

Clustering (IA)

Classificació conceptual (Michalski & Stepp 1983)

Lògica formal (instància vs objecte mostral)

Variables (atributs) qualitatiu

Input: K nro de classes

- Triar K instàncies (random) com a llavor de cada classe
 - Randomly
 - Strategically
- Associar a cada classe el concepte induït de cada llavor
- (A)ADetectar quin és l'atribut a generalitzar de cada classe tal que
 - Contingui el màxim nro d'objectes observats
 - Minimitzi l'sparsness
 - S'evitin sol.lapaments entre classes
- Avaluar els nous conceptes de classe sobre la matriu d'instàncies i recollir els que satisfan cada concepte
- Ajustar els conceptes de classe estrictament a les instàncies identificades en el pas anterior
- Repetir des d'A fins que totes les instàncies es col.loquin en alguna classe

Clustering (IA)

Conceptual clustering (Michalski & Stepp 1983)

	Age	Weight	cigarettes	Hard attacks
	<i>years</i>	<i>Kg</i>	<i>Pack/week</i>	<i>#</i>
A	30	Low	High	1
B	40	High	Moderate	1
C	30	Medium	Moderate	2
F	40	Low	Low	2
J	30	High	Low	2
M	30	Medium	Low	0
P	40	High	High	0
R	30	Low	Moderate	0
S	50	High	High	2
T	40	Medium	High	2

Choose 2 random seeds: M, R

Associate concepts:

M: [Age=30]&[Weight=Medium]&[Cigarettes=Low]&[H.Attack=0]

R: [Age=30]&[Weight=Low] &[Cigarettes=Moderate]&[H.Attack=0]

Conceptual clustering (Michalski & Stepp 1983)

Goal: Find K=2 classes

Initial random seeds:

M: [Age=30]&[Weight=Medium]&[Cigarettes=Low]&[H.Attack=0]

R: [Age=30]&[Weight=Low] &[Cigarettes=Moderate]&[H.Attack=0]

	Age	Weight	cigarettes	Hard attacks
	<i>years</i>	<i>Kg</i>	<i>Pack/week</i>	<i>#</i>
A	30	Low	High	1
B	40	High	Moderate	1
C	30	Medium	Moderate	2
F	40	Low	Low	2
J	30	High	Low	2
M	30	Medium	Low	0
P	40	High	High	0
R	30	Low	Moderate	0
S	50	High	High	2
T	40	Medium	High	2

Conceptual clustering (Michalski & Stepp 1983)

Goal: Find K=2 classes

Initial random seeds:

M: [Age=30]&[Weight=Medium]&[Cigarettes=Low]&[H.Atack=0]

R: [Age=30]&[Weight=Low] &[Cigarettes=Moderate]&[H.Atack=0]

Generalizations: variables with differences in M and R: Weight, Cigarettes

1) Grow M: weight ≠Low

Reevaluate and find class concept

M1: [Weight=Medium, High]

	Age	Weight	cigarettes	Hard attacks
	<i>years</i>	<i>Kg</i>	<i>Pack/week</i>	<i>#</i>
A	30	Low	High	1
B	40	High	Moderate	1
C	30	Medium	Moderate	2
F	40	Low	Low	2
J	30	High	Low	2
M	30	Medium	Low	0
P	40	High	High	0
R	30	Low	Moderate	0
S	50	High	High	2
T	40	Medium	High	2

Conceptual clustering (Michalski & Stepp 1983)

Goal: Find K=2 classes

Initial random seeds:

M: [Age=30]&[Weight=Medium]&[Cigarettes=Low]&[H.Atack=0]

R: [Age=30]&[Weight=Low] &[Cigarettes=Moderate]&[H.Atack=0]

Generalizations: variables with differences in M and R: Weight, Cigarettes

1) Grow M: weight ≠ Low

M1: [Weight=Medium, High]

2) Grow R: weight ≠ Medium

R2: [weight=Low, High]

	Age	Weight	cigarettes	Hard attacks
	years	Kg	Pack/week	#
A	30	Low	High	1
B	40	High	Moderate	1
C	30	Medium	Moderate	2
F	40	Low	Low	2
J	30	High	Low	2
M	30	Medium	Low	0
P	40	High	High	0
R	30	Low	Moderate	0
S	50	High	High	2
T	40	Medium	High	2

Conceptual clustering (Michalski & Stepp 1983)

Goal: Find K=2 classes

Initial random seeds:

M: [Age=30]&[Weight=Medium]&[Cigarettes=Low]&[H.Attack=0]

R: [Age=30]&[Weight=Low] &[Cigarettes=Moderate]&[H.Attack=0]

Generalizations: variables with differences in M and R: Weight, Cigarettes

1) Grow M: weight ≠ Low

M1: [Weight=Medium, High]

2) Grow R: weight ≠ Medium

R2: [weight=Low, High]

3) Grow M: Cigarettes ≠ Moderate

M3: [Cigarettes=Low, High]

	Age	Weight	cigarettes	Hard attacks
	years	Kg	Pack/week	#
A	30	Low	High	1
B	40	High	Moderate	1
C	30	Medium	Moderate	2
F	40	Low	Low	2
J	30	High	Low	2
M	30	Medium	Low	0
P	40	High	High	0
R	30	Low	Moderate	0
S	50	High	High	2
T	40	Medium	High	2

Conceptual clustering (Michalski & Stepp 1983)

Goal: Find K=2 classes

Initial random seeds:

M: [Age=30]&[Weight=Medium]&[Cigarettes=Low]&[H.Attack=0]

R: [Age=30]&[Weight=Low] &[Cigarettes=Moderate]&[H.Attack=0]

Generalizations: variables with differences in M and R: Weight, Cigarettes

1) Grow M: weight ≠ Low

M1: [Weight=Medium, High]

2) Grow R: weight ≠ Medium

R2: [weight=Medium]&[H.Att ≠ 1]

3) Grow M: Cigarettes ≠ Moderate

M3: [Weight=Low, High]

4) Grow R: Cigarettes ≠ Low

R4: [Cigarette=Moder, High]

	Age	Weight	cigarettes	Hard attacks
	years	Kg	Pack/week	#
A	30	Low	High	1
B	40	High	Moderate	1
C	30	Medium	Moderate	2
F	40	Low	Low	2
J	30	High	Low	2
M	30	Medium	Low	0
P	40	High	High	0
R	30	Low	Moderate	0
S	50	High	High	2
T	40	Medium	High	2

Conceptual clustering (Michalski & Stepp 1983)

Goal: Find $K=2$ classes

Initial random seeds:

M: [Age=30]&[Weight=Medium]&[Cigarettes=Low]&[H.Attack=0]

R: [Age=30]&[Weight=Low] &[Cigarettes=Moderate]&[H.Attack=0]

Generalizations: variables with differences in M and R: Weight, Cigarettes

A) Generalize M over R:

1) Grow M: weight \neq Low: M1: [Weight=Medium, High]

3) Grow M: Cigarettes \neq Moderate: M3: [Cigarettes=Low, High]

B) Generalize R over M:

2) Grow R: weight \neq Medium: R2: [weight=Low, High]

4) Grow R: Cigarettes \neq Low: R4: [Cigarette=Moder, High]

This produces 4 possible descriptions of the 2-class partition:

$P1=\{M1, R2\}$

$P2=\{M1, R4\}$

$P3=\{M3, R2\}$

$P4=\{M3, R4\}$

Some overlapp!!!!!!! Specification!!!!!!

Conceptual clustering (Michalski & Stepp 1983)

Goal: Find K=2 classes

Initial random seeds:

M: [Age=30]&[Weight=Medium]&[Cigarettes=Low]&[H.Atack=0]

R: [Age=30]&[Weight=Low] &[Cigarettes=Moderate]&[H.Atack=0]

Generalizations: variables with differences in M and R: Weight, Cigarettes

P1= {M1, R2}

M1: [Weight=Medium, High]

R2: [Weight=Low, High]

Overlap!!!! Specify R2, weight no high

Reevaluate and find class concept

M1: [Weight=Medium, High]

R1: [Age ≤40] &[Weight=Low]

Extension(P1)={M1={B,C,J,M,P,S,T}
R1={A,F,R}}

		Age	Weight	cigarettes	Hard attacks
		years	Kg	Pack/week	#
	A	30	Low	High	1
	B	40	High	Moderate	1
	C	30	Medium	Moderate	2
	F	40	Low	Low	2
	J	30	High	Low	2
	M	30	Medium	Low	0
	P	40	High	High	0
	R	30	Low	Moderate	0
	S	50	High	High	2
	T	40	Medium	High	2

Conceptual clustering (Michalski & Stepp 1983)

Goal: Find K=2 classes

Initial random seeds:

M: [Age=30]&[Weight=Medium]&[Cigarettes=Low]&[H.Attack=0]

R: [Age=30]&[Weight=Low] &[Cigarettes=Moderate]&[H.Attack=0]

Generalizations: variables with differences in M and R: Weight, Cigarettes

P1={M1,R1}

M1: [Weight=medium, high]

R1: [Age ≤40] &[Weight=Low]

P2= {M1, R4}

M1: [Weight=Low, High]

R4: [cigarettes=Moderate, High]

Overlapping.

Specify R4&weight=Medium

And expand weight=Medium

	Age	Weight	cigarettes	Hard attacks
	years	Kg	Pack/week	#
A	30	Low	High	1
B	40	High	Moderate	1
C	30	Medium	Moderate	2
F	40	Low	Low	2
J	30	High	Low	2
M	30	Medium	Low	0
P	40	High	High	0
R	30	Low	Moderate	0
S	50	High	High	2
T	40	Medium	High	2

Conceptual clustering (Michalski & Stepp 1983)

Goal: Find K=2 classes

Initial random seeds:

M: [Age=30]&[Weight=Medium]&[Cigarettes=Low]&[H.Atack=0]

R: [Age=30]&[Weight=Low] &[Cigarettes=Moderate]&[H.Atack=0]

Generalizations: variables with differences in M and R: Weight, Cigarettes

P1={M1,R1}

M1: [Weight=medium, high]

R1: [Age ≤40] &[Weight=Low]

P2= {M1, R4}

M1: [Weight=Low, High]

R4: [cigarettes=Moderate, High]

Overlapping.

Specify R4 and expand

Reevaluate concepts

M2: [Weight=Low, High]

R2: [Age ≤40]&[weight=Medium]

Extension(P2)={M2={C,MT}

R2={A,B,F,J,P,R,S}}

	Age	Weight	cigarettes	Hard attacks
	years	Kg	Pack/week	#
A	30	Low	High	1
B	40	High	Moderate	1
C	30	Medium	Moderate	2
F	40	Low	Low	2
J	30	High	Low	2
M	30	Medium	Low	0
P	40	High	High	0
R	30	Low	Moderate	0
S	50	High	High	2
T	40	Medium	High	2

Conceptual clustering (Michalski & Stepp 1983)

Goal: Find K=2 classes

Initial random seeds:

M: [Age=30]&[Weight=Medium]&[Cigarettes=Low]&[H.Attack=0]

R: [Age=30]&[Weight=Low] &[Cigarettes=Moderate]&[H.Attack=0]

Generalizations: variables with differences in M and R: Weight, Cigarettes

P3= M3, R2

M3=[Cigarettes=Low, High]

R2=[Weight=Low, High]

Overlapp, AND NO COVER!!!

R3:R2&[cigarette=moderate]

expand [cig=mod]

Reevaluate concepts and expand

M3: [Cigarettes=Low, High]

R3: [Age ≤40] &[Cigarette=Moderate]

Extension(P3)={M3={A,F,J,M,P,S,T}
R3={B,C,R}}

	Age	Weight	cigarettes	Hard attacks
	years	Kg	Pack/week	#
A	30	Low	High	1
B	40	High	Moderate	1
C	30	Medium	Moderate	2
F	40	Low	Low	2
J	30	High	Low	2
M	30	Medium	Low	0
P	40	High	High	0
R	30	Low	Moderate	0
S	50	High	High	2
T	40	Medium	High	2

Conceptual clustering (Michalski & Stepp 1983)

Goal: Find K=2 classes

Initial random seeds:

M: [Age=30]&[Weight=Medium]&[Cigarettes=Low]&[H.Attack=0]

R: [Age=30]&[Weight=Low] &[Cigarettes=Moderate]&[H.Attack=0]

Generalizations: variables with differences in M and R: Weight, Cigarettes

P3=

M3: [Weight=Low, High]

R3: [Age ≤40] &[Cigarette=Moder]

P4= M3, R4

M3: [Cigarette=Low, High]

R3: [Cigarette=Moder,high]

Overlapping. Specify R3, cig no high

Reevaluate concepts

M4: [Age ≤40] & [Cigarette=Low]&
[H.Att≠1]

R4: [Cigarette=Moder, High]

Extension(P4)={M4={F,J,M},
R4={A,B,C,P,R,S,T}}

	Age	Weight	cigarettes	Hard attacks
	years	Kg	Pack/week	#
A	30	Low	High	1
B	40	High	Moderate	1
C	30	Medium	Moderate	2
F	40	Low	Low	2
J	30	High	Low	2
M	30	Medium	Low	0
P	40	High	High	0
R	30	Low	Moderate	0
S	50	High	High	2
T	40	Medium	High	2

Conceptual clustering (Michalski & Stepp 1983)

Goal: Find K=2 classes

Initial random seeds:

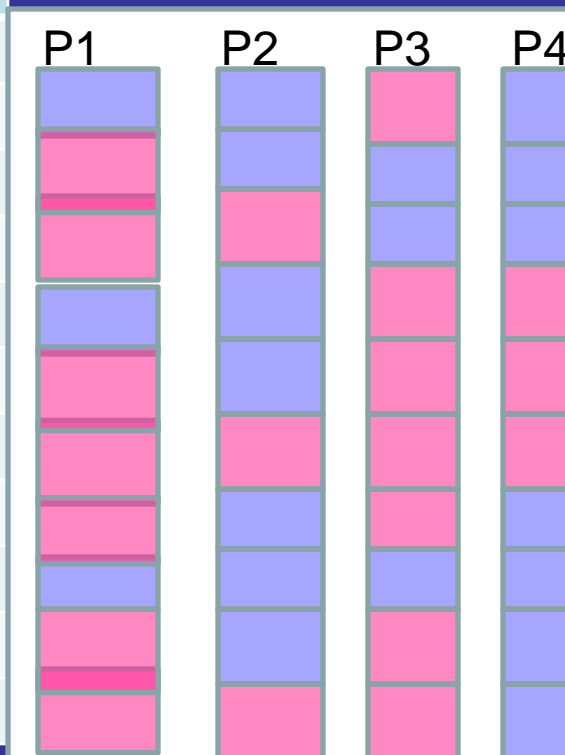
M: [Age=30]&[Weight=Medium]&[Cigarettes=Low]&[H.Atack=0]

R: [Age=30]&[Weight=Low] &[Cigarettes=Moderate]&[H.Atack=0]

Generalizations: variables with differences in M and R: Weight, Cigarettes

Four possible generalizations:

	Age	Weight	cigarettes	Hard attacks
	<i>years</i>	<i>Kg</i>	<i>Pack/week</i>	<i>#</i>
A	30	Low	High	1
B	40	High	Moderate	1
C	30	Medium	Moderate	2
F	40	Low	Low	2
J	30	High	Low	2
M	30	Medium	Low	0
P	40	High	High	0
R	30	Low	Moderate	0
S	50	High	High	2
T	40	Medium	High	2



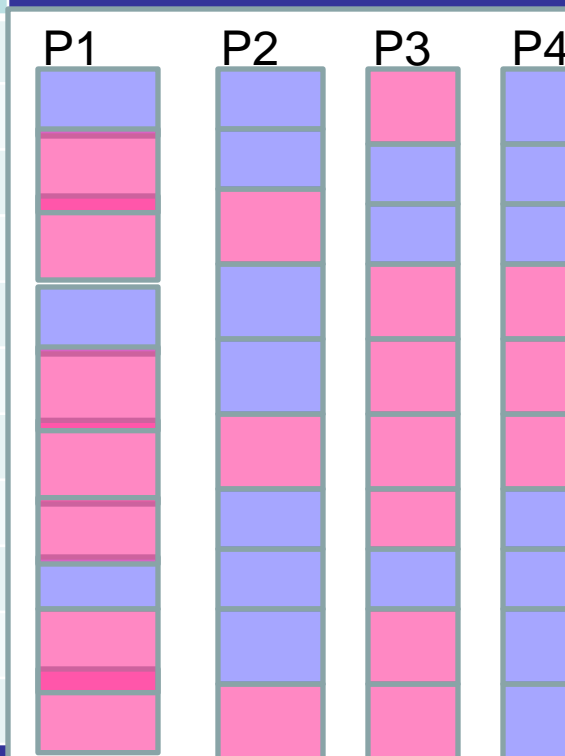
All directly understandable by user

Conceptual clustering (Michalski & Stepp 1983)

Goal: Find $K=2$ classes

Choose best generalization:

	Age	Weight	cigarettes	Hard attacks
	<i>years</i>	<i>Kg</i>	<i>Pack/week</i>	<i>#</i>
A	30	Low	High	1
B	40	High	Moderate	1
C	30	Medium	Moderate	2
F	40	Low	Low	2
J	30	High	Low	2
M	30	Medium	Low	0
P	40	High	High	0
R	30	Low	Moderate	0
S	50	High	High	2
T	40	Medium	High	2



Minimize over all clusterings P : $[\sum_{\alpha \in P} r(\alpha)]/\text{card}P$

Being $r(\alpha) = 1 - [p(\alpha)/t(\alpha)]$ the relative SPARSENESS of a concept α
 t : total events covered by concept α
 p : observed events covered by concept α

Clustering (IA)

Conceptual Clustering (Michalski & Stepp 1983)

Formal logics (instances vs sample objects)

Qualitative Variables (atributes)

Input: K number of classes

- Choose K instances as class seeds
 - Randomly
 - Strategically
- Associate the concept induced from each seed to each class
 - Detect the attributes for possible generalization
 - Generate all possible generalizations, recalculate class concepts
 - Correct overlappings by maximizing the coverage
 - Select the proposal with minimal sparseness
- Respecify the winer to cover strictly the observed cases
- Repeat generalization till sparseness do not improve anymore

Final result: $P = \{C1 = \{R, M, P, B, A\}, C2 = \{F, C, T, J, S\}\}$