

UNIVERSITÀ DI PISA



Dipartimento di Informatica

Semplice implementazione di una Echo State Network

Esperienze di programmazione

Professore:
Prof. Romani

Presentata da:
Eleonora Di Gregorio
(520655)

Sessione
Anno Accademico

Indice

1	Introduzione	1
2	Background	3
2.1	Notazione	3
2.2	Formulazione dei task	4
2.2.1	Task non temporali	4
2.2.2	Task temporali	4
2.3	Reti Neurali Ricorrenti	4
3	Echo State Network	8
3.1	Il modello di base	8
3.2	Le funzioni di attivazione	10
3.3	Leaking integration	10
3.4	Echo State Property	11
4	Implementazione	13
4.1	Inizializzazione del Reservoir	13
4.1.1	Dimensione del Reservoir	14
4.1.2	Sparsità del Reservoir	15
4.1.3	Raggio spettrale	16
4.1.4	Input scaling	16
4.1.5	Leaking Rate	17
4.2	Calcolo dello stato del reservoir	17

4.3	Allenamento del Readout	19
4.3.1	Regolarizzazione di Tikhonov	19
4.3.2	Soluzione con pseudoinversa	20
5	Test	22
5.1	10th order NARMA system	22
5.2	ESNtrain() ed ESNtest()	23
5.3	Risultati	24
5.4	Conclusioni	24
A	Codice	26

Elenco delle figure

2.1	Metodo di allenamento basato su discesa di gradiente che adatta tutti i pesi delle connessioni.	6
3.1	Reservoir Computing e allenamento del readout.	9
3.2	Funzioni di attivazione comuni.	10

Elenco delle tabelle

5.1	Parametri default <i>ESNtrain()</i>	23
-----	---	----

Capitolo 1

Introduzione

Il Reservoir Computing è un metodo di grande interesse nel campo della ricerca sulle reti neurali ricorrenti, grazie alla sua semplicità e alle sue buone performance su una certa varietà di task. Con Reservoir Computing ci si riferisce ad una classe di modelli di reti neurali ricorrenti che sono caratterizzati da una parte dinamica ricorrente e una parte di output semplice. L'idea alla base del Reservoir Computing è quella di prendere una rete ricorrente random di semplici nodi, ad esempio neuroni sigmoidali, chiamata appunto reservoir. I nodi sono collegati tra loro, associando dei pesi ai collegamenti che possono essere riscaldati per imporre il regime dinamico desiderato. Il reservoir permette dunque di mappare un input in uno spazio di dimensione maggiore che nel caso di task di classificazione può aumentare la probabilità di separazione lineare tra i dati. Oltre a questo permette anche di tenere traccia della storia degli input nel tempo, questi aspetti lo rendono ideale per molti task interessanti nel mondo reale che richiedono elaborazione temporale e mapping non lineare.

È proprio in questo contesto che si collocano le Echo State Network, affermandosi per la loro praticità e per essere concettualmente semplici. Per ottenere delle buone performance sui task vanno però tenuti in considerazione diversi aspetti e proprietà di cui si parlerà in seguito. Questa semplice implementazione delle Echo State Network vuole essere un primo approccio

guidato alla conoscenza dei meccanismi sui cui si basa questo modello e alla loro applicazione. In appendice è riportato il codice matlab della ESN sviluppata.

Capitolo 2

Background

Prima di iniziare la descrizione del modello delle ESN è necessario definire la notazione che verrà utilizzata in seguito e fissare alcuni concetti importanti riguardo alle reti neurali ricorrenti.

2.1 Notazione

Si utilizzano le lettere in corsivo maiuscolo e minuscolo per denotare gli scalari, le lettere in grassetto minuscolo per i vettori, quelle in grassetto maiuscolo per le matrici. Sia \mathbf{X} una matrice, allora $\mathbf{X}_{i,j}$ indica l'elemento alla i -esima riga e la j -esima colonna di \mathbf{X} ; inoltre $\mathbf{X}_{:,j}$ denota la j -esima colonna di \mathbf{X} e $\mathbf{X}_{i,:}$ la i -esima riga.

In particolare un vettore di input è rappresentato da \mathbf{u} , oppure da $\mathbf{u}(t)$ se si vuole indicare la posizione che occupa l'input all'interno di una sequenza. Quest'ultima viene rappresentata usando la notazione $s(\mathbf{u}) = [\mathbf{u}(1), \mathbf{u}(2) \dots \mathbf{u}(n)]$, dove n è il numero di elementi che appartengono ad una sequenza, definiti anche timesteps nel caso in cui questa rappresenti l'input di un task temporale. Analogamente per un elemento e una sequenza di output si utilizzano \mathbf{y} e $s(\mathbf{y})$, rispettivamente.

2.2 Formulazione dei task

Un task nel contesto del machine learning consiste nell'apprendere una relazione funzionale tra un dato input $\mathbf{u}(t) \in \mathbb{R}^{N_u}$ ed un output atteso $\hat{\mathbf{y}}(t) \in \mathbb{R}^{N_y}$, dove $t = 1, \dots, TeT$ è il numero di elementi $(\mathbf{u}(t), \hat{\mathbf{y}}(t))$ del dataset per il training, chiamati anche data points.

2.2.1 Task non temporali

Si ha un task non temporale quando i data points sono indipendenti gli uni dagli altri, l'obiettivo è approssimare una funzione $\mathbf{y}(t) = y(\mathbf{u}(t))$ che minimizzi la misura dell'errore $E(\mathbf{y}, \hat{\mathbf{y}})$.

2.2.2 Task temporali

Si ha un task temporale quando \mathbf{u} e $\hat{\mathbf{y}}$ sono segnali in un dominio discreto nel tempo $n = 1, \dots, T$, e l'obiettivo è approssimare una funzione $\mathbf{y}(t) = y(\dots, \mathbf{x}(t-1), \mathbf{x}(t))$ tale che minimizzi $E(\mathbf{y}, \hat{\mathbf{y}})$. A differenza dei task non temporali la funzione che bisogna apprendere ha memoria, ovvero dipende dalla storia degli input.

2.3 Reti Neurali Ricorrenti

Le reti neurali ricorrenti (RNN) sono una classe di modelli di reti neurali biologicamente ispirata e ampiamente utilizzata per il trattamento di dati sequenziali. In una RNN gli elementi che rappresentano i neuroni, chiamati unità, sono collegati tra loro a simulare le connessioni sinaptiche e sono l'origine delle attivazioni che percorrono la rete attraverso questi collegamenti. La principale caratteristica delle reti neurali ricorrenti, dalla quale ne deriva anche il nome, è che la topologia delle connessioni presenta dei cicli, ovvero sono ammesse connessioni sinaptiche ricorrenti (chiamate anche feedback).

Per tale ragione la RNN è il modello più simile al cervello umano. L'esistenza dei cicli ha profondi impatti come riportato in [1]:

- Una RNN può sviluppare autonomamente un sistema dinamico nel tempo lungo i suoi percorsi di connessione ricorrenti.
- Se guidata da un input, una RNN preserva nel suo stato interno una trasformazione non lineare della storia dell'input, ovvero ha una memoria dinamica.

Ogni unità della RNN rappresenta uno stato $x(t) \in \mathbb{R}$, la sua attivazione è calcolata in base al suo stato corrente ed al valore che riceve in input. Sia $\mathbf{x}(t-1)$ il vettore che rappresenta lo stato delle unità che al tempo t ricevono in input $\mathbf{u}(t)$, le loro nuove attivazioni saranno:

$$\mathbf{x}(t) = f_x(\mathbf{u}(t), \mathbf{x}(t-1)) = f(\mathbf{W}_{\text{in}}\mathbf{u}(t) + \mathbf{W}\mathbf{x}(t-1) + \mathbf{b}) \quad (2.1)$$

dove:

- f è la funzione di attivazione, generalmente la funzione simmetrica *tanh* o la funzione sigmoideale logistica,
- \mathbf{W}_{in} è la matrice dei pesi di ingresso, ovvero pesi associati alle connessioni tra l'input e lo stato interno,
- \mathbf{W} è la matrice dei pesi delle connessioni tra le diverse unità della rete,
- \mathbf{b} è un vettore di bias che vengono applicati alle unità.

L'output della rete viene poi calcolato secondo la seguente formula:

$$\mathbf{y}(t) = f_y(\mathbf{x}(t)) = f_{\text{out}}(\mathbf{W}_{\text{out}}\mathbf{x}(t)) \quad (2.2)$$

dove la f_{out} può essere una funzione di attivazione lineare, ad esempio la funzione identità, o una funzione non lineare. L'approccio classico per l'allenamento delle RNN è quello basato su discesa di gradiente, cioè vengono

adattati tutti i pesi \mathbf{W}_{in} , \mathbf{W} , \mathbf{W}_{out} rispettivamente in base a $\partial E/\partial \mathbf{W}_{\text{in}}$, $\partial E/\partial \mathbf{W}$, $\partial E/\partial \mathbf{W}_{\text{out}}$, con lo scopo di minimizzare $E(\mathbf{y}, \hat{\mathbf{y}})$.

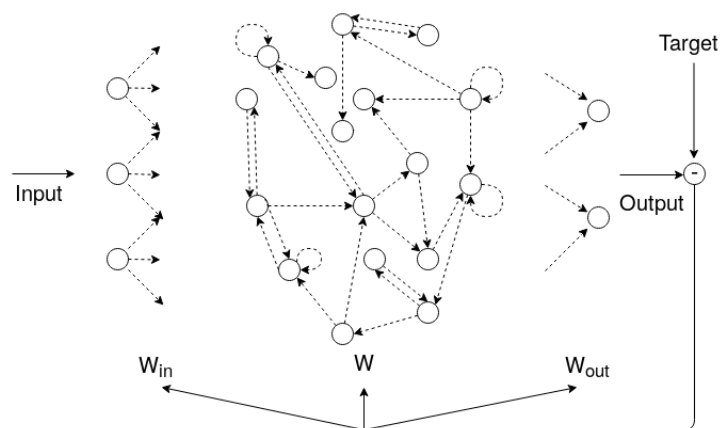


Figura 2.1: Metodo di allenamento basato su discesa di gradiente che adatta tutti i pesi delle connessioni.

Le reti neurali ricorrenti si presentano come strumenti molto promettenti per processare sequenze temporali non lineari per due ragioni. Innanzitutto rappresentano una approssimazione dei sistemi dinamici, in secondo luogo si avvicinano al modello del cervello biologico grazie alle connessioni ricorrenti. Tuttavia vi sono delle problematiche legate all'apprendimento evidenziate in [1] e riportate qui di seguito:

- Non può essere garantita la convergenza del metodo di allenamento a discesa di gradiente, vi sono dei punti in cui le informazioni sul gradiente possono essere mal definite.
- L'aggiornamento di un solo parametro può essere computazionalmente costoso, cioè possono essere necessari molti cicli di aggiornamento. Questo comporta lunghi tempi di addestramento che rendono le RNN modelli convenienti utilizzando solo poche unità.

- È molto difficile apprendere relazioni che richiedono molta memoria perché le informazioni fornite dal gradiente subiscono una dissolvenza esponenziale nel tempo.
- Gli algoritmi di apprendimento necessitano di profonde conoscenze ed esperienza per essere applicati con successo.

In questa situazione di lento progresso, nel 2001 un nuovo approccio al design e all'allenamento delle RNN, noto come Reservoir Computing, venne proposto indipendentemente da Wolfgang Maass [2] e da Herbert Jaeger [3].

Capitolo 3

Echo State Network

3.1 Il modello di base

Come evidenziato nella sezione 2.3 le diverse problematiche relative all'allenamento delle reti neurali ricorrenti hanno portato allo sviluppo del Reservoir Computing (RC). L'approccio del Reservoir Computing differisce dalla tecnica di apprendimento citata in precedenza perché crea una separazione concettuale e computazionale tra una parte dinamica chiamata *reservoir* ed una parte chiamata *readout*. Il *reservoir* è una rete neurale ricorrente che opera come funzione di espansione non lineare, mentre il *readout* produce l'output desiderato dall'espansione.

Una Echo state Network consiste in un livello di input di N_u unità, un insieme di N_r unità nascoste che compongono il suddetto *reservoir* e un livello di output di N_y unità tipicamente lineari e non ricorrenti che formano il *readout*. Le equazioni che descrivono la computazione sviluppata da una ESN sono:

$$\mathbf{x}(t) = f_x(\mathbf{u}(t), \mathbf{x}(t-1)) = f(\mathbf{W}_{\text{in}}\mathbf{u}(t) + \mathbf{W}\mathbf{x}(t-1) + \mathbf{b}) \quad (3.1)$$

$$\mathbf{y}(t) = f_y(\mathbf{x}(t)) = \mathbf{W}_{\text{out}}\mathbf{x}(t) \quad (3.2)$$

La separazione concettuale che caratterizza le echo state network è basata sulla comprensione dei diversi scopi che assumono f_x e f_y :

- f_x espande la storia dell'input $\mathbf{u}(t), \mathbf{u}(t-1), \dots$ in uno spazio degli stati del *reservoir* $\mathbf{x}(t) \in \mathbb{R}^{N_x}$,
- mentre f_y combina i segnali ricevuti dai neuroni $\mathbf{x}(t)$ per ottenere gli output attesi $\hat{\mathbf{y}}(t)$.

La principale caratteristica del RC e quindi delle ESN è che la parte ricorrente della rete può non essere allenata dopo essere inizializzata secondo alcune proprietà facili da soddisfare. L'allenamento si riduce dunque alla sola parte non ricorrente, si riducono così i costi computazionali che gravavano sul design delle RNN. Di seguito in figura 3.1 viene illustrato il design e l'allenamento di una ESN.

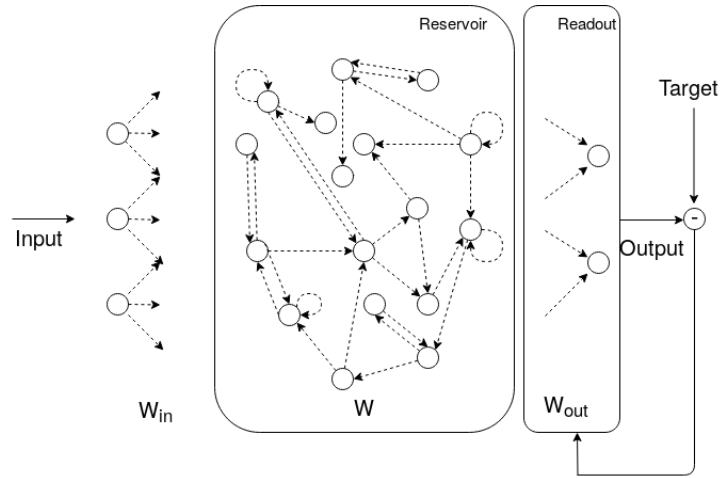


Figura 3.1: Reservoir Computing e allenamento del readout.

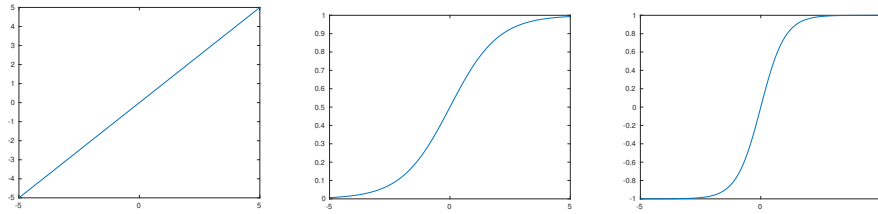
Ciò significa che \mathbf{W}_{in} e \mathbf{W} sono fisse e \mathbf{W}_{out} varia in base all'allenamento. Tuttavia non tutte le scelte di \mathbf{W}_{in} e \mathbf{W} producono una ESN valida, il

reservoir deve essere un eco ("echo") della storia dell'input, cioè le dinamiche del *reservoir* devono dipendere asintoticamente solo dai segnali di input.

3.2 Le funzioni di attivazione

Le funzioni di attivazioni alle quali si fa riferimento nella formula 3.1 possono essere diverse, tra queste:

- la funzione identità : $f_{id}(x) = x$ (3.2a)
- la funzione logistica : $f_{logistic}(x) = \frac{1}{1+e^{-x}}$ (3.2b)
- la tangente iperbolica : $f_{tanh}(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$ (3.2c)



(a) Funzione identità (b) Funzione logistica (c) Tangente iperbolica

Figura 3.2: Funzioni di attivazione comuni.

La funzione logistica e la tangente iperbolica sono le più usate, l'idea è che si può mappare ogni numero reale in un numero nell'intervallo $[0, 1]$ e rispettivamente in $[-1, 1]$, in questo modo si può dimostrare che una combinazione di queste funzioni può approssimare ogni funzione non lineare.

3.3 Leaking integration

Le unità nelle reti sigmoidali standard non hanno memoria; il loro valore ad un tempo $n + 1$ dipende solo parzialmente ed indirettamente dai valori

precedenti. Perciò queste reti si adattano al meglio per sistemi discreti nel tempo, d'altra parte è difficile l'apprendimento di dinamiche lente come onde sinusoidali molto lente. Per apprendere sistemi lenti e variabili in uno spazio continuo è più adeguato utilizzare delle reti con dinamiche continue. Per questo già in [3] viene introdotto un approccio ibrido che usa reti discrete nel tempo che approssimano delle reti continue grazie alle unità con *leaky integration*. In questo caso una semplificazione dell'equazione per il calcolo dello stato delle unità è la seguente:

$$\begin{aligned} \mathbf{x}(t) &= (1 - \alpha)\mathbf{x}(t - 1) + \alpha f_x(\mathbf{u}(t), \mathbf{x}(t - 1)) = \\ &(1 - \alpha)\mathbf{x}(t - 1) + \alpha f(\mathbf{W}_{\text{in}}\mathbf{u}(t) + \mathbf{W}\mathbf{x}(t - 1) + \mathbf{b}) \end{aligned} \quad (3.3)$$

dove $\alpha \in (0, 1]$ è il *leaking decay rate*, se si imposta $\alpha = 1$ si ottiene una ESN standard.

3.4 Echo State Property

Una ESN valida soddisfa la cosiddetta *Echo State Property* (ESP) (Jager, 2001). Questa dice che lo stato in cui la rete si trova dopo l'analisi di una lunga sequenza di input dipende solo dalla sequenza stessa. La dipendenza dallo stato iniziale della rete è progressivamente persa quando la lunghezza della sequenza di input tende all'infinito. Analogamente, lo stato corrente della rete $\mathbf{x}(t)$ è una funzione della storia dell'input indipendentemente dai valori iniziali dello stato. In formule la ESP può essere espressa come segue:

$$\begin{aligned} \forall s_n(\mathbf{u}) = [\mathbf{u}(1), \dots, \mathbf{u}(n)] \in (\mathbb{R}^{N_u})^n \text{ sequenza di input di lunghezza } n, \\ \forall \mathbf{x}, \mathbf{x}' \in \mathbb{R}^{N_r} : \end{aligned} \quad (3.4)$$

$$\|f_x(s_n(\mathbf{u}), \mathbf{x}) - f_x(s_n(\mathbf{u}), \mathbf{x}')\| \rightarrow 0 \text{ per } n \rightarrow \infty$$

che significa che la differenza tra gli stati che la rete assume dopo che viene valutata una sequenza di input di lunghezza n tende a zero al tendere di n

ad infinito, per ogni scelta iniziale dello stato.

Jaeger ha fornito due condizioni, rispettivamente necessaria e sufficiente, affinché una ESN goda di questa proprietà. La condizione necessaria è che il raggio spettrale, ovvero il maggiore degli autovalori in modulo, della matrice dei pesi del *reservoir* sia minore di uno.

$$\rho(\mathbf{W}) < 1 \quad (3.5)$$

Se questa condizione viene violata, il reservoir è localmente asintoticamente instabile nello stato $\mathbf{0} \in \mathbb{R}^{N_r}$ e la echo state property non può essere garantita se la sequenza nulla è un input ammissibile per il sistema.

La condizione sufficiente per la validità della proprietà è che il massimo valore singolare di \mathbf{W} sia minore di uno:

$$\sigma(\mathbf{W}) < 1 \quad (3.6)$$

Ciò assicura la stabilità globale del sistema e la presenza di "stati echo", garantendo così la echo state property.

Una caratteristica molto importante evidenziata in [4] è la contrattività delle dinamiche dello stato.

La contrattività della funzione di transizione dello stato f_x è una condizione sufficiente ma non necessaria per garantire la ESP.

Capitolo 4

Implementazione

Seguendo gli studi sviluppati in [5], vengono presentate le scelte progettuali fatte per l'implementazione e l'uso del *reservoir* e del *readout*.

4.1 Inizializzazione del Reservoir

Dato il modello delle RNN, il reservoir è definito dalla tupla $(\mathbf{W}_{\text{in}}, \mathbf{W}, \alpha)$. Le matrici dei pesi delle connessioni e dei pesi di input sono generate in modo random secondo dei parametri di cui si discute in seguito, e il leaking rate α è selezionato come parametro libero.

Analogamente a quanto viene fatto nella letteratura, quelli che per semplicità vengono chiamati parametri globali in questa sezione potrebbero essere chiamati iper-parametri, poiché non rappresentano direttamente i pesi delle connessioni ma regolano la loro distribuzione.

I parametri globali del reservoir sono: la dimensione N_r , la sparsità, ovvero la distribuzione degli elementi non nulli, e il raggio spettrale di \mathbf{W} ; lo *scaling* di \mathbf{W}_{in} ; ed il *leaking rate* α . I dettagli di ciascuno di queste scelte di design vengono illustrate in seguito. Per inizializzare un reservoir si invoca la funzione *weightMatrix()* con gli opportuni parametri.

```

1 function [win,w] = weightMatrix(varargin)
2
3 %initialization of default parameters
4 par.nr=100; %reservoir units
5 par.nu=1; %input dim
6 par.scale_in=0.1; %input scaling
7 par.rho=0.9; %spectral radius
8 par.alpha=1; %leaking integrator
9 par.dist='ud'; %type of distribution
10 par.density_con=1; %density of connections
11
12 %assignment of values passed as parameters
13 n_arg= length(varargin);
14 for iArg = 1:2:n_arg % considering couple of parameters
15     name_argument = varargin{iArg}; % arguments's name
16     value_argument = varargin{iArg+1}; % arguments's value
17     par.(name_argument) = value_argument;
18 end
19
20 switch par.dist
21     case 'ud'
22         %input weight matrix
23         win= (-1 + (2*rand(par.nr, par.nu+1)))*par.scale_in;
24         %connection weight matrix
25         aus= -1 + 2*rand(par.nr, par.nr);
26     case 'nd'
27         %input weight matrix
28         win= ((1/3)*randn(par.nr, par.nu+1))*par.scale_in;
29         %connection weight matrix
30         aus= (1/3)*randn(par.nr, par.nr);
31 end
32 w= (1-par.alpha).*eye(par.nr) + par.alpha.*aus;
33 w= w .* (par.rho/max(abs(eig(w)))));
34 w= (1/par.alpha) .* (w - (1-par.alpha).*eye(par.nr));
35 %setting element to zero
36 nrz= floor( (1-par.density_con)* par.nr);
37 if nrz > 0
38     for irow = 1:par.nr
39         indz= randperm(par.nr,nrz);
40         w(irow,indz) = 0;
41     end
42 end
43 end

```

/home/eleonora/Scrivania/tesi/weightMatrix.m

4.1.1 Dimensione del Reservoir

Un parametro fondamentale per il modello 3.1 è N_r , il numero delle unità del reservoir. In genere più è grande il reservoir, migliore è la performance ottenuta, tenendo conto delle tecniche di regolarizzazione appropriate per evitare l'overfitting. Poiché allenare ed applicare le ESN è computazionalmente economico rispetto alle altre reti neurali ricorrenti, non è insolito trovare dei reservoir di dimensioni dell'ordine di 10^4 .

Più è grande lo spazio dei segnali del reservoir $\mathbf{x}(n)$, più facile è trovare una combinazione lineare dei segnali per approssimare $\mathbf{y}_{target}(n)$. Il reservoir può essere troppo grande quando il task è banale oppure quando non si hanno

molti dati disponibili, ad esempio $T < 1 + N_u + N_r$.

Generalmente nell'ambito accademico la dimensione del reservoir viene limitata per convenienza e per compatibilità di risultati; inoltre buoni parametri sono trasferibili a reservoir di dimensioni maggiori, quindi può essere conveniente selezionare i parametri globali con reservoir ridotti e poi trasferirli a quelli più grandi.

Un limite inferiore per la dimensione del reservoir può essere considerato il numero di valori reali che il reservoir deve ricordare dall'input per realizzare con successo il task. Il massimo numero di valori conservati in una ESN non può essere superiore di N_r . il parametro che indica la dimensione del Reservoir passato alla funzione `weightMatrix()` è nr .

4.1.2 Sparsità del Reservoir

Nelle prime pubblicazioni riguardanti le ESN si raccomanda di realizzare reservoir con connessioni sparse, cioè con le matrici \mathbf{W}_{in} e \mathbf{W} aventi la maggior parte degli elementi uguali a zero. In generale questo parametro non influenza la performance così tanto ed ha una bassa priorità nell'ottimizzazione, come viene dimostrato anche in [4].

La sparsità fa sì che l'aggiornamento della matrice sia più rapido, velocizzando quindi la computazione. Matrici densamente connesse potrebbero rallentare la computazione, soprattutto in base al numero di unità di attivazione considerate, quando la dimensione del reservoir è grande si può decidere di generare \mathbf{W}_{in} e \mathbf{W} secondo lo stesso tipo di distribuzione ma una rispettivamente densa e sparsa. In particolare nell'implementazione fornita si può scegliere la distribuzione tra quella uniforme e quella normale associando, rispettivamente, le stringhe `'ud'` e `'nd'` al parametro `dist`. Inoltre si può esprimere la densità delle connessioni del reservoir attraverso il parametro `density_con`, il cui valore va da 0 a 1, dove con 1 si indica una matrice completamente connessa.

4.1.3 Raggio spettrale

Uno dei principali parametri delle ESN è il raggio spettrale della matrice dei pesi delle connessioni, cioè il massimo tra i valori assoluti degli autovalori della matrice. Questo parametro scala la matrice \mathbf{W} , o visto alternativamente scala l'ampiezza della distribuzione dei parametri diversi da zero.

Viene generata in modo random la matrice \mathbf{W} , viene calcolato il suo raggio spettrale $\rho(\mathbf{W})$, infine \mathbf{W} viene divisa per $\rho(\mathbf{W})$ così da ottenere una matrice con raggio spettrale unitario da poter moltiplicare per il raggio spettrale desiderato. Si ricorda che questo ultimo deve essere minore di uno affinché valga la *Echo State Property*.

Come principio guida, $\rho(\mathbf{W})$ dovrebbe essere scelto grande per i task per i quali è richiesto una storia dell'input vasta, piccolo per i task il cui output corrente dipende dalla storia recente. Il valore del raggio spettrale desiderato è espresso dal parametro *rho*.

4.1.4 Input scaling

Lo scaling dei pesi della matrice di input è un altro valore chiave da ottimizzare in una Echo State Network. Per \mathbf{W}_{in} generate con distribuzione uniforme generalmente ci si riferisce all'*input scaling a* come al modulo degli estremi dell'intervallo $[-a; a]$ dal quale vengono selezionati i valori di \mathbf{W}_{in} ; per le matrici generate con distribuzione normale si prende la deviazione standard come misura di scaling.

In questa implementazione, per avere un numero esiguo di parametri liberi, tutte le colonne di \mathbf{W}_{in} vengono scalate insieme in base ad un unico valore. In alternativa si può ottimizzare separatamente il valore di *input scaling* relativo alla colonna del *bias*, oppure si possono scalare le colonne di \mathbf{W}_{in} separatamente in base a come queste contribuiscono al task. A seconda delle impostazioni il numero di parametri liberi per \mathbf{W}_{in} varia da 1 a $N_u + 1$.

Nelle prime pubblicazioni venne suggerito di scalare e traslare i dati di input, ottimizzando queste operazioni. Lo stesso effetto può essere raggiunto sca-

lando i pesi di input e il *bias*, rispettivamente. Queste operazioni sono molto utili poiché permettono di avere dei valori limitati per i dati di input.

Per task molto lineari \mathbf{W}_{in} dovrebbe essere piccolo, permettendo alle unità di operare intorno allo zero dove la funzione di attivazione *tanh*, funzione di attivazione scelta in questa implementazione, è quasi lineare. Per \mathbf{W}_{in} con valori grandi, le unità si saturano facilmente vicino al loro valore -1 o 1 , agendo in modo non lineare. È chiaro che l'*input scaling*, insieme al valore $\rho(\mathbf{W})$, determina quanto lo stato corrente $\mathbf{x}(t)$ dipende dall'input corrente $\mathbf{u}(t)$ e quanto dallo stato precedente $\mathbf{x}(t-1)$. Nell'implementazione il valore dell'input scaling viene indicato da *scale_in*.

4.1.5 Leaking Rate

La variabile *alpha* esprime il valore del *leaking rate* come visto in 3.3 ed indica la velocità con cui vengono aggiornate le dinamiche del reservoir. Sebbene questo parametro viene maggiormente utilizzato nel calcolo dello stato del reservoir viene presentato qui poiché si deve tener conto di questo valore nel momento in cui viene riscalata la matrice dei pesi \mathbf{W} come dimostrato in [6].

4.2 Calcolo dello stato del reservoir

Per calcolare lo stato del reservoir è necessario invocare la funzione, riportata sotto, *setReservoir()*. Questa prende come parametri le informazioni relative al task e le matrici che rappresentano il reservoir, inizializzate come espresso sopra. Può prendere in ingresso dei parametri opzionali che rappresentano rispettivamente il valore del *leaking rate* e dell'*input bias*. Viene applicata la 3.3 per il calcolo dello stato del reservoir e viene eliminato il numero di stati transient come espresso nella definizione del task.

Ciò significa che i valori iniziali $\mathbf{x}(n)$ non vengono considerati per l'allenamento del readout e vengono scartati in quanto influenzati dallo stato iniziale $\mathbf{x}(0)$, che tipicamente è $\mathbf{x}(0) = 0$. Questo introduce uno stato iniziale non

naturale difficile da raggiungere quando la rete comincia ad adattarsi al task. Il numero di step da saltare dipende dalla memoria della rete e tipicamente è dell'ordine di decine di centinaia. Tuttavia se il task è composto da molte sequenze corte il transient iniziale è la normale modalità di lavoro della ESN ed in questi casi eliminare gli stati iniziali può essere pericoloso. Invece questo aspetto diventa fondamentale se si vogliono rendere indipendenti le sequenze, ad esempio quelle associate a fold diversi.

```

1 function [X] = setReservoir(task,win,w,varargin)
2
3     %initialization of default parameters
4     par.nr=size(w,1);           %reservoir units
5     par.alpha=1;                %leaking integrator
6     par.bias=0;                 %bias
7
8     %assignment of values passed as parameters
9     n_arg= length(varargin);
10    for iArg = 1:2:n_arg         % considering couple of parameters
11        name_argument = varargin{iArg}; % arguments's name
12        value_argument = varargin{iArg+1}; % arguments's value
13        par.(name_argument) = value_argument;
14    end
15
16    %training and test elements without transient for each fold
17    training= task.training -(ceil(task.training/task.timesteps)*task.transient);
18    design= task.design -(ceil(task.design/task.timesteps)*task.transient);
19    test= task.test -(ceil(task.test/task.timesteps)*task.transient);
20    % design + test element
21    len= task.design+task.test;
22    % reservoir states
23    X= zeros(par.nr,task.kfold*(design+test));
24    %target
25    T= cell(1,task.readouts(end));
26
27    % reservoir initialization
28    aus=1;
29    numseq=1;
30    for i = 1:task.kfold
31        taskin= task.in{i};
32        x = zeros(par.nr,1);
33        trans=task.transient;
34        for n = 1:len
35            u= taskin(:,n);
36            x= (1-par.alpha)* x + par.alpha*tanh(win*[u;par.bias] + w*x );
37            if n > task.timesteps*numseq
38                numseq=numseq+1;
39                trans=trans+task.timesteps;
40            end
41            %delete transient for each input sequence
42            if n>trans
43                %reservoir states
44                X(:,aus)=x;
45                aus=aus+1;
46            end
47        end
48    end
49
50
51 end

```

/home/eleonora/Scrivania/tesi/setReservoir.m

4.3 Allenamento del Readout

Per questa implementazione si prende in considerazione un readout lineare, senza connessioni di *feedback*.

Un readout lineare può essere descritto con la seguente equazione

$$\mathbf{Y} = \mathbf{W}_{\text{out}}\mathbf{X} \quad (4.1)$$

dove:

- $\mathbf{Y} \in \mathbb{R}^{N_y \times T}$ sono tutti gli $\mathbf{y}(n)$;
- $\mathbf{X} \in \mathbb{R}^{(1+N_r) \times T}$ sono tutti gli $[\mathbf{x}(n); 1]$
- \mathbf{W}_{out} sono i pesi ottimi che minimizzano l'errore tra $\mathbf{y}(n)$ e $\mathbf{y}_{\text{target}}(n)$

Nel caso in cui si volessero connettere i pesi di input al readout si avrebbe la matrice $\mathbf{X} \in \mathbb{R}^{(1+N_u+N_r) \times T}$, formata da tutti i $[\mathbf{x}(n); \mathbf{u}(n); 1]$.

Trovare i pesi ottimi \mathbf{W}_{out} che minimizzino l'errore tra $\mathbf{y}(n)$ e $\mathbf{y}_{\text{target}}(n)$ consiste nel risolvere un sistema di equazioni lineari del tipo:

$$\mathbf{Y}_{\text{target}} = \mathbf{W}_{\text{out}}\mathbf{X} \quad (4.2)$$

Per allenare il readout si deve invocare la funzione *train_readout()*, nelle sottosezioni successive vengono illustrati due approcci diversi per la risoluzione del sistema.

4.3.1 Regolarizzazione di Tikhonov

La più diffusa e stabile soluzione per 4.2 in questo contesto è la regolarizzazione di Tikhonov o Ridge Regression, in formula:

$$\mathbf{W}_{\text{out}} = \mathbf{Y}_{\text{target}}\mathbf{X}^T(\mathbf{X}\mathbf{X}^T + \beta\mathbf{I})^{-1} \quad (4.3)$$

dove β è il coefficiente di regolarizzazione.

Per misurare la qualità della soluzione prodotta dall'allenamento è consigliabile controllare i pesi ottenuti di \mathbf{W}_{out} , grandi pesi indicano che \mathbf{W}_{out}

amplifica piccole differenze tra le dimensioni di $\mathbf{x}(t)$ e può essere molto sensibile nelle situazioni diverse dalle esatte condizioni nelle quali la rete è stata allenata. Questo problema si accentua nel caso in cui la rete riceve il suo output come successivo input.

Per contrastare questi effetti viene introdotta la parte di regolarizzazione $\beta \mathbf{I}$. Per illustrare l'equazione risolta con la *ridge regression* viene considerato il RMSE (*root mean squared error*) come misura di errore:

$$\mathbf{W}_{\text{out}} = \arg \min_{\mathbf{W}_{\text{out}}} \frac{1}{N_y} \sum_{i=1}^{N_y} \left(\sum_{n=1}^T (y_i(n) - y_{i_{\text{target}}}(n))^2 + \beta \|w_{i_{\text{out}}}\|^2 \right) \quad (4.4)$$

con $w_{i_{\text{out}}}$ che indica la i -esima riga di \mathbf{W}_{out} e $\|\cdot\|$ la norma euclidea. È evidente il compromesso che si ha tra avere un basso *training error* e pesi di output piccoli, ed è regolato proprio dal parametro di regolarizzazione β . Il coefficiente di regolarizzazione ottimo β dipende dall'istanza di \mathbf{W} , quindi è buona norma scegliere questo parametro attraverso la validazione.

Settare β uguale a *zero* rimuove la regolarizzazione: la funzione in 4.4 diventa uguale al semplice calcolo del RMSE (*root mean squared error*), rendendo la regolarizzazione di Tikhononov una generalizzazione della regressione lineare. La soluzione con $\beta = 0$ diventa:

$$\mathbf{W}_{\text{out}} = \mathbf{Y}_{\text{target}} \mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^{-1} \quad (4.5)$$

Nella pratica, tuttavia, fissare $\beta = 0$ spesso causa instabilità numerica quando si inverte $(\mathbf{X} \mathbf{X}^T)$ in 4.5, è per questo che si raccomanda di usare per la selezione di β una scala logaritmica che non raggiunge lo zero o di usare la pseudoinversa invece che l'inversa come mostrato di seguito.

4.3.2 Soluzione con pseudoinversa

Una soluzione semplice per la risoluzione del sistema visto prima è:

$$\mathbf{W}_{\text{out}} = \mathbf{Y}_{\text{target}} \mathbf{X}^+ \quad (4.6)$$

dove \mathbf{X}^+ è la pseudoinversa di Moore-Penrose di \mathbf{X} . Se $\mathbf{X}\mathbf{X}^T$ è invertibile questa formula diventa uguale alla 4.5, ma funziona anche quando non lo è. Tuttavia, ha un elevato costo in memoria per grandi matrici \mathbf{X} , dunque si deve limitare la dimensione del reservoir e/o il numero di esempi di training T . Poichè non si ha la regolarizzazione il sistema di equazioni lineari 4.2 deve essere sovraddeterminato, cioè $1 + N_u + N_r \ll T$. In altre parole, il task deve essere difficile in relazione alla capacità del reservoir così che non avvenga l'*overfitting*.

In molte librerie viene fornita una implementazione della pseudoinversa, tuttavia ogni implementazione varia nella precisione, nell'efficienza computazionale e nella stabilità numerica. Per task di alta precisione, si deve controllare se la regressione $(\mathbf{Y}_{\text{target}} - \mathbf{W}_{\text{out}}\mathbf{X})\mathbf{X}^+$ sull'errore $\mathbf{Y}_{\text{target}} - \mathbf{W}_{\text{out}}\mathbf{X}$ è effettivamente uguale a zero, e aggiungerla a \mathbf{W}_{out} nel caso non lo sia. Questo trucco computazionale non dovrebbe funzionare in teoria ma a volte funziona in pratica in Matlab, probabilmente a causa di qualche ottimizzazione interna.

Capitolo 5

Test

In questo capitolo viene testata la rete neurale implementata attraverso il task NARMA e vengono confrontati i risultati ottenuti con quelli della letteratura.

5.1 10th order NARMA system

Questo task consiste nella predizione di un sistema di 10th ordine di media mobile autoregressivo non lineare. Il task è stato introdotto in Atiya e Palos (2000) ed è stato trattato con le ESN in Jaeger (2002) e in Cernanský e Tiño (2008). L'input del sistema è una sequenza di elementi $\mathbf{u}(n)$ scelti secondo una distribuzione uniforme in $[0, 0.5]$. L'output del sistema è calcolato come:

$$\bar{\mathbf{y}}(n) = 0.3\bar{\mathbf{y}}(n-1) + 0.05\bar{\mathbf{y}}(n-1) \left(\sum_{i=1}^{10} \bar{\mathbf{y}}(n-i) \right) + 1.5\mathbf{u}(n-10)\mathbf{u}(n-1) + 0.1. \quad (5.1)$$

Dato il valore di input $\mathbf{u}(n)$, il task è di predire il corrispondente valore di $\bar{\mathbf{y}}(n)$. Il training set è formato da $N_{train} = 2200$ esempi input-target, dei quali $N_{transient} = 200$ sono di transient iniziale. Una sequenza di lunghezza $N_{test} = 2000$ viene usata per il test.

Il task viene generato attraverso la funzione *genetateTask()*, passando come argomento *Tasks.narma*, questa funzione permette di generare tutti i task che

sono elencati nell'enumerazione *Tasks*, questi ultimi vengono creati secondo dei criteri ben precisi in modo da poter fornire alla rete tutte le informazioni necessarie. Per verificare che un task abbia tutti i campi necessari per essere ben definito è utilizzata la funzione *isTask()*.

5.2 ESNtrain() ed ESNtest()

Per facilitare l'uso della rete sviluppata viene fornita una funzione chiamata *ESNtrain()*. La funzione si occupa dell'inizializzazione delle matrici dei pesi, del calcolo dello stato del reservoir e del training del readout invocando le funzioni illustrate nel capitolo 5. A questa funzione devono essere passati come argomenti un task valido e tutti i parametri che si intendono usare nella rete, qualora questi non fossero presenti vengono utilizzati dei parametri standard illustrati nella tabella sottostante.

parametri	nome parametro	valore default
input scaling	<i>scale_in</i>	0.1
unità reservoir	<i>nr</i>	100
tipo di distribuzione	<i>dist</i>	'ud'
densità connessione	<i>density_con</i>	1
raggio spettrale	<i>rho</i>	0.9
leaky integrator	<i>alpha</i>	1
input bias	<i>bias</i>	1
par regolazizzazione	<i>lambda</i>	0
misura errore	<i>error</i>	'mse'
risultati test	<i>test</i>	true

Tabella 5.1: Parametri default *ESNtrain()*

In particolare si ha il parametro *test* che se settato a true permette di ottenere anche il risultato di test, eseguito dopo il training. Se viene usata questa opzione bisogna far attenzione a non usare il valore di test per la scelta del modello, in questo caso infatti la *model selection* verrebbe compromessa.

Per il test della rete deve essere usata la funzione `ESNtest()`. Questa funzione prende in ingresso il task ed eventualmente altri parametri per la realizzazione della rete, infatti si può decidere se passare o meno come parametri un reservoir ed un readout allenato. Nel primo caso viene eseguita solo una fase di test, altrimenti viene inizializzata un'altra rete secondo i parametri passati come argomento; quest'ultima operazione assume un certo valore soprattutto nel caso in cui i parametri passati come argomento sono quelli di una model selection precedente.

5.3 Risultati

Per poter ottenere dei risultati confrontabili con quelli in [4] viene fatta una model selection al variare del parametro ρ , ovvero del raggio spettrale. Nelle figure successive sono mostrati in ordine i valori dell'errore di train, dell'errore di test ottenuto con la funzione `ESNtrain()`, di quello di test ottenuto con la funzione `ESNtest()` avente come parametri rispettivamente quelli ottenuti della model selection ed il reservoir ed il readout prodotti in precedenza.

Si può notare come l'errore diminuisca al crescere del valore del raggio spettrale, andamento aspettato nel caso di task non lineare come questo.

5.4 Conclusioni

Bibliografia

- [1] Mantas Lukoševičius, Herbert Jaeger, Reservoir computing approaches to recurrent neural network training, *Computer Science Review*, Volume 3, Issue 3, 2009, Pages 127-149.
- [2] Wolfgang Maass, Thomas Natschläger, Henry Markram, Real-time computing without stable states: A new framework for neural computation based on perturbations, *Neural Computation* 14 (11) (2002) 2531–2560.
- [3] Herbert Jaeger, The “echo state” approach to analysing and training recurrent neural networks, Technical Report GMD Report 148, German National Research Center for Information Technology, 2001.
- [4] Claudio Gallicchio, Alessio Micheli. Architectural and Markovian factors of echo state networks, *Neural Networks*, Volume 24, Issue 5, 2011, Pages 440-456.
- [5] Lukoševičius M. (2012) A Practical Guide to Applying Echo State Networks. In: Montavon G., Orr G.B., Müller KR. (eds) *Neural Networks: Tricks of the Trade*. Lecture Notes in Computer Science, vol 7700. Springer, Berlin, Heidelberg.
- [6] Herbert Jaeger, Mantas Lukoševičius, Dan Popovici, Udo Siewert, Optimization and applications of echo state networks with leaky-integrator neurons, *Neural Networks*, Volume 20, Issue 3, 2007, Pages 335-352,

Appendice A

Codice

```
num(0).  
num(s(X)) :- num(X).
```