

# Shluková analýza – specializované algoritmy

---

**Jiří Kléma**

Katedra kybernetiky,  
FEL, ČVUT v Praze



<http://ida.felk.cvut.cz>











































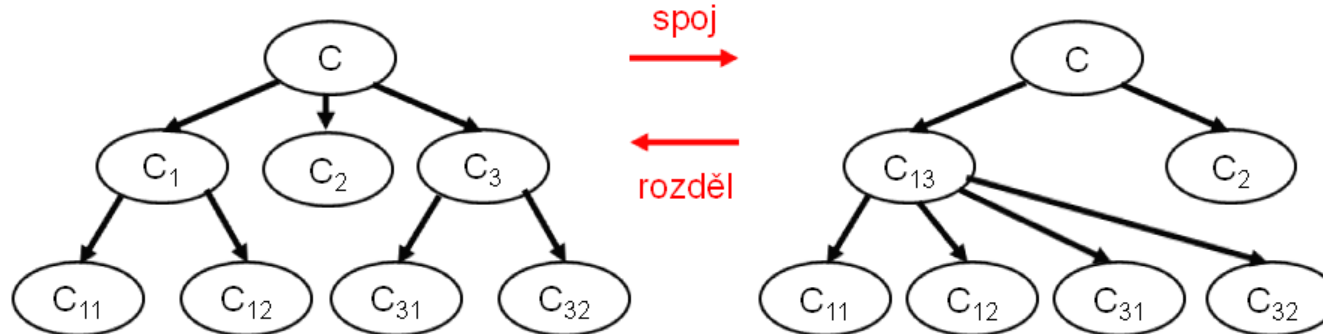






# COBWEB algoritmus

- příklad po příkladu tvoří (klasifikační) strom,
  - vnitřní uzel = (pravděpodobnostní) koncept, list = příklad,
- pro každý jednotlivý příklad volíme jednu z těchto operací
  - vytvoř novou třídu – najdi pozici v hierarchii pro novou podtřídu,
  - zařaď příklad do existující třídy – objekt je podobný objektům v jedné ze tříd,
- abychom minimalizovali vliv pořadí příkladů, zvážíme i operace
  - spojení tříd – dvě třídy nahrazeny jednou,
  - rozdělení třídy – jedna třída se rozpadá na podtřídy, resp. jednotlivé objekty,
- obousměrné gradientní prohledávání řízené funkcí kategoriální užitečnosti.















# Block clustering – US prezidentské volby

- procento hlasů odevzdaných pro republikánského kandidáta na prezidenta,
- jižanské státy v letech 1900-1968.

State	Year																	
	12	36	32	40	44	48	16	04	68	08	24	00	20	28	56	60	52	64
SC	1	1	2	4	4	4	2	5	39	6	2	7	4	9	25	49	49	59
MI	2	3	4	4	6	3	5	5	14	7	8	10	14	18	24	25	40	87
GA	4	13	8	15	18	18	7	18	30	31	18	29	29	45	33	37	30	54
LA	5	11	7	14	19	17	7	10	23	12	20	21	31	24	53	29	47	57
AA	8	13	14	14	18	19	22	21	14	24	27	35	31	48	39	42	35	70
TS	9	12	11	19	17	25	17	22	40	22	20	31	24	52	55	49	53	37
FA	8	24	25	26	30	34	18	21	41	22	28	19	31	57	57	52	55	48
AS	20	18	13	21	30	21	28	40	31	37	29	35	35	39	46	43	44	44
VA	17	29	30	32	37	41	32	37	43	38	33	44	38	54	55	52	56	46
NC	12	27	29	26	33	33	42	40	40	46	40	45	43	55	49	48	46	44
TE	24	31	32	33	39	37	43	43	38	46	44	45	51	54	49	53	50	44
KY	25	40	40	42	45	41	47	47	44	48	49	49	49	59	54	54	50	36
MD	24	37	36	41	48	49	45	49	42	49	45	52	55	57	60	46	55	35
MO	30	38	35	48	48	42	47	50	45	49	50	46	55	56	50	50	51	36
WV	21	39	44	43	45	42	49	55	40	53	49	54	55	58	54	47	48	32
DE	33	43	51	45	45	50	50	54	45	52	58	54	56	65	55	49	52	39

Hartigan: Direct Clustering of a Data Matrix.



# Dvojshlukování – příklad algoritmu

## ■ Cheng a Church

- první algoritmus použitý pro dvojshlukování dat z genových čipů,
- používá upravenou definici rozptylu (obecně reziduí)

$$\sum_{c=1}^k \frac{1}{|I||J|} \sum_{i \in I, j \in J} (x_{ij} - x_{Ij} - x_{iJ} + x_{IJ})^2$$

- \* superponujeme pozadí ( $x_{IJ}$ ), efekt genu ( $x_{iJ}$ ) a efekt biologických podmínek ( $x_{Ij}$ ),
- problém: některé triviální dvojshluky mají nulové reziduum ( $1 \times m$  nebo  $1 \times n$ )
  - \* zavádí prahové reziduum  $\delta$  a hledá největší dvojshluky s reziduem  $< \delta$ ,
- hladový algoritmus (hledá dvojshluky typu (i))
  1. začni s celou maticí,
  2. odstraň řádek či sloupec s největším poklesem rezidua,
  3. opakuj krok 2, skonči pokud reziduum nejde snížit nebo je  $< \delta$ ,
  4. hodnoty v nalezeném dvojshluku znáhodni, opakuj (1-3) pro  $k$  dvojshluků.

- obrázek: aditivní (vlevo) a multiplikativní (vpravo) dvojshluk.

1.0	2.0	5.0	0.0
2.0	3.0	6.0	1.0
4.0	5.0	8.0	3.0
5.0	6.0	9.0	4.0

1.0	2.0	0.5	1.5
2.0	4.0	1.0	3.0
4.0	8.0	2.0	6.0
3.0	6.0	1.5	4.5









