

# Flow Graphs - a New Paradigm for Intelligent Data Analysis

Zdzisław Pawlak

Institute of Theoretical and Applied Informatics

Polish Academy of Sciences, ul. Bałtycka 5, 44 100 Gliwice, Poland

and

Warsaw School of Information Technology, ul. Newelska 6, 01-447 Warsaw, Poland

Address for correspondence: ul. Zuga 29, 01 806 Warsaw, Poland

(48 22) 8345659, fax. (48 22) 8251635, e-mail: zpw@ii.pw.edu.pl

## Abstract

In this paper we propose a new approach to data (mining) and knowledge discovery based on information flow distribution study in a flow graph. Flow graphs introduced in this paper are different from those proposed by Ford and Fulkerson for optimal flow analysis and they model rather, e.g., flow *distribution* in a network, than the optimal flow. The flow graphs considered in this paper are not meant to physical media (e.g., water) flow analysis, but to information flow examination in decision algorithms. It is revealed that flow in the flow graph is governed by Bayes' rule, but the rule has entirely deterministic interpretation, not referring to its probabilistic roots. Besides, decision algorithm induced by the flow graph and dependency between conditions and decisions of decision rules are defined and studied. This idea is based on statistical concept of dependency but in our setting it has deterministic meaning.

*Keywords:* Flow graph; Data mining; Knowledge discovery.

## 1 Introduction

Searching for patterns in databases is of utmost importance in data mining in recent years [3]. Many methods have been developed and used in this domain, where statistical methods in particular Bayesian approach, play a

substantial role. However statistical method despite many advantages cause often problems due to probabilistic interpretation of obtained results.

Let us also observe that despite Bayes' rule fundamental role in statistical inference it has led to many philosophical discussions concerning its validity and meaning, and has caused much criticism [2], [3].

This drawback has deep roots related to understanding of probability, which will be discussed briefly next.

The concept of probability can be traced back to Laplace [8] who gave the definition of probability which is in use until now. But his idea of probability still causes many discussion and critics concerning its correctness and validity. One of the first who proposed the way out of the dilemma how free probability from its obscure meaning was Jan ukasiewicz [9], by suggesting to replace the concept of probability by truth values of propositional functions.

Similar ideas have been proposed independently many years later by Adams [1], Carnap [4], Ramsey [12] and Reichenbach [13].

In this paper we propose still another approach to solve this problem. Instead of using truth values in place of probability, stipulated by ukasiewicz, we propose, using of deterministic flow analysis in flow networks (graphs). However we analyze in the flow graph not the absolute value of the flow in each branch of the network but its relative value to the total flow expressed by a

fraction between 0 and 1. In this setting, flow is described by formulas which have formally probabilistic flavor (e.g., Bayes' rule), or by the corresponding logical calculus proposed by Lukasiewicz, though, the formulas have entirely deterministic meaning, and need neither probabilistic nor logical interpretation.

Flow graphs introduced in this paper are different from those proposed by Ford and Fulkerson [5] for optimal flow analysis and they model flow distribution in a network, than the optimal flow. More specifically they are used to information flow examination in decision algorithms. To this end branches of a flow graph can be interpreted as decision rules. With every decision rule (i.e. branch) three coefficients are associated, the *strength*, *certainty* and *coverage factors*.

These coefficients have been used under different names in data mining (see e.g., [4], [15]) but they were used first by ukasiewicz [9] in his study of logic and probability.

We start our consideration by defining basic concepts of the proposed approach, i.e., flow graph and its fundamental properties. Next decision algorithm induced by the flow graph and dependency between conditions and decisions of decision rules are defined and studied. This idea is based on statistical concept of dependency but in our setting it has deterministic meaning.

Simple tutorial examples is used to illustrate how the introduced ideas work in data mining. The presented ideas can be used, as a new tool for data mining, and knowledge discovery. Besides, it also throw a new light on the concept of probability.

This paper is a continuation of some authors' ideas presented in [11], where the relationship between Bayes' rule and flow graphs has been introduced and studied and is modified version of the plenary paper presented at KSS2004 [10].

## **2 Flow graphs**

### **2.1 Overview**

In this part the fundamental concepts of the proposed approach are defined and discussed. In particular flow graphs, certainty and coverage factors of branches of the flow graph are defined and studied. Next these coefficient are extended to paths and some classes of sub-graphs, called connections. Further a notion of a fusion of a flow graph is defined.

Further dependences of flow are introduced and examined. Finally dependency factor (correlation coefficient) is defined.

## 2.2 Basic concepts

A flow graph is a *directed, acyclic, finite* graph  $G = (N, \mathcal{B}, \varphi)$ , where  $N$  is a set of *nodes*,  $\mathcal{B} \subseteq N \times N$  is a set of *directed branches*,  $\varphi : \mathcal{B} \rightarrow R^+$  is a *flow function* and  $R^+$  is the set of non-negative reals. *Input* of a node  $x \in N$  is the set  $I(x) = \{y \in N : (y, x) \in \mathcal{B}\}$ ; *output* of a node  $x \in N$  is defined as  $O(x) = \{y \in N : (x, y) \in \mathcal{B}\}$ . We will also need the concept of *input* and *output* of a graph  $G$ , defined, respectively, as follows:  $I(G) = \{x \in N : I(x) = \emptyset\}$ ,  $O(G) = \{x \in N : O(x) = \emptyset\}$ . Inputs and outputs of  $G$  are *external nodes* of  $G$ ; other nodes are *internal nodes* of  $G$ .

If  $(x, y) \in \mathcal{B}$  then  $\varphi(x, y)$  is a *throughflow* from  $x$  to  $y$ . With every node  $x$  of a flow graph  $G$  we associate its *inflow*

$$\varphi_+(x) = \sum_{y \in I(x)} \varphi(y, x), \quad (1)$$

and *outflow*

$$\varphi_-(x) = \sum_{y \in O(x)} \varphi(x, y), \quad (2)$$

Similarly, we define an inflow and an outflow for the whole flow graph, which are defined as

$$\varphi_+(G) = \sum_{x \in I(G)} \varphi_-(x), \quad (3)$$

$$\varphi_-(G) = \sum_{x \in O(G)} \varphi_+(x). \quad (4)$$

We assume that for any internal node  $x$ ,  $\varphi_+(x) = \varphi_-(x)$ , where  $\varphi(x)$  is a *throughflow* of node  $x$ .

Obviously,  $\varphi_+(G) = \varphi_-(G) = \varphi(G)$ , where  $\varphi(G)$  is a *throughflow* of graph  $G$ .

The above formulas can be considered as *flow conservation equations* [5].

**Example.** We will illustrate basic concepts of flow graphs by an example of a group of 1000 patients put to the test for certain drug effectiveness.

Assume that patients are grouped according to presence of the disease, age and test results, as shown in Fig. 1.

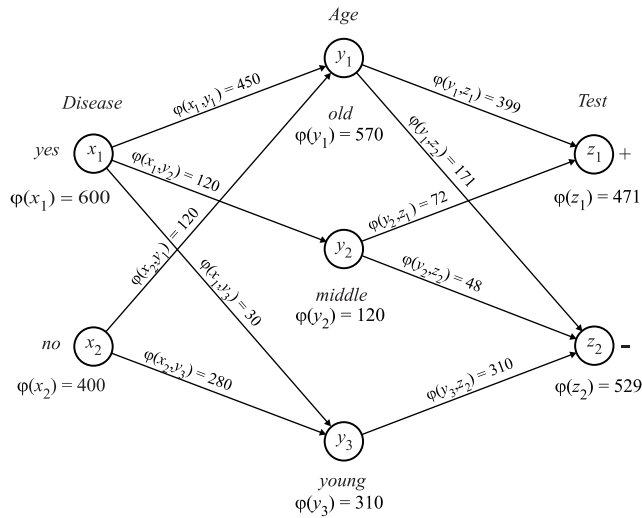


Figure 1: Flow graph

E.g.,  $\varphi(x_1) = 600$  means that these are 600 patients suffering from the

disease,  $\varphi(y_1) = 570$  means that these are 570 old patients  $\varphi(z_1) = 471$  means that 471 patients have positive test result;  $\varphi(x_1, y_1) = 450$  means that these are 450 old patients which suffer from disease etc.

Thus the flow graph gives clear insight into the relationship between different groups of patients.

Let us now explain the flow graph in more details.

Nodes of the flow graph are depicted by circles, labeled by  $x_1, x_2, y_1, y_2, y_3, z_1, z_2$ . A branch  $(x, y)$  is denoted by an arrow from node  $x$  to  $y$ . E.g., branch  $(x_1, z_1)$  is represented by an arrow from  $x_1$  to  $z_1$ , inputs of node  $y_1$  are nodes  $x_1$  and  $x_2$ , outputs of node  $x_1$  are nodes  $y_1, y_2$  and  $y_3$ .

Inputs of the flow graph are nodes  $x_1$  and  $x_2$ , whereas outputs of the flow graph are nodes  $z_1$  and  $z_2$ .

Nodes  $y_1, y_2$  and  $y_3$  are internal nodes of the flow graph. The throughflow of the branch  $(x_1, y_1)$  is  $\varphi(x_1, y_1) = 450$ . Inflow of node  $y_1$  is  $\varphi_+(y_1) = 450 + 120 = 750$ . Outflow of node  $y_1$  is  $\varphi_-(y_1) = 399 + 171 = 570$ .

Inflow of the flow graph is  $\varphi(x_1) + \varphi(x_2) = 600 + 400 = 1000$ , and outflow of the flow graph is  $\varphi(z_1) + \varphi(z_2) = 471 + 529 = 1000$ .

Throughflow of node  $y_1 = \varphi(y_1) = \varphi(x_1, y_1) + \varphi(x_2, y_1) = \varphi(y_1, z_1) + \varphi(y_1, z_2) = 570$ . □

We will define now a *normalized flow graph*. A normalized flow graph is



a *directed, acyclic, finite* graph  $G = (N, \mathcal{B}, \sigma)$ , where  $N$  is a set of *nodes*,  $\mathcal{B} \subseteq N \times N$  is a set of *directed branches* and  $\sigma : \mathcal{B} \rightarrow \langle 0, 1 \rangle$  is a *normalized flow* of  $(x, y)$  and

$$\sigma(x, y) = \varphi(x, y) \setminus \varphi(G) \quad (5)$$

is a *strength* of  $(x, y)$ . Obviously,  $0 \leq \sigma(x, y) \leq 1$ . The strength of the branch expresses simply the percentage of a total flow through the branch.

In what follows we will use normalized flow graphs only, therefore by flow graphs we will understand normalized flow graphs, unless stated otherwise.

With every node  $x$  of a flow graph  $G$  we associate its *inflow* and *outflow* defined as

$$\sigma_+(x) = \varphi_+(x) \setminus \varphi(G) = \sum_{y \in I(x)} \sigma(x, y), \quad (6)$$

$$\sigma_-(x) = \varphi_-(x) \setminus \varphi(G) = \sum_{y \in O(x)} \sigma(x, y). \quad (7)$$

Obviously for any internal node  $x$ , we have  $\sigma_+(x) = \sigma_-(x) = \sigma(x)$ , where  $\sigma(x)$  is a *normalized throughflow* of  $x$ .

Moreover, let

$$\sigma_+(G) = \varphi_+(G) \setminus \varphi(G) = \sum_{x \in I} \sigma_+(x), \quad (8)$$

$$\sigma_-(G) = \varphi_-(G) \setminus \varphi(G) = \sum_{x \in O(G)} \sigma_-(x). \quad (9)$$

Obviously,  $\sigma_+(G) = \sigma_-(G) = \sigma(G) = 1$  .

**Example (cont.).** The normalized flow graph of the flow graph presented in Fig. 1 is given in Fig. 2.

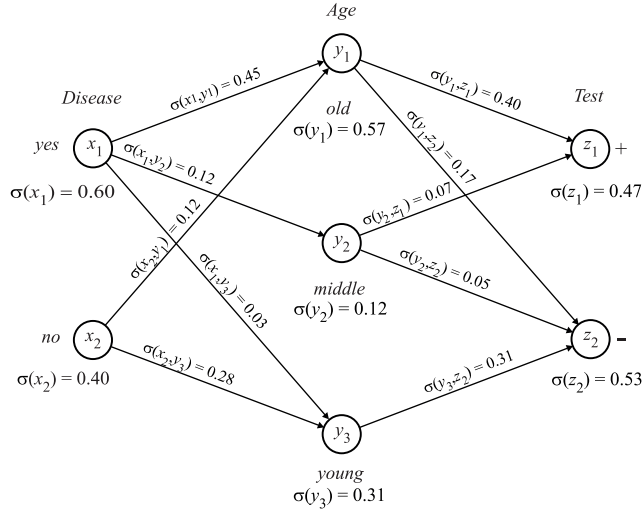


Figure 2: Normalized flow graph

In the flow graph e.g.,  $\sigma(x_1) = 0.60$ , that means that 60% of total inflow is associated with input  $x_1$ . The strength  $\sigma(x_1, y_1) = 45$  means that 45% of total flow flows through the branch  $(x_1, y_1)$  etc.  $\square$

Let  $G = (N, \mathcal{B}, \sigma)$  be a flow graph. If we invert direction of all branches in  $G$ , then the resulting graph  $G = (N, \mathcal{B}', \sigma')$  will be called an *inverted* graph of  $G$ . Of course the inverted graph  $G'$  is also a flow graph and all inputs and outputs of  $G$  become inputs and outputs of  $G'$  , respectively.

The inverted flow graph can be used to give reasons (explanation) for de-

cisions.

### 2.3 Certainty and Coverage Factors

With every branch  $(x, y)$  of a flow graph  $G$  we associate the *certainty* and the *coverage factors*.

The *certainty* and the *coverage* of  $(x, y)$  are defined as

$$cer(x, y) = \sigma(x, y) \setminus \sigma(x), \quad (10)$$

and

$$cov(x, y) = \sigma(x, y) \setminus \sigma(y). \quad (11)$$

respectively.

Evidently,  $cer(x, y) = cov(y, x)$ , where  $(x, y) \in \mathcal{B}$  and  $(y, x) \in \mathcal{B}$ .

**Example (cont.).** The certainty and the coverage factors for the flow graph presented in Fig. 2 are shown in Fig. 3.

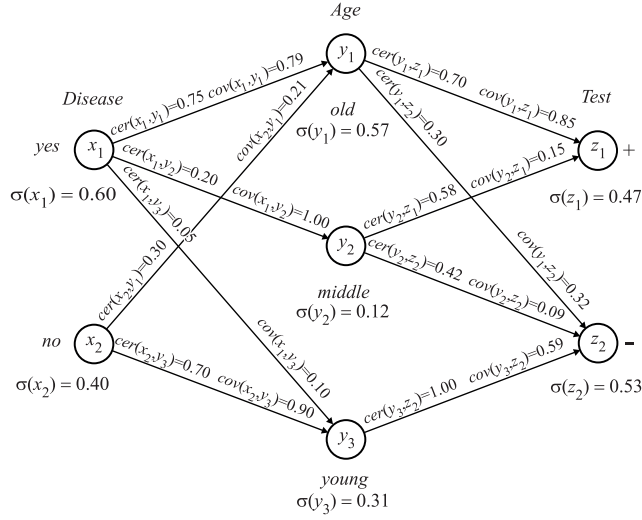


Figure 3: Certainty and coverage

E.g.,  $cer(x_1, y_1) = \sigma(x_1, y_1) \setminus \sigma(x_1) = 0.45 \setminus 0.60 = 0.75$ , and  $cov(x_1, y_1) = \sigma(x_1, y_1) \setminus \sigma(y_1) = 0.45 \setminus 0.57 \approx 0.21$ .  $\square$

**Example (cont.).** The inverted flow graph of the flow graph from Fig. 3 is shown in Fig. 4.

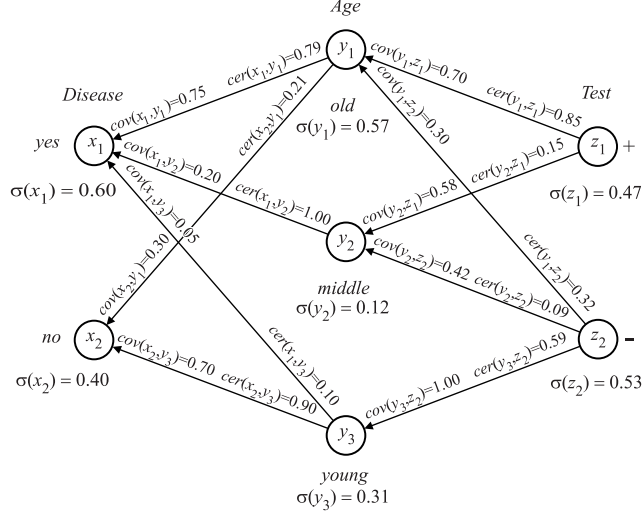


Figure 4: Inverted flow graph

□

Below some properties of certainty and coverage factors, which are immediate consequences of definitions given above, are presented:

$$\sum_{y \in O(x)} cer(x, y) = 1, \quad (12)$$

$$\sum_{x \in I(y)} cov(x, y) = 1, \quad (13)$$

$$\sigma(x) = \sum_{y \in O(x)} cer(x, y)\sigma(y) = \sum_{y \in O(x)} \sigma(x, y), \quad (14)$$

$$\sigma(y) = \sum_{x \in I(y)} cov(x, y)\sigma(x) = \sum_{x \in I(y)} \sigma(x, y), \quad (15)$$

$$cer(x, y) = cov(x, y)\sigma(y) \setminus \sigma(x), \quad (16)$$

$$cov(x, y) = cer(x, y)\sigma(x) \setminus \sigma(y). \quad (17)$$

Obviously the above properties have a probabilistic flavor, e.g., equations (14) and (15) have a form of total probability theorem, whereas formulas (16) and (17) are Bayes' rules. However, these properties in our approach are interpreted in a deterministic way and they describe flow distribution among branches in the network.

## 2.4 Paths, connections and fusion

A (*directed*) path from  $x$  to  $y$ ,  $x \neq y$  in  $G$  is a sequence of nodes  $x_1, \dots, x_n$  such that  $x_1 = x$ ,  $x_n = y$  and  $(x_i, x_{i+1}) \in \mathcal{B}$  for every  $i$ ,  $1 \leq i \leq n-1$ . A path from  $x$  to  $y$  is denoted by  $[x \dots y]$ .

The *certainty* of the path  $[x_1 \dots x_n]$  is defined as

$$cer[x_1 \dots x_n] = \prod_{i=1}^{n-1} cer(x_i, x_{i+1}), \quad (18)$$

the *coverage* of the path  $[x_1 \dots x_n]$  is

$$cov[x_1 \dots x_n] = \prod_{i=1}^{n-1} cov(x_i, x_{i+1}), \quad (19)$$

and the *strength* of the path  $[x \dots y]$  is

$$\sigma[x \dots y] = \sigma(x)cer[x \dots y] = \sigma(y)cov[x \dots y]. \quad (20)$$

The set of all paths from  $x$  to  $y$  ( $x \neq y$ ) in  $G$  denoted  $\langle x, y \rangle$ , will be called a

*connection* from  $x$  to  $y$  in  $G$ . In other words, connection  $\langle x, y \rangle$  is a sub-graph of  $G$  determined by nodes  $x$  and  $y$ .

The *certainty* of the connection  $\langle x, y \rangle$  is

$$cer\langle x, y \rangle = \sum_{[x\dots y] \in \langle x, y \rangle} cer[x\dots y], \quad (21)$$

the *coverage* of the connection  $\langle x, y \rangle$  is

$$cov\langle x, y \rangle = \sum_{[x\dots y] \in \langle x, y \rangle} cov[x\dots y], \quad (22)$$

and the *strength* of the connection  $\langle x, y \rangle$  is

$$\sigma\langle x, y \rangle = \sum_{[x\dots y] \in \langle x, y \rangle} \sigma[x\dots y] = \sigma(x)cer\langle x, y \rangle = \sigma(y)cov\langle x, y \rangle. \quad (23)$$

Let  $x, y (x \neq y)$  be nodes of  $G$ . If we substitute simultaneously every for the sub-graph  $\langle x, y \rangle$  of a given flow graph  $G$ , where  $x$  and  $y$  are input and output nodes of  $G$  respectively, by single branch  $(x, y)$  such that  $\sigma(x, y)$ , then in the resulting graph  $G'$ , called the *fusion* of  $G$ , we have  $cer(x, y) = cer\langle x, y \rangle$ ,  $cov(x, y) = cov\langle x, y \rangle$  and  $\sigma(G) = \sigma(G')$ .

**Example (cont.).** In the flow graph presented in Fig. 3 for the path  $p = [x_1, y_1, z_1]$  we have  $cer(p) = 0.75 \times 0.70 \approx 0.53$ ,  $cov(p) = 0.85 \times 0.79 \approx 0.67$ .

For example the connection  $\langle x_1, z_1 \rangle$  in the flow graph consists of paths  $[x_1, y_1, z_1]$  and  $[x_1, y_2, z_1]$ . This connection is shown in Fig. 5.

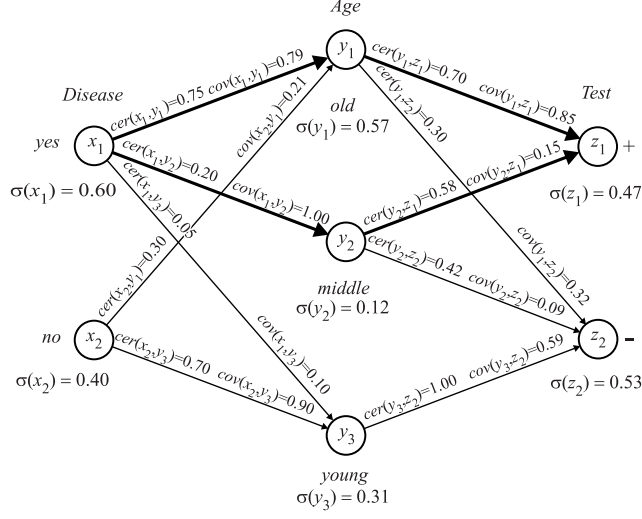


Figure 5: Connection

For this connection we have  $cer\langle x_1, z_1 \rangle = 0.75 \times 0.70 + 0.20 \times 0.58 \approx 0.64$ ;  
 $cov\langle x_1, z_1 \rangle = 0.85 \times 0.79 + 0.15 \times 1.00 \approx 0.82$ .

The strength of the connection  $x_1, z_1$  is  $0.64 \times 0.60 \approx 0.82 \times 0.47 \approx 0.38$ .

□

**Example (cont.).** The fusion of the flow graph shown in Fig. 3 is given in Fig. 6.

The fusion of a flow graph gives information about the flow distribution between input and output of the flow graph, i.e., it leads to the following conclusions:

- if the disease is present then the test result is positive with certainty 0.64,



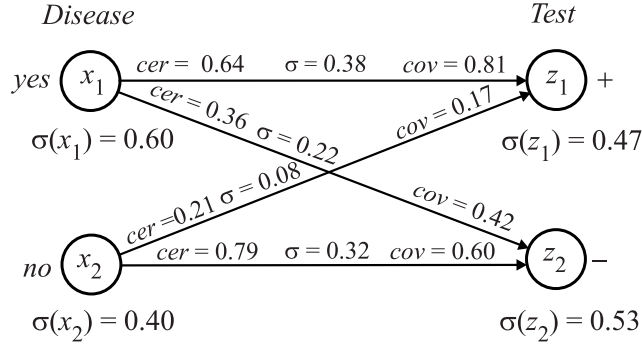


Figure 6: Fusion of a flow graph

- if the disease is absent then the test result is negative with certainty 0.79.

Explanation of test results is as follows:

- if the test result is positive then the disease is present with certainty 0.81,
- if the test result is negative then the disease is absent with certainty 0.60.

## 2.5 Dependences in flow graphs

Let  $x$  and  $y$  be nodes in a flow graph  $G = (N, \mathcal{B}, \sigma)$ , such that  $(x, y) \in \mathcal{B}$ .

Nodes  $x$  and  $y$  are *independent* in  $G$  if

$$\sigma(x, y) = \sigma(x)\sigma(y). \quad (24)$$

From (21) we get

$$\sigma(x, y)/\sigma(x) = cer(x, y) = \sigma(y), \quad (25)$$

and

$$\sigma(x, y)/\sigma(y) = cov(x, y) = \sigma(x) \quad (26)$$

If

$$cer(x, y) > \sigma(y), \quad (27)$$

or

$$cov(x, y) > \sigma(x), \quad (28)$$

then  $x$  and  $y$  are *positively dependent* on  $x$  in  $G$ . Similarly, if

$$cer(x, y) < \sigma(y), \quad (29)$$

or

$$cov(x, y) < \sigma(x), \quad (30)$$

then  $x$  and  $y$  are *negatively dependent* in  $G$ .

Let us observe that relations of independency and dependences are symmetric ones, and are analogous to those used in statistics.

For every branch  $(x, y) \in \mathcal{B}$  we define a *dependency(correlation) factor*  $\eta(x, y)$  defined as

$$\begin{aligned} \eta(x, y) = \\ & cer(x, y) - \sigma(y) / cer(x, y) + \sigma(y) = \\ & cov(x, y) - \sigma(x) / cov(x, y) + \sigma(x). \end{aligned} \tag{31}$$

Obviously  $-1 \leq \eta(x, y) \leq 1$ ;  $\eta(x, y) = 0$  if and only if  $cer(x, y)\sigma(y)$  and  $cov(x, y) = \sigma(x)$ ;  $\eta(x, y) = 1$  if and only if  $cer(x, y) = cov(x, y) = 0$ ;  $\eta(x, y) = -1$  if and only if  $\sigma(y) = \sigma(x) = 0$ . It is easy to check that if  $\eta(x, y) = 0$ , then  $x$  and  $y$  are independent, if  $-1 \leq \eta(x, y) < 0$  then  $x$  and  $y$  are negatively dependent and if  $0 < \eta(x, y) \leq 1$  then  $x$  and  $y$  are positively dependent. Thus the dependency factor expresses a degree of dependency, and can be seen as a counterpart of correlation coefficient used in statistics.

**Example (cont.).** Dependency factors for the flow graph shown in Fig. 6 are given Fig. 7.

Thus, there is positive dependency between presence of the disease and positive test result as well as between absence of disease and negative test result. However there is much stronger negative dependency between presence of the disease and negative test result or similarly – between absence of the

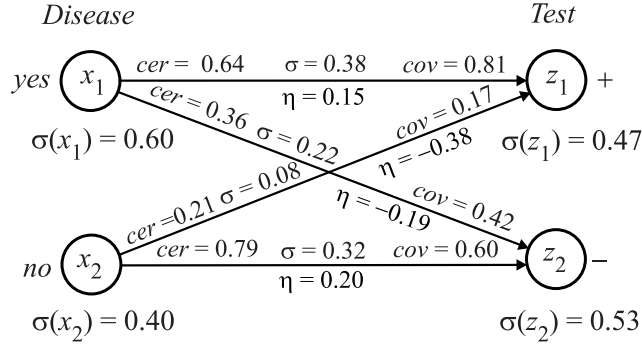


Figure 7: Dependencies in a flow graph

disease and positive test result. □

## 2.6 Flow graph and decision algorithms

Flow graphs can be interpreted as decision algorithms [6], [7], [11].

Let us assume that the set of nodes of a flow graph is interpreted as a set of logical formulas. The formulas are understood as propositional functions and if  $x$  is a formula, then  $\sigma(x)$  is to be interpreted as a truth value of the formula. Let us observe that the truth values are numbers from the closed interval  $\langle 0, 1 \rangle$ , i.e.,  $0 \leq \sigma(x) \leq 1$ , and can be also interpreted as probabilities.

With every branch  $(x, y)$  we associate a decision rule  $x \rightarrow y$ , read *if  $x$  then  $y$* ;  $x$  will be referred to as *condition*, whereas  $y$  – *decision* of the rule. Such a rule is characterized by three numbers,  $\sigma(x, y)$ ,  $cer(x, y)$  and  $cov(x, y)$ .

Thus every path  $[x_1 \dots x_n]$  determines a sequence of decision  $x_1 \rightarrow x_2$ ,

$x_2 \rightarrow x_3, \dots, x_{n-1} \rightarrow x_n$ .

From previous considerations it follows that this sequence of decision rules can be interpreted as a single decision rule  $x_1x_2 \dots x_{n-1} \rightarrow x_n$ , in short  $x^* \rightarrow x_n$ , where  $x^* = x_1x_2 \dots x_{n-1}$  characterized by

$$cer(x^*, x_n) = \sigma(x^*, x_n) / \sigma(x^*), \quad (32)$$

$$cov(x^*, x_n) = \sigma(x^*, x_n) / \sigma(x_n), \quad (33)$$

and

$$\sigma(x^*, x_n) = \sigma(x_1)cer[x_1 \dots x_n] = \sigma(x_n)cov[x_1 \dots x_n], \quad (34)$$

The set of all decision rules  $x_{i_1}x_{i_2} \dots x_{i_{n-1}} \rightarrow x_{i_n}$  associated with all paths  $[x_{i_1}, x_{i_n}]$  such that  $x_{i_1}$  and  $x_{i_n}$  are input and output of the graph respectively will be called a *decision algorithm* induced by the flow graph.

If  $x \rightarrow y$  is a decision rule then we say that condition and decision of the decision rule are independent if  $x$  and  $y$  are independent, otherwise condition and decision of the decision rule are dependent (positively or negatively).

To measure the degree of dependency between condition and decision of the decision rule  $x \rightarrow y$  we can use the dependency factor  $\eta(x, y)$ .

Thus every decision rule beside strength, certainty and coverage factor can be also characterized by the degree of dependency between its condition and

decision. This measure can be used as a new tool for data mining in pursuit of patterns in data.

**Example (cont.).** The decision algorithm induced by the flow graph shown in Fig. 3 is given below.

	certainty	coverage	strength
$x_1, y_1 \rightarrow z_1$	0.71	0.67	0.32
$x_1, y_1 \rightarrow z_2$	0.29	0.25	0.14
$x_1, y_2 \rightarrow z_1$	0.58	0.15	0.07
$x_1, y_2 \rightarrow z_2$	0.42	0.09	0.05
$x_1, y_3 \rightarrow z_2$	1.00	0.06	0.03
$x_2, y_1 \rightarrow z_1$	0.67	0.18	0.08
$x_2, y_1 \rightarrow z_2$	0.33	0.08	0.04
$x_2, y_3 \rightarrow z_2$	1.00	0.53	0.28

It can be easily seen that the decision rules  $x_1, y_3 \rightarrow z_2$  and  $x_1, y_3 \rightarrow z_2$  can be replaced by a single decisions rule  $y_3 \rightarrow z_2$ . Consequently the decision algorithm can be presented as

	certainty	coverage	strength
$x_1, y_1 \rightarrow z_1$	0.71	0.67	0.32
$x_1, y_1 \rightarrow z_2$	0.29	0.25	0.14
$x_1, y_2 \rightarrow z_1$	0.58	0.15	0.07
$x_1, y_2 \rightarrow z_2$	0.42	0.09	0.05
$x_2, y_1 \rightarrow z_1$	0.67	0.18	0.08
$x_2, y_1 \rightarrow z_2$	0.33	0.08	0.04
$y_3 \rightarrow z_2$	1.00	0.59	0.31

□

The corresponding flow graph is presented in Fig. 8.

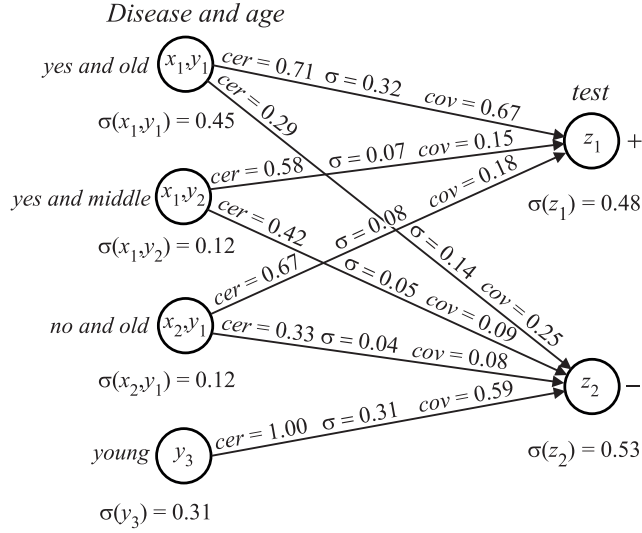


Figure 8: Flow graph for the decision algorithm

From the decision algorithm we can conclude that ill and old patients mostly (71%) have positive test result, ill and middle aged patients not display very clear difference between positive and negative test results (58% and 48% respectively), whereas healthy and old patients display mostly positive test results (67%) and young patients have always (100 %) negative test results.

The above conclusion can be expressed the following "qualitative" decision rules:

- (1) Ill old patients mostly (71%) have positive test results.
- (2) Ill middle aged patients have close positive and negative test results (58%,

42% respectively).

(3) Healthy old patients show mostly (67%) positive test results.

(4) Young patients display always (100%) negative test result.

The inverse decision algorithm yields the following explanation (reasons) for test results:

(i) Positive test result have mostly (67%) ill and old patients.

(ii) Negative test results have mostly (59%) young patients.

The dependency factor between ill and old patients and positive test results amounts to  $\eta \approx 0.19$ , whereas the dependency factor between young patients and negative test results equals  $\eta \approx 0.31$ .

That means that the relationship between young patients and negative test results is more substantial than – between ill old patients and positive test result. □

### 3 Conclusions

We propose in this paper a new approach to knowledge representation and data mining, based on flow analysis in a new kind of flow networks.

We advocate in this paper to represent relationships in data by means of flow graphs. Flow in the flow graph is meant to capture structure of data rather



than to describe any physical material flow in the network. It is revealed that information flow in the flow graph is governed by Bayes' formula, however the formula can be interpreted in entirely deterministic way, without referring to its probabilistic character. This representation allows us to study different relationships in data and can be used as a new mathematical tool for data mining.

## References

- [1] E. A. Adams, *The Logic of Conditionals, an Application of Probability to Deductive Logic*, D. Reidel Publishing Company, Dordrecht, Boston, 1975.
- [2] J. M. Bernardo, A. F. M. Smith, *Bayesian Theory*. Wiley series in probability and mathematical statistics. John Wiley & Sons, Chichester, New York, Brisbane, Toronto, Singapore, 1994.
- [3] M. Berthold, D. J. Hand, *Intelligent Data Analysis - An Introduction*, Springer-Verlag, Berlin, Heidelberg, New York, 1999.
- [4] R. Carnap, *Logical Foundation of Probability*, Routledge and Kegan Paul, London, 1950.

- [5] L. R. Ford, D. R. Fulkerson, *Flows in Networks*. Princeton University Press, Princeton. New Jersey, 1962.
- [6] S. Greco, Z. Pawlak, R. Sowiński, Generalized decision algorithms, rough inference rules and flow graphs, in: J.J. Alpigini, J. F. Peters, A. Skowron, N. Zhong (eds), *Rough Sets and Current Trends in Computing. Lecture Notes in Artificial Intelligence 2475*, Springer-Verlag, Berlin, 2002, pp. 93-104.
- [7] S. Greco, Z. Pawlak, R. Sowiński, Bayesian confirmation measures within rough set approach, *Proceedings of the RSCTS 2004*, Springer-Verlag, 2004, pp. 261-170.
- [8] P. S. Laplace, *Théorie Analytique des Probabilités*, Paris, 1812.
- [9] J. Łukasiewicz, *Die logischen Grundlagen der Wahrscheinlichkeitsrechnung*. Kraków (1913), in: L. Borkowski (Ed.), *Jan Łukasiewicz - Selected Works*, North Holland Publishing Company, Amsterdam, London, Polish Scientific Publishers, Warsaw, 1970.
- [10] Z. Pawlak, Flow graphs - a new paradigm for data mining and knowledge discovery. *JAIST Forum 2004 - Technology Creation Based on Knowledge Science: Theory and Practice*, jointly with KSS2004: 5th Inter-

national Symposium on Knowledge and Systems Science, Proceedings, JAIST, November 2004, pp. 147-153.

- [11] Z. Pawlak, Flow graphs and data mining, Transactions on Rough Sets III, 2005, pp. 1-36, Springer-Verlag, Berlin, Heidelberg.
- [12] F. P. Ramsey, Truth and Probability, In: H. E. Keyburg and H. E. Smokler (eds.) Studies in Subjective Probability, John Wiley and Sons, New York, 1965.
- [13] H. Reichenbach, Wahrscheinlichkeitslehre: eine Untersuchung ber die logischen und mathematischen Grundlagen der Wahrscheinlichkeitsrechnung, 1935, (English translation The theory of probability, an inquiry into the logical and mathematical foundations of the calculus of probability, Berkeley: University of California Press, 1948).
- [14] S. Tsumoto, H. Tanaka, Discovery of Functional Components of Proteins Based on PRIMEROSE and Domain Knowledge Hierarchy, Proceedings of the Workshop on Rough Sets and Soft Computing (RSSC-94), 1994: Lin, T.Y., and Wildberger, A.M. (Eds.), Soft Computing, SCS, 1995, pp. 280-285.

- [15] S. K. M. Wong, W. Ziarko, Algorithm for inductive learning. Bull. Polish Academy of Sciences 34 (5-6), 1986, pp. 271-276.