

# **In Pursuit of Patterns in Data**

**Reasoning from Data  
the Rough Set Way**

**Zdzisław Pawlak**

**Institute of Theoretical and Applied  
Informatics,  
Polish Academy of Sciences,**

**[zpw@ii.pw.edu.pl](mailto:zpw@ii.pw.edu.pl)**

***MOTTO:***

**„It is a capital mistake to theorise  
before one has data”**

**Sherlock Holmes**  
**In: A Scandal in Bohemia**

# AN EXAMPLE OF A DECISION TABLE

---

<i>Fact no.</i>	<i>Driving conditions</i>			<i>Consequence</i>	<i>N</i>
	<i>weather</i>	<i>road</i>	<i>time</i>	<i>accident</i>	
1	<i>misty</i>	<i>icy</i>	<i>day</i>	<i>yes</i>	80
2	<i>foggy</i>	<i>icy</i>	<i>night</i>	<i>yes</i>	140
3	<i>misty</i>	<i>not icy</i>	<i>night</i>	<i>yes</i>	40
4	<i>sunny</i>	<i>icy</i>	<i>day</i>	<i>no</i>	500
5	<i>foggy</i>	<i>icy</i>	<i>night</i>	<i>no</i>	20
6	<i>misty</i>	<i>not icy</i>	<i>night</i>	<i>no</i>	200

# DECISION RULES

---

Let  $S = (U, C, D)$  be a decision table.

Every  $x \in U$  determines a sequence

$c_1(x), \dots, c_n(x), d_1(x), \dots, d_m(x)$  where

$\{c_1, \dots, c_n\} = C$  and  $\{d_1, \dots, d_m\} = D$

- The sequence will be called a *decision rule induced by  $x$*  (in  $S$ ) and denoted by

$c_1(x), \dots, c_n(x) \rightarrow d_1(x), \dots, d_m(x)$  or in short  $C \xrightarrow{x} D$

- The number  $supp_x(C, D) = |C(x) \cap D(x)|$  will be called a *support* of the decision rule  $C \xrightarrow{x} D$

- The number

$$\sigma_x(C, D) = \frac{supp_x(C, D)}{|U|}$$

will be referred to as the *strength* of the decision rule  $C \xrightarrow{x} D$ , where  $|X|$  denotes the cardinality of  $X$

# CERTAINTY AND COVERAGE FACTORS

---

- A *certainty factor* of the decision rule, denoted  $cer_x(C, D)$  is defined as follows:

$$cer_x(C, D) = \frac{|C(x) \cap D(x)|}{|C(x)|} = \frac{\sigma_x(C, D)}{\pi(C(x))}$$

where  $C(x) \neq \emptyset$  and  $\pi(C(x)) = \frac{|C(x)|}{|U|}$

- A *coverage factor* of the decision rule, denoted  $cov_x(C, D)$  is defined as

$$cov_x(C, D) = \frac{|C(x) \cap D(x)|}{|D(x)|} = \frac{\sigma_x(C, D)}{\pi(D(x))}$$

where  $D(x) \neq \emptyset$  and  $\pi(D(x)) = \frac{|D(x)|}{|U|}$

# INVERSE DECISION RULES

---

- If  $C \xrightarrow{x} D$  is a decision rule then  $D \xrightarrow{x} C$  will be called an inverse decision rule
- The inverse decision rules can be used to give explanation (reason) for a decision

# CHARACTERIZATION OF DECISION RULES

---

<i>Fact no.</i>	<i>Strength</i>	<i>Certainty</i>	<i>Coverage</i>
1	0.082	1.000	0.308
2	0.143	0.877	0.538
3	0.041	1.167	0.154
4	0.510	1.000	0.695
5	0.020	0.123	0.027
6	0.204	0.833	0.278

# PROPERTIES OF DECISION RULES

---

Let  $C \xrightarrow{x} D$  be a decision rule. Then the following properties are valid:

$$\sum_{y \in C(x)} cer_y(C, D) = 1 \quad (1)$$

$$\sum_{y \in D(x)} cov_y(C, D) = 1 \quad (2)$$

$$\pi(D(x)) = \sum_{y \in C(x)} cer_y(C, D) \cdot \pi(C(y)) = \sum_{y \in C(x)} \sigma_y(C, D) \quad (3)$$

$$\pi(C(x)) = \sum_{y \in D(x)} cov_y(C, D) \cdot \pi(D(y)) = \sum_{y \in D(x)} \sigma_y(C, D) \quad (4)$$

$$cer_x(C, D) = \frac{cov_x(C, D) \cdot \pi(D(x))}{\pi(C(x))} = \frac{\sigma_x(C, D)}{\pi(C(x))} \quad (5)$$

$$cov_x(C, D) = \frac{cer_x(C, D) \cdot \pi(C(x))}{\pi(D(x))} = \frac{\sigma_x(C, D)}{\pi(D(x))} \quad (6)$$



# GRANULARITY OF DATA AND FLOW GRAPHS

---

- With every decision table we associate a flow graph
- To every decision rule  $C \xrightarrow{x} D$  we assign a directed branch  $x$  connecting the input node  $C(x)$  and the output node  $D(x)$
- Strength of the decision rule represents a throughflow of the corresponding branch.
- The throughflow of the graph is governed by formulas (1),..., (6)

# DECISIONS AND FLOW

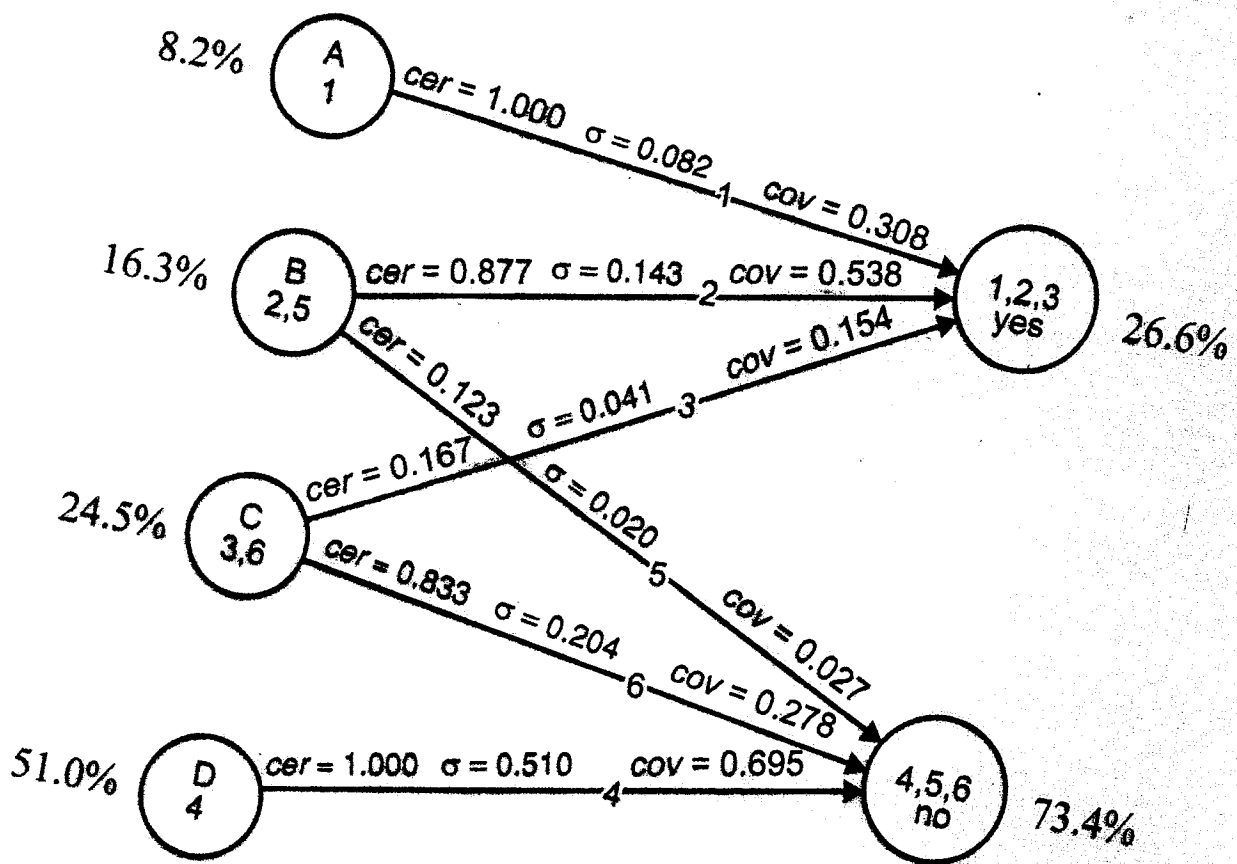
---

- **Classification of objects boils down to finding the maximal output flow in the flow graph**
- **Explanation of decisions is connected with the maximal input flow associated with the given decision**

# FLOW GRAPH

Driving  
Conditions

Accident



# DECISION SPACE

---

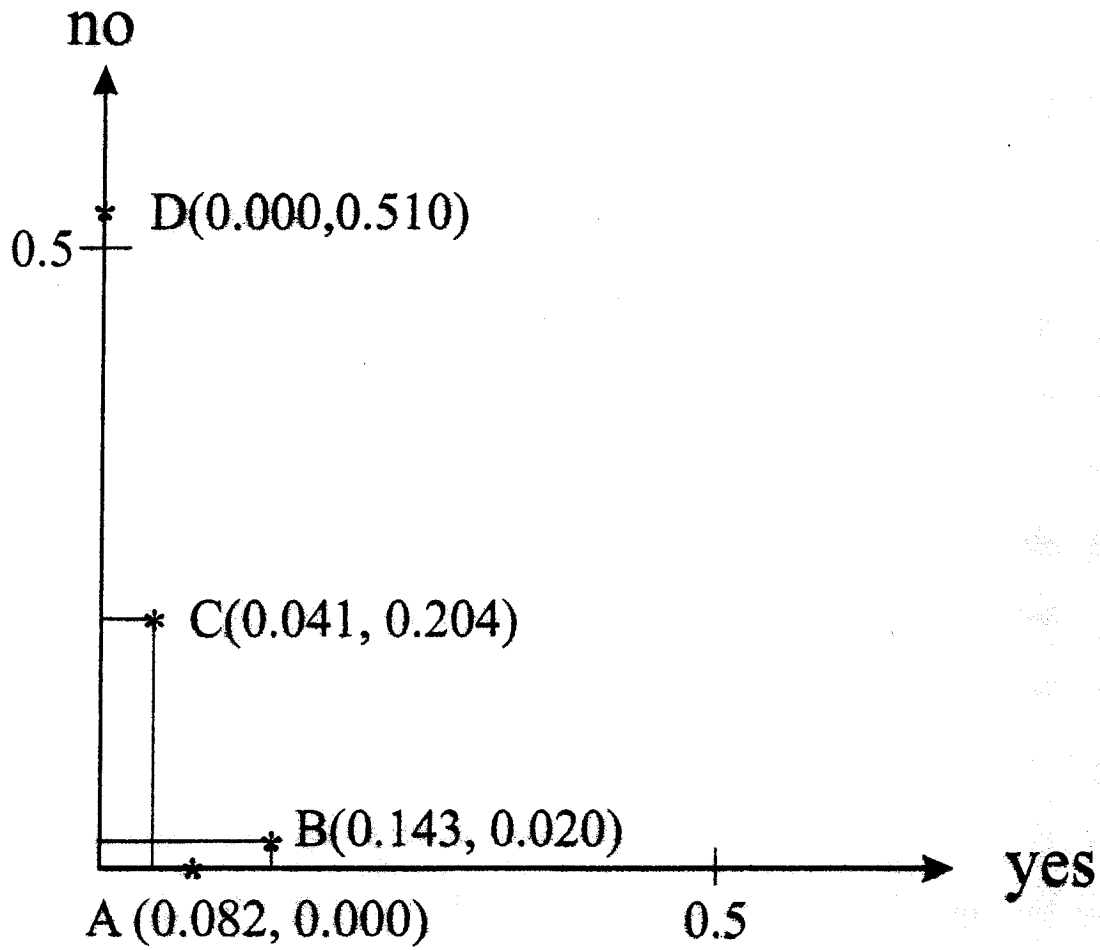
- With every decision table with one  $n$ -valued decision attribute we associate  $n$ -dimensional Euclidean space
- Decision granules determine  $n$  axis of the space
- Condition granules determine points of the space
- Strengths of decision rules are coordinates of granules
- Distance  $\delta(x, y)$  between granules  $x$  and  $y$  is defined as

$$\delta(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

where  $x = (x_1, \dots, x_n)$  and  $y = (y_1, \dots, y_n)$

# DECISION SPACE

---



# DISTANCE MATRIX

---

	<b>A</b>	<b>B</b>	<b>C</b>	<b>D</b>
<b>A</b>				
<b>B</b>	<b>0.064</b>			
<b>C</b>	<b>0.208</b>	<b>0.210</b>		
<b>D</b>	<b>0.517</b>	<b>0.510</b>	<b>0.309</b>	

# **SUPPLY – DEMAND**

---

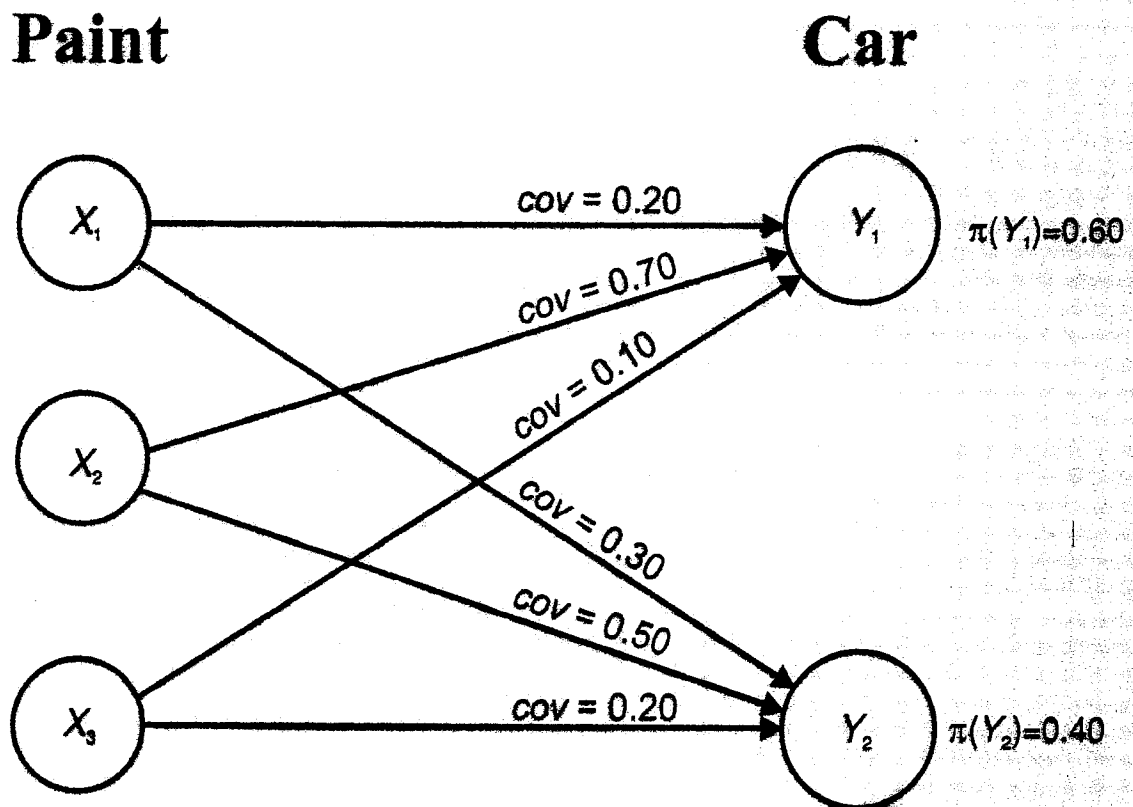
**Suppose that cars are painted into two colors  $Y_1$  and  $Y_2$  and that these colors can be obtained by mixing three paints  $X_1$ ,  $X_2$  and  $X_3$  in the following proportions:**

- $Y_1$  contains 20% of  $X_1$ , 70% of  $X_2$  and 10% of  $X_3$**
- $Y_2$  contains 30% of  $X_1$ , 50% of  $X_2$  and 20% of  $X_3$**

**We have to find demand of each paint and their distribution among colors  $Y_1$  and  $Y_2$**

# SUPPLY – DEMAND

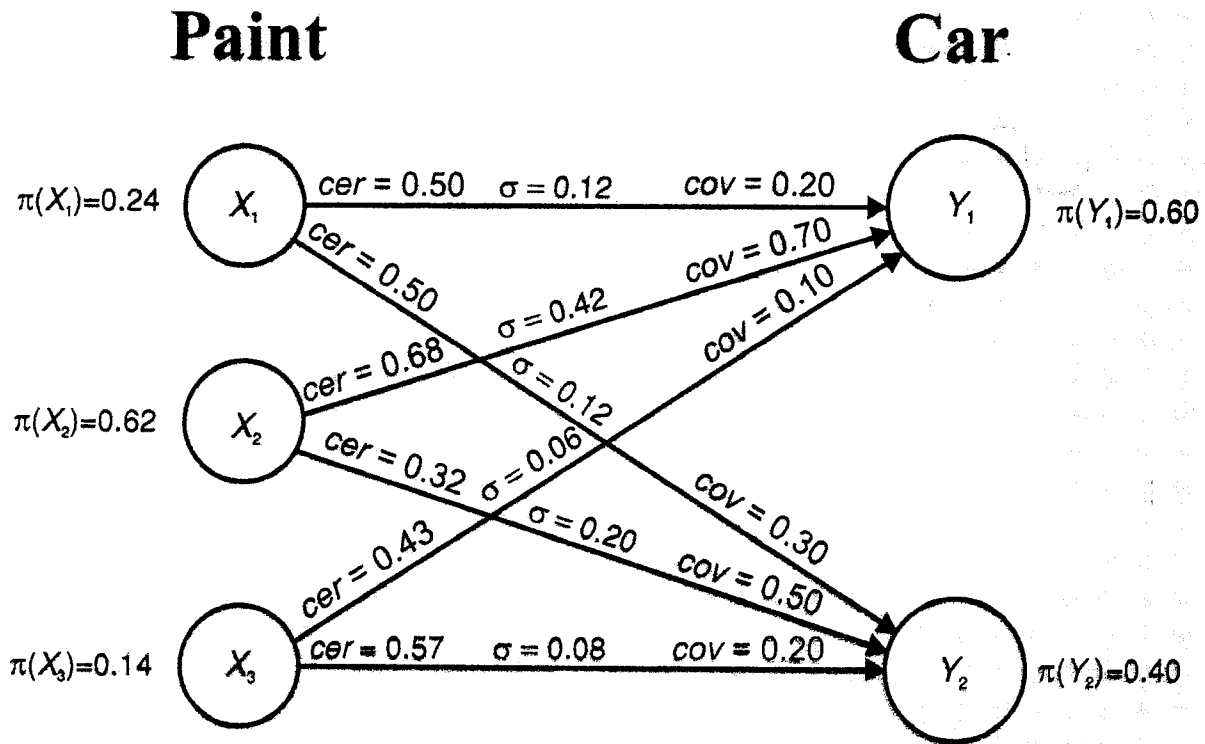
---





# FINAL RESULTS

---



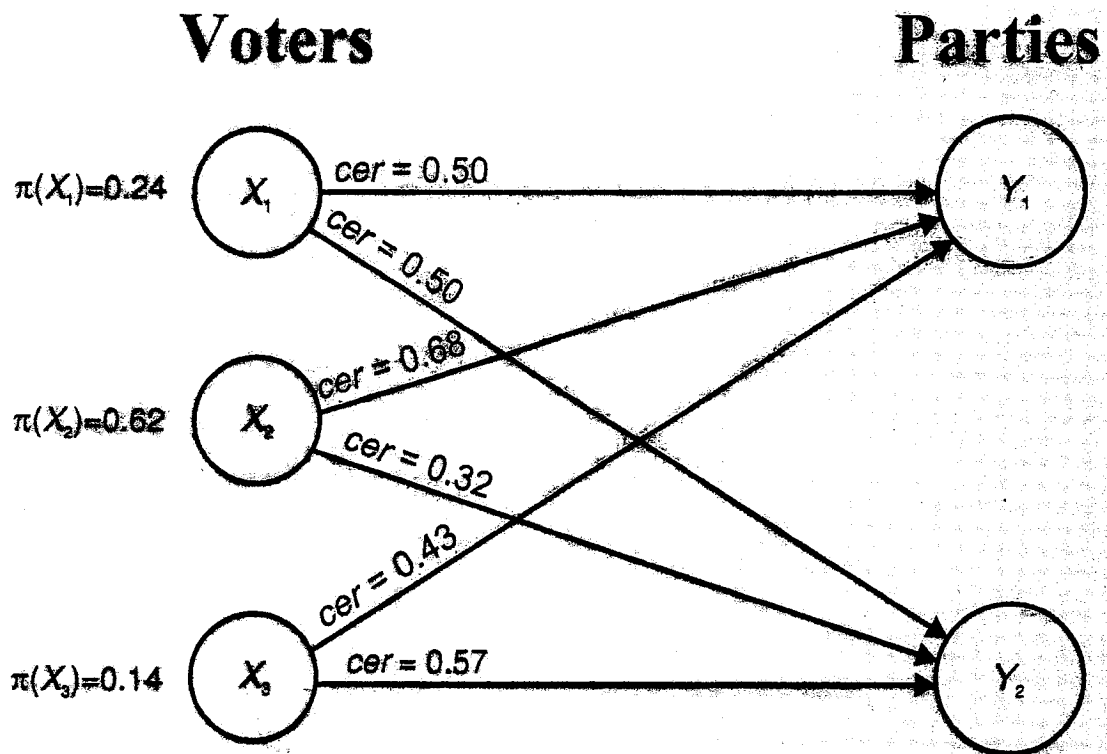
# INVERSE PROBLEM

---

- **Distribution of votes of three disjoint group  $X_1$ ,  $X_2$  and  $X_3$  of voters among two political parties  $Y_1$  and  $Y_2$**
- **$X_1$  consists of 24% of voters,  $X_2$  – 62% and  $X_3$  – 14%**
- **Votes distribution among parties is as follows:**
  - **group  $X_1$  gave 50% of its votes for each party**
  - **group  $X_2$  gave 68% of votes for party  $Y_1$  and 32% for party  $Y_2$**
  - **group  $X_3$  gave 43% votes for party  $Y_1$  and 57% votes for party  $Y_2$**

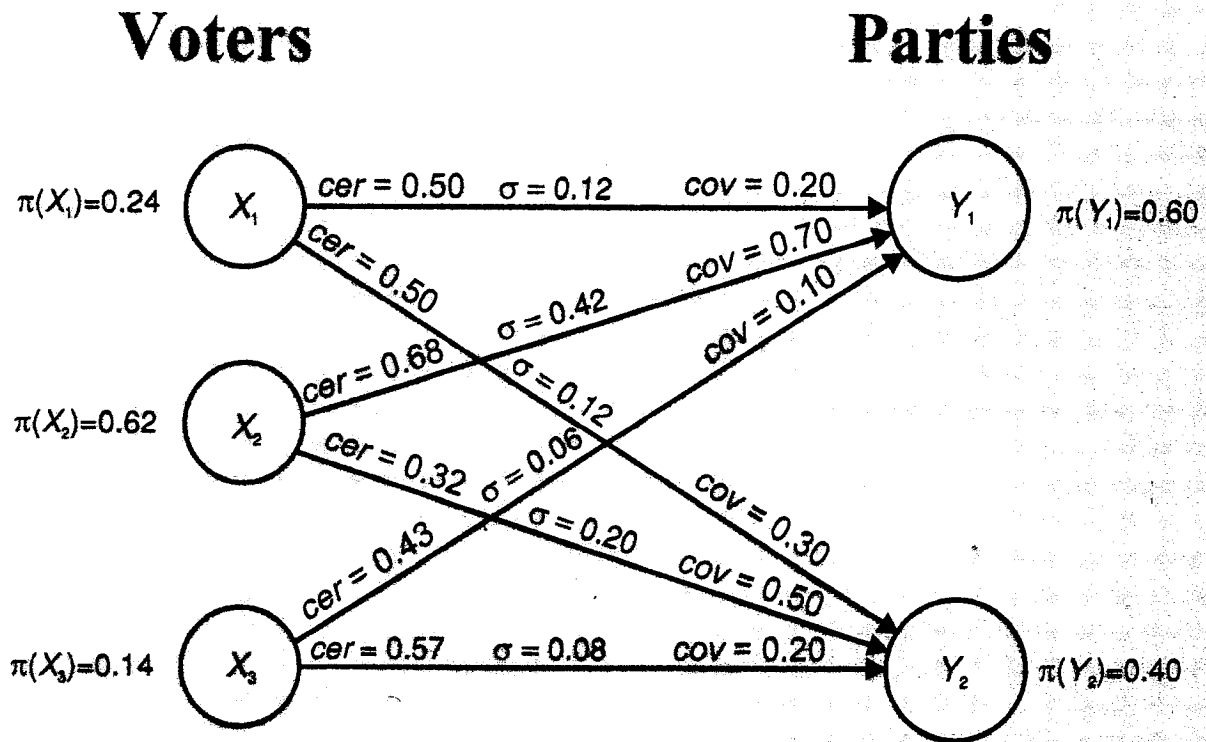
# VOTING ANALYSIS

---



# FINAL RESULTS

---



# FINAL RESULTS

---

- **Party  $Y_1$  obtained 60% votes**
- **Party  $Y_2$  obtained 40% votes**
- **Votes distribution for each party**
  - **Party  $Y_1$  obtained**
    - **20% votes from group  $X_1$ ,**
    - **70% from group  $X_2$  and**
    - **10% from group  $X_3$**
  - **Party  $Y_2$  obtained**
    - **30% votes from group  $X_1$ ,**
    - **50% from group  $X_2$  and**
    - **20% from group  $X_3$**

# FLOW GRAPH

---

A flow graph is a directed, acyclic finite graph

$$G = (N, \mathcal{B}, \varphi)$$

where

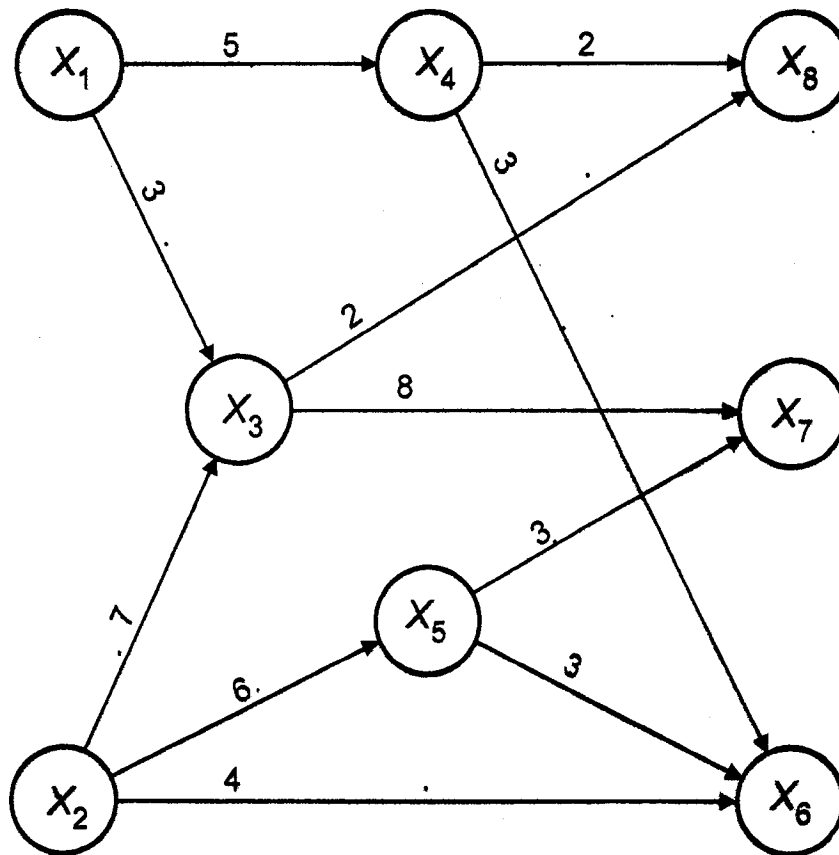
$N$  – set of nodes

$\mathcal{B} \subseteq N \times N$  set of directed branches

$\varphi : \mathcal{B} \rightarrow R^+$  – flow function

# A FLOW GRAPH

---



**flow conservation**  
**inflow = outflow**

# INPUTS AND OUTPUTS

---

- Input of  $x \in N$

$$I(x) = \{y \in N : (y, x) \in \mathcal{B}\}$$

- Output of  $x \in N$

$$O(x) = \{y \in N : (x, y) \in \mathcal{B}\}$$

- Input of  $G$

$$I(G) = \{x \in N : I(x) = \emptyset\}$$

- Output of  $G$

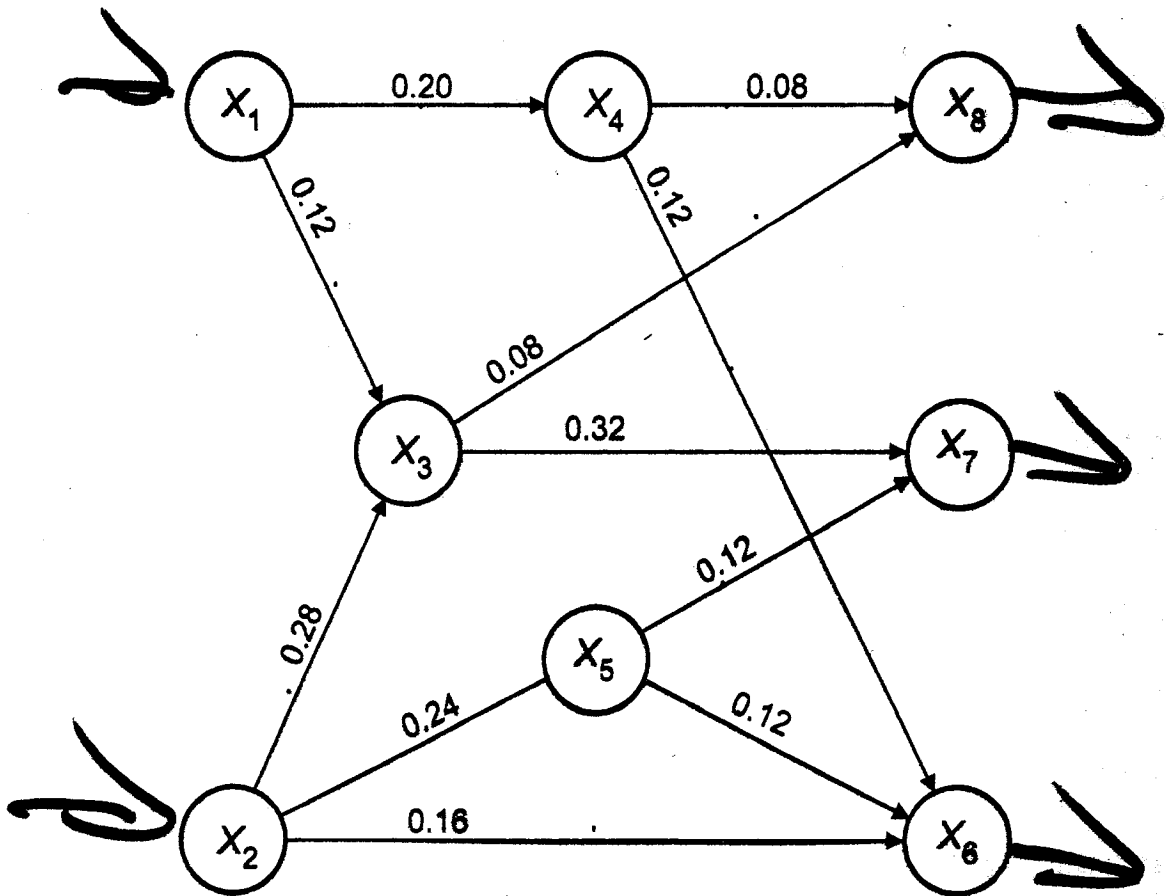
$$O(G) = \{x \in N : O(x) = \emptyset\}$$

- **Inputs and outputs of  $G$  are external nodes of  $G$ ; other nodes are internal nodes of  $G$**



# NORMALIZED FLOW GRAPH

---



$$\text{normalized flow} \equiv \frac{\text{flow}}{\text{total flow}}$$

# FLOW

---

- If  $(x, y) \in \mathcal{B}$  then  $\varphi(x, y)$  is throughflow from  $x$  to  $y$
- $\varphi_+(y) = \sum_{x \in I(y)} \varphi(x, y)$  is an inflow of  $y$
- $\varphi_-(x) = \sum_{y \in O(x)} \varphi(x, y)$  is an outflow of  $x$
- $\varphi_+(G) = \sum_{x \in I(G)} \varphi_-(x)$  is an inflow of  $G$
- $\varphi_-(G) = \sum_{x \in O(G)} \varphi_+(x)$  is an outflow of  $G$

# FLOW CONSERVATION

---

- We assume that for any internal node

$$\varphi_+(x) = \varphi_-(x) = \varphi(x)$$

$\varphi(x)$  – troughflow of  $x$

- Consequently

$$\varphi_+(G) = \varphi_-(G) = \varphi(G)$$

$\varphi(G)$  – troughflow of  $G$

# STRENGTH, CERTAINTY AND COVERAGE OF FLOW

---

- The strength of  $(x, y)$

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

- The certainty of  $(x, y)$

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

- The coverage of  $(x, y)$

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

- The normalized troughflow of  $x$

$$\sigma(x) = \sum_{y \in \theta(x)} \sigma(x, y) = \sum_{y \in I(x)} \sigma(y, x)$$

# PROPERTIES OF FLOW

---

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

$$\bullet \sum_{x \in I(y)} cov(x, y) = 1 \quad (2)$$

$$\bullet cer(x, y) = \frac{cov(x, y)\sigma(y)}{\sigma(x)} \quad (3)$$

$$\bullet cov(x, y) = \frac{cer(x, y)\sigma(x)}{\sigma(y)} \quad (4)$$

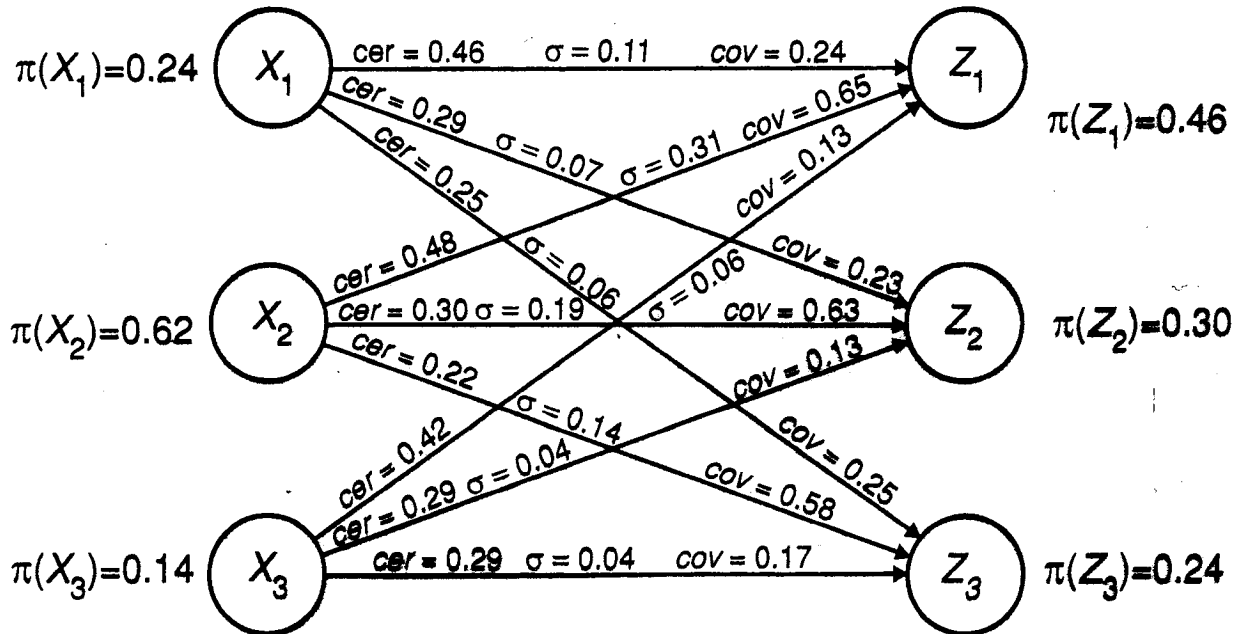
**(3) and (4) are Bayes' formulas**

# SUPPLY – DEMAND SIMPLIFIED GRAPH

---

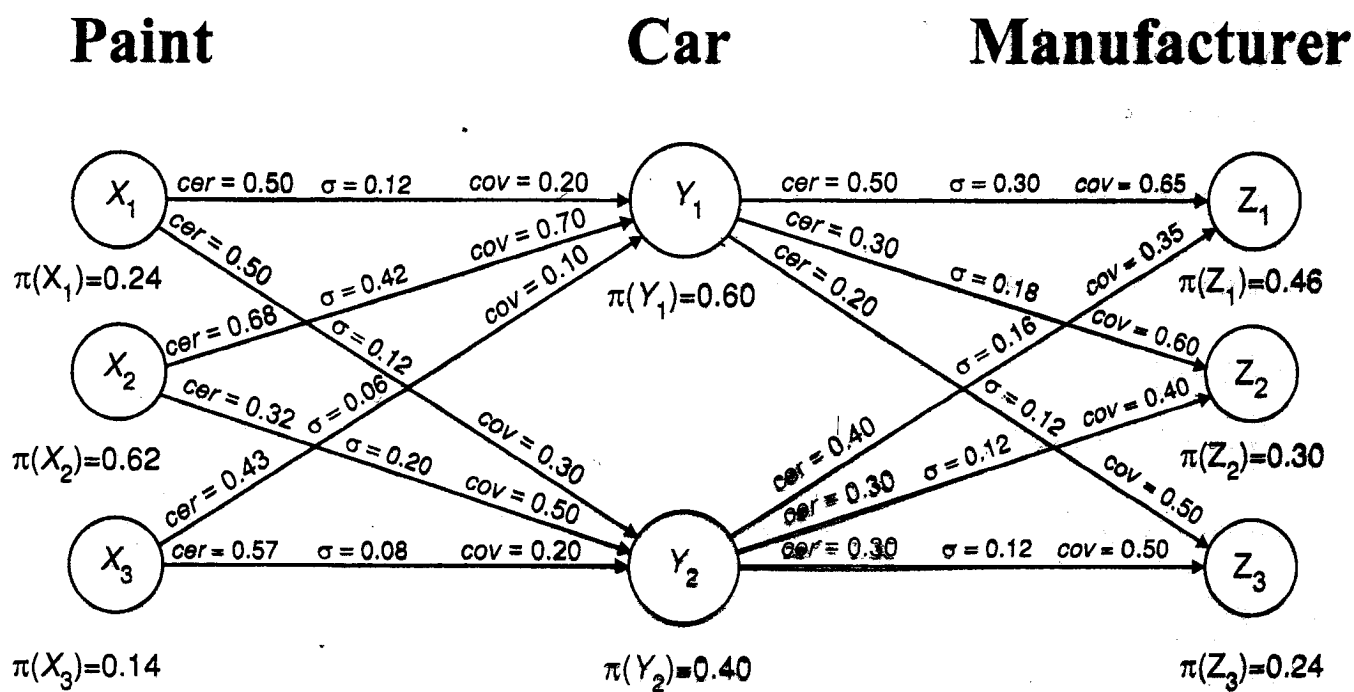
**Paint**

**Manufacturer**



# SUPPLY – DEMAND EXTENDED GRAPH

---



# PATH

---

- A (directed) path from  $x$  to  $y$ ,  $x \neq y$  denoted  $[x, y]$ , 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  $i$ ,  $1 \leq i \leq n-1$

- The certainty of  $[x_1, x_n]$

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

- The coverage of  $[x_1, x_n]$

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

- The strength of  $[x, y]$

$$\sigma[x, y] = \sigma(x) cer[x, y] = \sigma(y) cov[x, y]$$



# CONNECTIONS

---

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

- The certainty of  $\langle x, y \rangle$

$$cer \langle x, y \rangle = \sum_{[x,y] \in \langle x,y \rangle} cer[x,y]$$

- The coverage of  $\langle x, y \rangle$

$$cov \langle x, y \rangle = \sum_{[x,y] \in \langle x,y \rangle} cov[x,y]$$

- The strength of  $\langle x, y \rangle$

$$\sigma \langle x, y \rangle = \sum_{[x,y] \in \langle x,y \rangle} \sigma[x,y]$$

# THE RULE OF SUBSTITUTION

---

Let  $x, y$  ( $x \neq y$ ) be nodes of  $G$ . If we substitute the subgraph  $\langle x, y \rangle$  by a single branch  $(x, y)$  such that

$$\sigma(x, y) = \sigma \langle x, y \rangle$$

then

$$cer(x, y) = cer \langle x, y \rangle$$

$$cov(x, y) = cov \langle x, y \rangle$$

and

$$\varphi(G) = \varphi(G')$$

where  $G'$  is the graph obtained from  $G$  by substituting  $\langle x, y \rangle$  by  $(x, y)$

# DECISION TABLES

---

	<i>Paint</i>	<i>Car</i>	<i>Strength</i>
1	$X_1$	$Y_1$	0.12
2	$X_1$	$Y_2$	0.12
3	$X_2$	$Y_1$	0.42
4	$X_2$	$Y_2$	0.20
5	$X_3$	$Y_1$	0.06
6	$X_3$	$Y_2$	0.08

	<i>Car</i>	<i>Manu.</i>	<i>Strength</i>
1	$Y_1$	$Z_1$	0.30
2	$Y_1$	$Z_2$	0.18
3	$Y_1$	$Z_3$	0.12
4	$Y_2$	$Z_1$	0.16
5	$Y_2$	$Z_2$	0.12
6	$Y_2$	$Z_3$	0.12

	<i>Paint</i>	<i>Manu.</i>	<i>Strength</i>
1	$X_1$	$Z_1$	0.11
2	$X_1$	$Z_1$	0.08
3	$X_1$	$Z_3$	0.06
4	$X_2$	$Z_1$	0.29
5	$X_2$	$Z_2$	0.18
6	$X_2$	$Z_3$	0.14
7	$X_3$	$Z_1$	0.06
8	$X_3$	$Z_2$	0.04
9	$X_3$	$Z_3$	0.04

# SUMMARY

---

- **Flow graphs can be used to decision analysis**
- 
- **Flow in the graph represents strength of decisions**
- **Relation between decisions is expressed by Bayes' formula**
- **In this approach Bayes' formula has entirely deterministic character**
- **The presented approach leads to new computational algorithms and a new look on Bayesian methodology**