

## Data Mining

Appello del 9 settembre 2010

## Soluzioni

## Esercizio 1 - Sequential Patterns (4 punti)

Si consideri la seguente sequenza di input:

$$\begin{array}{cccccccc} < & \{A,C\} & \{C,D\} & \{F,H\} & \{A,B\} & \{B,C,D\} & \{E\} & \{A,B,D\} & \{F\} & > \\ & t=0 & t=1 & t=2 & t=3 & t=4 & t=5 & t=6 & t=7 & \end{array}$$

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  $min-gap = 1$  (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 $min-gap=1$
es.: $\langle \{C\} \{H\} \{C\} \rangle$	$\langle 0,2,4 \rangle \langle 1,2,4 \rangle$	$\langle 0,2,4 \rangle$
$w_1 = \langle \{A\} \{F\} \rangle$	$\langle 0,2 \rangle \langle 0,7 \rangle$ $\langle 3,7 \rangle \langle 6,7 \rangle$	$\langle 0,2 \rangle \langle 0,7 \rangle$ $\langle 3,7 \rangle$
$w_2 = \langle \{A\} \{A\} \{D\} \rangle$	$\langle 0,3,4 \rangle \langle 0,3,6 \rangle$	$\langle 0,3,6 \rangle$
$w_2 = \langle \{A\} \{A,B\} \{F\} \rangle$	$\langle 0,3,7 \rangle \langle 0,6,7 \rangle \langle 3,6,7 \rangle$	$\langle 0,3,7 \rangle$

## Esercizio 2 – Lift e confidenza di regole associative (2 punti)

---

La regola associativa “a → b” ha supporto pari a 0.6, confidenza = 0.8 e lift = 2. Si calcolino i valori del supporto di “a” e di “b”.

### Soluzione:

$$\text{confidenza} = s(a,b)/s(a) \quad \Rightarrow s(a) = s(a,b) / \text{confidenza} = 0.6 / 0.8 = 0.75$$

$$\text{lift} = s(a,b) / [s(a)s(b)] \quad \Rightarrow s(b) = s(a,b) / [s(a) \text{ lift}] = 0.6 / [0.75 \times 2] = 0.4$$

Si può notare, però, che il valore ottenuto per  $s(b) < s(a,b)$ , il che mostra che i valori di partenza del problema sono incoerenti, cioè non esiste alcun dataset D da cui essi possano emergere.

## Esercizio 3 – Itemset Frequenti (6 punti)

---

Considerare la seguente tabella di transazioni:

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

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

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

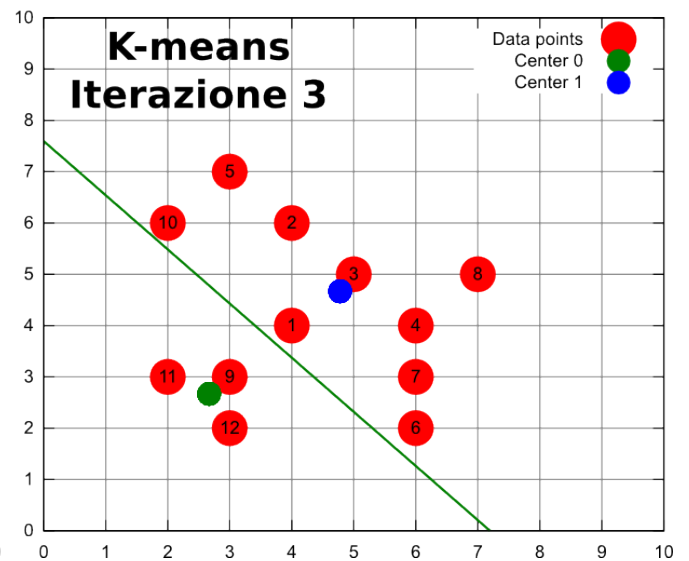
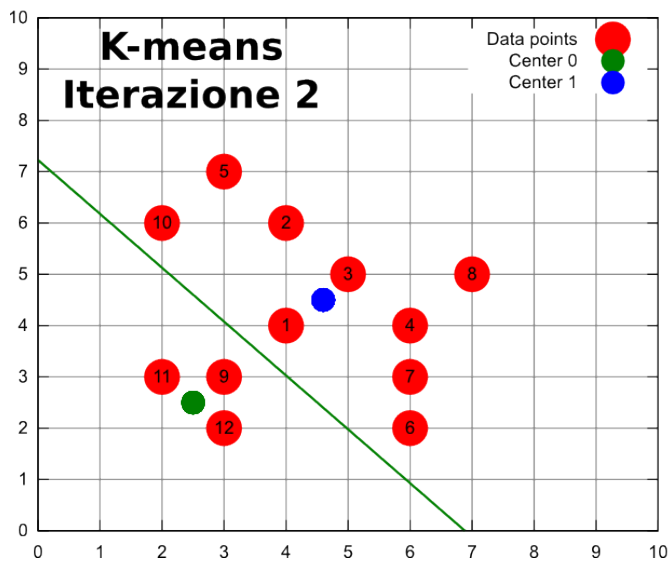
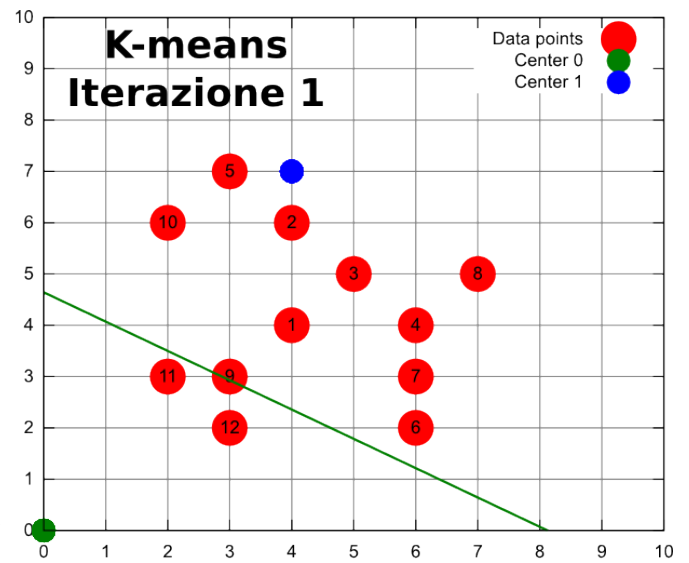
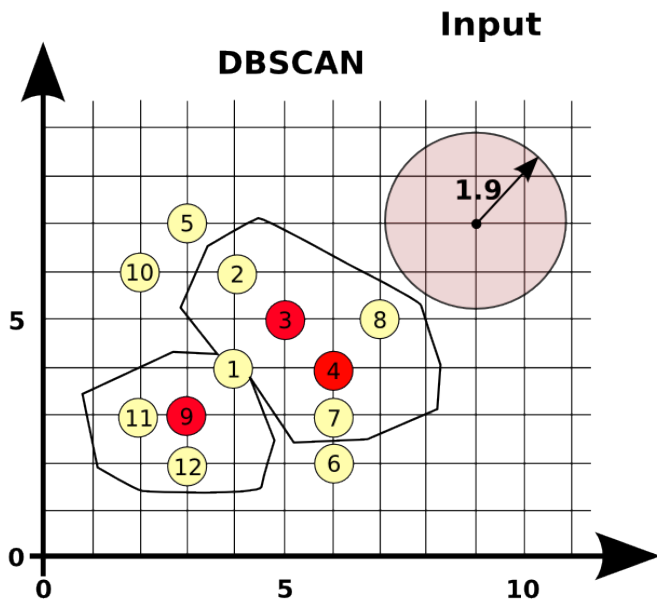
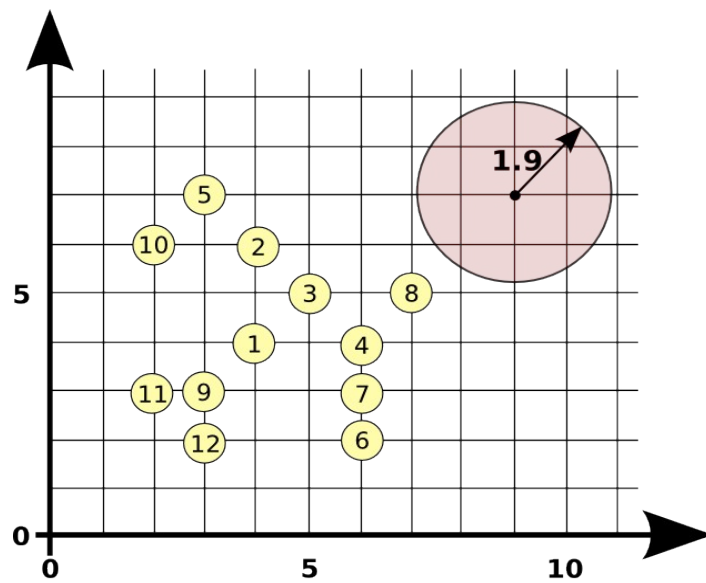
(m=massimale)
A (60.0)
B (70.0)
C (90.0)
D (40.0)
AB (30.0)
AC (60.0)
AD (20.0)
BC (60.0)
BD (30.0)
CD (40.0)
m ABC (30.0)
m ACD (20.0)
m BCD (30.0)

## Esercizio 4 - Clustering (10 punti)

---

Sul seguente dataset:

- A) Si utilizzi l'algoritmo di clustering density-based DBSCAN, con raggio ( $\epsilon$ ) pari a 1.9, e  $\text{minPts}$  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  $c_1=(0,0)$  e  $c_2=(4,7)$ . (5 punti)

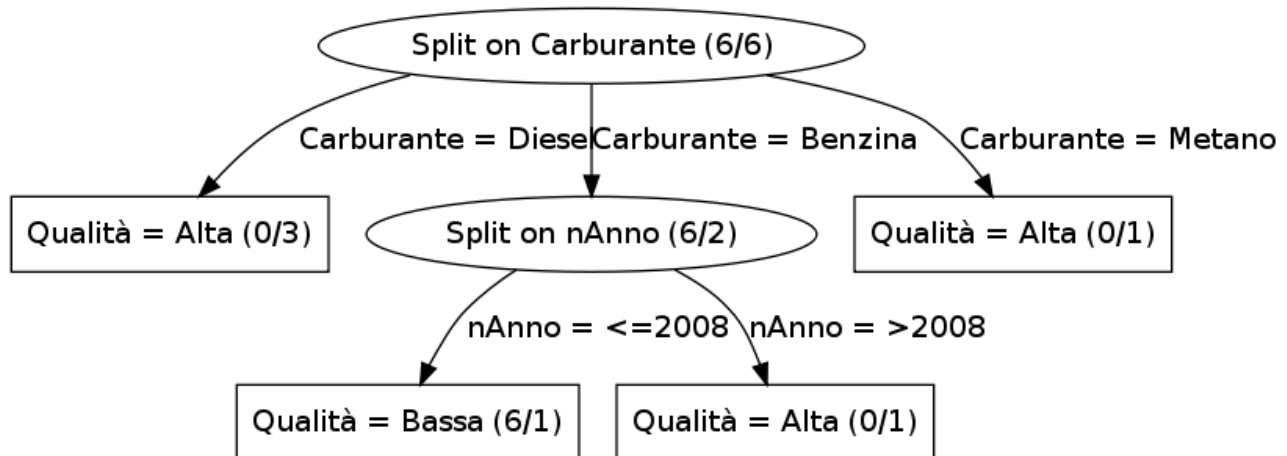


Esercizio 5 – Classificazione (10 punti)

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

nAnno	Usato	Carburante	Qualità
2004	Si	Benzina	Bassa
2008	No	Metano	Alta
2003	No	Diesel	Alta
2006	Si	Benzina	Alta
2001	No	Benzina	Bassa
2001	Si	Benzina	Bassa
2008	No	Benzina	Bassa
2006	Si	Diesel	Alta
2001	No	Diesel	Alta
2002	Si	Benzina	Bassa
2007	Si	Benzina	Bassa
2009	Si	Benzina	Alta

A) Si costruisca su tale dataset un albero di decisione per la variabile “Qualità”, 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)

	<b>nAnno</b>	<b>Usato</b>	<b>Carburante</b>	<b>Qualità</b>
Test set:	2001	No	Metano	Bassa
	2008	No	Diesel	Alta
	2010	No	Benzina	Bassa
	2006	Si	Diesel	Bassa
	2005	No	Metano	Alta
	2010	No	Benzina	Bassa
	2007	Si	Diesel	Bassa
	2001	No	Diesel	Bassa

Matrici di confusione e accuratezza:

Train:

	Alte	Basse	(Predetta)	Accuracy:	91,67%
Alte	5	1			
Basse	0	6			
(Reale)					

Test:

	Alta	Bassa	(Predetta)	Accuracy:	25,00%
Alta	0	0			
Bassa	6	2			
(Reale)					