

Machine Learning (part II)

Transformers

Angelo Ciaramella

Introduction

Transformers

- transform a set of vectors in some representation space into a corresponding set of vectors, having the same dimensionality, in some new space
- the fundamental concept that underpins a transformer is attention

Attention

- developed as an enhancement to RNNs for machine translation
- A transformer can be viewed as a richer form of embedding in which a given vector is mapped to a location that depends on the other vectors in the sequence



Attention



Figure 12.1 Schematic illustration of attention in which the interpretation of the word 'bank' is influenced by the words 'river' and 'swam', with the thickness of each line being indicative of the strength of its influence.



Introduction

Figure 12.3 The structure of the data matrix **X**, of dimension $N \times D$, in which row *n* represents the transposed data vector $\mathbf{x}_n^{\mathrm{T}}$.





$$\widetilde{\mathbf{X}} = \operatorname{TransformerLayer} [\mathbf{X}]$$

function that takes a data matrix as input and creates a transformed matrix of the same dimensionality as the output



Deep networks

- Multiple transformer layers
 - to construct deep networks capable of learning rich internal representations
 - each transformer layer contains its own weights and biases, which can be learned using gradient descent using appropriate cost function



Attention coefficients

attention weights (non-negative)

$$\mathbf{y}_n = \sum_{m=1}^N a_{nm} \mathbf{x}_m$$

$$a_{nm} \ge 0$$
$$\sum_{m=1}^{N} a_{nm} = 1.$$

The coefficients should be close to zero for input tokens that have little influence on the output y_n and largest for inputs that have most influence



Self-Attention

- choosing which movie to watch in an online movie streaming service
 - associate each movie with a list of attributes
 - attributes of each movie in a vector called the key
 - the corresponding movie file itself is called a value
 - personal vector of values for the desired attributes called query
- movie service could then compare the query vector with all the key vectors to find the best match and send the corresponding movie to the user in the form of the value file (hard attention)



Soft attention

- we use continuous variables to measure the degree of match between queries and keys
 - variables to weight the influence of the value vectors on the outputs



if all the input vectors are orthogonal

 $\mathbf{y}_m = \mathbf{x}_m$ for $m = 1, \dots, N$



Self attention



we are using the same sequence to determine the queries, keys, and values



Self attention



 $N \times D$

dot product self-attention

$$\begin{split} \mathbf{Q} &= \mathbf{X}\mathbf{W}^{(q)}\\ \mathbf{K} &= \mathbf{X}\mathbf{W}^{(k)}\\ \mathbf{V} &= \mathbf{X}\mathbf{W}^{(v)} \end{split}$$



Self attention



output from an attention layer



Scaled self-attention

- the gradients of the softmax function become exponentially small for inputs of high magnitude
- variance of the dot product would be D_K



Scaled dot-product self-attention

$$\mathbf{Y} = \operatorname{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) \equiv \operatorname{Softmax}\left[\frac{\mathbf{Q}\mathbf{K}^{\mathrm{T}}}{\sqrt{D_{k}}}\right]\mathbf{V}$$



in natural language some patterns might be relevant to tense whereas others might be associated with vocabulary

 $\mathbf{H}_h = \operatorname{Attention}(\mathbf{Q}_h, \mathbf{K}_h, \mathbf{V}_h)$

H heads

 $egin{aligned} \mathbf{Q}_h &= \mathbf{X}\mathbf{W}_h^{(\mathrm{q})}\ \mathbf{K}_h &= \mathbf{X}\mathbf{W}_h^{(\mathrm{k})}\ \mathbf{V}_h &= \mathbf{X}\mathbf{W}_h^{(\mathrm{v})} \end{aligned}$

 $\mathbf{H}_{1} \mathbf{H}_{2} \cdots \mathbf{H}_{H} \times \mathbf{W}^{(o)} = \mathbf{Y}$ $N \times HD_{v} \qquad HD_{v} \times D$

 $\mathbf{Y}(\mathbf{X}) = \operatorname{Concat} \left[\mathbf{H}_1, \dots, \mathbf{H}_H\right] \mathbf{W}^{(\mathrm{o})}$



Multi-head attention



Figure 12.8 Information flow in a multi-head attention layer. The associated computation, given by Algorithm 12.2, is illustrated in Figure 12.7.



Residual connections

- Improving the learning
 - $\mathbf{Z} = \text{LayerNorm}\left[\mathbf{Y}(\mathbf{X}) + \mathbf{X}\right]$
- Adding MLP
 - $\widetilde{\mathbf{X}} = \operatorname{LayerNorm}\left[\operatorname{MLP}\left[\mathbf{Z}\right] + \mathbf{Z}\right]$



H heads



Positional embeddings

$$PE(i, \delta) = \begin{cases} \sin(\frac{i}{10000^{2\delta'/d}}) & \text{if } \delta = 2\delta' \\ \cos(\frac{i}{10000^{2\delta'/d}}) & \text{if } \delta = 2\delta' + 1 \end{cases}$$



ML – Transformers

Attention is all You Need



Figure 1: Architecture of the standard Transformer (Vaswani et al., 2017)



Vision Transformers





An Image is Worth 16x16 Words

ViT

Vision Transformers

Transformers / Davide Concornini / 2023

33



Image captioning



A woman is throwing a <u>frisbee</u> in a park.



A $\underline{\text{dog}}$ is standing on a hardwood floor.



A <u>stop</u> sign is on a road with a mountain in the background.



A little <u>girl</u> sitting on a bed with a teddy bear.



A group of <u>people</u> sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.



ML – Transformers

ViT for video



ML – Transformers



ViT for video





ML – Transformers

Sample frames, extract 2D patches and linearly project (as in ViT). Consider a video as a "big image"



Visual tracking



Tracking vision transformer with class and regression tokens



Visual tracking



Outputs for sequence *BlurCar1* in OTB100 with the attention of two tokens overlaying the image. *First row:* <u>Bounding Box</u> outputs. Second row: Regression token attention. *Third row:* Classification token attention

Visual tracking



Outputs for sequence Group1_2 in UAV123 with the attention of two tokens overlaying the image. *First row:* <u>Bounding Box</u> outputs. Second row: Regression token attention. *Third row:* Classification token attention.