

Data Mining
 II verifica intermedia – 25 maggio 2010

Esercizio 1 - Sequential Patterns (6 punti)

Si consideri la seguente sequenza di input:

\langle {A} {A,B,C} {D} {H} {B, E} {A,B,D} \rangle
 t=0 t=1 t=2 t=3 t=4 t=5

Si indichi quali sono le occorrenze delle seguenti sotto-sequenze nella sequenza di input, senza considerare vincoli temporali (colonna sinistra) e considerando il vincolo temporale $max-gap = 1$ (colonna destra). Per brevità, si rappresenti ogni occorrenza tramite la corrispondente ennupla di tempi nella sequenza di input, es.: $\langle 0,2,3 \rangle = \langle t=0, t=2, t=3 \rangle$.

	Occorrenze	Occorrenze con $max-gap=1$
es.: $\langle \{A\} \{D\} \{H\} \rangle$	$\langle 0,2,3 \rangle \langle 1,2,3 \rangle$	$\langle 1,2,3 \rangle$
$w_1 = \langle \{A\} \{B\} \{D\} \rangle$		
$w_2 = \langle \{A\} \{H\} \{B\} \rangle$		
$w_2 = \langle \{A\} \{C\} \{E\} \rangle$		

Esercizio 2 – Itemset Frequenti (12 punti)

Considerare la seguente tabella di transazioni:

ID	ITEMS	ID	ITEMS
1	A B C	6	B C E
2	B	7	C D E
3	A B C D	8	A B
4	C	9	A B C D
5	A D	10	A B D

- A) Elencare gli itemset frequenti nel caso di un supporto minimo $\text{min_sup} = 25\%$ ed indicare il loro supporto.
- B) Quali itemset frequenti sono anche closed? Quali sono massimali?

Esercizio 3 – FP-tree (2 punti)

Si disegni l'FP-tree corrispondente al seguente dataset:

ID	Itemset
1	A
2	B D E
3	A C
4	B D
5	B C

Esercizio 4 - Clustering (12 punti)

Sul seguente dataset:

- A) Si utilizzi l'algoritmo di clustering density-based DBSCAN, con raggio (ϵ) pari a 1.9, e minPts pari a 4 (=3 vicini + il punto di cui si calcola la densità). Si richiede di (1) indicare il numero di cluster che si ottengono; (2) per ogni punto indicare il cluster di appartenenza; (3) per ogni punto dire se si tratta di un *core point*, *border point* o *rumore*. (8 punti)

- B) Si disegni il dendrogramma ottenuto con un algoritmo di clustering agglomerativo MIN-link (o *Single linkage*). (4 punti)

