

PROPERTIES OF PROBABILITY PRODUCTIVE DEPENDENCIES IN THE DATA ANALYSIS OF LARGE DATA VOLUMES

Oleksandr Pshenychnyi

Lviv Polytechnic National University

sasha.pshenychniy@gmail.com

Abstract: This paper describes findings in the area of aggregated associative dependencies detection. The work gives a method of building aggregated associative dependencies in large data volumes. This research can be applied to the wide range of data types and optimize data mining processes

Key words: aggregated associative dependency, data mining, large data volume, functional dependency.

Introduction

Nowadays computer systems operate with huge amount of data: there are plenty of terabyte storages and even a few petabyte-storages of structured data. Although hardware becomes more powerful every month, it isn't able to analyze these data storages.

This research discusses the problem of detection of associative relationships in large volumes of data, shows the results of investigation into the properties of associative dependencies and ways to optimize the search of such dependencies in large data volumes.

Data mining methods identifying dependencies and correlations are widely used in sociology, psychology, political science, physics, power industry, astronomy, computer science, and many other applied disciplines. The task of identifying the associative dependencies in polls was considered in [1]. This trend in data analysis is relatively old, but active research is still held in this area. For example, the research [2] describes a method of aggregate association rules of a construction based on simpler dependencies. This interest in the identification of data relationships can be explained by the increase in processing capacity of computer technologies and increasing volumes of collected data in many social areas. The data is collected in such amounts that its expert analysis is incomplete or may be even impossible. Modern computing hardware allows us to implement more intricate algorithms and to apply them to large amounts of data. It encourages researchers to develop such algorithms, and motivates the owners of large databases and data warehouses to develop data analysis software for the accumulated information.

Some branches of science and technology have already had powerful methods of data analysis, specifically designed to fit their needs and data

structures. The most notable ones are the following software products: CLASSIFI (Department of Pathology, UT Southwestern Medical Center) [3], BiNGO (Department of Plant Systems Biology, VIB / Ghent University) [4] and EASE (National Institute of Allergy and Infectious Diseases) [5]. However, most research institutions can't afford to develop such systems and require a common method, applicable to the wide range of data types and structures.

Thus, the effective detection of associative dependencies in multi-attribute data is an actual problem of modern data analysis.

It should be noted that data mining is a very broad area of data analysis, and the detection of associative relationships is only a part of it.

The goal of this research is to investigate the properties of aggregated associative dependencies that will help to implement effective data dependency identification algorithms, applicable to various data types and structures.

Relations and the work area of existing methods of data analysis

Data mining technologies are mostly used to identify data dependencies of the types "if ... then ..." or "for ... is true ...". Such dependencies are represented by implications, production rules or associative rules.

Data mining includes the wide range of mathematical and algorithmic tools, such as neural networks, evolutionary algorithms, decision trees etc. However, modern research is more and more accented on logical dependencies of the search in datasets. With their help the problems of classification, forecasting, creation of formal description of real objects and others are solved. [2]

It is possible to find an incredibly big amount of associative dependencies in large databases. It is impossible even to store all those dependencies in any data storage in the world, let alone analyzing them. Fortunately, it is not the goal of any business branch. Instead, there is a need for finding associative dependencies, which meet some predefined quality requirements (they are statistically substantiated).

The main problems of currently available data dependency detection methods are as follows:

- they can work only with binary properties of objects;
- they don't find dependencies with low support;
- they aren't able to effectively process addition of new data to data source;
- they are non-efficient in analyzing multi-attribute data;
- they are not flexible enough to satisfy the needs of dependency quality filtering.

One of the ways to solve some of these problems is building aggregated associative rules. The research [2] offers to use a system of 4 parameters of an associative dependency, which describe the rule. In general, such complex of criteria is more flexible and useful than a single association intensity attribute, proposed in [1], but sufficiency of this system has not yet been proven.

Taken into consideration the aforesaid, the given research discusses the urgent scientific and technical problems of the detection of aggregated associative dependencies and the development of valuation methods for the analysis of large databases.

Problem Statement

The search of arbitrary associative rules $P(x) \rightarrow Q(x), x \in r(R)$ in ratio $r(R)$ is a very broad task, the solution to which is still to be found and this problem is not the object of our research. The paper proposes studying the properties of associative dependencies, in which conditional and resulting predicates look like:

$$P = P_1^e \vee P_2^e \vee \dots \vee P_h^e = \bigvee_{k=1}^h P_k^e, \quad (1)$$

$$P_k^e = A_{i_1} \in \left\{ a_{(i_1)(j_1)} \right\} \wedge A_{i_2} \in \left\{ a_{(i_2)(j_2)} \right\} \wedge \dots \wedge A_{i_k} \in \left\{ a_{(i_k)(j_k)} \right\},$$

$$\forall l = \overline{1..k} : A_{i_l} \in R, \forall m = \overline{1..k} : \left\{ a_{(i_l)(j_m)} \right\} : a_{(i_l)(j_m)} \in \text{dom}(A_{i_l})$$

$$\forall i, j \in \{1..h\} : \arg(P_i^e) = \arg(P_j^e)$$

$$F_I : \bigvee_{k=1}^s P_k^e \rightarrow \bigvee_{l=1}^t Q_l^e \quad (2)$$

The denotation $\arg(P)$ is used for the operator returning a set of attribute-arguments of a predicate P .

We assign to such associative dependencies a separate term – probabilistic productive dependency (PPD). Various sources use different markings, interpretations and explanations of this concept. For example, in [2] a term "intensity of association" is used, in [1] it is "degree of confidence". However, these notions concern to a greater extent the existing data, and not the data to be worked with currently. Moreover, considering the system as static, not being renewed with new knowledge, we lose the sense of using such techniques as Laplacian smoothing [6] and other methods of protection against noise and data uncertainty.

Just for emphasising on the study of random processes, their dynamics and problems, the term "probability" has been added to the term of "dependency", which is being investigated. The second part of the term, namely "productive", does not need a special explanation, since the productive rule underlies a dependence.

Thus, the probabilistic productive dependence is a productional rule of a type Eq.(2) in the selection of the basic relation that holds for a significant number of objects of the selection. The threshold of significance should be determined in an expert way or based on calculations of the probability of a false selection of this relationship.

Let's write the Eqs.(1), (2) in terms of relational algebra:

$$P^e(x) = \pi_{A_{i_1} A_{i_2} \dots A_{i_k}}(x) \in \left\{ a_{(i_1)(j_1)} \right\} \times \left\{ a_{(i_2)(j_2)} \right\} \times \dots \times \left\{ a_{(i_k)(j_k)} \right\},$$

$$x \in r(R), \forall l = \overline{1..k} : A_{i_l} \in R, \forall m = \overline{1..k} : \left\{ a_{(i_l)(j_m)} \right\} : \quad (3)$$

$$a_{(i_l)(j_m)} \in \text{dom}(A_{i_l})$$

$$P(x) = \pi_{A_{i_1} A_{i_2} \dots A_{i_h}}(x) \in \left\{ a_{(i_1)(j_{1,1})} \right\} \times \left\{ a_{(i_2)(j_{1,2})} \right\} \times \dots \times \left\{ a_{(i_h)(j_{1,h})} \right\} \cup$$

$$\cup \left\{ a_{(i_1)(j_{2,1})} \right\} \times \left\{ a_{(i_2)(j_{2,2})} \right\} \times \dots \times \left\{ a_{(i_h)(j_{2,h})} \right\} \cup \dots \cup$$

$$\cup \left\{ a_{(i_1)(j_{h,1})} \right\} \times \left\{ a_{(i_2)(j_{h,2})} \right\} \times \dots \times \left\{ a_{(i_h)(j_{h,h})} \right\} \quad (4)$$

In other words, PPD predicates can be presented as predicates defined on processions of a ratio $r(R)$, not only on a set of attributes.

The PPD threshold of significance can be defined on the basis of an arbitrary function of assessing the importance of a detected dependence. However, the most often used indicators are those of a level of support and a level of confidence. In the source [2] it was shown that these parameters are insufficient to describe adequately the dependencies of a subject area and it was proposed to use the additional indicators: a level of improvement and a full mutual information.

Let us consider these notions in detail:

Level of support is a characteristic of a selection predicate in a relation, calculated as the ratio of the number of objects that satisfy the predicate P to the total number of objects in the relation:

$$Sup(P) = \frac{|\sigma_P(r)|}{|r|} \quad (5)$$

When calculating the level of support for PPD, conditional and resulting predicate are combined by a sign of conjunction:

$$Sup(S \rightarrow T) = Sup(S \wedge T) = \frac{|\sigma_{S \wedge T}(r)|}{|r|} \quad (6)$$

Level of confidence – the ratio of the number of objects, for which there is such a PPD to the number of objects in the selection:

$$Conf(S \rightarrow T) = P(S \rightarrow T) = \frac{|\sigma_{S \wedge T}(r)|}{|\sigma_S(r)|} \quad (7)$$

Using the concept of support level, the confidence level can be calculated as:

$$Conf(S \rightarrow T) = \frac{Sup(S \rightarrow T)}{Sup(S)} \quad (8)$$

Level of improvement is calculated as the ratio of the confidence levels and the support levels of PPD:

$$Imp(S \rightarrow T) = \frac{Conf(S \rightarrow T)}{Sup(T)} = \frac{Sup(S \wedge T)}{Sup(S) \cdot Sup(T)} \quad (9)$$

Full mutual information is calculated in general as:

$$I_{X \leftrightarrow Y} = \sum_{i=1}^n \sum_{j=1}^m p_{ij} \log_2 \frac{p_{ij}}{p_i r_j}, \quad (10)$$

where $p_{ij} = P(X = x_i \wedge Y = y_j)$ – the probability of the fact that X is in the condition x_i , and Y – in the condition y_j ; $p_i = P(X = x_i)$ – the probability of the fact that X is in the condition x_i ; $r_j = P(Y = y_j)$ – the probability of the fact that Y is in the condition y_j .

For associative rules the mutual information can be defined as:

$$I_{X \leftrightarrow Y} = \sum_{i=1}^n \sum_{j=1}^m Sup(x_i \rightarrow y_j) \log_2 Imp(x_i \rightarrow y_j) \quad (11)$$

Analogs of the rules of the functional dependencies derivation for PPD

As in the case with F-dependencies, the PPD set occurring in a given ratio can be represented by their certain subset, from which all the relevant PPD from the given ratio can be received by means of derivation rules. Since PPD is an extension of F-dependencies, we should consider the transformation of axioms of the functional dependencies F1-F6 derivation for PPD.

PPDs are characterized by many parameters, the most important of which are given in Eqs.(5)–(11). However, a parameter which is used the most often and is the easiest to understand is that of the confidence level. More sophisticated characteristics of dependencies are based on it, what has been found during the investigations.

Filtering by the level of support does not allow to derive PPDs from small partial dependencies – the level of support certainly increases when combining PPDs,

and therefore it is impossible to make a cut-off of some groups of dependencies based on this parameter.

The level of improvement is a nonlinear characteristic. That makes it impossible to perform a clipping of PPD generation based on the set of existing dependencies.

The full mutual information is also nonlinear by the power of the selection of both PPD predicates.

Thus, there is no sense in building derivation rules for the parameters of the improvement level and the full mutual information, because they are calculated from the level of support and the level of confidence; in addition, the given characteristics nonlinearly depend on the number of processions, which corresponds to the number of dependencies.

A strict proof of this fact is not given, because it is intuitively understood from the above considerations and a formal proof is rather cumbersome.

We resolve the problem of associative dependencies filtering flexibility using two-step filtering. It allows using arbitrary criteria of PPD quality at the second stage.

So, let us consider the transformation of the functional dependencies derivation for PPD:

Reflexivity of the confidence level.

$$Conf(s \in S \rightarrow s \in S) = 1 \text{ for any ratio } r(R).$$

Proof:

$$Conf(s \in S \rightarrow s \in S) = \frac{|\sigma_{s \in S \wedge s \in S}|}{|\sigma_{s \in S}|} = \frac{|\sigma_{s \in S}|}{|\sigma_{s \in S}|} = 1$$

Replenishment of the confidence level. If

$$Conf(s \in S \rightarrow t \in T) = p, \quad \text{then}$$

$$Conf(s \in S \wedge w \in D(W) \rightarrow t \in T) = p, \quad \text{where}$$

$D(W)$ – the domain of the attribute W of a ratio $r(R)$.

Proof:

$$\begin{aligned} Conf(s \in S \wedge w \in D(W) \rightarrow t \in T) &= \frac{|\sigma_{s \in S \wedge w \in D(W) \wedge t \in T}(R)|}{|\sigma_{s \in S \wedge w \in D(W)}(R)|} = \\ &= \frac{|\forall x \in r : q = \pi_{W=w}(x) \in D(W) \Rightarrow w \in D(W)|}{|\sigma_{s \in S \wedge t \in T}(R)|} = \\ &= \frac{|\sigma_{s \in S \wedge t \in T}(R)|}{|\sigma_{s \in S}(R)|} = Conf(s \in S \rightarrow t \in T) = p \end{aligned}$$

$$Conf(a \in \{a_1, a_2\} \rightarrow d \in \{d_1\}) = \frac{4}{8} = 0,5$$

$$\begin{aligned} Conf(a \in \{a_1, a_2\} \wedge b \in \{b_1, b_2, b_3, b_4\} \rightarrow d \in \{d_1\}) &= \\ &= \frac{|\sigma_{a \in \{a_1, a_2\} \wedge b \in \{b_1, b_2, b_3, b_4\} \wedge d \in \{d_1\}}(R)|}{|\sigma_{a \in \{a_1, a_2\} \wedge b \in \{b_1, b_2, b_3, b_4\}}(R)|} = \frac{4}{8} = 0,5 \end{aligned}$$

Table 1

Example 1

| A | B | C | D |
|----------------|----------------|----------------|----------------|
| a ₁ | b ₁ | c ₁ | d ₁ |
| a ₁ | b ₁ | c ₂ | d ₁ |
| a ₁ | b ₂ | c ₂ | d ₁ |
| a ₂ | b ₂ | c ₂ | d ₁ |
| a ₂ | b ₂ | c ₂ | d ₂ |
| a ₂ | b ₂ | c ₂ | d ₂ |
| a ₂ | b ₃ | c ₃ | d ₂ |
| a ₂ | b ₂ | c ₂ | d ₂ |
| a ₃ | b ₃ | c ₃ | d ₂ |

Additivity of the confidence level. If $Conf(s \in S \rightarrow t \in T) = p$ and

$Conf(s \in S \rightarrow w \in W) = 1$, then

$$Conf(s \in S \rightarrow t \in T \wedge w \in W) = p.$$

Proof:

$$\begin{aligned} Conf(s \in S \rightarrow t \in T \wedge w \in W) &= \frac{|\sigma_{s \in S \wedge t \in T \wedge w \in W}|}{|\sigma_{s \in S}|} = \\ &= |s \in S \rightarrow w \in W| = \frac{|\sigma_{s \in S \wedge t \in T}|}{|\sigma_{s \in S}|} = Conf(s \in S \rightarrow t \in T) = p \end{aligned}$$

Example 2:

Considering PPD from the example 1:

$$Conf(A \in \{a_1\} \rightarrow B \in \{b_1\}) = \frac{2}{3}$$

$$Conf(A \in \{a_1\} \rightarrow D \in \{d_1\}) = 1$$

Thereof, it can be concluded that

$$Conf(A \in \{a_1\} \rightarrow B \in \{b_1\} \wedge D \in \{d_1\}) = \frac{2}{3}. \quad \text{It is}$$

proved by the calculations with the Eq.(7):

$$Conf(A \in \{a_1\} \rightarrow B \in \{b_1\} \wedge D \in \{d_1\}) = \frac{|\sigma_{A \in \{a_1\} \wedge B \in \{b_1\} \wedge D \in \{d_1\}}|}{|\sigma_{A \in \{a_1\}}|} = \frac{2}{3}$$

Projectivity of the confidence level. If

$$Conf(s \in S \rightarrow t \in T \wedge w \in W) = p$$

and

$$Conf(s \in S \rightarrow w \in W) = 1,$$

then $Conf(s \in S \rightarrow t \in T) = p$.

Proof:

$$\begin{aligned} Conf(s \in S \rightarrow t \in T) &= \frac{|\sigma_{s \in S \wedge t \in T}|}{|\sigma_{s \in S}|} = \\ &= \frac{|\sigma_{s \in S \wedge t \in T \wedge w \in W}|}{|\sigma_{s \in S}|} |s \in S \rightarrow w \in W| = \\ &= Conf(s \in S \rightarrow t \in T \wedge w \in W) = p \end{aligned}$$

Example 3:

Let's consider the example, introduced in the previous paragraph, in a reverse variant:

From PPD

$$Conf(A \in \{a_1\} \rightarrow B \in \{b_1\} \wedge D \in \{d_1\}) = \frac{2}{3} \quad \text{and}$$

$Conf(A \in \{a_1\} \rightarrow D \in \{d_1\}) = 1$ it can be concluded

that $Conf(A \in \{a_1\} \rightarrow B \in \{b_1\}) = \frac{2}{3}$. The verification

can be done with the help of the Eq.(7):

$$Conf(A \in \{a_1\} \rightarrow B \in \{b_1\}) = \frac{|\sigma_{A \in \{a_1\} \wedge B \in \{b_1\}}|}{|\sigma_{A \in \{a_1\}}|} = \frac{2}{3}$$

Transitivity of the confidence level. If

$$Conf(s \in S \rightarrow t \in T) = p$$

and

$$Conf(t \in T \rightarrow w \in W) = 1,$$

then

$$Conf(s \in S \rightarrow w \in W) \geq p.$$

Proof:

$$Conf(s \in S \rightarrow w \in W) = \frac{|\sigma_{s \in S \wedge w \in W}|}{|\sigma_{s \in S}|}$$

$$Conf(s \in S \rightarrow t \in T) = \frac{|\sigma_{s \in S \wedge t \in T}|}{|\sigma_{s \in S}|} = p$$

Thus, since $|\sigma_{s \in S}| \geq 0$, $|\sigma_{s \in S \wedge w \in W}| \geq 0$ i

$|\sigma_{s \in S \wedge t \in T}| \geq 0$ (follows from the definition of a relational

selection operation), then in order to prove inequality

$Conf(s \in S \rightarrow w \in W) \geq Conf(s \in S \rightarrow t \in T)$ we

need to prove that $|\sigma_{s \in S \wedge w \in W}| \geq |\sigma_{s \in S \wedge t \in T}|$.

Let's consider a variable-procession x of a ratio $r(R)$, provided that $\pi_s(x) \in S$ and $\pi_t(x) \in T$. According to

the condition $Conf(t \in T \rightarrow w \in W) = 1$, if

$\pi_T(x) \in T$, then $\pi_w(x) \in W$. Thus

$$Conf(s \in S \rightarrow t \in T) = p \wedge Conf(t \in T \rightarrow w \in W) = 1:$$

$$\forall x \in r(R): \pi_s(x) \in S \wedge \pi_t(x) \in T \Rightarrow \pi_w(x) \in W$$

Hence we have got a number of consequences:

$$\sigma_{s \in S \wedge w \in W} \subseteq \sigma_{s \in S \wedge t \in T}$$

$$|\sigma_{s \in S \wedge w \in W}| \geq |\sigma_{s \in S \wedge t \in T}|$$

$$Conf(s \in S \rightarrow w \in W) \geq Conf(s \in S \rightarrow t \in T) = p$$

$$Conf(s \in S \rightarrow w \in W) \geq p$$

Hereby, the transitivity of the PPD confidence level has been proved.

Example 4:

Using the data of the example 2 the following PPD can be created:

$$\text{Conf}(C \in \{c_2\} \rightarrow B \in \{b_1\}) = \frac{1}{6}$$

$$\text{Conf}(B \in \{b_1\} \rightarrow D \in \{d_1\}) = 1$$

From the rule of the transitivity of the PPD confidence level we conclude that $\text{Conf}(C \in \{c_2\} \rightarrow D \in \{d_1\}) \geq \frac{1}{6}$.

Let us verify this by calculating $\text{Conf}(C \in \{c_2\} \rightarrow D \in \{d_1\})$ using the Eq.(7):

$$\text{Conf}(C \in \{c_2\} \rightarrow D \in \{d_1\}) = \frac{|\sigma_{C \in \{c_2\} \wedge D \in \{d_1\}}|}{|\sigma_{C \in \{c_2\}}|} = \frac{3}{6} = \frac{1}{2}$$

$$\begin{aligned} \text{Conf}(C \in \{c_2\} \rightarrow D \in \{d_1\}) &= \\ &= \frac{1}{2} \geq \frac{1}{6} = \text{Conf}(C \in \{c_2\} \rightarrow B \in \{b_1\}) \end{aligned}$$

Thus, the transitivity of the PPD confidence level has been confirmed.

The transitivity of the PPD confidence level is a powerful rule for making various assumptions and proofs.

Pseudotransitivity of the confidence level. This axiom of an F-dependencies derivation does not have a direct alternative for PPD, under the condition that only one restriction is imposed on the source dependencies. We are going to prove this statement.

Let us consider the dependencies

$$\text{Conf}(s \in S \rightarrow t \in T) = p \quad \text{and}$$

$\text{Conf}(t \in T \wedge q \in Q \rightarrow w \in W) = 1$, implying a limitation of genuineness on

$$\text{Conf}(t \in T \wedge q \in Q \rightarrow w \in W).$$

We denote $X = \sigma_{s \in S}(R)$, $Y = \sigma_{t \in T}(R)$, $Z = \sigma_{q \in Q}(R)$, $V = \sigma_{w \in W}(R)$. Then

$$\text{Conf}(t \in T \wedge q \in Q \rightarrow w \in W) = \frac{|Y \cap Z \cap V|}{|Y \cap Z|} = 1$$

$$|X \cap Z \cap V| = |X \cap Z|$$

$$\text{Conf}(s \in S \rightarrow t \in T) = \frac{|X \cap Y|}{|X|} = p$$

$$\begin{aligned} \text{Conf}(s \in S \wedge q \in Q \rightarrow w \in W) &= \frac{|X \cap Z \cap V|}{|X \cap Z|} = \\ &= \frac{|((X \cap Y) \cup (X \setminus Y)) \cap Z \cap V|}{|((X \cap Y) \cup (X \setminus Y)) \cap Z|} = \frac{|(X \cap Y \cap Z \cap V) \cup ((X \setminus Y) \cap Z \cap V)|}{|(X \cap Y \cap Z) \cup ((X \setminus Y) \cap Z)|} = \\ &= \frac{|(X \cap Y \cap Z) \cup ((X \setminus Y) \cap Z \cap V)|}{|(X \cap Y \cap Z) \cup ((X \setminus Y) \cap Z)|} \end{aligned}$$

$$\begin{aligned} &= \frac{|X \cap Y \cap Z| + |(X \setminus Y) \cap Z \cap V| - |X \cap Y \cap Z \cap (X \setminus Y) \cap Z \cap V|}{|X \cap Y \cap Z| + |(X \setminus Y) \cap Z|} = \\ &= \frac{|X \cap Y \cap Z| + |(X \setminus Y) \cap Z \cap V|}{|X \cap Y \cap Z| + |(X \setminus Y) \cap Z|} = (1) \end{aligned}$$

$$\begin{aligned} X \cap Y \cap Z &\subset X \cap Z \\ (X \setminus Y) \cap Z \cap V &\subset X \cap Z \\ (X \setminus Y) \cap Z &\subset X \cap Z \end{aligned}$$

Thereby, $X \cap Z$ is a universal set U of given expressions and the result of calculations Eq.(1) will not change, if we consider only the processions from $X \cap Z$.

We set $Y' = Y \cap (X \cap Z)$, $V' = V \cap (X \cap Z)$. Then

$$(1) = \frac{|Y'| + |\neg Y' \cap V'|}{|Y'| + |\neg Y'|} = \frac{|Y'| + |\neg Y' \cap V'|}{|U|}$$

$$|\neg Y' \cap V'| \in [0; |\neg Y'|]$$

$$\frac{|Y'| + |\neg Y' \cap V'|}{|U|} \in \left[\frac{|Y'|}{|U|}; 1 \right]$$

Let us return to the introduced settings: $Y' = Y \cap (X \cap Z)$. Initial conditions do not impose restrictions on the value of a given expression, so $|Y'| \in [0; |U|]$ and accordingly

$$\frac{|Y'| + |\neg Y' \cap V'|}{|U|} \in [0; 1]$$

$$\text{Conf}(s \in S \wedge q \in Q \rightarrow w \in W) \in [0; 1]$$

Thus, the dependencies $\text{Conf}(s \in S \rightarrow t \in T) = p$ i $\text{Conf}(t \in T \wedge q \in Q \rightarrow w \in W) = 1$ do not have any influence on the dependence

$$\text{Conf}(s \in S \wedge q \in Q \rightarrow w \in W).$$

Now let us consider the restrictions of the other dependence, combining $\text{Conf}(s \in S \rightarrow t \in T) = 1$ and $\text{Conf}(t \in T \wedge q \in Q \rightarrow w \in W) = p$. Using the settings given above, we obtain:

$$\text{Conf}(t \in T \wedge q \in Q \rightarrow w \in W) = \frac{|Y \cap Z \cap V|}{|Y \cap Z|} = p$$

$$\text{Conf}(s \in S \rightarrow t \in T) = \frac{|X \cap Y|}{|X|} = 1$$

$$|X \cap Y| = |X|$$

$$\text{Conf}(s \in S \wedge q \in Q \rightarrow w \in W) = \frac{|X \cap Z \cap V|}{|X \cap Z|}$$

In this case, the initial dependencies do not limit the resulting expression, and it is apparent just at the first step: a set V may not have mutual processions with X , immediately turning

$\text{Conf}(s \in S \wedge q \in Q \rightarrow w \in W)$ into zero. On the other hand, the variant $V \subset X \cap Z$ is also possible – then $\text{Conf}(s \in S \wedge q \in Q \rightarrow w \in W) = 1$. The function $\text{Conf}(s \in S \wedge q \in Q \rightarrow w \in W)$ linearly depends on $|X \cap Z \cap V|$, having a variable set V , and thus having a range of values $[0;1]$. So, the dependencies $\text{Conf}(s \in S \rightarrow t \in T) = 1$ and $\text{Conf}(t \in T \wedge q \in Q \rightarrow w \in W) = p$ do not limit the range of values $\text{Conf}(t \in T \wedge q \in Q \rightarrow w \in W)$.

Thus, it has been proved that the pseudotransitivity of functional dependencies does not have a direct analog among PPD.

As we can see from the above arguments, most axioms of the F-dependencies derivation may be transformed for PPD only with the significant limitations of the conditional part of one of the dependencies. In addition, some analogs do not give a precise formula for calculating the confidence level of a new dependence, but only limit it. Thus, the given set of PPD derivation rules is not complete. To ensure the completeness of derivation rules, it is necessary to introduce additional derivation rules, specific for PPD.

Operations on PPD

Factorization

We shall name the expansion of the dependence F_I $\pi_{A_1 A_2 \dots A_k}(s) \in \{s_1, s_2, \dots, s_m\} \rightarrow \pi_{A_1 A_2 \dots A_l}(s) \in \{t_1, t_2, \dots, t_n\}$ into the set of dependencies

$$\left\{ \pi_{A_1 A_2 \dots A_k}(s) = s_i \rightarrow \pi_{A_1 A_2 \dots A_l}(s) = t_j \right\}, i = \overline{1..m}, j = \overline{1..n}$$

a factorization and denote it as $F_I [Fact]$.

$$F_I = \sum_{i=1}^m \sum_{j=1}^n \left(\pi_{A_1 A_2 \dots A_k}(s) = s_i \rightarrow \pi_{A_1 A_2 \dots A_l}(s) = t_j \right) \quad (12)$$

Merging

Merging of PPD $s \in S_1 \rightarrow t \in T_1$ i $s \in S_2 \rightarrow t \in T_2$ is a new PPD $s \in S_1 \cup S_2 \rightarrow t \in T_1 \cup T_2$ (PPD denotations in terms of relational algebra are used – Eq. (4)).

$$\begin{aligned} & (s \in S_1 \rightarrow t \in T_1) + (s \in S_2 \rightarrow t \in T_2) = \\ & = s \in S_1 \cup S_2 \rightarrow t \in T_1 \cup T_2 \end{aligned} \quad (13)$$

Let us consider qualities of the PPD merging operation:

Commutativity

The operation of PPD merging has the quality of commutativity.

$$\begin{aligned} & (s \in S_1 \rightarrow t \in T_1) + (s \in S_2 \rightarrow t \in T_2) = \\ & = (s \in S_2 \rightarrow t \in T_2) + (s \in S_1 \rightarrow t \in T_1) \end{aligned} \quad (14)$$

Proof:

$$\begin{aligned} & (s \in S_1 \rightarrow t \in T_1) + (s \in S_2 \rightarrow t \in T_2) = \\ & = (s \in S_1 \cup S_2 \rightarrow t \in T_1 \cup T_2) = \\ & = (s \in S_2 \cup S_1 \rightarrow t \in T_2 \cup T_1) = \\ & = (s \in S_2 \rightarrow t \in T_2) + (s \in S_1 \rightarrow t \in T_1) \end{aligned}$$

Associativity

The operation of PPD merging has the quality of associativity.

$$\begin{aligned} & (s \in S_1 \rightarrow t \in T_1) + ((s \in S_2 \rightarrow t \in T_2) + (s \in S_3 \rightarrow t \in T_3)) = \\ & ((s \in S_1 \rightarrow t \in T_1) + (s \in S_2 \rightarrow t \in T_2)) + (s \in S_3 \rightarrow t \in T_3) \end{aligned} \quad (15)$$

Proof:

$$\begin{aligned} & (s \in S_1 \rightarrow t \in T_1) + ((s \in S_2 \rightarrow t \in T_2) + (s \in S_3 \rightarrow t \in T_3)) = \\ & = (s \in S_1 \rightarrow t \in T_1) + (s \in S_2 \cup S_3 \rightarrow t \in T_2 \cup T_3) = \\ & = s \in S_1 \cup S_2 \cup S_3 \rightarrow t \in T_1 \cup T_2 \cup T_3 \end{aligned}$$

We expand the right part of the associativity expression:

$$\begin{aligned} & ((s \in S_1 \rightarrow t \in T_1) + (s \in S_2 \rightarrow t \in T_2)) + (s \in S_3 \rightarrow t \in T_3) = \\ & = (s \in S_1 \cup S_2 \rightarrow t \in T_1 \cup T_2) + (s \in S_3 \rightarrow t \in T_3) = \\ & = s \in S_1 \cup S_2 \cup S_3 \rightarrow t \in T_1 \cup T_2 \cup T_3 = \\ & = (s \in S_1 \rightarrow t \in T_1) + ((s \in S_2 \rightarrow t \in T_2) + (s \in S_3 \rightarrow t \in T_3)) \end{aligned}$$

Rules of PPD Derivation

As it has been shown, the implementation of the transformed rules of the functional dependencies derivation is not enough to ensure the complete set of PPD derivation rules. We should consider the derivation rules, specific for PPDs, which will allow us to develop efficient algorithms for the search of these dependencies in databases.

Aggregation of a definition area. If there are PPD $s \in S_1 \rightarrow t \in T_1$ and $s \in S_2 \rightarrow t \in T_2$ and values $\sigma_{s=s_i}(r(R))$, $\sigma_{t=t_j}(r(R))$, $\sigma_{s=s_i \wedge t=t_j}(r(R))$ provided that $\bigcup_i S_i = S_1 \cup S_2$, $\bigcup_j T_j = T_1 \cup T_2$ then for the PPD

- $s \in S_1 \cup S_2 \rightarrow t \in T_1 \cup T_2$
- $s \in S_1 \cup S_2 \rightarrow t \in T_1 \cap T_2$

- c) $s \in S_1 \cap S_2 \rightarrow t \in T_1 \cup T_2$
- d) $s \in S_1 \cap S_2 \rightarrow t \in T_1 \cap T_2$
- e) $s \in S_1 \rightarrow t \in T_1 \cup T_2$
- f) $s \in S_1 \rightarrow t \in T_1 \cap T_2$
- g) $s \in S_2 \rightarrow t \in T_1 \cup T_2$
- h) $s \in S_2 \rightarrow t \in T_1 \cap T_2$
- i) $s \in S_1 \cup S_2 \rightarrow t \in T_1$
- j) $s \in S_1 \cap S_2 \rightarrow t \in T_1$
- k) $s \in S_1 \cup S_2 \rightarrow t \in T_2$
- l) $s \in S_1 \cap S_2 \rightarrow t \in T_2$

all the parameters from formulas (5)-(11) can be calculated.
Proof:

$$\sigma_{x \in X}(r) = \sum_{z \in X} \sigma_{x=z}(r) \quad (16)$$

Let us demonstrate the proving of the most complex (the first) of the given consequences. Others can be provided similarly and are not presented here for brevity sake.

$$\begin{aligned} & Sup(s \in S_1 \cup S_2 \rightarrow t \in T_1 \cup T_2) = \\ & = Sup(s \in S_1 \cup S_2 \wedge t \in T_1 \cup T_2) = \frac{|\sigma_{s \in S_1 \cup S_2 \wedge t \in T_1 \cup T_2}(r)|}{|r|} = \\ & = \frac{|\sigma_{s \in S_1 \wedge t \in T_1 \cup T_2}(r)| + |\sigma_{s \in S_2 \wedge t \in T_1 \cup T_2}(r)| - |\sigma_{s \in S_1 \cap S_2 \wedge t \in T_1 \cup T_2}(r)|}{|r|} \quad (17) \end{aligned}$$

$$\begin{aligned} & = \frac{1}{|r|} (|\sigma_{s \in S_1 \wedge t \in T_1}(r)| + |\sigma_{s \in S_1 \wedge t \in T_2}(r)| - |\sigma_{s \in S_1 \wedge t \in T_1 \cap T_2}(r)| + \\ & + |\sigma_{s \in S_2 \wedge t \in T_1}(r)| + |\sigma_{s \in S_2 \wedge t \in T_2}(r)| - |\sigma_{s \in S_2 \wedge t \in T_1 \cap T_2}(r)| - \\ & - |\sigma_{s \in S_1 \cap S_2 \wedge t \in T_1}(r)| - |\sigma_{s \in S_1 \cap S_2 \wedge t \in T_2}(r)| + |\sigma_{s \in S_1 \cap S_2 \wedge t \in T_1 \cap T_2}(r)|) \end{aligned}$$

$$\begin{aligned} & Conf(s \in S_1 \cup S_2 \rightarrow t \in T_1 \cup T_2) = \frac{|\sigma_{s \in S_1 \cup S_2 \wedge t \in T_1 \cup T_2}(r)|}{|\sigma_{s \in S_1 \cup S_2}(r)|} = \\ & = \frac{|r| \cdot Sup(s \in S_1 \cup S_2 \rightarrow t \in T_1 \cup T_2)}{|\sigma_{s \in S_1}(r)| + |\sigma_{s \in S_2}(r)| - |\sigma_{s \in S_1 \cap S_2}(r)|} \quad (18) \end{aligned}$$

$$\begin{aligned} & Imp(s \in S_1 \cup S_2 \rightarrow t \in T_1 \cup T_2) = \\ & = \frac{Conf(s \in S_1 \cup S_2 \rightarrow t \in T_1 \cup T_2)}{Sup(t \in T_1 \cup T_2)} = \\ & = \frac{|r| \cdot Conf(s \in S_1 \cup S_2 \rightarrow t \in T_1 \cup T_2)}{|\sigma_{t \in T_1}(r)| + |\sigma_{t \in T_2}(r)| - |\sigma_{t \in T_1 \cap T_2}(r)|} \quad (19) \end{aligned}$$

$$\begin{aligned} & I_{s \in S_1 \cup S_2 \rightarrow t \in T_1 \cup T_2} = \\ & = \sum_{x \in S_1 \cup S_2} \sum_{y \in T_1 \cup T_2} Sup(s = x \rightarrow t = y) \log_2 Imp(s = x \rightarrow t = y) \quad (20) \end{aligned}$$

The complexity of calculating the Eqs.(16)–(20) directly and linearly depends on the power of the sets S_1 and S_2 . In addition, there are data structures (Fibonacci heap, binomial heap [7]), allowing us to calculate the merging and the intersection of sets in a sublinear time.

The implementation of PPD derivation algorithms for the rule of aggregation can use subtypes and partial cases of the aggregation rule. For example, if $S_1 \cap S_2 = \emptyset$, Eqs.(16)-(20) become very simple and are calculated with the asymptotic complexity $O(1)$.

These rules of derivation (the given rule of PPD aggregation includes 12 subordinate rules, presented as consequences) are particularly effective for small sets of values, such as data from sociological and psychological interviews, weather observations, traffic studies and others.

PPD Derivation Rules Completeness

We shall consider an arbitrary PPD F_l :

$\bigvee_{k=1}^h P_k^e \rightarrow \bigvee_{l=1}^g Q_l^e$. It can be received by the aggregation of

the PPD $P_h^e \rightarrow Q_g^e$ $\exists \bigvee_{k=1}^{h-1} P_k^e \rightarrow \bigvee_{l=1}^{g-1} Q_l^e$ if $h > 1$ and

$g > 1$, $\exists P_h^e \rightarrow \bigvee_{l=1}^{g-1} Q_l^e$, if $h = 1$ and $g > 1$, \exists

$\bigvee_{k=1}^{h-1} P_k^e \rightarrow Q_g^e$, if $h > 1$ and $g = 1$.

Thus, we receive the factorization

$$\bigvee_{k=1}^h P_k^e \rightarrow \bigvee_{l=1}^g Q_l^e = \bigoplus_{i=0 \dots \max(h-1, g-1)} (P_{\max(h-i, 1)}^e \rightarrow Q_{\max(g-i, 1)}^e) \quad (21)$$

From Eq.(1)

$$P_k^e = A_{i_1} \in \{a_{(i_1)(j_1)}\} \wedge A_{i_2} \in \{a_{(i_2)(j_2)}\} \wedge A_{i_k} \in \{a_{(i_k)(j_k)}\}$$

So, the parameters of an arbitrary PPD can be calculated with the help of statistics

$$\sigma_{x=x_i}(r(R)), \sigma_{x=x_i \wedge y=y_j}(r(R)),$$

$\sigma_{x=x_i \wedge y=y_j \wedge z=z_k}(r(R))$ and so on. Using these means for presenting all PPD relations, which don't contain more than k parts of predicates conditions, we need

$$O\left(\max_{(i_1, i_2, \dots, i_k) \in Z^k} (|class(A_{i_1})| \cdot |class(A_{i_2})| \cdot \dots \cdot |class(A_{i_k})|)\right)$$

of memory, where $Z \subset R$ is a set of attributes, among which dependencies are searched, $|class(A_{i_j})|$ is a number of classification areas with the attribute A_{i_j} . The simplest

variant is: $class(A_{i_j}) = dom(A_{i_j})$, but for numerical and measurement data it is often convenient to divide them into

subordinate areas. This division increases the informational content and simplifies their search.

Keeping the full statistics with arbitrary deepness of enclosure is usually impossible due to limitations of available memory in a computing system, but in practice dependencies with more than 3-4 parts of the conditional predicate are not used. Accordingly, the presentation of all required data is quite possible even for quite large data sets.

Conclusions

This paper shows the results of research of building and merging aggregated associative dependencies. The investigated class of associative dependencies is widely used in data mining methods. The application is found in the branches of computer sciences, power engineering, physics, sociology etc.

As a result of the research the rules of expanding the full set of aggregated PPD from their certain set have been obtained. This fact allows storing only the minimal coverage of a data set with a selected class of dependencies, but not all the available dependencies in the context. This form of data presentation allows us to easily modify them (removing or adding processions, and also changing the attribute values of existing processions). The detection of this PPD property gives an important advantage over many other methods of static data analysis: data being changed, it is not necessary to recalculate all the data statistics, but only to update the required parameters.

The article provides the proof of effective calculation of such characteristics as a level of support, a level of confidence, a level of improvement and a full mutual information of PPD. However, there are more parameters which can be efficiently computed with the help of the studied properties and the developed rules of PPD derivation. During further research it is planned to examine the necessary conditions for the criteria of quality, so that they might be effectively calculated by applying the rules of PPD derivation.

The application of the PPD derivation rules described in this paper can reduce the required memory capacity of a computer system to

$$O \left(\max_{(i_1, i_2, \dots, i_k) \in Z^k} \left(\left| class(A_{i_1}) \right| \cdot \left| class(A_{i_2}) \right| \cdot \dots \cdot \left| class(A_{i_k}) \right| \right) \right),$$

where k is the maximum number of attributes that appears in the conditional and the resulting part of desired PPDs. Usually there is no need for values $k > 3$, but it may occur for some specific values of fixed attributes, and in this case it becomes possible to keep separate statistics for these attributes.

Thus, the derivation rules described in this paper allow us to store and find efficiently the PPDs in large data sets, based on its minimal coverage.

References:

1. Chesnokov S. Determinational analysis of socio-economic data.– M.: Nauka. – 1982.– 168p. [Rus]
2. Titova O. Methods of construction and evaluation of aggregate associative rules in intelligent databases. – Kharkiv – 2006.
3. Main site of the Pathology Department, UT Southwestern Medical Center . [Electronic resource] <http://pathcuric1.swmed.edu/pathdb/classifi.html>
4. Description of a utility BiNGO, site of the Ghent University, [Electronic resource] <http://www.psb.ugent.be/cbd/papers/BiNGO/Home.html>.
5. Official site of the National Institute of Allergy and Infectious Diseases (NIAID), NIH, [Electronic resource] <http://david.abcc.ncifcrf.gov/content.jsp?file=/ease/ease1.htm&type=1>
6. Norwig P., Thun S. Materials of online lecture “Machine Learning” of Stanford University – 2011. [Electronic resource] <https://www.ai-class.com/course/video/quizquestion/97>.
7. Llc Books, Heaps: Heapsort, Binary Heap, Smoothsort, Soft Heap, Fibonacci Heap, Treap, Binomial Heap, Pairing Heap, Leftist Tree, Skew Heap. Memphis, Tennessee, General Books LLC. – 2010. – 74p.

ВЛАСТИВОСТІ АСОЦІАТИВНИХ ЗАЛЕЖНОСТЕЙ У АНАЛІЗІ ДАНИХ

О. Пшеничний

У статті наведено результати дослідження властивостей асоціативних залежностей та можливостей їх ефективного агрегування. Розроблено метод виявлення асоціативних залежностей широкого класу у великих наборах даних.



Oleksandr Pshenychnyi – M.Sc., graduated in 2009 from Lviv Polytechnic National University in speciality “Intelligent Decision Support Systems”. Currently is a post-graduate of Lviv Polytechnic National University, works as a software developer for Eleks Software Ltd.

Scientific interests area: data analysis methods, optimization and effective computation, parallel programming, grid- and cloud-based calculations, distributed data processing.