

Intelligenza Artificiale II

Self-Organizing Systems

Introduzione

Marco Piastra

Self-Organization

- (*Wikipedia*)

“Self-organization is a process in which the internal organization of a system, normally an open system, increases in complexity without being guided or managed by an outside source.

Self-organizing systems typically (though not always) display *emergent properties*.”

Emergence

- (Corning P.A., 2002, in *Wikipedia*)

“ Perhaps the most elaborate recent definition of **emergence** was provided by Jeffrey Goldstein in the inaugural issue of “Emergence”. (Goldstein 1999)

To Goldstein, emergence refers to “the arising of novel and coherent structures, patterns and properties during the process of self-organization in complex systems.”

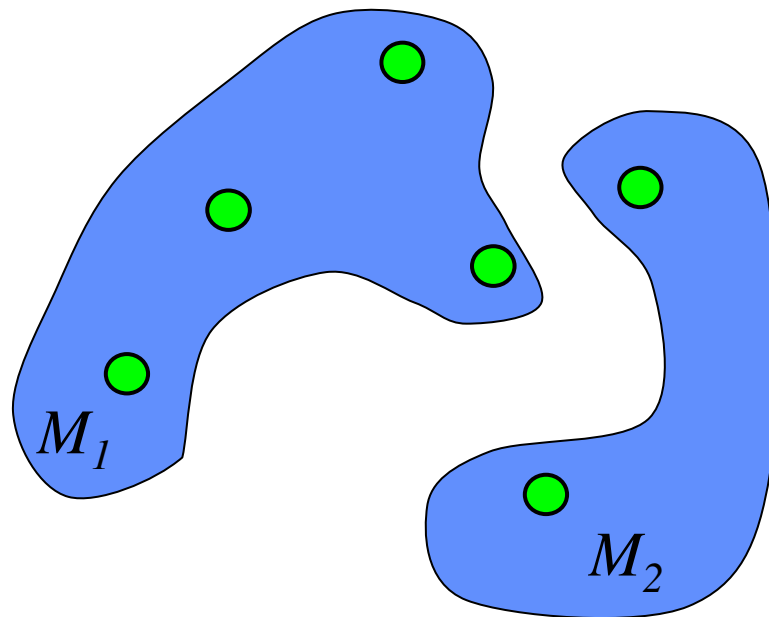
The common characteristics are:

- 1) radical novelty (features not previously observed in systems)
- 2) coherence or correlation (meaning integrated wholes that maintain themselves over some period of time)
- 3) a global or macro "level" (i.e. there is some property of "wholeness")
- 4) it is the product of a dynamical process (it evolves)
- 5) it is "ostensive" - it can be perceived ”

Preludio topologico-geometrico (1)

- Tessellazione di Voronoi

I punti verdi sono *campioni* tratti da due 'oggetti' M_1 e M_2

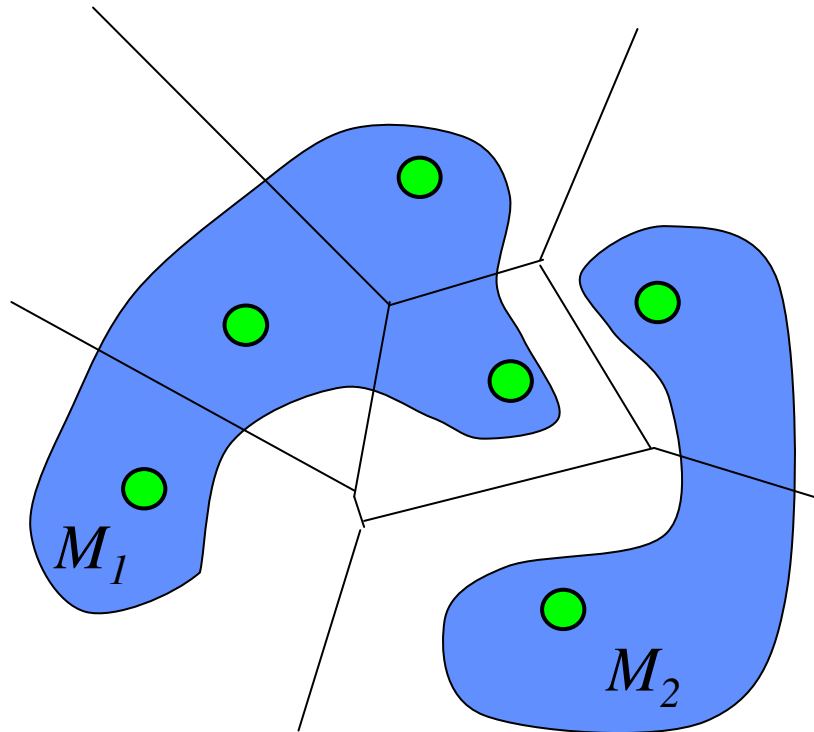


Preludio topologico-geometrico (1)

▪ Tessellazione di Voronoi

I punti verdi sono *campioni* tratti da due 'oggetti' M_1 e M_2

La *regione di Voronoi* associata a ciascun campione è formata da tutti i punti che sono più vicini al campione che a qualsiasi altro

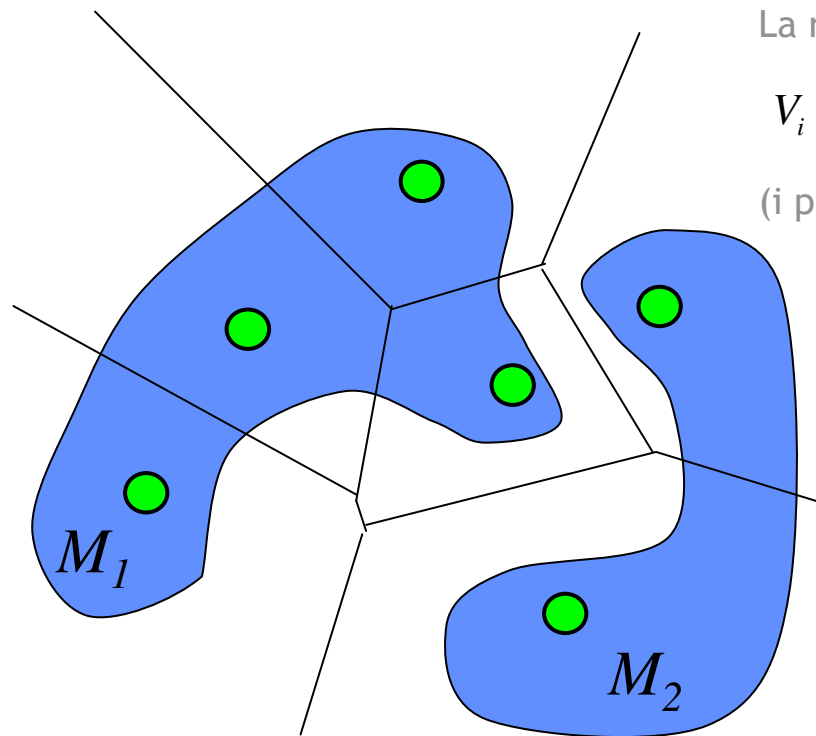


Preludio topologico-geometrico (1)

▪ Tessellazione di Voronoi

I punti verdi sono *campioni* tratti da due 'oggetti' M_1 e M_2

La *regione di Voronoi* associata a ciascun campione è formata da tutti i punti che sono più vicini al campione che a qualsiasi altro



La regione di Voronoi associata a ciascun p_i

$$V_i = \left\{ \xi \in \mathbb{R}^d \mid i = \arg \min_j \|\xi - p_j\| \right\}$$

(i punti sul confine sono *condivisi*)

Preludio topologico-geometrico (1)

▪ Tessellazione di Voronoi

I punti verdi sono *campioni* tratti da due 'oggetti' M_1 e M_2

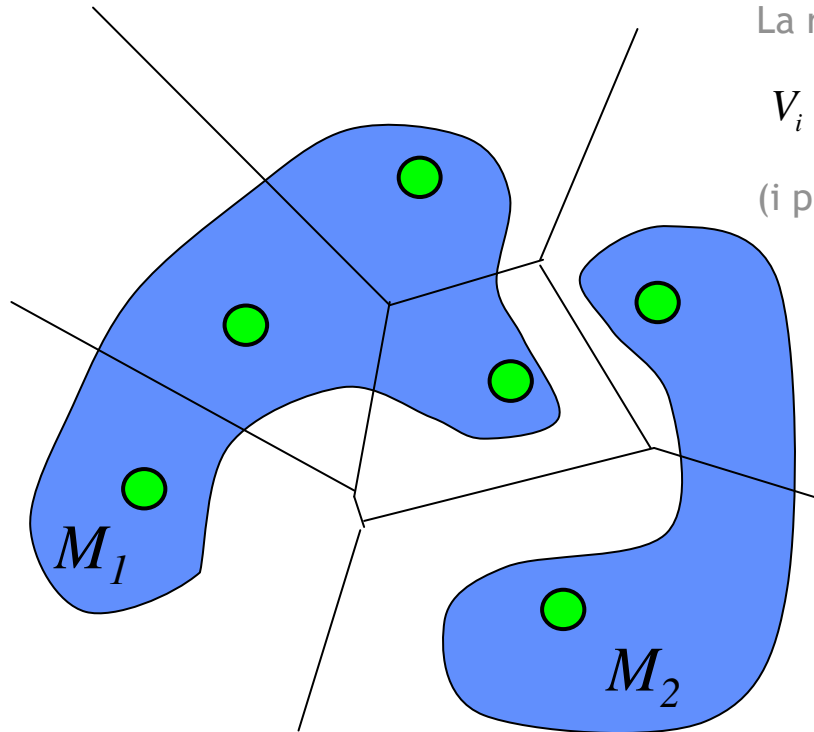
La **regione di Voronoi** associata a ciascun campione è formata da tutti i punti che sono più vicini al campione che a qualsiasi altro

La **tessellazione di Voronoi** (di \mathbf{R}^d) è formata dall'insieme delle regioni di Voronoi

La regione di Voronoi associata a ciascun p_i

$$V_i = \left\{ \xi \in \mathbf{R}^d \mid i = \arg \min_j \|\xi - p_j\| \right\}$$

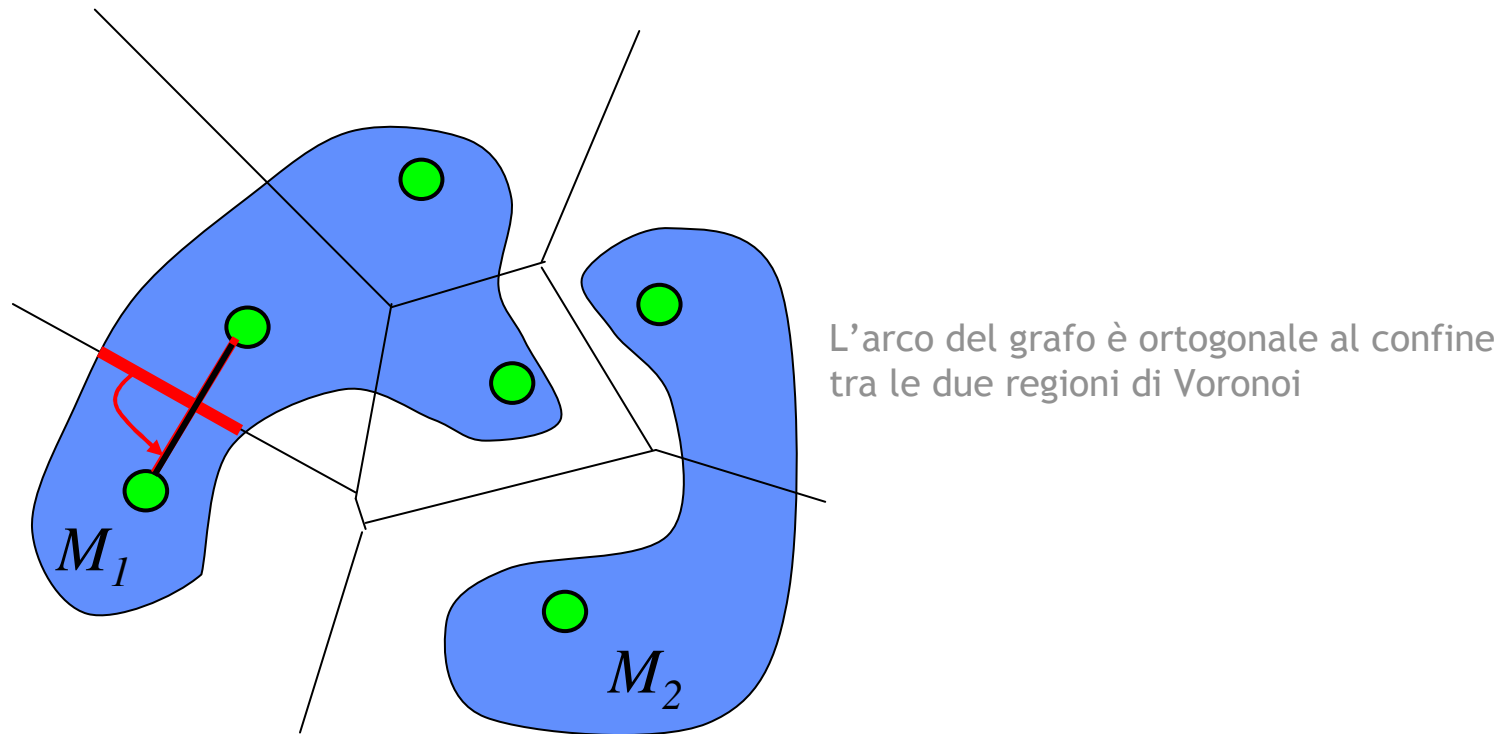
(i punti sul confine sono *condivisi*)



Preludio topologico-geometrico (2)

▪ Grafo di Delaunay

Due campioni le cui regioni di Voronoi hanno punti in comune (= sono confinanti) sono uniti da un arco

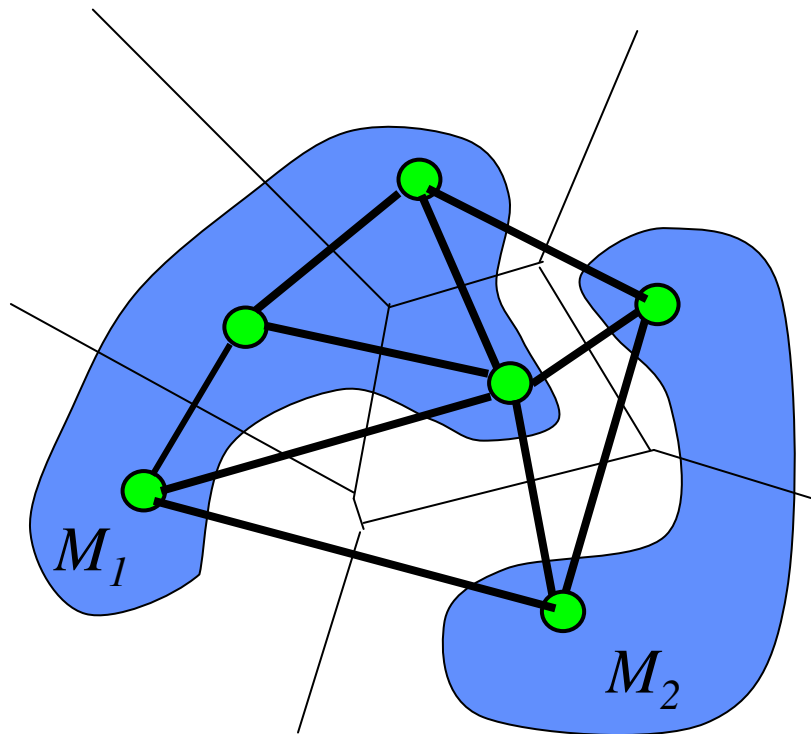


Preludio topologico-geometrico (2)

▪ Grafo di Delaunay

Due campioni le cui regioni di Voronoi hanno punti in comune (= sono confinanti) sono uniti da un arco

Il **grafo di Delaunay** è formato dai campioni e l'insieme di tutti gli archi che uniscono regioni Voronoi confinanti



Il grafo di Delaunay forma un *inviluppo convesso* dei campioni, diviso all'interno in regioni triangolari

La triangolazione di Delaunay ha proprietà notevoli, p.es. massimizza l'angolo minimo dei triangoli

k-means (*Generalized Lloyd's Algorithm*)

- Descrizione

Dato un set di dati D (appartenenti ad uno spazio metrico), si costruisce un insieme di k vettori (o campioni) A che *rappresenta* D

- 1) La posizione iniziale dei vettori di A viene scelta a caso in \mathbf{R}^d (o in un sotto-spazio opportuno)

$$A = \{ w_1, w_2, \dots, w_k \}$$

- 2) Si calcola la regione di Voronoi V_i associata a ciascun vettore w_i
- 3) Si sposta ciascun w_i nel punto medio del sotto-insieme di D che appartiene a V_i

$$w_i = \frac{1}{|\{\xi \in V_i\}|} \sum_{\{\xi \in V_i\}} \xi \quad (\xi \in D)$$

- 4) L'algoritmo termina se nessun w_i è stato mosso, altrimenti torna a 2)

Dimostrabilmente, l'algoritmo converge ad un minimo *locale* dell'errore quadratico medio:

$$MSE(D, A) = \frac{1}{|D|} \sum_{w_i \in A} \sum_{\xi \in V_i} \|w_i - \xi\|^2$$

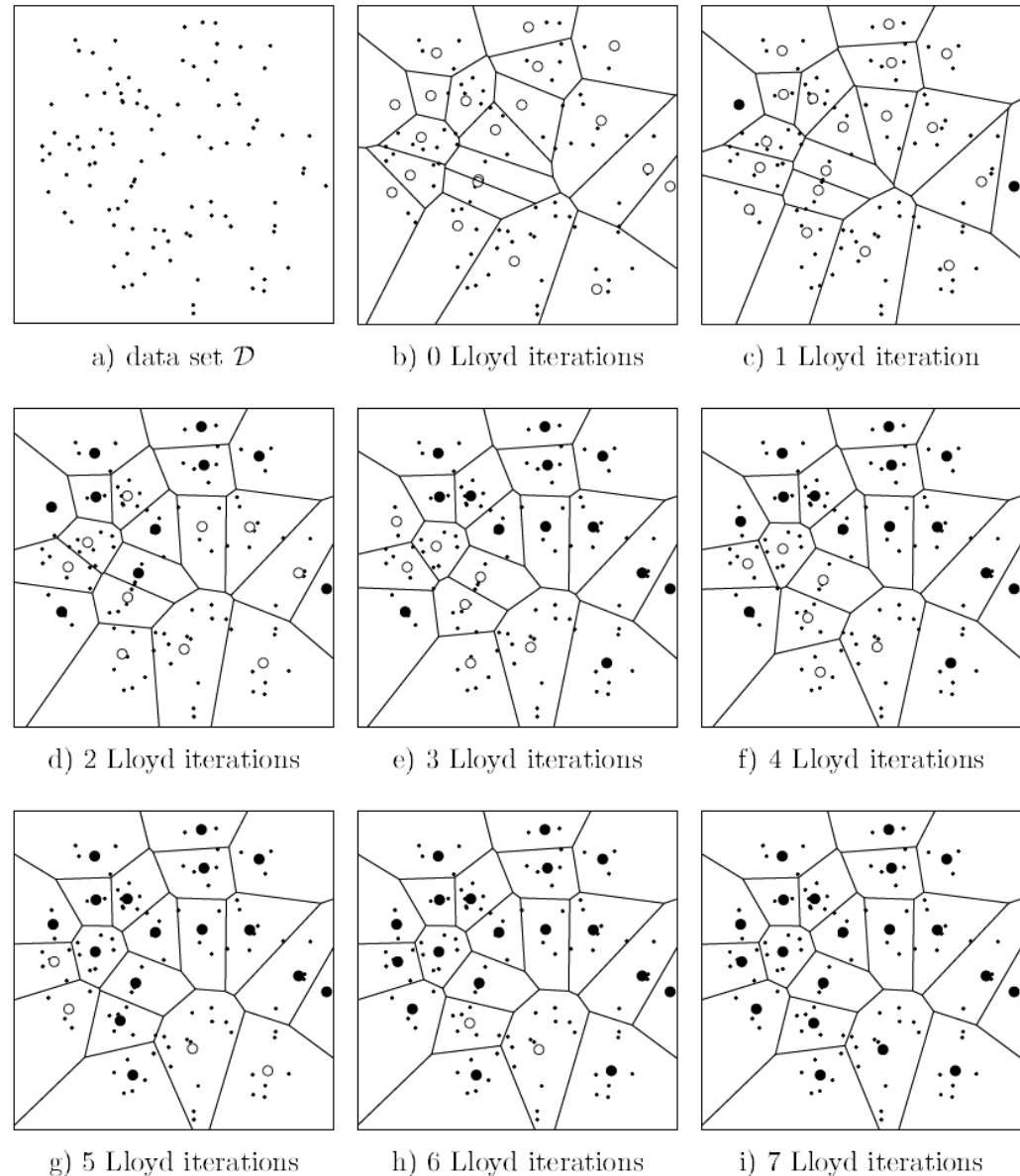
k-means

L'algoritmo di Lloyd opera in modo *batch*:
in ciascuna iterazione si considerano *tutti*
i vettori in \mathbf{A} e i dati in \mathbf{D}

L'insieme di vettori in \mathbf{A} è talvolta chiamato
codebook : se si usa \mathbf{A} per codificare \mathbf{D} ,
si ha un errore $MSE(\mathbf{D}, \mathbf{A})$

Il minimo di $MSE(\mathbf{D}, \mathbf{A})$ cui converge è locale,
quindi dipende dalla scelta iniziale di \mathbf{A}

L'algoritmo adatta la tessellazione di Voronoi
indotta da \mathbf{A} al set di dati \mathbf{D}



k-means e apprendimento probabilistico

- Il *k-means* come metodo di ottimizzazione

Ottimizzazione iterativa ed alternata della funzione obiettivo:

$$MSE(\mathbf{D}, \mathbf{A}) = \frac{1}{|\mathbf{D}|} \sum_{w_i \in \mathbf{A}} \sum_{\xi \in V_i} \|w_i - \xi\|^2 \quad \text{Mean Square Error}$$

- Si minimizza $MSE(\mathbf{D}, \mathbf{A})$ rispetto a ξ assegnando i dati al cluster del w_i più vicino
- Si minimizza $MSE(\mathbf{D}, \mathbf{A})$ rispetto ai w_i spostandoli nel punto medio del cluster

$$w_i = E(\xi \mid \xi \in V_i)$$

- Analogia: algoritmo *EM* (*Expectation-Maximization*)

Modello grafico: n cluster, ciascuno con parametri θ_i :

ciascun cluster genera i propri dati $X \in V_i$ con distribuzione $P(X \mid \theta_i)$

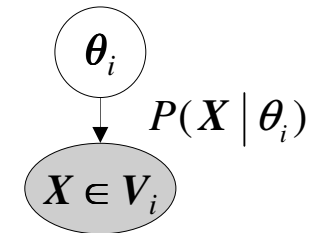
Sono disponibili solo le osservazioni $\xi \in \mathbf{D}$

Ottimizzazione iterativa ed alternata della funzione obiettivo:

$$MLE = P(\mathbf{D} \mid \theta) \quad \text{Maximum Likelihood Estimator}$$

- Calcolo della distribuzione $P(X \mid \xi, \{\theta_i\})$
- MLE* del valor medio di $L(\{\theta_i\} \mid X, \xi)$ rispetto a X , data la distribuzione $P(X \mid \xi, \{\theta_i\})$

$$\{\theta_i^*\} = \arg \max_{\{\theta_i\}} \sum_{\{X_i\}} P(X \mid \xi, \{\theta_i\}) P(X, \xi \mid \{\theta_i\})$$



Hard Competitive Learning

(anche *winner-takes-all*)

■ Un algoritmo *incrementale*

In questo caso non si ha un set di dati D , ma un *flusso* di dati in input un insieme di n vettori A che lo *rappresenta*

Si assume che i dati in input arrivino in sequenza, con probabilità $P(\xi)$

- 1) Ciascun vettore w_i viene scelto a caso in R^d (o in un sotto-spazio opportuno)

$$A = \{ w_1, w_2, \dots, w_n \}$$

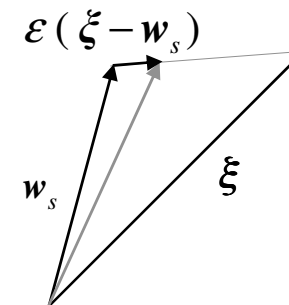
- 2) Generazione di un segnale ξ secondo $P(\xi)$

- 3) Si determina il *winner*

$$w_s = w_i \mid i = \arg \min_j \| \xi - w_j \|$$

- 4) Si adatta il *winner* al segnale

$$\Delta w_s = \varepsilon (\xi - w_s) \quad \varepsilon \text{ è il learning rate}$$

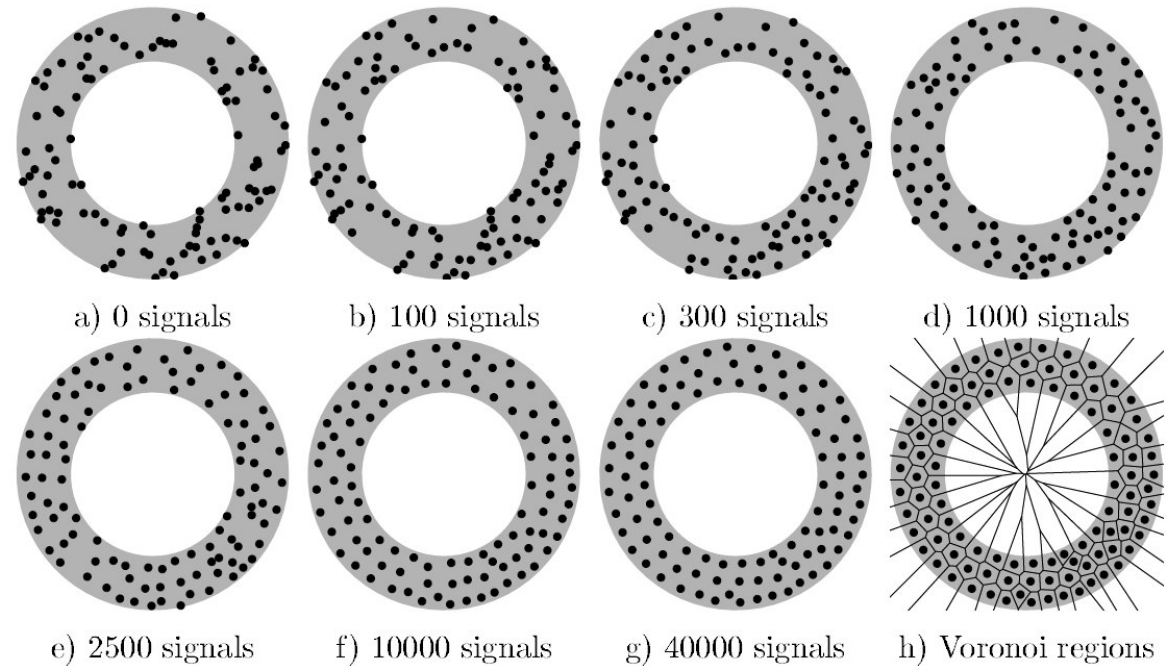


- 5) L'algoritmo termina se il flusso è esaurito, altrimenti torna a 2)

Hard Competitive Learning

A prima vista, il comportamento è simile al *k*-means

In realtà l'algorithmo NON CONVERGE: il *winner* si adatta sempre al segnale e quindi si sposta



Per forzare la convergenza, si può far 'raffreddare' (diminuire) il *learning rate* ε nel tempo

Con:

$$\varepsilon = \frac{1}{t} \quad (t \text{ numero di segnali})$$

ciascun w_i converge *asintoticamente* al valor medio di $P(\xi)$ nella regione di Voronoi

$$w_i \rightarrow E(\xi \mid \xi \in V_i)$$

Adattare vettori o reti?

- Adattamento di vettori (*clustering*)

Ai vettori di A è associata una regione (di Voronoi)

Adattare un insieme di vettori A ad un flusso di dati in input significa adattare la tessellazione complessiva

Non si ha alcuna rappresentazione delle connessioni tra regioni diverse

Le connessioni che interessano sono quelle del flusso di dati (es. modello *hidden* composto da una curva con rumore)

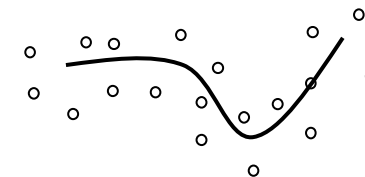
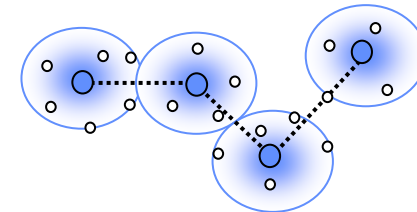
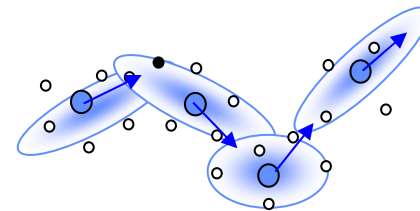
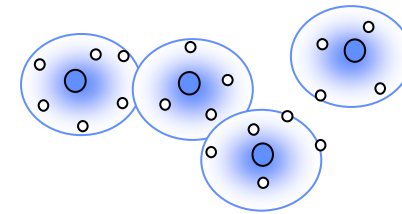
- Struttura a rete

Formata da un'insieme di unità (nodi)

A ciascuna unità è associato un vettore

Le unità possono essere connesse

Adattare una rete significa adattare la tessellazione (indotta dai vettori) e le connessioni



Self-organizing maps (SOM) (Kohonen, 1985)

- Struttura (tipica) a due livelli

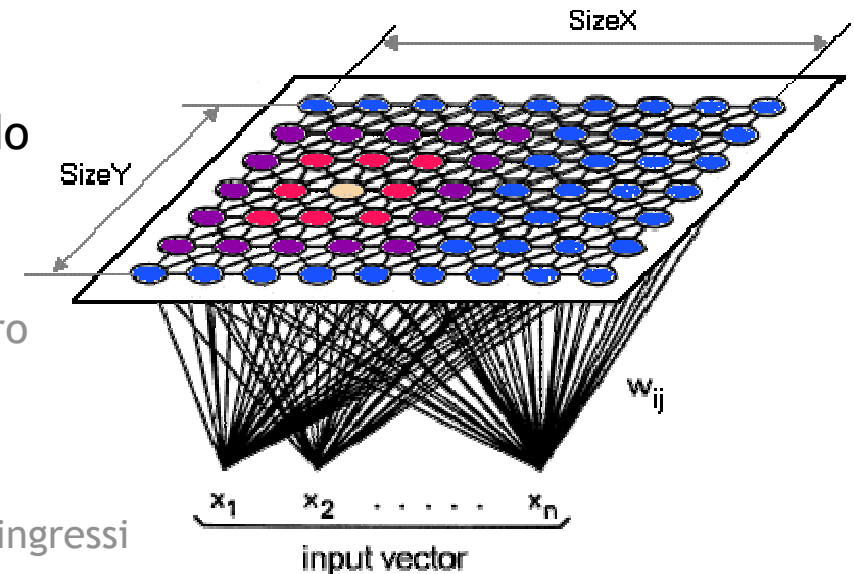
Livello di mappa, con unità organizzate secondo una topologia prestabilita

A ciascuna unità c_i è associato un vettore w_i

Le unità del livello di mappa sono connesse tra loro

Livello di input, o ingressi

Il livello di mappa è completamente connesso agli ingressi (non sono queste le connessioni che ci interessano)



La struttura delle connessioni tra le unità al livello di mappa è prestabilita e non varia nel corso dell'adattamento

Ad esempio, potrebbe trattarsi di un reticolo a base quadrata

Self-organizing maps (SOM)

■ Adattamento incrementale

Un insieme di unità $A = \{ c_1, c_2, \dots, c_n \}$ a ciascuna delle quali è associato un vettore $w_i \in R^d$

Un insieme di connessioni $C : A \times A$ prefissato (p.es. un reticolo quadrato)

Un flusso di dati $\xi \in R^d$ in input, con probabilità $P(\xi)$

- 1) Ciascun vettore w_i viene scelto a caso in R^d
- 2) Generazione di un segnale ξ secondo $P(\xi)$
- 3) Si determina il *winner*

$$c_s = c_i \mid i = \arg \min_j \|\xi - w_j\|$$

- 4) Si adattano tutti i vettori di A al segnale

$$\Delta w_i = \varepsilon(t) h(i, s) (\xi - w_i)$$

$\varepsilon(t)$ è un *learning rate* che varia nel tempo

$h(i, s)$ è una funzione che varia con la *distanza nella rete* tra c_i e c_s

- 5) L'algoritmo termina se il flusso è esaurito, altrimenti torna a 2)

Self-organizing maps (SOM)

- Scelta delle funzioni

Learning rate variabile

$$\varepsilon(t) = \varepsilon_{\max} \left(\varepsilon_{\min} / \varepsilon_{\max} \right)^{t/t_{\max}}$$

$\varepsilon(t)$ decresce da un valore massimo ad uno minimo

Adattamento delle unità

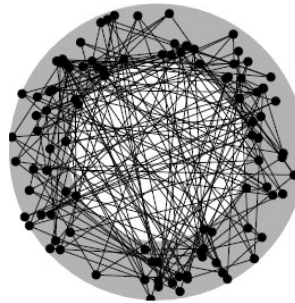
Una gaussiana che definisce un *intorno* del *winner*

$$h(i, s) = \exp\left(\frac{d(c_i, c_s)}{2\sigma(t)^2}\right)$$

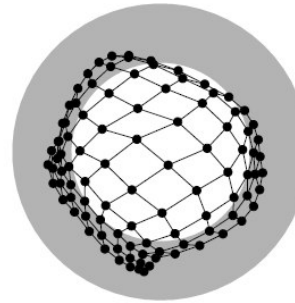
avente una varianza $\sigma(t)$ decresce nel tempo: $\sigma(t) = \sigma_{\max} \left(\sigma_{\min} / \sigma_{\max} \right)^{t/t_{\max}}$

dove $d(c_i, c_s)$ è la distanza tra le unità nel reticolo (p.es. una *manhattan distance*)

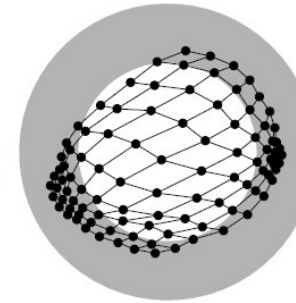
In pratica, l'*intorno* del *winner* si restringe nel tempo



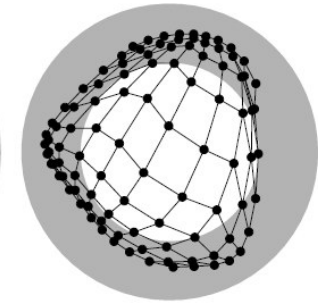
a) 0 signals



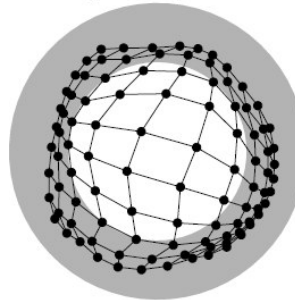
b) 100 signals



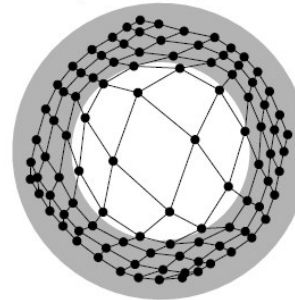
c) 300 signals



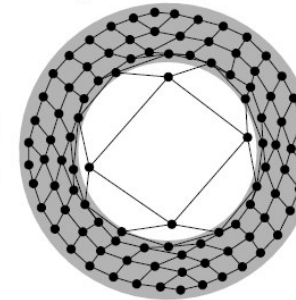
d) 1000 signals



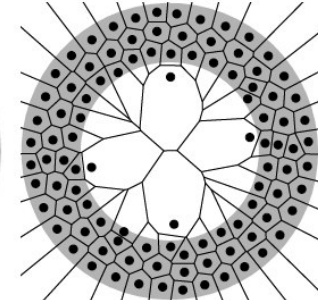
e) 2500 signals



f) 10000 signals



g) 40000 signals



h) Voronoi regions

Hebbian learning e grafo di Delaunay (Martinetz e Schulten, 1993)

- Un algoritmo incrementale per costruire il grafo di Delaunay (o un suo sottoinsieme)

Un insieme di unità $A = \{ c_1, c_2, \dots, c_n \}$ a ciascuna delle quali è associato un vettore $w_i \in R^d$

Un insieme di connessioni C inizialmente vuoto

Un flusso di dati $\xi \in R^d$ in input, con probabilità $P(\xi)$

1) I due vettori w_i vengono scelti a caso in R^d

2) Generazione di un segnale ξ secondo $P(\xi)$

3) Si determina il *winner* ed il *second winner*

$$c_s = c_i \mid i = \arg \min_j \|\xi - w_j\| \quad c_r = c_i \mid i = \arg \min_{j \neq s} \|\xi - w_j\|$$

3) Si aggiunge a C la connessione (c_s, c_r) , se non la contiene già

4) L'algoritmo termina se il flusso è esaurito, altrimenti torna a 2)

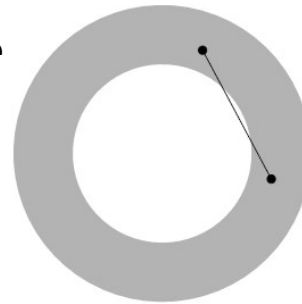
Growing Neural Gas (Fritzke, 1995)

Un algoritmo incrementale per adattare unità e connessioni

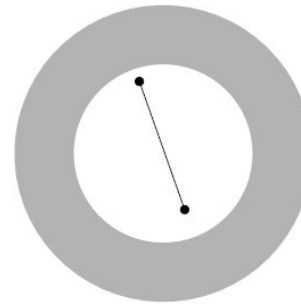
Utilizza lo *Hebbian learning* per le connessioni

Introduce nuove unità in base all'*errore accumulato*

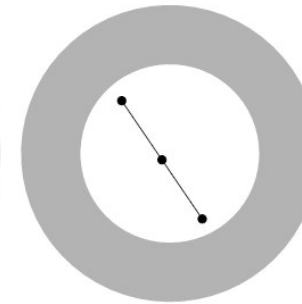
Utilizza un meccanismo di *aging* per eliminare connessioni e unità non più necessarie



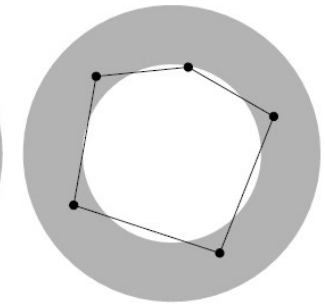
a) 0 signals



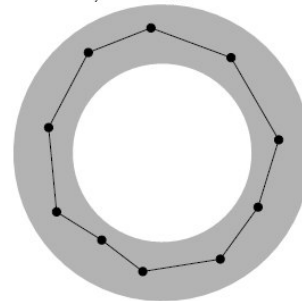
b) 100 signals



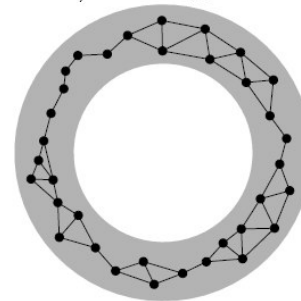
c) 300 signals



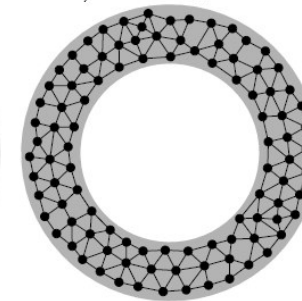
d) 1000 signals



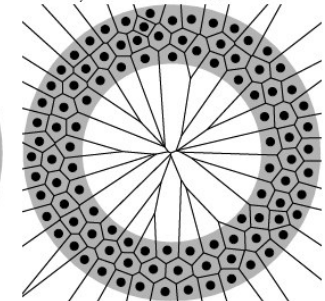
e) 2500 signals



f) 10000 signals



g) 40000 signals



h) Voronoi regions

Growing Neural Gas (GNG)

- Un algoritmo incrementale per adattare unità e connessioni

Un insieme che inizialmente contiene solo due unità $A = \{ c_1, c_2 \}$
a ciascuna delle quali è associato un vettore $w_i \in R^d$

Un insieme di connessioni C inizialmente vuoto

Un flusso di dati $\xi \in R^d$ in input, con probabilità $P(\xi)$

1) I due vettori w_1 e w_2 vengono scelti a caso in R^d

2) Generazione di un segnale ξ secondo $P(\xi)$

3) Si determina il *winner* ed il *second winner*

$$c_s = c_i \mid i = \arg \min_j \|\xi - w_j\| \quad c_r = c_i \mid i = \arg \min_{j \neq s} \|\xi - w_j\|$$

4) Si aggiunge a C la connessione (c_s, c_r) , se C non la contiene già
Alla connessione (c_s, c_r) viene assegnata $age = 0$

5) Si incrementa l'errore accumulato dall'unità *winner*

$$\Delta E_s = \|\xi - w_s\|^2$$

6) Si adattano i vettori dell'unità *winner* e delle unità direttamente connesse

$$\Delta w_s = \varepsilon_b (\xi - w_s) \quad \Delta w_i = \varepsilon_n (\xi - w_i), \quad i \mid (c_i, c_s) \in C$$

(continua)

Growing Neural Gas

7) Si incrementa il valore di *age* di tutte le connessioni dell'unità *winner*

$$age_{(c_i, c_s)} = age_{(c_i, c_s)} + 1, \quad i | (c_i, c_s) \in C$$

8) Si rimuovono tutte le connessioni con valore di *age* maggiore di una soglia. Vengono rimosse anche le unità che rimangono prive di connessioni.

9) Se il numero di segnali generati è multiplo di un valore prestabilito, si inserisce una nuova unità

9.1) Si individua l'unità c_q con il massimo errore accumulato E_q e l'unità connessa c_f con il valore di E più elevato

9.2) Si rimuove da C la connessione (c_q, c_f)

9.3) Si aggiunge ad A una nuova unità c_r con un vettore associato che è la media dei due

$$\mathbf{w}_r = (\mathbf{w}_q + \mathbf{w}_f) / 2$$

9.4) Si aggiungono a C le connessioni (c_q, c_r) e (c_r, c_f) , con $age = 0$

9.5) Si ripartiscono i valori di E tra c_q , c_f e c_r

10) Si decrementano tutti i valori di errore accumulato

$$\Delta E_i = -\beta E_i$$

11) L'algoritmo termina se il flusso è esaurito, altrimenti torna a 1)

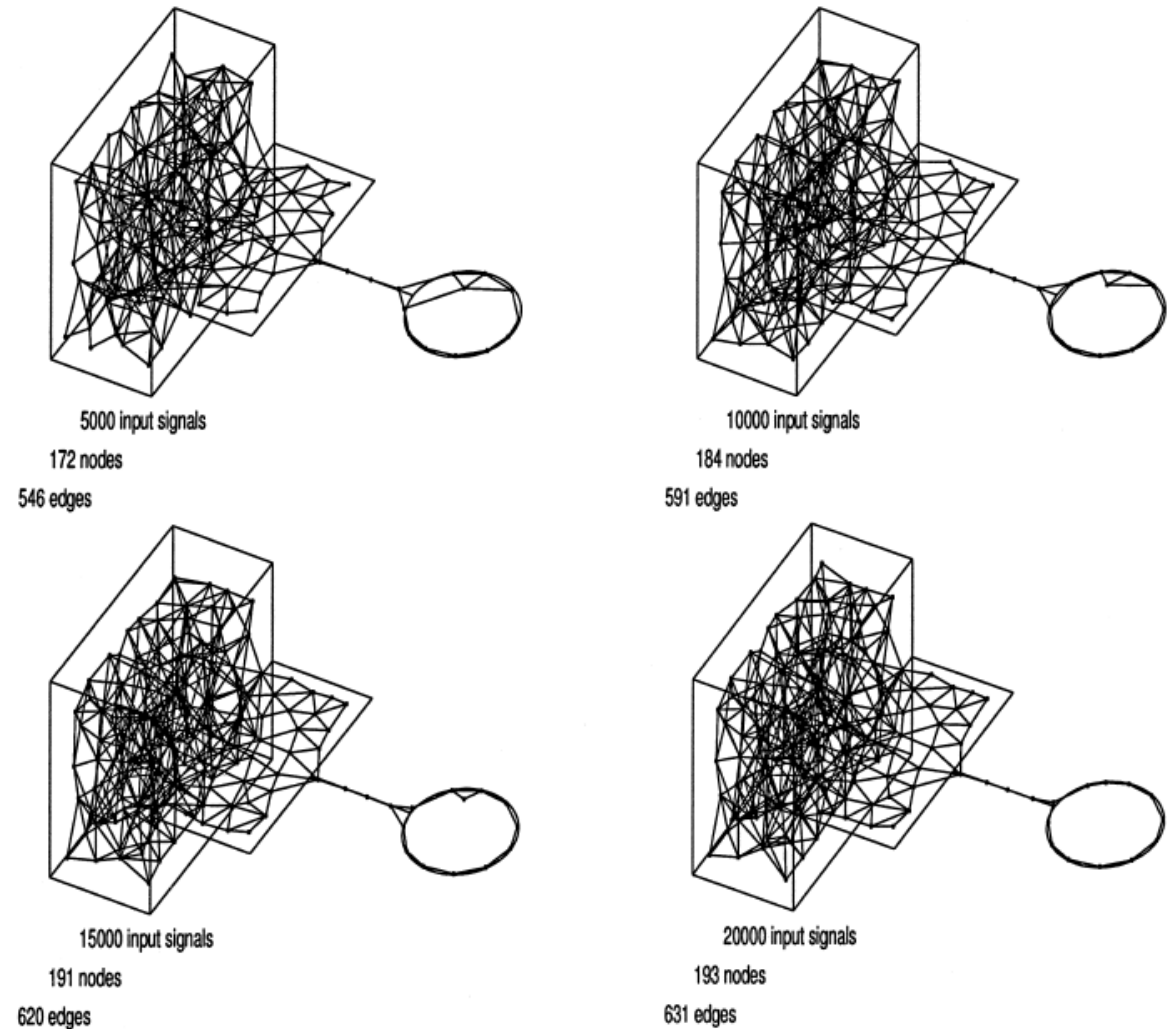
Adattamento e topologia variabile

Rispetto alle *SOM* (topologia fissa) i *GNNG* possono adattare la propria topologia a quella del flusso di dati

Qui a destra: il flusso di dati proviene da una struttura composta da parti a dimensione intrinseca diversa

La rete si adatta e l'organizzazione risultante rappresenta la topologia della struttura originaria

Il calcolo è distribuito: ogni unità 'vede' soltanto i propri vicini



Grow-when-required networks (Marsland, 2002)

- Un limite dei *Growing Neural Gas*

La rete cresce aggiungendo un'unità a periodi stabiliti

La crescita non ha limiti, salvo ovviamente un massimo prestabilito

Ciò rappresenta un'ulteriore elemento di perturbazione

- *Grow-when-required networks* (GWR)

Molto simile ai *GNG*

Due differenze principali:

a) Le unità hanno un valore di *habituation* (assuefazione) che decresce esponenzialmente, con costante τ , tutte le volte che l'unità è *winner*

Sotto una certa soglia di *habituation*, l'unità riduce la 'propensione all'adattamento'

b) Sotto la soglia di *habituation*, ciascuna unità ha un raggio prestabilito di *activity* (che definisce una *d*-palla in \mathbf{R}^d)

I segnali in input per cui l'unità è *winner* vengono accettati solo se cadono all'interno del raggio di *activity*

In caso contrario, si crea una nuova unità