

Sistemi Intelligenti On-line learning

Alberto Borghese
Università degli Studi di Milano
Laboratorio di Sistemi Intelligenti Applicati (AIS-Lab)
Dipartimento di Informatica
Alberto.borghese@unimi.it



A.A. 2023-2024

1/43

<http://borghese.di.unimi.it/>



Riassunto



- **Valutazione di un modello**
- Modelli multi-scala on-line
- Introduzione alle reti neurali

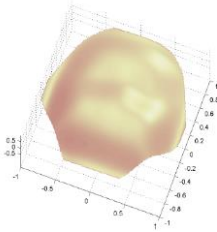
A.A. 2023-2024

2/43

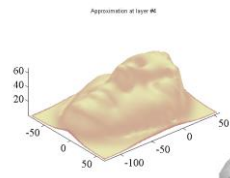
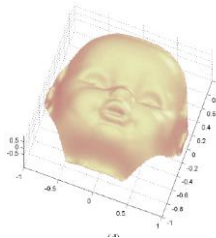
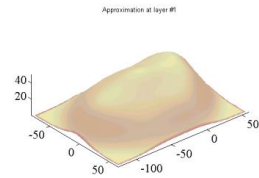
<http://borghese.di.unimi.it/>



Underfitting e overfitting



Quanti parametri?



Quante unità?

Come valutiamo?



How to classify the error introduced by a model?



Does it cover the input domain (in all dimensions – **dimensionality discovery**)?

This is not enough to obtain a good model!!! Is the model good enough?


Is a model that produces 0 error on the data a good model?

Total quality is an inspection system that detects all the defects. No defected objects are left behind.


The model should be properly tuned to the data. Errors can be classified in:

- Bias
- Variability

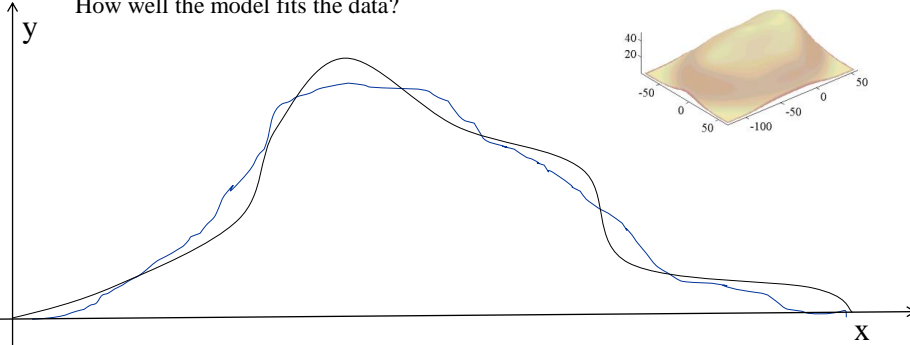
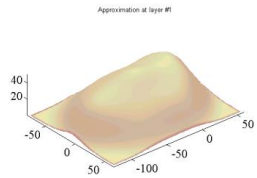
Relationship between number of parameters and bias/variability



Bias



How well the model fits the data?


Blue represents the curve of the real data $\{x_{true}, y_{true}\}$
 Black represents the curve produce by the model: $y = f(x)$

How good is the model?


Model output

We can compute the error: $dist(y_{true}, y=f(x_{true}))$, for instance Euclidean distance.
 As such it is the **bias of the model**.

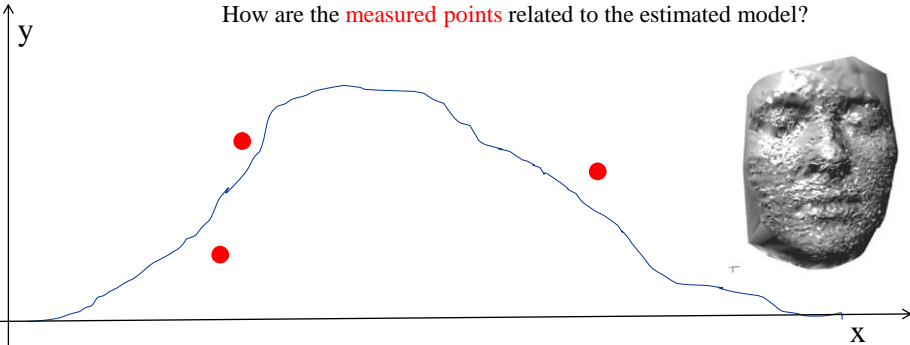

A.A. 2023-2024
5/43
<http://borghese.di.unimi.it/>



Variability



How are the **measured points** related to the estimated model?

Given $P_m(x_m, y_m)$ and $y = f^*(x)$, the true data behavior, the error is:
 $dist(y_m, y=f^*(x_m))$, for instance Euclidean distance.
 As such **variability** is the **measurement error**.

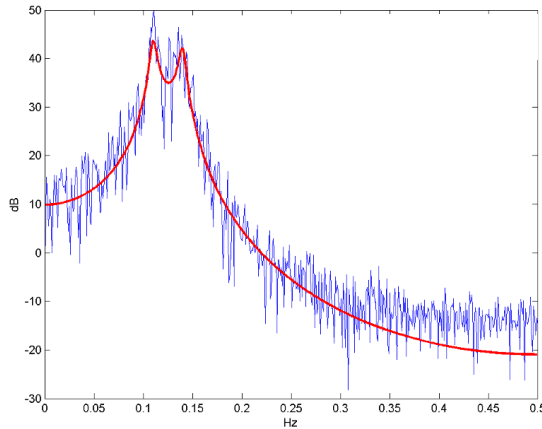
If variability goes to zero, bias increases and overfitting arises (model fits data and the noise too).
In a good model, variability tends to the statistics of the measurement noise
 (cf. regularization parameter setting).

A.A. 2023-2024
6/43
<http://borghese.di.unimi.it/>



Variability

How are the measured points related to the estimated model?



Given $P_m(x_m, y_{mes})$ and $y = f(x)$, the true data behavior, the error is measured as: $dist(y_m, f(x_m))$, for instance Euclidean distance. It is associated to measurement error.

If variability goes to zero, bias increases and overfitting arises.
In a good model, variability tends to the statistics of the data noise.

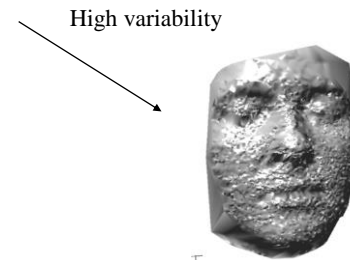
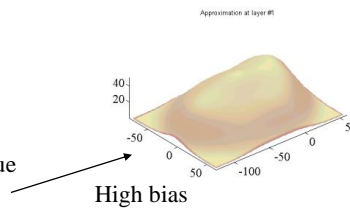


Bias and variability

Bias and variability trade-off

Bias is the distance of the model curve from the true curve, **that is unknown**. It is the model error.

Variability is the distance of the true curve, **that is unknown**, from the **measured data**. It is the measurement error.





Parametric model Error vs number of parameters



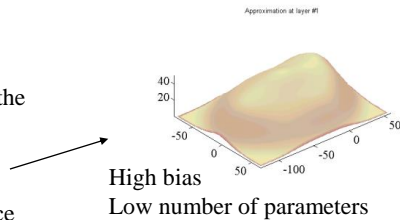
Bias and variability trade-off

Bias is the distance of the model curve from the true curve, **that is unknown**. It is the model error.

If we have **few parameters**, we can reproduce only the outline of the data (under-fitting -> bias).

Variability is the distance of the true curve, **that is unknown**, from the **measured data**. It is the measurement error.

If we have many parameters we can reproduce the fastest variations, that are due to noise (over-fitting -> variability).



High variability
High number of parameters



Cross validation



- I dati vengono suddivisi in due sotto-insiemi: training (costruzione del modello) e test.
- Errore sull'insieme di training = Errore sull'insieme di test.
- L'insieme di test è rappresentativo dei dati che potranno essere presentati in futuro al modello.

Si vuole evitare che il modello si specializzi troppo sui pattern di training e non sia in grado di interpolare correttamente su altri dati (e.g. dati di test).

- La procedura viene ripetuta k volte su sottoinsiemi diversi estratti a caso: k-fold cross-validation.

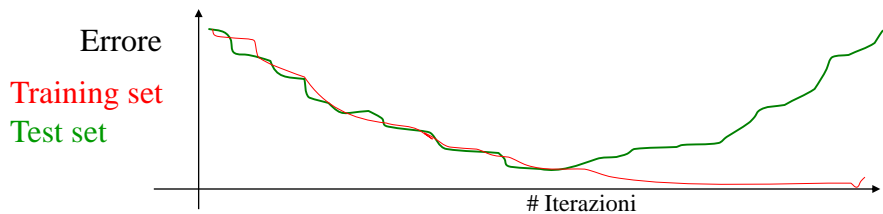


Semi-convergenza

Ipotesi: ho un insieme di parametri sufficientemente ampio.

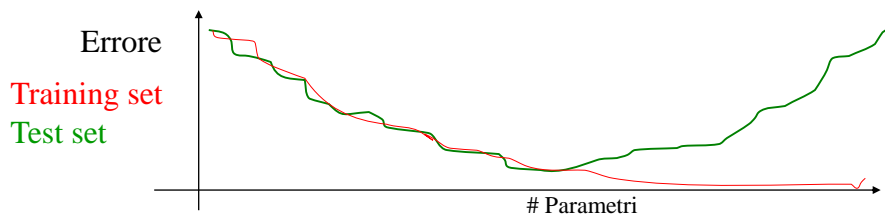
Semi-convergenza: non porto l'algoritmo fino alla convergenza nel punto di ottimo ma arresto le iterazioni prima ("early stopping")

Il modello non sarà perfettamente aderente ai dati, ma il residuo sarà tendenzialmente l'errore di misura.



Scelta empirica del numero di parametri

Il numero di parametri viene aumentato fino a quando **entrambi** gli errori diminuiscono.



Parto da un modello con pochi parametri e aumento il numero di parametri o viceversa.

Approccio simile a quello utilizzato nel criterio di discrepanza.

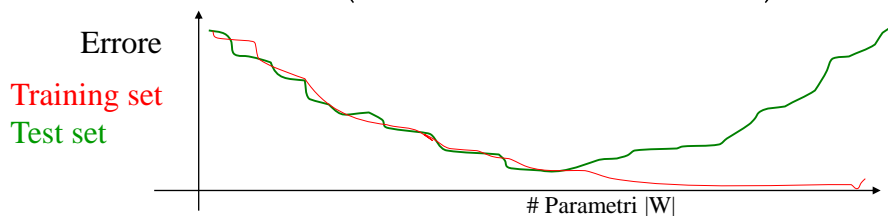


Scelta teorica

Quale funzione costo minimizzo? Come posso inserire l'informazione di complessità del modello nella funzione costo?

Penalizzo i modelli con tanti parametri. Aggiungo nel calcolo della distanza tra dati e modello (funzione costo) un termine che cresce con il numero dei parametri -> Regularization with Reproducible Hilbert Kernels as regularizers.

$$w = \underset{w}{\operatorname{argmin}} \left(\sum_i \|f(x_m) - y_m\|^2 + \lambda g(|W|) \right)$$



A.A. 2023-2024

13/43

<http://borghese.di.unimi.it/>



I vari tipi di apprendimento

$$x(t+1) = f[x(t), a(t)] \quad \text{Ambiente}$$

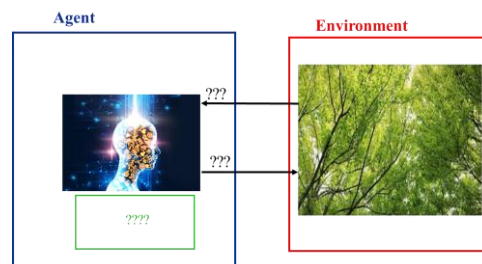
$$a(t) = g[x(t)] \quad \text{Agente}$$

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

Semi-Supervisionato. Il pattern desiderato in output viene specificato solamente per **alcuni** pattern di input. Gli altri pattern contribuiscono a definire la forma del manifold dei dati.

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

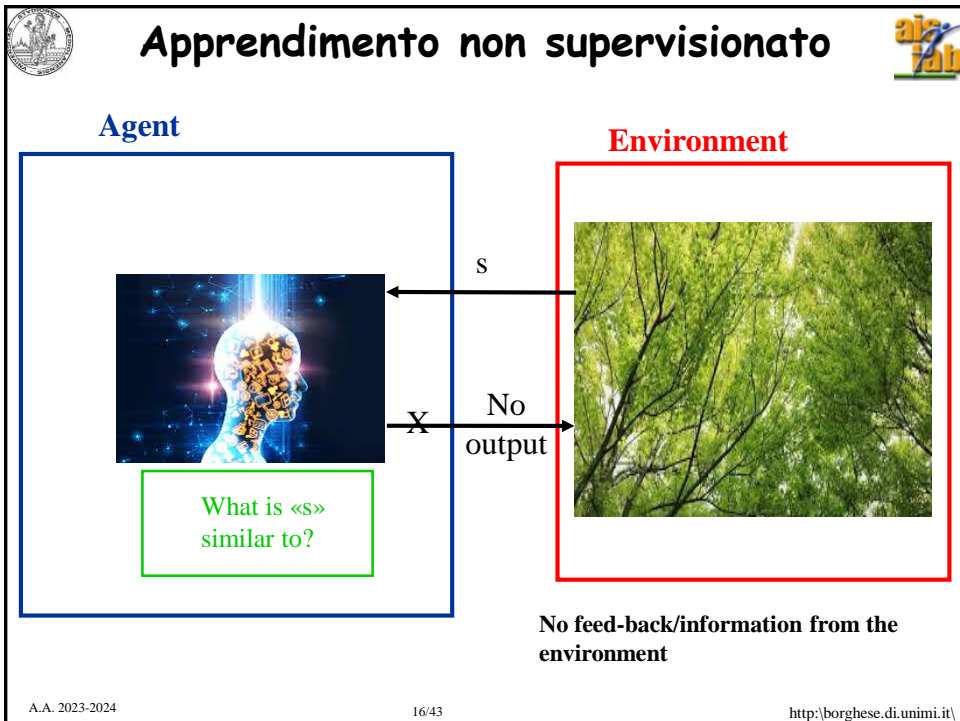
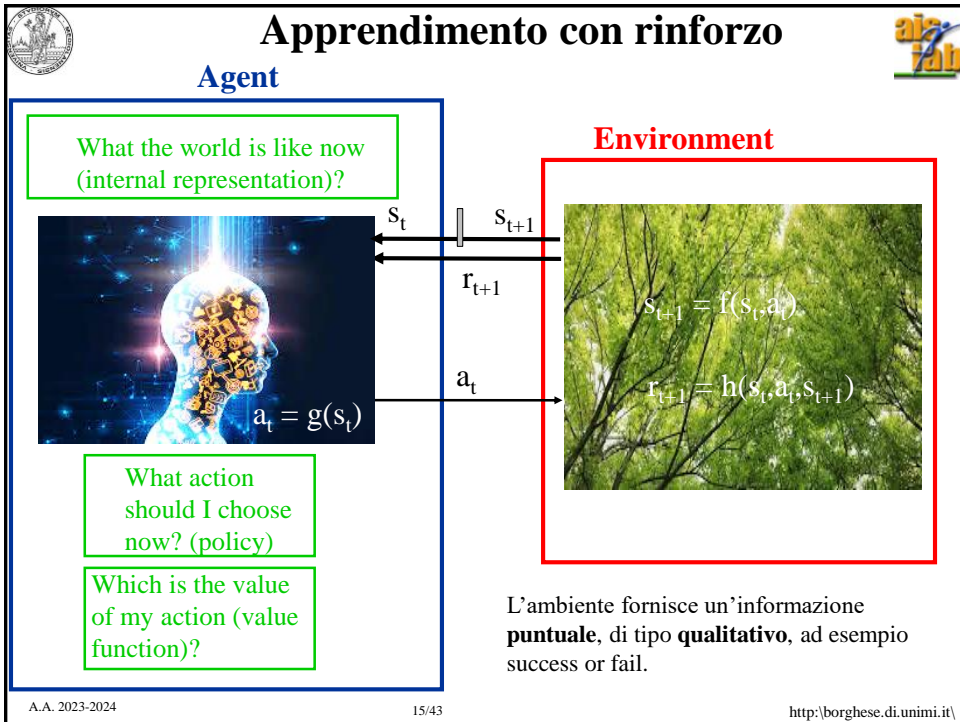
Apprendimento con rinforzo (reinforcement learning, learning with a critic). L'ambiente fornisce un'informazione puntuale, di tipo qualitativo, ad esempio success or fail.

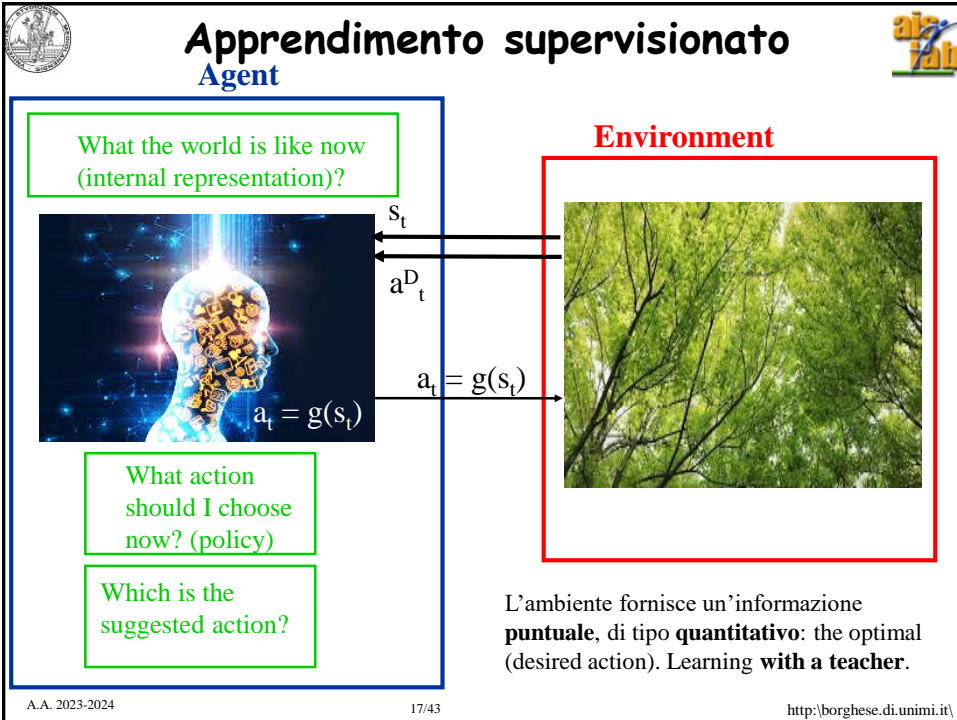


A.A. 2023-2024

14/43

<http://borghese.di.unimi.it/>





- ## Riassunto
- Valutazione di un modello
 - **Modelli multi-scala on-line**
 - Introduzione alle reti neurali
- A.A. 2023-2024 18/43 <http://borghese.di.unimi.it/>



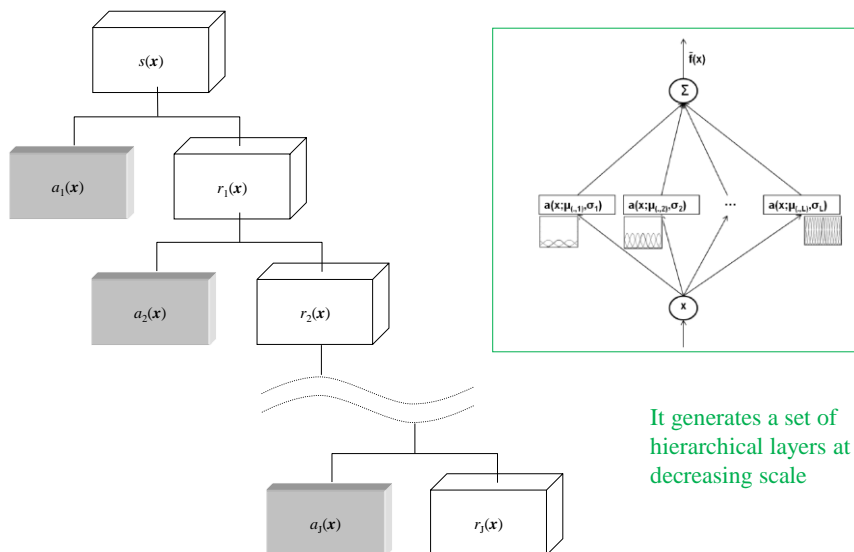
Dati e apprendimento

- Apprendimento batch. Ho a disposizione tutti i dati o lotti (batch) di dati. Costruzione del modello una-tantum.
- Apprendimento on-line. Ho a disposizione un dato alla volta. Costruzione del modello incrementale.

Active sampling. Dirigo l'attenzione in certe regioni dello spazio dei dati e li raccolgo più dati -> modulo di pianificazione dell'acquisizione dei dati.



HRBF networks



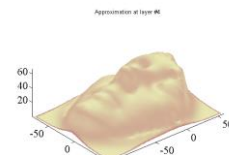
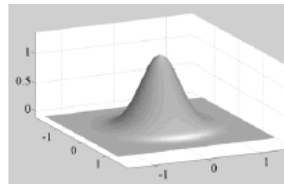
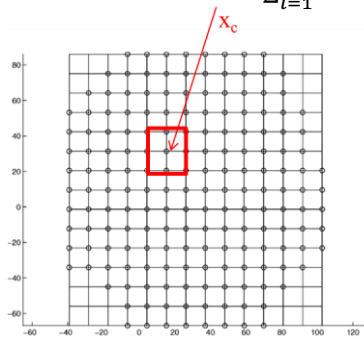
It generates a set of hierarchical layers at decreasing scale



HRBF layer

- Convolutional layer, 1, with a Gaussian at a given scale σ_1 .
- Units are equally spaced with large overlap.
- Weights are computed gridding the data.

$$f(x_{c_k}) = \frac{\sum_{i=1}^{N_c} \left(f(x_i) e^{-\frac{(x_i - x_{c_k})^2}{\sigma^2}} \right)}{\sum_{i=1}^{N_c} \left(e^{-\frac{(x_i - x_{c_k})^2}{\sigma^2}} \right)} \hat{f}(x) = \sum_{k=1}^N f(x_{c_k}) G(x; x_{c_k}, \sigma) \Delta x$$



21/43

<http://borgnese.di.unimi.it/>



Batch regression

- All data are present at start.
$$\hat{f}(x) = \sum_{k=1}^N f(x_{c_k}) G(x; x_{c_k}, \sigma) \Delta x$$
- The $w_{c_k} = f(x_{c_k})$ can be computed from all the local data (those inside the receptive field of each unit).
- The local residual can be computed from all the local data.
- Decision for splitting a quad into four quads at half resolution can be taken from all the data and it is taken in parallel for all quads of that layer
- Efficient data support through in-place ordering of the data.

Here we have 1 datum at a time....

A.A. 2023-2024

22/43

<http://borgnese.di.unimi.it/>



On-line single layer

- Each new point, $\{x_k, y_k\}$, contributes to the estimate of the function height inside the receptive fields of the associated Gaussians.
- The estimate of $f(x_c)$ has to be recomputed:

$$f(x_c) = \frac{\sum_{i=1}^{N_c} \left(f(x_i) e^{-\frac{(x_i - x_c)^2}{\sigma^2}} \right)}{\sum_{i=1}^{N_c} \left(e^{-\frac{(x_i - x_c)^2}{\sigma^2}} \right)}$$

- On-line strategy: numerator and denominator are updated separately.



On-line estimate of $f(x_c)$

$$Num(x_c) = \sum_{i=1}^{N_c} \left(f(x_i) e^{-\frac{(x_i - x_c)^2}{\sigma^2}} \right) + \left(f(x_{new}) e^{-\frac{(x_{new} - x_c)^2}{\sigma^2}} \right) =$$

$$Num'(x_c) = \sum_{i=1}^{N_c+1} \left(f(x_i) e^{-\frac{(x_i - x_c)^2}{\sigma^2}} \right)$$

$$Den(x_c) = \sum_{i=1}^{N_c} \left(e^{-\frac{(x_i - x_c)^2}{\sigma^2}} \right) + \left(e^{-\frac{(x_{new} - x_c)^2}{\sigma^2}} \right) =$$

$$Den'(x_c) = \sum_{i=1}^{N_c+1} \left(e^{-\frac{(x_i - x_c)^2}{\sigma^2}} \right)$$

$$f(x_c) = \frac{Num'(x_c)}{Den'(x_c)}$$

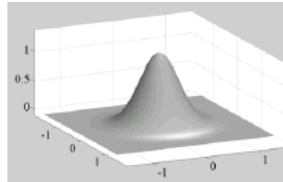
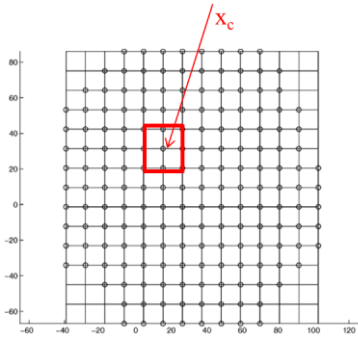
Few computations are required for any number of point: recursive estimate of $f(x_c)$

For each new point a new term is added and the ratio is recomputed only for the Gaussians whose receptive field contains that point.



For $N_c \rightarrow \infty$

$$\lim_{N_c \rightarrow \infty} \frac{Num'(x_c)}{Den'(x_c)} = E(f(x_c))$$

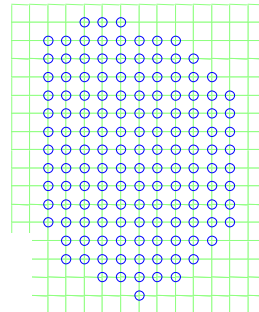
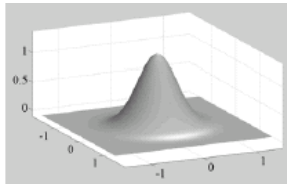


How many points are required to get a good estimate? Experimental answer.

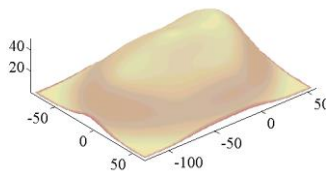
Is it sufficient to obtain a good reconstruction?



First layer



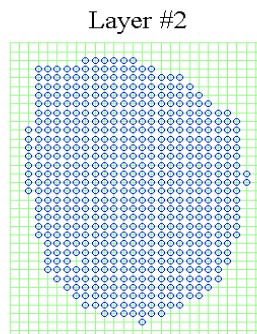
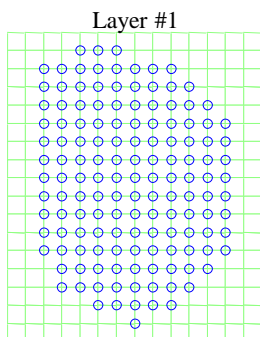
Approximation at layer #1



Asymptotically, we cannot obtain anything better than this.
Few Gaussians, large scale.



How to move to next layer?



A reliable estimate of f_c on Layer 1 ~~X~~ a reliable estimate of f_c on Layer 2.

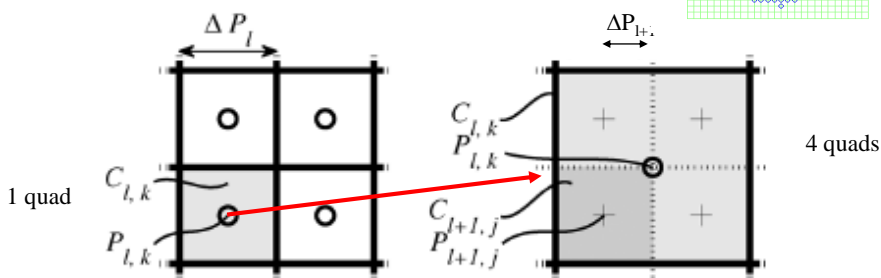
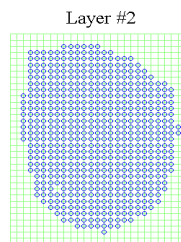
The points that belong to the receptive field of f_c on Layer 1 are split over the different quads in Layer 2.

When shall we start to estimate f_c in Layer 2?



Local operations

- Local split of each quad is achieved when:
 - Residual is higher than threshold
 - **Enough points have been sampled**
- **4 new Gaussians are generated at the higher level**





On-line estimate of $f(x_c)$ on new layer



$$Num'(x_c) = \sum_{i=1}^{N_c+1} \left(r(x_i) e^{-\frac{(x_i-x_c)^2}{\sigma^2}} \right)$$

$$Den'(x_c) = \sum_{i=1}^{N_c+1} \left(e^{-\frac{(x_i-x_c)^2}{\sigma^2}} \right)$$

This requires that the points $\{x_i\}$ inside the receptive field of the 4 Gaussians created are extracted from the $\{x_i\}$ of the Gaussian of the current layer that we have split.

How?

In-place ordering of the $\{x_i\}$ associated to the quad of the Gaussian of the current layer such that they are distributed in the quads of the four Gaussians created.

The approximation of the residual, $a(x_c)$ is initialized with few points. Computation of $Num'(x_c)$ and $Den'(x_c)$ for the four new Gaussians requires little effort.

$$r(x_c) = \frac{Num'(x_c)}{Den'(x_c)}$$

A.A. 2023-2024

29/43

<http://borghese.di.unimi.it/>



On-line version



- Data do not arrive all together (batch)
- One data at a time.
- **Growing while scanning**



2 min video



A.A. 2023-2024

30/43

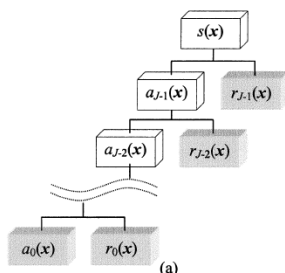
<http://borghese.di.unimi.it/>



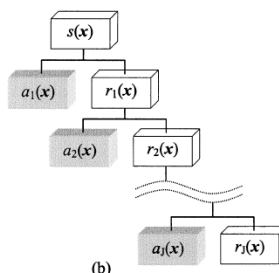
Comparison with Multi-Resolution Analysis (wavelet)



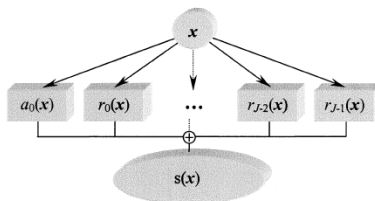
MRA – Coefficients determination



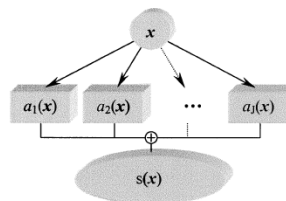
HRBF – Parameters determination



MRA – Reconstruction



HRBF – Reconstruction



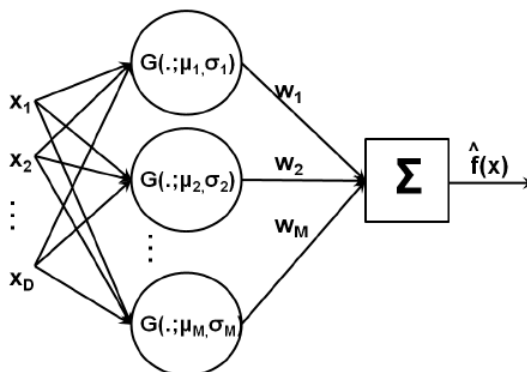
Riassunto



- Valutazione di un modello
- Modelli multi-scala on-line
- **Introduzione alle reti neurali**



Verso le reti neurali artificiali



- $G(\cdot)$ può essere sostituita da un'altra funzione di base
- Le unità **non devono essere necessariamente equispaziate**
- I parametri sono **personalizzati su ciascuna unità**

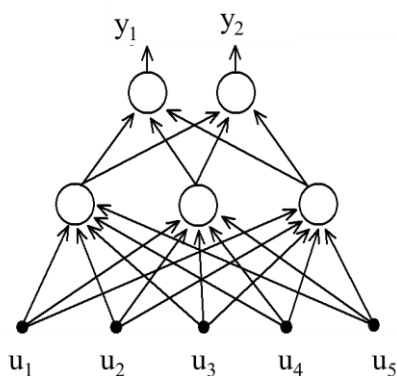
A.A. 2023-2024

33/43

<http://borgnese.di.unimi.it/>



Neuroni artificiali «moderni»



Connessionismo classico. Uscita compresa tra min – Max. Tra 0 e 1 (o tra -1 e 1): $y_i = [0 \ 1]$.

L'uscita può essere binarizzata definendo una soglia.

L'input può essere compreso tra $-\infty$ e $+\infty$ o limitato a seconda delle unità utilizzate in input.

A.A. 2023-2024

34/43

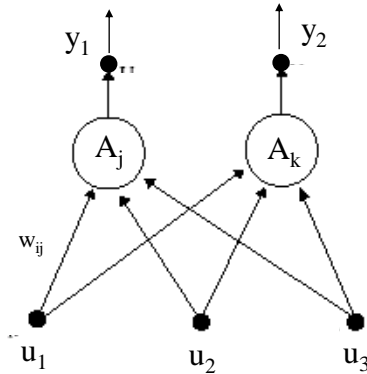
<http://borgnese.di.unimi.it/>



La rete neurale a un livello



La rete opera una trasformazione dallo spazio di input allo spazio di output.



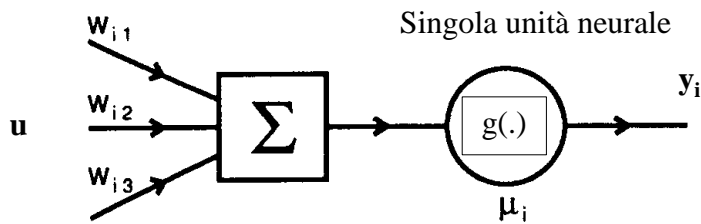
$$y_i = \mathbf{g}(\sum_j (w_{ij}u_j) - \mu_i)$$

La trasformazione o mappatura dipende dai parametri $\{w_{ij}\}$ e $\{\mu_i\}$ in modo tale che la rete neurale approssimi la trasformazione tra i pattern di input e di output.

Se $\mathbf{g}(\cdot) = 1$, la rete diventa un **modello lineare**: $y_i = \sum_j (w_{ij}u_j) - \mu_i$

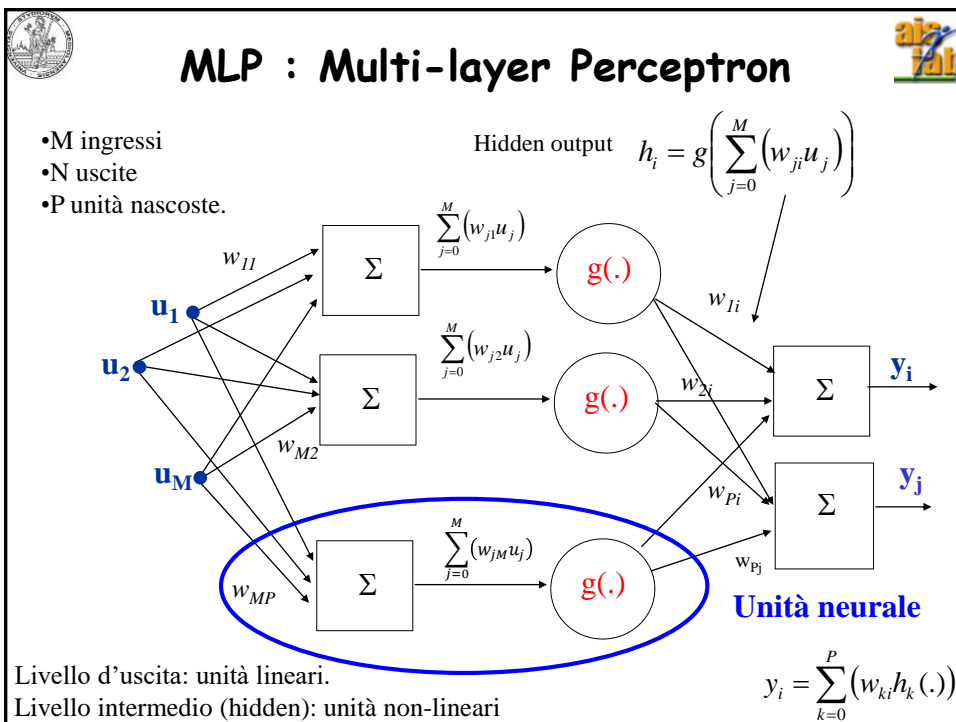
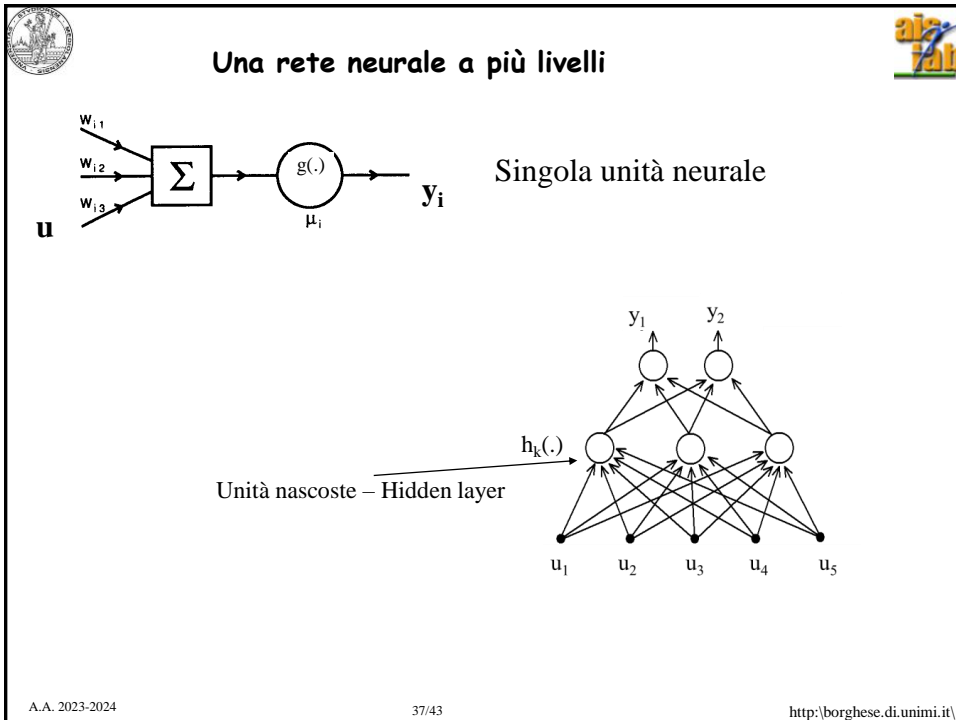


Una rete neurale a più livelli



y_i viene attivata quando la somma pesata degli ingressi supera una certa soglia.

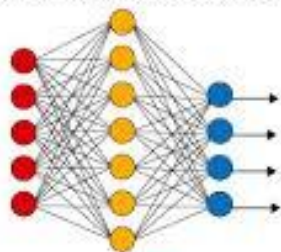
Perceptrone lineare o non-lineare



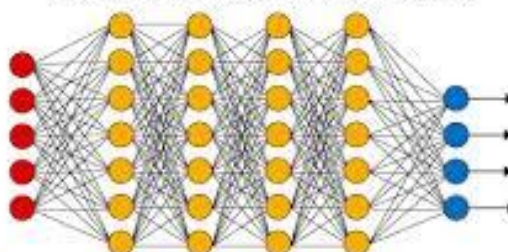


Deep Neural Networks

Simple Neural Network



Deep Learning Neural Network



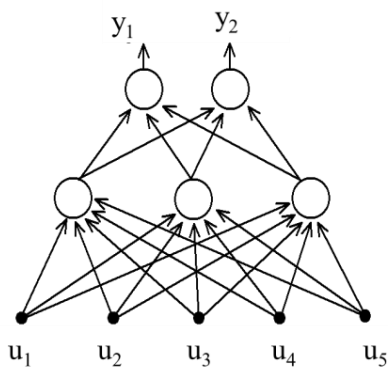
● Input Layer ● Hidden Layer ● Output Layer

Convolutional layers (e.g. RBF networks) - filters
Pooling
Subsampling
Gather-scatter...

Operazioni semplici e parallele -> match perfetto per le GPU



Caratteristiche



Livelli di unità di attivazione

Collegamento in cascata

Input convergenti, output divergenti.

Capacità di approssimazione universale

Perceptrone: layered networks, flusso unidirezionale dell'elaborazione.

L'output viene interpretato come frequenza di scarica del neurone d'uscita della rete.



Complessità della funzione realizzabile



Quanti più neuroni artificiali vengono connessi tanto più la funzione complessiva approssimabile diviene più complessa

$$Y = |y_1, y_2, y_3, \dots, y_n|^T$$

$$X = |x_1, x_2, x_3, \dots, x_m|^T$$

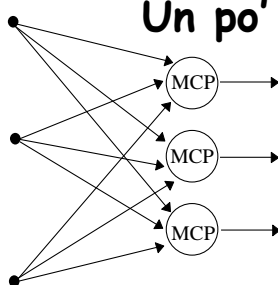
$$y = F(X)$$



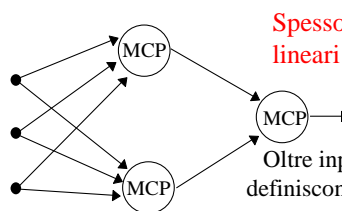
Reti neurali = approssimatori universali.



Un po' di tassonomia



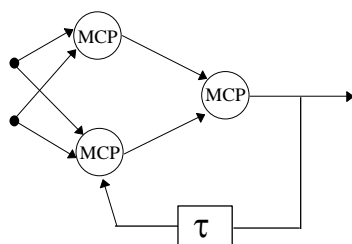
Perceptrone semplice



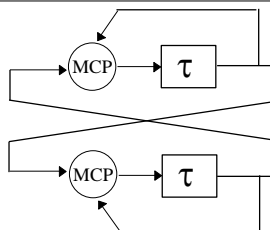
Spesso unità lineari

Oltre input/output si definiscono anche unità nascoste (**hidden units**)

Perceptrone multistrato



Ricorrente



Ricorrente completamente connessa: autoassociativa (ingresso=stato)



Riassunto



- Valutazione di un modello
- Modelli multi-scala on-line
- Introduzione alle reti neurali