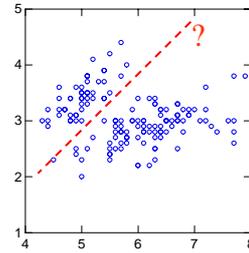


Classificazione (aka Cluster Analysis)

- Classificazione non gerarchica
 - es. *k*-means
- Classificazione gerarchica divisiva
- Classificazione gerarchica agglomerativa
 - Legame: singolo, completo, medio, ...
 - Coefficiente di correlazione cofenetica
- Classificazione con vincoli
- Metodi innovativi (machine learning)
 - es. Self-organizing maps

Classificazione

(= *Cluster Analysis*)



Obiettivo:

massimizzare l'omogeneità dei gruppi (o classi o clusters)

(cioè: *gli oggetti simili devono essere nello stesso cluster*)

Problema generale:

cercare le discontinuità nello spazio dei dati da classificare

Discontinuità

Reali o “naturali”:

es. tassonomia

Arbitrarie:

es. ecologia delle comunità

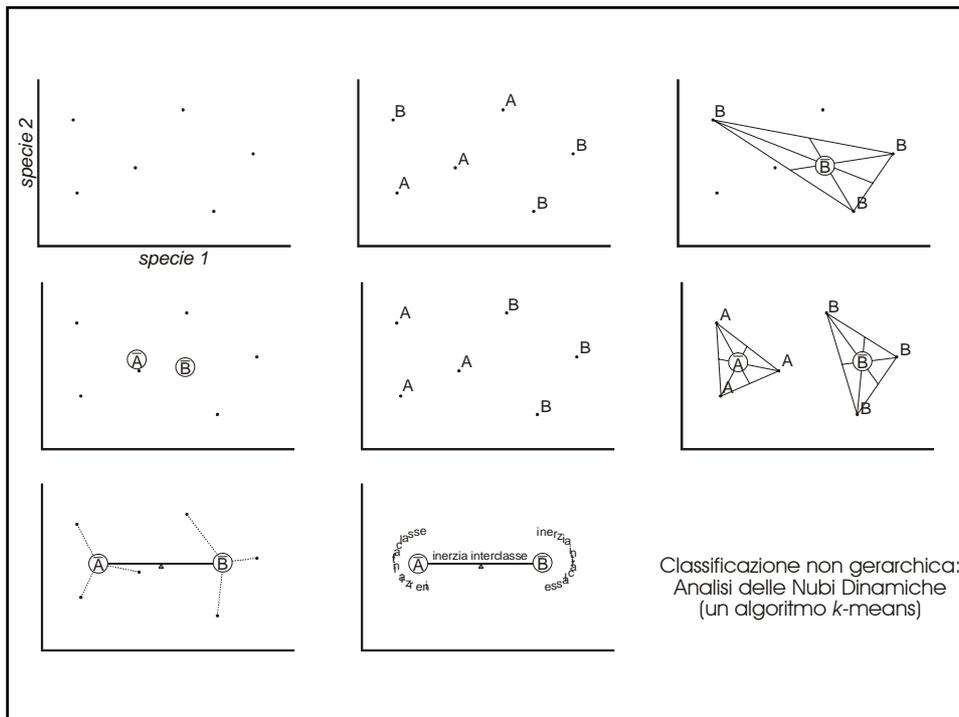
Clustering non-gerarchico: k -means

- Si utilizza direttamente la matrice dei dati calcolando distanze euclidee
- Si massimizza la varianza inter-classe per un numero dato *a priori* di classi
- In pratica, equivale a cercare le classi che massimizzano il rapporto F di una MANOVA a una via

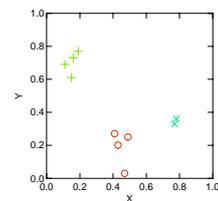
k -means

Procedura iterativa:

1. Scegli un numero di classi
2. Assegna gli oggetti alle classi
(a caso o in base ad un'altra classificazione)
3. Sposta gli oggetti nelle classi il cui centroide è più vicino (la varianza intra-classe diminuisce)
4. Ripeti lo step 3 finchè non c'è più nessun cambiamento nella composizione delle classi

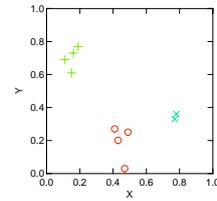


k-means con 3 classi



Variabile	InterSQ	gdl	IntraSQ	gdl	rapporto F
X	0.536	2	0.007	7	256.163
Y	0.541	2	0.050	7	37.566
** TOTALE **	1.078	4	0.058	14	

k-means con 3 classi



Classe 1 di 3 (contiene 4 casi)

Membri			Statistiche			
Caso	Distanza	Variabile	Min	Media	Max	Dev.St.
Caso 1	0.02	X	0.41	0.45	0.49	0.04
Caso 2	0.11	Y	0.03	0.19	0.27	0.11
Caso 3	0.06					
Caso 4	0.05					

Classe 2 di 3 (contiene 4 casi)

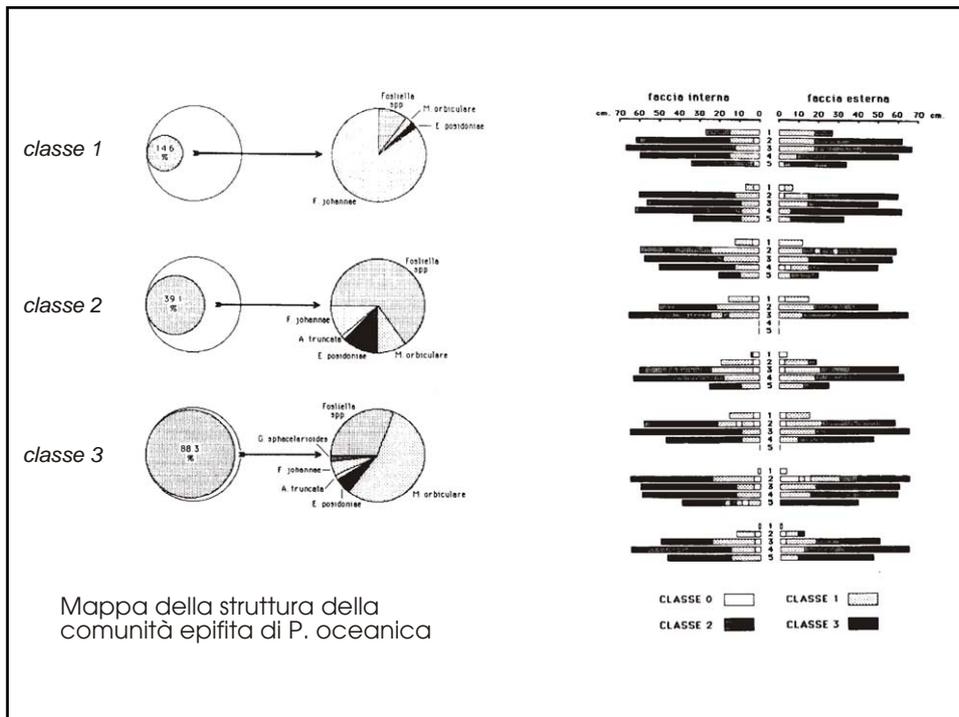
Membri			Statistiche			
Caso	Distanza	Variabile	Min	Media	Max	Dev.St.
Caso 7	0.06	X	0.11	0.15	0.19	0.03
Caso 8	0.03	Y	0.61	0.70	0.77	0.07
Caso 9	0.02					
Caso 10	0.06					

Classe 3 di 3 (contiene 2 casi)

Membri			Statistiche			
Caso	Distanza	Variabile	Min	Media	Max	Dev.St.
Caso 5	0.01	X	0.77	0.77	0.78	0.01
Caso 6	0.01	Y	0.33	0.35	0.36	0.02

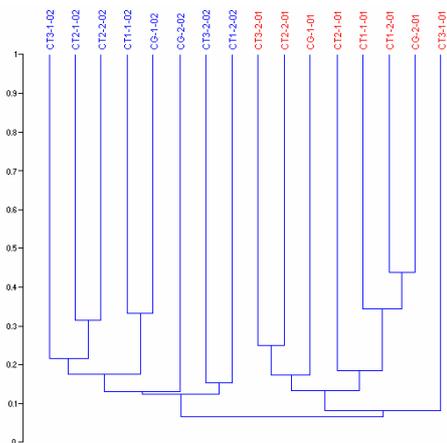
Svantaggi della classificazione k-means

- Ottimizza le classi, ma non in modo globale (iterando si può ottenere una partizione più efficace)
- Bisogna scegliere il numero di classi *a priori*
 - Ma quante classi devono essere scelte?



Classificazione gerarchica

- Solitamente si rappresenta con un dendrogramma

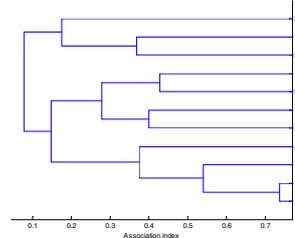


Classificazione gerarchica

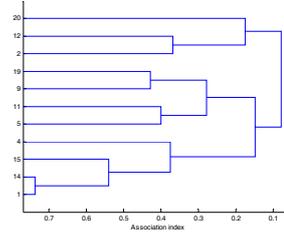
- Divisiva
 - Si cercano partizioni successive di un insieme di oggetti in due sotto-insiemi in modo da soddisfare un criterio predefinito
- Agglomerativa
 - Si aggregano gli oggetti in gruppi ed i gruppi fra loro secondo un criterio predefinito

Classificazione gerarchica divisiva

- Si inizia con tutti gli oggetti in una sola classe, poi:
 - si dividono mediante classificazioni k-means o altre tecniche divisive con partizioni in due classi;
 - si divide ogni volta la classe più eterogenea...
 - Oppure si dividono simultaneamente tutte le classi esistenti.
- Si tratta di tecniche efficienti, ma poco usate!



Classificazione gerarchica agglomerativa



- Si inizia con un oggetto per classe
- Le classi vengono fuse finchè ne rimane una soltanto
- E' la classe di metodi di gran lunga più ampiamente diffusa

Classificazione gerarchica agglomerativa

- Si lavora su una matrice di dissimilarità o distanza
- Procedura generale:
 1. Si calcola una matrice simmetrica di dissimilarità/distanza fra le classi
 2. Le classi fra cui la dissimilarità/distanza è minima vengono fuse fra loro
 3. Si calcola la dissimilarità/distanza fra la nuova classe ottenuta e tutte le altre

N.B. I diversi algoritmi differiscono nel punto 3

Classificazione gerarchica agglomerativa

	A	B	C	D	E		AD	B	C	E
A	0	AD	0	.	.	.
B	0.35	0	.	.	.	B	?	0	.	.
C	0.45	0.67	0	.	.	C	?	0.67	0	.
D	0.11	0.45	0.57	0	.	E	?	0.56	0.78	0
E	0.22	0.56	0.78	0.19	0					

Prima fusione: A e D

Come calcolare le nuove distanze?

Classificazione gerarchica agglomerativa

Legame semplice

	A	B	C	D	E		AD	B	C	E
A	0	AD	0	.	.	.
B	0.35	0	.	.	.	B	0.35	0	.	.
C	0.45	0.67	0	.	.	C	?	0.67	0	.
D	0.11	0.45	0.57	0	.	E	?	0.56	0.78	0
E	0.22	0.56	0.78	0.19	0					

$$d(AD,B) = \text{Min}\{d(A,B), d(D,B)\}$$

Classificazione gerarchica agglomerativa
Legame completo

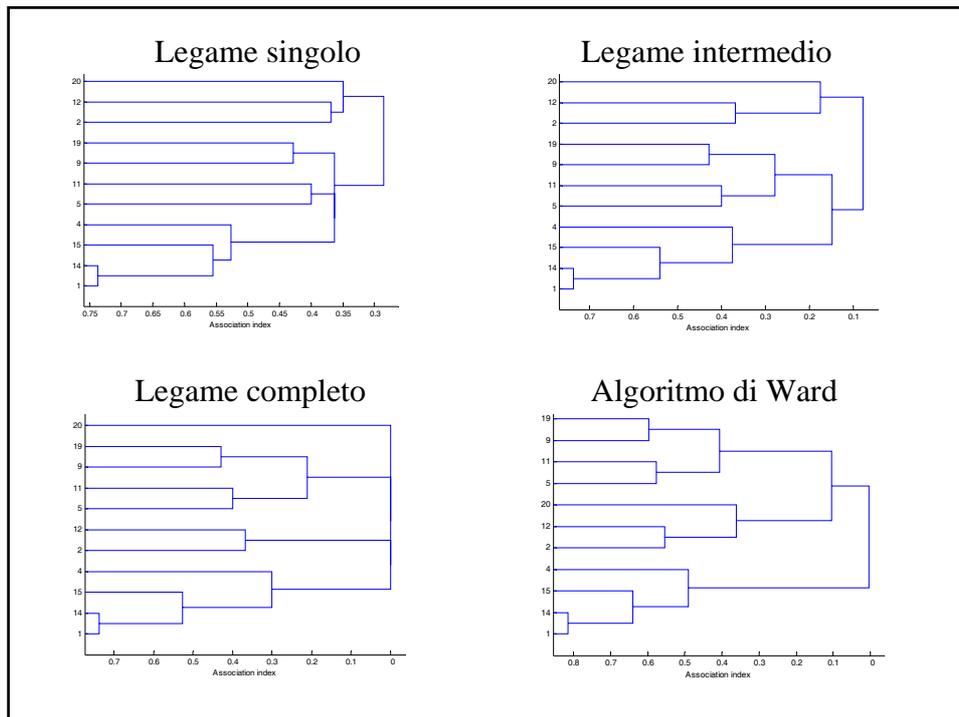
	A	B	C	D	E		AD	B	C	E
A	0	AD	0	.	.	.
B	0.35	0	.	.	.	B	0.45	0	.	.
C	0.45	0.67	0	.	.	C	?	0.67	0	.
D	0.11	0.45	0.57	0	.	E	?	0.56	0.78	0
E	0.22	0.56	0.78	0.19	0					

$$d(AD,B)=Max\{d(A,B), d(D,B)\}$$

Classificazione gerarchica agglomerativa
Legame intermedio

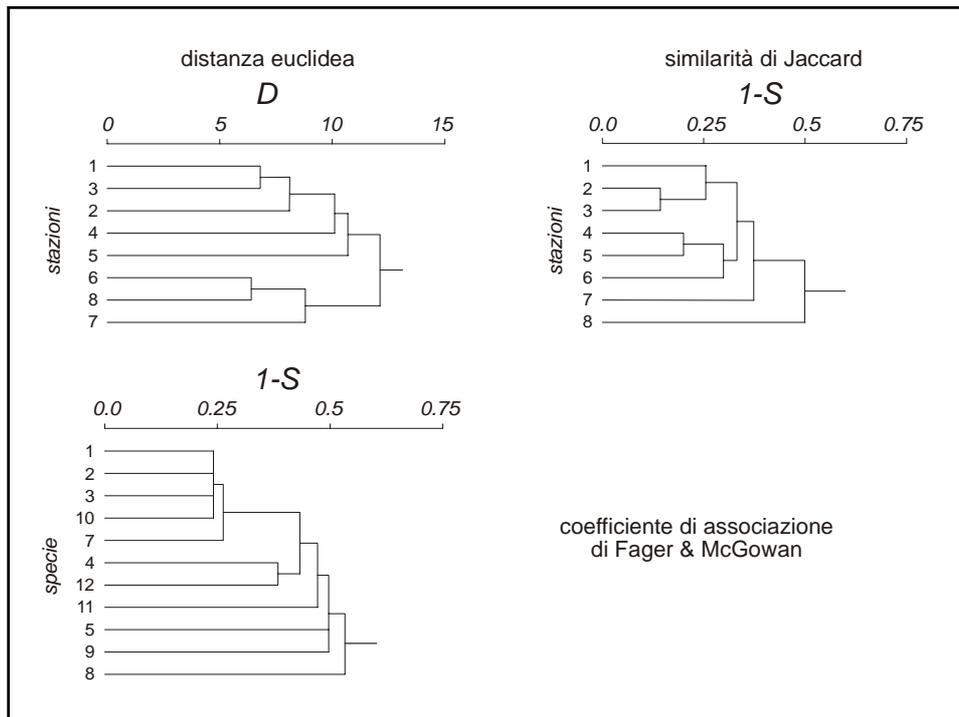
	A	B	C	D	E		AD	B	C	E
A	0	AD	0	.	.	.
B	0.35	0	.	.	.	B	0.40	0	.	.
C	0.45	0.67	0	.	.	C	?	0.67	0	.
D	0.11	0.45	0.57	0	.	E	?	0.56	0.78	0
E	0.22	0.56	0.78	0.19	0					

$$d(AD,B)=Mean\{d(A,B), d(D,B)\}$$



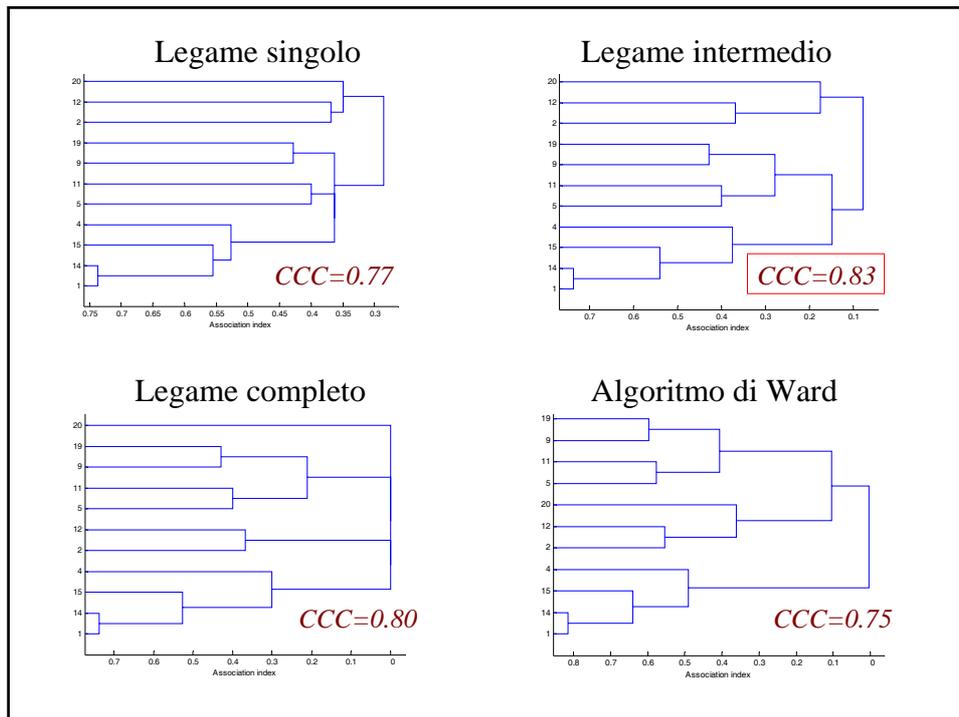
Metodi di classificazione gerarchica agglomerativa

- Legame semplice
 - Produce “catene” di oggetti, “contrae” lo spazio intorno alle classi
 - Non appropriato se l’errore di campionamento è grande
 - Usato in tassonomia
 - Invariante rispetto alle trasformazioni dei dati
- Legame completo
 - Produce classi compatte e “dilata” lo spazio intorno a esse
 - Non appropriato se l’errore di campionamento è grande
 - Invariante rispetto alle trasformazioni dei dati
- Legame intermedio, centroide, di Ward
 - Maggiore probabilità di riconoscere partizioni “naturalì”
 - Non invarianti rispetto alle trasformazioni dei dati



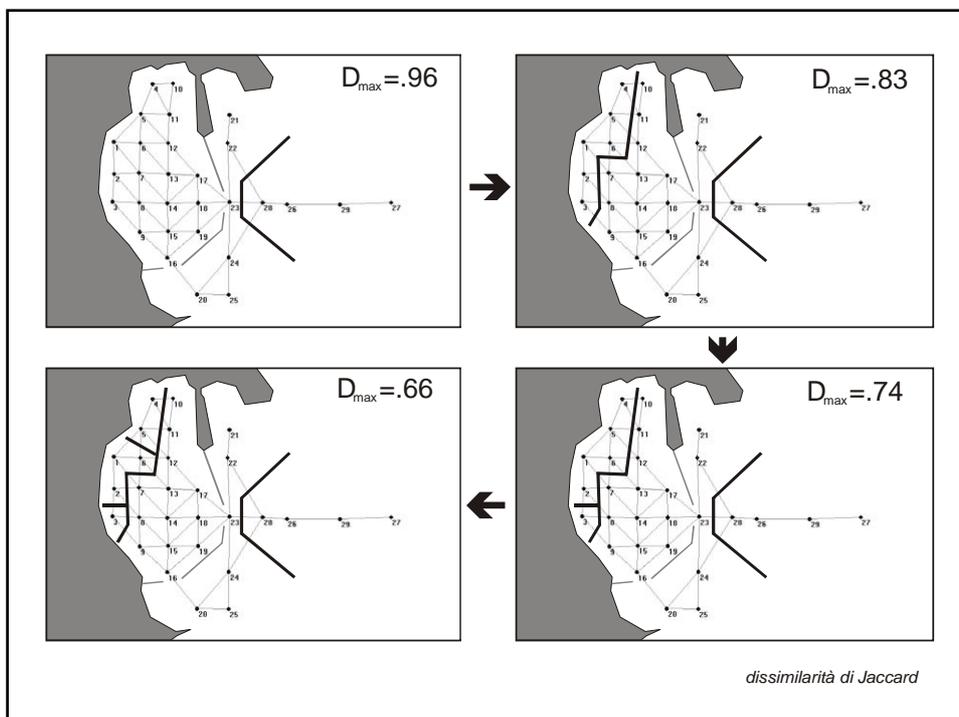
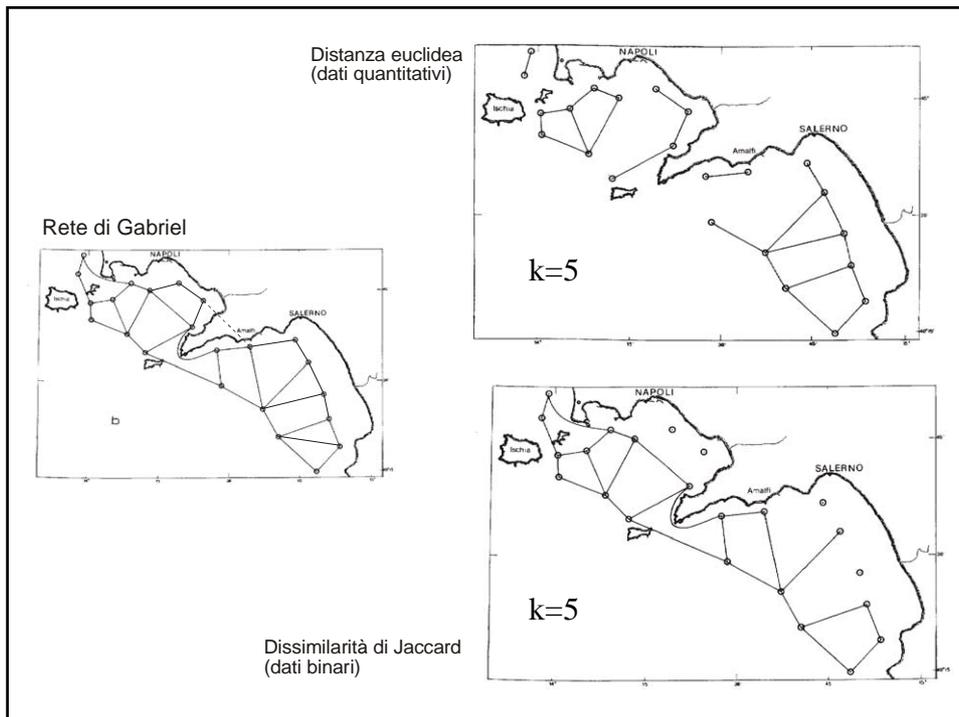
Coefficiente di correlazione cofenetica (CCC)

- Correlazione fra la matrice di dissimilarità originale e quella inferita in base alla classificazione
- CCC $> \sim 0.8$ indica un buon accordo
- CCC $< \sim 0.8$ indica che il dendrogramma non è una buona rappresentazione delle relazioni fra oggetti
- Si può usare il CCC per scegliere l'algoritmo ottimale



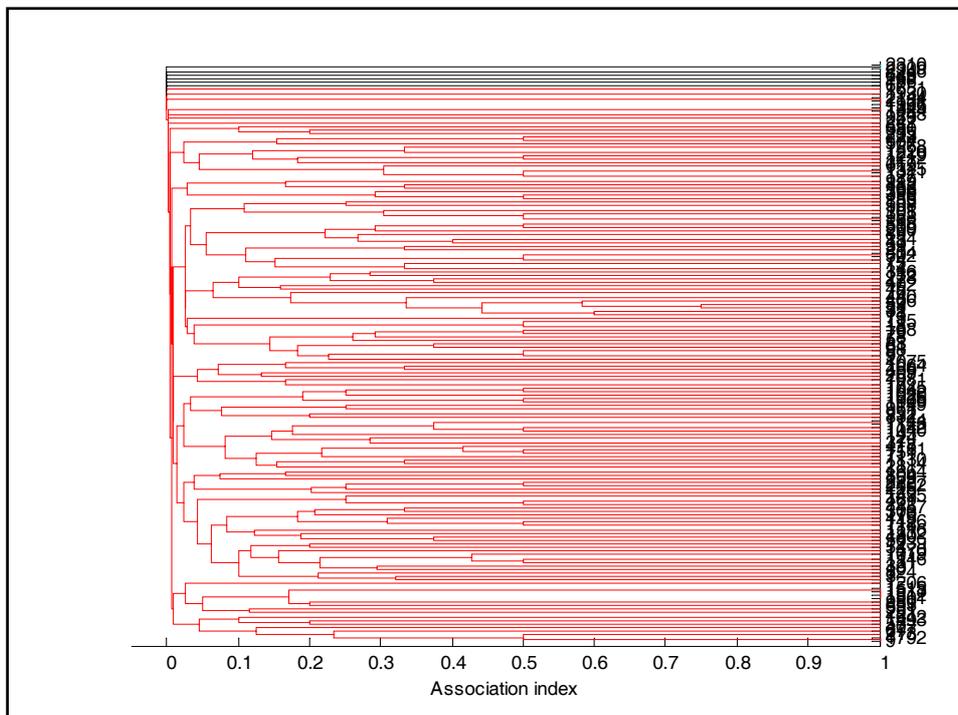
Classificazione vincolata

- Si definisce un vincolo di natura spaziale (es. contiguità fra oggetti) o di natura temporale (es. successione)
- Si definisce un metodo per verificare che il vincolo sia rispettato (es. rete di Gabriel per rappresentare la contiguità)
- Si applica il vincolo ad un qualsiasi algoritmo di classificazione
- Le classi ottenute possono essere più vicine alla nostra percezione della realtà
- Il fuoco si sposta dalla composizione delle classi alle discontinuità fra di esse



Problemi aperti

- Esistono realmente classi biologicamente o ecologicamente rilevanti?
- Il dendrogramma rappresenta la realtà biologica (web-of-life vs. tree-of-life)?
- Quante classi usare?
 - le regole per la selezione sono arbitrarie!
- Quale metodo usare?
 - la tecnica ottimale dipende dai dati!
- I dendrogrammi non sono efficienti nel visualizzare classificazioni complesse...

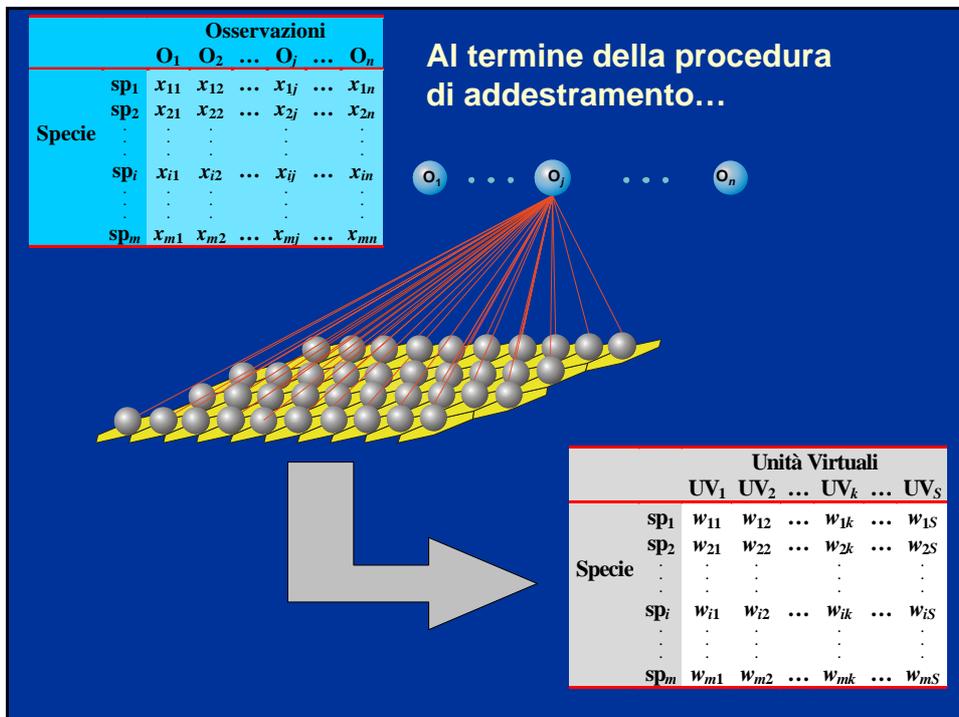
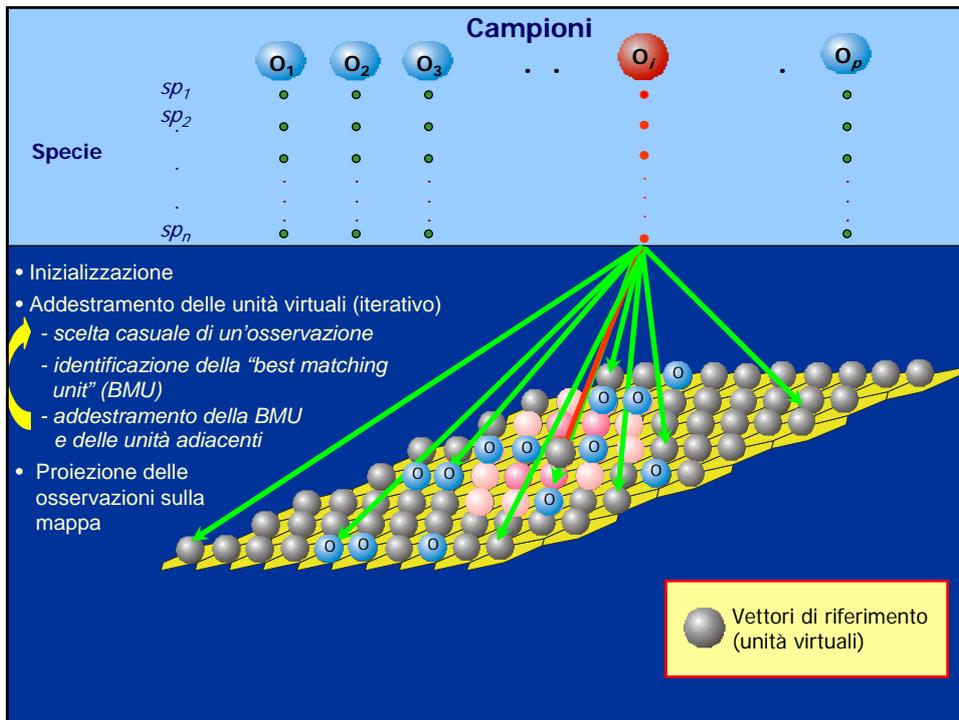


C'è qualcosa di meglio per casi complessi?

- Esistono metodi di nuova generazione, che non sono ancora molto diffusi in Ecologia
- La loro diffusione sta crescendo, ma bisogna che siano accettati dagli ecologi
- In molti casi, sfruttano le capacità di calcolo attualmente disponibili per tutti
- Si tratta soprattutto di tecniche mutuatae dal campo del Machine Learning, dell'Intelligenza Artificiale e del Data Mining

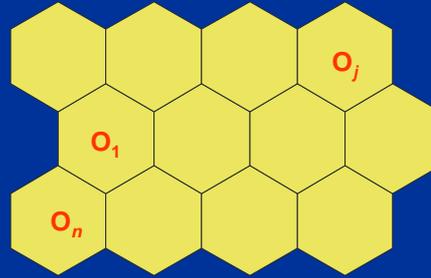
Ordinamento e classificazione:

Self-Organizing Maps (SOMs) o Mappe di Kohonen



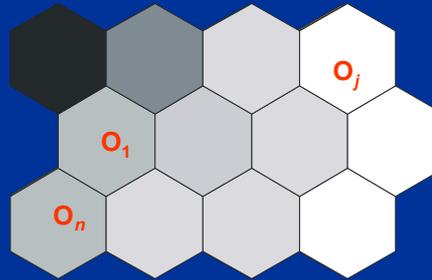
• Visualizzazione delle osservazioni

	Osservazioni					
	O_1	O_2	...	O_j	...	O_n
SP1	x_{11}	x_{12}	...	x_{1j}	...	x_{1n}
SP2	x_{21}	x_{22}	...	x_{2j}	...	x_{2n}
...
SPi	x_{i1}	x_{i2}	...	x_{ij}	...	x_{in}
...
SPm	x_{m1}	x_{m2}	...	x_{mj}	...	x_{mn}



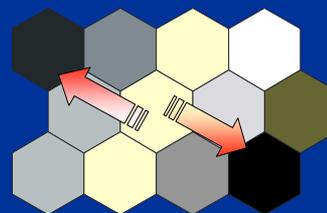
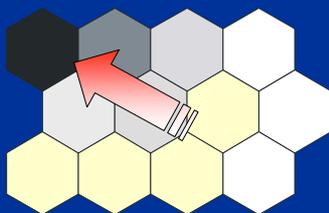
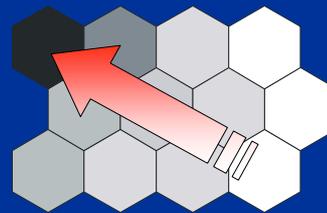
• Visualizzazione delle Unità Virtuali

	Unità Virtuali					
	UV_1	UV_2	...	UV_k	...	UV_s
SP1	w_{11}	w_{12}	...	w_{1k}	...	w_{1s}
SP2	w_{21}	w_{22}	...	w_{2k}	...	w_{2s}
...
SPi	w_{i1}	w_{i2}	...	w_{ik}	...	w_{is}
...
SPm	w_{m1}	w_{m2}	...	w_{mk}	...	w_{ms}

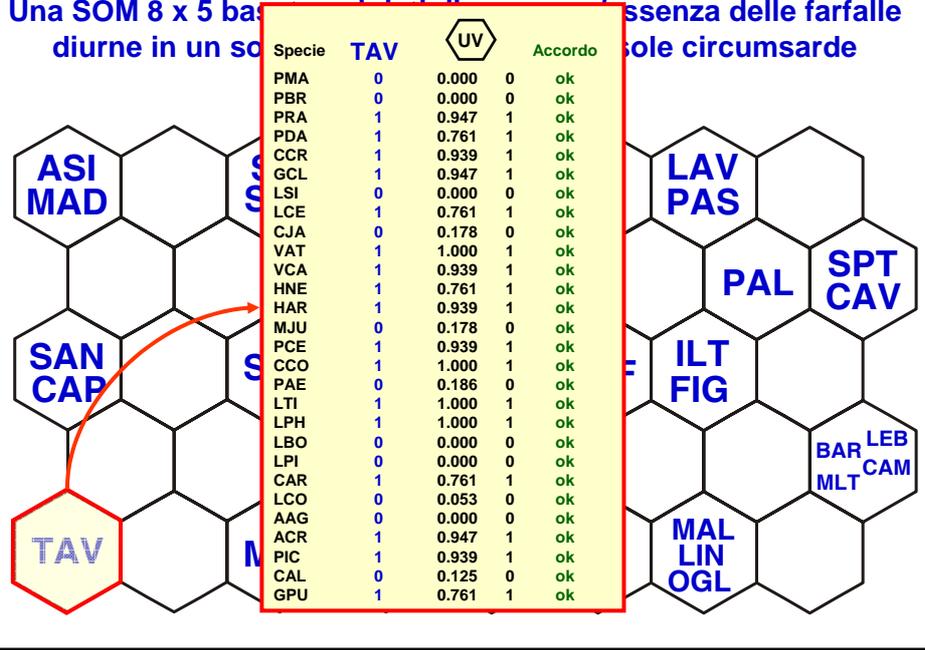


Ogni specie ha una diversa distribuzione...

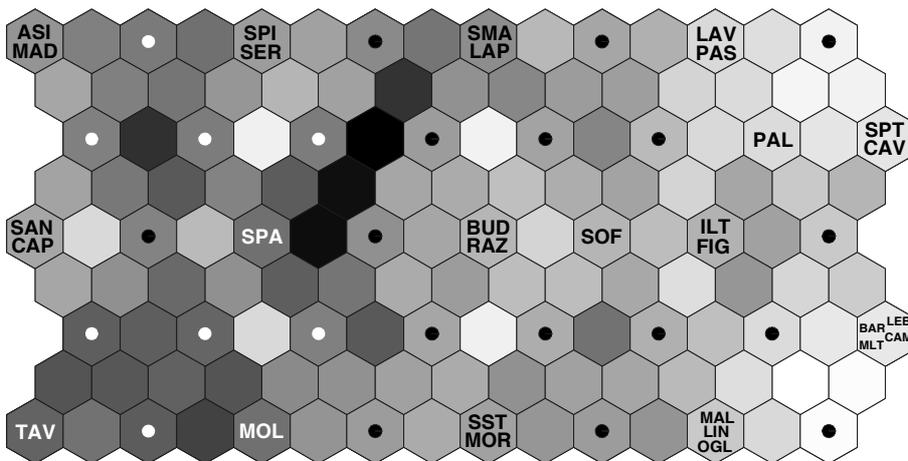
	Unità Virtuali					
	UV_1	UV_2	...	UV_k	...	UV_s
SP1	w_{11}	w_{12}	...	w_{1k}	...	w_{1s}
SP2	w_{21}	w_{22}	...	w_{2k}	...	w_{2s}
...
SPi	w_{i1}	w_{i2}	...	w_{ik}	...	w_{is}
...
SPm	w_{m1}	w_{m2}	...	w_{mk}	...	w_{ms}



Una SOM 8 x 5 basata sulla presenza delle farfalle diurne in un solo mese di sole circumsardegna



La matrice U consente di visualizzare le distanze fra gli elementi di una SOM.



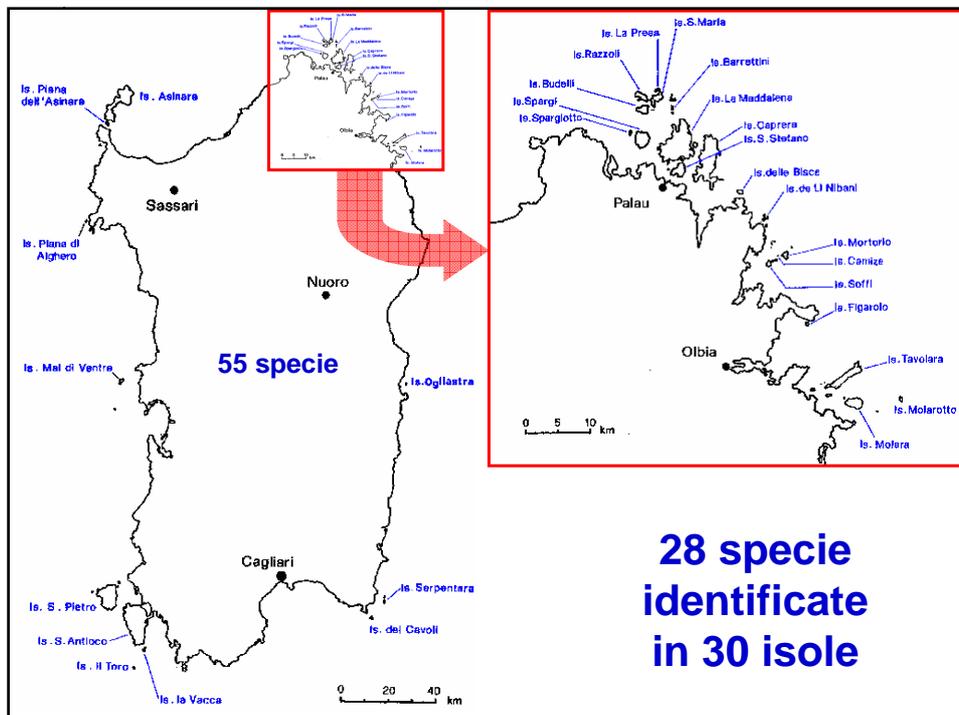
Un'applicazione delle **Self-Organizing Maps** (SOMs)

Dati estratti da:

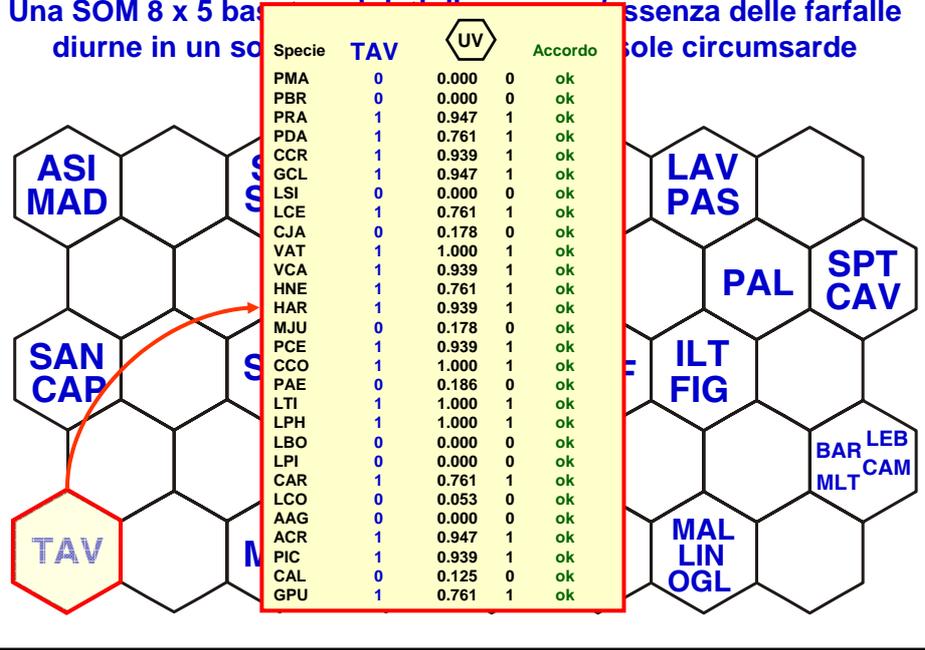
Marina Cobolli, Marco Lucrelli e Valerio
Sbordoni (1996).

Le farfalle diurne delle piccole isole
circumsarde.

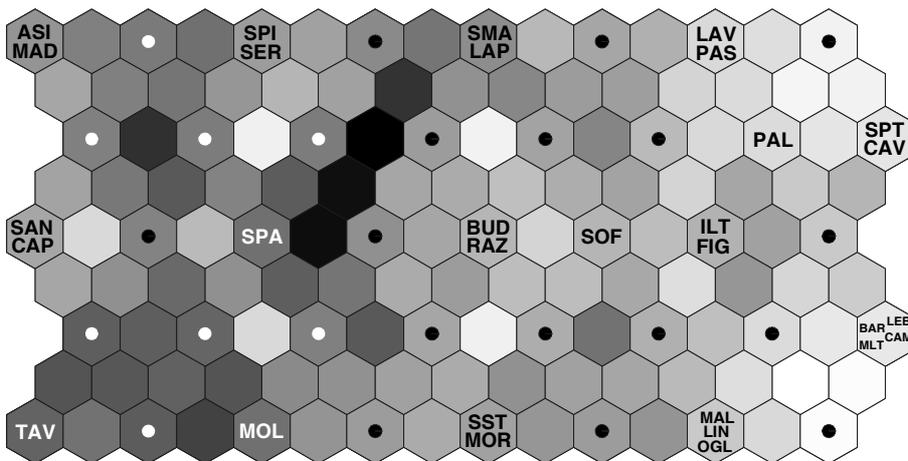
Biogeographia, 18: 569-582



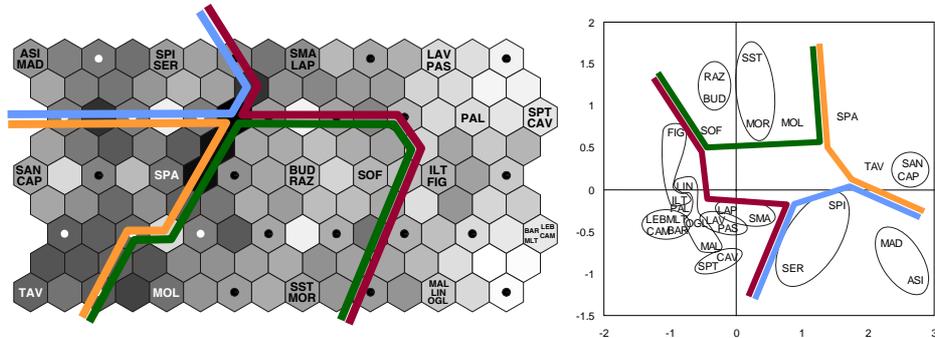
Una SOM 8 x 5 basata sulla presenza delle farfalle diurne in un solo mese di sole circumsardegna



La matrice U consente di visualizzare le distanze fra gli elementi di una SOM.



Self-Organizing Map vs. Analisi delle Coordinate Principali



Gonepteryx cleopatra (L., 1767)

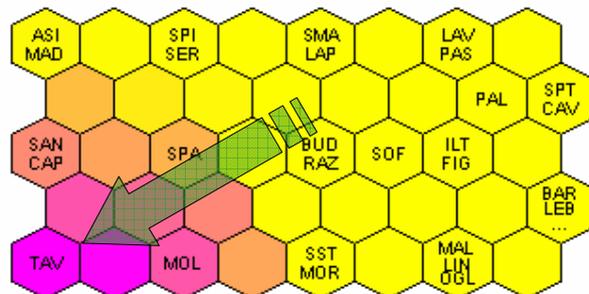
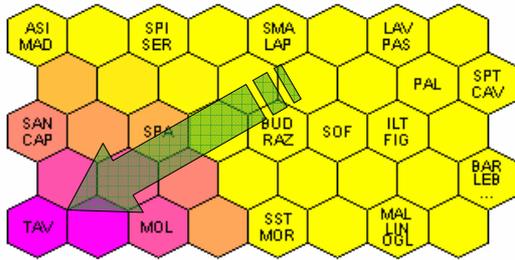
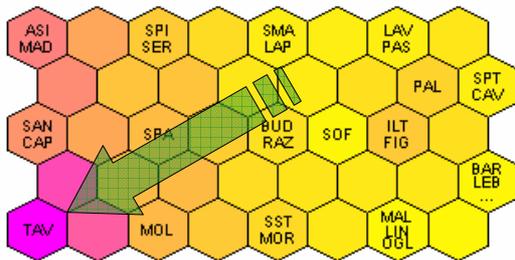


Foto: www.leps.it



Gonepteryx cleopatra

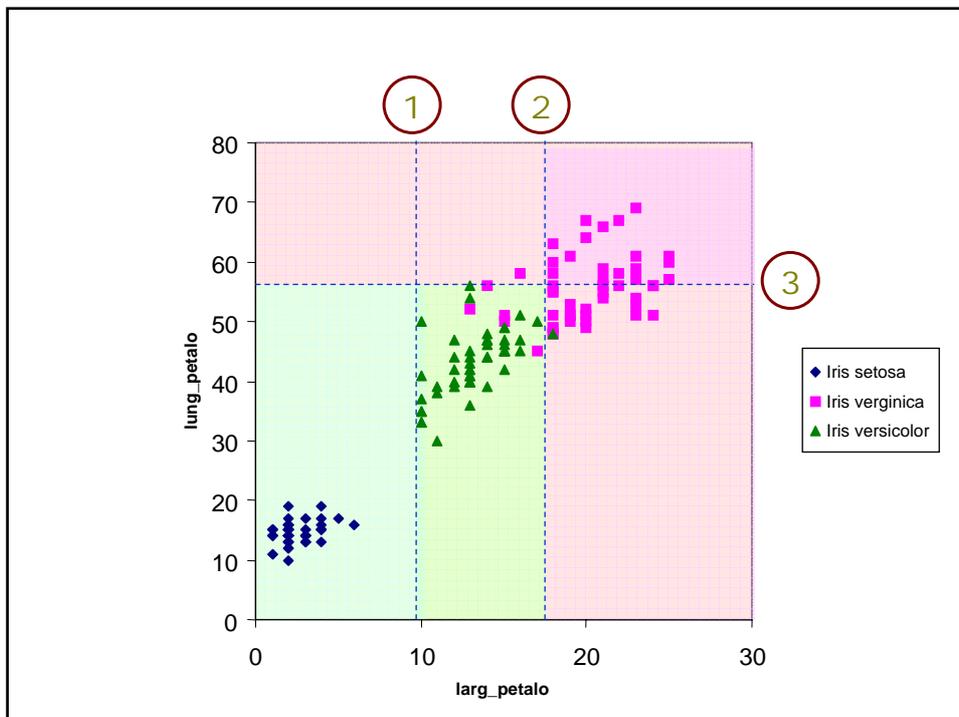
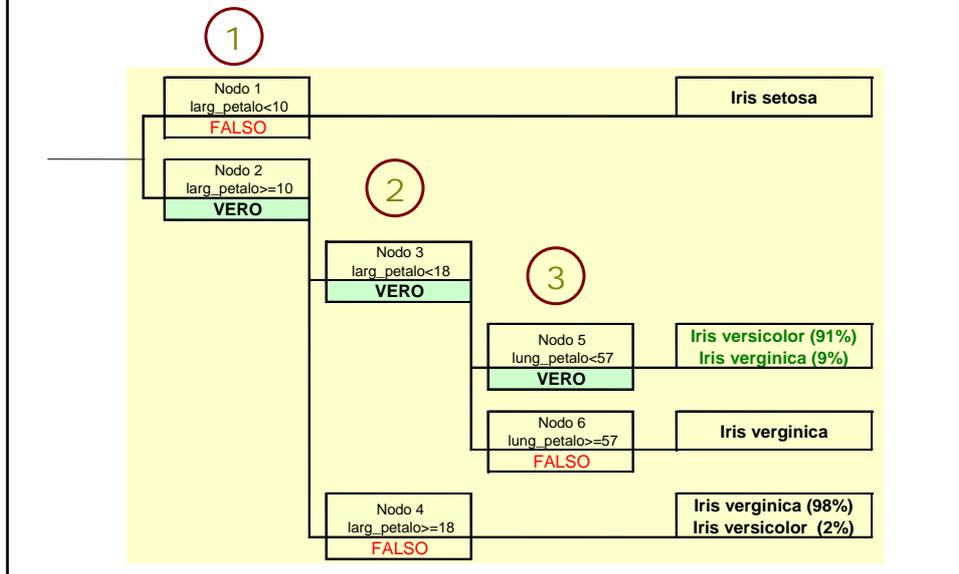


Altitudine

Foto: www.leps.it

Classification Trees

Il metodo: un esempio classico...



Modelli strutturali: classification tree

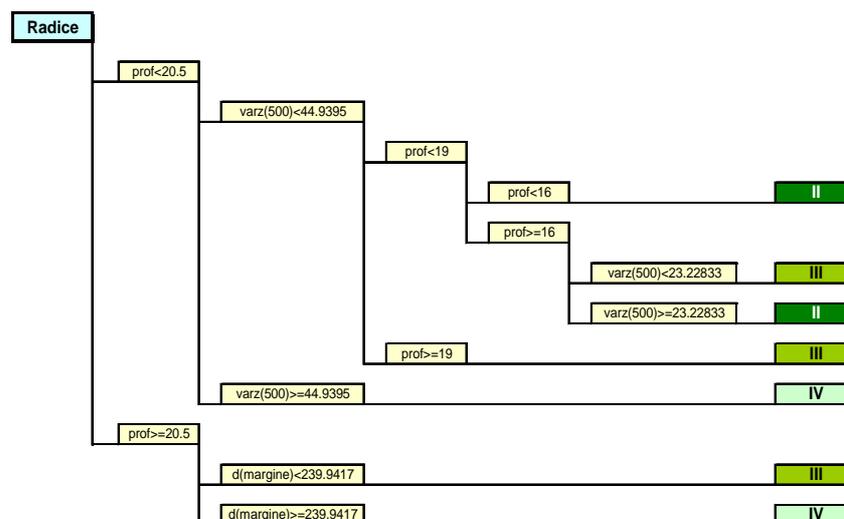
Variabili predittive:

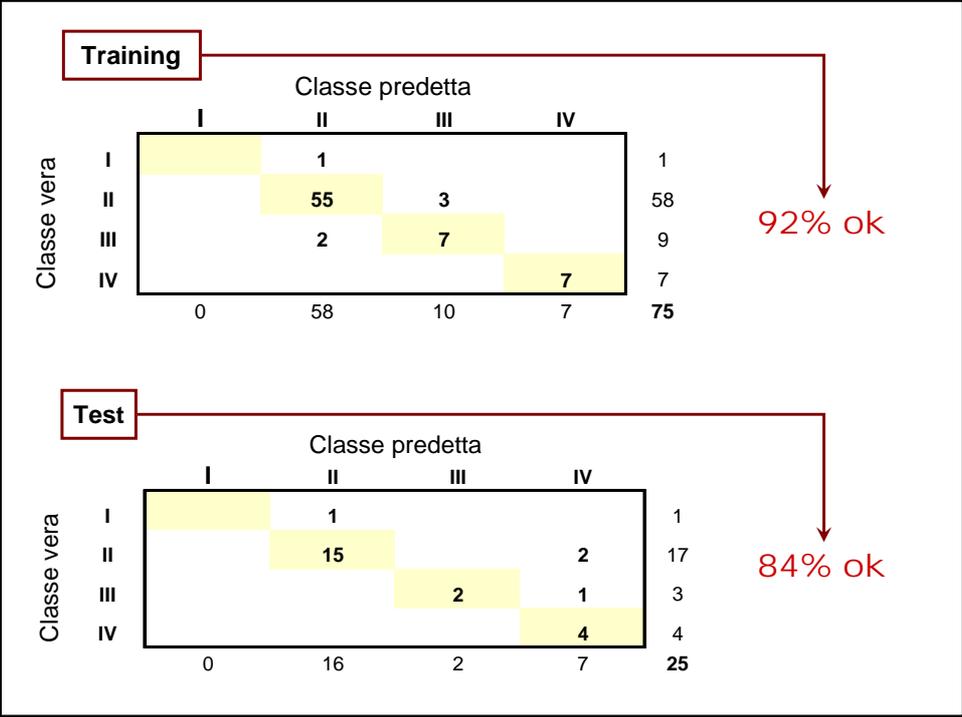
- Profondità
- Varianza della profondità (raggio=500 m)
- Distanza dal punto noto senza *Posidonia* più vicino (distanza dal margine della prateria stimata per eccesso)

Variabile da predire:

- Classe di densità della prateria (secondo Giraud)

Un modello per Rosignano:





In sintesi...

Alcuni metodi di classificazione

Tipo	Tecnica	Uso
<i>Non gerarchico</i>	<i>k</i> -means	Partizionare insiemi di dati
<i>Gerarchico divisivo</i>	<i>k</i> -means ripetuto	Buona tecnica su piccoli insiemi di dati
<i>Gerarchico agglomerativo</i>	Legame semplice	Tassonomia
	Legame completo	Classi omogenee
	Legame intermedio, centroide, di Ward	Scelta più frequente
<i>Machine learning</i>	SOM	Visualizzazione di insiemi di dati complessi
	Classification trees	Previsione e pattern recognition