

Transformers

Transformers map sequences of vectors to sequences of vectors — in contrast to standard NNs, which map a fixed vec to fixed out.

Input: sequence of N "tokens" $x_n^{(0)}$ of dimension D .

Collected into input matrix $X^{(0)} \in \mathbb{R}^{D \times N}$

$$X^{(0)} = (x_1^{(0)} \dots x_N^{(0)})$$

- Image: pixels are vectors in $[255]^3$. You can imagine taking 4×4 patches, each patch has $16 \cdot 3$ elements in $\{0, 1, \dots, 255\}$. Stack these as columns
- Language: text can be broken up into words ("tokens"), each word has a fixed vector which represents it. eg., for "vocab" of V words, could use one-hot encoding for each word/token in the vocab. Could learn embeddings, or keep fixed.
- Supervised learning: if you have (x_i, y_i) pairs, you could let tokens be $\begin{pmatrix} x_i \\ y_i \end{pmatrix} \in \mathbb{R}^{d_H}$.

Transformer blocks:

- Each layer is mapped to next layer through "transformer block":

$$X^{(m)} = \text{xformer-block}(X^{(m-1)}) \in \mathbb{R}^{D \times N}.$$

- Two main components:

- Interactions between tokens: generally requires N^2 time, allows for understanding how similar/dissimilar i^{th} & j^{th} tokens are
- Per-token feature refinement: apply nonlinear NN to token's representation

Stage 1: self-attention

We will map $X^{(m)} \rightarrow Y^{(m)} \in \mathbb{R}^{D \times N}$

$$Y^{(m)} = X^{(m-1)} A^{(m)} \Rightarrow y_n^{(m)} = \sum_{n'=1}^N X_{n'}^{(m-1)} A_{n',n}^{(m)}, \text{ where}$$

columns $y_n^{(m)}$.

- $A_{n',n}^{(m)} \in \mathbb{R}^{N \times N}$

is attention matrix, satisfies

$$\sum_{n'=1}^N A_{n',n}^{(m)} = 1, A_{n',n} \geq 0.$$

- large vals of $A_{n',n}^{(m)} \leftrightarrow n'$ token highly relevant to n token.

- $y_n^{(m)}$ is convex combo of tokens $\{X_i^{(m-1)}\}$, weighted by attn.

Self-attn uses the $X_i^{(m-1)}$ to generate $A^{(m)}$.

Many variants of attn. One simple one:

$$A_{n,n'} := \frac{\exp(\langle x_n, x_{n'} \rangle)}{\sum_{\ell=1}^N \exp(\langle x_n, x_{\ell} \rangle)}$$

→ each dot product of D dim vectors.
→ N^2 entries $\Rightarrow N^2 D$ time.

A more flexible form of attn would be to introduce learnable matrix $U \in \mathbb{R}^{K \times D}$ so that

$$A_{n,n'} = \frac{\exp(\langle Ux_n, Ux_{n'} \rangle)}{\sum_{\ell=1}^N \exp(\langle Ux_{n'}, Ux_{\ell} \rangle)}$$

→ Computing $Ux_n \in \mathbb{R}^K$ requires computing K dot products of D dim vecs;
 KD time; compute for all N vecs, NKD time

→ then computing N^2 dot prods of K dim vectors,
 $N^2 K$ time.

→ total of $NKD + N^2 K$ time, $N^2 D$ if $K < D/2$

Similarly could use two learnable matrices U_q, U_k & compute

$$A_{n,n'} = \frac{\exp(\langle U_k x_n, U_q x_{n'} \rangle)}{\sum_{\ell} \exp(\langle U_k x_{n'}, U_q x_{\ell} \rangle)}$$

$q_n := U_q x_n$, $k_n := U_k x_n$ are called query & key vectors, resp.

Multi-head self attention (MHSA):

- Above allows for one notion of similarity - $\langle U_k x_n, U_q x_{n'} \rangle$.
would be better to allow for multiple similarity metrics.

- Introduce H heads: each "head" has params $U_{q,h}^{(m)}, U_{k,h}^{(m)}, U_{v,h}^{(m)} \in \mathbb{R}^{K \times D}$,
 $U_{p,h}^{(m)} \in \mathbb{R}^{D \times K}$

$$q_{h,n}^{(m)} := U_{q,h}^{(m)} x_n^{(m-1)}, \quad k_{h,n}^{(m)} := U_{k,h}^{(m)} x_n^{(m-1)}$$

$$[A_h^{(m)}]_{n,n'} := \frac{\exp(\langle k_{h,n}^{(m)}, q_{h,n'}^{(m)} \rangle)}{\sum_{\ell} \exp(\langle k_{h,\ell}^{(m)}, q_{h,n'}^{(m)} \rangle)}$$

$$y_h^{(m)} := \underbrace{U_{v,h}^{(m)}}_{K \times D} \underbrace{x^{(m-1)}}_{D \times N} \underbrace{A_h^{(m)}}_{N \times N}; \quad \in \mathbb{R}^{K \times N}$$

$$y^{(m)} = \text{MHSA}_{\Theta}(x^{(m-1)}) = \sum_{h=1}^H \underbrace{U_{p,h}^{(m)}}_{D \times K} \underbrace{y_h^{(m)}}_{K \times N}$$

Θ = params. For each MHSA layer, have H heads, each of which has
key/query/value/projection matrices
 $\uparrow U_{v,h} \quad \hookrightarrow U_{p,h}$

Second stage: after getting $y^{(m)}$, we pass each column (token representation) through an MLP, $\tilde{x}_n^{(m)} = \text{MLP}_\theta(y_n^{(m)})$.

- MLP_θ is usually a 2-3 layer F.C. network
- Can often dominate computation time in transformers

Two more ingredients before can introduce full xformer block.

- (1) Residual connection (2) Token normalization.

(1) we parameterize $x^{(m)} = x^{(m-1)} + \text{res}_\theta(x^{(m-1)}) \rightarrow$
 just learning $x^{(m)} - x^{(m-1)}$.

(2) Layer norm is default technique.

$$[X_n]_l \mapsto \frac{1}{\sqrt{\text{var}(X_n)}} ([X_n]_l - \text{mean}(X_n))$$

$$\text{mean}(X_n) = \frac{1}{D} \sum_i^D [X_n]_l,$$

$$\text{var}(X_n) = \frac{1}{D} \sum_i^D ([X_n]_l - \text{mean}(X_n))^2.$$

Final recipe:

$$x^{(m)} \mapsto \text{LayerNorm}(x^{(m-1)}) = \bar{x}^{(m-1)} \mapsto y^{(m)} = x^{(m-1)} + \text{MHSA}(\bar{x}^{(m-1)}) \mapsto \bar{y}^{(m)} = \text{LayerNorm}(y^{(m)})$$

$$\downarrow$$

$$x^{(m)} = \bar{y}^{(m)} + \text{MLP}(y^{(m)})$$

Position encoding:

- Transformer as we've described doesn't depend on order of tokens
- OK in some contexts, not for language
- The fix is to use "position encoders": if token appears in i th position, we add to token embedding $x_n^{(0)}$ some value $p_n^{(0)}$.
- eg in language, input to transformer is,

$$x_n^{(0)} = \underbrace{e_n^{(0)}}_{\text{embedding corresp. to token}} + \underbrace{p_n^{(0)}}_{\text{embedding corresp. to "n"th token in seq.}}$$

Since transformers can map sequences to sequences, they can operate as learning algorithms at test time:

after being trained, they can take as input (x_i, y_i) sequences, then formulate predictions for query examples.

We'll now discuss one setting where we can understand this behavior precisely:

Linear transformers trained on random linear regression tests.

Model: A single layer, single-head softmax-based attn has form,
 for input matrix $E \in \mathbb{R}^{D \times N}$ (no normalization layer)

$$f(E; W^K, W^P, W^V, W^Q) = E + W^P W^V E \operatorname{softmax} \left(\frac{(W^K E)^T W^Q E}{\rho} \right)$$

ρ : normaliz. factor

We will look at following simplified model: for $\Theta = (W^{PV}, W^{KQ})$,

$$f_{LSA}(E; \Theta) = E + W^{PV} E \cdot \frac{E^T W^{KQ} E}{\rho}$$

Linear Self Attention: no softmax. Also merged $W^K, W^Q \in W^P, W^V$.

Setting: We suppose we have data as follows. Let $\lambda > 0, \Lambda \in \mathbb{R}^{d \times d}$.

$w^{(z)} \stackrel{iid}{\sim} N(0, I_d)$

$x_i^{(z)} \stackrel{iid}{\sim} N(0, \Lambda)$

$y_i^{(z)} \stackrel{iid}{\sim} \langle x_i^{(z)}, w^{(z)} \rangle$.

Receive datasets

$$D^{(z)} = \left\{ (x_i^{(z)}, y_i^{(z)}) \right\}_{i=1}^{N+1}, \quad z=1, \dots, B$$

Define token embedding matrices

$$E^{(z)} = \begin{pmatrix} x_1^{(z)} & \dots & x_N^{(z)} & x_{N+1}^{(z)} \\ y_1^{(z)} & \dots & y_N^{(z)} & 0 \end{pmatrix} \in \mathbb{R}^{d+1 \times N+1}$$

Prediction for $y_{N+1}^{(z)}$: $f_{LSA}(E^{(z)}; \Theta) \in \mathbb{R}^{d+1 \times N+1}$
 $[f_{LSA}(E^{(z)}; \Theta)]_{d+1, N+1} = \hat{y}_{query}^{(z)}$

Define loss function

$$\hat{L}(\Theta) = \frac{1}{2B} \sum_{z=1}^B \left(\hat{y}_{query}^{(z)} - y_{N+1}^{(z)} \right)^2$$

comes from f_{LSA} ; depends on $\Theta = (W^{KQ}, W^{PV})$

Consider limit as $B \rightarrow \infty$: this is "infinite pretraining data" limit. Then,

$$L(\theta) = \lim_{B \rightarrow \infty} \hat{L}(\theta) = \frac{1}{2} \mathbb{E}_{w^c, x_i^c} \left[(y_{\text{query}}^c - y_{\text{NH}}^{(c)})^2 \right].$$

Let's consider gradient flow on this: $\frac{d\theta}{dt} = -\nabla L(\theta)$.

Theorem [Zhang-Frei-Bartlett JMLR '24]:

For a suitable random init, GF converges to a global min of $L(\theta)$. Moreover,

if $\Gamma_N := (1 + \frac{1}{N})\Lambda + \frac{1}{N} \text{tr}(\Lambda) \mathbb{I}_d$, then the transformer converges to a network $\hat{\theta}$ which satisfies, for any $(x_i, y_i)_{i=1}^M$,

$$\begin{aligned} \left[\text{first} \left(\begin{pmatrix} x_1 & \dots & x_M & x \\ y_1 & \dots & y_M & 0 \end{pmatrix}, \hat{\theta} \right) \right]_{d_H, M+1} &= x^T \Gamma_N^{-1} \left(\frac{1}{M} \sum_{i=1}^M y_i x_i \right) \\ &= x^T \left[\left(1 + \frac{1}{N}\right) \Lambda + \frac{1}{N} \text{tr}(\Lambda) \mathbb{I}_d \right]^{-1} \cdot \left(\frac{1}{M} \sum_{i=1}^M y_i x_i \right). \end{aligned}$$

$$=: \hat{y}_{\text{query}}(x)$$

Remarks. If $y_i = \langle w, x_i \rangle \forall i$, then: $\frac{1}{M} \sum_{i=1}^M y_i x_i = \frac{1}{M} \sum_{i=1}^M x_i x_i^T w = \left(\frac{1}{M} \sum_{i=1}^M x_i x_i^T \right) w$.

• If $N \rightarrow \infty$, $\Gamma_N \rightarrow \Lambda$, so if \uparrow holds too, we get

$$\hat{y}_{\text{query}}(x) \rightarrow x^T \Lambda^{-1} \cdot \left(\frac{1}{M} \sum_{i=1}^M x_i x_i^T \right) w. \quad \Gamma_N \text{ is a regularized version of } \Lambda.$$

• If x_i are iid, $\frac{1}{M} \sum_{i=1}^M x_i x_i^T \rightarrow \mathbb{E} x x^T$. If $\mathbb{E} x x^T = \Lambda$, then $\hat{y}_{\text{query}}(x) \rightarrow x^T w$.

- If x_i are not iid, or if $\mathbb{E}xx^T \neq \Lambda$ (distribution shift), then we should not expect the trained transformer to do well for new linear regression tasks.

What happens if $y_i \neq \langle w, x_i \rangle$?

Theorem (Zhang-Forei-Bartlett) Suppose $(x_i, y_i) \stackrel{iid}{\sim} \mathcal{D}$, where $x_i \sim N(0, \Lambda)$. Suppose $\mathbb{E}y$, $\mathbb{E}xy$, $\mathbb{E}y^2xx^T$ exist & are finite. Then, for \hat{y} : prediction for trained transformer, if we let

$$a := \Lambda^{-1} \mathbb{E}xy, \quad \Sigma := \mathbb{E}[(xy - \mathbb{E}xy)(xy - \mathbb{E}xy)^T], \quad \text{then:}$$

$$\mathbb{E}[(\hat{y} - y)^2] = \min_{w \in \mathbb{R}^d} \mathbb{E}[(\langle w, x \rangle - y)^2] + \frac{1}{M} \text{tr}(\Sigma \Lambda^{-2} \Lambda) + \frac{1}{N^2} \left(\|a\|_{\Lambda^{-2} \Lambda}^2 + 2 \text{tr}(\Lambda) \|a\|_{\Lambda^{-2} \Lambda}^2 + \text{tr}(\Lambda)^2 \|a\|_{\Lambda^{-2} \Lambda}^2 \right)$$

Remarks

- A type of "emergent behavior": trained on noisy linear regression, but achieves best linear prediction error. So does well on noisy regression, nonlinear, ...
- Training contexts have different effect than test-time: $\frac{1}{N^2}$ vs $\frac{1}{M}$.

Proof ingredients : Fairly involved ...

① Show $L(\theta) \Leftrightarrow L(u) = \mathbb{E} \left[(u^T H u - y)^2 \right]$ for some non-psd H ,
so L is non-conv.

② Establish ρ -inequality:

$$\|\nabla L(u)\|^2 \geq c \cdot (L(u) - \min_u L(u))^2 \quad \text{for some } c > 0.$$

This is 90% of proof. Has some analogues to prod ideas in training
two-layer linear networks (recal: $f(E; \theta) = E + W^{IV} E \cdot \frac{E^T W^{KR} E}{\rho}$. Two linear layers,
 w^{KR}, w^{IV} appear.)

③ Identify limit of gradient flow

④ Characterize statistical properties of the limit.