

Sistemi Intelligenti Learning and Clustering

Alberto Borghese and Iuri Frosio

Università degli Studi di Milano
Laboratorio di Sistemi Intelligenti Applicati (AIS-Lab)
Dipartimento di Scienze dell'Informazione
alberto.borghese@unimi.it



A.A. 2012-2013

1/35

<http://homes.dsi.unimi.it/~borghese/>



Riassunto



- I tipi di apprendimento
- Il clustering

A.A. 2012-2013

2/35

<http://homes.dsi.unimi.it/~borghese/>



I vari tipi di apprendimento

$$\begin{aligned}x(t+1) &= f[x(t), a(t)] && \text{Ambiente} \\ a(t) &= g[x(t)] && \text{Agente}\end{aligned}$$

Supervisionato (learning with a teacher). Viene specificato per ogni pattern di input, il pattern desiderato in output.

Semi-Supervisionato. Viene specificato solamente per **alcuni** pattern di input, il pattern desiderato in output.

Non-supervisionato (learning without a teacher). Estrazione di similitudine statistiche tra pattern di input. Clustering. Mappe neurali.

Apprendimento con rinforzo (reinforcement learning, learning with a critic). L'ambiente fornisce un'informazione puntuale, di tipo qualitativo, ad esempio success or fail.

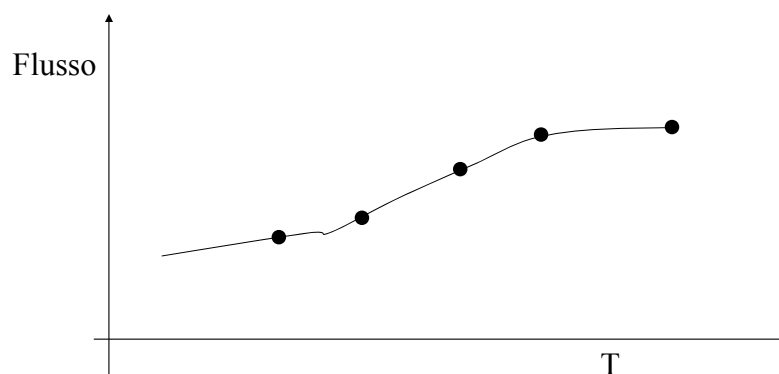
A.A. 2012-2013

3/35

<http://homes.dsi.unimi.it/~borghese/>



Apprendimento supervisionato: regressione = predictive learning



Controllo della portata di un condizionatore in funzione della temperatura. "Imparo" una funzione continua a partire da alcuni campioni: devo imparare ad **interpolare**.

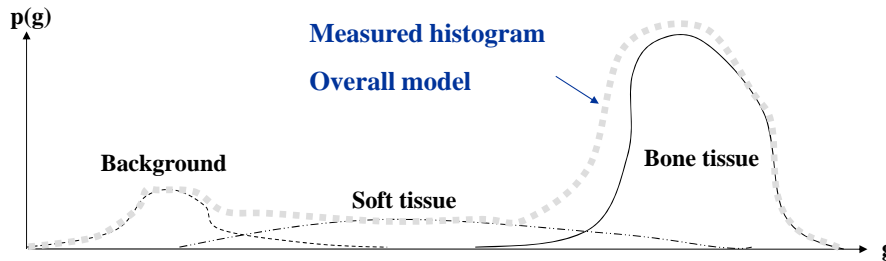
A.A. 2012-2013

4/35

<http://homes.dsi.unimi.it/~borghese/>



I modelli parametrici



$$p(g) = \sum_{j=1}^M P(j) \cdot p(g | j) = \sum_{j=1}^M w_j \cdot p_j(g)$$

La probabilità di avere un livello di grigio g è la somma pesata delle tre probabilità di avere background, $p_1(g)$, tessuto molle, $p_2(g)$ o tessuto osseo, $p_3(g)$.

A.A. 2012-2013

5/35

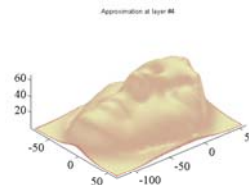
<http://homes.dsi.unimi.it/~borghese/>



I modelli semi-parametrici

- L'approssimazione è ottenuta mediante funzioni “generiche”, dette di **base**, soluzione molto utilizzata nelle NN e in Machine learning.
- (Il concetto di Base in matematica è definito mediante certe proprietà di approssimazione che qui non consideriamo, consideriamo solo l'idea intuitiva).

$$z(p(x, y)) = \sum_i w_i G(p, p_i; \sigma)$$



Combinazione
lineare di funzioni
di base

Da calcolare

Funzione di base (fissate)

A.A. 2012-2013

6/35

<http://homes.dsi.unimi.it/~borghese/>



Modelli lineari e non lineari



Classificazione alternativa dei modelli. Vengono utilizzate classi molto diversi di algoritmi per stimare i parametri di questi due tipi di modelli.

$$z(p(x, y)) = f(x) = \sum_i w_i x \qquad z(p(x, y)) = \sum_i f_i(G(p, p_i; \sigma))$$

$f(.) = w_i$ è funzione lineare

$f(.)$ è funzione non lineare

e.g. $f(.) = e^{G(.)}$

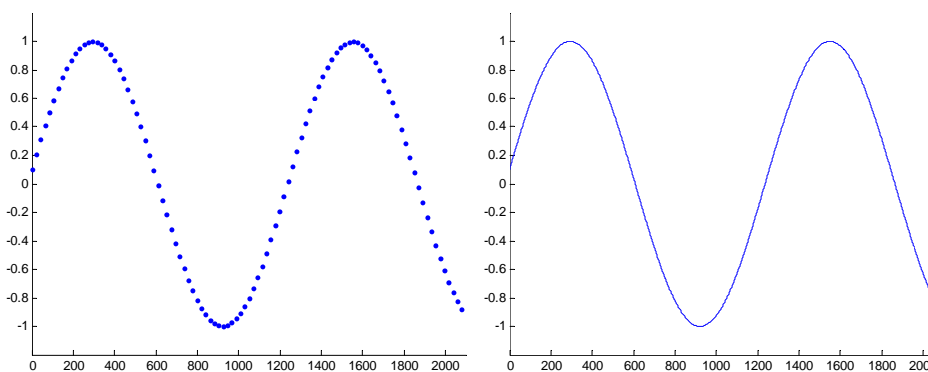
A.A. 2012-2013

7/35

<http://homes.dsi.unimi.it/~borghese/>



Funzionamento di un modello parametrico (non-lineare)



I punti vengono fittati perfettamente da una sinusoida: $y = A \sin(\omega x + \phi)$.
Devo determinare i parametri della sinusoida (non lineare), i cui valori ottimali sono: $\omega = 1/200$, $\phi = 0.1$.

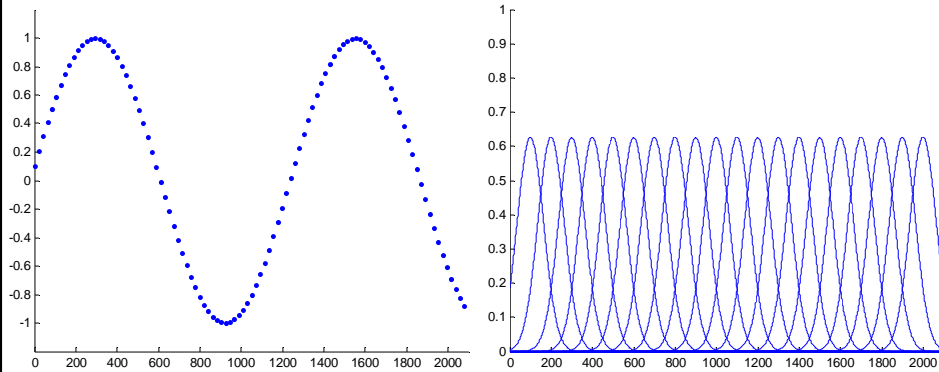
A.A. 2012-2013

8/35

<http://homes.dsi.unimi.it/~borghese/>



Approssimazione mediante un modello semi-parametrico (lineare)



Vogliamo fittare i punti con l'insieme di Gaussiane riportate sulla dx. In questo caso hanno tutte $\sigma = 90$. Come le utilizzo?

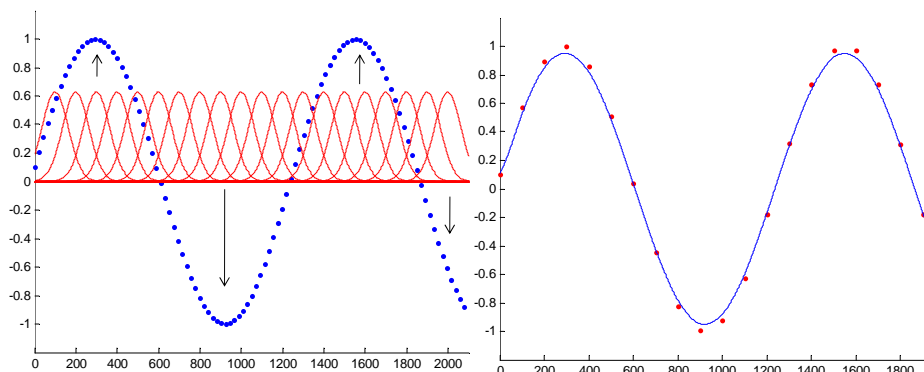
A.A. 2012-2013

9/35

<http://homes.dsi.unimi.it/~borghese/>



Funzionamento di un modello semi-parametrico (lineare)



$$y(x) = \sum_{i=1}^{20} w_i G(x - x_{o_i}; 90)$$

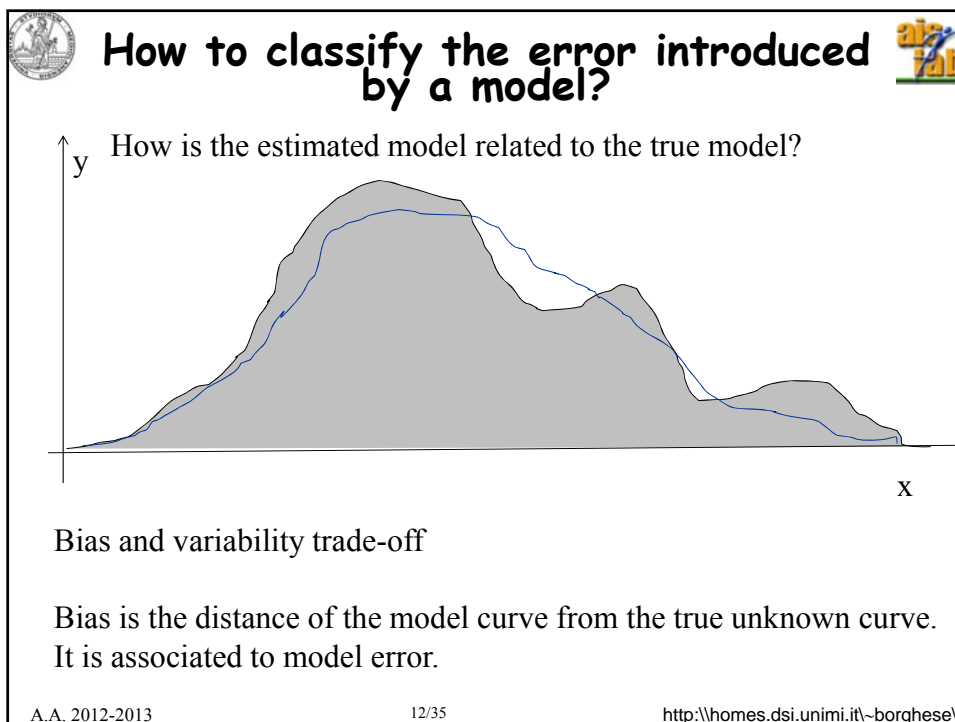
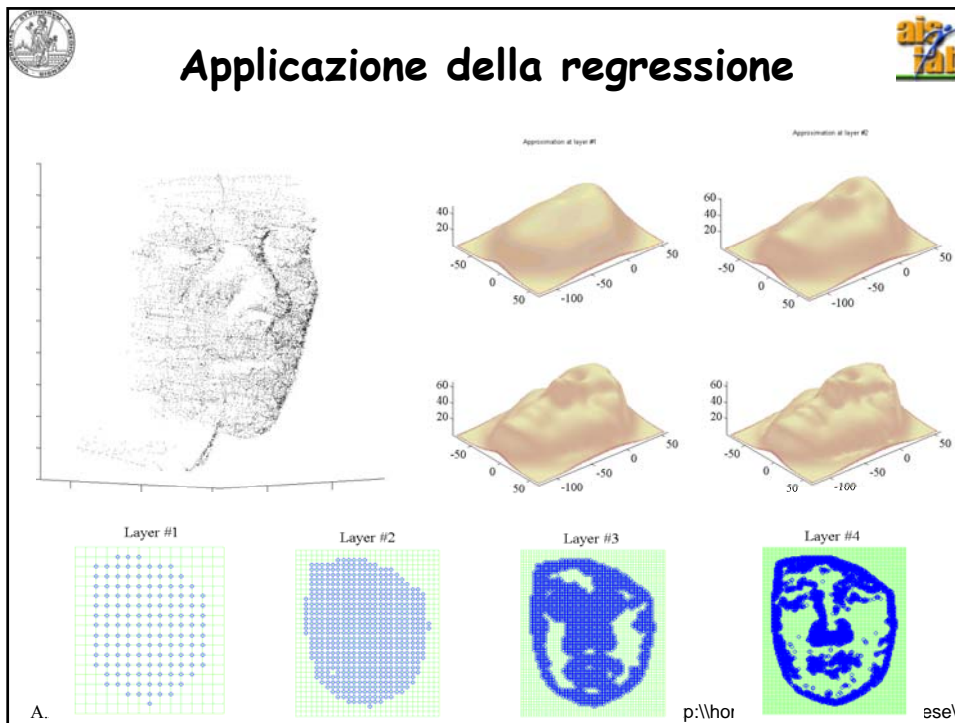
Devo definire, i $\{w_i\}$

I σ sono tutti uguali ed uguali a 90, le Gaussiane sono equispaziate.
Le Gaussiane sono note tutte a priori, devono essere definiti i pesi.

A.A. 2012-2013

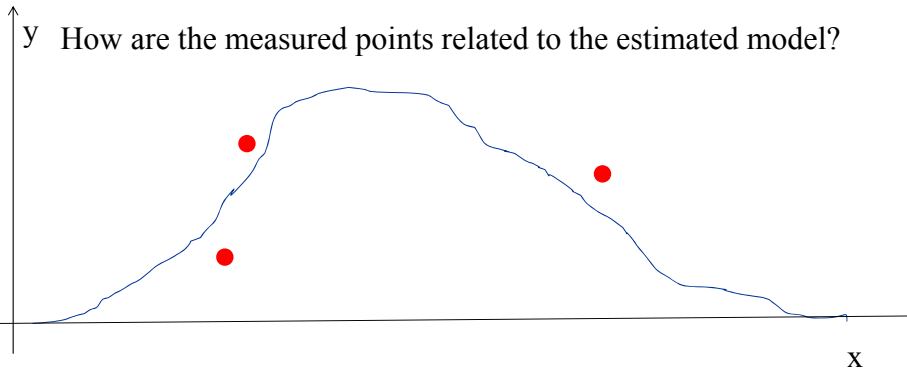
10/35

<http://homes.dsi.unimi.it/~borghese/>





Variability



Given $P_{mes}(x,y)$ and $y = y(x)$, the error is measured as:

- $\min \| P - y(x) \|$
- $\| y_{mes}(x_{mes}) - y(x_{mes}) \|$

It is associated to measurement error.

If variability goes to zero, bias increases and overfitting arises.

A.A. 2012-2013

13/35

<http://homes.dsi.unimi.it/~borgnese/>



Problemi

Quando si termina l'algoritmo di apprendimento?

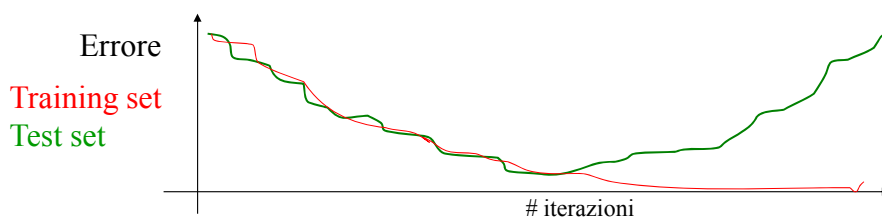
Bootstrap – Vengono estratti pattern con ripetizioni.

Cross-Validation - Errore sull'insieme di training =

Errore sull'insieme di test.

Utilizzare lo “structural risk” invece dell’”empirical risk”.

Si vuole evitare che la rete si specializzi troppo sui pattern di training e non sia in grado di interpolare.



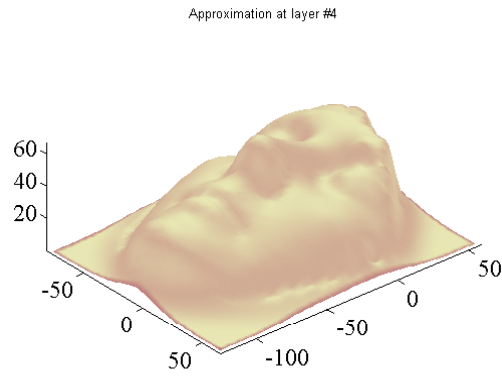
A.A. 2012-2013

14/35

<http://homes.dsi.unimi.it/~borgnese/>



Problema dell'overfitting dovuto a sovrapparametrizzazione



Quante unità?

A.A. 2012-2013

15/35

<http://homes.dsi.unimi.it/~borghese>

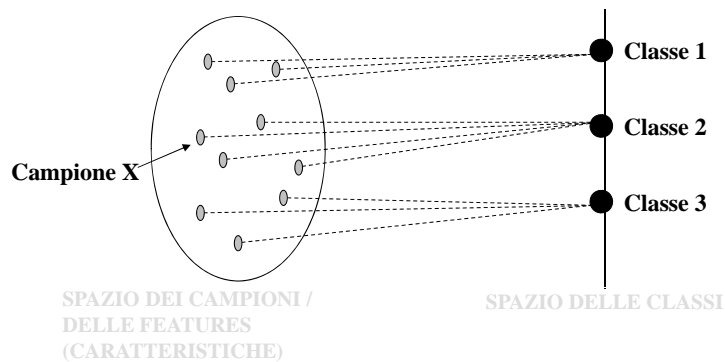


Classificazione



Un'interpretazione geometrica:



Mappatura dello spazio dei campioni nello spazio delle classi.



A.A. 2012-2013

16/35

<http://homes.dsi.unimi.it/~borghese>



Apprendimento Supervisionato:
Classificazione

b
 b
 b
 b
 b
 → B

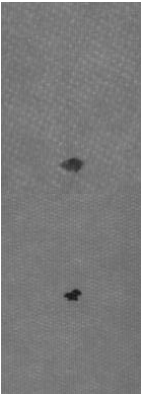
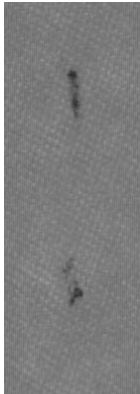
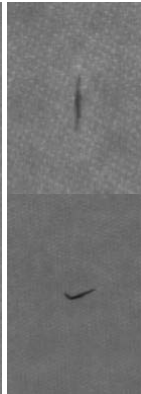
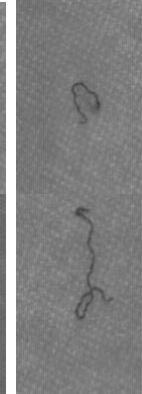
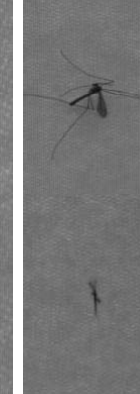
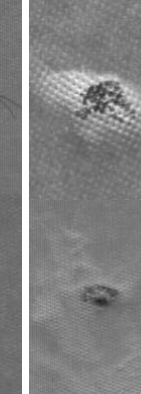
a
 a
 a
 a
 a
 → A

Task di classificazione
 Uscita intera (etichetta o label della classe)

A.A. 2012-20 <http://homes.dsi.unimi.it/~borghese/>

CLASSIFICAZIONE: Riconoscimento difetti in linee di produzione
 (progetto finanziato da Electronic Systems: 2006-2007)

					
regolari	irregolari	allungati	fili	insetti	macchie su denso

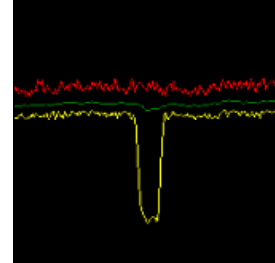
Difetti – Classificazione real-time e apprendimento mediante **boosting**.
 Committee (linear combination) of weak (binary) classifiers.

A.A. 2012-2013 <http://homes.dsi.unimi.it/~borghese/>



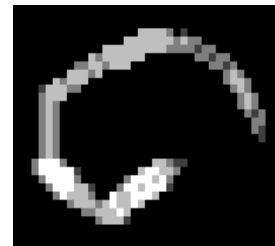
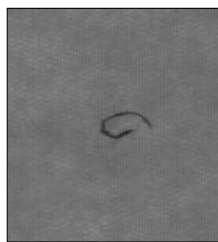
Features

Macchie
dense




- *Località.*
- *Significatività.*
- *Rinoscibilità.*

Fili




Riassunto

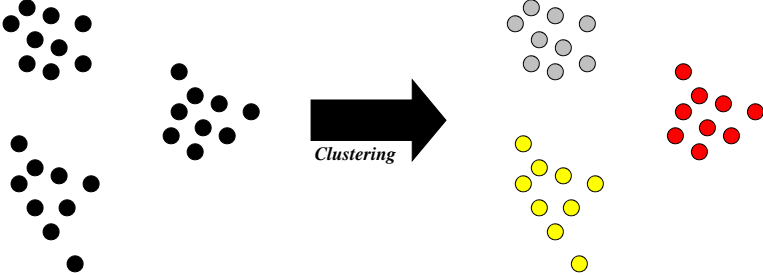
- I tipi di apprendimento
- **Il clustering**



Clustering



- Clustering: raggruppamento degli “oggetti” in classi omogenee tra loro.
 - ◆ Raggruppamento per colore
 - ◆ Raggruppamento per forme
 - ◆ Raggruppamento per tipi
 - ◆



Novel name: **data mining**

A.A. 2012-2013
21/35
<http://homes.dsi.unimi.it/~borghese/>



Esempio di clustering





Clustering -> Indicizzazione

Ricerca immagini su WEB.



A.A. 2012-2013
22/35
<http://homes.dsi.unimi.it/~borghese/>



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.”).
- ... “Ragionare” sui dati o elaborare i dati in modo stereotipato.

A.A. 2012-2013

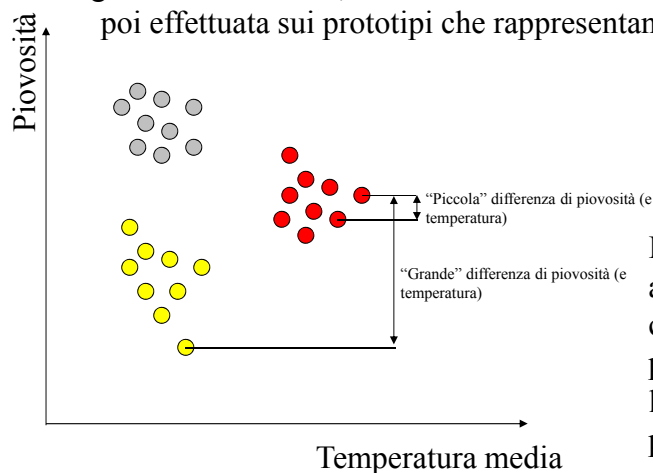
23/35



Clustering



Processo attraverso il quale i dati (pattern, vettori) vengono organizzati in cluster, basata sulla similarità. L’elaborazione verrà poi effettuata sui prototipi che rappresentano ciascun cluster.



I pattern appartenenti ad un cluster valido sono più simili l’uno con l’altro rispetto ai pattern appartenenti ad un cluster differente.

A.A. 2012-2013

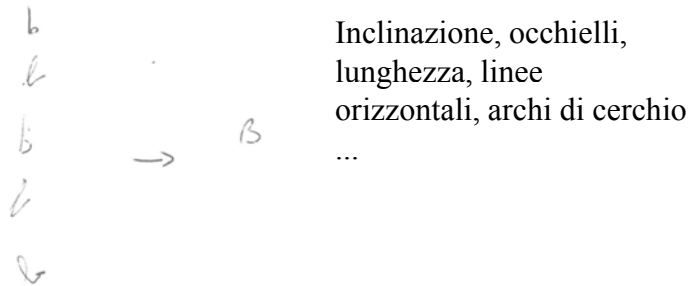
24/35



Clustering: definizioni



- **Pattern:** un singolo dato $\mathbf{p} = [p_1, p_2, \dots, p_n]$. Il dato appartiene quindi ad uno spazio multi-dimensionale, solitamente eterogeneo.
- **Feature:** le caratteristiche dei dati significative per il clustering, possono costituire anch'esso un vettore, il vettore delle feature: x_1, x_2, \dots, x_d . Questo vettore costituisce l'input agli algoritmi di clustering.



A.A. 2012-2013



Clustering: definizioni



- **n:** dimensione dello spazio dei pattern;
- **d:** dimensione dello spazio delle feature;
- **Cluster:** in generale, insieme che raggruppa dati simili tra loro, valutati in base alle feature;
- **Funzione di similarità o distanza:** una metrica (o quasi metrica) nello spazio delle feature, usata per quantificare la similarità tra due pattern.
- **Algoritmo:** scelta di come effettuare il clustering (motore di clustering).

A.A. 2012-2013

26/35

Analisi mediante clustering

Da Xu and Wunsch, 2005

I cluster ottenuti sono significativi?
 Il clustering ha operato con successo?

NB i cammini all'indietro consentono di fare la sintonizzazione dei diversi passi.

A.A. 2012-2013 27/35 <http://homes.dsi.unimi.it/~borghese/>

Il clustering

Per una buona review: Xu and Wunsch, IEEE Transactions on Neural Networks, vol. 16, no. 3, 2005.

Il clustering non è di per sé un problema ben posto. Ci sono diversi gradi di libertà da fissare su come effettuare un clustering.

- Rappresentazione dei pattern;
- Calcolo delle feature;
- Definizione di una misura di prossimità dei pattern attraverso le feature;
- Tipo di algoritmo di clustering (gerarchico o partizionale)
- Validazione dell'output (se necessario) -> Testing.

Problema a cui non risponderemo: **quanti cluster**? Soluzione teorica (criterio di Akaike), soluzione empirica (growing networks di Fritzke).

A.A. 2012-2013 28/35



Rappresentazione dei pattern



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

Esempio: descrizione del clima e della città di Roma.

Roma: [17°; 500mm; 1.500.000 ab.]

- Come scegliere i pattern?
 - ◆ Vicinanza ai bordi di ciascun cluster (Support Vector Machines)
 - ◆ Tutti i pattern
- Come valutare le feature?
 - ◆ Analisi statistica del potere discriminante: correlazione tra feature e loro significatività.



Feature & feature



- Feature extraction: trasformazione delle feature per creare nuove, significative feature;
- Elaborazione di primo livello, per ottenere informazioni caratteristiche del fenomeno che, ad esempio, siano invarianti.

Esempio: descrizione di oggetti circolari.

Posso misurare l'area e il perimetro, ma il loro rapporto è più significativo.

Esempio: descrizione del clima.

Milano: [13°; 900mm; 265 giorni sole; 100 giorni pioggia]

oppure

Milano: [13°; 900mm / 100 giorni pioggia; 265 giorni sole]



Similarità tra feature



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

Esempio:

Distanza euclidea...

dist (Roma, Milano) = ...

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

= ... Distanza euclidea? = ...

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

Ha senso?



Normalizzazione feature



Att.ne!

dist (Roma, Milano) = ...

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

= ... Distanza euclidea? = ...

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

La distanza tra le due città in termini di gradi è insignificante nel nostro calcolo... **E' necessario trovare una metrica corretta per la rappresentazione dei dati. Ad esempio, normalizzare le feature!**

$T_{Max} = 20^\circ$ $T_{Min} = 5^\circ \rightarrow T_{Norm} = (T - T_{Min}) / (T_{Max} - T_{Min})$

$P_{Max} = 1000\text{mm}$ $P_{Min} = 0\text{mm} \rightarrow P_{Norm} = (P - P_{Min}) / (P_{Max} - P_{Min})$

$Roma_{Norm} = [0.8 \ 0.5]$

$Milano_{Norm} = [0.53 \ 0.9]$

dist($Roma_{Norm}$, $Milano_{Norm}$) = $((0.8-0.53)^2+(0.5-0.9)^2)^{1/2} = 0.4826$



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:
 $\text{dist}(x,y)=(x_k-y_k)S^{-1}(x_k-y_k)$, con S matrice di covarianza.
- Context dependent:
 $\text{dist}(x,y)=f(x, y, \text{context})$



Tassonomia (sintetica) degli algoritmi di clustering



- Algoritmi gerarchici (agglomerativi, divisivi), e.g. **Hierarchical clustering**.
- Algoritmi partizionali, hard: **K-means, quad-tree decomposition**.
- Algoritmi partizionali, soft: fuzzy c-mean, neural-gas, enhanced vector quantization, **mappe di Kohonen**.
- Algoritmi statistici: **mixture models**.



Riassunto



- I tipi di apprendimento
- Il clustering