

L'intelligenza biologica

Reti Neurali con funzione di attivazione a base radiale

Alberto Borghese
Università degli Studi di Milano
Laboratorio di Motion Analysis and Virtual Reality (MAVR)
Dipartimento di Scienze dell'Informazione
borghese@dsi.unimi.it



Sommario



RBF: reti neurali con neuroni a base radiale.

Struttura della rete.

Apprendimento ibrido.

Teoria della regolarizzazione.

Stima Bayesiana.

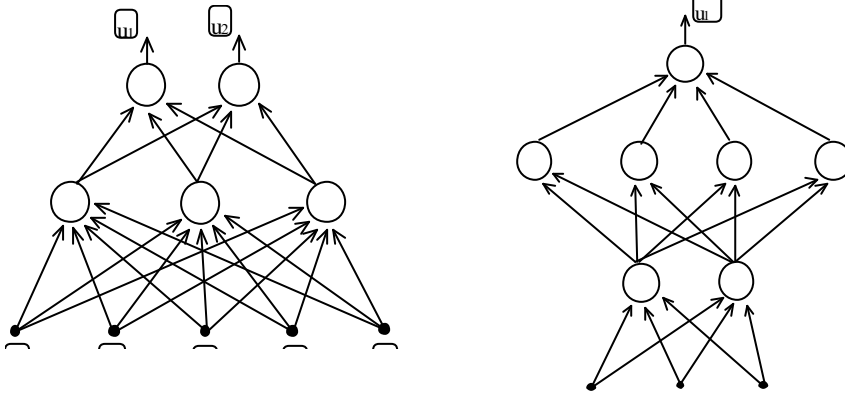
Teoria del filtraggio.

Approccio gerarchico.

Approccio gerarchico locale.



Problema con MLP



Difficult to make it learn!



Approssimazione e RBF networks: ipotesi



Apprendimento da esempi è in molti casi equivalente ad approssimare una funzione multi-variabile.

Dati $\{P(x_1, x_2, x_3, \dots, x_M, x_{M+1})\}$ posso scrivere il mio set di esempi come: $\{P(\mathbf{x}, y(\mathbf{x}))\}$ con $y = x_{M+1}$ e $\mathbf{x} \in \mathbb{R}^M$.

Ad esempio nello spazio 3D la mia funzione rappresenterà l'altezza della superficie $z = f(x,y)$.

E' una ricostruzione $2\frac{1}{2}$ D.



Alcuni Richiami sulla Teoria del Filtraggio



Prodotto di convoluzione (visualizzazione grafica)



Input

(Linear) System
= Filter – $h(x)$

Output

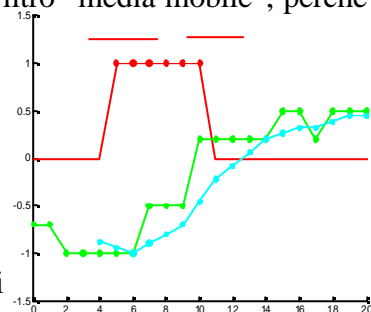
$$out(x_k) = 1/M \sum_{n=0}^M h(x_n) inpu(x_n - x_k) = h(x) * inpu(x)$$

Il segnale “scorre” verso destra e sinistra rispetto al filtro al variare del valore di x_n ”.

In ogni posizione viene effettuata la somma dei prodotti dei campioni del segnale e del fiotor.

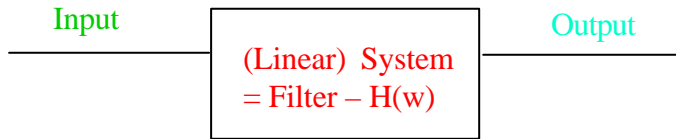
Il primo valore di uscita si ha per $x_n = 5$ pari all’ampiezza del filtro.

Filtro “media mobile”, perchè?

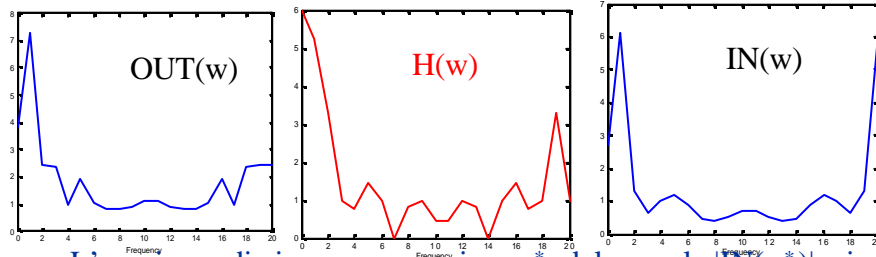




Rappresentazione nel dominio delle trasformate



La convoluzione nel dominio dello spazio equivale ad un prodotto nel dominio delle frequenze $OUT(x) = \int h(\tilde{x})IN(x-\tilde{x})d\tilde{x}$ $OUT(w) = H(w) * IN(w)$



L'ampiezza di ciascuna armonica, ω^* , del segnale $|IN(\omega^*)|$, viene modificata selettivamente a seconda dell'ampiezza di $|H(\omega^*)|$.

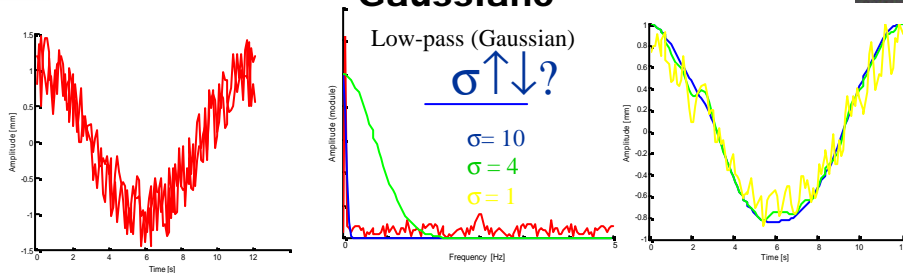
A.A. 2003-2004

7/74

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



Filtraggio reale passa-basso Gaussiano



Consideriamo il filtro Gaussiano: $h(x; x_c) = g(x - x_c; \sigma) = e^{-\frac{(x-x_c)^2}{2\sigma^2}}$

La posizione della Gaussiana è definita da x_c , la sua ampiezza da σ .

La sua Trasformata di Fourier è ancora una Gaussiana: $H(w) = e^{-1/4s^2w^2} e^{-2pjwx_c}$

$|H(w)| = e^{-1/4s^2w^2}$ A parità di ω , l'ampiezza cresce con il decresce di σ .
 σ regola l'ampiezza della banda passante del filtro.

La Gaussiana attenua in modo progressivamente maggiore le armoniche con frequenze più elevate (filtraggio passa-basso)

se



Struttura delle RBF networks



$$S(P) = \sum_{k=1}^M w_k G(P - P_k | \mathbf{s}_k)$$

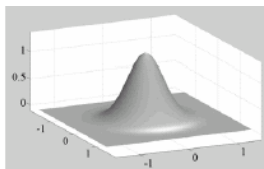


Funzioni quasi-locali



Ispirazione biologica. Campo recettivo di un neurone sensoriale (e.g. I neuroni della coclea sono sensibili ad una banda di frequenza, i neuroni della corteccia somatosensoriale primaria ad una regione del corpo, alcuni neuroni della corteccia visiva sono sensibili all'orientamento dei contorni...).

Ciascuna unità risponde ad input in una regione limitata, ovvero ciascun input attiva un certo numero di unità.



- Output della Gaussiana va velocemente a 0.
- Spline.
- In generale, funzioni a base radiale, in modulo decrescenti.



Funzioni di attivazione a simmetria radiale



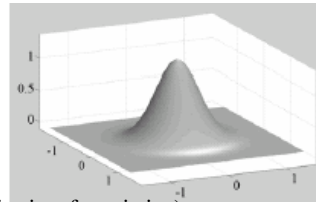
Supponiamo che vale: $f(x) = f(Rx)$, con R matrice di rotazione. Allora, si dice che $f(x)$ è invariante a rotazione, ovvero sia ha simmetria radiale.

Per le funzioni di attivazione RBF, questo vuol dire passare dalla generica formulazione:

$$\mathbf{z} = f(\mathbf{x}) = \sum_{k=1}^M \mathbf{w}_k e^{-(\mathbf{x}-\mathbf{x}_k)\Sigma_k(\mathbf{x}-\mathbf{x}_k)^T}$$

Alla formulazione a simmetria radiale:

$$z = f(\mathbf{x}) = \frac{1}{\sqrt{\pi s^2}} e^{-\frac{(\mathbf{x}-\mathbf{x}_k)^2}{s^2}} \quad (\text{ha norma } L^1 \text{ unitaria, cf. statistica})$$



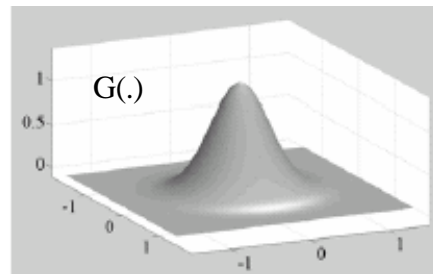
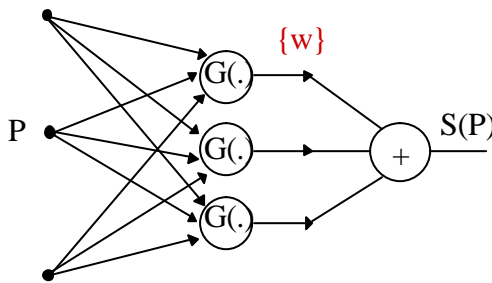
Questo riflette il fatto che non si hanno assunzioni a-priori sull'importanza delle variabili: tutte le direzioni hanno la stessa importanza.



Reti con elementi di attivazione radiali (RBF)



Perceptrone con unità quasi-locali. $S(P) = \sum_{k=1}^M w_k G(P - P_k | s_k)$

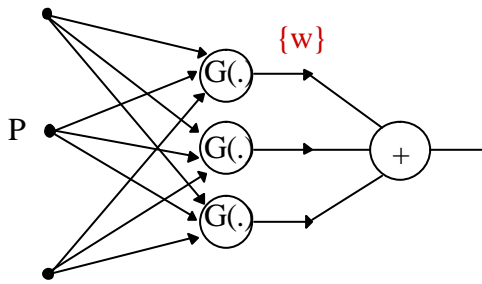


Distribuiamo le Gaussiane nell'insieme di definizione e calcoliamo la superficie come somma pesata.

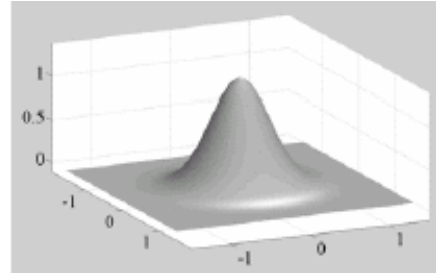
$P, P_k \in \mathbb{R}^m$ $S(P), w_k, \sigma_k$ (simmetria radiale) $\in \mathbb{R}$



Parametri delle RBF



$$S(P | \mathbf{a}) = \sum_{k=1}^M w_k G(P - P_k | \mathbf{s}_k)$$



$\{P_k\}$, M , $\{\sigma_k\}$ – parametri strutturali.
 $\{w_k\}$ – pesi sinaptici.

Per ogni Gaussiana si può definire un campo recettivo (regione di influenza). L'ampiezza è determinata da σ_k .

Al variare di α , varia la superficie ricostruita $S(P | \alpha)$.



Determinazione dei parametri (supervised learning)



Definisco una funzione costo $J(\cdot)$, funzione dei parametri α .

$$S(P | \mathbf{a}) = \sum_{k=1}^M w_k G(P - P_k | \mathbf{s}_k) \quad \min_{\{\mathbf{a}\}} J(\mathbf{a}) = \min_{\{\mathbf{a}\}} (S[(P | \mathbf{a}) - y]^2)$$

Il minimo di $J(\cdot)$ si ha quando le derivate rispetto ai parametri si annullano:

$$\frac{\partial J(\cdot)}{\partial w_i} = -2 * (\sum_k w_k g(P; P_k | \mathbf{s}_k) - y) * g(P; P_i | \mathbf{s}_i) = 0$$

$$\frac{\partial J(\cdot)}{\partial P_i} = -2 * (\sum_k w_k g(P; P_k | \mathbf{s}_k) - y) * w_i \frac{\partial g(P; P_i | \mathbf{s}_i)}{\partial P_i} = 0$$

$$\frac{\partial J(\cdot)}{\partial \mathbf{s}_i} = -2 * (\sum_k w_k g(P; P_k | \mathbf{s}_k) - y) * w_i \frac{\partial g(P; P_i | \mathbf{s}_i)}{\partial \mathbf{s}_i} = 0$$

Non ci sono vincoli sui parametri. I valori di σ tendono a crescere ed i centri a concentrarsi al centro o a respingersi fuori dal campo recettivo.

Minimi locali.



Sommario



- RBF: reti neurali con neuroni a base radiale.
- Struttura della rete.
- Apprendimento ibrido.
- Teoria della regolarizzazione.
- Stima Bayesiana.
- Teoria del filtraggio.
- Approccio gerarchico.
- Approccio gerarchico locale.



Learning ibrido



- 1) Determinazione di posizione e standard deviation delle unità.
- 2) Determinazione dei pesi sinaptici (perceptrone con unità di attivazione Gaussiane).

$$S(P | \mathbf{a}) = \sum_{k=1}^M w_k G(P - P_k | \mathbf{s}_k)$$



Determinazione di posizione ed ampiezza delle unità



Per determinare la posizione delle unità: $S(P | \alpha) = \sum_{k=1}^M w_k g(P - P_k | \sigma_k)$

- Definizione del numero di unità.
- Posizionamento delle unità mediante clustering: minimizzazione di un errore di rappresentazione (k-means, fuzzy-clustering,...).

Vado ad inserire unità dove ci sono dati.

Per determinarne l'ampiezza:

- Si utilizzano euristiche di tipo P-nearest-neighbour.

E.g. primo ordine: $\sigma = \langle \Delta x_{\alpha\beta} \rangle$

Dove α, β sono tutte le coppie di centri ottenute associando ad α il centro più vicino.

σ controlla il grado di sovrapposizione tra due Gaussiane vicine.

La media può essere calcolata localmente: P-loc-nearest-neighbor.



Determinazione dei pesi sinaptici



$$S(P | \mathbf{a}) = \sum_{k=1}^M w_k G(P - P_k | \mathbf{s}_k)$$

Utilizzo la tecnica utilizzata per il perceptrone. Quale?

Ricordiamo che le reti RBF sono chiamate anche perceptroni non lineari o Gaussiani.

NB rilevanza dal punto di vista biologico e computazionale. Ho una struttura fissa: funzioni con campo recettivo locale, con queste unità posso realizzare funzioni complesse semplicemente agendo sui pesi. Le unità mi forniscono già una pre-elaborazione dell'input locale in spazio ma anche in frequenza.



Determinazione dei pesi sinaptici



$$S(P | \mathbf{a}) = \sum_{k=1}^M w_k G(P - P_k | \mathbf{s}_k)$$

Definisco una funzione costo e calcolo derivate rispetto ai **pesi**:

$$J(\{w_k\}) = (S(P_m) - S(P_m | \mathbf{a}))^2 = \left(S(P_m) - \sum_{k=1}^M w_k G(P_m - P_k | \mathbf{s}_k) \right)^2$$

Indico i nostri esempi come $\{P_m, S(P_m)\}$

$$\frac{\partial J(\cdot)}{\partial w_i} = -2 * \left(\sum_k w_k G((P - P_k) | \mathbf{s}_k) - S(P) \right) * G((P - P_i) | \mathbf{s}_i) = 0$$

Scrivo l'equazione di cui sopra per ogni punto misurato (esempio) e avrò un sistema lineare nei pesi \Rightarrow "facilmente" risolvibile.



Risultati da Moody & Darken

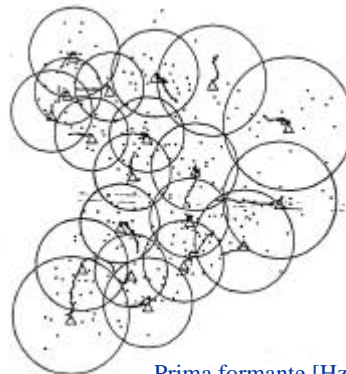


Esperimento di riconoscimento automatico del parlato.

Analizzo in tempo reale lo spettro per estrarre le formanti.

Dalle formanti voglio riconoscere le vocali (sillabe).

Seconda formante [Hz]



Prima formante [Hz]

Posizione e standard deviation delle unità della rete (20 unità). Queste unità mi ricostruiscono una superficie (pre-processing) su cui poi andrò ad applicare tecniche di clustering per individuare le diverse vocali.



Sommario



RBF: reti neurali con neuroni a base radiale.

Struttura della rete.

Apprendimento ibrido.

Teoria della regolarizzazione.

Stima Bayesiana.

Teoria del filtraggio.

Approccio gerarchico.

Approccio gerarchico locale.



Approssimazione e RBF networks: ipotesi



Apprendimento da esempi è in molti casi equivalente ad approssimare una funzione multi-variabile.

Dati $\{P(x_1, x_2, x_3, \dots, x_M, x_{M+1})\}$ posso scrivere il mio set di esempi come: $\{P(\mathbf{x}, y(\mathbf{x}))\}$ con $y = x_{M+1}$ e $\mathbf{x} \in \mathbb{R}^M$.

Ad esempio nello spazio 3D la mia funzione rappresenterà l'altezza della superficie $z = f(x,y)$.

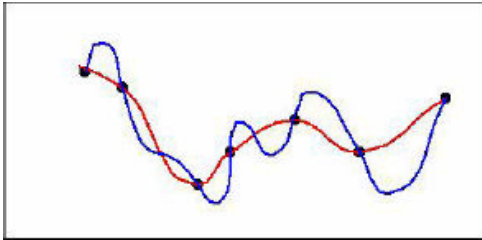
E' una ricostruzione $2\frac{1}{2}$ D.



Interpolazione dei dati



Dato un insieme $\{y_i = f(\mathbf{x}_i)\} \Rightarrow$ trovare $f(\mathbf{x})$.



Esistono infinite soluzioni.
E' un problema *mal-posto*.

Occorre introdurre delle ipotesi sull'andamento della funzione per "scegliere" quale funzione meglio rappresenta i dati.

Difficilmente siamo in grado di ricostruire andamenti complicati tra i punti, privilegeremo soluzioni "smooth".



RBF come soluzione di un problema di regolarizzazione



Dato un insieme $\{y_i = f(\mathbf{x}_i)\} \Rightarrow$ trovare $f(\mathbf{x})$ non è sufficiente, si aggiungono ipotesi sulla funzione (stabilizzatori della soluzione).

Regolarizzazione \rightarrow Problema variazionale:

$$H[f] = \min_{\{f(\cdot)\}} \left(\sum_i (f(\mathbf{x}_i) - y_i)^2 + \lambda \Phi[f] \right)$$

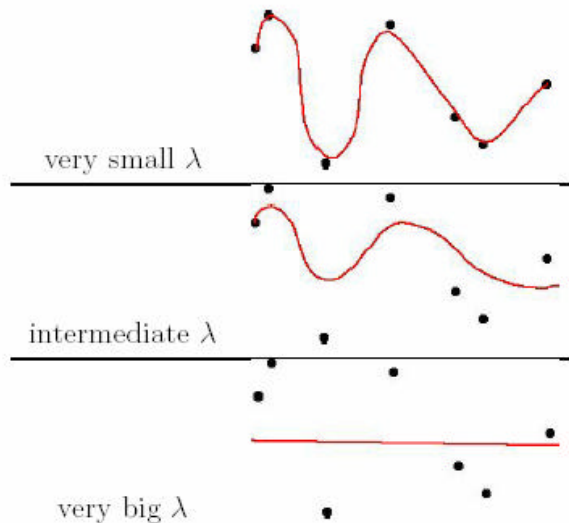
Misura la fedeltà della soluzione $f(\cdot)$ ai campioni, grado di interpolazione.

Stabilizzatore (e.g. Penalizza le variazioni brusche), grado di smoothness.

λ è un parametro scelto dall'utente.



Ruolo del parametro λ



Soluzione di un problema di regolarizzazione



$$H[f] = \min_{\{f(\cdot)\}} \left(\frac{1}{2} \sum_i (f(\mathbf{x}_i) - y_i)^2 + \lambda \Phi[f] \right) \implies \frac{dH[f]}{df} = 0$$

Trovo $f(\cdot)$ tale per cui $H[f(\cdot)]$ è minimo. E' una minimizzazione funzionale. Gestibile in pochi casi in forma analitica.

Esempio (interpolazione mediante spline cubica).

Utilizzo come stabilizzatore: $\lambda \Phi[f(\cdot)] = \int \frac{d^2 f(\cdot)}{dx^2} dt$

$$\frac{dH[f]}{df} = 0 \quad \frac{d}{df} \left(\frac{1}{2} \sum_i (f(\mathbf{x}_i) - y_i)^2 + \lambda \int \frac{d^2 f(\cdot)}{dx^2} dt \right) = 0$$



RBF come soluzione di un problema di regolarizzazione



$$H[f] = \min_{(f(\cdot))} \left(\frac{1}{2} \sum_i (f(\mathbf{x}_i) - y_i)^2 + \mathbf{I} \Phi[f] \right) \quad \frac{dH[f]}{df} = 0$$

Scrivo il funzionale come:
$$\Phi[f] = \int_{R^D} \frac{|\tilde{f}(\mathbf{w})|^2}{\tilde{G}(\mathbf{w}; \mathbf{s})} d\mathbf{w}$$

$\tilde{G}(\mathbf{w}; \mathbf{s})$ è la trasformata di Fourier di una funzione a base radiale decrescente a 0 (ad esempio, una Gaussiana). In generale σ descrive il potere filtrante di $G(\cdot)$

L'ampiezza di ogni armonica $\bar{\mathbf{w}}$ di $f(\mathbf{x})$, $\tilde{f}(\bar{\mathbf{w}})$ viene divisa per $\tilde{G}(\bar{\mathbf{w}}; \mathbf{s})$ l'ampiezza dell'armonica corrispondente del filtro. Ne risulta che armoniche ad alta frequenza producono un costo elevato.

La soluzione $f(\mathbf{x})$ ha uno spazio nullo. Sono le funzioni $p(\mathbf{x})$, che non alterano l'integrale, e si annullano nei punti \mathbf{x}_i .



Sviluppo della soluzione (I)



Trasformiamo secondo Fourier anche la funzione $f(\mathbf{x})$:

$$f(\mathbf{x}) = C \int_{R^D} d\mathbf{w} \tilde{f}(\mathbf{w}) e^{i\mathbf{x}\mathbf{w}}$$

Il funzionale diventa:

$$H[\tilde{f}] = 1/2 \sum_{i=1}^N \left(y_i - C \int_{R^D} d\mathbf{w} \tilde{f}(\mathbf{w}) e^{i\mathbf{x}_i\mathbf{w}} \right)^2 + \mathbf{I} \int_{R^D} d\mathbf{w} \frac{|\tilde{f}(\mathbf{w})|^2}{\tilde{G}(\mathbf{w}; \mathbf{s})}$$

espresso nel dominio delle frequenze.

La norma di $\tilde{f}(\mathbf{w})$ coincide con la norma di $f(\mathbf{t})$ negli spazi di Hilbert. La norma della trasformata coincide con la norma della funzione.



Sviluppo della soluzione (II)



$$H[\tilde{f}] = 1/2 \sum_{i=1}^N \left(y_i - C \int_{R^D} d\mathbf{w} \tilde{f}(\mathbf{w}) e^{i\mathbf{x}\mathbf{w}} \right)^2 + I \int_{R^D} d\mathbf{w} \frac{|\tilde{f}(\mathbf{w})|^2}{\tilde{G}(\mathbf{w}; \mathbf{s})}$$

Minimizzo rispetto alla trasformata di $f(\mathbf{x})$:

$$\frac{dH[\tilde{f}(\mathbf{w})]}{d\tilde{f}(\cdot)} = \sum_{i=1}^N \left(y_i - C \int_{R^D} d\mathbf{w} \tilde{f}(\mathbf{w}) e^{i\mathbf{x}\mathbf{w}} \right)^2 + I \int_{R^D} d\mathbf{w} \frac{|\tilde{f}(\mathbf{w})|^2}{\tilde{G}(\mathbf{w}; \mathbf{s})}$$

Dato che $f(\mathbf{x})$ è reale, la sua trasformata sarà pari (simmetrica) e quindi: $\tilde{f}^*(\mathbf{w}) = \tilde{f}(-\mathbf{w})$

$\tilde{G}(\mathbf{w}; \mathbf{s})$ Viene supposta simmetrica (cf. Gaussiana), per cui la sua (anti)trasformata sarà reale.



Sviluppo della soluzione (III)



Consideriamo il primo termine: $\frac{d}{d\tilde{f}(v)} \left\{ 1/2 \sum_{i=1}^N \left(y_i - C \int_{R^D} d\mathbf{w} \tilde{f}(\mathbf{w}) e^{i\mathbf{x}\mathbf{w}} \right)^2 \right\}$

$$= \left\{ \sum_{i=1}^N \left(y_i - C \int_{R^D} d\mathbf{w} \tilde{f}(\mathbf{w}) e^{i\mathbf{x}\mathbf{w}} \right) \int_{R^D} d\mathbf{w} \frac{d\tilde{f}(\mathbf{w})}{d\tilde{f}(v)} e^{i\mathbf{x}\mathbf{w}} \right\}$$

Antitrasformo

$$= \left\{ \sum_{i=1}^N (y_i - f(x_i)) \int_{R^D} d(\mathbf{w}-v) e^{i\mathbf{x}\mathbf{w}} d\mathbf{w} \right\}$$

$$= \left\{ \sum_{i=1}^N (y_i - f(x_i)) e^{i\mathbf{x}\mathbf{v}} \right\}$$



Sviluppo della soluzione (IV)



Consideriamo il secondo termine: $\frac{d}{df(v)} \left\{ \int_{\mathbb{R}^d} \frac{\tilde{f}(-\mathbf{w})\tilde{f}(\mathbf{w})}{\tilde{G}(\mathbf{w};\mathbf{s})} d\mathbf{w} \right\}$

$$= \int_{\mathbb{R}^d} \frac{\tilde{f}(-\mathbf{w})}{\tilde{G}(\mathbf{w};\mathbf{s})} \frac{d\tilde{f}(\mathbf{w})}{df(v)} d\mathbf{w}$$

$$\stackrel{2}{=} \int_{\mathbb{R}^d} \frac{\tilde{f}(-\mathbf{w})}{\tilde{G}(\mathbf{w};\mathbf{s})} d(\mathbf{w}-v) d\mathbf{w}$$

$$= 2 \frac{\tilde{f}(-v)}{\tilde{G}(v;\mathbf{s})}$$



Sviluppo della soluzione (V)



Ricapitolando, nel dominio delle frequenze:

$$\frac{dH[\tilde{f}(\cdot)]}{df(v)} = 0 \quad ? \quad \left\{ \sum_{i=1}^N (y_i - f(x_i)) e^{ixv} \right\} + \lambda \frac{\tilde{f}(-v)}{\tilde{G}(v;\mathbf{s})} = 0$$

Sostituiamo v a $-v$ (le trasformate non cambiano perchè sono simmetriche per ipotesi). Moltiplicando entrambi i membri per $\tilde{G}(w)$:

$$\tilde{f}(v) = \tilde{G}(-v;\mathbf{s}) \frac{\sum_{i=1}^N (y_i - f(x_i)) e^{ixv}}{1} e^{ixv}$$



Sviluppo della soluzione (VI)



$$\tilde{f}(v) = \frac{\sum_{i=1}^N (y_i - f(x_i))}{I} \tilde{G}(-v; \mathbf{s}) e^{ix_i v}$$

Non dipende da w

Prodotto nel dominio delle frequenze = convoluzione nel dominio dello spazio.

Definisco i coefficienti: $w_i = \frac{y_i - f(x_i)}{I}$

Termini ortogonali al funzionale (polonomi), L dimensione dello spazio nullo.

Da cui antitrasformando si ottiene:

$$f(x) = \sum_{i=1}^N w_i \mathbf{d}(x - x_i) * G(x; \mathbf{s}) = \sum_{i=1}^N w_i G(x - x_i; \mathbf{s}) + \sum_{l=1}^L d_l Q_l(\mathbf{x})$$



Analisi della soluzione



$$f(x) = \sum_{i=1}^N c_i \mathbf{d}(x - x_i) * G(x; \mathbf{s}) = \sum_{i=1}^N w_i G(x - x_i; \mathbf{s}) + \sum_{l=1}^L d_l Q_l(\mathbf{x})$$

$$w_i = \frac{y_i - f(x_i)}{I}$$

$$d_l = 0$$

Tante unità (N) quanti sono il numero di punti.

Se G(.) è una Gaussiana otteniamo RBF Gaussiane.

Due elementi che regolano la smoothness della soluzione

- Standard deviation, σ , di G(.), uguale per tutte le unità.
- Grado di fedeltà ai dati λ .



Come determino il valore dei pesi $\{w_i\}$



$$f(x) = \sum_{i=1}^N w_i G(x - x_i; \mathbf{s}) \quad w_i = \frac{y_i - f(x_i)}{I} \quad f(.) \text{ dipende dai } w_i$$

Ma $f(x_k)$ è una combinazione lineare di funzioni radiali.

$$f(x_j) = \sum_i w_i G(x_j - x_i; \mathbf{s}) + \sum_{l=1}^L d_l Q_l(x_j)$$

Per ogni punto m misurato, possiamo scrivere:

$$y_m = f(\mathbf{x}_m) + \lambda w_m$$

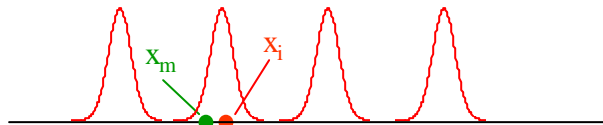
Abbiamo tanti w_m quanti sono gli esempi $([\mathbf{x}_m, y_m])$



Determinazione dei $\{w_i\}$ (I)



$$f(x) = \sum_{i=1}^N w_i G(x - x_i; \mathbf{s}) \quad w_i = \frac{y_i - f(x_i)}{I} \quad f(.) \text{ dipende dai } w_i$$



Se le Gaussiane hanno domini separati, preso un punto, x_m , in un intorno del centro della i -esima Gaussianiana, x_i ,

$$f(x) = \sum_{i=1}^N w_i G(x - x_i; \mathbf{s}) \text{ avrà come componente ? } 0 \text{ solamente: } w_i G(x - x_i; \mathbf{s})$$

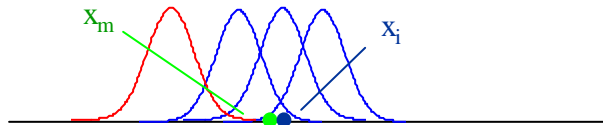
→ Ciascun peso, w_i , può essere determinato solamente guardando alla coppia di valori $[y_m, w_i G(x_m - x_i; \sigma)]$, per tutti i punti che stanno nel campo recettivo della Gaussianiana i -esima. Problemi disaccoppiati. $G(x - x_i; \sigma)$ non estende la sua influenza sugli altri punti (è utile?).



Determinazione dei $\{w_i\}$ (II)



$$f(x) = \sum_{i=1}^N w_i G(x - x_i; \mathbf{s}) \quad w_i = \frac{y_i - f(x_i)}{I} \quad f(.) \text{ dipende dai } w_i$$



Se le Gaussianne hanno domini parzialmente sovrapposti, preso un punto, x_m , in un intorno del centro della i -esima Gaussianne, x_i , un certo numero di Gaussianne daranno un contributo non nullo a $f(x) = \sum_{i=1}^N w_i G(x - x_i)$ quali?

→ Ciascun peso non può essere determinato solamente guardando alla coppia di valori $[y_m, w_i G(x_m - x_i; \sigma)]$ per quei punti che cadono nel campo recettivo della Gaussianne i -esima.

$G(x - x_i; \sigma)$ estende l'influenza sugli altri punti; occorre tenerne conto. Come?



Determinazione dei $\{w_i\}$: soluzione matriciale (I)



$$f(x) = \sum_{i=1}^N w_i G(x - x_i; \mathbf{s}) \quad y_i = f(x_i) + I w_i \quad f(.) \text{ dipende da tutti i } w_i$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \dots \\ y_N \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^N w_k G(x_1 - x_i; \mathbf{s}) \\ \sum_{i=1}^N w_k G(x_2 - x_i; \mathbf{s}) \\ \sum_{i=1}^N w_k G(x_3 - x_i; \mathbf{s}) \\ \dots \\ \sum_{i=1}^N w_k G(x_N - x_i; \mathbf{s}) \end{bmatrix} + \lambda \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ \dots \\ w_N \end{bmatrix}$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \dots \\ y_N \end{bmatrix} = \begin{bmatrix} G(x_1 - x_1; \mathbf{s}_1) & G(x_1 - x_2; \mathbf{s}_2) & G(x_1 - x_3; \mathbf{s}_3) & G(x_1 - x_N; \mathbf{s}_N) \\ G(x_2 - x_1; \mathbf{s}_1) & G(x_2 - x_2; \mathbf{s}_2) & G(x_2 - x_3; \mathbf{s}_3) & G(x_2 - x_N; \mathbf{s}_N) \\ G(x_3 - x_1; \mathbf{s}_1) & G(x_3 - x_2; \mathbf{s}_2) & G(x_3 - x_3; \mathbf{s}_3) & G(x_3 - x_N; \mathbf{s}_N) \\ \dots & \dots & \dots & \dots \\ G(x_N - x_1; \mathbf{s}_1) & G(x_N - x_2; \mathbf{s}_2) & G(x_N - x_3; \mathbf{s}_3) & G(x_N - x_N; \mathbf{s}_N) \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ \dots \\ w_N \end{bmatrix} + \lambda \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ \dots \\ w_N \end{bmatrix}$$



Determinazione dei $\{w_i\}$: soluzione matriciale (II)



$$f(x) = \sum_{i=1}^N w_i G(x - x_i; \mathbf{s}) \quad y_i = f(x_i) + I w_i \quad f(.) \text{ dipende da tutti i } w_i$$

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \dots \\ y_N \end{bmatrix} = \begin{bmatrix} G(x_1 - x_1; \mathbf{s}_1) & G(x_1 - x_2; \mathbf{s}_2) & G(x_1 - x_3; \mathbf{s}_3) & \dots & G(x_1 - x_N; \mathbf{s}_N) \\ G(x_2 - x_1; \mathbf{s}_1) & G(x_2 - x_2; \mathbf{s}_2) & G(x_2 - x_3; \mathbf{s}_3) & \dots & G(x_2 - x_N; \mathbf{s}_N) \\ G(x_3 - x_1; \mathbf{s}_1) & G(x_3 - x_2; \mathbf{s}_2) & G(x_3 - x_3; \mathbf{s}_3) & \dots & G(x_3 - x_N; \mathbf{s}_N) \\ \dots & \dots & \dots & \dots & \dots \\ G(x_N - x_1; \mathbf{s}_1) & G(x_N - x_2; \mathbf{s}_2) & G(x_N - x_3; \mathbf{s}_3) & \dots & G(x_N - x_N; \mathbf{s}_N) \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ \dots \\ w_N \end{bmatrix} + \lambda \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ \dots \\ w_N \end{bmatrix}$$

$$\mathbf{y} = (\mathbf{G} + \lambda \mathbf{I}) \mathbf{w} \quad \mathbf{y}, \mathbf{w} : N \times 1; \lambda \text{ scalare, } \mathbf{G}, \mathbf{I} : N \times N$$

$$\rightarrow \mathbf{w} = (\mathbf{G} + \lambda \mathbf{I})^{-1} \mathbf{y}$$

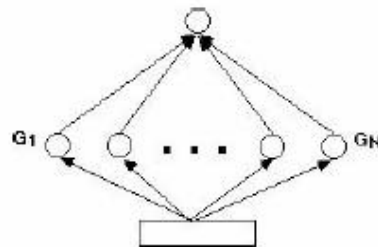
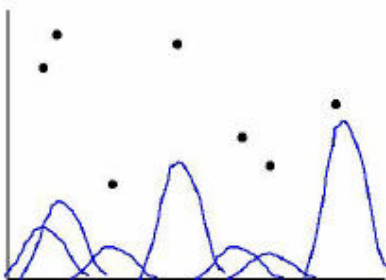
- $\lambda \rightarrow ? \infty$ il peso i -esimo è una frazione (che tende a 0) dell'altezza misurata per il campione i -esimo.
- $\lambda \rightarrow ? 0$ il peso i -esimo dipende dalla forma di G . Se s è molto piccola coinciderà con y_m , altrimenti, se più campioni cadono sotto la stessa Gaussiana, sarà una versione filtrata degli stessi.



Come agisce la regolarizzazione?



Ciascuna funzione radiale viene posta in corrispondenza di un campione e cerca quindi di interpolare in modo smooth tra gli stessi.



La struttura della rete è data dal numero di punti e dalla loro posizione, i pesi vengono calcolati tramite un sistema lineare che tiene conto dei vincoli di interpolazione e smoothness.



Sommario



RBF: reti neurali con neuroni a base radiale.

Struttura della rete.

Apprendimento ibrido.

Teoria della regolarizzazione.

Stima Bayesiana.

Teoria del filtraggio.

Approccio gerarchico.

Approccio gerarchico locale.



Approccio Bayesiano



Supponiamo che i nostri esempi: $\{\mathbf{x}_i; y_i\}$ siano ottenuti campionando in modo randomico una funzione $f(\cdot)$ non nota, in presenza di rumore: $y_i = f(\mathbf{x}_i) + \varepsilon_i$ ($i=1, \dots, N$). Gli ε_i sono indipendenti, sono funzioni randomiche che rappresentano il rumore con una sua distribuzione statistica.

Supponiamo di considerare la funzione $f(\cdot)$ come una realizzazione di un campo randomico con una certa distribuzione a-priori, che indichiamo con $P[f]$. Questo rappresenta la conoscenza a-priori sul campo di funzioni da cui è estratta e si può utilizzare per dare un vincolo soft alla soluzione, favorendo quelle funzioni di probabilità che meglio soddisfano le ipotesi a priori sul campo di funzioni.



Approccio Bayesiano: formalizzazione



Definisco $f(\cdot)$ la mia funzione probabilità da determinare e $g(\cdot) = \{\mathbf{x}_i; y_i\}$ l'insieme dei dati.

- $P[f|g]$ – Probabilità **condizionata** della funzione f , dati i campioni g .
- $P[g|f]$ – Probabilità **condizionata** di ottenere i dati g da un modello f : probabilità di ottenere i valori $\{y_i\}$ campionando la funzione nei $\{\mathbf{x}_i\}$: rappresenta un modello dell'errore.
- $P[f]$ – Probabilità **a-priori** sul campo randomico f .

Dal teorema di Bayes si ottiene: $P[f|g] = P[g|f] P[f]$



Approccio Bayesiano: le ipotesi



- $P[g|f]$ – Assumiamo che il rumore di misura (ε_i) sia normalmente distribuito con varianza σ .

$$P[g | f] = e^{-\frac{1}{2\sigma^2} \sum_{k=1}^N (y_i - f(\mathbf{x}_i))^2}$$

- $P[f]$ – Viene scelto in modo analogo al termine di regolarizzazione: diamo probabilità alta a funzioni che hanno costo basso:

$$P(f(\cdot)) = e^{-\alpha f[f]}$$

Tramite la regola di Bayes, otteniamo:

$$P(f | g) = e^{-\frac{1}{2\sigma^2} \sum_{k=1}^N (y_i - f(\mathbf{x}_i))^2} e^{-\alpha f[f]} = e^{-\frac{1}{2\sigma^2} \left[\sum_{k=1}^N (y_i - f(\mathbf{x}_i))^2 + \alpha \sigma^2 f[f] \right]}$$

Un modo classico di stimare f è mediante la MAP, $f(\cdot)$ che minimizza l'esponente \rightarrow problema di regolarizzazione con $\lambda = 2\sigma^2\alpha$.



Sommario



RBF: reti neurali con neuroni a base radiale.

Struttura della rete.

Apprendimento ibrido.

Teoria della regolarizzazione.

Stima Bayesiana.

Teoria del filtraggio.

Approccio gerarchico.

Approccio gerarchico locale.



Autoscan

$$\{z_i = S(x_i, y_i)\}$$



- Pair of video-cameras + standard laser pointer.
- The range data are obtained by “painting” the surface manually.
- Set of range data, which is denser where required.
- High precision in spot localization (cross-correlation, bright image).

Drawback: High scanning time.



RBF networks e filtri: ipotesi



Apprendimento da esempi è in molti casi equivalente ad approssimare una funzione multi-variabile.

Dati $\{P(x_1, x_2, x_3, \dots, x_M, x_{M+1})\}$ posso scrivere il mio set di esempi come: $\{P(\mathbf{x}, y(\mathbf{x}))\}$ con $y = x_{M+1}$ e $\mathbf{x} \in \mathbb{R}^M$.

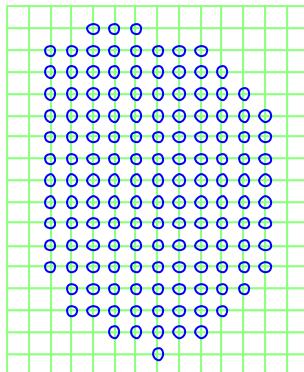
Nel caso dello scanner, i punti 3D rappresentano dei punti nello spazio geometrico 3D (e la ricostruzione è 2½ D), ma possono essere punti in uno spazio di qualsiasi dimensionalità: possono essere punti in uno spazio cognitivo, di controllo.....



The RBF model



$$S(P) = \sum_k^M w_k G(P - P_k | \sigma_k)$$



Grid support



Simple computation

Structural parameters to be set:

- M – Number of units.
- $\{P_k\}$ – Position of the units.
- σ_k – Gaussian width.

Linear parameters:

- $\{w_k\}$ – “Weights”.



RBF e Filtri



$$S(P) = \sum_k^M w_k G(P - P_k | \mathbf{s}_k)$$

Gaussiane equispaziate su un grid.

$$|P_{k+1} - P_k| = \Delta P_k$$

$$\sigma_k = \sigma \quad \forall k$$

$$w_k = S(P_k)$$

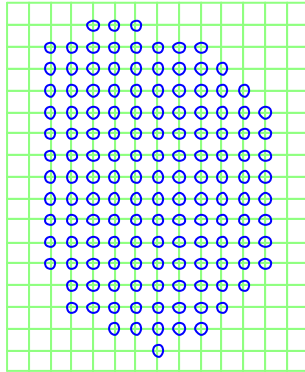
Gaussiane normalizzate $\|g(\cdot)\| = 1$

$$S(P) = \sum_{k=1}^M S_{m_k} G(P - P_k | \mathbf{s}) =$$

$$S_{m_1} G(P - P_1 | \mathbf{s}) + S_{m_2} G(P - P_2 | \mathbf{s}) + \dots =$$

$$S_m(\cdot) * G(P - \cdot | \mathbf{s})$$

$G(\cdot)$ costituisce un filtro passa-basso: $S(P) = S_m(\cdot) * G(P - \cdot | \mathbf{s})$

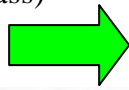


Role of the scale s

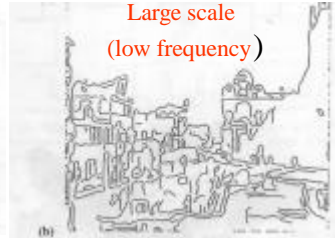
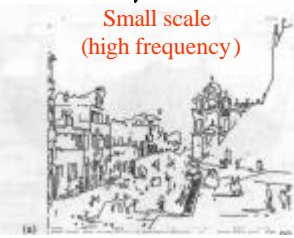


Linear filter
(low-pass)

$$S(P) = \sum_k^M S_k G(P - P_k | \mathbf{s})$$



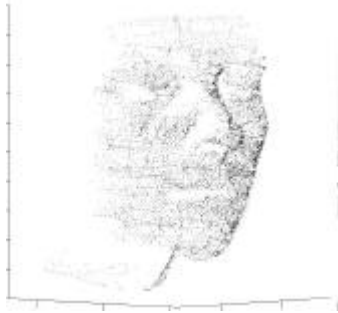
σ - Scale of the filter (bandwidth)



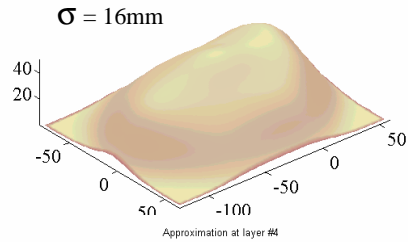
Setting s , we set the spatial frequency of the reconstruction.



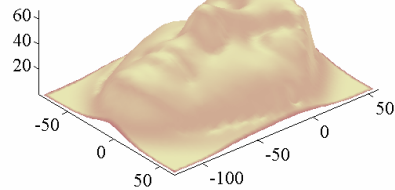
Ricostruzione da renga data a scala diversa



Approximation at layer #1



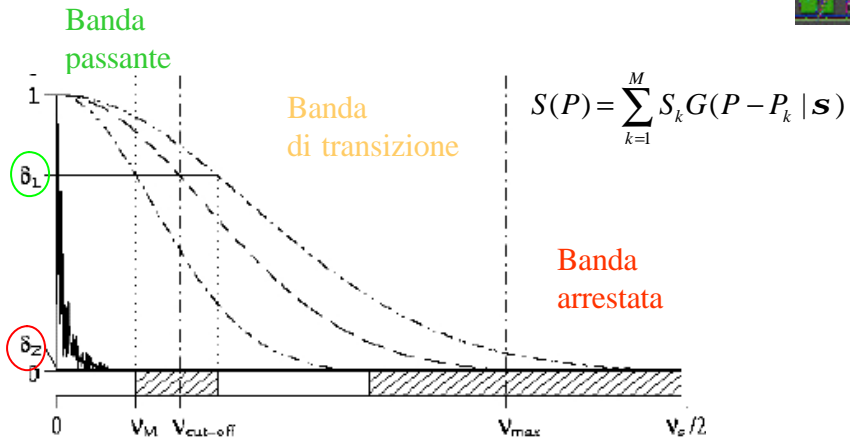
$\sigma = 2\text{mm}$



$$S(P) = \sum_k^M S_k G(P - P_k | \mathbf{s})$$



RBF come filtro -> criterio per s



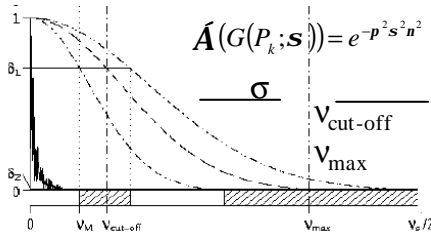
v_M Massima frequenza fatta passare inalterata (e.g. attenuazione: -3dB).

v_{max} Massima frequenza del filtro (e.g. ampiezza filtro: -40dB).

Al diminuire di σ , aumenta la frequenza -> aumenta la frequenza di campionamento (vicinanza tra due Gaussiane. $\sigma \rightarrow v_{max} \rightarrow v_s \rightarrow \Delta P_k$)

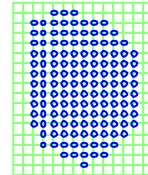


Parameters computation



Con $\delta_1 = -3\text{dB}$ e $\delta_2 = -40\text{dB}$
 $\rightarrow s_1 = 1.465 DP_1$ equivalente ad un 75% di overlap tra Gaussiani adiacenti.

Quasi-finite support



Conservative with respect to Parzen window estimators (65% of overlap)

$$S(P) = \sum_k^M S_k G(P - P_k | s) \quad \sigma \rightarrow \Delta P \{P_k\} \rightarrow M$$

All the structural parameters have been set.

With $\{S_k\}$ the surface could be reconstructed. But we do not have them



Computation of the $\{S_k\}$



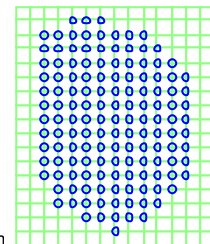
$$S(P) = \sum_k^M S_k G(P - P_k | s)$$

We do not have the height of the surface in the grid crossing, $\{S_k\}$, and we must estimate it.

$$S_k = \frac{\sum_m S_m(P_m) f(|P_m - P_k|)}{\sum_m f(|P_m - P_k|)}$$



We can apply weighted mean estimate where each point is estimated with a function, $f(|\cdot|)$, of its distance from P_k .



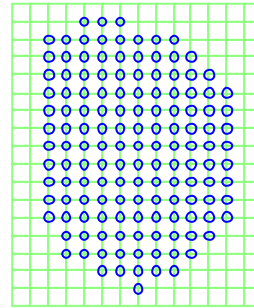


Riassunto sulla configurazione delle RBF network



$$S(P) = \sum_k^M S_k G(P - P_k | \sigma)$$

- Scelta di un valore di σ adeguato.
- Calcolo dell'”impaccamento” delle Gaussiane: $\sigma \rightarrow v_s \rightarrow \Delta P_k$ (teoria del filtraggio).
- Stima dell'altezza della superficie nei grid crossings.



$$S_k = \frac{\sum_m S_m(P_m) f(|P_m - P_k|)}{\sum_m f(|P_m - P_k|)}$$



Problems



$$S(P) = \sum_{k=1}^M S_k g(P - P_k | \sigma) \quad S_k = \frac{\sum_m S_m(P_m) f(|P_m - P_k|)}{\sum_m f(|P_m - P_k|)}$$

Related to the choice of σ :

- No information on spatial frequency content.
- Frequency variability. σ should be chosen small enough to reconstruct the finest details (high packing). Few points close to P_k to estimate $S_k = S(P_k)$.

No control on the reconstruction error.

$\{S_k\}$ are computed using all the data points. No locality assumption.

The approach as it is has little adaptivity and little hope to be real-time.



Sommario



RBF: reti neurali con neuroni a base radiale.

Struttura della rete.

Apprendimento ibrido.

Teoria della regolarizzazione.

Stima Bayesiana.

Teoria del filtraggio.

Approccio gerarchico.

Approccio gerarchico locale.



Strategia costruttiva di $s(P)$



Inizio con una scala molto ampia e valuto il residuo.

Approssimazione del residuo, dove presente, ad una scala più piccola.

Calcolo di un secondo residuo

.....

Fino a che l'errore di approssimazione soddisfa il mio criterio di approssimazione (cross-validation, contenuto in frequenza del rumore, rischio empirico o strutturale...).

NB Le Gaussianne vengono inserite negli incroci della griglia **solamente** dove l'errore (residuo) è sopra soglia (*calcolo del residuo locale*).



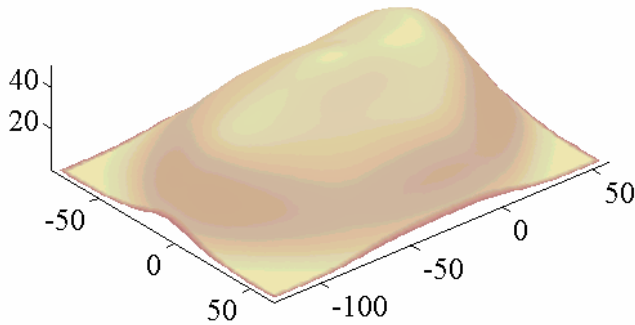
Costruzione di un primo livello a scala ampia



$\sigma_1 = 16\text{mm}$

Approximation at layer #1

$$s(P) = \sum_{k=1}^M S_k G(P - P_k | \mathbf{s}_1)$$



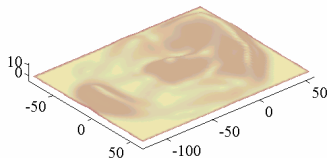
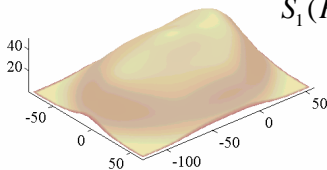
Calcolo del residuo



Residuo della ricostruzione: $r_1 = \{z_m - S_1(P)\} \forall P$ misurato, P_m .

Approximation at layer #1

$$S_1(P) = \sum_{k=1}^M S_k G(P - P_k | \mathbf{s}_1)$$



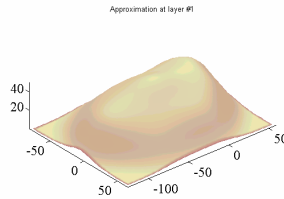
Residuo: $\{z_m - S_1(P_m)\}$



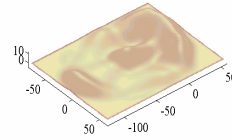
Error evaluation



Key element is the residual: $r(P_m) = z_m - S_1(P_m)$.



Continuous surface at a low scale ($\sigma = 16\text{mm}$)



Residual at the sampled points

The residual is defined only in the surface sampled points (range data).

The residual is evaluated using an **integral metric** to avoid outliers (spike reproduction). The error can be computed as:
$$e = \frac{\sum_m |r_1(P_m)|}{N}$$



Costruzione di un secondo livello gerarchico



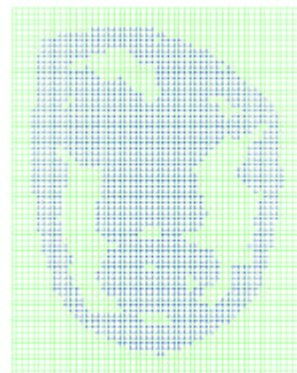
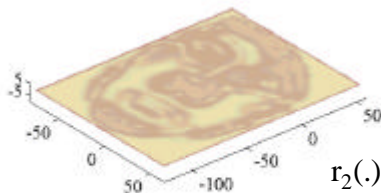
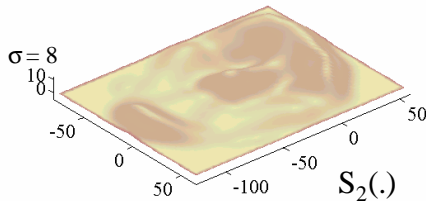
At layer 1:

$$\text{Error}(P_k) = \frac{\sum_m |r_1(P_m) - S_1(P_m)|}{N_k}$$

Output of layer #2

NB nei livelli successivi si approssima il residuo = si aggiunge dettaglio.

Layer #2





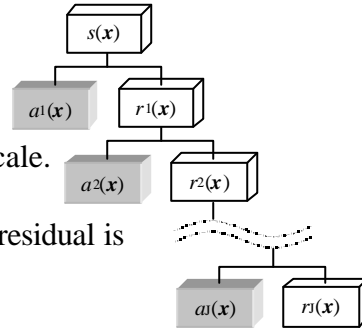
The HRBF configuration algorithm



- 1) Start with a large scale and reconstruct the surface.
- 2) Compute the residual.

Do:

- 2) Create a denser grid (lower scale).
- 3) Compute the residual at this lower scale.
- 4) Evaluate the residual.
- 5) Insert the Gaussians only where the residual is under-threshold.



Until:

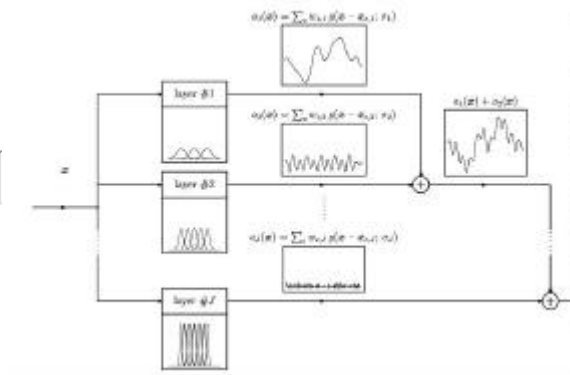
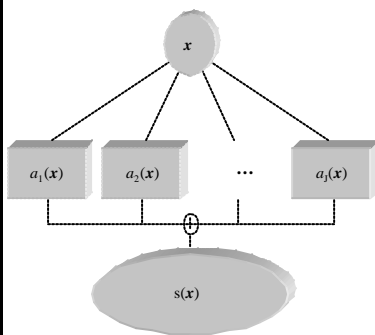
- 6) Until the residual goes under threshold.

Reconstruction is obtained by adding the approximations: $s(\mathbf{x}) = \sum_j a_j(\mathbf{x})$

Each approximation is obtained as: $S(\mathbf{x}) = \sum_k^M S_k G(\mathbf{x} - \mathbf{x}_k | S)$



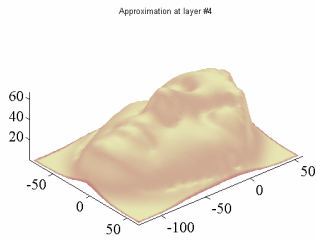
Sintesi della superficie



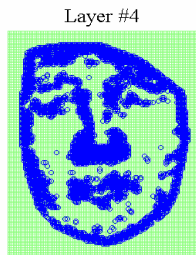
- Diversi strati (reti RBF).
- Ciascuno strato opera ad una certa scala.
- Gli strati non sono completi ma allocano Gaussiani solo dove l'errore locale è sopra-soglia.



HRBF Networks



Noise



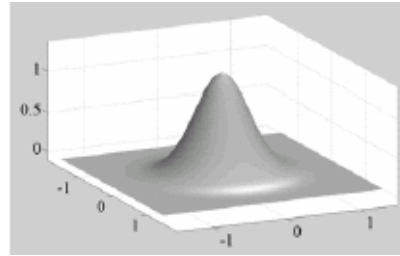
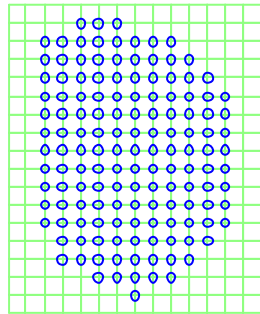
Sommario



- RBF: reti neurali con neuroni a base radiale.
 - Struttura della rete.
 - Apprendimento ibrido.
 - Teoria della regolarizzazione.
 - Stima Bayesiana.
 - Teoria del filtraggio.
 - Approccio gerarchico.
 - Approccio gerarchico locale.**



Sfruttamento della quasi-localita'



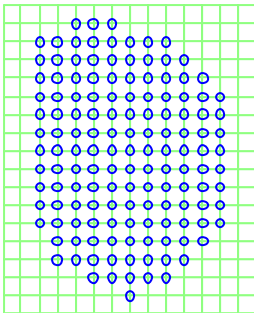
I campi recettivi sono parzialmente sovrapposti.

Operazioni quasi-locali sui punti all'interno dei campi recettivi:

- Calcolo dell'altezza della funzione nell'incrocio.
- Calcolo dell'errore locale.



Where locality is interesting?



$$S(P) = \sum_{k=1}^M S_k g(P - P_k | \sigma)$$

Estimate of the Surface in the grid crossings:

$$S_k = \frac{\sum_m S_m(P_m) f(|P_m - P_k|)}{\sum_m f(|P_m - P_k|)}$$

Computation of the error:
$$e = \frac{\sum_m |r(P_m)|}{N}$$

and its association to the grid crossings ["5) Insert the Gaussians only where the residual is under-threshold."]

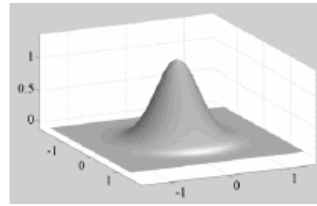
Costly operation because $f(\cdot)$ has to be computed for all the data points and there is no criterion on how to evaluate locally e .



Receptive field



$$S(P) = \sum_{k=1}^M S_k G(P - P_k | \mathbf{s})$$



- $S(P_k)$ is estimated as a weighted (with the distance) mean of the points which lie inside the receptive field of the Gaussian:

$$S_k = S(P_k) = \frac{\sum_m S(P_m) g(P_m - P_k | \mathbf{s})}{\sum_m g(P_m - P_k | \mathbf{s})} \quad P_m \in \text{RF}(P_k)$$

- An error measure can be associated to each Gaussian as the residual computed to all the points inside the receptive field.

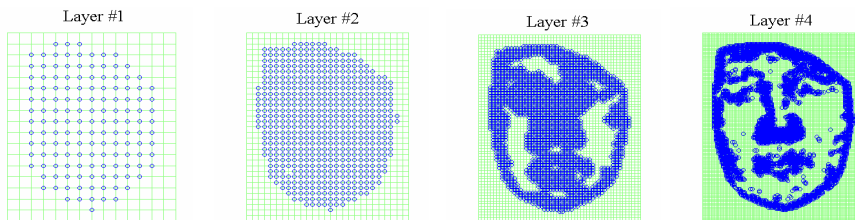
$$e(P_k) = \frac{\sum_m |r(P_m)|}{N(P_k)} \quad P_m \in \text{RF}(P_k)$$



Sparse approximation



Grids do not need to be complete, but Gaussians can be inserted only in those crossings where the error is over threshold.



We obtain a sparse approximation (less dense units, where less dense sampling and larger details).

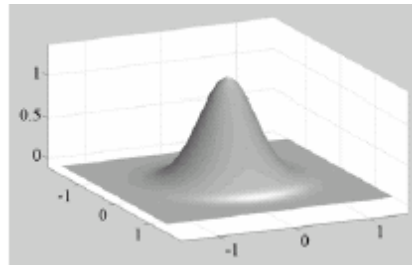
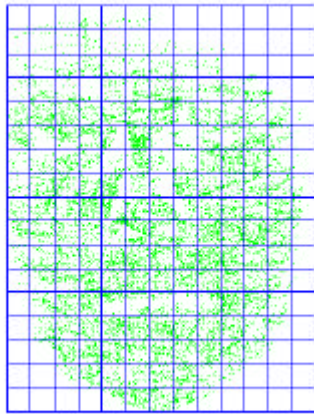
Local control of the reconstruction error. Error globally under threshold.



Representation of data locality



Approximate the Receptive field with squares (Manhattan distance).

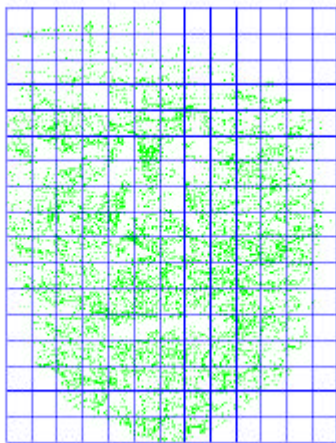


 Quad-tree subdivision

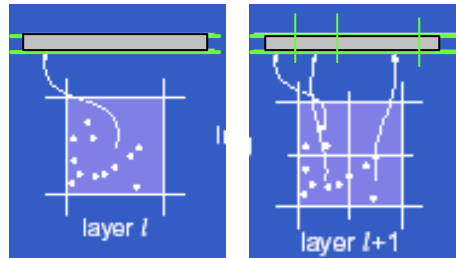
Data are partitioned into parallelepipeds supported by the grid at the lowest (largest) scale.




Efficient data storage



Quad-tree subdivision by means of in-place sorting.



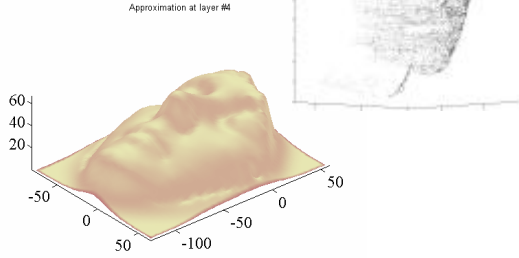
Block of range data associated to the grid square 

Data can be retrieved directly.

At every layer, only data inside the receptive field have to be sorted.



Quantitative results



#Layers	Mean error [mm]	Error std [mm]	# of gaussians	σ [mm]
1	2.33	4.38	141/270	16
2	0.03	2.06	573/980	8
3	0.04	1.02	1,778/3,795	4
4	0.06	0.77	3,078/14,933	2

Data set cardinality: 12,471. Processing time on a Pentium III, 800Mhz, was < 1s.

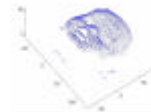


Summary of HRBF

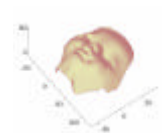


Real object

- Incremental Reconstruction (multi-scale).
- Local computation.
- Error-driven (local adaptation of the scale).
- Local control on the error.



Range data



3D surface