

Intro aux Transformers

28/11/25
Simon Gafre

I Intro ~~à~~ Historique

Test part du papier "Attention is all you need" 2017, Google BOK.cit.

Remarque notable Transf: BERT, RoBERTa, Camembert → classif

Ensuite: Generative Pretrained Transformers
↳ puis à BERT qui fait classif.

D'autres: modèles open source

Rq: OpenAI GPT, OpenAI GPT-4, LLaMA 1, DEEPSEEK, QWEN, MISTRAL

Ceci sont des modèles de langage (LLM).
En 2021 on se rend compte que transformers peut remplacer les convolutions

de les modèles d'image "An image is worth 1000 words" ce qui a conduit

aux ViT (Vision Transformers). C'est bcp + fort aujourd'hui les

modèles image gen n'utilisent plus les convs, seulement le transformer.
C'est un peu vrai aussi par la classif.

II Tokens

Dans les transformers les données sont des suites de tokens.

Un transformer est une famille de fonctions $f: \mathbb{R}^{T \times d} \rightarrow \mathbb{R}^{T \times d}$

où $T = \text{nb de tokens}$ et $d = \text{dim du modèle}$. On manipule des suites de longueur T de tokens de dim d .

Δ d est fixe par emb. T peut être amené à changer.

L'idée (qui vient du langage) est que 1 token = 1 unit sémantique.

Ex: langage: syllabe

vision: patch image 8x8

audis: morceaux de spectrogramme

2) Tokenisation

Problème transformer une suite de caractères en une suite de tokens.
Nb de caractères: ~ 2500 Nbd tokens: 3000 (conjugaisons inclus) mais

en vrai juste 25k. Une lettre = hpp pphr, pas d'inf sémantique pas de hien.

Ex le site "Tribkenize" recense les bigrammes utilisés

en pratique, le + populaire est byte-pair encoding (BPE)

Ex: math + emphique. Tém si "en" "ar" "ge" existait dans la base!

Ex: les nombres sont reliés. sans nul tous les chiffres sont décomposés. 10000? 10000?

Ex: Histoir quand deux bair sr linguistique. En fait on sé emb-lyne et

achomah séle process. Par BPE prend la paire de lettre les +

caractères, puis les triples etc. C'est data-driven.

III ~~Arch~~ L'architecture du transformeur

$$f = f_1 \circ \dots \circ f_L \quad \infty \quad L = \text{nb de couches}$$

f_i = bloc de transformeur.

Ex: $L = 120$ pour GPT-4, $L = 70$ pour LLaMA 4

Le point clé est que les f_i ont des mécanismes d'attention. Il y a plusieurs.

1) Self-attention

On a un input $x \in \mathbb{R}^d$ de matrice $Q, K, V \in \mathbb{R}^{d \times d}$

$$\begin{array}{l|l} d^s \text{ par } & q = xQ \rightarrow q \in \mathbb{R}^q \\ & k = xK \rightarrow k \in \mathbb{R}^k \\ & v = xV \rightarrow \text{valeur} \end{array}$$

On définit une matrice de score $S = qk^T \sim [TT^T]$ (gross!)
 $S_{ij} = \langle q_i, k_j \rangle$

Matrice d'attention qui joue le rôle de matrice de transition / Markov

$$A_{ij} = \frac{e^{S_{ij}}}{\sum_{k=1}^n e^{S_{ik}}} \quad (\text{cas où } A = e^S \text{ normaliser la ligne})$$

Output: $o = AV$

$$\text{tel que } (AV)_i = \sum_{j=1}^n A_{ij} V_j$$

Rq: cela se fait de façon canonique:

$$x \mapsto \text{softmax}(xQK^T) xV \quad \text{équivalent à } \text{SATTN}(x)$$

Rq: Le + important est S : qui dit comment les tokens interagissent

entre eux. On pourrait appeler S directement, mais elle doit dépendre de x

On doit se que $S = qk^T = xQK^T$. On pourrait appeler

QK^T direct ~~mais~~ ?

2) Cross-Attention

$$\text{XATTN}(x) = \text{softmax}(xQK^T) yV$$

par $x \in \mathbb{R}^d$ $y \in \mathbb{R}^n$

∞ $x = \text{prompt}$ $y = \text{image}$

3) Full-head Attention

$$\text{MHATTN}(x) = (\text{SATTN}(k))_i$$

ou relier

Rappel: Un transformeur est un f de \mathbb{R}^d à \mathbb{R}^d $T = \text{nb tokens}$
 $d = \text{dim modèle}$

$f = f_1 \circ \dots \circ f_L$ où chaque f_i est un bloc de n architectures dit "bloc transformeur"

chaque bloc f_i est lui-même composé de couches, typiquement

$$f_i = \text{"MLP"} \circ \text{"Attention"}$$

les signaux sont: Feed Forward Network (q, k, v)

la couche de self-attention est paramétrisée par des matrices $Q, K, V \in \mathbb{R}^{d \times d}$
 de se résumer en $x \mapsto \text{softmax}(xQK^T)xV$

Quelle interprétation?

La matrice de score $S_{ij} = \langle q_i, k_j \rangle$ + calcul grad, + la base q_j ,
 donne de l'importance à la key k_j .

Ex: didactique si Key = "not", value = "day"

Exprimez une query "magesis". Il n'y a pas Hd qui dit les dicos.

Fais il y a "mays", "mays". Parce que 2 keys pertinents pour le query.

Raisons on extrair, combiner les valeurs correspondantes.

Historiquement en info, base de données, ce sont des pbs / notes classiques.

Ensuite la matrice de transition est $P_{ij} = \frac{e^{S_{ij}}}{\sum_{k=1}^n e^{S_{ik}}}$

elle que softmax $(X \Omega K^T X) x^T V = P V, v = x V$

Rq: On fait le choix que K, q, v sont linéaires ex. C'est pas obligé
| et des valeurs normales exulter. Fais le lien est suffisant à expliquer.

Il y a beaucoup de variables, et expliquer ça un peu beaucoup Nott.Hend Allah's

Dq: On a des matrices $\Omega^i K^i V^i \sim [d, \frac{d}{4}]$ pour $i=1..H$ et n de mots

On définit des $q^i = x \Omega^i K^i v^i$ etc

output $O^i = \text{softmax}(x \Omega^i K^i v^i)$

la seule manière la concaténation $(O^1, \dots, O^H) \sim [T, d]$

Ex: Group-Org attention: idem mais avec $K^i \equiv K, V^i \equiv V$

seuls les Ω^i sont différents.
Parce que la V est commune pour tous les groupes liés à la structure grammaticale
mais plein de Ω^i différents: qui est le sujet, le verbe etc.

4) Complexité de l'orthogonalité

Complexité de calcul? Évaluer la fonction c'est évaluer la matrice de transition P, donc $O(T^2)$ car il faut appliquer le softmax à la matrice T^2 .

Beaucoup essayer de faire de calculs différents qui coûte moins que T^2

En mémoire c'est un peu pareil T^2 mais + facile à améliorer

Ex: linéar attention Objectif: réduire le nbps.

Rapports que $P_{ij} = \frac{e^{\langle q_i, k_j \rangle}}{\sum_l e^{\langle q_i, k_l \rangle}}$. Pourquoi $\langle q_i, k_j \rangle$ est une similarité, on peut prendre n'importe quelle autre fonction de similarité $\rho(q_i, k_j)$, $\rho: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$.

Idee à la RKHS: $\rho(x, y) = \langle \sigma(x), \sigma(y) \rangle$ en prenant Reu...

Alors $P_{ij} = \frac{\langle \sigma(q_i), \sigma(k_j) \rangle}{\sum_l \langle \sigma(q_i), \sigma(k_l) \rangle} = \frac{\langle \sigma(q_i), \sum_l \sigma(k_l) \rangle}{\langle \sigma(q_i), \sum_l \sigma(k_l) \rangle}$

Donc $P = \frac{1}{\langle \sigma(q_i), \sum_l \sigma(k_l) \rangle} \cdot \sigma(q_i) \sigma(k_l)^T x V$

Où, bien que P est $T \times T$ ici on peut réduire: évaluer en manipulant que des matrices $d \times d$ et $T \times T$.

Est-ce que ça répond ?

Rapports: l'expression de la ressource et la n°

+ du à évaluer (calcul l'inférence)

réduire les 2 pour gagner du temps à l'inférence

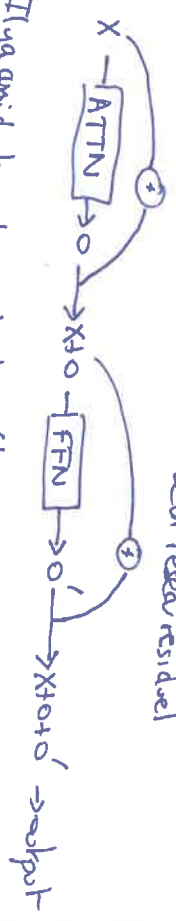
les calculer ailleurs...

5) Le bloc transforme

On rappelle que $f_i = \text{FFN} \circ \text{ATTN}$

Le FFN est un $\text{Feed} \rightarrow$ qui s'applique individuellement à chaque token, c'est $\text{FFN}(x) = (\rho(\beta), \dots, \rho(\beta_T))$ où ρ est un NN

Rq: En pratique ρ est un ReLU et le réseau résiduel



Il y a ainsi des layers de normalisation (layer norm, add norm, rmsnorm) qui sont des $\text{rms}(x) = \frac{1}{\sqrt{N}} \sum x^2$ et γ est un paramètre à apprendre.

C'est un gain de savoir où le mettre. Rq ex: $\frac{x}{\sqrt{N}}$ est une moyenne, $\frac{x}{\sqrt{N}}$ est une variance.

$X \rightarrow \text{rms} \rightarrow \text{ATTN} \rightarrow \text{rms} \rightarrow \text{FFN} \rightarrow \text{output}$

La normalisation permet que le réseau voit toujours la même distribution de données et stabilise le système.

Le résidu dit qu'on peut avoir un message de déplaçant à peu près disons $x \rightarrow x + 0$.

6) Le FFN

a) La tradition: FFN est un MLP: $\rho(x) = A_2 \sigma(A_1 x)$

Multilayer Perceptron. Utilisé à l'époque (GPT-2)

b) Gated Activation (SWIGLU)

$\text{SWIGLU}(x) = A_3 (A_2 x \odot \sigma(A_1 x))$ ← c'est quelque chose $\times x$!

swiglu gated linear unit

$\sigma(x) = \frac{x^2 e^x}{1 + e^x}$ ← "Mish" ou "Swish"

c) Mixture d'Experts (SoDA 85)

On a une famille de MLP "experts" ρ_1, \dots, ρ_n

On a un classifieur "router" $I: x \rightarrow I(x) \in \{1, \dots, n\}$

Puis $\text{MOE}(x) = \rho_{I(x)}(x)$

Rq: Exemple inspiré des 3 dernières années.

Rq ex: la taille a minimalement 8 experts (8E) mais à l'infini ça est un + vite.

Ex: Mistral 8x7B avec 8 experts. Deepseek: 32

Rq: c'est important à avoir: si un expert devient redondant on peut le supprimer. Il y a plein de recettes. Rq ex: on s'entraîne avec un certain ρ KL(ρ_I, unf) \sim petit.

7) Flash attention Ti Dao 222

C'est une implémentation du layer attention avec un gain K/D , sur GPU.

Il observe que le GPU utilise 2 mémoires : HBM et SDRAM
Le défaut du GPU se situe sur HBM (~10Gb/s) qui a un bande passante
pas si élevé de 1-2 TB/s \rightarrow 1 bit émis / lin
La SDRAM est ~100KB mais a un bus de 8 TB/s.

Processus : le GPU lit & écrit dans la HBM, l'écrit dans la SDRAM
la SDRAM est une unité de traitement, qui calcule. Puis se réécrit
dans la HBM (bit)

Après le self attention demand de manipuler $P_n [T, T]$ donc des produits
de $T^2 \rightarrow T^2$ accès à la HBM \rightarrow dr.

L'idée serait de décomposer cette matrice en ~~de~~
ou gagner un facteur $\frac{d^2}{m}$ où $M =$ taille de la SDRAM. Équivalent de $\frac{1}{D}$.