



Clustering

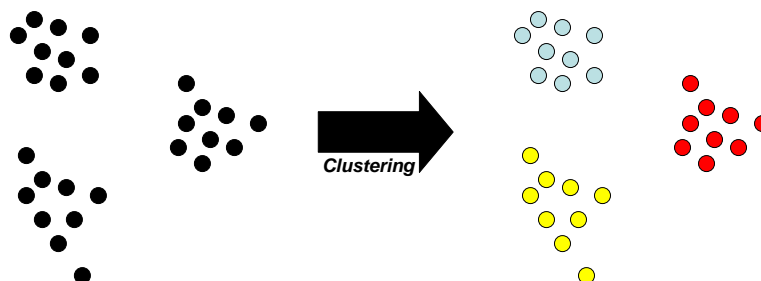
Iuri Frosio
frosio@dsi.unimi.it

Approfondimenti in A.K. Jan, M. N. Murty, P. J. Flynn, "Data clustering: a review", ACM Computing Surveys, Vol. 31, No. 3, September 1999, ref. pp. 265-290, disponibile in <http://citeseer.ist.psu.edu/jain99data.html>



Clustering

- Intuitivamente...





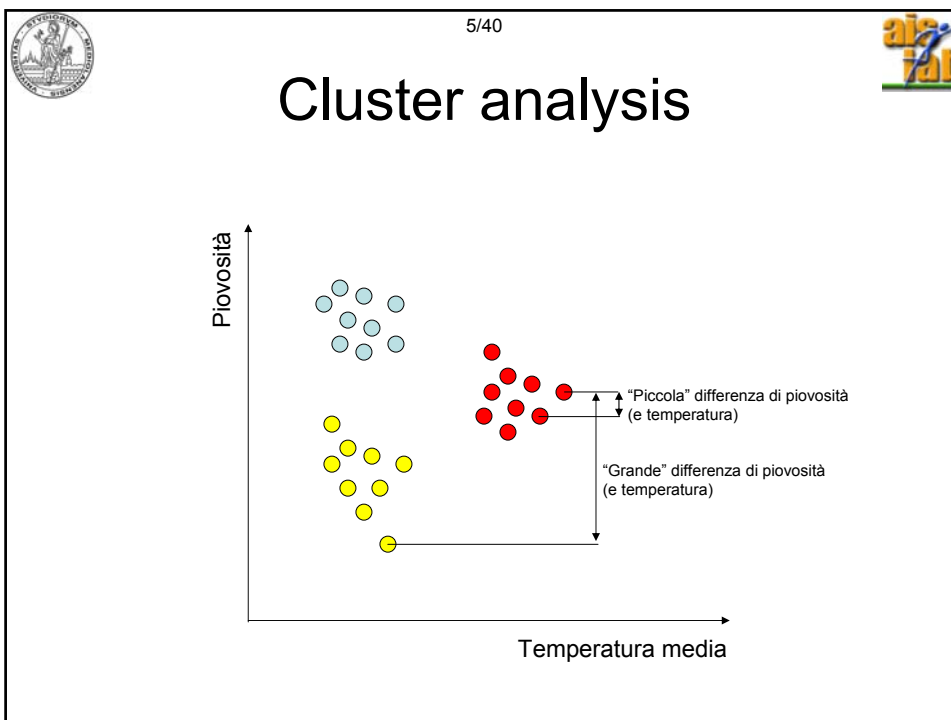
Il clustering per...

- ... Confermare ipotesi sui dati (es. “E’ possibile identificare tre diversi tipi di clima in Italia: mediterraneo, continentale, alpino...”);
- ... Esplorare lo spazio dei dati (es. “Quanti tipi diversi di clima sono presenti in Italia?”);
- ... Semplificare l’interpretazione dei dati (“Il clima di ogni città d’Italia è approssimativamente mediterraneo, continentale o alpino.”).

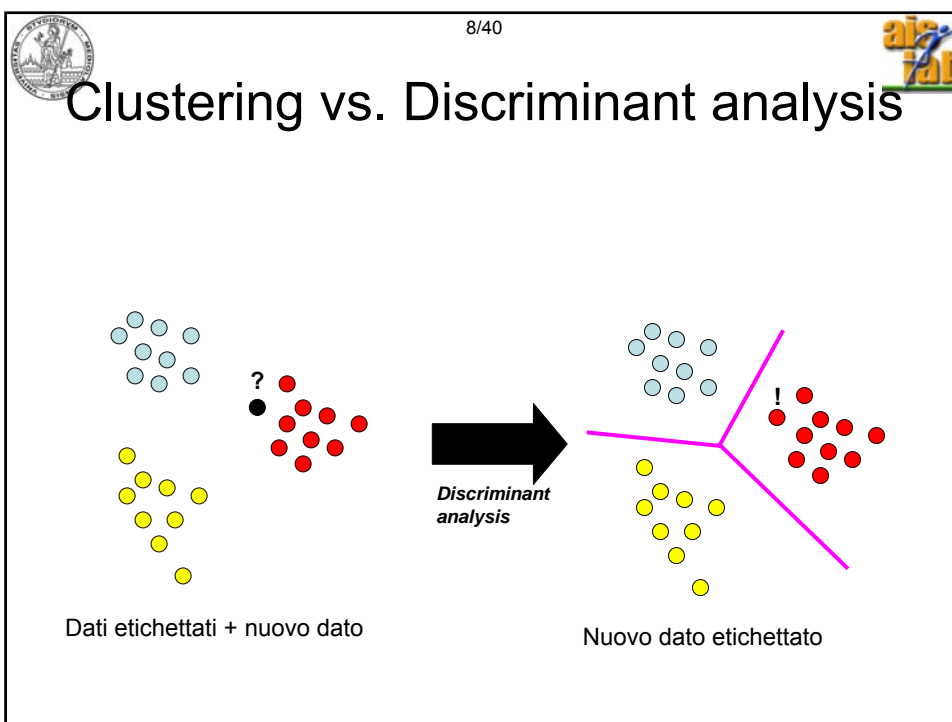
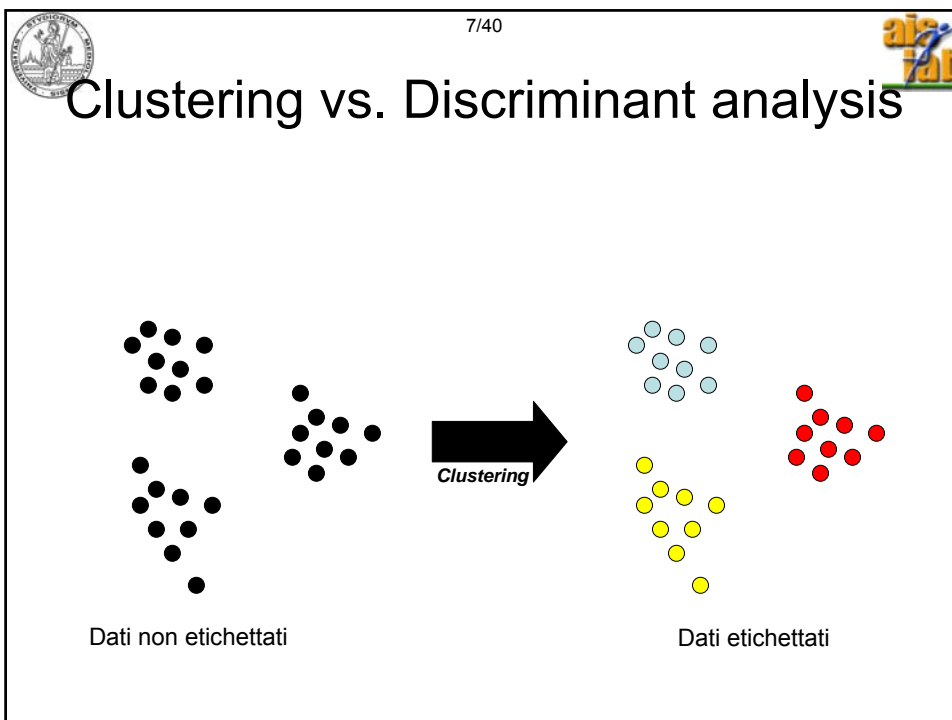


Cluster analysis

- Cluster analysis: organizzazione di una collezione di dati (pattern, vettori) in cluster, basata sulla similarità.
- I pattern appartenenti ad un cluster valido sono più simili l’uno con l’altro rispetto ai pattern appartenenti ad un cluster differente.



- 6/40
-
- Clustering vs. Discriminant analysis**
- Si definisce clustering un procedimento **non supervisionato** di **classificazione di pattern non precedentemente etichettati** in cluster.
 - Nel caso in cui i pattern siano stati precedentemente etichettati, si parla di discriminant analysis (determinazione della classe di appartenenza di un nuovo pattern, data in partenza una serie di pattern già etichettati).





Passi del clustering

- 1) Rappresentazione dei pattern;
- 2) Definizione di una misura di prossimità dei pattern;
- 3) Clustering propriamente detto;
- 4) Data abstraction (se necessario);
- 5) Validazione dell'output (se necessario).



Rappresentazione dei pattern

- Numero di classi, numero di pattern disponibili;
- Numero, tipo e scala delle **feature** (caratteristiche) utilizzate dall'algoritmo di clustering.

Roma : [17°; 500mm] Pattern

Milano : [13°; 900mm]

Feature 1 Feature 2



Rappresentazione dei pattern

- Feature selection: identificazione delle feature più significative per la descrizione dei pattern.

esempio:

Roma: [17°; 500mm; ~~1.500.000 ab.~~]

- Feature extraction: trasformazione delle feature per creare nuove, significative feature;

esempio:

Milano: [13°; 900mm; 255 giorni sole; 100 giorni pioggia]

oppure

Milano: [13°; 900mm / 100 giorni pioggia; 255 giorni sole]



Prossimità dei pattern

- Definizione di una **misura di distanza tra due patterns**;

esempio:

Distanza euclidea...

dist (Roma, Milano) = ...

dist ([17°; 500mm], [13°; 900mm]) = ...

= ... Distanza euclidea? = ...

= $((17-13)^2 + (500-900)^2)^{1/2} = 400.02 \sim 400$



Normalizzazione feature

Att.ne!

$$\begin{aligned} \text{dist}(\text{Roma}, \text{Milano}) &= \dots \\ \text{dist}([17^\circ; 500\text{mm}], [13^\circ; 900\text{mm}]) &= \dots \\ &= \dots \text{ Distanza euclidea?} = \dots \\ &= ((17-13)^2 + (500-900)^2)^{1/2} = 400.02 \sim 400 \end{aligned}$$

La distanza tra le due città in termini di gradi è insignificante nel nostro calcolo... **E' necessario normalizzare i dati!**

Es.

$$\begin{aligned} T_{\text{Max}} &= 20^\circ \quad T_{\text{Min}} = 5^\circ \rightarrow T_{\text{Norm}} = (T - T_{\text{Min}}) / (T_{\text{Max}} - T_{\text{Min}}) \\ P_{\text{Max}} &= 1000\text{mm} \quad P_{\text{Min}} = 0\text{mm} \rightarrow P_{\text{Norm}} = (P - P_{\text{Min}}) / (P_{\text{Max}} - P_{\text{Min}}) \\ \text{Roma}_{\text{Norm}} &= [0.8 \quad 0.5] \\ \text{Milano}_{\text{Norm}} &= [0.53 \quad 0.9] \\ \text{dist}(\text{Roma}_{\text{Norm}}, \text{Milano}_{\text{Norm}}) &= ((0.8-0.53)^2 + (0.5-0.9)^2)^{1/2} = 0.4826 \end{aligned}$$



Altre funzioni di distanza

- Distanza euclidea:
 $\text{dist}(x,y) = [\sum_{k=1..d} (x_k - y_k)^2]^{1/2}$
- Minkowski:
 $\text{dist}(x,y) = [\sum_{k=1..d} (x_k - y_k)^p]^{1/p}$
- Mahalanobis:
 $\text{dist}(x,y) = (x_k - y_k) S^{-1} (x_k - y_k)$
- Context dependent:
 $\text{dist}(x,y) = f(x, y, \text{context})^*$

*Approfondimenti in A.K. Jan, M. N. Murty, P. J. Flynn, "Data clustering: a review", ACM Computing Surveys, Vol. 31, No. 3, September 1999, ref. pp. 272-273.



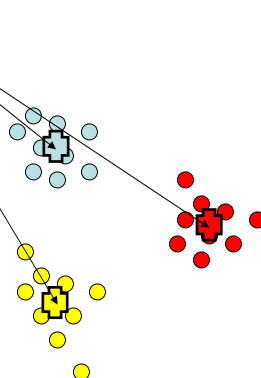
Clustering

- Hard vs. fuzzy clustering;
- Hierarchical vs. partitional clustering;
- Agglomerative vs. divisive;
- Deterministic vs. stochastic...



Data abstraction

- **Prototipi** per la descrizione di ciascuna classe.





Validazione dell'output

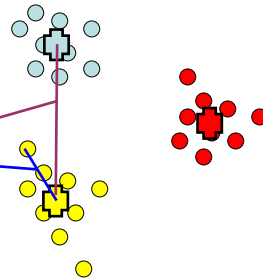
- Valutazione qualitativa (es. operatore umano) della classificazione;
- Valutazione quantitativa della classificazione;

es. minimo di:

$$k_1 \cdot \sum_j \sum_i [\text{dist}(x_i^j, \text{prot}_j)] / N_j$$

$$k_2 \cdot \sum_j \sum_h [\text{dist}(\text{prot}_j, \text{prot}_h)]$$

La valutazione quantitativa è più semplice se ci si basa su metodi statistici.



Clustering: definizioni

- **Pattern:** un singolo dato $x = [x_1, x_2, \dots, x_d]$;
- **Feature:** ogni componente x_1, x_2, \dots, x_d ;
- **d:** dimensione dello spazio delle feature;
- **Classe:** in generale, un processo che governa la generazione dei pattern di un determinato cluster (in particolare, la ddp dei dati di un cluster);
- **Hard vs. fuzzy clustering:** ogni pattern è etichettato con una sola label nel primo caso o da un grado di appartenenza ad un cluster nel secondo caso;
- **Funzione di distanza:** una metrica (o quasi metrica) nello spazio delle feature, usata per quantificare la similarità tra due pattern.



Altre funzioni di distanza

- Distanza euclidea:
 $\text{dist}(x,y)=[\sum_{k=1..d}(x_k-y_k)^2]^{1/2}$
- Minkowski:
 $\text{dist}(x,y)=[\sum_{k=1..d}(x_k-y_k)^p]^{1/p}$
- Mahalanobis:
 $\text{dist}(x,y)= (x_k-y_k)S^{-1}(x_k-y_k)$
- Context dependent:
 $\text{dist}(x,y)= f(x, y, \text{context}) *$


*Approfondimenti in A.K. Jan, M. N. Murty, P. J. Flynn, "Data clustering: a review", ACM Computing Surveys, Vol. 31, No. 3, September 1999, ref. pp. 272-273.



Tecniche per il clustering

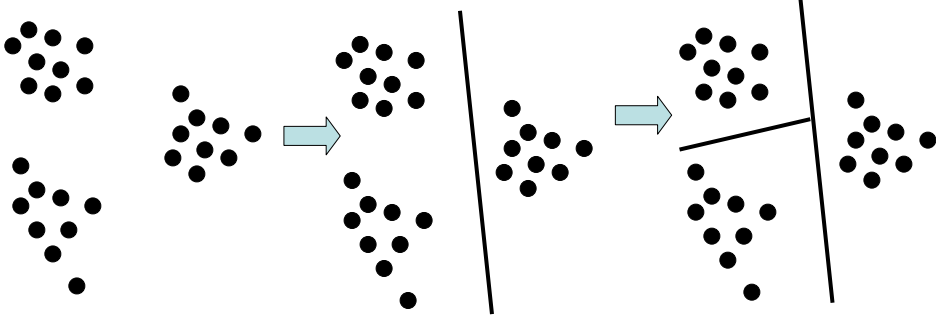
- **Hierarchical vs. partitional:** Gli algoritmi gerarchici producono una serie di partizioni al contrario degli algoritmi partitional, che lavorano sempre con una sola partizione.

21/40




Tecniche per il clustering

Hierarchical:



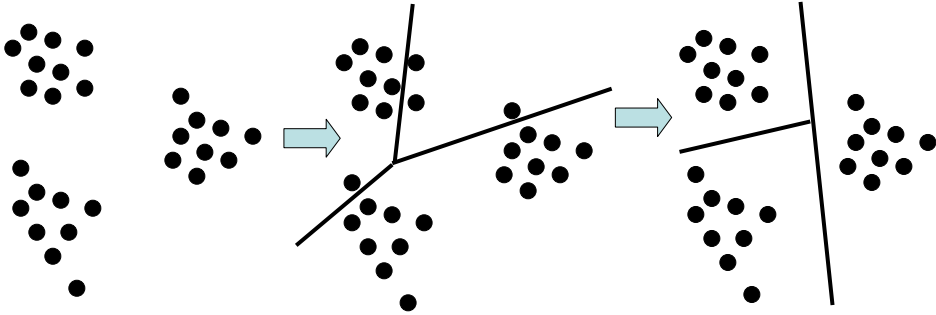
The diagram illustrates the hierarchical clustering process. It starts with two distinct clusters of points. An arrow indicates the first step, where a vertical line is drawn to separate the two clusters. A second arrow indicates the next step, where a horizontal line is drawn to further divide the top cluster into two sub-clusters, resulting in three clusters in total.

22/40





Tecniche per il clustering

Partitional:



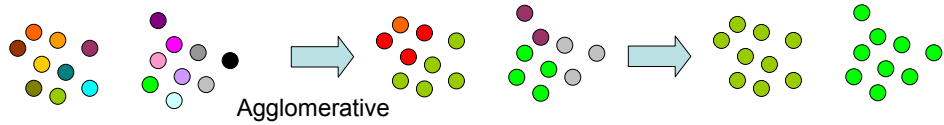
The diagram illustrates the partitional clustering process. It starts with two distinct clusters of points. An arrow indicates the first step, where a vertical line is drawn to separate the two clusters. A second arrow indicates the next step, where a horizontal line is drawn to further divide the top cluster into two sub-clusters, resulting in three clusters in total.

23/40

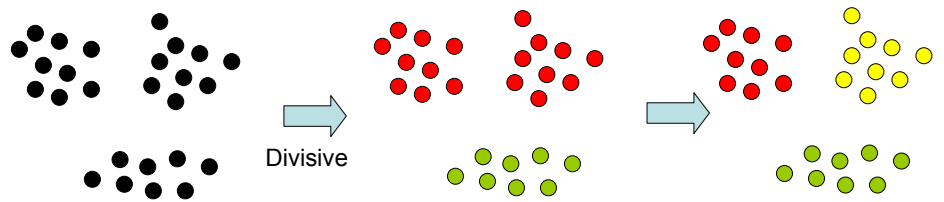



Tecniche per il clustering

- **Agglomerative vs. divisive:** i dati vengono accorpati oppure divisi nel processo di clustering.





Agglomerative



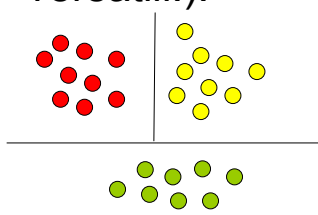
Divisive

24/40

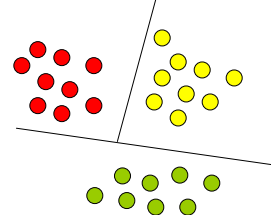



Tecniche per il clustering

- **Monothetic vs. polythetic:** si considera una sola feature (monothetic) oppure tutte le feature (polythetic) per volta per effettuare il clustering... La maggior parte degli algoritmi sono di tipo polythetic (più versatili!).



Monothetic



Polythetic



Tecniche per il clustering

- **Hard vs. fuzzy:** label vs. grado di appartenenza. Un algoritmo fuzzy può sempre essere defuzzyficato (classificazione dell'i-esimo pattern nella classe con grado di appartenenza più alto).



Tecniche per il clustering

- **Deterministic vs. stochastic:** introduzione di elementi random nell'algoritmo di clusterizzazione (es. GA) che rendono l'output dell'algoritmo non predicibile sulla base dei dati di partenza (stochastic).

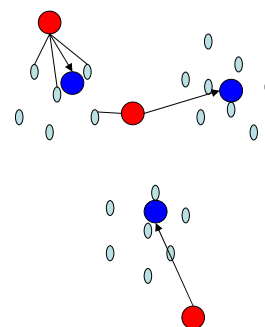
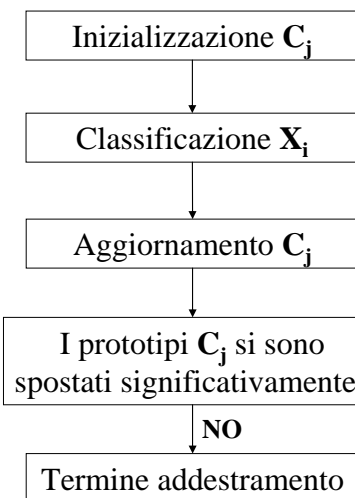


K-means (partitional): framework

- Siano $\mathbf{X}_1, \dots, \mathbf{X}_D$ i dati di addestramento (per semplicità, definiti in \mathbb{R}^2);
- siano $\mathbf{C}_1, \dots, \mathbf{C}_K$ i *prototipi* di K classi, definiti anch'essi in \mathbb{R}^2 ; ogni *prototipo* identifica il baricentro della classe corrispondente;
- lo schema di classificazione adottato sia il seguente: " \mathbf{X}_i appartiene a \mathbf{C}_j se e solo se \mathbf{C}_j è il *prototipo* più vicino a \mathbf{X}_i (distanza euclidea)";
- l'algoritmo di addestramento permette di determinare le posizioni dei *prototipi* \mathbf{C}_j mediante successive approssimazioni.



K-means: addestramento



Aggiornamento \mathbf{C}_j : baricentro degli \mathbf{X}_i classificati da \mathbf{C}_j .



K-means: limiti

- Partitional, polythetic, hard, deterministic;
- Veloce, semplice da implementare;
- Trova un minimo locale della funzione $f = \sum_j \sum_i [\text{dist}(x_i, \text{prot}_j)] / N_j$;
- Il risultato dipende dall'inizializzazione!
- Possono essere usati altri metodi (es. GA) per inizializzare K-means... es. GA per la minimizzazione di f , effettuano una ricerca globale, ma sono lenti!



K-Means e EM

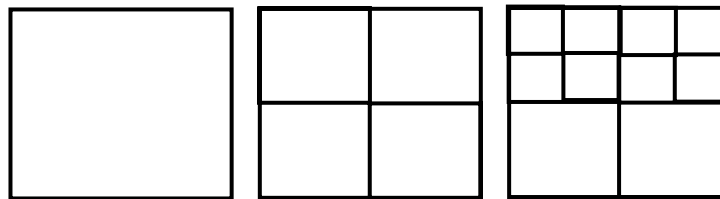
- K-Means come versione “limite” di EM (assegnazione “hard” delle responsabilità delle diverse componenti per ciascun dato);*
- Expectation → assegnazione (hard) dei dati ai cluster;
- Maximization → aggiornamento dei prototipi.

*Approfondimenti in Christopher M. Bishop, Pattern Recognition and Machine Learning, Capitolo 9.1, (K-means vs. mixture models, EM).



Algoritmi hierarchical: QTD

- Quad Tree Decomposition;
- Suddivisione gerarchica dello spazio delle feature, mediante splitting dei cluster;
- Criterio di splitting (\sim distanza tra cluster).



Algoritmi hierarchical: QTD

- Clusterizzazione immagini RGB, 512x512;
- Pattern: pixel (x,y);
- Feature: canali R, G, B.
- Distanza tra due pattern (non euclidea):

$$\text{dist}(p_1, p_2) =$$

$$\text{dist}([R_1 \ G_1 \ B_1], [R_2 \ G_2 \ B_2]) =$$

$$\max(|R_1 - R_2|, |G_1 - G_2|, |B_1 - B_2|).$$



33/40



Algoritmi hierarchical: QTD

$$p1 = [0 \ 100 \ 250]$$

$$p2 = [50 \ 100 \ 200]$$

$$p3 = [255 \ 150 \ 50]$$

$$\text{dist}(p1, p2) = \text{dist}([R1 \ G1 \ B1], [R2 \ G2 \ B2]) = \\ \max(|R1-R2|, |G1-G2|, |B1-B2|) = \max([50 \ 0 \ 50]) = 50.$$

$$\text{dist}(p2, p3) = 205.$$

$$\text{dist}(p3, p1) = 255.$$



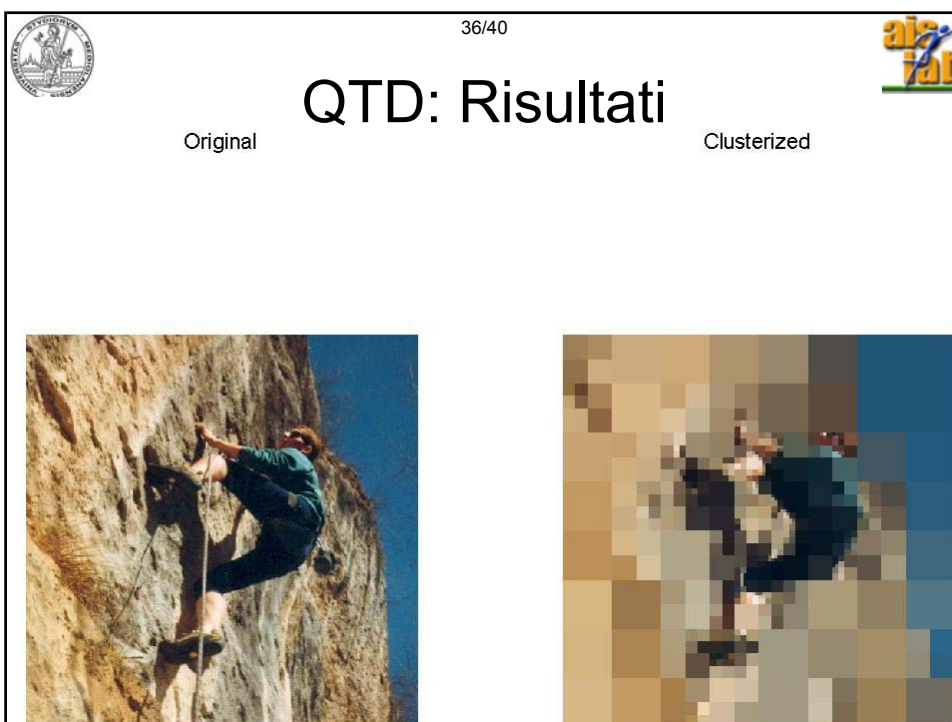
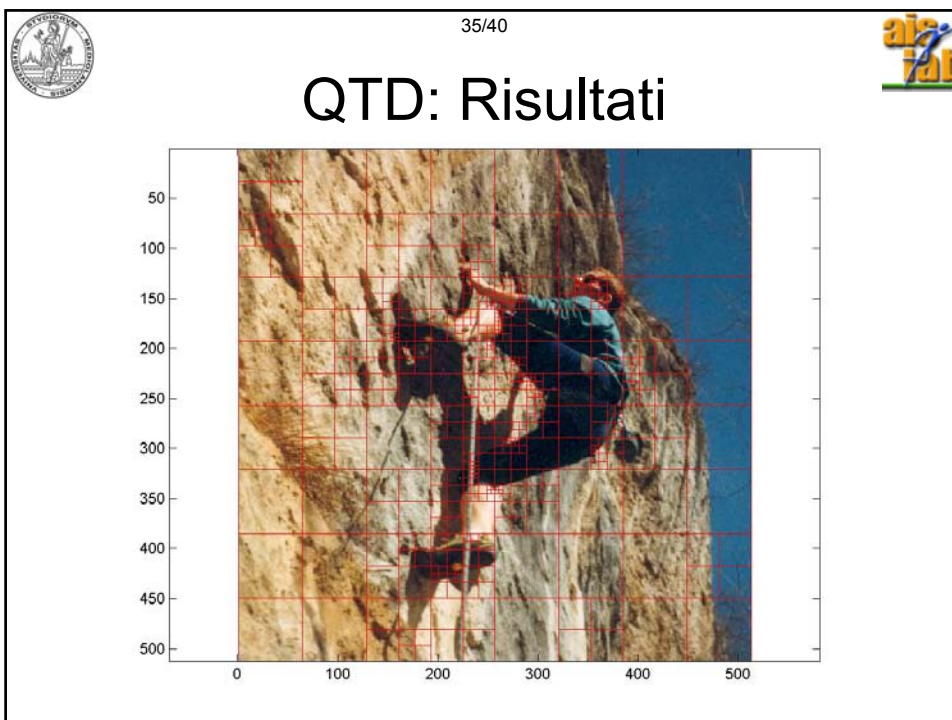
34/40

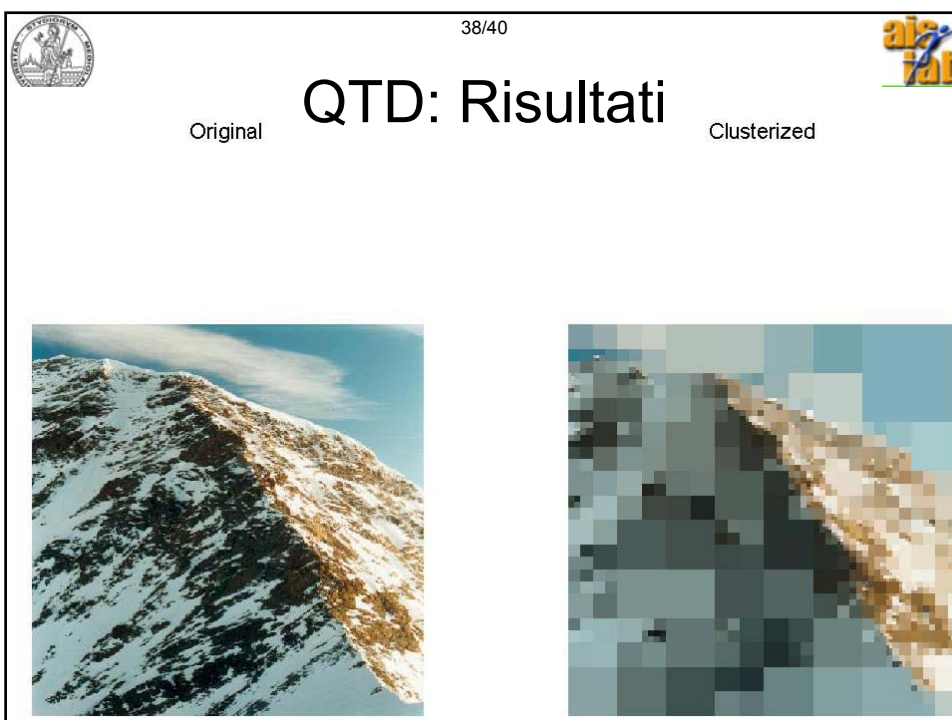
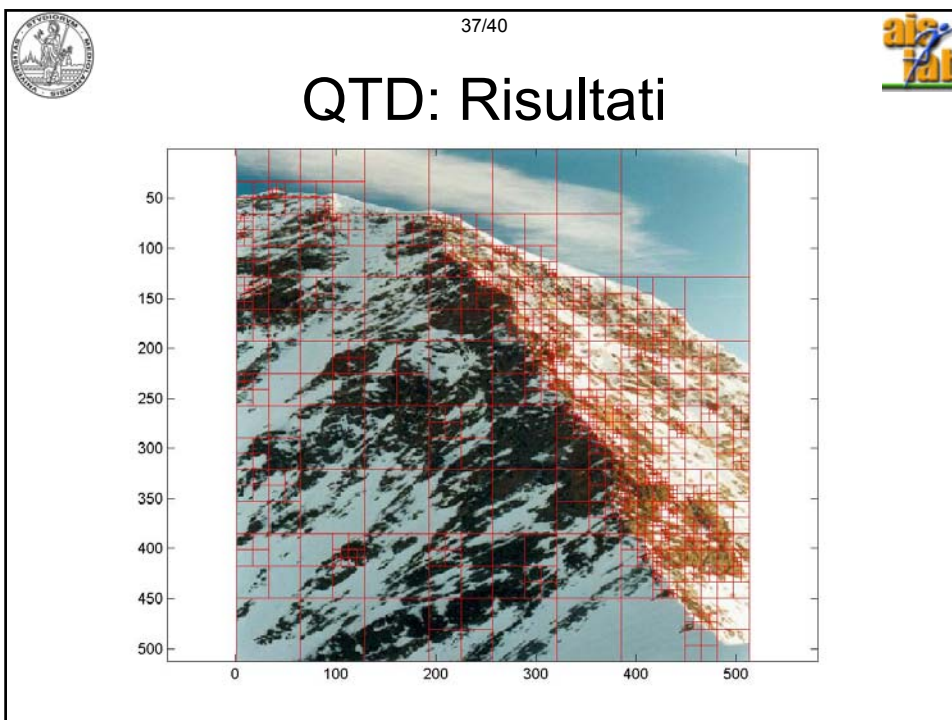


Algoritmi hierarchical: QTD

Criterio di splitting: se due pixel all'interno dello stesso cluster distano più di una determinata soglia, il cluster viene diviso in 4 cluster.

Esempio applicazione: segmentazione immagini, compressione immagini, analisi locale frequenze immagini...



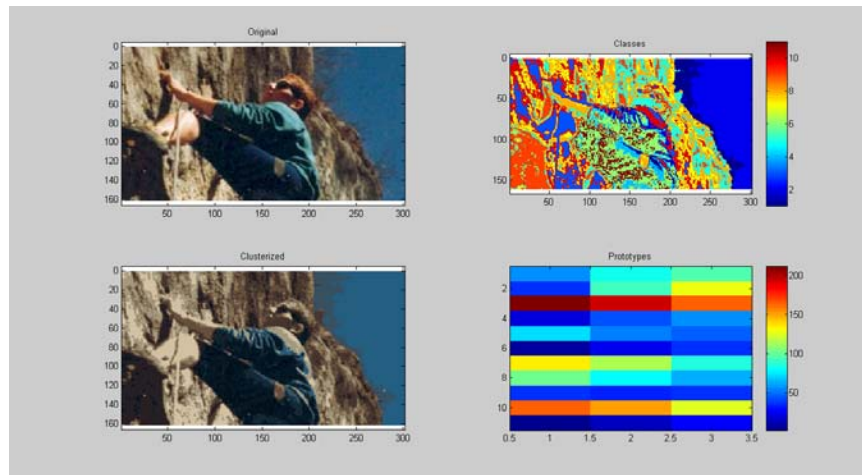




39/40



K-Means per immagine RGB



40/40



M-Files

- M-files (disponibili sul sito <http://homes.dsi.unimi.it/~frosio/>) e documenti allegati:
 - KMeans;
 - KMeans per clusterizzazione immagini RGB;
 - QTD Clustering.