# Transactions on
# Rough Sets III

James F. Peters · Andrzej Skowron

**Editors-in-Chief**

**Springer**

James F. Peters   Andrzej Skowron (Eds.)

# Transactions on
# Rough Sets III

≜ Springer

Volume Editors

James F. Peters
University of Manitoba
Department of Electrical and Computer Engineering
Winnipeg, Manitoba R3T 5V6, Canada
E-mail: jfpeters@ee.umanitoba.ca

Andrzej Skowron
Warsaw University
Institute of Mathematics
Banacha 2, 02-097 Warsaw, Poland
E-mail: skowron@mimuw.edu.pl

# Lecture Notes in Computer Science 3400

*Commenced Publication in 1973*
Founding and Former Series Editors:
Gerhard Goos, Juris Hartmanis, and Jan van Leeuwen

## Editorial Board

# Preface

Volume III of the Transactions on Rough Sets (TRS) introduces advances in the theory and application of rough sets. These advances have far-reaching implications in a number of research areas such as approximate reasoning, bioinformatics, computer science, data mining, engineering (especially, computer engineering and signal analysis), intelligent systems, knowledge discovery, pattern recognition, machine intelligence, and various forms of learning. This volume reveals the vigor, breadth and depth in research either directly or indirectly related to the rough sets theory introduced by Prof. Zdzisław Pawlak more than three decades ago. Evidence of this can be found in the seminal paper on data mining by Prof. Pawlak included in this volume. In addition, there are eight papers on the theory and application of rough sets as well as a presentation of a new version of the Rough Set Exploration System (RSES) tool set and an introduction to the Rough Set Database System (RSDS).

Prof. Pawlak has contributed a pioneering paper on data mining to this volume. In this paper, it is shown that information flow in a flow graph is governed by Bayes' rule with a deterministic rather than a probabilistic interpretation. A cardinal feature of this paper is that it is self-contained inasmuch as it not only introduces a new view of information flow but also provides an introduction to the basic concepts of flow graphs. The representation of information flow introduced in this paper makes it possible to study different relationships in data and establishes a basis for a new mathematical tool for data mining.

In addition to the paper by Prof. Pawlak, new developments in rough set theory are represented by five papers that investigate the validity, confidence and coverage of rules in approximation spaces (Anna Gomolińska), decision trees considered in the context of rough sets (Mikhail Ju. Moshkov), study of approximation spaces and information granulation (Andrzej Skowron, Roman Świniarski and Piotr Synak), a new interpretation of rough sets based on inverse probabilities and the foundations for a rough Bayesian model (Dominik Ślęzak), and formal concept analysis and rough set theory considered relative to topological approximations (Marcin Wolski). The theory papers in this volume are accompanied by four papers on applications of rough sets: knowledge extraction from electronic devices for power system substation event analysis and decision support (Ching-Lai Hor and Peter Crossley), processing of musical data using rough set methods, RSES and neural computing (Bożena Kostek, Piotr Szczuko, Paweł Żwan and Piotr Dalka), computational intelligence in bioinformatics (Sushmita Mitra), and an introduction to rough ethology, which is based on a biologically inspired study of collective behavior and reinforcement learning in intelligent systems using approximation spaces (James Peters).

This volume also celebrates two landmark events: a new version of RSES and the availability of a Rough Set Database System (RSDS). The introduction of a new version of the Rough Set Exploration System (RSES 2.2) is given in a

paper by Jan G. Bazan and Marcin Szczuka. This paper gives an overview of the basic features of the new version of RSES: improved graphical user interface as well as production of decomposition trees and rules based on training samples. The decomposition tree and rules resulting from training can be used to classify unseen cases. The paper by Zbigniew Suraj and Piotr Grochowalski gives an overview of RSDS, which now includes over 1900 entries and over 800 authors. RSDS includes a number of useful utilities that make it possible for authors to update the database via the web, namely, append, search, download, statistics and help. In addition, RSDS provides access to biographies of researchers in the rough set community.

This issue of the TRS has been made possible thanks to the efforts of a great many generous persons and organizations. We express our thanks to the many anonymous reviewers for their heroic efforts in providing detailed reviews of the articles in this issue of the TRS. The editors and authors of this volume also extend an expression of gratitude to Alfred Hofmann, Ursula Barth, Christine Günther and the other LNCS staff members at Springer for their support in making this volume of the TRS possible. The Editors of this volume have been supported by the Ministry of Scientific Research and Information Technology of the Republic of Poland, Research Grant No. 3T11C00226, and the Natural Sciences and Engineering Research Council of Canada (NSERC), Research Grant No. 185986.

January 2005

James F. Peters
Andrzej Skowron

# LNCS Transactions on Rough Sets

This journal subline has as its principal aim the fostering of professional exchanges between scientists and practitioners who are interested in the foundations and applications of rough sets. Topics include foundations and applications of rough sets as well as foundations and applications of hybrid methods combining rough sets with other approaches important for the development of intelligent systems.

The journal includes high-quality research articles accepted for publication on the basis of thorough peer reviews. Dissertations and monographs up to 250 pages that include new research results can also be considered as regular papers. Extended and revised versions of selected papers from conferences can also be included in regular or special issues of the journal.

# Lecture Notes in Computer Science

For information about Vols. 1–3382

please contact your bookseller or Springer

¥679.68元

# Table of Contents

# Flow Graphs and Data Mining

Zdzisław Pawlak[1,2]

[1] Institute for Theoretical and Applied Informatics,
Polish Academy of Sciences,
ul. Bałtycka 5, 44-100 Gliwice, Poland
[2] Warsaw School of Information Technology,
ul. Newelska 6, 01-447 Warsaw, Poland
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 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 flow distribution in a network rather than the optimal flow which is used for information flow examination in decision algorithms. It is revealed that flow in a flow graph is governed by Bayes' rule, but the rule has an entirely deterministic interpretation without referring to its probabilistic roots. Besides, a decision algorithm induced by a flow graph and dependency between conditions and decisions of decision rules is introduced and studied, which is used next to simplify decision algorithms.

**Keywords:** flow graph, data mining, knowledge discovery, decision algorithms.

## Introduction

In this paper we propose a new approach to data analysis (mining) 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 [4] for optimal flow analysis and they model rather flow *distribution* in a network, than the optimal flow.

The flow graphs considered in this paper are not meant to model physical media (e.g., water) flow analysis, but to model 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., [14, 15]) but they were used first by Łukasiewicz [8] in his study of logic and probability.

This interpretation, in particular, leads to a new look at Bayes' theorem. 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 [1, 3, 13].

This paper is a continuation of some of the authors' ideas presented in [10, 11], where the relationship between Bayes' rule and flow graphs has been introduced and studied (see also [6, 7]).

This paper consists of two parts. Part one introduces basic concepts of the proposed approach, i.e., flow graph and its fundamental properties. It is revealed that flow in a flow graph is governed by Bayes' rule, but the rule has an entirely deterministic interpretation that does not refer to its probabilistic roots. In addition, dependency of flow is defined and studied. This idea is based on the statistical concept of dependency but in our setting it has a deterministic meaning.

In part two many tutorial examples are given to illustrate how the introduced ideas work in data mining. These examples clearly show the difference between classical Bayesian inference methodology and the proposed one.

The presented ideas can be used, among others, as a new tool for data mining, and knowledge representation. Besides, the proposed approach throws new light on the concept of probability.

# 1    Flow Graphs

## 1.1    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 coefficients are extended to paths and some classes of sub-graphs called connections. Further a notion of fusion of a flow graph is defined.

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

Observe that in many cases the data flow order, represented in flow graphs, explicitly follows from the problem specification. However, in other cases the relevant order should be discovered from data. This latter issue will be discussed elsewhere.

## 1.2    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*, $\varphi : \mathcal{B} \to R^+$ 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 by $O(x) = \{y \in N : (x, y) \in \mathcal{B}\}$.

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*

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

and *outflow*

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

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

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

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

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

Then, 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* [4].

**Example**

We will illustrate the 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.

For example, $\varphi(x_1) = 600$ means that these are 600 patients suffering from the disease, $\varphi(y_1) = 570$ means that there are 570 old patients $\varphi(z_1) = 471$ means that 471 patients have a positive test result; $\varphi(x_1, y_1) = 450$ means that there 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 detail.

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 $y$. For example, branch $(x_1, z_1)$ is represented by an arrow from $x_1$ to $z_1$.

For example, 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 the 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 $\varphi_+(y_1) = 450 + 120 = 570$. 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 $\varphi(z_1) + \varphi(z_2) = 471 + 529 = 1000$.

**Fig. 1.** Flow graph.

Throughflow of node $y_1$ is equal to $\varphi(y_1) = \varphi(x_1, y_1) + \varphi(x_2, y_1) = \varphi(y_1, z_1) + \varphi(y_2, z_2) = 570$. □

We will now define 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} \to \ <0,1> \ $ is a *normalized flow* of $(x, y)$ and

$$\sigma(x, y) = \frac{\varphi(x, y)}{\varphi(G)} \tag{5}$$

is a *strength* of $(x, y)$. Obviously, $0 \leq \sigma(x, y) \leq 1$. The strength of the branch (multiplied by 100) 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 by

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

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

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

$$\sigma_+(G) = \frac{\varphi_+(G)}{\varphi(G)} = \sum_{x \in I(G)} \sigma_-(x), \qquad (8)$$

$$\sigma_-(G) = \frac{\varphi_-(G)}{\varphi(G)} = \sum_{x \in O(G)} \sigma_+(x). \qquad (9)$$

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.

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 $\sigma(x_1, y_1) = 0.45$ means that 45% of total flow of $x_1$ flows through the branch $(x_1, y_1)$ etc. □

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.

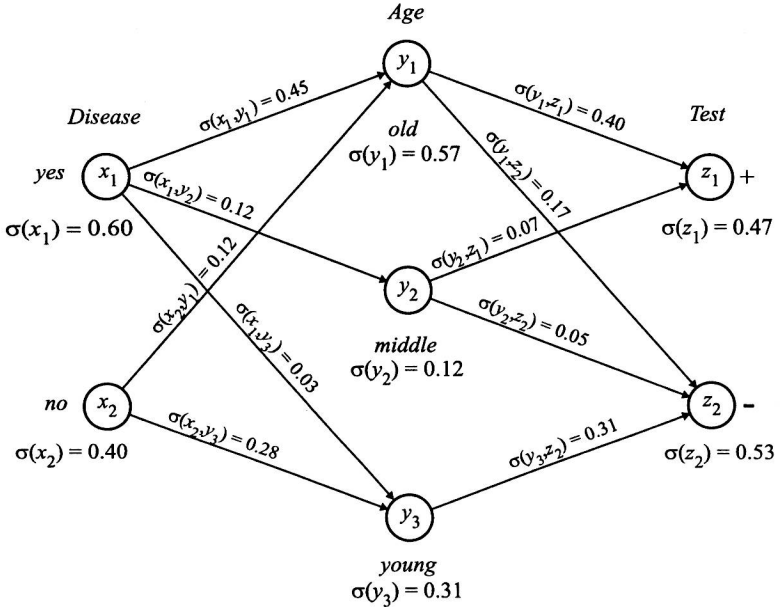**Example (cont.)** The inverted flow graph of the flow graph from Fig. 2 is shown in Fig. 3. □
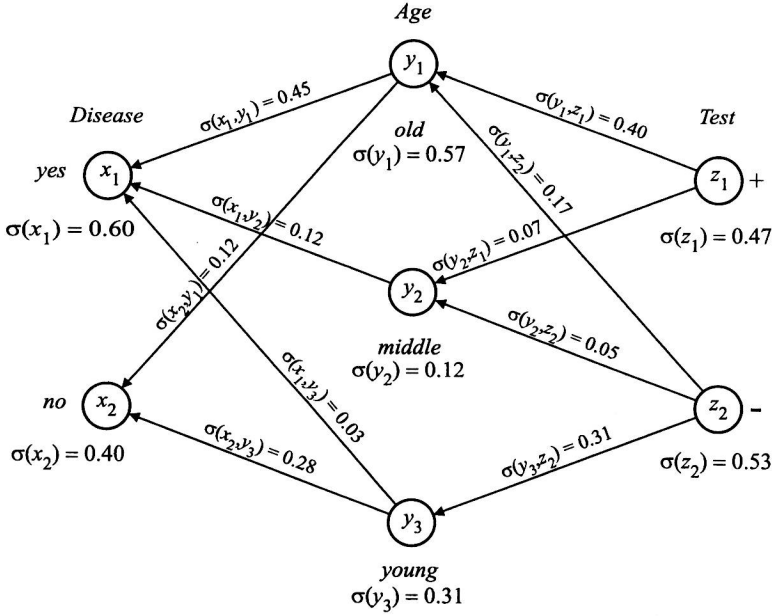


**Fig. 2.** Normalized flow graph.

**Fig. 3.** Inverted flow graph.

## 1.3   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 by

$$cer(x, y) = \frac{\sigma(x, y)}{\sigma(x)}, \tag{10}$$

and

$$cov(x, y) = \frac{\sigma(x, y)}{\sigma(y)}. \tag{11}$$

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. 4.

For example, $cer(x_1, y_1) = \frac{\sigma(x_1, y_1)}{\sigma(x_1)} = \frac{0.45}{0.60} = 0.75$, and $cov(x_1, y_1) = \frac{\sigma(x_1, y_1)}{\sigma(y_1)} = \frac{0.45}{0.57} \approx 0.79$. □

Below some properties of certainty and coverage factors, which are immediate consequences of definitions given above, are presented:
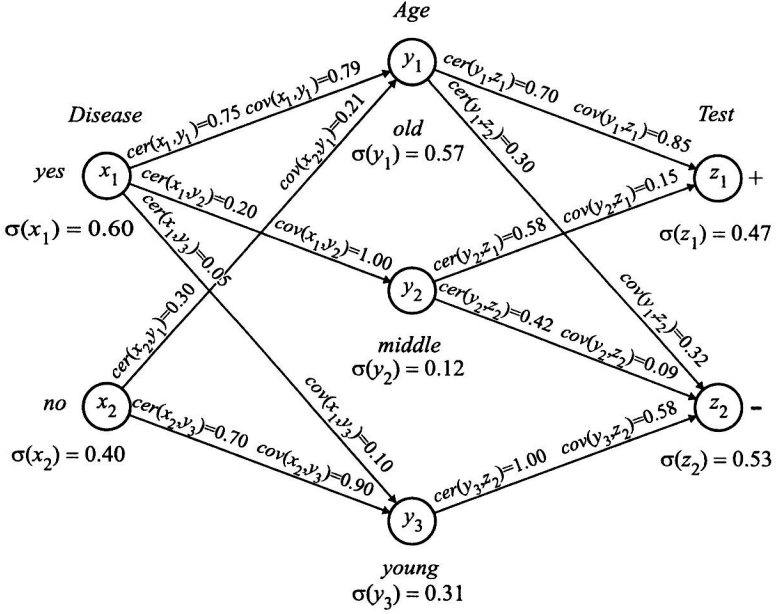
$$\sum_{y \in O(x)} cer(x, y) = 1, \tag{12}$$

**Fig. 4.** Certainty and coverage.

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

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

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

$$cer(x,y) = \frac{cov(x,y)\sigma(y)}{\sigma(x)}, \tag{16}$$

$$cov(x,y) = \frac{cer(x,y)\sigma(x)}{\sigma(y)}. \tag{17}$$

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.

### 1.4   Paths, Connections and Fusion

A (*directed*) *path* from $x$ to $y$, $x \neq y$ in $G$ is a sequence of nodes $x_1, \ldots, x_n$ 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 $[x \ldots y]$.