

Transformers e la tutela della proprietà dei dati: Un'analisi matematica e applicativa delle Deep Neural Networks

Tesi di Laurea Magistrale in Matematica

Candidato: Bernardo Valente

Relatore : Prof. Vittorio Colao

Correlatori : Ph.D. Angelica Liguori, Ph.D. Ettore Ritacco

23/04/2024



Università della Calabria

Dipartimento di Matematica e Informatica

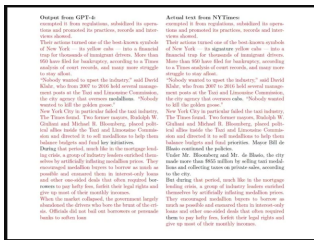
Introduzione

- **Generative AI**
 - La Generative AI è un campo dell'Intelligenza Artificiale che si occupa della creazione di Modelli capaci di generare dati (testo, immagini, audio, video, ...) che siano:
 - **Nuovi**: i dati generati non devono essere una copia di dati esistenti;
 - **Coerenti**: i dati generati devono esibire le caratteristiche del dominio di riferimento
- **Problema**:
 - Gli attuali modelli generativi sono davvero capaci di generare contenuti del tutto **nuovi**?

Introduzione

Perchè è importante risolvere il problema?

- Oggi giorno il problema della violazione dei diritti di copyright è al centro di numerose controversie tra diverse società, l'esempio più recente è la disputa sorta tra il NYT e OpenAI. Pertanto, riconoscere una violazione di dati sensibili diventa una sfida importante da affrontare.



Actual text from NYTimes:
concerned in these regulations, subordinated its operations and promoted its practices, records and interviews shared.

These actions turned one of the best-known symbols of New York – its signature yellow cabs – into a financial trap for thousands of immigrant drivers. More than 900 have filed for bankruptcy, according to a Times analysis of court records, and some were unable to stay afloat.

"Nobody wanted to spend the industry," said David Khale, who from 2007 to 2016 held several management posts at the Taxi and Limousine Commission, the city agency that oversees cabs. "Nobody wanted to kill the golden goose."

New York City is particularly failed the taxi industry. The Times found. Two former mayors, Rudolph W. Giuliani and Michael R. Bloomberg, played pivotal roles after taking the Taxi and Limousine Commission and directed it to self-sufficiency to help them balance budgets and fund priorities. Mayor Bill de Blasio continued the policies.

Under Mr. Bloomberg and Mr. de Blasio, the city made more than \$600 million by selling taxi medallions and collecting taxes on private cabs, according to the city.

But during that period, much like in the mortgage lending crisis, a group of industry leaders enriched themselves by artificially inflating market prices. They encouraged speculation buyers to borrow as much as possible and consumed them in interest-only loans and other one-sided deals that often required them to pay hefty fees, further their legal rights and give up most of their monthly incomes.

(b) Immagine generata da AI



Introduzione

State of Art

Modelli Generativi

Problema: Come avviene la generazione dei dati?

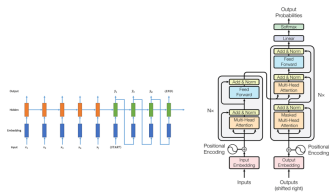
Data una sequenza di elementi $x_1, x_2 \dots x_n \in \mathcal{X}$ un modello generativo p_θ massimizza la *likelihood*:

$$\mathbb{L}(\theta|x) = \sum_i p_\theta(x_i|x_{i-k}, \dots x_{i-1}) \quad (4)$$

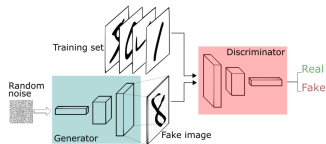
dove k è la dimensione della finestra di contesto. L'obiettivo è dunque risolvere il seguente problema di massimizzazione:

$$\hat{\theta} = \operatorname{argmax}_{\theta \in \Theta} \mathbb{L}(\theta|x) \quad (5)$$

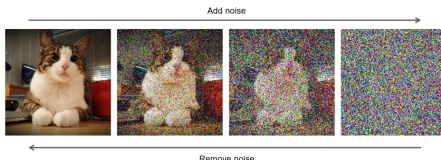
Modelli Generativi



(a) Modello Transformer



(b) GAN



(c) Diffusion Models

Introduzione

Transformers

Descrizione matematica dei transformers

Open Questions

Transformers

Descrizione matematica dei transformers

Transformers come flow maps

Idea: Neural ordinal differential equations

$$\begin{cases} \dot{x}(t) = w(t)\sigma(a(t)x(t) + b(t)), & t \in (0, T), \\ x(0) = x. \end{cases} \quad (6)$$

Che si sono dimostrati un metodo efficace per lo studio di reti neurali residuali con L hidden layers. Sia $x \in \mathbb{R}^d$, una residual network può essere vista come:

$$\begin{cases} x(k) = x(k-1) + f_{\theta}(x(k-1)) & \text{con } k = 1, \dots, L-1 \\ x(0) = x \end{cases} \quad (7)$$

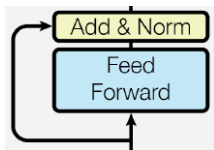
$f_\theta(x(k)) = w(k)\sigma(a(k)x(s)+b(k))$ Lipschitziana e $\theta(\cdot) = (w(\cdot), a(\cdot), b(\cdot)) \in \mathbb{R}^{d \times d} \times \mathbb{R}^{d \times d} \times \mathbb{R}^d$ sono i parametri del modello.

Transformers come flow maps

I Transformers processano *sequenze di input*

- $x_i(0)$ è chiamato *token* o *particella*
- $(x_i(0))_{i=1,\dots,n}$ è chiamato *prompt*

Le implementazioni pratiche del transformer utilizzano un layer di normalizzazione,



che costringono la dinamica delle particelle ad evolversi su un ellissoide allineato con gli assi che variano nel tempo. Assumeremo che $x_i(t) \in \mathbb{S}^{d-1}$ e dunque il *prompt* $(x_i(t))_{i=1\dots n} \in (\mathbb{S}^{d-1})^n$

Transformers come flow maps

Un transformer può essere visto, dunque, come una mappa di flusso (flow map): dato un prompt come input $(x_i(0))_{i=1,\dots,n} \in (\mathbb{S}^{d-1})^n$, questo evolve nel tempo secondo la dinamica

$$\dot{x}_i(t) = \mathbb{P}_{x_i}(t) \left(\frac{1}{Z_{\beta,i}(t)} \sum_{j=1}^n e^{\beta \langle Q(t)x_i(t), K(t)x_j(t) \rangle} V(t)x_j(t) \right) \quad (8)$$

per ogni $i \in \{1, \dots, n\}$ e $t \geq 0$. Con $\mathbb{P}_{x_i}(t)$ ci riferiamo alla mappa di proiezione sullo spazio tangente alla sfera unitaria d – *dimensionale* in x , $\mathbb{P}_x : S^{d-1} \rightarrow T_x S^{d-1}$ così definita:

$$\mathbb{P}_x y = y - \langle x, y \rangle x \quad (9)$$

$Q(\cdot), K(\cdot), V(\cdot)$ sono le matrici parametriche del modello, mentre $\beta > 0$ rappresenta un parametro intrinseco al modello.

Transformers come flow maps

Multi-head-attention:

- Invece di proiettare *query, key, values* in uno spazio d_{model} dimensionale, verranno proiettate h volte tramite matrici di parametri, negli spazi di dimensione d_k, d_k, d_v rispettivamente. Per ognuna di queste versioni proiettate si applica la self-attention. Questo permette di ricavare informazioni da diversi sottospazi di rappresentazione.

Scaled Dot-Product Attention



Multi-Head Attention

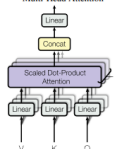


Figura 6: Self attention VS Multi-head attention

Transformers come flow maps

Dinamica transformers con Multi-Head-Attention:

$$\dot{x}_i(t) = P_{x_i(t)} \left(\sum_{h=1}^H \sum_{j=1}^n \frac{e^{\beta \langle Q_h x_i(t), K_h x_j(t) \rangle}}{Z_{\beta, i, h}(t)} V_h x_j(t) \right) \quad (12)$$

con $Z_{\beta, i, h}(t)$ definita con le matrici Q_h, K_h . Qui $H \geq 1$ è un intero. L'analisi del meccanismo della multi-head-attention ci porta a formulare un problema di approssimazione ancora aperto.

Definition 1

Una funzione $f: (\mathbb{S}^{d-1})^n \rightarrow (\mathbb{S}^{d-1})^n$ è detta permutation-equivariant se $f(\pi X) = \pi(f_1(X), f_2(X), \dots, f_n(X))$ per ogni $X \in (\mathbb{R}^d)^n$ e per ogni $\pi \in S_n$.

Transformers come flow maps tra misure

Problema:

- Non essendoci ricorrenze, abbiamo bisogno di un meccanismo che codifica la posizione relativa o assoluta di un token in una sequenza.

Soluzione proposta:

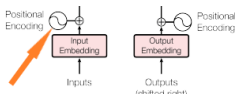


Figura 7

- **Positional Encoding.** Sia $(w_i)_{i=1\dots n} \in (\mathbb{R}^d)^n$, il positional encoding del vettore w_i è definito come il vettore:

$$(p_i)_{2k} = \sin\left(\frac{i}{M^{2k/d}}\right), \quad (p_i)_{2k+1} = \cos\left(\frac{i}{M^{2k/d}}\right),$$

con $k \in \{1, \dots, \frac{d}{2} - 1\}$ e $M > 0$. Infine l' i -esimo token è rappresentato da $x_i(0) = p_i + w_i$.

Transformers come flow maps tra misure

Osservazione: L'output di un transformer è una misura di probabilità. Quindi possiamo rappresentare ogni sequenza di input $\{x_1(0), \dots, x_n(0)\}$ tramite la misura empirica dei tokens costituenti la sequenza di input:

$$\frac{1}{n} \sum_{i=1}^n \delta_{x_i(0)}.$$

Cambio prospettiva: Consideriamo ora i transformers come delle mappe di flusso tra spazi di probabilità sulla sfera $\mathcal{P}(\mathbb{S}^{d-1})$. Dunque possiamo riscrivere la dinamica (10) come segue:

$$\dot{x}_i(t) = \mathcal{X}[\mu(t)](x_i(t)) \quad (13)$$

per ogni $i \in \{1, \dots, n\}$, e per $t \geq 0$. dove:

$$\mu(t, \cdot) = \frac{1}{n} \sum_{i=1}^n \delta_{x_i(t)}(\cdot) \quad (14)$$

è la misura empirica che codifica una sequenza.

Funzionale di interazione di energia

Questo solleva la questione di caratterizzare i minimi e i massimi globali di E_β , che è l'obiettivo del risultato seguente.

Proposizione 4.1

Sia $\beta > 0$ e $d \geq 2$. L'unico minimizzatore globale di E_β su $\mathcal{P}(\mathbb{S}^{d-1})$ è la misura uniforme σ_d . Ogni massimizzatore globale E_β su $\mathcal{P}(\mathbb{S}^{d-1})$ è una massa di Dirac δ_{x^} centrata in qualche punto $x^* \in \mathbb{S}^{d-1}$.*

Transformers

Open Questions

Introduzione

State of Art

Modello realizzato

Transformers

Descrizione matematica dei transformers

Open Questions

Conclusioni

Appendice

Modelli Ricorrenti

Architettura Encoder-Decoder

Transformer

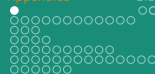
Dinamica dei Transformers

Fenomeno di clustering

Matrici Generali

Bibliografia

GRAZIE PER L'ATTENZIONE



Introduzione

State of Art

Modello realizzato

Transformers

Descrizione matematica dei transformers

Open Questions

Conclusioni

Appendice

Modelli Ricorrenti

Architettura Encoder-Decoder

Transformer

Dinamica dei Transformers

Fenomeno di clustering

Matrici Generali

Bibliografia

Appendice

Modelli Ricorrenti

Modelli ricorrenti

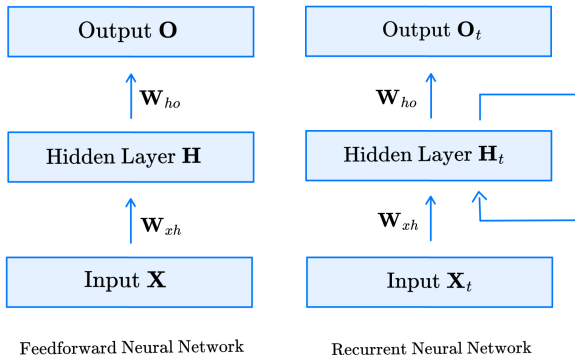


Figura 9: Differenza tra RRN e FFN

Modelli Ricorrenti

- (Forget Gate) Decidere quanta informazione selezionare dalla cella di memoria.

$$f_t = \sigma(\theta_{fx}x_t + \theta_{hf}h_{t-1} + b_f) \quad (25)$$

- (Input Gate) Inserimento della nuova informazione processata nella cella di memoria.

$$i_t = \sigma(\theta_{ix}x_t + \theta_{hi}h_{t-1} + b_i) \quad (26)$$

Modelli Ricorrenti

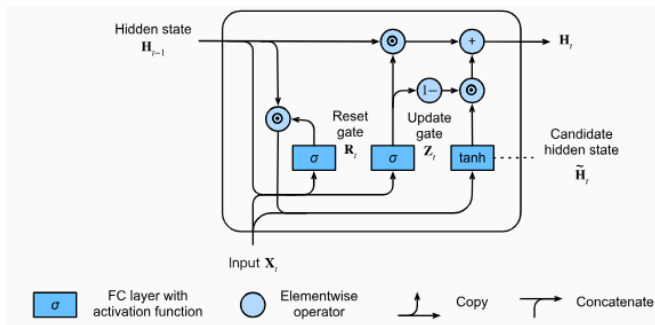


Figura 11: Calcolo dell'hidden state in una GRU

Appendice

Architettura Encoder-Decoder

Architettura Encoder-Decoder

I modelli Seq2Seq processano delle sequenze di dati per produrre nuove sequenze. L'architettura encoder-decoder è un tipo di architettura neurale che viene utilizzata per implementare modelli Sequence2Sequence. L'architettura è formata principalmente da due parti : un encoder che processa una sequenza di input producendo un insieme di vettori detti *context vectors* che vengono poi utilizzati dal decoder per produrre degli output. Il context vector riassume in un certo modo tutte le proprietà delle parole in quella frase e verrà poi passata in input al decoder.

Architettura Encoder-Decoder

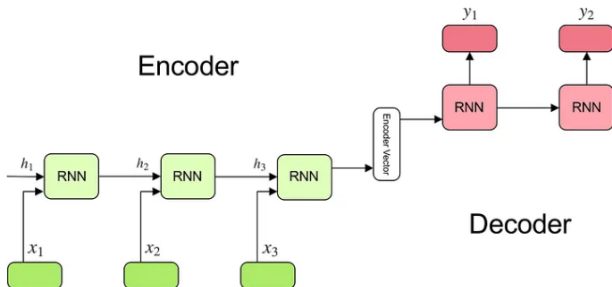


Figura 12: Architettura encoder-decoder

Appendice

Transformer

Multi-head attention

Multi-head-attention:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O \quad (31)$$

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \quad (32)$$

dove $W^O \in \mathbb{R}^{hd_v \times d_{\text{model}}}$. e $W_i^Q \in \mathbb{R}^{d_{\text{model}} \times d_k}$ $W_i^K \in \mathbb{R}^{d_{\text{model}} \times d_k}$ $W_i^V \in \mathbb{R}^{d_{\text{model}} \times d_v}$ e $h = 8$, $d_k = d_v = \frac{d_{\text{model}}}{h}$.

Appendice

Dinamica dei Transformers

Approssimazione del flusso del gradiente di Wasserstein

Introduciamo gli spazi di *Polish*.

Definition 2

Sia (X, τ) uno spazio topologico, diremo che è *completamente metrizzabile* se ammette una metrica $d : X \times X \rightarrow \mathbb{R}$ tale che induce la topologia τ su X e (X, d) è uno spazio metrico completo.

Definition 3

Uno spazio topologico (X, τ) è detto di *Polish* se è completamente metrizzabile ed è anche separabile.

62 / 84

Approssimazione del flusso del gradiente di Wasserstein

Poniamo

$$E_{\beta}(X) = \frac{1}{2\beta} \sum_{i=1}^n \sum_{j=1}^n e^{\beta \langle Vx_i, x_j \rangle}. \quad (39)$$

Si dimostra che (10) può essere equivalentemente scritta come

$$\dot{X} = \nabla E_{\beta}(X(t)) \quad (40)$$

dove il gradiente è calcolato rispetto alla metrica definita in (38) su $(\mathbb{S}^{d-1})^n$.

Appendice

Fenomeno di clustering

Cluster singolo in regime di alta dimensionalità

Osservazione: Il prodotto interno non è altro che una misura della similarità di due particelle. Dunque la formazione dei clusters avviene quando $\langle x_i(t), x_j(t) \rangle \rightarrow 1 \forall i, j \in \{1, \dots, n\}$ e per $t \rightarrow +\infty$.

Teorema 7

Siano $\beta \geq 0$, $d, n \geq 2$ arbitrarie. Consideriamo una sequenza iniziale $(x_i(0))_{i \in \{1, \dots, n\}} \in (\mathbb{S}^{d-1})^n$ di n punti ortogonali a due a due: $\langle x_i(0), x_j(0) \rangle = 0$ per $i \neq j$, e con $(x_i(\cdot))_{i \in \{1, \dots, n\}} \in C^0(\mathbb{R}_{\geq 0}; (\mathbb{S}^{d-1})^n)$ denotiamo l'unica soluzione del corrispondente problema di Cauchy per (10) (rispettivamente per (19)). Allora l'angolo $\angle(x_i(t), x_j(t))$ è lo stesso per tutti $i, j \in \{1, \dots, n\}$ distinti:

$$\angle(x_i(t), x_j(t)) = \theta_\beta(t)$$

per $t \geq 0$ e qualche $\theta_\beta \in C^0(\mathbb{R}_{\geq 0}; \mathbb{T})$.

Appendice

Matrici Generali

77 / 84

Pure Self Attention

Per esempio quando $V = I_d$ tranne in casi eccezionali [2] l'involuppo convesso (*convex hull*) \mathcal{K} delle particelle $z_i(t)$, si restringe e converge a qualche politopo convesso.

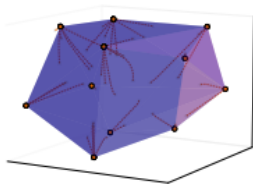


Figura 20: Quando $V = I_3$

Problema 14

Possiamo estendere i risultati di clustering ottenuti nella tabella 19 ad altre matrici Q, K, V ? Inoltre quali sono le forme limite risultante?

Introduzione

State of Art

Modello realizzato

Transformers

Descrizione matematica dei transformers

Open Questions

Conclusioni

Appendice

Modelli Ricorrenti

Architettura Encoder-Decoder

Transformer

Dinamica dei Transformers

Fenomeno di clustering

Matrici Generali

Bibliografia

