

1 Storia

Uno dei primi esempi di rete neurale fu il percettrone di Rosenbatt (1958). Di seguito fu creata ADELIN (Widrow & Hoff, 1960), composta da un singolo layer con un singolo neurone e utilizzava il *LMS* (least mean square) insieme ad un algoritmo di discesa stocastica del gradiente per la configurazione dei pesi. Ne descriviamo adesso le caratteristiche. Siano $x \in R^N$ gli input e $w \in R^N$ i pesi; dal modello abbiamo:

$$\hat{y} = x \cdot w. \quad (1)$$

Dato l'output desiderato $y \in R$ possiamo definire la funzione errore e e la *loss function* L . L'obiettivo é quello di trovare il valore dei pesi che minimizzi L .

$$e = y - \hat{y} \quad (2)$$

$$L = e^2 \quad (3)$$

$$w \leftarrow \arg \min(L). \quad (4)$$

L'algoritmo *LMS* ottimizza L usando il gradiente discendente:

$$\nabla L_w := -2ex \quad (5)$$

$$w \leftarrow w + \alpha ex \quad (6)$$

dove α é il tasso di apprendimento.

Widrow, McCool and Ball (1975) hanno esteso l'algoritmo *LMS* al dominio complesso fornendo la derivazione delle parti reali e immaginarie. Brandwood (1983) generalizzò la teoria applicando il gradiente al numero complesso, senza separarlo in parte reale e parte immaginaria attraverso il gradiente di Wirtinger (1927).

Aggiorniamo quindi il problema nel dominio complesso:

$$L = e\bar{e} \quad (7)$$

$$\nabla L_w := -2e\bar{x} \quad (8)$$

$$w \leftarrow w + \alpha e\bar{x}. \quad (9)$$

Nonostante i risultati di Brandwood, fino a non molto tempo fa la letteratura non applicava il calcolo di Wirtinger, a favore della derivazione separata della parte reale da quella immaginaria.

2 Funzione d'attivazione

Definiamo l'input $x \in C^M$ e i pesi $\mathbf{W} \in C^{N \times M}$, con M e N rispettivamente le dimensioni di input e output; di conseguenza l'output $y \in C^N$ di un qualsiasi layer é determinato da:

$$z = \mathbf{W}x \quad (10)$$

$$y = f(z) \quad (11)$$

dove f solitamente é una funzione di attivazione non lineare.

3 Loss function

La maggior parte della letteratura attuale utilizza come funzione di costo la *mean squarred error*. Dato un target y e l'output ottenuto \hat{y} , entrambi in C^N e l'errore

$$e := y - \hat{y} \quad (12)$$

la *complex mean squared loss function* é definita come segue:

$$L(e) = \sum_{i=0}^{N-1} |e_i|^2 \quad (13)$$

$$= \sum_{i=0}^{N-1} e_i \overline{e_i}. \quad (14)$$

La (??) é una funzione a valori reali scalari non negativi, che tende a zero insieme al modulo dell'errore. Savitha, Suresh e Sundararajan proposero di sostituire l'errore (??) con:

$$e := \log \hat{y} - \log y. \quad (15)$$

La *loss function* diventerebbe quindi:

$$L(e) = (\log |\hat{y}_i| - \log |y_i|)^2 + (\arg \hat{y}_i - \arg y_i)^2 \quad (16)$$

L'equazione (??) ha la proprietà di rappresentare esplicitamente l'ampiezza e la fase.

Le equazioni (??) e (??) possono essere *error function* appropriate per reti neurali complesse.

4 Ottimizzazione

4.1 Il gradiente complesso ed il calcolo di Wirtinger

Brandwood e Van den Bos formularono le prime derivazioni del gradiente complesso. Wirtinger (1927) fornì un formalismo equivalente che rese il calcolo della derivata di funzioni con valori complessi meno oneroso rispetto a funzioni olomorfe e non analitiche, agendo interamente nel campo complesso. Nonostante tale apparente comodità, solo recentemente si ricominciò ad utilizzare il calcolo di Wirtinger per la *backpropagation* di reti neurali complesse.

Definiamo:

$$f(z) := f(z, \bar{z}) \quad (17)$$

$$= g(x, y) \quad (18)$$

$$= u(x, y) + iv(x, y) \quad (19)$$

con $z \in C$, $x, y \in R$ e $z = x + iy$.

Usando la prima definizione avremo le derivate in z e \bar{z} date da:

$$\left. \frac{\partial f}{\partial z} \right|_{\bar{z} \text{ costante}} \quad (20)$$

$$\left. \frac{\partial f}{\partial \bar{z}} \right|_{z \text{ costante}} \quad (21)$$

le quali, espresse in funzione di x e y , diventano:

$$\frac{\partial f}{\partial z} = \frac{1}{2} \left(\frac{\partial f}{\partial x} - i \frac{\partial f}{\partial y} \right) \quad (22)$$

$$\frac{\partial f}{\partial \bar{z}} = \frac{1}{2} \left(\frac{\partial f}{\partial x} + i \frac{\partial f}{\partial y} \right). \quad (23)$$

Si osservi che la derivata parziale rispetto a \bar{z} è nulla per ogni funzione olomorfa. Richiamiamo le condizioni di esistenza di Cauchy-Riemann per la derivata complessa della funzione $f(z, \bar{z})$ presa in considerazione:

$$\frac{\partial u}{\partial x} = \frac{\partial v}{\partial y} \quad (24)$$

$$\frac{\partial v}{\partial x} = -\frac{\partial u}{\partial y}. \quad (25)$$

Se applichiamo le condizioni di Cauchy-Riemann alla derivata di f in \bar{z} (??) notiamo che effettivamente essa si annulla. Le funzioni olomorfe quindi non dipendono esplicitamente da \bar{z} . Brandwood dimostrò che l'annullarsi di (??) o (??) per una generica $f : C \rightarrow R$ è condizione sufficiente e necessaria affinché f abbia un punto stazionario. Per estensione se $f : C^N \rightarrow R$ è una funzione a valori reali di un certo vettore $z = [z_0, z_1, \dots, z_{N-1}]^T \in C^N$ e definiamo il cogradiente e il gradiente coniugato come:

$$\frac{\partial}{\partial z} := \left[\frac{\partial}{\partial z_0}, \frac{\partial}{\partial z_1}, \dots, \frac{\partial}{\partial z_{N-1}} \right] \quad (26)$$

$$\frac{\partial}{\partial \bar{z}} := \left[\frac{\partial}{\partial \bar{z}_0}, \frac{\partial}{\partial \bar{z}_1}, \dots, \frac{\partial}{\partial \bar{z}_{N-1}} \right] \quad (27)$$

allora $\frac{\partial f}{\partial z} = 0$ o $\frac{\partial f}{\partial \bar{z}} = 0$ sono condizioni sufficienti e necessarie per determinare un punto di stazionarietà.

Se f è una funzione di un vettore complesso z , la sua derivata totale è:

$$df = \frac{\partial f}{\partial z} dz + \frac{\partial f}{\partial \bar{z}} d\bar{z}. \quad (28)$$

Se f è reale allora avremmo:

$$df = 2 \operatorname{Re} \left\{ \frac{\partial f}{\partial z} dz \right\}. \quad (29)$$

Definendo ora l'operatore gradiente come:

$$\nabla_z := \left(\frac{\partial}{\partial \bar{z}} \right)^T \quad (30)$$

$$= \left(\frac{\partial}{\partial z} \right)^* \quad (31)$$

può essere dimostrato, usando la disuguaglianza di Cauchy-Schwarz, che f ha il maggior tasso di cambiamento lungo il gradiente. Grazie a queste definizioni possiamo costruire una *cost function* reale con argomenti complessi anche se

alcuni elementi della funzione non sono olomorfi. In generale, date le funzioni arbitrarie f e g , combinando gli jacobiani come segue

$$J_{f \circ g} = J_f J_g + J_f^{(c)} \overline{J_g^{(c)}} \quad (32)$$

$$J_{f \circ g}^{(c)} = J_f J_g^{(c)} + J_f^{(c)} \overline{J_g}. \quad (33)$$

Supponiamo di avere una funzione composta $(f \circ g \circ h)(z, \bar{z})$, con f la *cost function* reale, g una funzione complessa non olomorfa e h una funzione olomorfa. tenendo a mente che f é una funzione a valori reali di variabile complessa e che h é olomorfa ($\frac{\partial h}{\partial \bar{z}} = 0$), applichiamo la *chain rule*:

$$J_f = \frac{\partial f}{\partial g} \quad (34)$$

$$J_f^{(c)} = \frac{\partial f}{\partial \bar{g}} \quad (35)$$

$$J_{f \circ g} = J_f \frac{\partial g}{\partial h} + J_f^{(c)} \overline{\left(\frac{\partial g}{\partial \bar{h}} \right)} \quad (36)$$

$$J_{f \circ g \circ h} = J_{f \circ g} \frac{\partial h}{\partial z} \quad (37)$$

$$\nabla_z f = (J_{f \circ g \circ h})^* \quad (38)$$

Il Calcolo di Wirtinger rende un po' piú semplice costruire un grafico computazionale per reti complesse aventi composizioni miste di operazioni olomorfe e non olomorfe.

4.2 Inputs

La funzione che la rete neurale deve apprendere ha come argomenti i vettori reali α e φ che rappresentano rispettivamente le componenti di ampiezza e fase. Di seguito mostriamo alcune modalità per rappresentare l'input:

Ampiezza-Fase I parametri di ampiezza e fase vengono concatenati nel vettore reale:

$$x^{(ap)} := [\dots, \alpha_i, \varphi_i, \dots]^T \in R^{2K} \quad (39)$$

Complessa I parametri vengono rappresentati all'interno del vettore complesso:

$$x^{(c)} := [\dots, \alpha_i e^{i\varphi_i}, \dots]^T \in C^K \quad (40)$$

Reale-Immaginaria Possiamo scomporre il vettore (??) in parte reale e parte immaginaria e ottenere un input nella forma:

$$x^{(ri)} := \left[Re \left\{ x^{(c)} \right\}^T, Im \left\{ x^{(c)} \right\}^T \right] \quad (41)$$

$$= [\dots, \alpha_i \cos \varphi_i, \dots, \alpha_i \sin \varphi_i, \dots]^T \quad (42)$$

Complesso aumentata L'input é un vettore complesso contenente sia (??) che il suo coniugato:

$$x^{(ca)} := \left[\left(x^{(c)} \right)^T, \left(x^{(c)} \right)^* \right]^T \quad (43)$$

$$= [\dots, \alpha_i e^{i\varphi_i}, \dots, \alpha_i e^{-i\varphi_i}, \dots]^T \quad (44)$$

Il *mapping* definito dal complesso coniugato é antilineare e di conseguenza non può essere calcolato attraverso la moltiplicazione di matrici complesse. ciò potrebbe richiedere alla rete l'apprendimento di un ulteriore parametro. In effetti, l'utilizzo del vettore complesso aumentato fornisce un *path* più semplice per la modellizzazione della distribuzione completa al secondo ordine dei dati. Gli inut aumentati sono utilizzati in una stima lineare.

4.3 Architettura

Esaminiamo un modello *feedforward* avente un unico hidden layer; abbiamo di conseguenza:

$$h = f \left(W^{(h)} x + b^{(h)} \right) \quad (45)$$

$$\hat{y} = W^{(o)} x + b^{(o)} \quad (46)$$

con $\theta = \{W^{(h)}, W^{(o)}, b^{(h)}, b^{(o)}\}$ sono i parametri del modello. Definendo M il numero dei nodi dell'input e N il numero dei nodi dell'output, avremo:

$$\begin{array}{ll} W^{(h)} \in R^{M \times M} & or \quad W^{(o)} \in C^{M \times M} \\ b^{(h)} \in R^M & or \quad b^{(o)} \in C^M \\ W^{(o)} \in R^{N \times M} & or \quad W^{(o)} \in C^{N \times M} \\ b^{(o)} \in R^N & or \quad b^{(o)} \in C^N \end{array} \quad (47)$$

In tutti i casi gli inut, i pesi e gli output appartengono allo stesso dominio numerico.

4.4 Funzioni di attivazione

Identità Permette una modellizzazione lineare

$$f^{(-)}(x) := x \quad (48)$$

Tangente Iperbolica É una funzione sigmoidale e differenziabile. Permette una non linearità ampiamente utilizzata per l'apprendimento delle reti neurali

$$f^{(s)} := \tanh(x) \quad (49)$$

Split reale-immaginario Applica la tangente iperbolica separatamente alla parte reale e alla parte immaginaria:

$$f^{(ri)}(x) := \tanh(\operatorname{Re} x) + i \tanh(\operatorname{Im} x) \quad (50)$$

Split ampiezza-fase Applica la tangente iperbolica al modulo del numero complesso, senza modificarne la fase (questa funzione non é differenziabile nel campo complesso)

$$f^{(ap)}(x) := \tanh(|x|) e^{i \arg x} \quad (51)$$

ModReLU Arjosky nel 2015 propose una variazione alla classica ReLU utilizzata nelle reti neurali reali, definita come segue:

$$ModReLU(z) = ReLU(|z| + b) e^{i\theta_z} = \begin{cases} (|z| + b) \frac{z}{|z|} & \text{se } |z| + b \geq 0 \\ 0 & \text{altrimenti} \end{cases} \quad (52)$$

dove θ_z é la fase di z e b il bias, apprendibile dalla rete neurale, e necessario per creare una *dead zone* di raggio b attorno all'origine dove il neurone é inattivo. Tale funzione non soddisfa le equazioni di Cauchy-Riemann e quindi non é olomorfa.

CReLU Consiste in una applicazione della funzione ReLU individualmente alla parte reale e alla parte immaginaria:

$$CReLU(z) = ReLU(Re(z)) + i ReLU(Im(z)). \quad (53)$$

Questa soddisfa le condizioni di Cauchy-Riemann solo se sia la parte reale che la parte immaginaria sono contemporaneamente strettamente positive o strettamente negative, quindi nell'intervallo $\theta_z \in (0, \frac{\pi}{2})$ oppure $\theta_z \in (\pi, \frac{3\pi}{2})$.

zReLU Proposta nel 2016 da Guberman e basata anch'essa sulla ReLU, é definita come:

$$zReLU(z) = \begin{cases} z & \text{se } \theta_z \in [0, \frac{\pi}{2}] \\ 0 & \text{altrimenti} \end{cases} \quad (54)$$