

Deep Learning HW3

Yuanhe Guo

February 2024

1 Theory

1.1 Attention

- (a) $\mathbf{H} = \mathbf{V} \text{softargmax}_\beta(\mathbf{K}^T \mathbf{Q}) \in \mathbb{R}^{t \times n}$
- (b) The higher β , the more focused the attention will be, vice versa.
- (c) When β goes to $+\infty$ (temperature gets cold), the output of softargmax_β goes to onehot encoding, thus preserving one value vector v . Here we are referring to hard attention. In fully connected layer, we can use a linear transformation with no bias and weight being identity matrix.
- (d) When β goes to 0 (temperature goes hot), the output of softargmax_β goes to the average value. Here we are referring to soft attention. In fully connected layer, we can use a linear transformation with no bias and weight being a scalar matrix.
- (e) First let's write out the $\hat{\mathbf{H}}$ after perturbation.
 $\hat{\mathbf{H}} = \mathbf{V} \text{softargmax}_\beta([\mathbf{k}_1^T \mathbf{Q}, \dots, (\mathbf{k}_i + \epsilon)^T \mathbf{Q}, \dots, \mathbf{k}_m^T \mathbf{Q}])$
Given $\mathbb{E}(\epsilon) = 0$, we know $\mathbb{E}[(\mathbf{k}_i + \epsilon)^T \mathbf{Q}] = \mathbf{k}_i^T \mathbf{Q}$. Therefore, a small perturbation to one of the \mathbf{k} won't change \mathbf{H} .
- (f) Suppose we change the i -th \mathbf{k} . Now the new $\hat{\mathbf{H}}$ is
 $\hat{\mathbf{H}} = \mathbf{V} \text{softargmax}_\beta([\mathbf{k}_1^T \mathbf{Q}, \dots, \mathbf{k}_i^T \mathbf{Q}, \dots, \mathbf{k}_m^T \mathbf{Q}])$
Let's consider the softargmax_β function:

$$\text{softargmax}_\beta(\mathbf{K}^T \mathbf{Q}) \begin{cases} \frac{\exp(\beta \mathbf{k}_i^T \mathbf{Q})}{\sum_{j=1}^{i-1} \exp(\beta \mathbf{k}_j^T \mathbf{Q}) + \sum_{j=i+1}^m \exp(\beta \mathbf{k}_j^T \mathbf{Q}) + \exp(\beta \mathbf{a} \mathbf{k}_i^T \mathbf{Q})} & \text{if } t \neq i \\ \frac{\exp(\beta \mathbf{a} \mathbf{k}_i^T \mathbf{Q})}{\sum_{j=1}^{i-1} \exp(\beta \mathbf{k}_j^T \mathbf{Q}) + \sum_{j=i+1}^m \exp(\beta \mathbf{k}_j^T \mathbf{Q}) + \exp(\beta \mathbf{a} \mathbf{k}_i^T \mathbf{Q})} & \text{if } t = i \end{cases}$$

As a result, the \mathbf{h} value for the i -th \mathbf{k} becomes larger, while all other values become smaller.

1.2 Multi-headed Attention

- (a) Let the concatted matrix of attentions to be S .

$$S = \begin{pmatrix} V W_V^1 \text{softargmax}_\beta((K W_K^1)^T Q W_Q^1) \\ \dots \\ V W_V^h \text{softargmax}_\beta((K W_K^h)^T Q W_Q^h) \end{pmatrix}$$

Then we apply an linear transformation to get H

$$H = S W^0$$

- (b) Both conv1d and multi-headed attention applies weight to input in order to project it into a new matrix, then connects the result to a linear layer.

1.3 Self Attention

- (a) Consider the i -th head. $Q^i = C W_Q^i, K^i = C W_K^i, V^i = C V_K^i$.
 $S^i = \text{Attention}(Q^i, K^i, V^i) = V^i \text{softargmax}_\beta((K^i)^T Q^i)$
 $H = S W^0$
- (b) Since self-attention handles data simultaneously, we need positional encoding when the input is sequential. For example, the position of each word in a text input matters, so we want to add positional encoding.
- (c) If $Q = K$, and Q is an orthogonal matrix, then $K^T Q = Q^{-1} Q = I$. Now the attention behaves like an identity layer.
- (d) When the result of $\text{softargmax}_\beta(K^T Q)$ is a diagonal matrix, then the self attention just does linear projection to each location. The proper name for “running” linear layer should be “1Dconvolution”, which applies linear projection one element at a time.
- (e) Let’s consider a special situation. The attention score calculated by $\text{softargmax}_\beta(K^T Q)$ for each location (as well as its relative neighbors) are the same. In this way, each self attention head acts like a kernel in a convolution layer, which applies a shared weight to one location and its neighbors.

1.4 Transformer

- (a) The way of handling input sequence is different between Transformer and previous models. Previous models knows the neighbors of current input by using the latent of previous inputs, resulting in sequential processing. Transformer takes in sequence in parallel, and gets to know the neighbors of one element by using attention mechanism, while the position information is preserved by positional encoding. Compared with previous models, Transformer is able to handle a large input sequence efficiently.

- (b) Self-attention is an attention mechanism that relates different positions of an input sequence to output a representation. It has three significant usages. First, self-attention reduces computational complexity per layer when compared with recurrent layers. Second, self-attention can parallelized a large amount of computation. Third, self-attention enables the network to better learn long-range dependencies.
- (c) Multi-head attention learns to linearly projects queries, keys values h times, and apply attention to each version of queries, keys and values to get h heads. After that, all heads are concat together. The benefit is that multi-head attention allows the model to learn from different representations in parallel.
- (d) Feed-forward networks used in transformer are position-wise. One feed-forward network consists of two linear transformations and a ReLU in the middle. This enables the model to utilize different parameters across layers at the same position.
- (e) Layer normalization normalize the layer with means and variance of all hidden units in that layer. This technique reduces the undesired correlated changes caused by the output of one layer. It also overcomes the constraints on mini-batch size of batch normalization. Transformer architecture applies layer normalization after residual connection in both “encoder” and “decoder”.

1.5 Vision Transformer

- (a) Vision Transformer turns image into sequence by embedding image patches, while CNN turns image into latent by weight sharing. In ViT, the process of patching images is a convolution layer whose kernel size and stride is the same.
- (b) There are three main differences. First, ViT takes images as input, and Transformer takes words as input. Second, ViT has “encoder” but doesn’t have “decoder”, so there’s only one part of attention in ViT model. Third, the output of ViT is the predicted class, which is in different modality of the input.
- (c) In ViT, positional embeddings are used for preserving locality information in the image. Positional embeddings in ViT are learned 1D embeddings, while in Transformer, the embedding is a combination of sine and cosine functions.
- (d) ViT connects a one hidden layer MLP to the output of multi-head attention as classification head during pre-training. During fine-tuning, the classification head is just a one layer MLP.