



# *Získavanie znalostí z databáz ak sú znalosťami fuzzy pravidlá*

Systemová integrácia 2008  
13. a 14. november

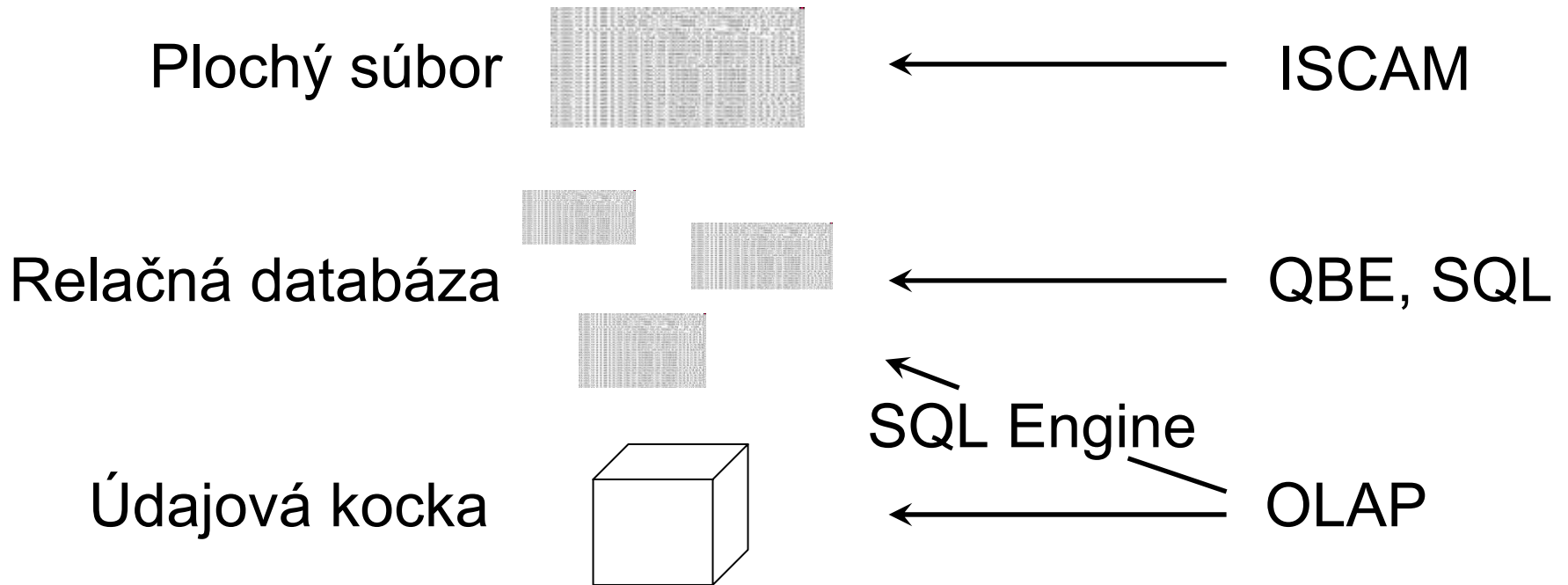
© Ján Boháčik  
Jan.Bohacik@gmail.com

# O čom budeme hovoriť?

---

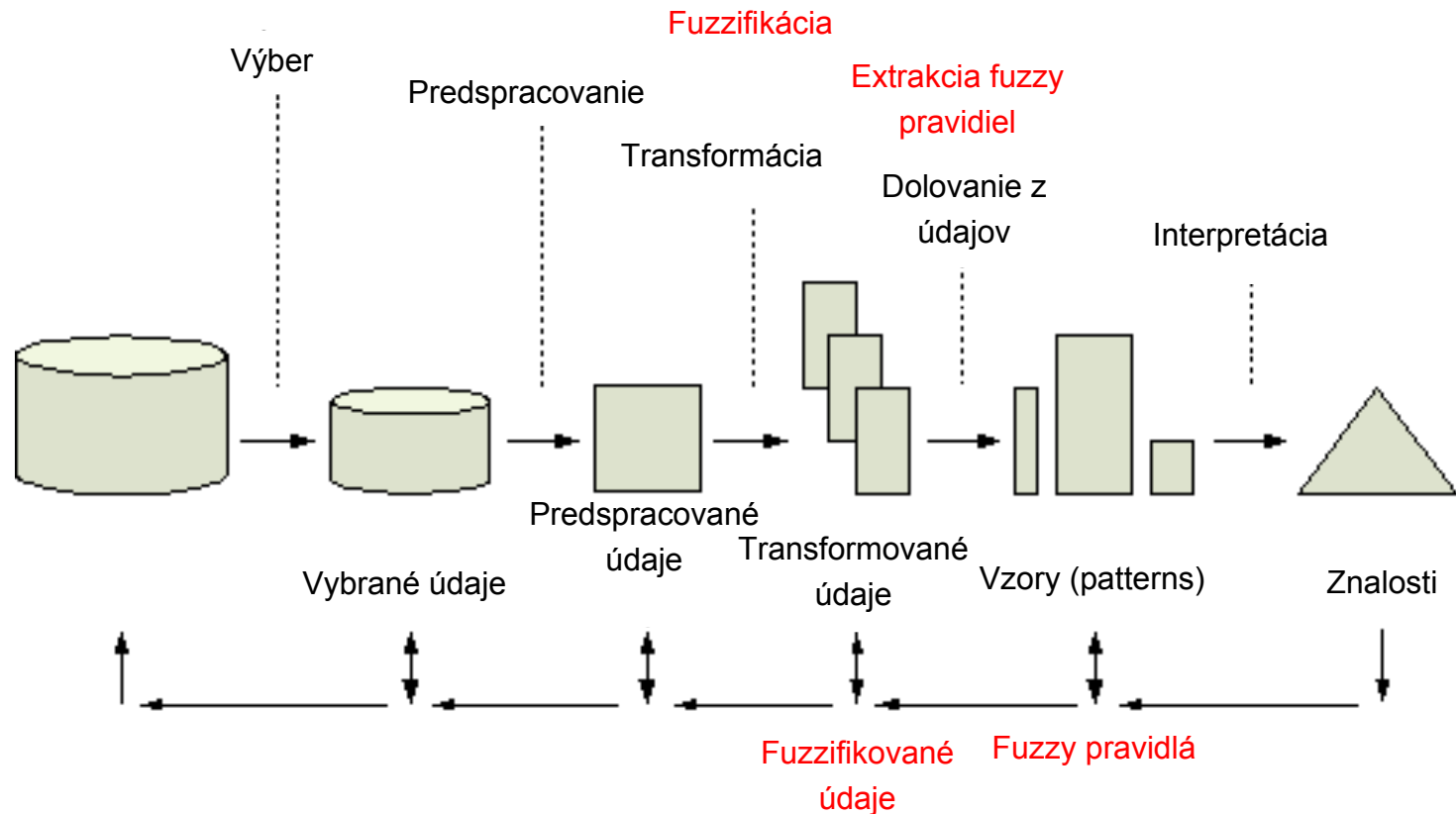
- Vymedzenie a definícia problému.
  - ✓ Kladenie dotazov. Získavanie znalostí z databáz.
  - ✓ Prečo používať fuzzy logiku?
  - ✓ Čo sú znalosti? Klasifikácia pravidiel.
  
- Získavanie fuzzy pravidiel z databáz.
  - ✓ Výber a predspracovanie údajov.
  - ✓ Fuzzifikácia údajov (transformácia údajov).
  - ✓ Interpretácia a vyhodnotenie fuzzy pravidiel.
  
- Zhrnutie a ďalšie plány.

# Databázy a kladenie dotazov



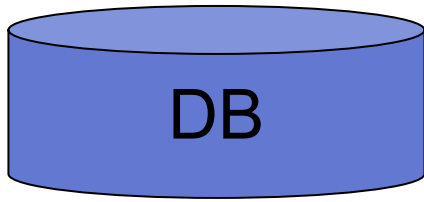
Použiteľné iba ak sa vopred vie, čo sa hľadá!

# Získavanie znalostí z databáz [Fayyad et al., 1996]



[Fayyad et al., 1996] Fayyad, U., Piatetsky-Shapiro, G., Smyth, P., Uthurusamy, R.: Advances in Knowledge Discovery and Data Mining, AAAI Press/MIT Press, Menlo Park, 1996, pp. 37-54.

# Spôsoby reprezentácie znalostí



atd'.

Štatistické metódy

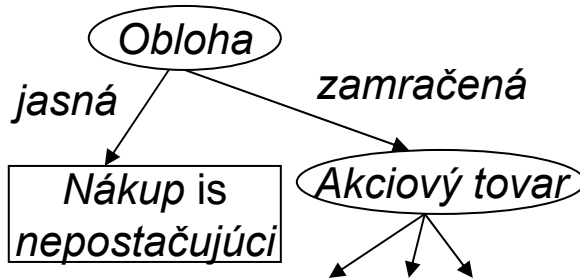
Rozhodovacie tabuľky

Pravidlá IF-THEN

Neurónová sieť



Rozhodovací strom



IF *Obloha is jasná* THEN *Nákup is postačujúci*

IF *Obloha is zamračená AND Akciový tovar is žiadny* THEN *Nákup is postačujúci*

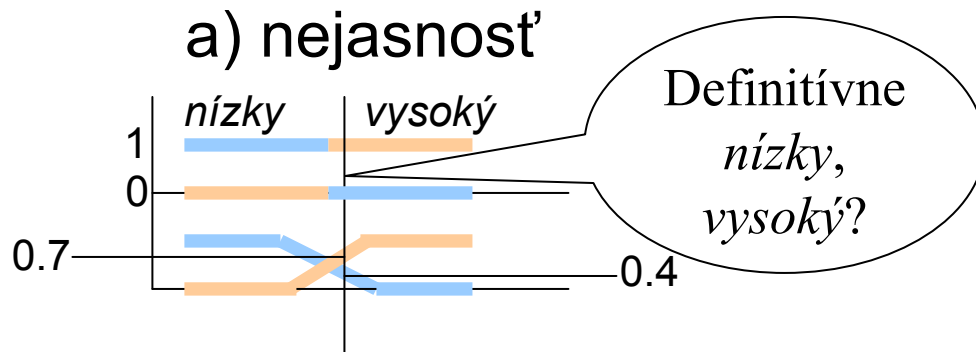
- Klasifikácia
- Nezrozumiteľné znalosti
- Preučenie
- Numerické dáta

- Klasifikácia
- Dôležité atribúty
- Na pravidlá
- Kategoriálne triedy

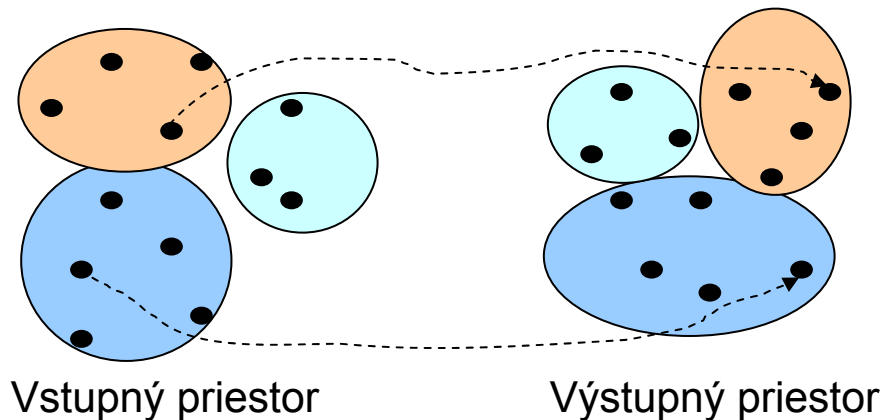
- Klasifikácia, závislosti
- Blízke ľudskému uvažovaniu
- Vytvoriteľné a použiteľné tak ľuďmi ako aj počítačmi
- Mnoho použiteľných variant
- V zásade kategoriálne dáta

# Prečo používať fuzzy logiku?

- Neurčitost' poznávacieho procesu

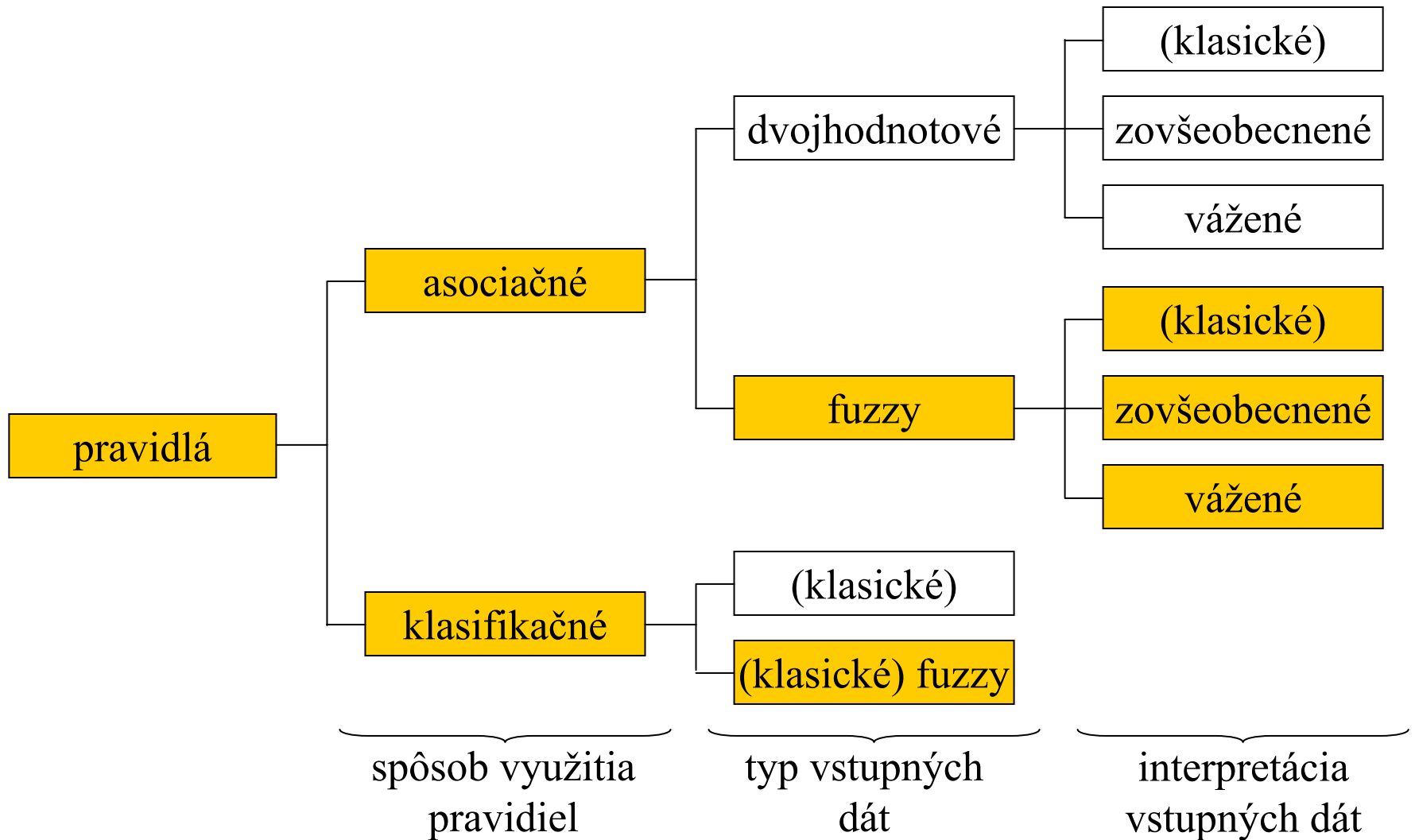


- Redukcia zložitosti



- Ľudia často vôbec nepoužívajú presné údaje

# Klasifikácia pravidiel



# Výber a predspracovanie údajov (1)

---

1. Vytvorenie jednej tabuľky,
2. Odstránenie nadbytočných inštancií a atribútov,
3. Kontrola formátu údajov,
4. Identifikácia odchýlok,
5. Odstránenie chýbajúcich hodnôt,
6. Identifikácia chybných kategoriálnych hodnôt,
7. Tvorba odvodených atribútov.

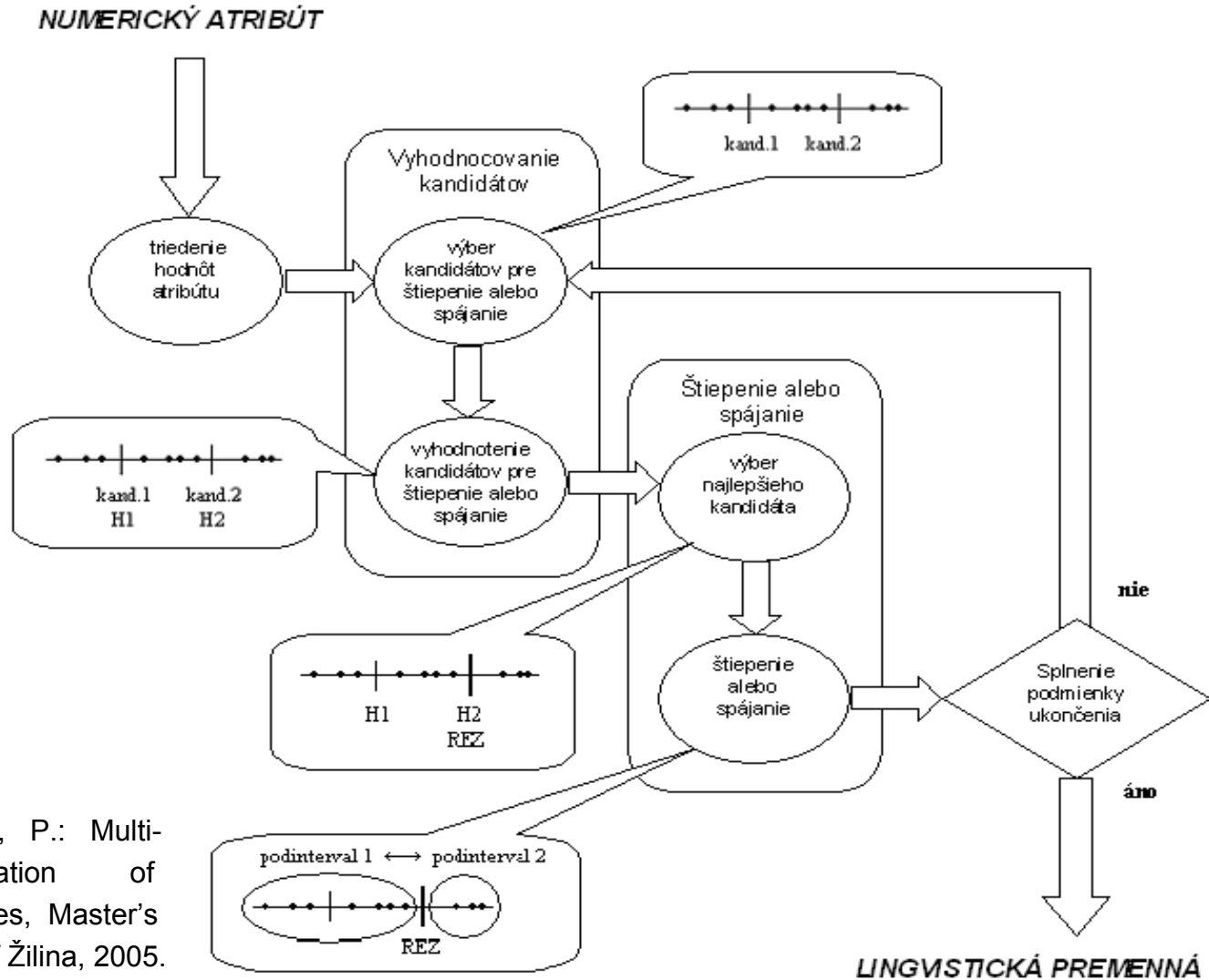


# Výber a predspracovanie údajov (2)

<i>Inšancia</i>	<i>Teplota</i>	<i>Obloha</i>			<i>Novinka</i>	<i>Nákup</i>
		<i>daždivá</i>	<i>zamračená</i>	<i>jasná</i>		
$(e_1)$	15,4543	0,1	0,8	0,1	<i>áno</i>	20 000
$(e_2)$	16,0138	0,7	0,3	0,0	<i>nie</i>	15 000
$(e_3)$	11,9943	0,9	0,1	0,0	<i>áno</i>	14 000
$(e_4)$	14,2764	0,2	0,2	0,6	<i>nie</i>	14 800

Fragment údajov po výbere a predspracovaní

# Fuzzifikácia údajov (1)



[Pazúrik., 2005] Pazúrik, P.: Multi-Interval Discretization of Continuous Attributes, Master's Thesis, University of Žilina, 2005.

# Fuzzifikácia údajov (2)

<i>Inšancia</i>	<i>Teplota</i>			<i>Obloha</i>			<i>Novinka</i>		<i>Nákup</i>	
	<i>nízka</i>	<i>stredná</i>	<i>vysoká</i>	<i>daždivá</i>	<i>zamračená</i>	<i>jasná</i>	<i>nie</i>	<i>áno</i>	<i>nepostačujúci</i>	<i>postačujúci</i>
$(e_1)$	0,0	0,9	0,1	0,1	0,8	0,1	0,0	1,0	0,1	0,9
$(e_2)$	0,2	0,7	0,1	0,7	0,3	0,0	1,0	0,0	0,2	0,8
$(e_3)$	0,5	0,5	0,0	0,9	0,1	0,0	1,0	0,0	0,9	0,1
$(e_4)$	0,1	0,6	0,3	0,2	0,2	0,6	1,0	0,0	0,9	0,1

Fragment údajov po fuzzifikácii

# Fragment databázy web obchodu

V	A <sub>1</sub>			A <sub>2</sub>			A <sub>3</sub>			A <sub>4</sub>		C	
	a <sub>1,1</sub>	a <sub>1,2</sub>	a <sub>1,3</sub>	a <sub>2,1</sub>	a <sub>2,2</sub>	a <sub>2,3</sub>	a <sub>3,1</sub>	a <sub>3,2</sub>	a <sub>3,3</sub>	a <sub>4,1</sub>	a <sub>4,2</sub>	c <sub>1</sub>	c <sub>2</sub>
e <sub>1</sub>	0.2	0.7	0.1	0.3	0.7	0.0	0.2	0.8	0.0	0.0	1.0	0.4	0.6
e <sub>2</sub>	0.9	0.1	0.0	1.0	0.0	0.0	0.8	0.1	0.1	0.6	0.5	0.2	0.8
e <sub>3</sub>	0.8	0.2	0.0	0.6	0.4	0.0	0.0	1.0	0.0	1.0	1.0	0.0	0.7
e <sub>4</sub>	0.0	0.7	0.3	0.8	0.2	0.0	0.1	0.9	0.0	0.8	0.2	0.1	0.9
e <sub>5</sub>	0.0	0.1	0.9	0.7	0.3	0.0	0.3	0.4	0.3	0.5	0.5	1.0	0.0
e <sub>6</sub>	0.0	0.7	0.3	0.0	0.3	0.7	0.7	0.3	0.0	0.8	0.2	0.8	0.2
e <sub>7</sub>	0.9	0.1	0.0	0.2	0.8	0.0	0.1	0.9	0.0	0.0	1.0	1.0	0.0
e <sub>8</sub>	0.0	0.9	0.1	0.0	0.9	0.1	0.1	0.9	0.0	0.7	0.0	1.0	0.0
e <sub>9</sub>	0.0	0.0	1.0	0.0	0.0	1.0	0.6	0.0	0.4	0.8	0.2	1.0	0.0
e <sub>10</sub>	1.0	0.0	0.0	0.5	0.5	0.0	0.0	1.0	0.0	1.0	0.0	0.0	1.0
e <sub>11</sub>	0.0	0.3	0.7	0.0	0.0	1.0	0.0	1.0	0.1	0.9	0.1	1.0	0.0
e <sub>12</sub>	0.0	1.0	0.0	0.0	0.2	0.8	0.2	0.8	0.0	1.0	0.0	0.7	0.3
e <sub>13</sub>	1.0	0.0	0.0	1.0	0.0	0.0	0.3	0.0	0.7	0.6	0.4	0.2	0.8
e <sub>14</sub>	0.9	0.1	0.0	0.0	0.3	0.7	0.0	1.0	0.9	0.1	0.9	0.7	0.3
e <sub>15</sub>	0.7	0.3	0.0	1.0	0.0	0.0	1.0	0.0	0.0	0.8	0.2	0.3	0.7
e <sub>16</sub>	0.2	0.6	0.2	0.0	1.0	0.0	0.0	0.7	0.3	0.7	0.3	0.4	0.6

$$e_i \in V \subseteq U$$

V – známe inštanície

U – úplná množina

$A = \{A_1; A_2; A_3; A_4\} = \{\text{Teplota}; \text{Obloha}; \text{Akciový tovar}; \text{Návštevnosť}\}$

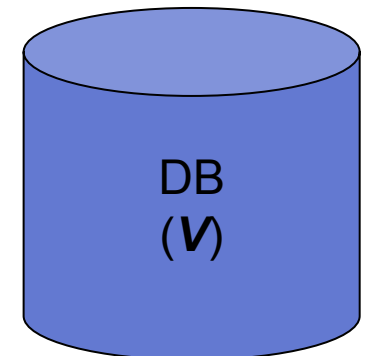
$A_1 = \text{Teplota} = \{a_{1,1}; a_{1,2}; a_{1,3}\} = \{\text{vysoká}; \text{stredná}; \text{nízka}\}$

$A_2 = \text{Obloha} = \{a_{2,1}; a_{2,2}; a_{2,3}\} = \{\text{jasná}; \text{zamračená}; \text{daždivá}\}$

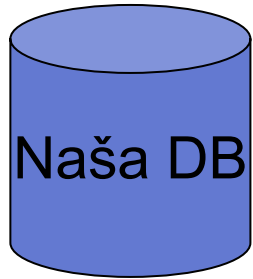
$A_3 = \text{Akciový tovar} = \{a_{3,1}; a_{3,2}; a_{3,3}\} = \{\text{žiaden}; \text{slabý}; \text{bohatý}\}$

$A_4 = \text{Návštevnosť} = \{a_{4,1}; a_{4,2}\} = \{\text{malá}; \text{vysoká}\}$

$C = \text{Nákup} = \{c_1; c_2\} = \{\text{postačujúci}; \text{nepostačujúci}\}$



# Interpretácia a vyhodnotenie (1)



Určujú sa triedne hodnoty novej inštancie  $e_{nová}$ :

U	A <sub>1</sub>			A <sub>2</sub>			A <sub>3</sub>			A <sub>4</sub>		Nákup (C)	
	a <sub>1,1</sub>	a <sub>1,2</sub>	a <sub>1,3</sub>	a <sub>2,1</sub>	a <sub>2,2</sub>	a <sub>2,3</sub>	a <sub>3,1</sub>	a <sub>3,2</sub>	a <sub>3,3</sub>	a <sub>4,1</sub>	a <sub>4,2</sub>	postaču- júci (c <sub>1</sub> )	neposta- čujúci (c <sub>2</sub> )
$e_{nová}$	0.2	0.7	0.1	0.3	0.7	0.0	0.2	0.8	0.0	0.0	1.0	?	?

**IF** *Obloha* is *jasná* **THEN** *Nákup* is *nepostačujúci*

**IF** *Obloha* is *zamračená* **AND** *Akciový tovar* is *žiaden* **THEN** *Nákup* is *postačujúci*

**IF** *Obloha* is *zamračená* **AND** *Teplota* is *vysoká* **AND** *Návštevnosť* is *malá*

**THEN** *Nákup* is *nepostačujúci*

**IF** *Obloha* is *zamračená* **AND** *Teplota* is *vysoká* **AND** *Návštevnosť* is *vysoká*

**THEN** *Nákup* is *postačujúci*

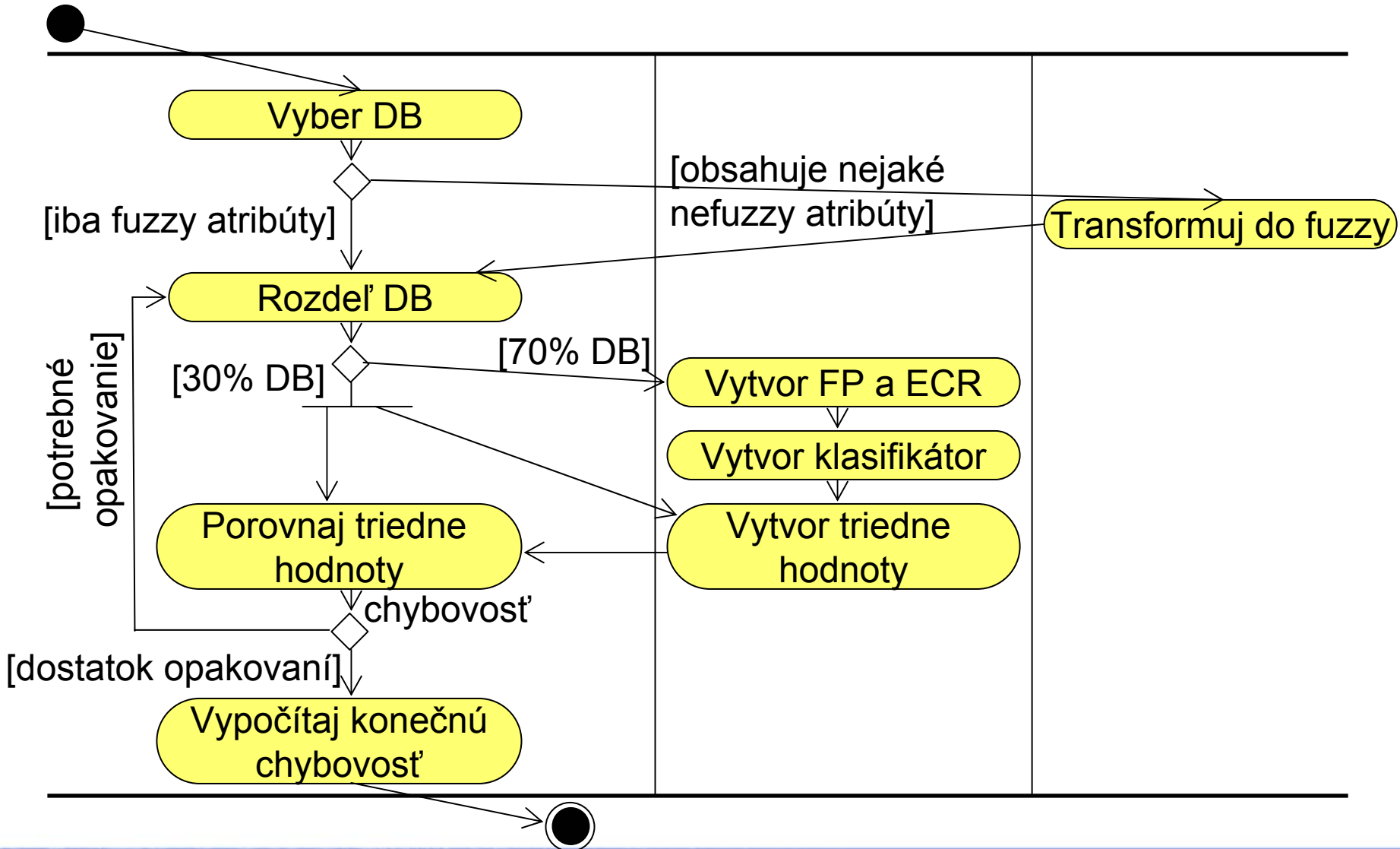
**IF** *Obloha* is *zamračená* **AND** *Teplota* is *stredná* **THEN** *Nákup* is *postačujúci*

**IF** *Teplota* is *nízka* **THEN** *Nákup* is *postačujúci*

**IF** *Obloha* is *zamračená* **AND** *Akciový tovar* is *slabý* **THEN** *Nákup* is *postačujúci*

**IF** *Obloha* is *daždivá* **THEN** *Nákup* is *postačujúci*

# Interpretácia a vyhodnotenie (2)



# Interpretácia a vyhodnotenie (3)

---

**IF** *Návštevnosť is vysoká* **THEN** *Nákup is postačujúci* **AND** *Teplota is nízka*

**IF** *Obloha is daždivá* **THEN** *Nákup is postačujúci*

·  
·  
·

**IF** *Nákup is postačujúci* **AND**

*Akciový tovar is žiaden* **THEN** *Obloha is zamračená*

# Interpretácia a vyhodnotenie (4)

- Jednoduchosť
- Užitočnosť

$$\text{Užitočnosť}(\mathbf{R}) = \frac{1}{\mathbf{M}(\mathbf{V})} \sum_{e \in \mathbf{V}} \mathbf{T}(\mu_{\text{Predpoklad}}(e); \mu_{\text{Záver}}(e)) \in \langle 0; 1 \rangle$$

- Určitosť

$$\text{Určitosť}(\mathbf{R}) = \frac{\sum_{e \in \mathbf{V}} \mathbf{T}(\mu_{\text{Predpoklad}}(e); \mu_{\text{Záver}}(e))}{\sum_{e \in \mathbf{V}} \mu_{\text{Predpoklad}}(e)} \in \langle 0; 1 \rangle$$

- Originalita
- Zisk

$$\text{Zisk}(\mathbf{R}) = \text{Určitosť}(\mathbf{R}) - \frac{1}{\mathbf{M}(\mathbf{V})} \sum_{e \in \mathbf{V}} \mu_{\text{Záver}}(e)$$



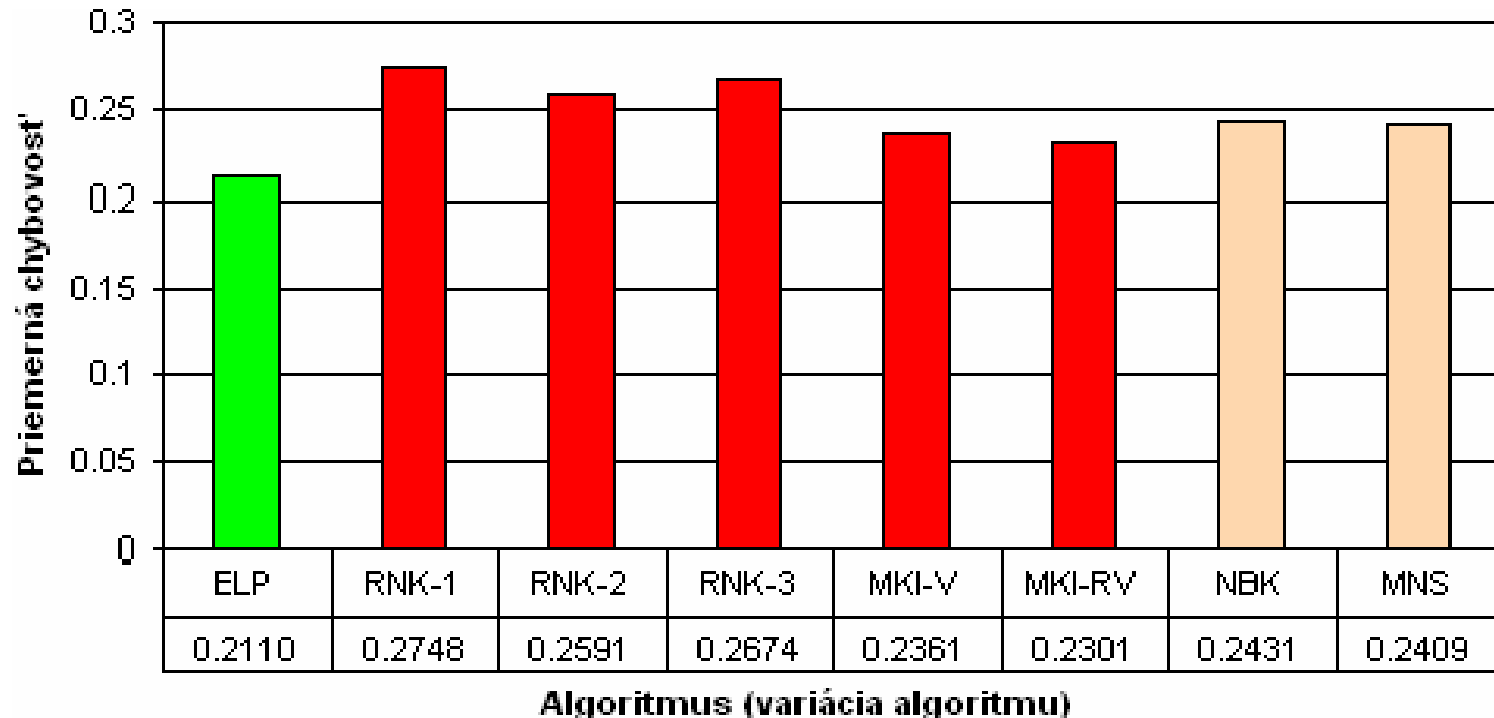
# Eliminácia fuzzy atribútov [Boháčik, 2007]

{fuzzy pravidlo} = vytvorPravidlá( $\alpha; \beta; \mathbf{A}; C; V$ )	
Krok 1	Inicializuj dočasne vytvorené pravidlá $R := \theta$ , dočasné tréningové inštancie $I := V$ , momentálne uvažované atribúty $P := \mathbf{A}$ , maximálny počet povolených atribútov v predpoklade pravidla $maxDĺžka := \mathbf{M}(\mathbf{A})$ , $dĺžka := 1$ a $ponechanýAtribút := false$ .
Krok 2	Nastav eliminovaný atribút $A_{eliminovaný} := \operatorname{argmax}\{\mathbf{NEJED}_\alpha(A_k; C; I) \mid A_k \in P\}$ , $I_1 := \theta$ , $P_1 := P - A_{eliminovaný}$ , $I_2 := \theta$ , $P_2 := P$ . Ak $ponechanýAtribút = false$ , potom nastav $Z := V$ . Inak polož $Z := I$ .
Krok 3	Pre každé $e \in I$ urob: ak neexistuje žiadne $f \in Z$ , $f \neq e$ , také, že pre všetky $A_k \in P_1$ sa $\operatorname{argmax}\{\mu_{a_{k,l}}(f) \mid a_{k,l} \in A_k\} = \operatorname{argmax}\{\mu_{a_{k,l}}(e) \mid a_{k,l} \in A_k\}$ a $\operatorname{argmax}\{\mu_{c_j}(e) \mid c_j \in C\}$ , potom $I_1 := I_1 \cup \{e\}$ . Inak, $I_2 := I_2 \cup \{e\}$ .
Krok 4	Ak $dĺžka + 1 < maxDĺžka$ , choď na Krok 5. Ak nie je, vytvor jedno pravidlo pre každé $e \in I$ a vlož ich bezduplicitne do $R$ . Forma pravidiel je: IF $A_{k_1}$ is $\operatorname{argmax}\{\mu_{a_{k_1,l}}(e) \mid a_{k_1,l} \in A_{k_1}\}$ AND $A_{k_2}$ is $\operatorname{argmax}\{\mu_{a_{k_2,l}}(e) \mid a_{k_2,l} \in A_{k_2}\}$ AND ... AND $A_{k_q}$ is $\operatorname{argmax}\{\mu_{a_{k_q,l}}(e) \mid a_{k_q,l} \in A_{k_q}\}$ THEN $C$ is $\operatorname{argmax}\{\mu_{c_j}(e) \mid c_j \in C\}$ , kde $\{A_{k_1}, A_{k_2}, \dots, A_{k_q}\} = P_1$ ak $e \in I_1$ a $\{A_{k_1}, A_{k_2}, \dots, A_{k_q}\} = P_2$ ak $e \in I_2$ .
Krok 5	Vykonaj Kroky 2-5 jeden krát s $I := I_1$ , $P := P_1$ , $dĺžka := dĺžka + 1$ , $ponechanýAtribút := false$ a jeden krát s $I := I_2$ , $P := P_2$ , $dĺžka := dĺžka + 1$ , $ponechanýAtribút := true$ .
Krok 6	Vypočítaj $UP_0(E_i; c; V)$ pre každé pravidlo IF $E_i$ THEN $C$ is $c_j$ v $R$ . Ako výsledok algoritmu vráť pravidlá, pre ktoré $UP_0(E_i; c; V) \geq \beta$ .

[Boháčik., 2007] Boháčik, J.: Induction by fuzzy attribute elimination, Journal of Information, Control and Management Systems, Vol. 5, No. 2, 2007, pp. 291-301.

# Porovnanie algoritmov

- UCI Repository of ML Databases [Asuncion and Newman, 2007],
- Fuzzifikácia algoritmom [Lee et al., 2001].



[Asuncion and Newman, 2007] UCI Machine Learning Repository [<http://www.ics.uci.edu/~mllearn/MLRepository.html>], Irvine, CA: University of California, School of Information and Computer Science, 2007.

[Lee et al., 2001] Lee, H. M., Chen, C. M., Chen, J. M., JOU, Y. L.: An Efficient Fuzzy Classifier with Feature Selection Based on Fuzzy Entropy. In: IEEE Transactions on Systems, Man, and Cybernetics – Cybernetics 3, 2001, pp.426-432.

# Zhrnutie a záver

---

- Proces Získavania znalostí z databáz:
  - Definovaný v 1996.
  - Odvtedy vzniklo viacero systémov na jeho riešenie, napr. Clementine, DB2 Intelligent Miner, alebo Enterprise Miner.
- Momentálne v tejto oblasti:
  - Vedci sa sústreďujú na metódy pre výber a predspracovanie údajov, interpretáciu či vyhodnotenie získaných znalostí, a metódy pre zdokonalenie existujúcich algoritmov.
  - Ukazuje sa, že využitie fuzzy množín a ďalších princípov fuzzy logiky pomáha odstrániť ich nedostatky.
- Hlavné prínosy:
  - Súhrnné rozpracovanie procesu Získavania znalostí z databáz s ohľadom na vytváranie fuzzy IF-THEN pravidiel vo fáze dolovania.
  - Poukázanie na jednotlivé čiastkové výskumné úlohy, ktoré treba v jednotlivých fázach tohto procesu riešiť.