

Towards a Categorical Foundation of Deep Learning: A Survey

Una rassegna di approcci categorici al *deep learning*

Francesco Riccardo Crescenzi

11 settembre 2024

Alma mater studiorum - Università di Bologna
CdL in Matematica

**We are in an AI summer,
but is winter coming?**

Problemi con il deep learning

Mancano fondamenta teoriche:

- approcci *ad hoc*

Mancano fondamenta teoriche:

- approcci *ad hoc*
- complessità fine a se stessa

Mancano fondamenta teoriche:

- approcci *ad hoc*
- complessità fine a se stessa
- assenza di garanzie di correttezza

La ricerca viene rallentata da:

- *research debt*

La ricerca viene rallentata da:

- *research debt*
- mancata replicabilità

Teoria delle categorie:

una lingua franca della matematica

La teoria delle categorie studia strutture e relazioni, e può essere vista come un'estensione del celebre Erlangen Programme.

Teoria delle categorie:

una lingua franca delle scienze

La teoria delle categorie può essere applicata con successo anche in fisica, informatica, chimica... ovunque ci sia **composizionalità**.

- ottiche parametriche

- ottiche parametriche
- (co)algebre categoriche

- ottiche parametriche
- (co)algebre categoriche
- *integral transforms*

- ottiche parametriche
- (co)algebre categoriche
- *integral transforms*
- *functor learning*

- ottiche parametriche
- (co)algebre categoriche
- *integral transforms*
- *functor learning*
- *string diagrams*

Lenti parametriche

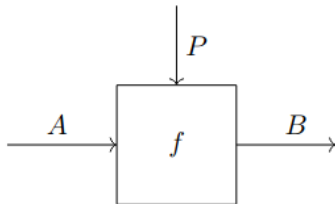
per modellare il gradient-based learning

DEFINIZIONE: Il costrutto Para

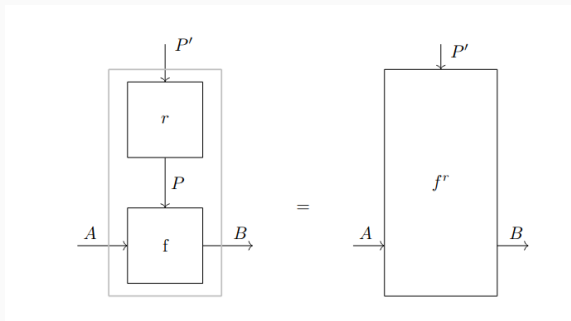
Sia $(\mathcal{C}, I, \otimes)$ una categoria monoidale strettamente simmetrica. Allora, $\mathbf{Para}_{\otimes}(\mathcal{C})$ è la 2-categoria definita come segue.

- Le 0-celle sono oggetti di \mathcal{C} .
- Le 1-cells sono coppie $(P, f) : A \rightarrow B$, dove $P : \mathcal{C}$ e $f : P \otimes A \rightarrow B$.
- The 2-celle sono $r : (P, f) \Rightarrow (Q, g)$, dove $r : P \rightarrow Q$ è un morfismo in \mathcal{C} che rispetta certe condizioni di naturalità.

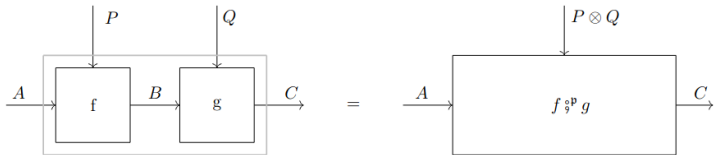
Gradient-based learning con lenti parametriche



Gradient-based learning con lenti parametriche



Gradient-based learning con lenti parametriche

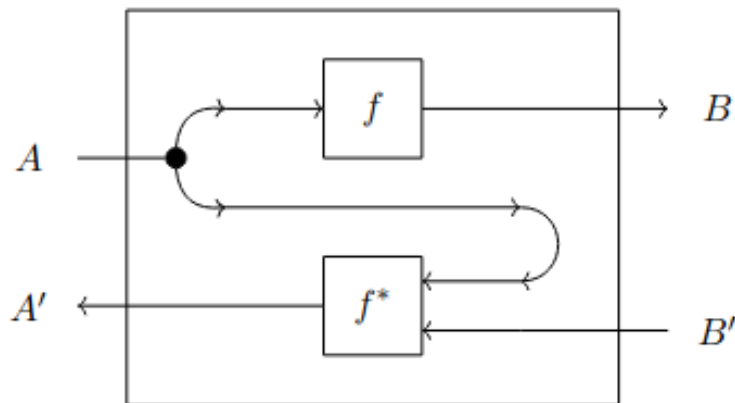


DEFINIZIONE: Il costrutto Lens

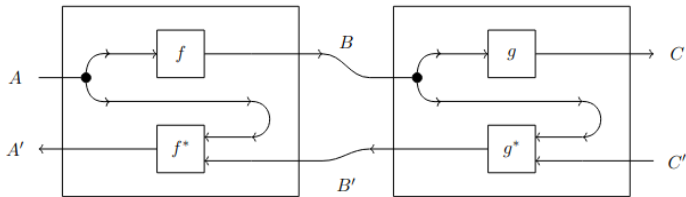
Sia $(\mathcal{C}, 1, \times)$ una categoria Cartesiana. Allora, **Lens** (\mathcal{C}) è la categoria definita come segue.

- Un oggetto di **Lens** (\mathcal{C}) è una coppia $\left(\begin{smallmatrix} A \\ A' \end{smallmatrix}\right)$ di oggetti di \mathcal{C} .
- Un morfismo $\left(\begin{smallmatrix} A \\ A' \end{smallmatrix}\right) \rightarrow \left(\begin{smallmatrix} B \\ B' \end{smallmatrix}\right)$ (anche chiamato lente) è una coppia $\left(\begin{smallmatrix} f \\ f' \end{smallmatrix}\right)$ di morfismi di \mathcal{C} tali che $f : A \rightarrow B$ and $f' : A \times B' \rightarrow A'$. La mappa f è nota come *forward pass* della lente, mentre la mappa f' è nota come *backward pass*.

Gradient-based learning con lenti parametriche



Gradient-based learning con lenti parametriche



DEFINIZIONE: Cartesian reverse differential category

Una *Cartesian reverse differential category* (CRDC) \mathcal{C} è una categoria Cartesiana con una struttura additiva dove è definito un operatore differenziale R che ha le proprietà di una *reverse derivative*.

ESEMPIO: Smooth

Consideriamo **Smooth**, ovvero la categoria degli spazi Euclidei e delle funzioni lisce. **Smooth** è una CRDC rispetto all'operatore

$$R[f] : (x, y) \mapsto \mathcal{J}_f(x)^T y.$$

DEFINIZIONE: Lenti con backward pass additivo

Sia \mathcal{C} una CRDC. Allora, definiamo la sottocategoria $\mathbf{Lens}_A(\mathcal{C})$ di $\mathbf{Lens}_A(\mathcal{C})$, i cui oggetti sono coppie $\begin{pmatrix} A \\ A' \end{pmatrix}$ e i cui morfismi hanno la forma $\begin{pmatrix} f \\ \mathbf{R}[f] \end{pmatrix}$.

TEOREMA: Struttura cartesiana di $\mathbf{Lens}_A(\mathcal{C})$

La struttura

$$I = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad \begin{pmatrix} A \\ A \end{pmatrix} \otimes \begin{pmatrix} B \\ B \end{pmatrix} = \begin{pmatrix} A \times B \\ A \times B \end{pmatrix}$$

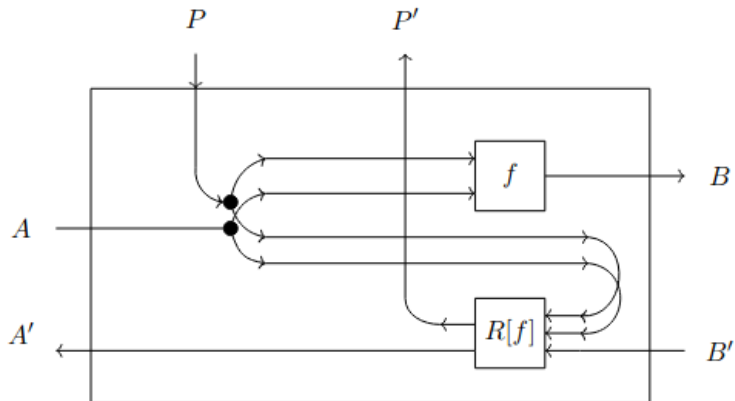
definita su $\mathbf{Lens}_A(\mathcal{C})$ è Cartesiana.

DEFINIZIONE: Lenti parametriche

Sia \mathcal{C} una CRDC. Allora, definiamo la categoria delle lenti parametriche su \mathcal{C} come

$$\mathbf{Para}_{\otimes}(\mathbf{Lens}_A(\mathcal{C})).$$

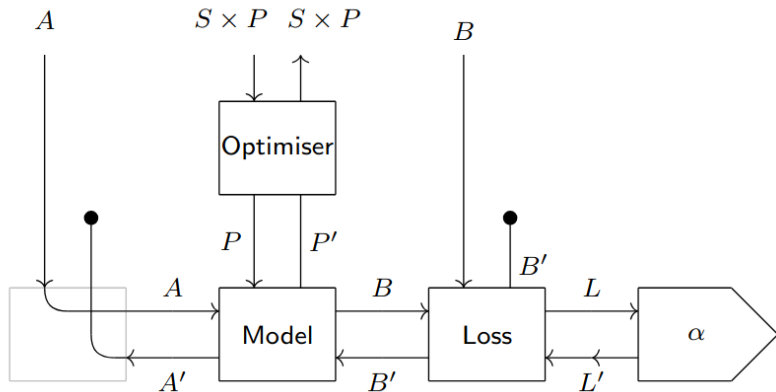
Gradient-based learning con lenti parametriche



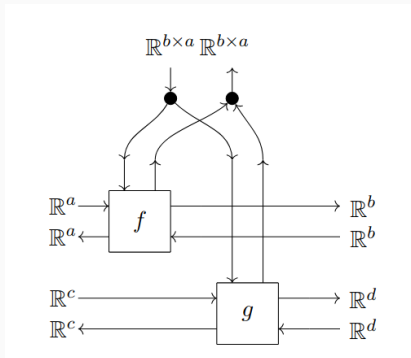
Gradient-based learning con lenti parametriche

Le lenti parametriche in $\mathbf{Para}_{\otimes}(\mathbf{Lens}_A(\mathcal{C}))$ supportano la *automatic differentiation* e possono essere utilizzate per implementare il *gradient-based learning*.

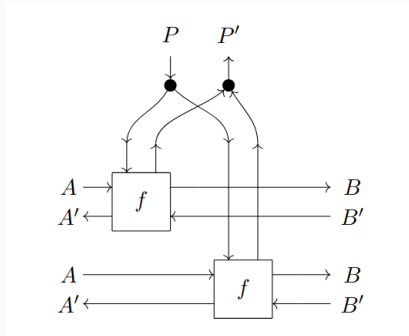
Gradient-based learning con lenti parametriche



Gradient-based learning con lenti parametriche



Gradient-based learning con lenti parametriche



Categorical deep learning:

(co)algebre categoriche

come teoria delle architetture

Geometric deep learning

Il *geometric deep learning* è una teoria delle architetture di reti neurali che imita l'*Erlangen Programme*, organizzando le architetture in base al concetto di equivarianza rispetto ad azioni di gruppi.

DEFINIZIONE: Funzione equivariante

Sia G un gruppo e siano (S, \cdot) e $(T, *)$ azioni di G . Una funzione $f : S \rightarrow T$ è equivariante rispetto a tali azioni se

$$f(g \cdot s) = g * f(s),$$

per ogni $s \in S$ e per ogni $g \in G$.

ESEMPIO

I convolutional layers delle reti neurali rappresentano mappe invarianti rispetto a traslazioni.

Categorical deep learning

Il *categorical deep learning* è una teoria delle architetture di reti neurali che generalizza il GDL, organizzando le architetture in base al concetto di omomorfismo di (co)algebre categoriche.

DEFINIZIONE: Algebra su un endofuntore

Sia $F : \mathcal{C} \rightarrow \mathcal{C}$ un endofuntore. Un'algebra su F è una coppia (A, a) dove A è un oggetto di \mathcal{C} e $a : F(A) \rightarrow A$ è un morfismo in \mathcal{C} .

DEFINIZIONE: Omomorfismo di algebre

Siano (A, a) e (B, b) algebre sullo stesso endofuntore $F : \mathcal{C} \rightarrow \mathcal{C}$. Un omomorfismo di algebre $(A, a) \rightarrow (B, b)$ è un morfismo $f : A \rightarrow B$ in \mathcal{C} tale che $F(f) \circ b = a \circ f$.

$$\begin{array}{ccc} T(A) & \xrightarrow{T(f)} & T(B) \\ \downarrow a & & \downarrow b \\ A & \xrightarrow{f} & B \end{array}$$

Il CDL generalizza il GDL poiché le azioni di un gruppo G si possono definire come algebre su una monade, e le mappe invarianti si recuperano come omomorfismi tra queste algebre.

DEFINIZIONE: Monade delle azioni di G

Consideriamo l'endofunttore $G \times - : \mathbf{Set} \rightarrow \mathbf{Set}$ che mappa $A \mapsto G \times A$ e $f \mapsto G \times f$. La monade delle azioni di G è definita dotando l'endofunttore delle trasformazioni naturali di $\mu_A : (g, h, a) \mapsto (gh, a)$ e $\eta_A : a \mapsto (e, a)$.

Il CDL costruisce collega algoritmi e strutture dell'informatica classica con le reti neurali.

ESEMPIO: Liste

Sia A un insieme. Consideriamo l'endofuntore $1 + A \times -$ su **Set**.
Sia $\text{List}(A)$ l'insieme delle liste di elementi di A . Allora, se
 $\text{Nil} : 1 \rightarrow \text{List}(A)$ mappa l'unico oggetto di 1 alla lista vuota e
 $\text{Cons} : A \times \text{List}(A) \rightarrow \text{List}(A)$ aggiunge un elemento a una lista,
 $(\text{List}(A), [\text{Nil}, \text{Cons}])$, è un algebra su $1 + A \times -$.

ESEMPIO: List folds

Consideriamo due algebre $(\text{List}(A), [\text{Nil}, \text{Cons}])$ e $(Z, [r_0, r_1])$ su $1 + A \times -$. Un omomorfismo $f : \text{List}(A) \rightarrow Z$ tra queste due algebre deve soddisfare

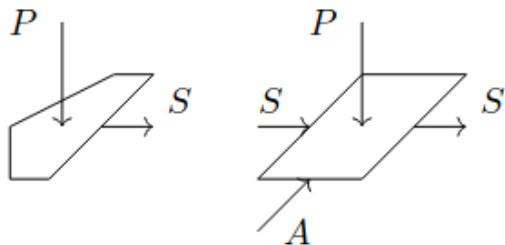
$$\begin{aligned}f(\text{Nil}) &= r_0, \\f(\text{Cons}(a, l)) &= r_1(a, f(l)).\end{aligned}$$

Hence, f è necessariamente un *fold* che riduce liste di elementi di A a singoli elementi di Z .

ESEMPIO: Una cella di un folding RNN

Consideriamo l'endofuntore $1 + A \times X$: e la struttura cartesiana $(1, \times)$ su **Set**. Su questo funtore, può essere costruito un 2-funtore $\mathbf{Para}(1 + A \times X) : \mathbf{Para}_\bullet(\mathbf{Set}) \rightarrow \mathbf{Para}_\bullet(\mathbf{Set})$.

Consideriamo un'algebra $(S, (P, \text{Cell}))$ su tale funtore. Tramite l'isomorfismo $P \times (1 + A \times X) \cong P + P \times A \times X$, deduciamo che $\text{Cell} = [\text{Cell}_0, \text{Cell}_1]$, dove $\text{Cell}_0 : P \rightarrow S$ e $\text{Cell}_1 : P \times A \times S \rightarrow S$. Le funzioni Cell_0 e Cell_1 si possono interpretare come celle di un folding recurrent neural network.

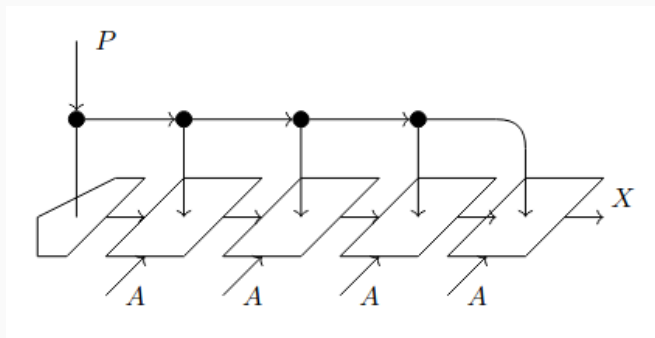


ESEMPIO: Unrolling di un folding RNN

Consideriamo le due algebre $(\text{List}(A), [\text{Nil}, \text{Cons}])$ e $(S, (P, \text{Cell}))$ sull'endofuntore **Para** $(1 + A \times X)$. Ora consideriamo un omomorfismo di algebre

$(P, f, \Delta_P) : (\text{List}(A), [\text{Nil}, \text{Cons}]) \rightarrow (S, (P, \text{Cell}))$. Si può dimostrare che una funzione f così definita è l'unrolling di un folding recurrent neural network. L'algebra $(\text{List}(A), [\text{Nil}, \text{Cons}])$ fornisce gli input della rete neurale, mentre l'algebra $(S, (P, \text{Cell}))$ fornisce le celle.

$$\begin{array}{ccc} T(A) & \xrightarrow{T(f)} & T(B) \\ a \downarrow & \swarrow \kappa & \downarrow b \\ A & \xrightarrow{f} & B \end{array}$$

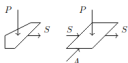


Folding recurrent
neural network

$$1 + A \times S$$

$$\downarrow (P, \text{cell}^{\text{rct}})$$

$$S$$



Unfolding recurrent
neural network

$$S$$

$$\downarrow (P, \langle \text{cell}_o, \text{cell}_n \rangle)$$

$$O \times S$$

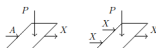


Recursive
neural network

$$A + S^2$$

$$\downarrow (P, \text{cell}^{\text{rcsv}})$$

$$S$$



Full recurrent
neural network

$$S$$

$$\downarrow (P, \text{cell}^{\text{Mealy}})$$

$$(I \rightarrow O \times S)$$

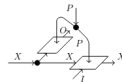


"Moore machine"
neural network

$$S$$

$$\downarrow (P, \text{cell}^{\text{Moore}})$$

$$O \times (I \rightarrow S)$$



Functor learning:

modelli come funtori che sfruttano la struttura dei dati

Categorical representation learning

Il *categorical representation learning* consiste nell'immergere funtorialmente una categoria \mathcal{C} di dati in una categoria \mathcal{R} di vettori latenti.

Grazie alla funtorialità dell'embedding, il *categorical representation learning* consente di preservare la struttura dei dati nello spazio latente.

DEFINIZIONE: Struttura categorica dello spazio latente

tegoria i cui oggetti sono i vettori di \mathbb{R}^n and tale che, per ogni u, v in \mathcal{R} , i morfismi $u \rightarrow v$ sono le matrici $M \in \mathbb{R}^{n \times n}$ tali che $v = Mu$. La composizione è l'ordinario prodotto riga per colonna e l'identità relativa a un generico $v \neq 0$ è $\text{id}_v = \frac{vv^T}{|v|^2}$. L'identità relativa al vettore nullo è la matrice nulla.

Data una categoria \mathcal{C} di dati, il funtore di *embedding* $\mathcal{C} \rightarrow \mathcal{R}$ è realizzato da due **neural embedding layers**: il primo produce rappresentazioni degli elementi di \mathcal{C} come vettori, mentre il secondo produce rappresentazioni dei morfismi di \mathcal{C} come matrici.

DEFINIZIONE: Negative sampling loss

La *objective function* utilizzata per addestrare i due *layers* è la *negative sampling loss*

$$\mathcal{L} = -\mathbb{E}_{(a,b) \sim p(a,b)} (\log P(a \rightarrow b)) + \mathbb{E}_{b' \sim p(b')} \log(1 - P(a \rightarrow b')) ,$$

dove la probabilità che sussista una relazione $a \rightarrow b$ è misurata come

$$P(a \rightarrow b) = \text{sigmoid} \left(F \left(\bigoplus_f v_a^T M_f v_b \right) \right) .$$

ESEMPIO: Traduzione non supervisionata

Se \mathcal{C} e \mathcal{D} sono database di formule chimiche in inglese e cinese, rispettivamente, possiamo usare il CRL per immergere le due categorie in \mathcal{R} funtorialmente. Poi si può imparare un funtore $\mathcal{F} : \mathcal{R} \rightarrow \mathcal{R}$ che preservi la struttura categorica. Tale funtore opererà la traduzione.

ESEMPIO: Traduzione non supervisionata

Se \mathcal{C} e \mathcal{D} sono database di formule chimiche in inglese e cinese, rispettivamente, possiamo usare il categorical representation learning per immergere le due categorie in \mathcal{R} funtorialmente. Poi si può imparare un funtore $\mathcal{F} : \mathcal{R} \rightarrow \mathcal{R}$ che preservi la struttura categorica. Tale funtore effettuerà la traduzione. \mathcal{F} può essere implementato come una matrice $V_{\mathcal{F}}$ che mappa

$$\begin{aligned}v &\mapsto V_{\mathcal{F}}v, \\M_f &\mapsto V_{\mathcal{F}}M_f.\end{aligned}$$

ESEMPIO: Traduzione non supervisionata

La matrice può essere imparata minimizzando la *loss function*

$$\mathcal{L}_{\text{struc}} = \sum_f \|V_{\mathcal{F}} M_f - M_{\mathcal{F}(f)} V_{\mathcal{F}}\|^2,$$

a cui si può aggiungere anche la *loss*

$$\mathcal{L}_{\text{align}} = \sum_{a \in A} \|V_{\mathcal{F}} v_a - v_{\mathcal{F}(a)}\|,$$

per dare parziale supervisione.

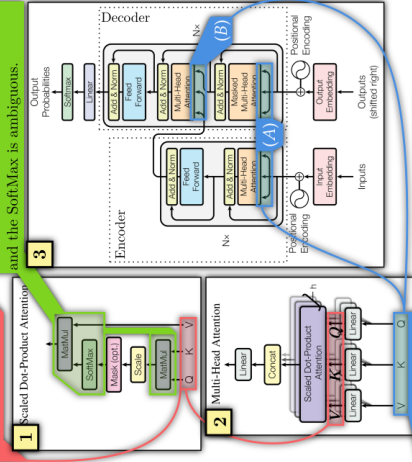
Neural circuit diagrams:

rappresentazioni dettagliate di architetture neurali

Neural circuit diagrams

The Q , K , and V inputs to each attention head come from the learned linear projections of multi-head attention.

The dimension of matrix multiplication and the SoftMax is ambiguous.



These V , K , and Q values are copies in situation (A); while Q is separate in situation (B)

I diagrammi di *deep learning models* generalmente utilizzati nella letteratura scientifica sono inadeguati in quanto non riportano molti dettagli importanti ai fini dell'implementazione.

I *monoidal string diagrams* usati in teoria delle categorie applicata sono inadeguati poiché non riescono a rappresentare funtori e trasformazioni naturali.

I *functor string diagrams* prendono ispirazione dai *monoidal string diagrams* e li adattano per rappresentare funtori e trasformazioni naturali.

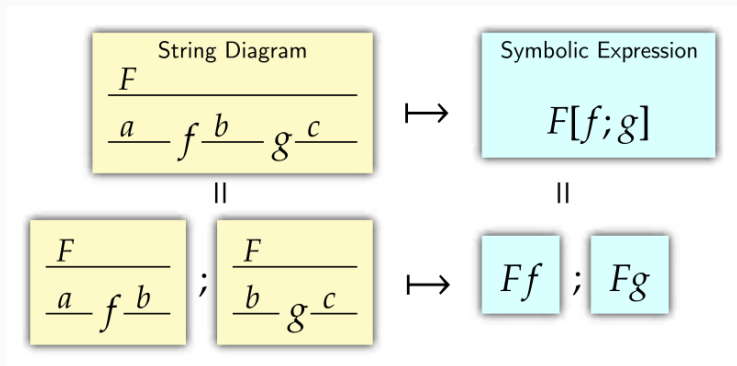
PRINCIPIO di decomposizione verticale

Uno *string diagram* può essere diviso in colonne verticali. Una singola colonna deve contenere solo oggetti o solo morfismi. Colonne con oggetti e colonne con morfismi devono alternarsi.

PRINCIPIO dell'espressione equivalente

Deve essere possibile sostituire ogni diagramma che segue una nuova notazione grafica con uno diagramma equivalente che segue la vecchia notazione.

Neural circuit diagrams

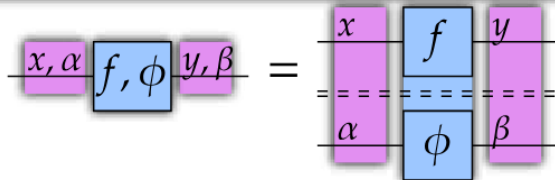


Neural circuit diagrams

$$\begin{array}{c} \text{Functor Diagram} \\ \hline \frac{F}{a} \boxed{\eta_a} \frac{G}{a} \frac{G}{f} \frac{G}{b} \end{array} = \begin{array}{c} \text{Functor Diagram} \\ \hline \frac{F}{a} \eta \frac{G}{a} \frac{G}{f} \frac{G}{b} \end{array} = \begin{array}{c} \text{Functor Diagram} \\ \hline \frac{F}{a} \frac{F}{f} \eta \frac{G}{b} \frac{G}{b} \end{array}$$

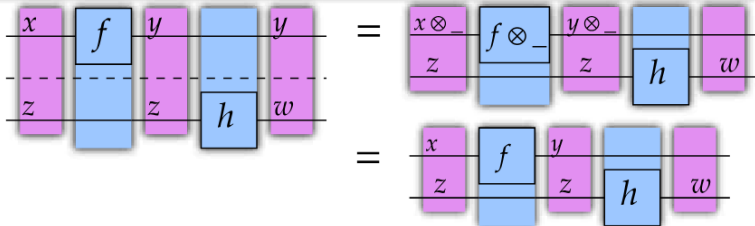
Neural circuit diagrams

We develop double-dashed lines as an equivalent expression to represent product categories $\mathcal{C} \times \mathcal{D}$, where objects and morphisms are tuples.



Neural circuit diagrams

We can reexpress monoidal products using monoidal string diagrams wherein we develop equivalent expressions to view $x \otimes _$ and $f \otimes _$ as functors and natural transformations.



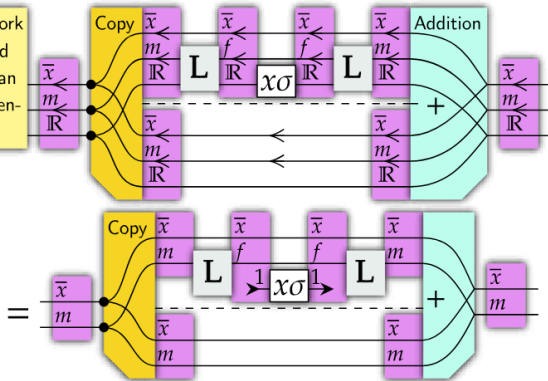
I *neural circuit diagrams* sono *functor string diagrams* specializzati nel rappresentare architetture di reti neurali.

Neural circuit diagrams

A section of a neural network is an expression in **Set**, and can be re-expressed using an equivalent expression for tensor objects.

Equivalent Expression for Tensor Objects

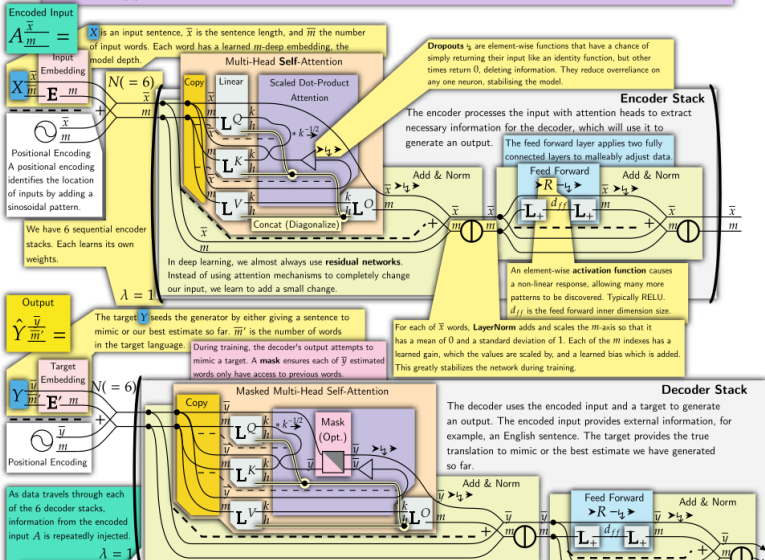
$$\begin{array}{|c|} \hline \bar{x} \\ \hline m \\ \hline \mathbb{R} \end{array} = \begin{array}{|c|} \hline \bar{x} \\ \hline m \\ \hline \mathbb{R} \end{array}$$



Neural circuit diagrams

Neural Circuit Diagram for Transformers

Neural circuit diagrams are a visual and explicit framework for representing deep learning models. Transformer architectures have changed the world, and we provide a novel and comprehensive diagram for the original architecture from *Attention is All You Need*. We describe all necessary components, enabling technically proficient novices who have read our paper to understand the transformer architecture.



Prospettive future

C'è una competizione in corso tra varie discipline, che puntano a spiegare le reti neurali utilizzando ciascuna i propri strumenti:

- fisica matematica,

C'è una competizione in corso tra varie discipline, che puntano a spiegare le reti neurali utilizzando ciascuna i propri strumenti:

- fisica matematica,
- topologia,

C'è una competizione in corso tra varie discipline, che puntano a spiegare le reti neurali utilizzando ciascuna i propri strumenti:

- fisica matematica,
- topologia,
- probabilità,

C'è una competizione in corso tra varie discipline, che puntano a spiegare le reti neurali utilizzando ciascuna i propri strumenti:

- fisica matematica,
- topologia,
- probabilità,
- e così via...

La teoria delle categorie, oltre a offrire strumenti propri, potrebbe creare un ponte tra queste discipline e potrebbe unificare i loro approcci in una teoria generale del deep learning.