

Corso di Laurea Specialistica in Ingegneria Gestionale
a.a. 2005-06

Sistemi Informativi per le Decisioni LS

Clustering Association Rules

Brian Lent - Arun Swami - Jennifer Widom

Prof. Marco Patella

GRUPPO IX

Angelini Andrea

Menenti Alice

Della Vittoria Lucia

Il Framework ARCS

- **Ambiti applicativi:**

- Segmentazione
- Classificazione

- **Input:**

- Large Data Base
- Valore dell'attributo di segmentazione (es.: G = 'BUONA')
- 2 attributi significativi per analizzare le caratteristiche del DB (X, Y)
- Numero di bin (parametro)

- **Output:**

- Cluster di regole associative costruiti su X e Y

X	Y		G
ETA'	SALARIO	FIGLI	SPESA
25	10000	0	BUONA
...
50	35000	2	SCARSA

Cos'è un cluster di regole associative

Dalle seguenti regole associative:

$$(Et\grave{a} = 40) \wedge (\text{Salario} = \$42,350) \Rightarrow (\text{Spesa} = \text{BUONA})$$

$$(Et\grave{a} = 41) \wedge (\text{Salario} = \$57,000) \Rightarrow (\text{Spesa} = \text{BUONA})$$

$$(Et\grave{a} = 41) \wedge (\text{Salario} = \$48,750) \Rightarrow (\text{Spesa} = \text{BUONA})$$

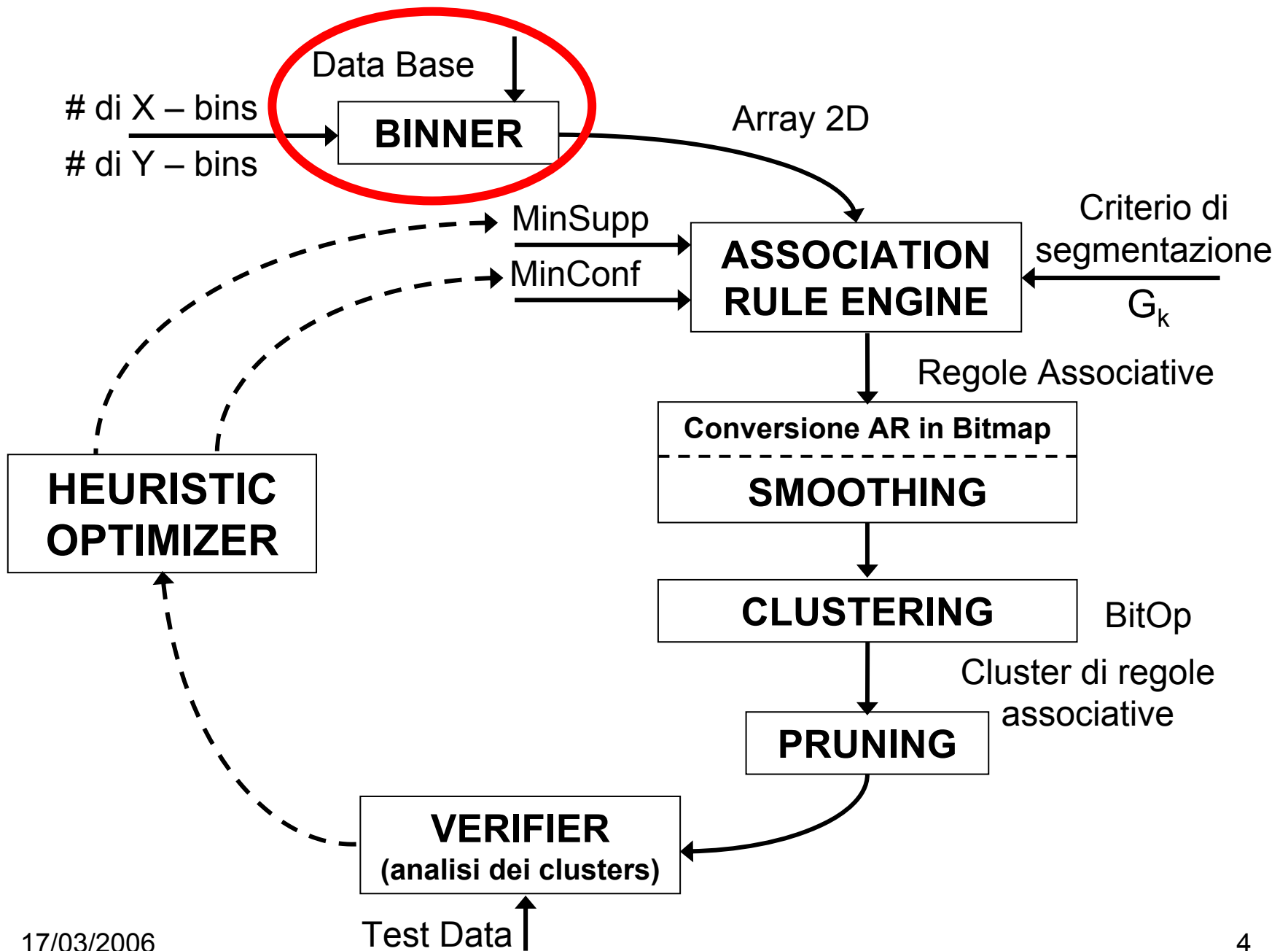
$$(Et\grave{a} = 40) \wedge (\text{Salario} = \$52,600) \Rightarrow (\text{Spesa} = \text{BUONA})$$

Si determina il cluster di regole associative:

$$(40 \leq Et\grave{a} < 42) \wedge (\$40,000 \leq \text{Salario} < \$60,000) \Rightarrow (\text{Spesa} = \text{BUONA})$$

Generalizzando:

$$(\mathbf{X}_1 \leq \mathbf{X} < \mathbf{X}_2) \wedge (\mathbf{Y}_1 \leq \mathbf{Y} < \mathbf{Y}_2) \Rightarrow (\mathbf{G} = \mathbf{G}_K)$$



Binner

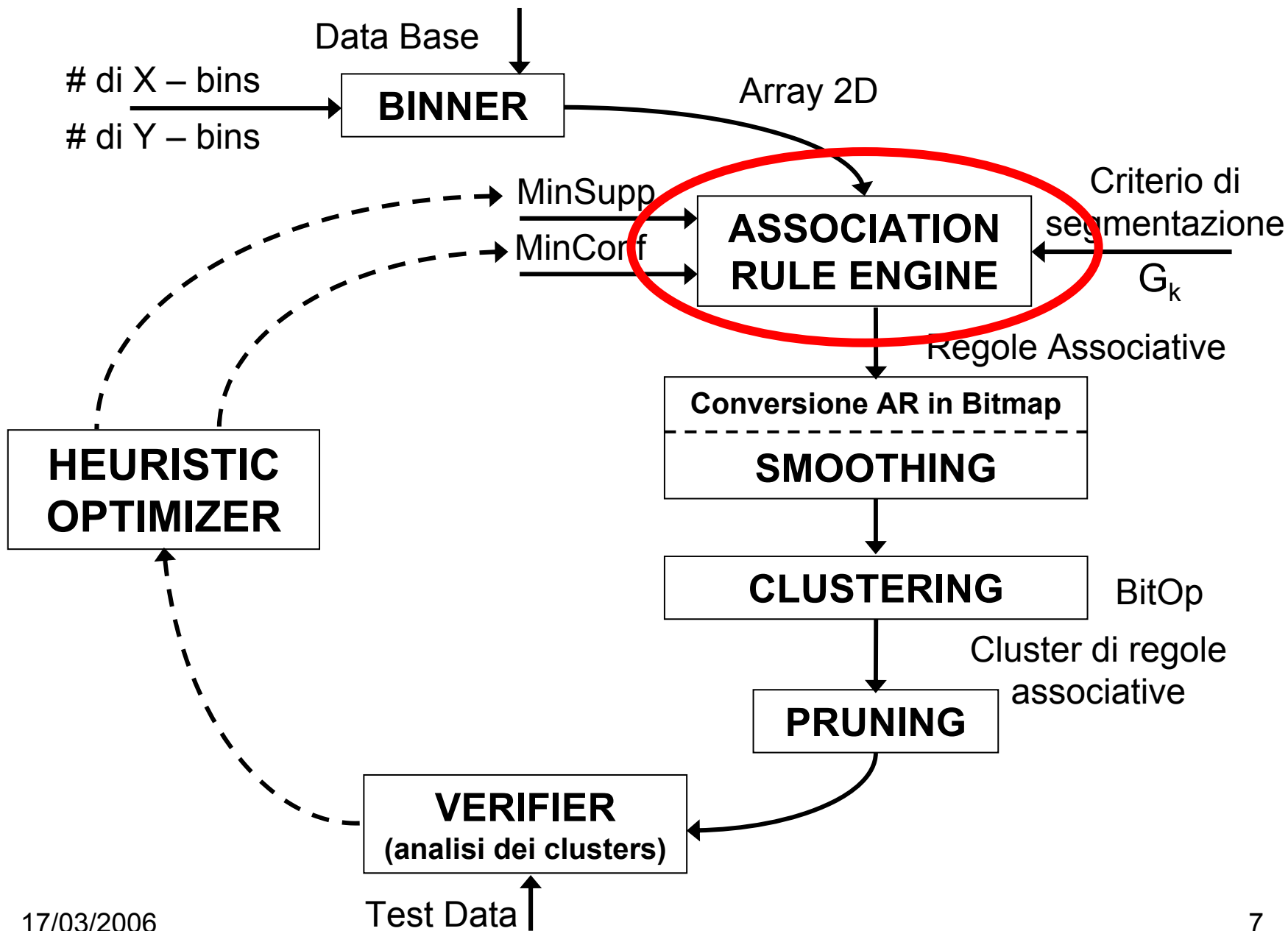
- Gli attributi X e Y sono quantitativi
→ è necessario scegliere il numero di bin con cui effettuare il partizionamento del loro dominio
- Il numero di bin è un **parametro** scelto dall'utente!!
- Partizionamento domini con criterio equi-width
- Successiva costruzione dell'**array 2D**, rappresentazione sintetica delle caratteristiche del DB

Array 2D

(bin_x, bin_y)	n°tot tuple	G_1	G_2	...	G_k
(1,1)					
(1,2)					
...					
(i,j)	500	50	10		440
...					
(n _x ,n _y)					

$$|G|=n_{seg}$$

- Ogni bin viene associato ad un numero intero progressivo
- La dimensione dell'array 2D è $n_x * n_y * (n_{seg}+1)$
- L'array 2D può essere memorizzato in memoria centrale
 - cambiando il valore di G_k per effettuare una nuova segmentazione, non è necessario rileggere i dati da DB




Association rule engine (i)

Array 2D

(bin _x ,bin _y)	n°tot tuple	G ₁	G ₂	...	G _k
(1,1)					
(1,2)					
...					
(i,j)	500	50	10		440
...					
(n _x ,n _y)					

MinSupp 25%

MinConf 80%  Supporto $| (i, j, G_k) | / N = 440/1500 \approx 30\%$

N 1500

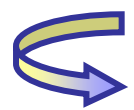
Confidenza $| (i, j, G_k) | / | (i, j) | = 440/500 \approx 88\%$

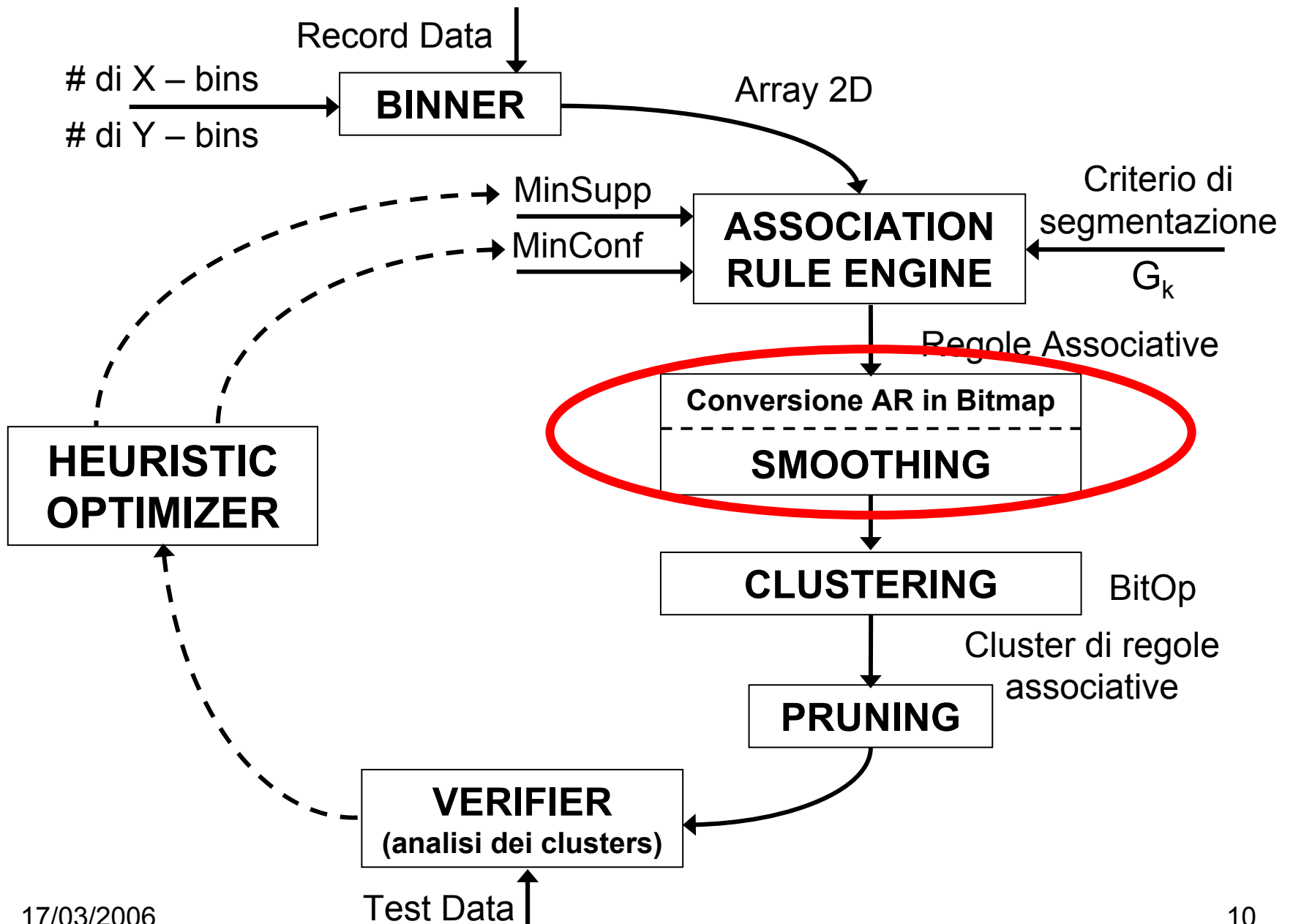
GenAssociationRules() restituisce la coppia (i,j) che identifica la regola associativa $(X = i) \wedge (Y = j) \Rightarrow G_k$

Association rule engine (ii)

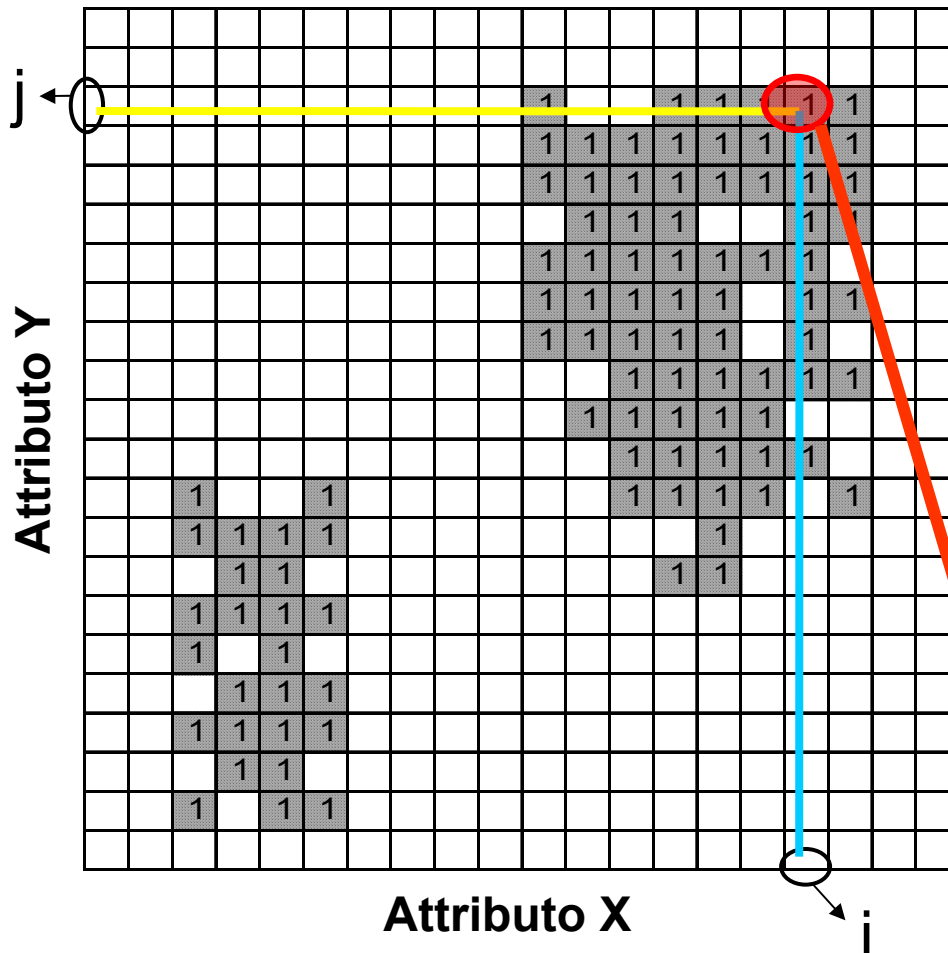
Quali **vantaggi** offre l'Association Rule Engine applicato all'array 2D?

- Tutte le regole associative per G_k in un'unica scansione
- Array 2D memorizzabile in memoria centrale

 parametri MinSupp e MinConf modificabili senza riesaminare le tuple del DB



Conversione AR in Bitmap Grid



Che cos'è una **Bitmap Grid**?

- Ad ogni coppia (i,j) restituita dalla procedura `GenAssociationRules()` viene assegnato il valore "1"
- Le altre caselle contengono valore "0"

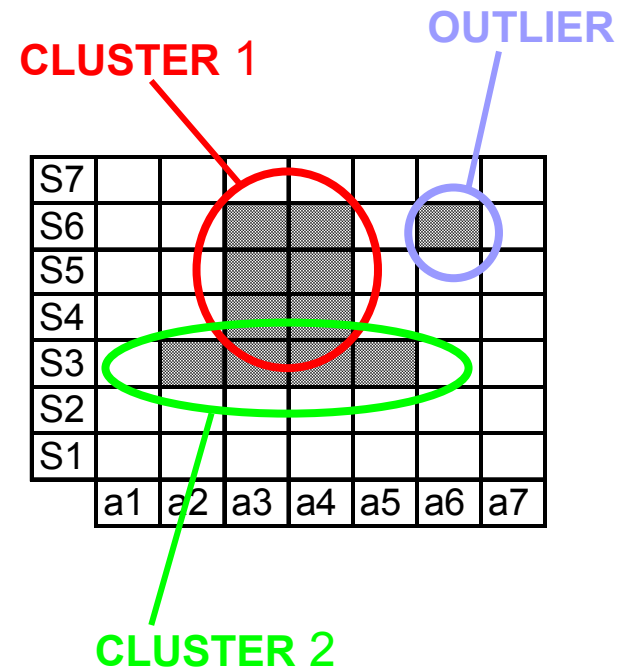
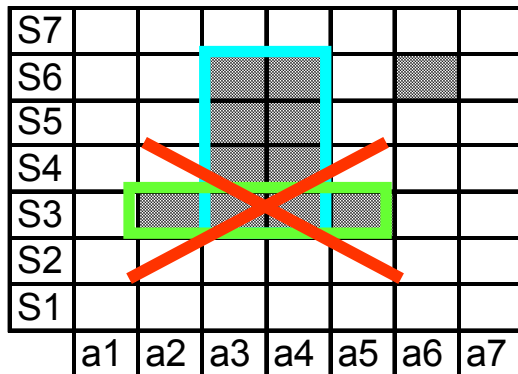
$$(X = i) \wedge (Y = j) \Rightarrow G_k$$

Come utilizzare la Bitmap Grid?

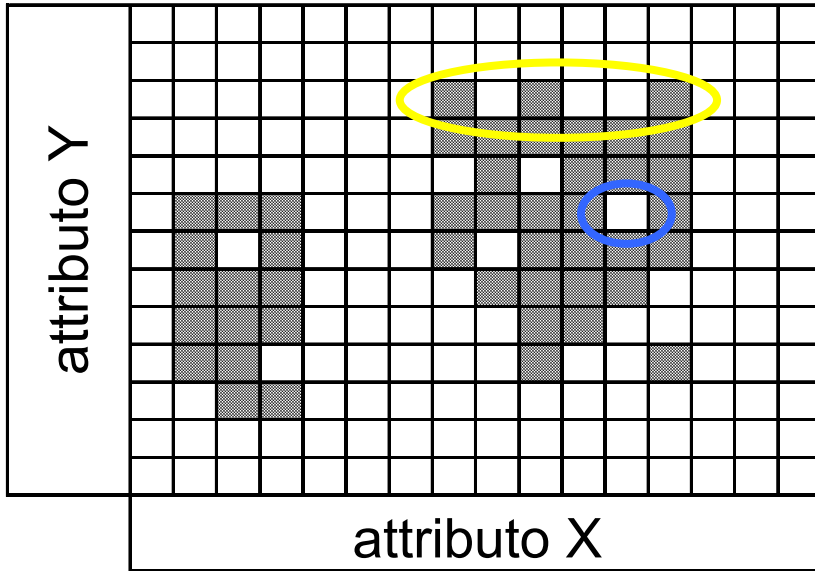
Dalla Bitmap Grid si vuole individuare il minimo numero di cluster:

- di forma rettangolare
- non sovrapposti

$$(X = i) \wedge (Y = j) \Rightarrow G_k$$

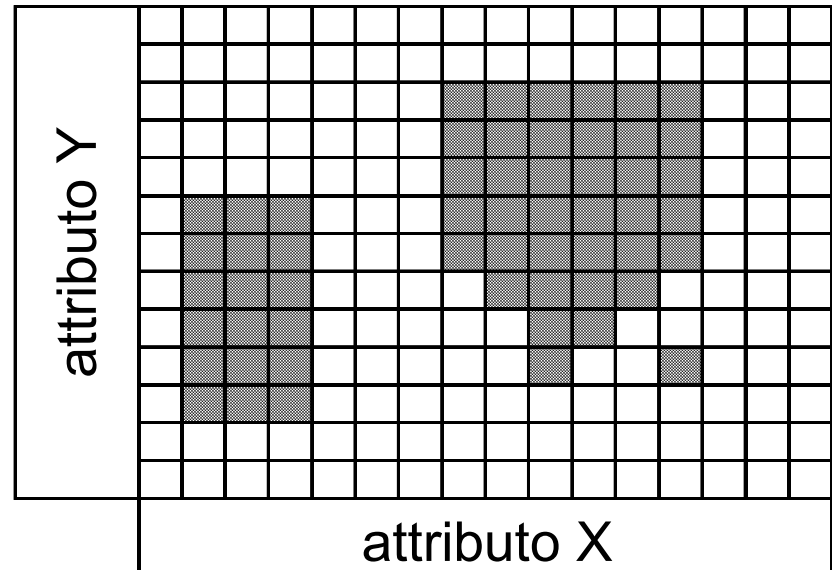
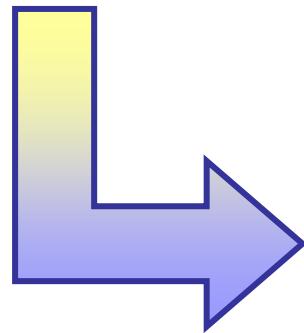


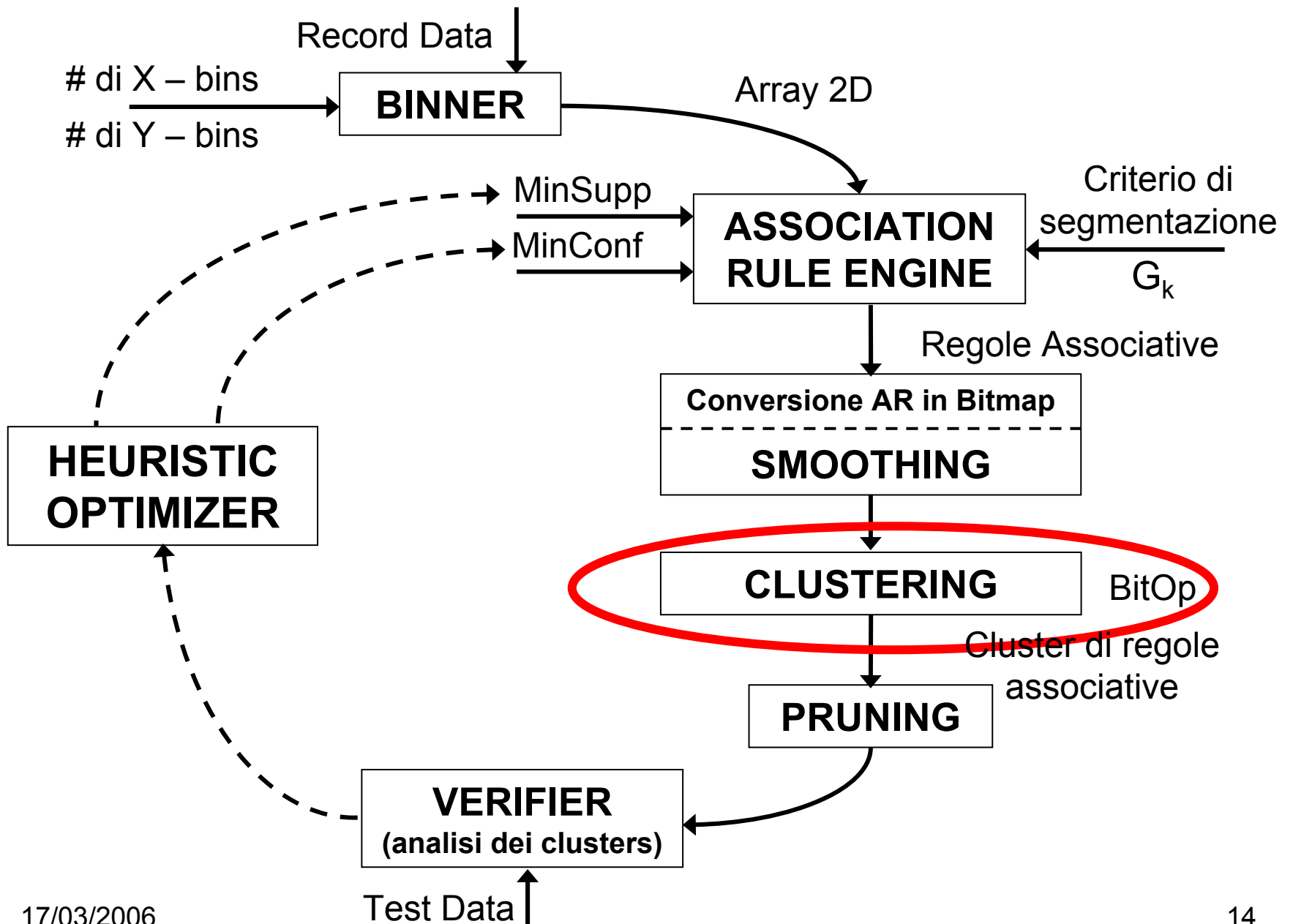
Smoothing



Tipicamente, le celle significative nella Bitmap Grid tendono a formare regioni con:

- confini frastagliati
- vuoti all'interno (rumore)

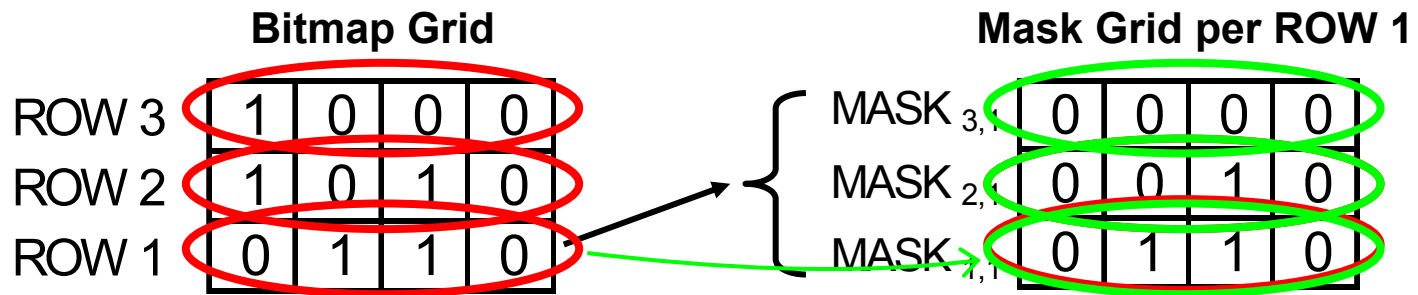




Clustering (i)

Algoritmo BitOp

Input:



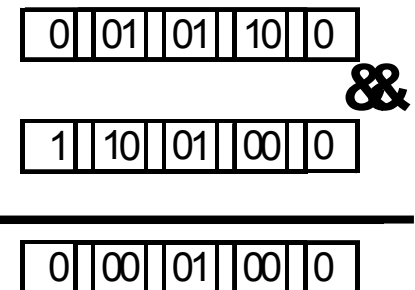
Procedura:

- Per ogni riga della Bitmap Grid si crea una Mask Grid
- Costruzione Mask Grid per riga 1:

1) $MASK_{1,1} = ROW\ 1$

2) $MASK_{2,1} = MASK_{1,1} \& ROW\ 2$

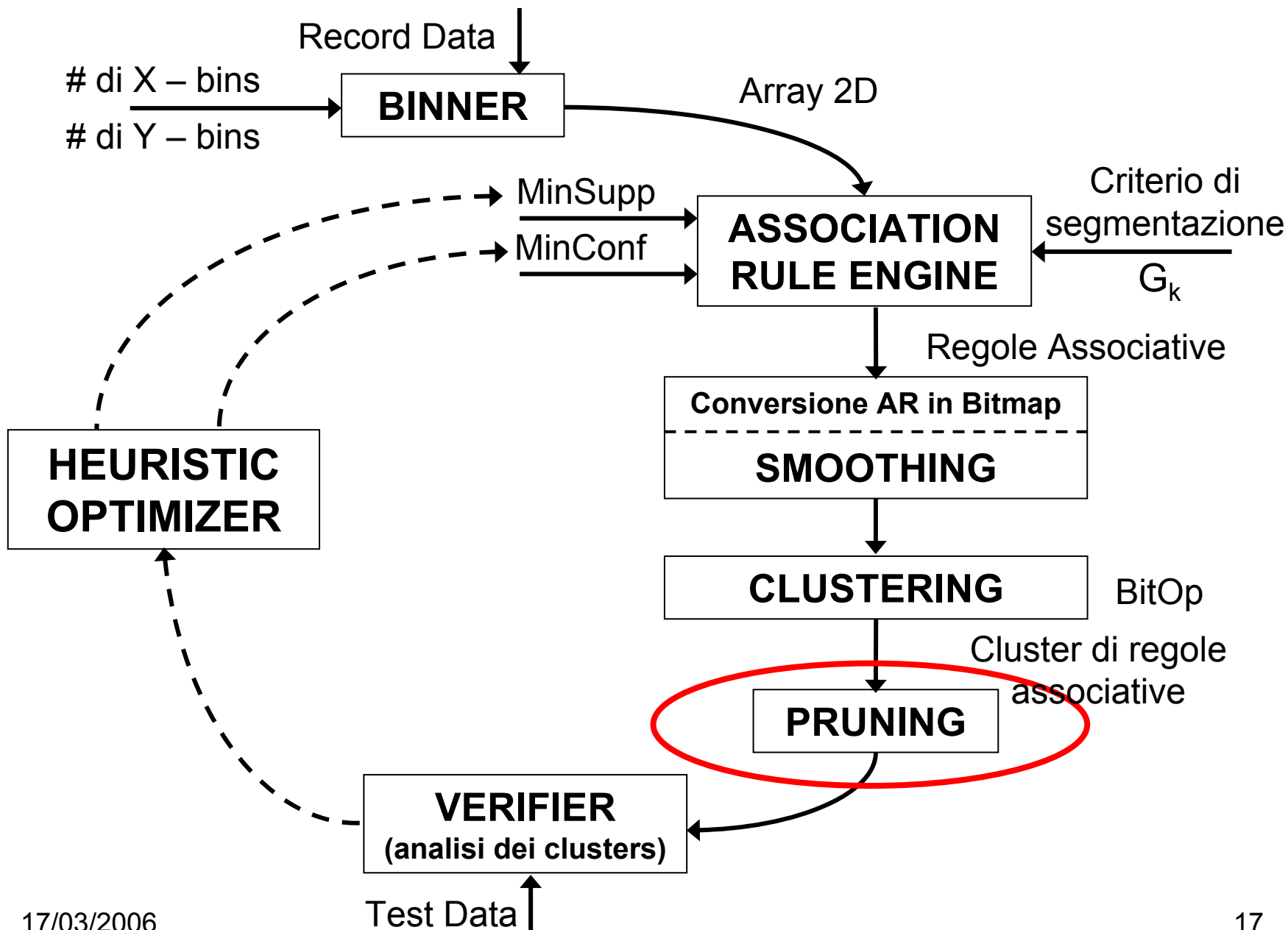
3) $MASK_{n,1} = MASK_{(n-1),1} \& ROW\ n$



Clustering (ii)

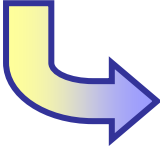
MASK _{3,1}	0	0	0	0
MASK _{2,1}	0	0	1	0
MASK _{1,1}	0	1	1	0

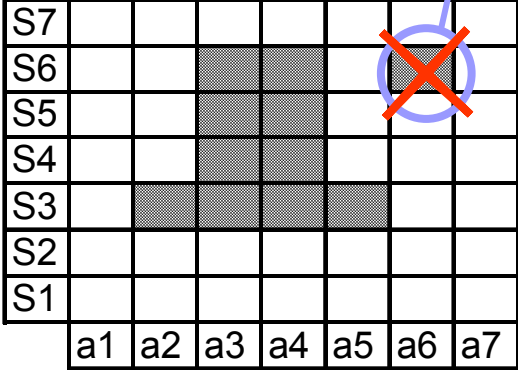
- Minimo numero di cluster di massima dimensione
- Complessità $O(\sum_{S \in C} |S|)$,
C = set finale di clusters individuati
- Approccio euristico
- Approssimazione efficiente di un clustering ottimale



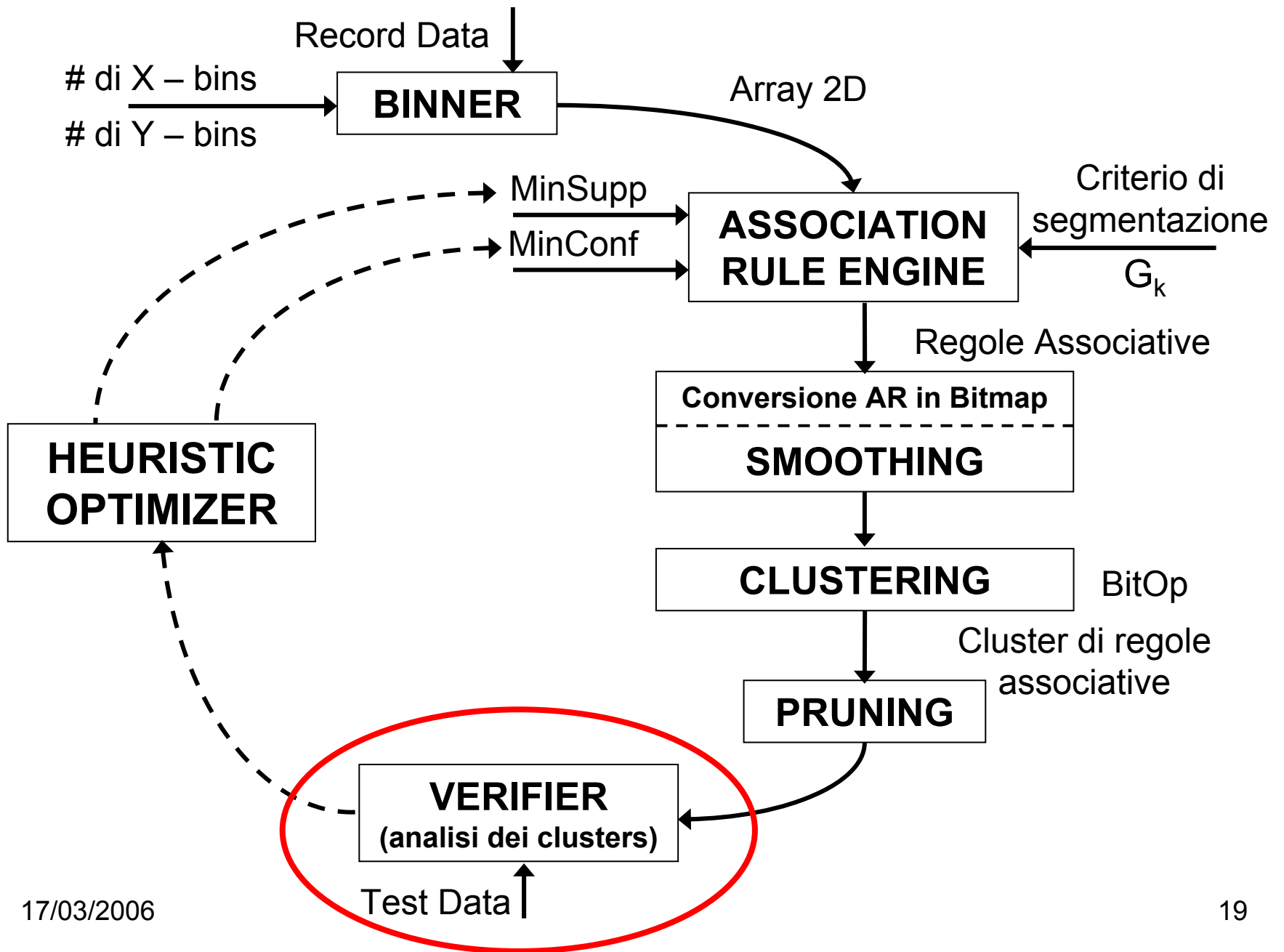
Pruning

Tipicamente i clusters di dimensioni <1% della superficie della Bitmap Grid risultano non rilevanti per una segmentazione significativa

 Eliminazione degli **outliers**



S7							
S6							
S5							
S4							
S3							
S2							
S1							
	a1	a2	a3	a4	a5	a6	a7



Verifier - qualità di segmentazione

Utilizzo del principio **MDL** (Minimum Description Length):

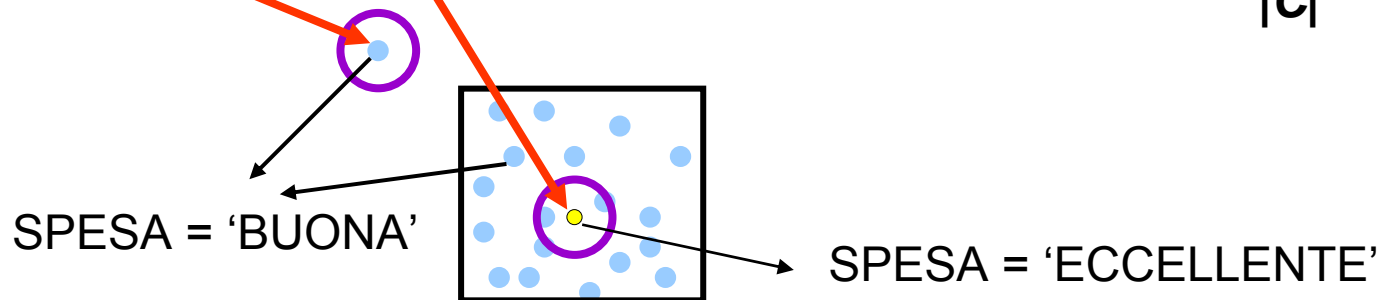
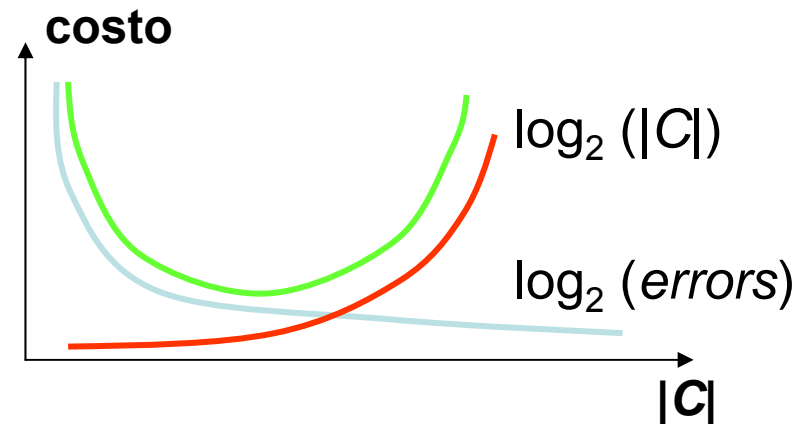
Obiettivo: min costo totale di descrizione dei clusters

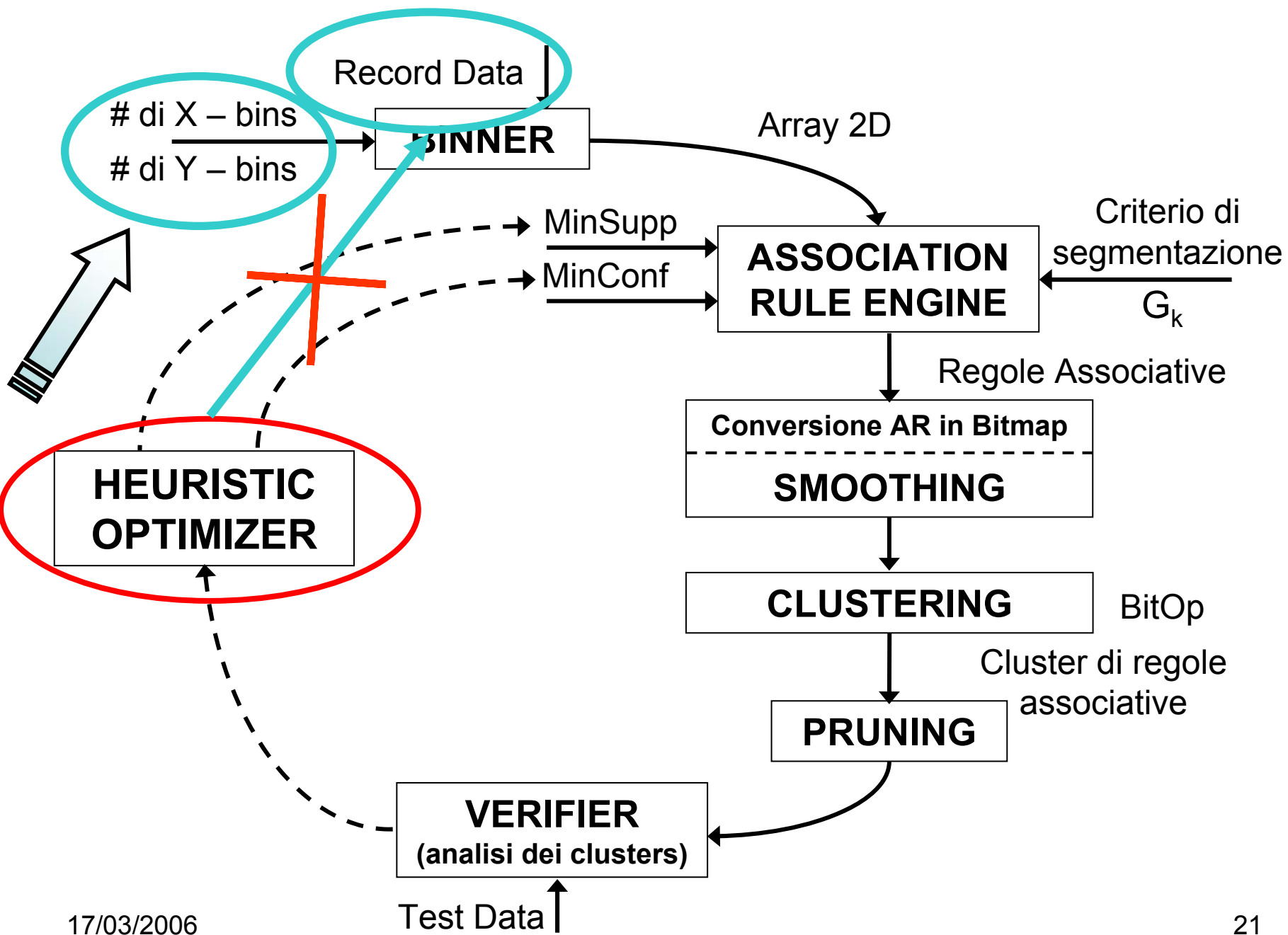
$$\text{costo} = w_c \log_2 (|\mathbf{C}|) + w_e \log_2 (\text{errors})$$

$|\mathbf{C}|$ = numero di clusters

w_c, w_e = pesi (scelti dall'utente)

errors = somma di falsi-positivi
e falsi-negativi



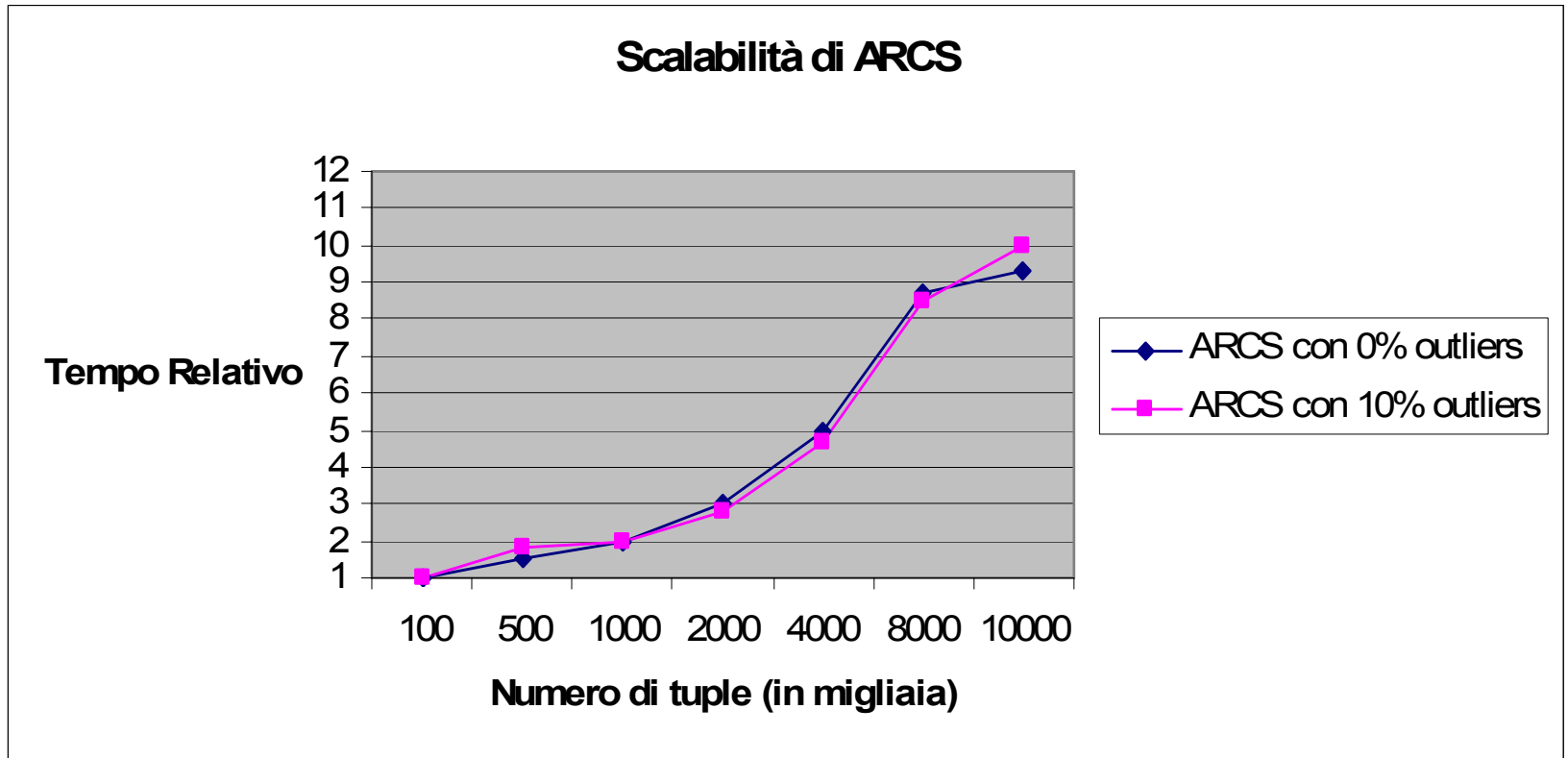


Confronto ARCS e C4.5(i)

- Comparazione dei risultati ottenuti su entrambi i datasets con ARCS e C4.5 su un Data Set sintetico con aggiunta di 5% di rumore e outlier
- C4.5 permette di costruire accurati decision trees per la classificazione dei dati

		ARCS		C4.5	
	D	$\sigma = 0\%$	$\sigma = 10\%$	$\sigma = 0\%$	$\sigma = 10\%$
	20000	27	28	14	27
	50000	37	43	81	190
	100000	42	41	210	893
	200000	45	47	650	3K
	500000	56	62	4K	20K
	1 milione	80	82	15K	86K
	2 milioni	123	117	n/a	n/a
	4 milioni	203	192	n/a	n/a
	8 milioni	367	349	n/a	n/a
	10 milioni	420	426	n/a	n/a

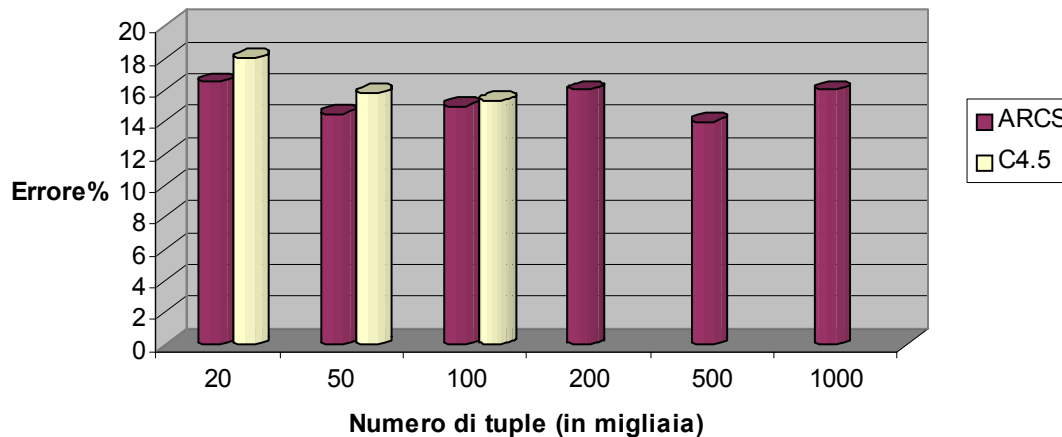
Confronto ARCS e C4.5 (iii)



Il tempo di esecuzione di ARCS aumenta al massimo linearmente con la dimensione del DB

Confronto ARCS e C4.5 (ii)

Errore con rumore al 10%

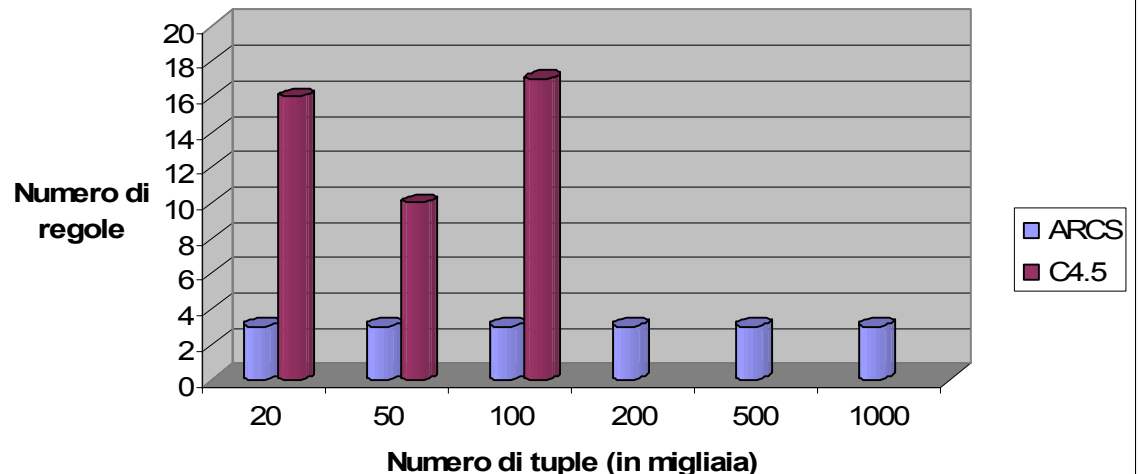


ARCS garantisce migliori prestazioni in presenza di rumore grazie a:

- smoothing
- pruning
- ricerca di MinSupp e MinConf ottimali

Riuscire ad ottenere poche regole associative clusterizzate è fondamentale per l'utente che dovrà poi utilizzarle

Numero di regole generate con rumore al 10%



Criticita'

1) Dimensione dei bin

- Ogni volta in cui il numero di bin viene cambiato è necessario rileggere tutti i dati

2) Cluster di regole associative solo



- a 2 dimensioni
- con attributi quantitativi

$$(X = i) \wedge (Y = j) \Rightarrow G_k$$

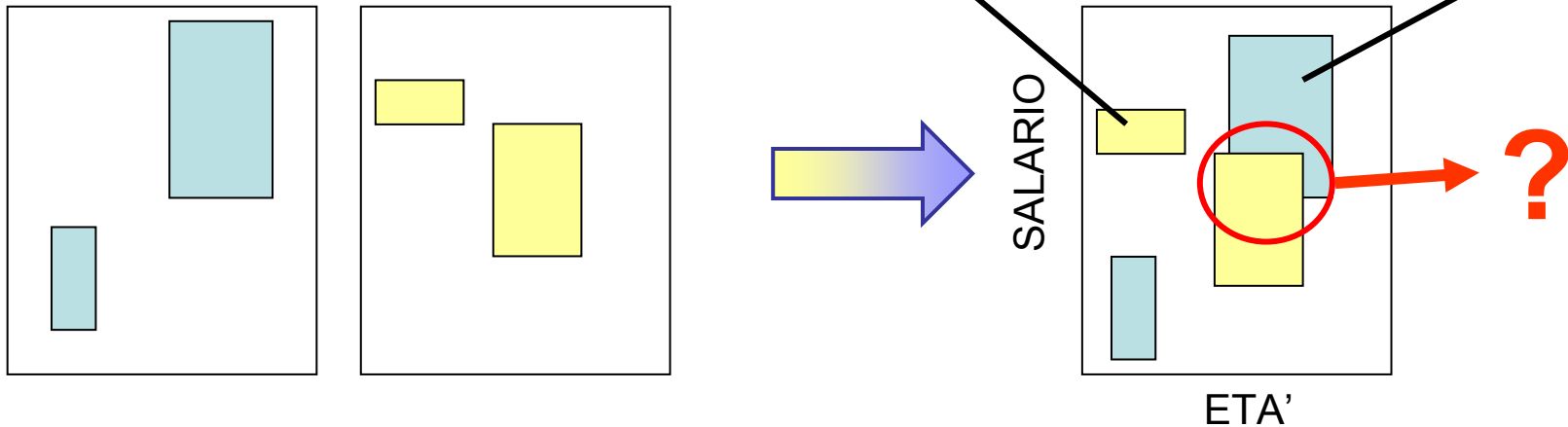
3) Overlapping dei segmenti

- Sovrapponendo griglie bitmap costruite per valori diversi dell'attributo di label, potrebbero sovrapporsi celle di "1"

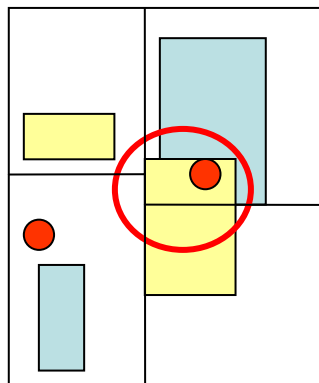
→ classificazione non univoca

Criticita' (3)

Segmenti ottenuti per due diversi valori dell'attributo di segmentazione in ARCS: Spesa = 'ECCELLENTE' Spesa = 'BUONA'



DT divide in sequenza, ortogonalmente i due attributi



Conclusioni

- Il data set su cui è stato valutato ARCS ne favorisce l'efficienza
- Un'attenta analisi del paper ha evidenziato le criticità precedentemente riportate
- Una nostra proposta di confronto fra ARCS e DT

Conclusioni

Caratteristiche del DT proposto:

- Partiziona il dominio degli attributi in bin
- Memorizza la tabella di contingenza 3D sugli attributi partizionati in bin

