k\}\} = 45k.$$

Therefore,  $\langle \text{Mgr}, [32k, 40k, 45k] \rangle \in \min(\text{Employees.salary by Employees.pos})$ . For ‘Engr’,

$$\min_{t \in \text{Employees} \wedge t.pos = \text{'Engr'}} \max \nu(t.salary) = \min\{\max\{20k, 32k\}\} = 32k.$$

Therefore,  $\langle \text{Engr}, [20k, 32k, 32k] \rangle \in \min(\text{Employees.salary by Employees.pos})$ .

For the time complexity of *maximum* and *minimum*, let  $n_\Sigma = \sum_{t \in T} |\nu(t.A_y)|$ . By Theorems 4.2 and 4.3, we can compute *maximum* and *minimum* in  $O(n_\Sigma)$  for each  $a_x \in \bigcup_{t \in T} \nu(t.A_x)$ . Therefore, the time complexity for our *minimum* and *maximum* are  $O(m_\Sigma \cdot n_\Sigma)$ , where  $m_\Sigma = |\bigcup_{t \in T} \nu(t.A_x)|$ .

## 5 Discussion and Future Studies

Partial values have been used to represent imprecise data in databases. Referring to the work of DeMichiel [11] on resolving domain mismatch problems in multidatabase systems by partial values, data imprecision may come from their unavailability or data/schema incompatibilities in a multidatabase system. Therefore, in addition to imprecise data manipulation, research work on partial values is also important for the interoperability of a multidatabase system.

Since, in practice, relational algebra or calculus are inadequate for many important applications involving statistical information or aggregations, modern query languages are equipped with some useful aggregate operations. In this paper, we define a set of extended aggregate operations, namely *sum*, *average*, *count*, *maximum*, and *minimum*, which can be applied to attributes containing partial values.

Since the number of the interpretations of a set of partial values can be very large, the properties presented in this paper can be used to speed up the computations of *count*, *maximum* and *minimum*. By the properties presented in Section 3.2, we know that *maximum* and *minimum* can be solved in polynomial time. For *count*, we show



that the lower bound is equal to the cardinality of the minimum cover set and the upper bound is equal to the cardinality of the maximum matching. The maximum matching can be found in polynomial time while the minimum cover problem is NP-complete. A greedy algorithm is given to find a minimal cover set. For *sum* and *average*, we point out that the cardinalities of  $\nu(\text{sum}(T.A))$  and  $\nu(\text{avg}(T.A))$  are equal to the cardinality of  $\prod_{t \in T} |\nu(t.A)|$  in the worst case and no polynomial time algorithms can be found.

The *count* operation actually corresponds to those in SQL with the keyword DISTINCT. For those in SQL with the keyword ALL, *count* can be obtained by directly counting the number of tuples in the involved relation.

Although Imieliński and Vadaparty [18, 19] pointed out that if partial values are allowed to occur in databases, the data complexity of query processing jumps from polynomial time to CoNP [13], there also exist some types of queries which have polynomial time complexity [18]. Our studies on the query processing over partial values intend to discover more polynomial time algorithms from algebraic point of view. In our recent work, we found the *division* (by restricting the divisor to be definite) [31] and the *projection with redundancy elimination* [29] over partial values can be done in polynomial time. Besides, we generalized partial values into *probabilistic partial values* in [30]. Similar properties for probabilistic partial values are also being studied.

### **Acknowledgement**

The authors would like to thank the anonymous referees for their valuable suggestions and comments.

## References

- [1] D. Barbará, H. Garcia-Molina and D. Porter, "The Management of Probabilistic Data," *IEEE Trans. Knowledge and Data Engineering*, vol. 4, no. 4, pp. 487-501, 1992.
- [2] J.A. Bondy and U.S.R. Murty, *Graph Theory with Applications*, Macmillan Press, NY, 1976.
- [3] P. Bosc, M. Galibourg and G. Hamon, "Fuzzy Querying with SQL: Extensions and Implementation Aspects," *Fuzzy Sets and Systems*, vol. 28, pp. 333-349, 1988.
- [4] B.P. Buckles and F.E. Petry, "A Fuzzy Representation of Data for Relational Databases," *Fuzzy Sets and Systems*, vol. 7, pp. 213-226, 1982.
- [5] R. Cavallo and M. Pittarelli, "The Theory of Probabilistic Databases," in *Proc. VLDB Conf.*, 1987, pp. 71-81.
- [6] E.F. Codd, "Understanding Relations," Installment no. 7, *ACM SIGMOD Record FDT Bulletin*, vol. 7, no. 3-4, pp. 23-28, 1975.
- [7] E.F. Codd, "Extending the Database Relational Model to Capture More Meaning," *ACM Trans. Database Systems*, vol. 4, no. 4, pp. 397-434, 1979.
- [8] E.F. Codd, "Missing Information (Applicable and Inapplicable) in Relational Databases," *ACM SIGMOD Record*, vol. 15, no. 4, pp. 53-78, 1986.
- [9] C.J. Date, *A Guide to the SQL Standard*, Addison-Wesley, MA, 1989.
- [10] L.G. DeMichiel, "Performing Database Operations over Mismatched Domains," Ph.D. dissertation, Stanford Univ., Stanford, CA, 1989.
- [11] L.G. DeMichiel, "Resolving Database Incompatibility: An Approach to Performing Relational Operations over Mismatched Domains," *IEEE Trans. Knowledge and Data Engineering*, vol. 1, no. 4, pp. 485-493, 1989.
- [12] D. Dubois and H. Prade, *Possibility Theory: An Approach to Computerized Processing*, Plenum Press, NY, 1986.
- [13] M.R. Garey and D.S. Johnson, *Computers and Intractability: A Guide to the Theory of NP-Completeness*, Freeman, SF, 1979.
- [14] J. Grant, "Partial Values in a Tabular Database Model," *Information Processing Letters*, vol. 9, no. 2, pp. 97-99, 1979.
- [15] J. Grant and J. Minker, "Answering Queries in Indefinite Databases and the Null Value Problem," in *Advances in Computing Research*, P. Kanellakis, Eds., vol. 3, *The Theory of Databases*, JAI Press, 1986, pp. 247-267.
- [16] P. Hall, "On Representatives of Subsets," *J. London Math. Soc.*, vol. 10, pp. 26-30, 1935.
- [17] J.E. Hopcroft and R.M. Karp, "An  $n^{5/2}$  Algorithm for Maximum Matching in Bipartite Graphs," *SIAM J. Computing*, vol. 2, no. 4, pp. 225-231, 1973.
- [18] T. Imieliński and K. Vadaparty, "Complexity of Query Processing in Databases with OR-Objects," in *Proc. ACM Symp. on Principles of Database Systems*, 1989, pp. 51-65.

- [19] T. Imieliński, “Incomplete Deductive Databases,” *Annals of Mathematics and Artificial Intelligence*, vol. 3, no. 2-4, pp. 259-294, 1984.
- [20] W. Lipski, “On Semantic Issues Connected with Incomplete Information Systems,” *ACM Trans. Database Systems*, vol. 4 no. 3, pp. 262-296, 1979.
- [21] D. Maier, *The Theory of Relational Databases*, Computer Science Press, Rockville, MD, 1983.
- [22] A. Motro, “Accommodating Imprecision in Database Systems: Issues and Solutions,” *ACM SIGMOD RECORD*, vol. 19, no. 4, pp. 69-74, 1990.
- [23] A. Ola, “Relational Databases with Exclusive Disjunctions,” in *Proc. 8th IEEE Int. Conf. Data Engineering*, 1992, pp. 328-336.
- [24] G. Özsoyoğlu, Z.M. Özsoyoğlu, and V. Matos, “Extending Relational Algebra and Relational Calculus with Set-Valued Attributes and Aggregate Functions,” *ACM Trans. Database Systems*, vol. 12, no. 4, pp. 566-592, 1987.
- [25] H. Prade and C. Testemale, “Generalizing Database Relational Algebra for the Treatment of Incomplete/Uncertain Information and Vague Queries,” *Information Sciences*, vol. 34, pp. 115-143, 1984.
- [26] E.A. Rundensteiner and L. Bic, “Aggregates in Possibilistic Databases,” in *Proc. 15th Int. Conf. Very Large Data Bases*, 1989, pp. 287-294.
- [27] E.A. Rundensteiner and L. Bic, “Evaluating Aggregates in Possibilistic Relational Databases,” *Data & Knowledge Engineering*, vol. 7, no. 3, pp. 239-267, 1992.
- [28] P.S.M. Tsai and A.L.P. Chen, “Querying Uncertain Data in Heterogeneous Databases,” in *Proc. IEEE Int’l Workshop on Research Issues on Data Engineering (RIDE)*, 1993, pp. 161-168.
- [29] F.S.C. Tseng, A.L.P. Chen and W.P. Yang, “Searching a Minimal Semantically-Equivalent Subset of a Set of Partial Values,” *The VLDB Journal*, vol. 2, no. 4, pp. 489-512, 1993.
- [30] F.S.C. Tseng, A.L.P. Chen, and W.P. Yang, “Answering Heterogeneous Database Queries with Degrees of Uncertainty,” *Distributed and Parallel Databases: An International Journal*, vol. 1, no. 3, pp. 281-302, 1993.
- [31] F.S.C. Tseng, A.L.P. Chen, and W.P. Yang, “Generalizing the Division Operation on Indefinite Databases,” in *Proc. 2nd Far-East Workshop on Future Database Systems*, Kyoto, Japan, 1992, pp. 347-354.
- [32] F.S.C. Tseng, A.L.P. Chen, and W.P. Yang, “Refining Imprecise Data by Integrity Constraints,” *Data and Knowledge Engineering Journal*, vol. 11, no. 3, pp. 299-316, 1993.
- [33] L.A. Zadeh, “Fuzzy Sets as a basis for a theory of possibility,” *Fuzzy Sets and Systems*, vol. 1, no. 1, pp. 3-28, 1978.
- [34] M. Zemankova and A. Kandel, “Implementing Imprecision in Information Systems,” *Information Science*, vol. 37, pp. 107-141, 1985.

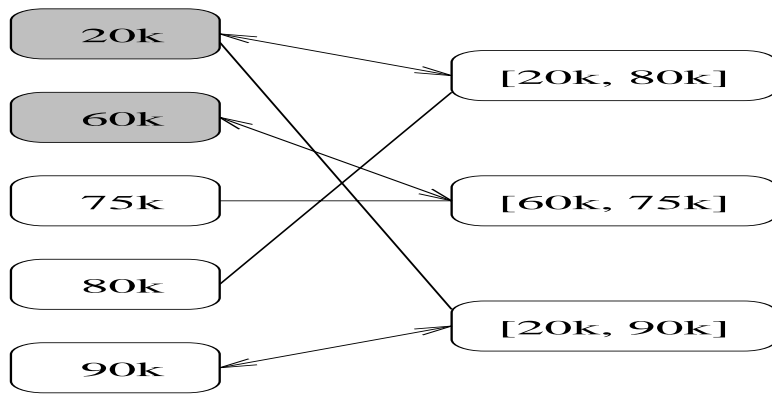


Figure 1:

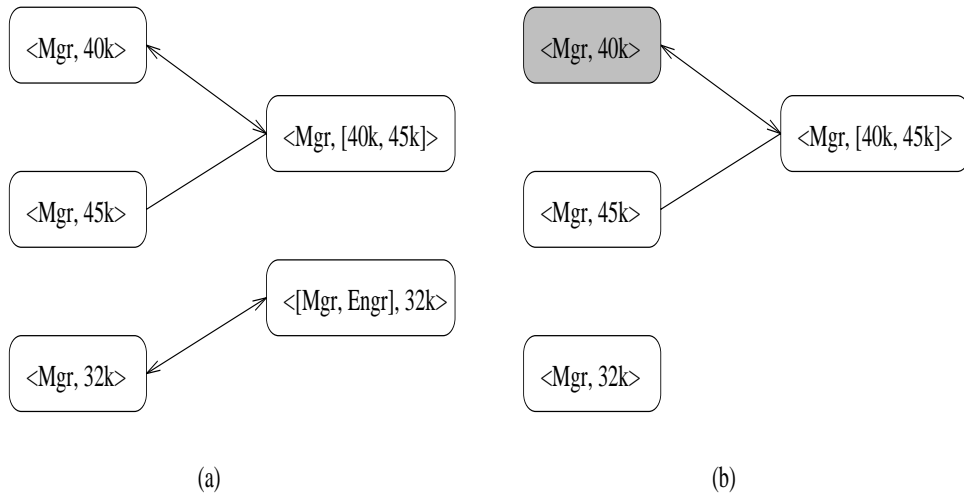


Figure 2: