

Adding Context to Preferences*

Kostas Stefanidis, Evaggelia Pitoura and Panos Vassiliadis
Computer Science Department, University of Ioannina, GR-45110 Ioannina, Greece
{kstef, pitoura, pvassil}@cs.uoi.gr

Abstract

To handle the overwhelming amount of information currently available, personalization systems allow users to specify the information that interests them through preferences. Most often, users have different preferences depending on context. In this paper, we introduce a model for expressing such contextual preferences. Context is modeled as a set of multidimensional attributes. We formulate the context resolution problem as the problem of (a) identifying those preferences that qualify to encompass the context state of a query and (b) selecting the most appropriate among them. We also propose an algorithm for context resolution that uses a data structure, called the profile tree, that indexes preferences based on their associated context. Finally, we evaluate our approach from two perspectives: usability and performance.

1. Introduction

Today, a very large and steadily increasing amount of information is available to a wide spectrum of users, thus creating the need for personalized information processing. Instead of overwhelming the users with all available data, a personalized query returns only the information that is of interest to them [10]. In general, to achieve personalization, users express their preferences on specific pieces of data either explicitly or implicitly. However, often users may have different preferences under different circumstances. For instance, a user visiting Athens may prefer to visit *Acropolis* in a nice sunny summer day and the *archaeological museum* in a cold and rainy winter afternoon. In other words, the results of a preference query may depend on context.

Context is a general term used to capture any information that can be used to characterize the situations of an entity [4]. Common types of context include the *computing context* (e.g., network connectivity, resources), the *user context* (e.g., profile, location), the *physical context* (e.g., noise levels, temperature) and *time* [2]. We model context as a set of multidimensional context parameters. A context

state corresponds to an assignment of values to context parameters. Different levels of abstraction for the captured context data are introduced by allowing context parameters to take values from hierarchical domains. For instance, the context parameter *location* may take values from a *region*, *city* or *country* domain. Users employ context descriptors to express preferences on specific database instances for a variety of context states expressed with varying levels of detail.

Each query is associated with one or more context state. The context state of a query may, for example, be the current state at the time of its submission. Furthermore, a query may be explicitly enhanced with context descriptors to allow exploratory queries about hypothetical context states. We formulate the *context resolution problem* that refers to the problem of identifying those preferences whose context states are the most relevant to the context state of the query. The problem can be divided into two steps: (a) the identification of all the candidate context states that encompass the query state and (b) the selection of the most appropriate state among these candidates. The first subproblem is resolved through the notion of the “covers” partial order between states that relates context states expressed at different levels of abstraction. For instance, the notion of coverage allows relating a context state in which location is expressed at the level of a *city* and a context state in which location is expressed at the level of a *country*. To resolve the second subproblem, we consider two distance metrics that capture similarity between context states.

Our algorithm for context resolution uses a *profile tree* that indexes user preferences based on their associated context. Intuitively, the algorithm starts from the query context and incrementally “extends” its coverage until a matching state is found in the profile tree. Finally, we evaluate our approach from two perspectives: usability and performance.

The rest of this paper is organized as follows. In Section 2, we present our reference example. In Section 3, we introduce our context and preference model and the profile tree. In Section 4, we focus on processing contextual queries, while in Section 5, we present our evaluation results. Section 6 describes related work. Finally, Section 7 concludes the paper with a summary of our contributions.

*Work partially supported by the Greek General Secretariat for Research and Technology through PENED 03-ED-591

2. Reference Example

We consider a simple database that maintains information about *points_of_interest*. *Points_of_interest* may for example be museums, monuments, archaeological places or zoos. The database schema consists of a single database relation: *Points_of_Interest*(*pid*, *name*, *type*, *location*, *open-air*, *hours_of_operation*, *admission_cost*). We consider three context parameters as relevant: *location*, *temperature* and *accompanying_people*. Depending on context, users prefer *points_of_interest* that have specific attribute values. Such preferences are expressed by providing a numeric score between 0 and 1. For instance, the interest score of a preference that is related to the *type* of the visiting place depends on the *accompanying_people* that might be *friends*, *family*, or *alone*. For example, a *zoo* may be a preferred place to visit than a *brewery* in the context of *family*.

3. Context and Preference Model

First, we present the fundamental concepts related to context modeling, and then, define user preferences.

3.1. Modeling Context

Context is modeled through a finite set of special-purpose attributes, called *context parameters* (C_i). In particular, for a given application X , we define its context environment CE_X as a set of n context parameters $\{C_1, C_2, \dots, C_n\}$. For instance, the context environment of our example is $\{location, temperature, accompanying_people\}$. Each context parameter C_i is characterized by a *context domain*, $dom(C_i)$. As usual, a *domain* is an infinitely countable set of values. A *context state* w is a n -tuple of the form (c_1, c_2, \dots, c_n) , where $c_i \in dom(C_i)$. For instance, a context state in our example may be: $(Plaka, warm, friends)$. The set of all possible context states called *world*, W , is the Cartesian product of the domains of the context attributes: $W = dom(C_1) \times dom(C_2) \times \dots \times dom(C_n)$.

To allow more flexibility in defining preferences, we model context parameters as multidimensional attributes. In particular, we assume that each context parameter participates in an associated *hierarchy of levels* of aggregated data, i.e., it can be viewed from different levels of detail. Formally, an *attribute hierarchy* is a lattice (L, \prec) : $L = (L_1, \dots, L_{m-1}, ALL)$ of m levels and \prec is a partial order among the levels of L such that $L_1 \prec L_i \prec ALL$, for every $1 < i < m$. We require that the upper bound of the lattice is always the level ALL , so that we can group all values into the single value ‘all’. The lower bound of the lattice is called the *detailed level* of the context parameter. We use the notation $dom_{L_j}(C_i)$ for the domain of level L_j of parameter C_i . For the domain of the detailed level, we shall use both $dom_{L_1}(C_i)$ and $dom(C_i)$ interchangeably.

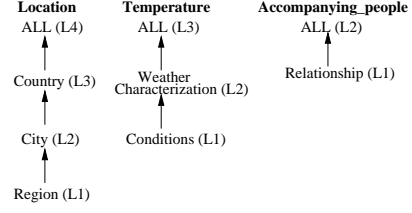


Figure 1. Example Hierarchies.

Regarding our running example, levels of location are *Region* (L_1), *City* (L_2), *Country* (L_3) and *ALL* (L_4). For example, a value of *City* is *Athens* and two *Regions* in *Athens* are *Plaka* and *Kifisia*. For weather, there are three levels: the detailed level *Conditions* (L_1) whose domain includes *freezing*, *cold*, *mild*, *warm* and *hot*, the level *Weather Characterization* (L_2) which just refers to whether the weather is *good* (grouping *mild*, *warm* and *hot*) or *bad* (grouping *freezing* and *cold*) and the level *ALL* (L_3). Finally, the context parameter *accompanying_people* has the lower level *Relationship* (L_1) which consists of the values *friends*, *family*, *alone* and the level *ALL* (L_2). Figure 1 depicts these hierarchies.

The relationship between the values of the context levels is achieved through the use of the set of $anc_{L_i}^{L_j}$, $L_i \prec L_j$, functions [17]. A function $anc_{L_i}^{L_j}$ assigns a value of the domain of L_i to a value of the domain of L_j . For instance, $anc_{Region}^{City}(Plaka) = Athens$. Formally, the set of functions $anc_{L_i}^{L_j}$ satisfies the following conditions:

1. For each pair of levels L_1 and L_2 such that $L_1 \prec L_2$, the function $anc_{L_1}^{L_2}$ maps each element of $dom_{L_1}(C_i)$ to an element of $dom_{L_2}(C_i)$.
2. Given levels L_1, L_2 and L_3 such that $L_1 \prec L_2 \prec L_3$, the function $anc_{L_1}^{L_3}$ is equal to the composition $anc_{L_1}^{L_2} \circ anc_{L_2}^{L_3}$.
3. For each pair of levels L_1 and L_2 such that $L_1 \prec L_2$, the function $anc_{L_1}^{L_2}$ is monotone, i.e., $\forall x, y \in dom_{L_1}(C_i)$, $L_1 \prec L_2$, $x < y \Rightarrow anc_{L_1}^{L_2}(x) \leq anc_{L_1}^{L_2}(y)$.

The function $desc_{L_1}^{L_2}$ is the inverse of $anc_{L_1}^{L_2}$, that is $desc_{L_1}^{L_2}(v) = \{x \in dom_{L_1}(C_i) : anc_{L_1}^{L_2}(x) = v\}$. For instance, $desc_{Region}^{City}(Athens) = \{Plaka, Kifisia\}$. We use $L_1 \preceq L_2$ between two levels to mean $L_1 \prec L_2$ or $L_1 = L_2$.

We define the extended domain for a parameter C_i with m levels as $edom(C_i) = \cup_{j=1}^m dom_{L_j}(C_i)$. Then, an (*extended*) *context state* is an assignment of values to context parameters from their extended domain. In particular, an extended context state s is a n -tuple of the form (c_1, c_2, \dots, c_n) , where $c_i \in edom(C_i)$. For instance, a context state in our example may be $(Greece, warm, friends)$ or $(Greece, good, all)$. The set of all possible extended

context states called *extended world*, EW , is the Cartesian product of the extended domains of the context attributes: $EW = \text{edom}(C_1) \times \text{edom}(C_2) \times \dots \times \text{edom}(C_n)$.

Users can express conditions regarding the values of a context parameter through *context descriptors*.

Definition 1 (Context parameter descriptor) A context parameter descriptor $\text{cod}(C_i)$ for a parameter C_i is an expression of the form:

1. $C_i = v$, where $v \in \text{edom}(C_i)$, or
2. $C_i \in \{\text{value}_1, \dots, \text{value}_m\}$, where $\text{value}_k \in \text{edom}(C_i)$, $1 \leq k \leq m$, or
3. $C_i \in [\text{value}_1, \text{value}_m]$, where $[\text{value}_1, \text{value}_m]$ denotes a range of values $x \in \text{edom}(C_i)$, such that $\text{value}_1 \leq x \leq \text{value}_m$.

For example, given a context parameter *location*, a context parameter descriptor can be of the form $\text{location} = \text{Plaka}$, or $\text{location} \in \{\text{Plaka}, \text{Acropolis}\}$. Given a context parameter *temperature*, a range-based context parameter descriptor can be of the form $\text{temperature} \in [\text{mild}, \text{hot}]$, signifying thus the set of values $\{\text{mild}, \text{warm}, \text{hot}\}$. There is a straightforward way to translate context parameter descriptors to sets of values. Practically, this involves translating range descriptors to sets of values (recall that all domains are infinitely countable, hence, they are not dense and all ranges can be translated to finite sets of values).

Definition 2 (Context of a context parameter descriptor) Given a context parameter descriptor $c = \text{cod}(C_i)$ for a parameter C_i , its context is a finite set of values, computed as follows:

$$\text{Context}(c) = \begin{cases} \{v\} & \text{if } c \text{ of the form } C_i = v \\ \{v_1, \dots, v_m\} & \text{if } c \text{ of the form } \\ & C_i \in \{v_1, \dots, v_m\} \\ \{v_1, \dots, v_m\} & \text{if } c \text{ of the form } \\ & C_i \in [v_1, v_m] \end{cases}$$

A context descriptor is a specification that a user can make for a set of context parameters, through the combination of simple parameter descriptors.

Definition 3 (Composite context descriptor) A (composite) context descriptor cod is a formula $\text{cod}(C_{i_1}) \wedge \text{cod}(C_{i_2}) \wedge \dots \wedge \text{cod}(C_{i_k})$ where each C_{i_j} , $1 \leq j \leq k$ is a context parameter and there is at most one parameter descriptor per context parameter C_{i_j} .

Given a set of context parameters C_1, \dots, C_n , a composite context descriptor describes a set of possible context states, with each state having a specific value for each parameter. Clearly, one context descriptor can produce more than one state. The production of these states can be performed by computing the Cartesian product of the context states of all the individual parameter descriptors of a context descriptor. If a context descriptor does not contain all context parameters, that means that the absent context parameters have irrelevant values. This is equivalent to a condition

$C_i = \text{all}$. Observe that the set of produced states is finite, due to the finite character of the context of the parameter descriptors.

Definition 4 (Context of a context descriptor) Assume a set of context parameters C_1, \dots, C_n and a context descriptor $\text{cod} = \text{cod}(C_{i_1}) \wedge \text{cod}(C_{i_2}) \wedge \dots \wedge \text{cod}(C_{i_k})$, $0 \leq k \leq n$. Without loss of generality, we assume that the parameters without a parameter descriptor are the last $n - k$ ones. The context states of a context descriptor, called $\text{Context}(\text{cod})$ are defined as:

$$\text{Context}(\text{cod}(C_{i_1})) \times \dots \times \text{Context}(\text{cod}(C_{i_k})) \times \{\text{all}\} \times \dots \times \{\text{all}\}$$

Suppose for instance, the context descriptor ($\text{location} = \text{Plaka} \wedge \text{temperature} \in \{\text{warm}, \text{hot}\} \wedge \text{accompanying_people} = \text{friends}$). This descriptor corresponds to the following two context states: $(\text{Plaka}, \text{warm}, \text{friends})$ and $(\text{Plaka}, \text{hot}, \text{friends})$.

3.2. Contextual Preferences

In this section, we define how context affects the results of queries, so that the same query returns different results based on the context of its execution. Such context-aware personalization is achieved through the use of preferences. In general, there are two different approaches for expressing preferences: a quantitative and a qualitative one. With the *quantitative approach* (e.g., [1]), preferences are expressed indirectly by using scoring functions that associate a numeric score with every tuple of the query answer. In the *qualitative approach* (such as the work in [3]), preferences between tuples in the query answer are specified directly, typically using binary preference relations.

Although, our context model can be used with both quantitative and qualitative approaches, we use a simple quantitative preference model to demonstrate the basic issues underlying contextualization. In particular, users express their preference for specific database instances by providing a numeric score which is a real number between 0 and 1. This score expresses their degree of interest. Value 1 indicates extreme interest, while value 0 indicates no interest. Interest is expressed for specific values of non context attributes of a database relation, for instance for the various attributes (e.g., *type*, *location*) of our *Point_of_Interest* database relation. In particular, a *contextual preference* is defined as follows.

Definition 5 (Contextual preference) A contextual preference is a triple of the form $\text{contextual_preference} = (\text{cod}, \text{attributes_clause}, \text{interest_score})$, where cod is a context descriptor, the *attributes_clause* $\{A_1\theta_1a_1, A_2\theta_2a_2, \dots, A_k\theta_ka_k\}$ specifies a set of attributes A_1, A_2, \dots, A_k with their values a_1, a_2, \dots, a_k with $a_i \in \text{dom}(A_i)$, $\theta_i \in \{=, <, >, \leq, \geq, \neq\}$ and *interest_score* is a real number between 0 and 1.

The meaning is that in the set of context states specified by cod , all database tuples (instances) for which

the attributes A_1, A_2, \dots, A_m have respectively values a_1, a_2, \dots, a_m are assigned the indicated interest score. Since our focus in this paper is on context descriptors, we further simplify our model, so that in the following, we shall use *attributes_clauses* with a single attribute A of the form $A = a$, for $a \in \text{dom}(A)$. Further, we assume that for tuples for which more than one preference applies, appropriate combining preference functions exist [1, 14].

In our reference example, there are three context parameters *location*, *temperature* and *accompanying_people*. As non-context parameters, we use the attributes of the relation *Points_of_Interest*. For example, consider that a user wants to express the fact that, when she is at *Plaka* and the weather is *warm*, she likes to visit *Acropolis*. This may be expressed through the following contextual preference: $\text{contextual_preference}_1 = ((\text{location} = \text{Plaka} \wedge \text{temperature} = \text{warm}), (\text{name} = \text{Acropolis}), 0.8)$. Similarly, she may also express the fact that when she is with friends, she likes to visit breweries through preference $\text{contextual_preference}_2 = ((\text{accompanying_people} = \text{friends}), (\text{type} = \text{brewery}), 0.9)$. More involved context descriptors may be used as well, for example, $\text{contextual_preference}_3 = ((\text{location} = \text{Plaka} \wedge \text{temperature} \in \{\text{warm}, \text{hot}\}), (\text{name} = \text{Acropolis}), 0.8)$, where the context descriptor is $\text{cod} = (\text{location} = \text{Plaka} \wedge \text{temperature} \in \{\text{warm}, \text{hot}\})$.

A contextual preference may conflict with another one. For example, assume that a user defines that she prefers to visit *Acropolis* in a nice sunny day, giving a high score of 0.8 to this preference. If, later on, she gives the interest score 0.3 to the same preference, this will cause a conflict.

Definition 6 (Conflicting preferences) A

$\text{contextual_preference}_i = (\text{cod}_i, (A_i = a_i), \text{interest_score}_i)$ conflicts with a $\text{contextual_preference}_j = (\text{cod}_j, (A_j = a_j), \text{interest_score}_j)$ if and only if:

1. $\text{Context}(\text{cod}_i) \cap \text{Context}(\text{cod}_j) \neq \emptyset$, and
2. $A_i = A_j$ and $a_i = a_j$, and
3. $\text{interest_score}_i \neq \text{interest_score}_j$.

Such conflicting preferences are detected when users enter their preferences. Finally, we define *profile P* as:

Definition 7 (Profile) A *profile P* is a set of non-conflicting contextual preferences.

3.3. The Profile Tree

Contextual preferences are stored in a hierarchical data structure called *profile tree*. Let P be a profile and S be the set of context states of all context descriptors that appeared in P . The basic idea is to store in the profile tree the context states in S such that there is exactly one path in the tree for each context state $s \in S$. Specifically, the profile tree for P is a directed acyclic graph with a single root node and at most $n+1$ levels. Each one of the first n levels corresponds to a context parameter and the last one

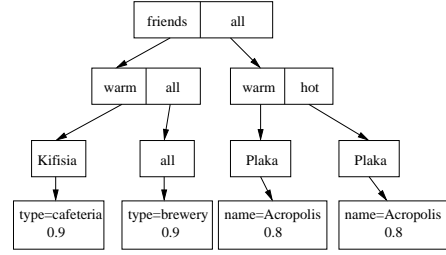


Figure 2. An instance of a profile tree.

to the leaf nodes. For simplicity, assume that context parameter C_i is mapped to level i . Each non-leaf node at level k maintains cells of the form $[key, pointer]$, where $key \in \text{edom}(C_k)$ for some value of c_k that appeared in a state $s \in S$. No two cells within the same node contain the same key value. The pointer points to a node at level $k+1$ having cells with key values in $\text{edom}(C_{k+1})$ which appeared in the same state s with the *key*. Each leaf node is a set of cells of the form $[attribute = value, interest_score]$, where $[attribute = value, interest_score]$ is the one associated with the path (state) leading to it.

For example, assume an instance of a profile P consisting of the following preferences: $\text{contextual_preference}_1 = ((\text{location} = \text{Kifisia} \wedge \text{temperature} = \text{warm} \wedge \text{accompanying_people} = \text{friends}), (\text{type} = \text{cafeteria}), 0.9)$, $\text{contextual_preference}_2 = ((\text{accompanying_people} = \text{friends}), (\text{type} = \text{brewery}), 0.9)$, and $\text{contextual_preference}_3 = ((\text{location} = \text{Plaka} \wedge \text{temperature} \in \{\text{warm}, \text{hot}\}), (\text{name} = \text{Acropolis}), 0.8)$. Assume further that *accompanying_people* is assigned to the first level of the tree, *temperature* to the second and *location* to the third one. Fig. 2 depicts the profile tree for P .

The way that the context parameters are assigned to the levels of the context tree affects its size. Let $m_i, 1 \leq i \leq n$, be the cardinality of the domain, then the maximum number of cells is $m_1 * (1 + m_2 * (1 + \dots (1 + m_n)))$. The above number is as small as possible, when $m_1 \leq m_2 \leq \dots \leq m_n$, thus, it is better to place context parameters with domains with higher cardinalities lower in the context tree.

4. Contextual Preference Queries

In this section, we define contextual queries, formulate the problem of identifying the preferences that are most relevant to a query and present an algorithm that locates them.

4.1. Contextual Queries

A contextual query is a query enhanced with information regarding context. Implicitly, the context associated with a contextual preference query is the current context, that is, the context surrounding the user at the time of the submission of the query. The current context should correspond to a single context state, where each of the values of the

context parameters takes a specific value from its most detailed domain. However, in some cases, it may be possible to specify the current context using only rough values, for example, when the values of some context parameters are provided by sensor devices with limited accuracy. In this case, a context parameter may take a single value from a higher level of the hierarchy or even more than one value. Besides the implicit context, we also consider queries that are explicitly augmented with an extended context descriptor. For example, a user may want to pose an exploratory query of the form: “When I travel to Athens with my family this summer (implying good weather), what places should I visit?”. Formally,

Definition 8 (Contextual query) A contextual query CQ is a query Q enhanced with a context descriptor denoted cod^Q .

Now, the problem is: given the cod^Q of a contextual query CQ and a user profile P , identify the contextual preferences that are the most relevant to the context states specified by cod^Q . Next, we first formalize the problem and then, provide a procedure for locating such preferences.

4.2. Context State of a Query

Assume a contextual query CQ enhanced with a context descriptor of the form $cod^Q = (location = Athens \wedge weather = warm)$ and a simple profile $P = \{((location = Greece \wedge weather = warm), attributes_clause, interest_score_1), ((location = Europe \wedge weather = warm), attributes_clause, interest_score_2)\}$. Intuitively, we are seeking for a context descriptor in P that is more general than the query descriptor. Both context descriptors in P satisfy this requirement, however, the first one is more “specific” and should be the one used. Next, we formalize the notion of a set of states covering another one.

Definition 9 (Covering context state) An extended context state $s^1 = (c_1^1, c_2^1, \dots, c_n^1) \in EW$ covers an extended context state $s^2 = (c_1^2, c_2^2, \dots, c_n^2) \in EW$, iff $\forall k, 1 \leq k \leq n, c_k^1 = c_k^2$, or $c_k^1 = anc_{L_i}^{L_j}(c_k^2)$ for some levels $L_i \prec L_j$.

It can be shown (proof in [15]) that:

Theorem 1 The covers relationship between context states is a partial order relationship.

Definition 10 (Covering set) A set S_i of extended context states, $S_i \subseteq EW$ covers a set S_j of extended context states, $S_j \subseteq EW$, iff $\forall s \in S_j, \exists s' \in S_i$, such that s' covers s .

Now, we define formally, which context descriptor matches the state of a query.

Definition 11 (Matching context) Let P be a profile, cod a context descriptor and C_P the set of context descriptors appearing in the contextual preferences of P . We say that a context descriptor $cod' \in C_P$ is a match for cod iff

- (i) $Context(cod')$ covers $Context(cod)$, and

- (ii) $\neg \exists cod'' \in C_P, cod'' \neq cod',$ such that $Context(cod')$ covers $Context(cod'')$ and $Context(cod'')$ covers $Context(cod)$.

There are two issues, one is whether there is at least one context preference that matches a given cod and the other one is what happens if there are more than one match. Regarding the first issue, if there is no matching context, we consider that there is no context associated with the query. In this case, the query is executed as a normal (i.e., non contextual) preference query. Note that the user can define non contextual preference queries, by using empty context descriptors which correspond to the (all, all, \dots, all) state (see Def. 4).

As an example for the case of more than one match, consider again the $cod = (location = Athens \wedge weather = warm)$ and the profile $P = \{((location = Greece \wedge weather = warm), attributes_clause, interest_score_1), ((location = Athens \wedge weather = good), attributes_clause, interest_score_2)\}$. Both context descriptors in P satisfy the first condition of Def. 11 (i.e., it holds $Context(location = Greece \wedge weather = warm)$ covers $Context(location = Athens \wedge weather = warm)$ and $(location = Athens \wedge weather = good)$ covers $Context(location = Athens \wedge weather = warm)$), but none of them covers the other. In this case, it is necessary to define which one is the most closely related state, i.e., a better match. There are many ways to handle such ties, including letting the user decide. In the next section, we propose two ways of defining similarity among context states.

4.3. State Similarity

To express similarity between two context states, we introduce a distance function named *hierarchy distance*. Using the hierarchy distance leads to choosing the preference that refers to the most specific context state, that is the state that is defined in the most detailed hierarchy level. To define the hierarchy distance, we define first the level of a state as follows.

Definition 12 (Levels of a state) Given a state $s = (c_1, c_2, \dots, c_n)$, the hierarchy levels that correspond to this state are $levels(s) = [L_{j_1}, L_{j_2}, \dots, L_{j_n}]$, such that, $c_i \in dom_{L_{j_i}}(C_i), i = 1, \dots, n$.

The distance between two levels is defined as:

Definition 13 (Level distance) Given two levels L_1 and L_2 , their distance $dist_H(L_1, L_2)$ is defined as follows:

1. if a path exists in a hierarchy between L_1 and L_2 , then $dist_H(L_1, L_2)$ is the minimum number of edges that connect L_1 and L_2 ;
2. otherwise $dist_H(L_1, L_2) = \infty$.

Having defined the distance between two levels, we can now define the level-based distance between two states.

Definition 14 (Hierarchy state distance) Given two states $s^1 = (c_1^1, c_2^1, \dots, c_n^1)$ and $s^2 = (c_1^2, c_2^2, \dots, c_n^2)$, the hierarchy distance $dist_H(s^1, s^2)$ is defined as:

$$dist_H(s^1, s^2) = \sum_{i=1}^n |dist_H(L_i^1, L_i^2)|.$$

The hierarchy state distance produces an ordering of states that is compatible with the covers partial order in the sense that between two covering states s^2 and s^3 of s^1 , the matching one is the one with the smallest hierarchy distance (proof in [15]):

Property 1 Assume a state $s^1 = (c_1^1, c_2^1, \dots, c_n^1)$. For any two different states $s^2 = (c_1^2, c_2^2, \dots, c_n^2)$ and $s^3 = (c_1^3, c_2^3, \dots, c_n^3)$, $s^2 \neq s^3$, that both cover s^1 , that is s^2 covers s^1 and s^3 covers s^1 , if s^3 covers s^2 , then $dist_H(s^3, s^1) > dist_H(s^2, s^1)$.

A second way to define the distance between two states is to use the Jaccard distance function. In this case, we compute all the descendants of each value of a state. For two values of two states corresponding to the same context parameter, we measure the fraction of the intersection of their corresponding lowest level value sets over the union of this two sets. In this case, we consider as a better match, the “smallest” state in terms of cardinality. Next, we define the Jaccard distance of two values v_1 and v_2 and use it to define the Jaccard state distance.

Definition 15 (Jaccard distance) The Jaccard distance of two values v_1 and v_2 , belonging to levels L_i and L_j of the same hierarchy H that has level L_1 as the most detailed level, is defined as:

$$dist_J(v_1, v_2) = 1 - \frac{desc_{L_1}^{L_i}(v_1) \cap desc_{L_1}^{L_j}(v_2)}{desc_{L_1}^{L_i}(v_1) \cup desc_{L_1}^{L_j}(v_2)}.$$

The Jaccard distance is consistent with the ordering produced by the level distance (proof in [15]):

Property 2 Assume three values, v_i, v_j, v_k , defined at different levels $L_i \prec L_j \prec L_k$ of the same hierarchy having L_1 as the most detailed level, such that $v_k = anc_{L_j}^{L_k}(v_j) = anc_{L_i}^{L_k}(v_i)$. Then, $dist_J(v_k, v_i) \geq dist_J(v_j, v_i)$.

Now, we define the Jaccard distance between states.

Definition 16 (Jaccard state distance) Given two states $s^1 = (c_1^1, c_2^1, \dots, c_n^1)$ and $s^2 = (c_1^2, c_2^2, \dots, c_n^2)$, the Jaccard state distance $dist_J(s^1, s^2)$ is defined as

$$dist_J(s^1, s^2) = \sum_{i=1}^n |dist_J(c_i^1, c_i^2)|.$$

Due to Properties 1 and 2, it holds that:

Property 3 Assume a state $s^1 = (c_1^1, c_2^1, \dots, c_n^1)$. For any two different states $s^2 = (c_1^2, c_2^2, \dots, c_n^2)$ and $s^3 = (c_1^3, c_2^3, \dots, c_n^3)$, $s^2 \neq s^3$, that both cover s^1 , that is s^2 covers s^1 and s^3 covers s^1 , if s^3 covers s^2 , then $dist_J(s^3, s^1) > dist_J(s^2, s^1)$.

4.4. A Context Resolution Algorithm

Given a database and a certain context descriptor (that characterizes either the current or a hypothetical context), the problem is to locate the tuples of the database that correspond to the given context descriptor and score them appropriately. The problem is further divided in two parts:

1. (Context Resolution) Locate in the profile tree the paths (context states) that correspond to the given context descriptor (in an exact or approximate fashion).
2. On the basis of the leaves of these paths (i.e., expressions of the form $A_i = value, score$), determine the corresponding tuples in the underlying database and annotate them with the appropriate score.

In the following, we detail each of these steps.

Context Resolution. Given a contextual query CQ with an extended context descriptor, for each context state $s = (c_1, c_2, \dots, c_n)$ in the context of the descriptor, we search the contextual preferences in the profile to locate a state that matches it. To this end, we use the profile tree. If there is a state that exactly matches it, that is a state (c_1, c_2, \dots, c_n) , then the associated preference is returned to the user. Note, that this state is easily located, by a single depth-first-search traversal of the profile tree. Starting from the root of the tree (level 1), at each level i , we follow the pointer associated with $key = c_i$. If such a state does not exist, we search for a state s' that matches s . If more than one such state exists, we select the one with the smallest distance, using either the hierarchy or the Jaccard distance.

Algorithm 1 presents the *Search_CS* algorithm that implements context resolution. Given a *Profile tree* whose root node is R_P , the algorithm returns all paths whose context state is either the same or covers the searching context state (c_1, c_2, \dots, c_n) . Each candidate path counts the distance from the searching path. To search an extended context state, at first we invoke *Search_CS*($R_P, \{c_1, c_2, \dots, c_n\}, 0$). At the end of the execution of this call, we can sort all results on the basis of their distances and select the one with the minimum distance, i.e., the one that differs the least from the searched path based on one of the distances. Clearly the last step can be easily replaced by a simple runtime check that keeps the current closest leaf if its distance is smaller than the one currently tested. Still, we prefer to keep this variant of the algorithm to cover the general case where more than one candidate can be selected by the system or the user.

It can be easily proved that the algorithm is correct, i.e., if applied for all extended context states specified by the extended context descriptor of the query, it leads to the desired set of states according to Def. 11 (proof in [15]).

For estimating the complexity of context resolution, we consider the number of cells accessed. In the case of an exact match, locating the related preferences requires just a simple root-to-leaf traversal of the profile tree. At level i , we search for the cell having as key the i^{th} value in the query and descend to the next level, following the pointer of the corresponding cell. For a profile tree with n context parameters (C_1, C_2, \dots, C_n) , in the worst case, we need to access $\sum_{i=1}^n |edom(C_i)|$ cells. In the case of a non-exact match, at each level i , for the i^{th} query value, we

also need to consider all its ancestors in the hierarchy. If each parameter C_i has h_i hierarchy levels, then in the worst case, we need to access $|edom(C_1)| + |edom(C_2)| \times h_1 + |edom(C_3)| \times h_2 \times h_1 + \dots + |edom(C_n)| \times h_{n-1} \times \dots \times h_1$ cells. Note, that in the case of a sequential scan, we need to access all cells in the profile, even for an exact match, since there may be more than one context state exactly matching the query (referring to a different non context-aware attribute). Thus, a sequential scan needs to visit $|edom(C_1)| \times |edom(C_2)| \times \dots \times |edom(C_n)|$ cells.

Algorithm 1 Search_CS Algorithm

Input: A node R_P of the *Profile tree*, the searching context state (c_1, c_2, \dots, c_n) , the current distance of each candidate path.

Output: A *ResultSet* of tuples of the form (Attribute name = attribute value, interest score, distance) characterizing a candidate path whose context state is either the same or best covers the searching context state.

Begin

if R_P is a non leaf node **then**

$\forall x \in R_P$ such that $(x = c_i)$ or $(x = anc_{L_i}^{L_j}(c_i))$
 $Search_CS(x \rightarrow child, \{c_{i+1}, \dots, c_n\}, dist(x, c_i) + distance)$

else if R_P is a leaf of the form $(A_i = value, score)$ **then**

$ResultSet = ResultSet \cup (A_i = value, score, distance)$

end if

End

Determination of the database tuples that correspond to the identified states. Assume a relation $R(A_1, A_2, \dots, A_n)$ and a profile tree P with leaves containing expressions of the form $(A_i = value, score)$. The problem now is that given a context descriptor cod , we need to rank the tuples of relation R with respect to cod . A simple algorithm is employed for this task.

The algorithm *Search_CS* is invoked for all extended context states specified by the query descriptor. Each invocation returns an expression that characterizes one or more tuples of the underlying relation. Then, we perform all the produced expressions as selections of the relational algebra over the underlying relation. It is straightforward (and practically orthogonal to our problem) to add (a) ranking of the expressions by their score (and consequently, ranking of the results of the queries over the relation) and (b) removal of duplicate tuples produced by these selection queries by keeping the *max* (equivalently, *avg*, *min*, or some weighted average) for the score of tuples appearing more than once in the *ResultSet*.

Extension to arbitrary queries. Algorithm 2 can be straightforwardly extended to capture context-sensitive database queries. Assume a context-independent expression E in relational algebra, and its context-dependent extension E^C

Algorithm 2 Rank_CS Algorithm

Input: A profile tree P , a relation $R(A_1, A_2, \dots, A_n)$ and a context descriptor cod

Output: A *TupleResultSet* of tuples of R ranked by the appropriate score.

Variables: An (initially empty) *ExprResultSet* of expressions of *Search_CS* results.

Begin

\forall state $s \in context(cod)$ {

Pick minimum distance tuple t from the result of $Search_CS(P, s, 0)$

$ExprResultSet = ExprResultSet \cup \{t\}$

\forall expression $e : (A_i = value, score) \in ExprResultSet$ {

$ResultSet = ResultSet \cup \sigma_{A_i=value}(R)$,
with the latter annotated with score. }

End

(that incorporates the context C of its run-time). Assume that Algorithm 1 returns a set of tuples t_1, \dots, t_n with expressions $\{\phi_1 : A_{i1} = value_1, \dots, \phi_n : A_{in} = value_n\}$. The evaluation of an expression $\sigma_{\phi_i}(E)$ returns the database tuples that correspond to preference t_i . Finally, the answer to the query E^C is $\bigcup_i \sigma_{\phi_i}(E)$. Again, the ranking of the results and the removal of duplicates is straightforward.

5. Evaluation

We evaluate our approach along two perspectives: usability and performance.

5.1. Usability Evaluation

We use a real database of points-of-interest of the two largest cities in Greece, namely Athens and Thessaloniki. To ease the specification of contextual preferences, we created a number of default profiles based on the (a) age (below 30, between 30-50, above 50), (b) sex (male or female) and (c) taste (broadly categorized as mainstream or out-of-the-beaten track). Based on the above three characteristics, users were assigned one of the 12 available profiles. Each of these profiles has 650 user preferences. Each preference consists of three context values (*accompanying_people*, *time*, *location*), an *attribute_name*, an *attribute_value* and an *interest score*. The active domains of the context parameters have 4, 17, 100 values, respectively.

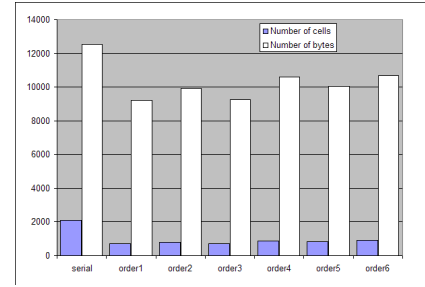
Users were allowed to modify the default profiles assigned to them by adding, deleting or updating preferences. We evaluated the system along two lines: easy of profile specification and quality of results. We run our prototype implementation for 10 users; the results are summarized in Table 1. For all users, it was the first time that they used the system.

Table 1. User Study Results

	User 1	User 2	User 3	User 4	User 5	User 6	User 7	User 8	User 9	User 10
Num of updates	22	31	12	28	24	32	38	13	18	25
Update time (mins)	30	45	20	30	30	40	45	15	20	25
Exact match	100%	90%	90%	95%	90%	100%	100%	85%	100%	100%
1 cover state	100%	95%	90%	85%	90%	100%	100%	85%	90%	100%
More cover states										
<i>Hierarchy</i>	90%	85%	80%	80%	90%	90%	90%	70%	85%	85%
<i>Jaccard</i>	95%	90%	85%	100%	95%	90%	100%	75%	85%	95%

With regards to profile specification, we count the number of modifications (insertions, deletions, updates) of preferences of the default profile that was originally assigned to the users. In addition, we report how long (in minutes), it takes users to specify/modify their profile. This also includes the time it took users to understand how profile specification works. The results are reported in the first two lines of Table 1. The general impression was that predefined profiles save time in specifying user preferences. Furthermore, having default profiles makes it easier for someone to understand the main idea behind the system, since the preferences in the profile act as examples. With regards to time, as expected, there is deviation among the time users spent on specifying profiles: some users were more meticulous than others, spending more time in adjusting the profiles assigned to them.

With regards to the quality of the results, users were asked to rank the results of each contextual query manually. Then, we compare the ranking specified by the users with what was recommended by the system, for the following three cases: (i) when there is an exact match (ii) when there is exactly one cover, and (iii) when there is more than one cover, and the *hierarchy* or the *Jaccard* distance functions are used. For each case, we consider the best 20 results, i.e., the 20 points-of-interest that were ranked higher. When there are ties in the ranking, we consider all results with the same score. We report the percentage of the results returned that belong to the results given by the user. As shown in Table 1, this percentage is generally high. Surprisingly, sometimes users do not conform even to their own preferences as shown by the results for exact match queries. In this case, although the context state of the preferences used was an exact match of the context state in the query, still some users ranked their results differently than the system. In such cases, traceability helps a lot, since users can track back which preferences were used to attain the results and either modify the preferences or reconsider their ranking. Note that users that customized their profile by making more modifications got more satisfactory results than those that spent less time during profile specification. Finally, it seems that the *Jaccard* distance produces more ac-


Figure 3. Size using real profiles.

curate results than the *hierarchy* distance mainly because the *hierarchy* distance produces rankings with many ties.

5.2. Performance Evaluation

To evaluate performance, we run a set of experiments using both real and synthetic profiles. The *real profile* is the one used for the usability study. We consider: (a) the space required to store preferences when using a profile tree as opposed to storing them sequentially and (b) the complexity of context resolution.

Size of the Profile Tree. In this set of experiments, we evaluate the size of the context tree for different mappings of the context parameters to the levels of the tree.

Using a Real Profile. We count the total number of cells and the total number of bytes of the context tree that is created for each ordering of the context parameters. Let A stand for *accompanying_people*, T for *time* and L for *location*. We call *order 1* the ordering in which A is assigned to the first level of the tree, T to the second and L to the third one, that is the ordering (A, T, L) . *Order 2* is the ordering (A, L, T) , *order 3* is (T, A, L) , *order 4* is (T, L, A) , *order 5* is (L, A, T) and *order 6* is (L, T, A) . As shown in Fig. 3, the orderings that result in trees with smaller sizes are the ones that map the context parameter with large domains lower in the tree. In addition, all trees occupy less space than storing preferences sequentially.

Using Synthetic Profiles. We study the size of the tree as a function of the size of the profile (i.e., number of user preferences). Synthetic profiles have three context parameters,

and thus, the profile tree has three levels (plus one for the leaves). There are three different types regarding the cardinality of the domains of the context parameters: a domain with 50 values, a domain with 100 values and a domain with 1000 values and profiles with various numbers (500, 1000, 5000 and 10000) of user preferences. Context values are drawn from their corresponding domain, either using a uniform data distribution, or a zipf data distribution with $a = 1.5$. The size of the tree depends on the ordering of context parameters. We call *order 1* the ordering in which the parameter whose domain has 50 values is assigned to the first level, the parameter with 100 values to the second one, and the parameter with 1000 values to the last one. *Order 2* is the ordering (50, 1000, 100), *order 3* is (100, 50, 1000), *order 4* is (100, 1000, 50), *order 5* is (1000, 50, 100) and *order 6* is (1000, 100, 50). As expected, storage is minimized when the parameters with large domains are placed lower in the tree (Fig. 4 (left, center)). For the zipf distribution (Fig. 4, center), the total number of cells is smaller than for the uniform distribution (Fig. 4, left), because “hot” values appear more frequently in preferences, i.e., more context values are the same. We also consider the case in which one of the parameter is highly skewed. In this case, it may be more space efficient to map it higher in the tree, even if its domain is large (Fig. 4, right). In this experiment, the profile has 5000 preferences, and the context parameters have domains with 50, 100 and 200 values. The values of the parameters with domains with 50 and 100 values are selected using a *uniform* data distribution and the values of the parameter with 200 values using a *zipf* data distribution with various values for the parameter a , varying from 0 (corresponding to the *uniform* distribution) to 3.5 (corresponding to a very high skew). *Order 1* is the ordering (50, 100, 200), *order 2* is (50, 200, 100) and *order 3* is (200, 50, 100).

Context Resolution. To study the usefulness of the profile tree in answering preference queries, and in particular for finding the appropriate preferences, we performed a set of experiments in which we count the number of cell accesses during context resolution. We run this experiment using both the real (Fig. 5, left) and synthetic profiles (Fig. 5, center and right). We use synthetic profiles with 500, 1000, 5000, and 10000 preferences. In all cases, the profile tree is the one in which the larger domains are mapped in lower levels. In synthetic profiles, the context values are selected from the corresponding domain, using both a uniform and a zipf data distribution with $a = 1.5$. We report results where the context parameters have hierarchies with (i) 2, 3 and 3 levels and (ii) 4, 6 and 6 levels, and for 50 randomly generated queries, where the context parameters take values from different hierarchy levels. With the profile tree, exact match queries are resolved by a simple root-to-leaf traversal, while non exact matches need to consider multi-

ple candidate paths. In the case of the sequential scan, for both the exact and the non exact matches the whole profile needs to be scanned.

6. Related Work

There is little work on context-aware preference queries. In our previous research [14, 16], we have addressed the same problem of expressing contextual preferences. However, the model used there for defining preferences includes only a *single* context parameter. Interest scores of preferences involving more than one context parameter are computed by a simple weighted sum of the preferences of single context parameters. Here, we allow contextual preferences that involve more than one context parameter and also associate context with queries. The problem of context state resolution and its development is also new here. Contextual preferences, called situated preferences, are also discussed in [9]. Situations (i.e., context states) are uniquely linked through an N:M relationship with preferences expressed using the quantitative approach. Our model is compatible with this approach and further supports multidimensional attributes and context resolution. Next, we discuss related work on context queries and on preference queries.

Context and Queries. Although, there is much research on location-aware query processing in the area of spatio-temporal databases, integrating other forms of context in query processing is less explored. In the context-aware querying processing framework of [7], there is no notion of preferences, instead context attributes are treated as normal attributes of relations. Storing context data using data cubes, called context cubes, is proposed in [8] for developing context-aware applications that use archive sensor data. The Context Relational Model ([12]) is an extended relational model that allows attributes to exist under some contexts or to have different values under different contexts. Context as a set of dimensions (e.g., context parameters) is also considered in [13] where the problem of representing context-dependent semistructured data is studied. A similar context model is deployed in [6] for enhancing web service discovery with contextual parameters. Recently, context has been used in information filtering to define context-aware filters which are filters that have attributes whose values change frequently [5]. Finally, in [11], the current contextual state of a system is represented as a multidimensional subspace within or near other situation subspaces.

Preference Queries. The research literature on preferences is extensive. In the context of database queries, there are two different approaches for expressing preferences: a quantitative and a qualitative one. With the *quantitative approach*, preferences are expressed indirectly by using scoring functions that associate a numeric score with every tuple of the query answer. In our work, we have adapted the general quantitative framework of [1]. In the quantitative framework of [10], user preferences are stored as degrees

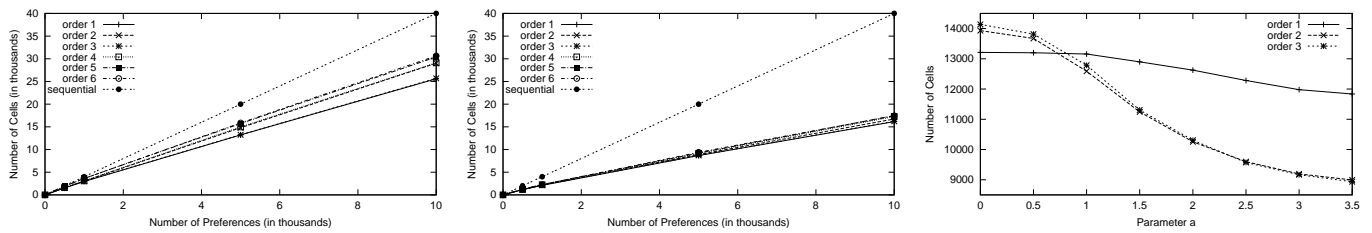


Figure 4. Uniform (left), zipf with $a=1.5$ (center) and combined (right) data distribution.

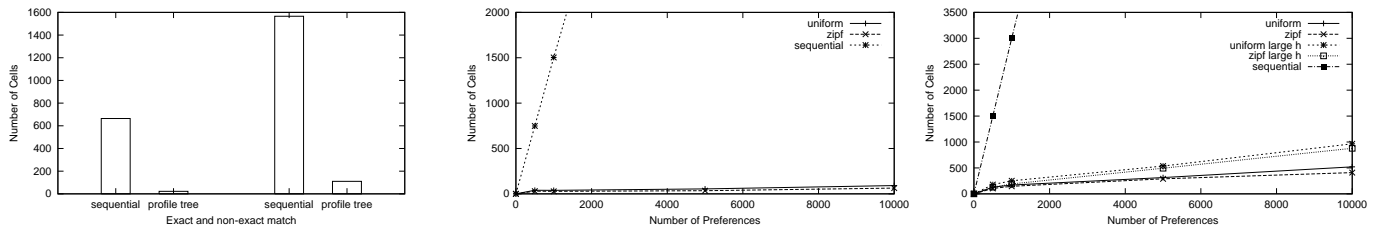


Figure 5. Number of cells accessed to find related preferences to queries, in real profiles (left) and in synthetic profiles, for an exact (center) and non exact match (right).

of interest in *atomic query elements* (such as individual selection or join conditions) instead of interests in specific attribute values. Our approach can be generalized for this framework as well, by making the degree of interest for each atomic query element depend on context. In the *qualitative approach* (i.e., [3]), the preferences between tuples in the answer to a query are specified directly, typically using binary preference relations. This framework can also be readily extended to include context.

7. Summary

In this paper, we focus on handling contextual preferences. We define context descriptors for specifying conditions on context parameters that allow the specification of context states at various levels of detail. Preferences are augmented with context descriptors that specify their scope of applicability. Similarly, queries are enhanced with context descriptors. We formulate the context resolution problem of identifying the preferences that are most relevant to the context of a query. To address this problem, we develop the notion of cover between states as well as appropriate distance functions. We also present an algorithm that locates relevant preferences. Finally, we evaluate both the usability of our model and the performance of context resolution.

References

- [1] R. Agrawal and E. L. Wimmers. A framework for expressing and combining preferences. In *ACM SIGMOD*, 2000.
- [2] G. Chen and D. Kotz. A Survey of Context-Aware Mobile Computing Research. Technical Report TR2000-381, Dartmouth College, Computer Science, November 2000.
- [3] J. Chomicki. Preference formulas in relational queries. *ACM Trans. Database Syst.*, 2003.
- [4] A. K. Dey. Understanding and Using Context. *Personal and Ubiquitous Computing*, 2001.
- [5] J.-P. Dittrich, P. M. Fischer, and D. Kossmann. Agile: Adaptive indexing for context-aware information filters. In *ACM SIGMOD*, 2005.
- [6] C. Doukeridis and M. Vazirgiannis. Querying and Updating a Context-Aware Service Directory in Mobile Environments. *Web Intelligence*, 2004.
- [7] L. Feng, P. Apers, and W. Jonker. Towards Context-Aware Data Management for Ambient Intelligence. In *DEXA*, 2004.
- [8] L. D. Harvel, L. Liu, G. D. Abowd, Y.-X. Lim, C. Scheibe, and C. Chatham. Context cube: Flexible and effective manipulation of sensed context data. In *Pervasive*, 2004.
- [9] S. Holland and W. Kiessling. Situated preferences and preference repositories for personalized database applications. In *ER*, 2004.
- [10] G. Koutrika and Y. Ioannidis. Personalization of queries in database systems. *ICDE*, 2004.
- [11] A. Padovitz, S. W. Loke, and A. Zaslavsky. Towards a theory of context spaces. *PerCom*, 2004.
- [12] Y. Roussos, Y. Stavarakas, and V. Pavlaki. Towards a Context-Aware Relational Model. *CRR*, 2005.
- [13] Y. Stavarakas and M. Gergatsoulis. Multidimensional Semistructured Data: Representing Context-Dependent Information on the Web. *CAiSE*, 2002.
- [14] K. Stefanidis, E. Pitoura, and P. Vassiliadis. On Supporting Context-Aware Preferences in Relational Database Systems. *MCMP*, 2005. (Extended version to appear in Intl Journal of Pervasive Computing and Communications, 2006).
- [15] K. Stefanidis, E. Pitoura, and P. Vassiliadis. Adding Context to Preferences (extended version). *University of Ioannina, Computer Science Department, TR 2006-07*, 2006.
- [16] K. Stefanidis, E. Pitoura, and P. Vassiliadis. Modeling and Storing Context-Aware Preferences. In *ADBIS*, 2006.
- [17] P. Vassiliadis and S. Skiadopoulos. Modelling and Optimisation Issues for Multidimensional Databases. In *CAiSE*, 2000.