

Data Mining

Appello del 1 giugno 2010

Soluzioni

Esercizio 1 - Sequential Patterns (4 punti)

Si consideri la seguente sequenza di input:

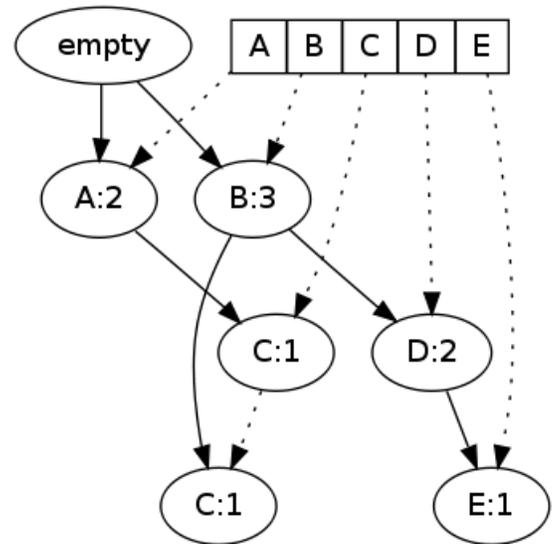
$$\begin{array}{cccccc} < & \{A\} & \{A,B,C\} & \{C,D,E\} & \{A,E,H\} & \{B\} & \{A,B,D\} & > \\ & t=0 & t=1 & t=2 & t=3 & t=4 & t=5 & \end{array}$$

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.: $\langle \{A\} \{D\} \{H\} \rangle$	$\langle 0,2,3 \rangle \langle 1,2,3 \rangle$	$\langle 1,2,3 \rangle$
$w_1 = \langle \{B\} \{E\} \{D\} \rangle$	$\langle 1,2,5 \rangle \langle 1,3,5 \rangle$	nessuna
$w_2 = \langle \{A\} \{A,B\} \rangle$	$\langle 0,1 \rangle \langle 0,5 \rangle \langle 1,5 \rangle \langle 3,5 \rangle$	$\langle 0,1 \rangle$
$w_2 = \langle \{A\} \{C\} \{E\} \rangle$	$\langle 0,1,2 \rangle \langle 0,1,3 \rangle \langle 0,2,3 \rangle \langle 1,2,3 \rangle$	$\langle 0,1,2 \rangle \langle 1,2,3 \rangle$

Esercizio 2 – FP-tree (2 punti)

Si ricostruisca il dataset di transazioni (itemset) da cui il seguente FP-tree è stato ottenuto.



Soluzione:

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

Esercizio 3 – Itemset Frequenti (6 punti)

Considerare la seguente tabella di transazioni:

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

ID	ITEMS
6	B C
7	C D E
8	A E
9	A B D E
10	A C D

- A) Elencare gli itemset frequenti nel caso di $\text{min_sup} = 20\%$ ed indicare il loro supporto.
- B) Quali itemset frequenti sono anche massimali?

(m=massimale)

- A (60.0)
- m B (30.0)
- C (60.0)
- D (60.0)
- E (40.0)
- AC (30.0)
- AD (40.0)
- m AE (30.0)
- CD (40.0)
- m CE (20.0)
- m DE (20.0)
- m ACD (20.0)

Esercizio 4 - Clustering (10 punti)

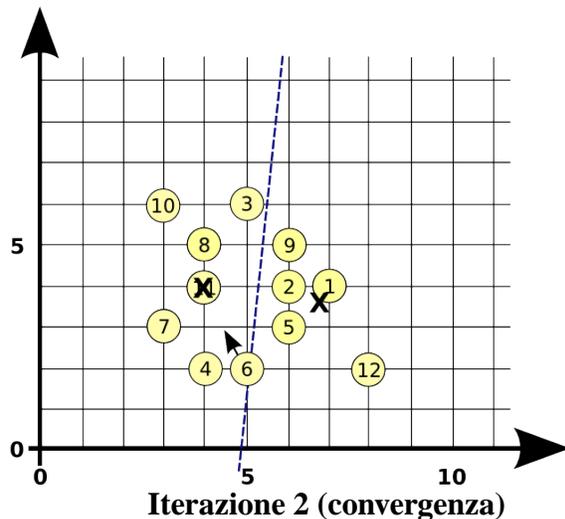
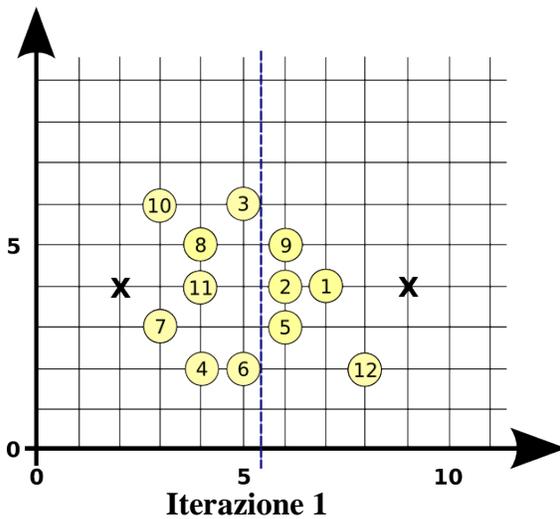
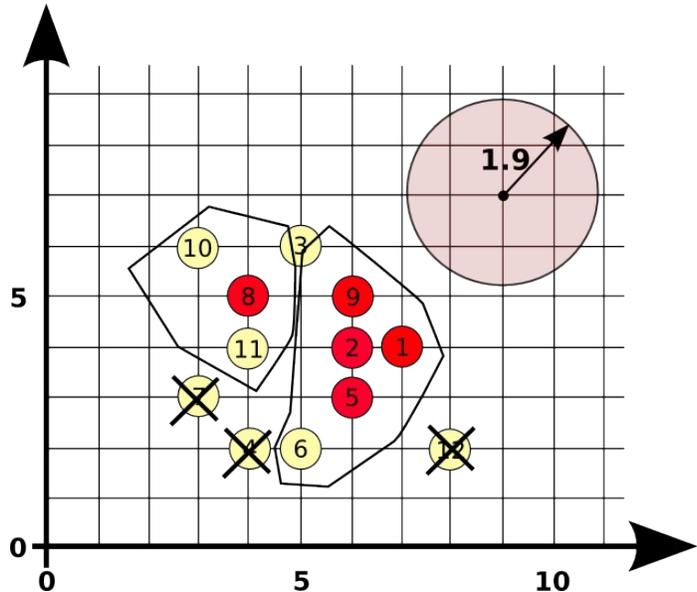
Sul seguente dataset:

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).

(1) per ogni punto dire se si tratta di un *core point*, *border point* o *rumore*;

(2) indicare la composizione dei cluster ottenuti. (5 punti)

B) Simulare l'esecuzione dell'algoritmo k-means sullo stesso insieme di punti, con $k=2$ e centri iniziali $c_1=(2,4)$ e $c_2=(9,4)$. (5 punti)



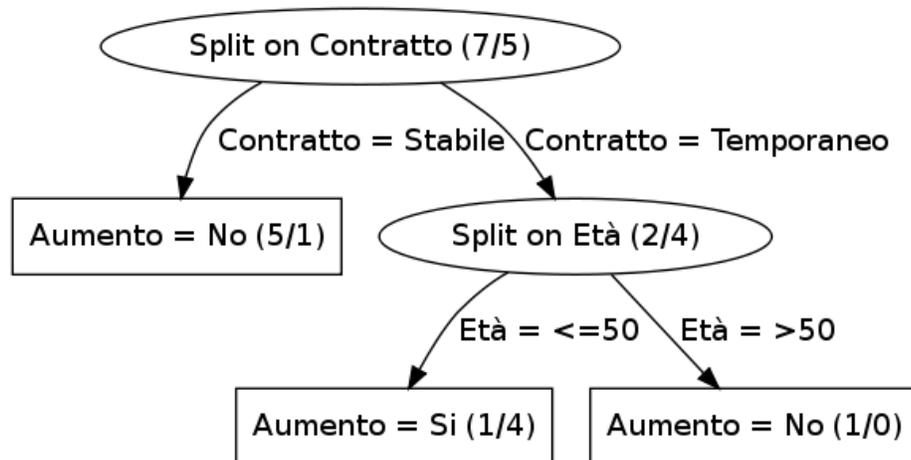
Esercizio 4 – Classificazione (10 punti)

Si consideri il seguente insieme di transazioni (*training set*).

Training set:

Età	Contratto	Sesso	Aumento
50	Temporaneo	F	Si
35	Stabile	M	No
50	Temporaneo	F	Si
40	Stabile	M	No
25	Temporaneo	F	Si
30	Temporaneo	F	No
50	Stabile	F	No
40	Stabile	F	Si
55	Temporaneo	M	No
55	Stabile	F	No
25	Temporaneo	M	Si
40	Stabile	M	No

A) Si costruisca su tale dataset un albero di decisione per la variabile “Aumento”, utilizzando il criterio di split basato su “misclassification rate”, espandendo i nodi dell'albero fino a che la precisione non è più migliorabile localmente (ovvero nessuno split da' un guadagno). (7 punti)



B) Si mostrino accuratezza e matrice di confusione dell'albero ottenuto al punto A), calcolati sia sul training set che sul test set riportato qui sotto. Confrontare i risultati. **(3 punti)**

	Età	Contratto	Sesso	Aumento
Test set:	35	Stabile	F	No
	45	Temporaneo	M	No
	30	Stabile	M	No
	25	Stabile	F	Si
	40	Temporaneo	M	No
	55	Stabile	F	Si
	30	Temporaneo	F	Si
	35	Temporaneo	M	Si

Matrici di confusione e accuratezza:

Train:

	Si	No
Si	4	1
No	1	6

(Reale)

(Predetta)

Accuracy:

83,33%

Test:

	Si	No
Si	2	2
No	2	2

(Reale)

(Predetta)

Accuracy:

50,00%