

GDT: White-box transformers

Claim: l'architecture assez spécifique des transformers décèle "naturellement" de certains choix de modélisation sur la manière d'encoder des données / tokens.

→ meilleure interprétabilité (\neq black box)

→ avec la porte à des architectures plus adaptées à un pb.

I Rappels sur l'architecture des transformers

1) Architecture à haut niveau

Un transformer manipule des tokens $x \in \mathbb{R}^d$, obtenus après tokenisation de données (ici d est la dimension du modèle ~ 64)

Plus précisément on traite une suite de tokens, de longueur T

$$X = \begin{bmatrix} x_1 & \dots & x_T \end{bmatrix} \in \mathbb{R}^{d \times T}$$

Rq: Ici les vecteurs sont des colonnes, contrairement à l'exposé précédent.

Un transformer est une fonction $f: \mathbb{R}^{d \times T} \rightarrow \mathbb{R}^{d \times T}$ qui est la composée

de L couches de transformer blocs $f = f_L \circ \dots \circ f_1$ [$L \sim 10^2$]

Chaque bloc est lui-même composé de deux blocs:

$$f_i = \text{FFN} \circ \text{ATT}$$

↙
petit réseau de neurones

↘ bloc d'attention: essentiel

2) Bloc d'attention

- Couche d'attention classique SATT: $\mathbb{R}^{d \times T} \rightarrow \mathbb{R}^{d \times T}$

$$\text{SATT}(X) = (VX) \text{softmax} \left[(KX)^T (QX) \right]$$

où $V, K, Q \in \mathbb{R}^{d \times d}$ sont les paramètres de value, key, query.

- Couche d'attention multi-tête MHATT: $\mathbb{R}^{d \times T} \rightarrow \mathbb{R}^{d \times T}$

$$\text{MHATT}(X) = A \begin{bmatrix} \text{SATT}_1(X) \\ \vdots \\ \text{SATT}_H(X) \end{bmatrix}$$

$$\text{où } \text{SATT}_h(X) = (V_h X) \text{softmax} \left[(K_h X)^T (Q_h X) \right]$$

avec $V_h, K_h, Q_h \in \mathbb{R}^{p \times d}$, $p \ll d$ (typiquement $p = \frac{d}{H}$)

et $A \in \mathbb{R}^{d \times Hp}$ une matrice d'agglomération qui combine les infos de toutes les heads.

- En général la couche d'attention ne se résume pas à MHATT.

Il y a également :

* la skip connection / residual connection, qui consiste à ajouter tid :

$$X \mapsto X + \text{MHATT}(X)$$

* une couche de normalisation

$$X \mapsto \gamma \frac{X}{|X|} \quad (=: \text{RMS}_\gamma(X))$$

En conclusion le bloc d'attention s'écrit :

$$X \mapsto \text{RMS}_\gamma \left(X + \text{MHATT}(X) \right)$$

3) FFN

la couche Fully connected Feedforward Network la plus simple est le Multi Layer Perceptron (MLP) [Désuet]

$$\text{MLP: } \mathbb{R}^{dxT} \longrightarrow \mathbb{R}^{dxT}$$
$$X \longmapsto (\varphi(x_1), \dots, \varphi(x_T))$$

$$\text{ou } \varphi: \mathbb{R}^d \longrightarrow \mathbb{R}^d$$

$$x \longmapsto A_2 \sigma(A_1 x)$$

$$\text{avec } A_1 \in \mathbb{R}^{d' \times d} \text{ et } A_2 \in \mathbb{R}^{d \times d'}$$

II Un paradigme white box (et de belles promesses)

1) Idée générale

On veut trouver l'apprendre une fonction $f: \mathbb{R}^{dxT} \longrightarrow \mathbb{R}^{dxT}$ qui soit capable de transformer des données d'entrée X en Y , de telle façon que même si X suit une distribution non linéaire et multi modale, alors Y est essentiellement linéaire (par morceaux) et compacte.

Rq: Ici $Y = (y_1, \dots, y_T)$. Il est vraisemblable que la loi ~~jointe~~ ^{jointe} sur Y soit compliquée, mais on s'attend à ce que les marginales soient simples &

Ce genre d'hypothèse peut se modéliser par le fait que f doit minimiser un certain critère, c'est qu'elle soit solution

d'un problème $\min_{f \in \mathcal{F}} \mathbb{E}_Y \Phi(Y, \text{param})$
 $Y = f(x)$

les auteurs proposent alors l'idée suivante: une bonne fonction candidate est la transformation $Y \mapsto \text{algo}(Y)$ où "algo" est une étape d'un algo de minimisation par la fonction $Y \mapsto \Phi(Y, \text{param})$

Rq: Pourquoi passe-t-on de la minimisation sur f à une minimisation sur Y ???

L'avantage de cette approche est que l'on peut voir la succession de couches comme le fait que les données passent successivement les itérés de cet algo d'optimisation.

2) Sparse Rate Reduction

a) Taux d'encodage

les auteurs se tournent vers des outils de la théorie de l'information pour obtenir une fonction $R(Y)$ qui mesure à quel point des données sont bien "compressées", représentées sous forme "compacte".

Rq: Je ne suis pas sûr mais je crois que par des distributions continues on a

$$R_\varepsilon(Y) = \min_{\mathbb{E} \|Y - \hat{Y}\|^2 \leq \varepsilon} I(Y; \hat{Y}) \leftarrow \text{information mutuelle}$$

\leftarrow taux d'erreur admis lors de la quantification cf Cover Thomas chap 10

En tout cas c'est vrai pour les Gaussiennes.

~~###~~

b) ~~Gain d'information~~ Mixtures de Gaussiennes

* On va faire l'hypothèse que les tokens y_i suivent une distribut^o qui est une mixture de Gaussiennes de basse dimension

↳ hypothèse de modélisat^o centrale du papier

Plus précisément on va parler de K gaussiennes, centrées ($\mu_k = 0$) avec ~~la~~ covariance $\Sigma_k \in \mathbb{R}^{d \times d}$ de faible rang, rang $\Sigma_k = r$. ↙ p dans le papier.
On note $U_k \in \mathbb{R}^{d \times r}$ la famille orthonormale qui engendre le support.

c) Gain d'information

Par mesure qu'une distribution Y est bien ~~proche~~ proche d'une mixture de Gaussienne, on va vouloir que

- le taux d'encodage de Y w.r.t. ce mixture est bas (= bon)

[" " " " autre chose (ex: $N(0, \Sigma)$) est haut]

On va donc vouloir minimiser $R_{\text{mix}}(Y) \left[-R_{\text{gauss}}(Y) \right]$ ← - taux réduction - information gain

Lem: 1) $R_{\text{gauss}}(Y) \approx \frac{1}{2} \log \det \left[I + \frac{d}{T \Sigma^2} Y^T Y \right]$

2) $R_{\text{mix}}(Y) \approx \sum_{k=1}^K R_{\text{gauss}}(U_k^T Y)$ ← cf a). Pas discuté ds le papier.

Dem 2): $U_k^T Y$ projette Y sur le ~~support~~ support de Σ_k , ce qui donne une borne supérieure naïve de R_{mix}

Minimiser cette fct va promouvoir des Y qui st des mix Gauss.
↳ $R_{\text{mix}}(Y)$

d) Sparse encoding

La fonction $R_{\text{mix}}[\cancel{R_{\text{gauss}}}]$ est invariante par changement de base car c'est une mesure "intrinsèque" de la représentation.

De coup le pb est mal posé et a peut encore spécifier des caractéristiques pr notre distrib^o. (K)

cf + hard

~~Tant qu'à faire on va demander que la représentation Y soit sparse.~~

~~Donc on va aussi rajouter un terme $+\lambda \|Y\|_1$ au pb. Ah et $Y \geq 0$ aussi car ça arrange...~~

Rg: Arguant un peu confus, par moi sparse voudrait dire jouer sur la taille du rang r ..?

e) le problème final

$$\min_{\substack{Y \geq 0 \\ Y=f(x)}} \mathbb{E}_Y \underbrace{\lambda \|Y\|_1 + R_{\text{mix}}(Y) - \cancel{R_{\text{gauss}}(Y)}}_{\Phi(Y)}$$

possiblement remplacé?! oui

ici les paramètres sont
 $U_k \in \mathbb{R}^{d \times r}$
 $\lambda, \epsilon \in \mathbb{R}$

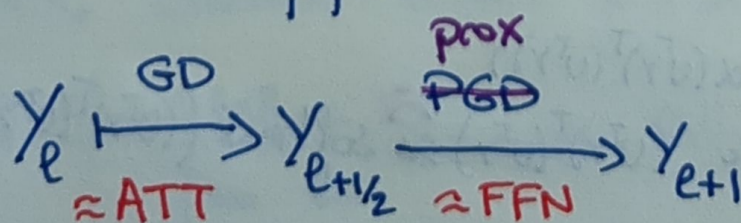
3) Dén algo ou transformer

On va minimiser $\Phi(Y)$ en deux étapes:

- descente de gradient (GD) sur $R_{\text{mix}}(Y)$

- ~~prox-gradient / Forward-Backward~~ (PGD) sur $\lambda \|Y\|_1 + \delta_{Y \geq 0} + R_{\text{gauss}}(Y)$

Ainsi on aura une pipe line:



a) GD sur R_{mix}

On veut pour $Y_{l+\frac{1}{2}} = Y_l - \gamma \nabla R_{mix}(Y_l)$

Pro: $Y - \gamma \nabla R_{mix}(Y) \approx (1 - \hat{\gamma})Y + \hat{\gamma} \widehat{MHATT}(Y)$

où $\widehat{MHATT}(Y) = A \begin{bmatrix} \widehat{SATT}_1(Y) \\ \vdots \\ \widehat{SATT}_K(Y) \end{bmatrix}$

avec $\widehat{SATT}_K(Y) = (U_K^T Y) \text{softmax} \left[(U_K^T Y)^T (U_K^T Y) \right] \in \mathbb{R}^{r \times T}$

et $A = [U_1 \dots U_K] \in \mathbb{R}^{d \times Kr}$

On voit que

- les Q_k, K_k, V_k sont tous égaux à $U_k^T \rightarrow$ vague justifié bio
- la matrice d'agglomération A est aussi paramétrée par les U_k
- nb de gaussiennes $K =$ nb de têtes H
- ~~support~~ rang des gaussiennes $r =$ projet^o des têtes p
- on retrouve à peu près la skip connection
 - $\hookrightarrow \hat{\gamma}$ peut être absorbé par A
 - $\hookrightarrow 1 - \hat{\gamma}$ dans les U_k via drift voir?
- on n'a pas la norm layer. (purement technique?)

Dem: $R_{mix}(Y) = \sum R(U_k^T Y)$ où $\nabla R(Y) = \alpha Y (I + \alpha Y^T Y)^{-1}$, $\alpha = \frac{d}{T \epsilon^2}$

$\nabla R_{mix}(Y) = \alpha \sum U_k U_k^T Y (I + \alpha (U_k^T Y)^T (U_k^T Y))^{-1}$

Or $(I + A)^{-1} \approx I - A$

$\hookrightarrow \approx \alpha \sum U U^T Y (I - \alpha (U^T Y)^T (U^T Y))$

et $(U^T Y) =$ proj^o de Y sur U donc $(U^T Y)^T (U^T Y) \approx \text{softmax}((U^T Y)^T (U^T Y))$

Après c'est en jeu de réécriture.

b) PSD sur $\|\cdot\|_1$ - R_{gauss}

On veut poser $y_{l+1} = \underset{\lambda \|\cdot\|_1 + \delta_{\geq 0}}{\text{prox}}_{\gamma \phi} (y_{l+\frac{1}{2}} + \gamma \nabla R_{\text{gauss}}(y_{l+\frac{1}{2}}))$

Pro: En fait c'est complètement con, donc on va s'en passer à la fin.

(*) 2. d la suite

Rq: Ici j'improvise en m'inspirant du papier afin d'avoir un truc qui tienne la route.

Minimiser R_{mix} est invariant par transformⁿ orthogonale, donc on peut rajouter un prior. Ici les auteurs suggèrent que l'output soit, si possible, sparse dans un dictionnaire.

On se donne donc $D = [d_1 | \dots | d_d]$ qu'on suppose orthogonale (les auteurs disent que avec famille incohérente $D^T D \approx I$ ça passe)

et on voudrait que le $Y = [y_1 | \dots | y_T]$ en sortie soit tel que

$y_T \approx d_{i_T}$ c'ad $y_T \in D$. Ceci équivaut à $y_T = D e_{i_T} \Leftrightarrow D^* y_T = e_{i_T}$

$\Leftrightarrow D^* y_T$ sparse $\forall T \Rightarrow D^* Y$ sparse $\forall T > 0$

Rq: Bon ici c'est pas top car la sparité de Y n'équivaut pas à la sparité des y_T . Il faudrait une notion de sparité par groupes...

Donc on va rajouter un terme de pénalisation $+ \lambda \|D^* Y\|_1^*$ où

$\|\cdot\|_1^* \sim \|\cdot\|$. (Standard). Et aussi $D^* Y \geq 0$ de $+ \lambda \|D^* Y\|_1^*$

où $\|x\|_1^* = \begin{cases} \|x\|, & \text{si } x \geq 0 \\ \infty, & \text{sin} \end{cases}$ c'ad $\|x\|_1 + \delta_{\mathbb{R}_+^n}(x)$.

b) Prox sur $\|D^* Y\|_1^+$

On veut pour $Y_{t+1} = \text{prox}_{\gamma \lambda \|D^* \cdot\|_1^+}(Y_{t+\frac{1}{2}})$

lem: $\text{prox}_{\gamma \lambda \|D^* \cdot\|_1^+}(Y) = D \sigma(D^* Y - \gamma \lambda \mathbb{1})$
↑ ReLU & d'où le $D^* Y \geq 0 \dots$ BoF
↑ matricielle que des 1

On voit que

- On retrouve une architecture MLP avec $A_1 = D$ $A_2 = D^*$
- Il y a un terme de biais (normal? non? pas discuté dans le papier)
- La ReLU est un peu forcée selon moi.

4) Bonus: ce que les auteurs font vraiment

* Dans un premier temps les auteurs proposent de minimiser

$$R_{\text{mix}}(Y) - R_{\text{gauss}}(Y) + \lambda \|Y\|_1^+$$

en disant qu'ils veulent Y sparse (pe pas..) et que la différence

$R_{\text{mix}} - R_{\text{gauss}}$ a du sens (information gain pr discriminer mix vs gauss).

Donc l'étape (b) consiste à minimiser $\lambda \|Y\|_1^+ - R_{\text{gauss}}(Y)$

* Ils font le calcul exact en annexe, avec un pas prox-gradient:

$$Y_{t+1} = \text{prox}_{\gamma \lambda \|Y\|_1^+}(Y_{t+\frac{1}{2}} + \gamma \nabla R_{\text{gauss}}(Y_{t+\frac{1}{2}}))$$

où le prox va donner un ReLU et

$$\nabla R(Y) = \alpha Y (I + \alpha Y^T Y)^{-1}, \quad \alpha = \frac{d}{T \epsilon^2}$$

Là ils font l'hypothèse que $Y_{t+\frac{1}{2}}$ est orthogonale et donc que $\nabla R(Y) = \frac{\alpha}{1+\alpha} Y$

Dans une update $Y_{t+1} = \text{prox} \left(Y + \frac{\gamma d}{1+\alpha} Y \right) = \text{prox} \left(\left(1 + \frac{\gamma d}{1+\alpha}\right) Y \right)$
 $= \sigma \left[\left(1 + \frac{\gamma d}{1+\alpha}\right) Y - \gamma \lambda \mathbb{1} \right]$ (en vrai ils soustraient le γ une ligne sur deux...)

Bon là ils sortent du chapeau que oui R est invariant par transformation orthogonale, et avec un changement de variable illégal fait apparaître un D :

$$Y_{t+1} = \sigma \left[\left(1 + \frac{\gamma \lambda}{1+\alpha}\right) D^* Y_{t+\frac{1}{2}} - \gamma \lambda \mathbb{1} \right]$$

* Dans le corps du papier ils disent que $\nabla R_{\text{gauss}}(Y)$ trop cher à calculer en général et donc laissent tomber l'analyse.

Ils sortent ad hoc le fait qu'à t on a $Y_{t+\frac{1}{2}}$ et qu'il faut le sparsifier, à savoir trouver Y_{t+1} tq $D Y_{t+1} = Y_{t+\frac{1}{2}}$, donc il faut maintenant

résoudre un nouveau pb $\underset{Y \geq 0}{\text{arg min}} \lambda \|Y\|_1 + \frac{1}{2} \|D Y - Y_{t+\frac{1}{2}}\|_F^2$

Ceci donne $Y_{t+1} = \sigma \left[Y + Y D^* (Y - D Y) - \gamma \lambda \mathbb{1} \right]$

- Dans les deux cas ils veulent vraiment sortir les paramètres sparse et pas les représentations ... (et donc perd le A_i)

- Casse totalement l'histoire du papier

- Ils disent vaguement que c'est de l'optim alternée, mais

$$\min_{Y, Z} \lambda \|Y\|_1 + \frac{1}{2} \|D Y - Z\|_F^2 + R_{\text{mix}}(Z)$$

ferait que la 1^{re} étape devrait dériver ici

~~Step~~

III Quelques points positifs

1) XPs

- Les auteurs vérifiait que leur architecture fait bien ce qu'elle prétend
→ pas A mais appris
- compression ($R_{mix} \downarrow$)
- sparsité ($|| \cdot || \downarrow$)
- Ils ne testent pas si les transforme standard le font?
- À nb de paramètres \approx , \hat{m} perf que transforme standard (ViT)
→ possible avantage de partage les param entre ViT & QA
→ ne testent "que" sur $\begin{cases} \text{Im Net} \\ \text{CIFAR finetuning} \end{cases}$? pas de texte?
- Bonnes lois de scaling par les hyperparamètres.

2) Un dernier nugget théorique

Les auteurs ont une autre histoire pour justifier le softmax dans le transformer.

Supposons $N=1$

On veut que on ait $x \rightarrow y$ avec y mixte de Gaussiennes

On va faire comme si x était une moins bonne approximation de y ,
càd $x = z + \sigma w$ où $w \sim \mathcal{N}(0, I)$.

Retrouver z à partir de x peut se faire en calculant

$$z_{\text{opt}} = \mathbb{E}[z | z_e] = \underset{z}{\operatorname{argmin}} \mathbb{E}_{z, w} \| f(z + \sigma w) - z \|_2^2$$

Pro (Tweedie) $z_{e+1} = z_e + \sigma_e^2 \nabla \log q_e(z_e)$ as q_e densité de z^e

Thm: Sans hypothèses,

$$z_{e+1} \approx A \begin{bmatrix} \hat{\text{SATT}}_1(z_e) \\ \vdots \\ \hat{\text{SATT}}_k(z_e) \end{bmatrix}, \quad \hat{\text{SATT}}(z) = \underset{\uparrow}{\text{softmax}} \left(\frac{1}{2\sigma_e^2} (U_k^T z)^T (U_k^T z) \right)$$

pas de V! ni skip connexion

$$A = [U_1 \dots U_k]$$