

MSc. Data Science & AI

Class: Advanced Deep Learning

Professors: Frederic Precioso, Rémy Sun, Pierre-Alexandre Mattei

Student: Aglind Reka

Mail: aglind.reka@etu.univ-cotedazur.fr

Contents

1	Introduction	1
1.1	Context	1
2	Transformers	1
2.1	Scaled Dot Product	1
2.2	Multi-Head Attention	3
2.3	Encoder Block	7
2.4	Transformer Predictor	9
3	Questions	13
3.1	How do we solve the sequence reversal problem at the end of the notebook? Discuss both the method and results.	13
3.2	What would be different if we wanted to predict something about the sequence as a whole?	14

1. Introduction

1.1 Context

Transformers, introduced in the 2017 paper "Attention is All You Need" [1], have become a dominant architecture in natural language processing and machine learning. They leverage attention mechanisms for parallelized and scalable training, distinguishing them from traditional sequence models. The self-attention mechanism allows elements in a sequence to prioritize information, enabling the model to capture long-range dependencies efficiently. Transformers commonly employ an encoder-decoder architecture, incorporating positional encoding to handle the sequence order. Their impact extends to various tasks, making them a foundational technology in contemporary deep learning applications.

2. Transformers

2.1 Scaled Dot Product

2.1.1 Code

Code extraction 1: Scaled Dot-Product Attention

```
1 import torch
2 import torch.nn.functional as F
3 import math
4
5 def scaled_dot_product(q, k, v, mask=None):
6     d_k = q.size()[-1]
7
8     #####
9     ### YOUR CODE HERE! ###
10    #####
11
12    # Compute attn_logits
13    attn_logits = torch.matmul(q, k.transpose(-2, -1))
14    attn_logits /= math.sqrt(d_k)
15
16    # Apply mask if not None
17    if mask is not None:
18        attn_logits = attn_logits.masked_fill(mask == 0, -
19        1e14)
20
21    # Pass through softmax
22    attention = F.softmax(attn_logits, dim=-1)
23
24    # Weight values accordingly
25    output_values = torch.matmul(attention, v)
```

```

26 #####
27 ###      END      ###
28 #####
29
30 return output_values, attention

```

The purpose of the function is to implement the scaled dot-product attention mechanism, which is commonly used in the context of attention mechanisms in neural network architectures, especially in transformer models.

The input parameters that the function takes:

- **q**: Query vectors.
- **k**: Key vectors.
- **v**: Value vectors.
- **mask**: An optional mask to control which positions should be ignored in the attention computation, by default `None`.

2.1.2 The steps that the function follows (Explanation)

Step 1: Compute Attention Logits

Use matrix multiplication to compute the raw attention scores between the query (**q**) and key (**k**) vectors. This step measures the compatibility between each query and key.

Step 2: Scale Logits

Divide the attention logits by the square root of the dimension of the key vectors (**d_k**). This scaling helps stabilize the training process and prevents the gradients from becoming too small during backpropagation.

Step 3: Apply Mask (Optional)

If a mask is provided, use it to mask certain positions in the attention logits. This step is useful, for example, in sequence-to-sequence tasks where certain positions in the input sequence should be ignored during attention computation.

Step 4: Softmax

Apply the softmax function to obtain normalized attention weights. Softmax converts the attention logits into a probability distribution, ensuring that the weights sum to 1.

Step 5: Weight Values

Use the computed attention weights to weight the values (**v**). This step combines information from the values based on the attention weights. Finally, the function returns the weighted values (**output_values**) and the attention weights (**attention**).

The function is a crucial part of the attention mechanism in transformers, allowing the model to focus on different parts of the input sequence during processing. The attention weights indicate how much attention each element in the input sequence should receive. This mechanism is used to capture long-range dependencies and improve the model's ability to handle sequential data.

2.2 Multi-Head Attention

The scaled dot product attention enables a network to focus on different aspects of a sequence. To address the need for attending to multiple aspects of a sequence element, attention mechanisms are extended to multiple heads. This involves transforming query, key, and value matrices into sub-queries, sub-keys, and sub-values. These are independently processed through scaled dot product attention, and the resulting heads are concatenated and combined using a final weight matrix. This process enhances the network's ability to capture diverse patterns in the input sequence.

2.2.1 Code

Code extraction 2: Multihead Attention Class

```

1 import torch
2 import torch.nn as nn
3
4 class MultiheadAttention(nn.Module):
5     def __init__(self, input_dim, embed_dim, num_heads):
6         super().__init__()
7         assert embed_dim % num_heads == 0, "Embedding
            dimension must be 0 modulo number of heads."
8
9         self.embed_dim = embed_dim # dimension of
            concatenated heads
10        self.num_heads = num_heads
11        self.head_dim = embed_dim // num_heads
12
13        #####
14        ### YOUR CODE HERE! ###
15        #####
16
17        # Create linear layers for both qkv and output
18        # TIP: Stack all weight matrices 1...h together for
            efficiency
19        self.o_proj = nn.Linear(embed_dim, embed_dim)
20        self.qkv_proj = nn.Linear(input_dim, embed_dim * 3)
21
22        #####
23        ###          END          ###
24        #####
25
26        self._reset_parameters()
27

```

```

28     def _reset_parameters(self):
29         # Original Transformer initialization, see PyTorch
           documentation
30         nn.init.xavier_uniform_(self.qkv_proj.weight)
31         self.qkv_proj.bias.data.fill_(0)
32         nn.init.xavier_uniform_(self.o_proj.weight)
33         self.o_proj.bias.data.fill_(0)
34
35     def forward(self, x, mask=None, return_attention=False):
36
37         #####
38         ### YOUR CODE HERE! ###
39         #####
40
41         batch_dim, seq_length, input_dim = x.shape
42
43         # Compute linear projection for qkv and separate heads
44         # QKV: [Batch, Head, SeqLen, Dims]
45         qkv = self.qkv_proj(x)
46         qkv = qkv.reshape(batch_dim, seq_length,
           self.num_heads, 3 * self.head_dim)
47         qkv = qkv.permute(0, 2, 1, 3)
48         q, k, v = torch.chunk(qkv, 3, dim=-1)
49
50         # Apply Dot Product Attention to qkv ()
51         attention_values, attention = scaled_dot_product(q, k,
           v)
52
53         # Concatenate heads to [Batch, SeqLen, Embed Dim]
54         attention_values = attention_values.permute(0, 2, 1, 3)
55         attention_values = attention_values.reshape(batch_dim,
           seq_length, self.embed_dim)
56
57         # Output projection
58         o = self.o_proj(attention_values)
59
60         if return_attention:
61             return o, attention
62         else:
63             return o
64
65         #####
66         ###          END          ###
67         #####

```

2.2.2 Class Explanation

The class `MultiheadAttention` is a neural network module representing a multihead attention mechanism.

Initialization

The constructor initializes the module's parameters, including the embedding dimension, the number of attention heads, and the input dimension. It ensures that the embedding dimension is divisible evenly by the number of heads.

Linear Layers

Two linear layers are created – `o_proj` for output projection and `qkv_proj` for the linear projection of queries, keys, and values.

Parameter Initialization

The `_reset_parameters` method initializes the linear layer weights and biases using Xavier initialization. Xavier uniform initialization for the weight parameters of two linear layers (`self.qkv_proj` and `self.o_proj`). Additionally, it sets the bias parameters of these layers to zero. This choice aims to establish a neutral starting point for learning and simplifies the initial learning dynamics.

Forward Method

The `forward` method executes the forward pass of the attention mechanism.

Linear Projection and Head Separation

The input sequence `x` undergoes linear projection (`self.qkv_proj`) to obtain the query (`q`), key (`k`), and value (`v`) matrices. These are then reshaped and permuted to separate the heads.

Dot Product Attention

The scaled dot product attention is applied to the separated query, key, and value matrices, producing attention values and attention weights.

Concatenation and Output Projection

The attention values from different heads are concatenated and reshaped to the desired output dimensions. The concatenated attention values are linearly projected using `self.o_proj` to produce the final output.

Return Statement

The method returns either the output values or both the output values and attention weights based on the `return_attention` parameter.

Overall Purpose

The class provides a modular implementation of a multihead attention module, following the design principles of the original Transformer architecture. Its primary utility lies in integration within neural network models, especially those dealing with sequential data. The code prioritizes modularity for flexibility and efficiency in handling attention computations across multiple heads. One crucial property of multi-head attention is its permutation-equivariance with respect to inputs. This implies that when swapping two elements in the input sequence (e.g., `X1-X2`, ignoring the batch dimension), the output remains the same except for the interchange of elements 1 and

2. Consequently, multi-head attention processes the input as a set of elements rather than a fixed sequence.

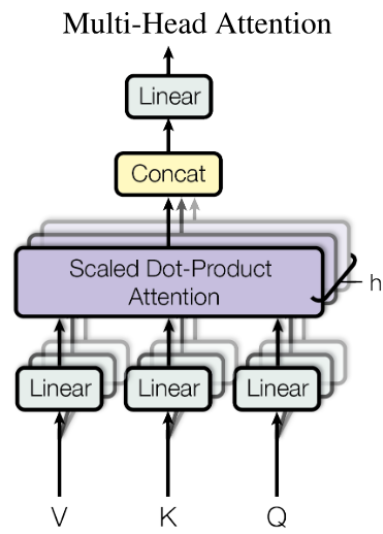


Figure 1: Multi-Head Attention

2.3 Encoder Block

The encoder is composed of N identical blocks applied sequentially. Given input x , it undergoes a Multi-Head Attention block, as implemented previously. The output is then combined with the original input using a residual connection, and a subsequent Layer Normalization is applied to the sum. In summary, it computes $\text{LayerNorm}(x + \text{Multihead}(x, x, x))$, with x as the input to the attention layer (Q, K, and V).

2.3.1 Code

Code extraction 3: EncoderBlock Class

```
1 import torch.nn as nn
2
3 class EncoderBlock(nn.Module):
4     def __init__(self, input_dim, num_heads, dim_feedforward,
5         dropout=0.0):
6         """
7         Args:
8             input_dim: Dimensionality of the input
9             num_heads: Number of heads to use in the attention
10                block
11             dim_feedforward: Dimensionality of the hidden
12                layer in the MLP
13             dropout: Dropout probability to use in the dropout
14                layers
15         """
16         super().__init__()
17
18         #####
19         ### YOUR CODE HERE! ###
20         #####
21
22         # Create Attention layer
23         self.self_attn = MultiheadAttention(input_dim,
24             input_dim, num_heads)
25
26         # Create Two-layer MLP with dropout
27         self.mlp = nn.Sequential(
28             nn.Linear(input_dim, input_dim*2),
29             nn.ReLU(),
30             nn.Dropout(dropout),
31             nn.Linear(2*input_dim, input_dim)
32         )
33
34         # Layers to apply in between the main layers (Layer
35             Norm and Dropout)
36         self.norm = nn.Sequential(
37             nn.LayerNorm(input_dim),
38             nn.Dropout(dropout)
39         )
```

```

34
35 #####
36 ###      END      ###
37 #####
38
39 def forward(self, x, mask=None):
40     # Compute Attention part
41     attn = self.self_attn(x)
42     x = self.norm(x + attn)
43
44     # Compute MLP part
45     x = self.norm(x + self.mlp(x))
46
47     return x

```

2.3.2 Explanation

Class Definition

The `EncoderBlock` class is a neural network module representing an encoder block in a transformer architecture.

Initialization

The constructor initializes the encoder block's parameters, including the input dimension, the number of attention heads, the dimensionality of the hidden layer in the MLP (multi-layer perceptron), and an optional dropout probability.

Attention Layer

An attention layer (`self_attn`) is created using the `MultiheadAttention` class with the specified input dimension, acting as a self-attention mechanism.

MLP (Multi-Layer Perceptron)

A two-layer MLP (`mlp`) is constructed with an initial linear layer, ReLU activation, dropout, and a final linear layer.

Layer Normalization and Dropout

Two sequential layers (`norm`) are defined, consisting of layer normalization followed by dropout. These layers are applied in between the main layers (attention and MLP) to improve training stability.

Forward Method

The `forward` method executes the forward pass of the encoder block.

Attention Part

The attention mechanism is computed by applying the self-attention layer to the input. The result is normalized and added to the original input.

MLP Part

The MLP is applied to the normalized input, and the result is again normalized and added to the original input.

Return Statement

The method returns the final output of the encoder block.

Overall Purpose

The `EncoderBlock` class serves as a container for a transformer encoder block, amalgamating a self-attention mechanism and a multi-layer perceptron. It is engineered to capture intricate patterns and dependencies within input sequences, augmenting the overall capabilities of the transformer model. The code prioritizes modularity and flexibility, emphasizing seamless integration of attention and MLP components.

2.4 Transformer Predictor

The provided template outlines a classifier built upon the Transformer encoder for sequence prediction. In addition to the Transformer architecture, it incorporates:

- A compact input network, responsible for mapping input dimensions to model dimensions.
- The inclusion of positional encoding.
- An output network tasked with transforming output encodings into predictions.

It's essential to note that the output network operates on a 3D tensor (batch samples, sequence length, model dimension) and produces a 2D tensor (batch samples, sequence length). Each output value signifies the prediction for the corresponding reversed number.

2.4.1 Code

Code extraction 4: TransformerPredictor Class

```
1 import torch.nn as nn
2
3 class TransformerPredictor(nn.Module):
4     def __init__(
5         self,
6         input_dim,
7         model_dim,
8         num_classes,
9         num_heads,
10        num_layers,
11        dropout=0.0,
12        input_dropout=0.0,
13    ):
14        """
15        Args:
16            input_dim: Hidden dimensionality of the input
17            model_dim: Hidden dimensionality to use inside the
                        Transformer
```

```

18         num_classes: Number of classes to predict per
           sequence element
19         num_heads: Number of heads to use in the
           Multi-Head Attention blocks
20         num_layers: Number of encoder blocks to use.
21         lr: Learning rate in the optimizer
22         warmup: Number of warmup steps. Usually between 50
           and 500
23         max_iters: Number of maximum iterations the model
           is trained for. This is needed for the
           CosineWarmup scheduler
24         dropout: Dropout to apply inside the model
25         input_dropout: Dropout to apply on the input
           features
26     """
27     super().__init__()
28     self.input_dim = input_dim
29     self.model_dim = model_dim
30     self.num_classes = num_classes
31     self.num_heads = num_heads
32     self.num_layers = num_layers
33     self.dropout = dropout
34     self.input_dropout = input_dropout
35
36     #####
37     ### YOUR CODE HERE! ###
38     #####
39
40     # Create a Generic Input Encoder Input dim -> Model
       dim with input dropout
41     self.input_net = nn.Sequential(
42         nn.Linear(input_dim, model_dim),
43         nn.Dropout(input_dropout)
44     )
45
46     # Create positional encoding for sequences
47     self.positional_encoding =
       PositionalEncoding(model_dim, max_len=16)
48
49     # Create transformer Encoder
50     self.transformer = TransformerEncoder(num_layers,
       input_dim=model_dim, dim_feedforward=model_dim*2,
       num_heads=num_heads, dropout=dropout)
51
52     # Create output classifier per sequence element
       Model_dim -> num_classes
53     self.output_net = nn.Linear(model_dim, num_classes)
54
55     def forward(self, x, mask=None,
       add_positional_encoding=True):

```

```

56     """
57     Args:
58         x: Input features of shape [Batch, SeqLen,
59           input_dim]
60         mask: Mask to apply on the attention outputs
61           (optional)
62         add_positional_encoding: If True, we add the
63           positional encoding to the input.
64                                     Might not be desired for
65                                     some tasks.
66     """
67
68     x = self.input_net(x)
69     if add_positional_encoding:
70         x = self.positional_encoding(x)
71     x = self.transformer(x, mask=mask)
72     x = self.output_net(x)
73
74     #####
75     ###         END         ###
76     #####
77
78     return x
79
80 @torch.no_grad()
81 def get_attention_maps(self, x, mask=None,
82   add_positional_encoding=True):
83     """Function for extracting the attention matrices of
84       the whole Transformer for a single batch.
85
86       Input arguments same as the forward pass.
87     """
88
89     x = self.input_net(x)
90     if add_positional_encoding:
91         x = self.positional_encoding(x)
92     attention_maps =
93         self.transformer.get_attention_maps(x, mask=mask)
94     return attention_maps

```

2.4.2 Explanation

Class Definition

The `TransformerPredictor` class is a neural network module designed for sequence prediction tasks using a transformer architecture.

Initialization

The constructor initializes the predictor's parameters, including the input dimension, model dimension, number of classes to predict per sequence element, number of attention heads, number of encoder blocks, and optional dropout probabilities for the model and input features.

Input Network

A generic input encoder (`input_net`) is created, consisting of a linear layer mapping input dimension to model dimension, followed by dropout on the input features.

Positional Encoding

A positional encoding module is instantiated (`positional_encoding`) to provide positional information to the input sequences.

Transformer Encoder

A transformer encoder (`transformer`) is created with the specified number of layers, input dimension, hidden layer dimension, number of attention heads, and dropout. This component captures complex dependencies in the input sequences.

Output Network

An output classifier (`output_net`) is defined, mapping the model dimension to the number of classes. This layer produces predictions for each sequence element.

Forward Method

The `forward` method executes the forward pass of the predictor.

Input Processing

The input features go through the input encoder (`input_net`) and optionally receive positional encoding.

Transformer Processing

The processed input is fed into the transformer encoder (`transformer`), capturing hierarchical dependencies in the sequence.

Output Generation

The final output is obtained by passing the transformer's output through the output classifier (`output_net`).

Attention Maps Method (`get_attention_maps`)

`get_attention_maps` is defined to extract the attention matrices of the entire transformer for a single batch. It applies the input encoder and positional encoding (if specified) to the input and retrieves attention maps using the `get_attention_maps` method of the transformer.

Overall Purpose

The `TransformerPredictor` class is designed for sequence prediction tasks, employing a transformer architecture to capture both local and global dependencies in input sequences. It prioritizes modularity, enabling users to adjust parameters for diverse tasks. Additionally, the class features attention map extraction, offering insights into the learned attention patterns of the transformer.

3. Questions

3.1 How do we solve the sequence reversal problem at the end of the notebook? Discuss both the method and results.

The sequence reversal problem is tackled through the utilization of a Transformer architecture, employing the following key components:

Dataset Handling

- The `ReverseDataset` class is designed to generate random sequences for training, where labels represent the reversed sequences.

Model Configuration

- The `TransformerPredictor` class is employed, featuring a single encoder block and a single head in the Multi-Head Attention for simplicity.

Training Approach

- Training involves processing sequences through the Transformer encoder.
- Cross-Entropy loss is applied, treating each number as a one-hot vector.
- The AdamW optimizer is utilized.

Training Outcome

- Achieves 100% accuracy on both training and validation sets, as well as on the test set.

Attention Visualization

- Attention maps are extracted post-training, unveiling the model's focus during the sequence reversal task.
- Attention patterns are visually represented using a dedicated plotting function.

Results and Discussion

- The Transformer model effectively solves the sequence reversal task, demonstrating perfect accuracy.
- Attention maps show the model's ability to concentrate on the token at the flipped index, aligning with the task requirements.
- The model adapts to positional encoding patterns, showcasing its flexibility in handling sequential dependencies.

In summary, this solution highlights the Transformer's proficiency in sequential tasks, emphasizing accuracy and providing insight into the model's attention mechanisms through attention map visualization.

3.2 What would be different if we wanted to predict something about the sequence as a whole?

To shift the model's focus from predicting individual sequence elements to making predictions about the sequence as a whole, several adjustments are required:

1. **Output Layer Modification:** The output layer needs to be adapted to consider the entire sequence rather than individual elements. This adjustment involves using appropriate layers or mechanisms, such as pooling layers, recurrent structures, or global attention mechanisms, to aggregate information across the entire sequence.
2. **Choice of Loss Function:** Depending on the nature of the sequence-level prediction task, the loss function must align with the desired output. Common choices include Binary Cross-Entropy, Categorical Cross-Entropy, or Mean Squared Error for tasks involving the sequence as a cohesive entity.
3. **Model Architecture Adjustments:** The Transformer architecture may undergo modifications to better capture global dependencies within the sequence. This could entail changes in the arrangement of Transformer layers or the inclusion of additional components tailored for capturing holistic sequence characteristics.
4. **Input Representation Reevaluation:** The input representation should be scrutinized to ensure it encapsulates comprehensive features relevant to the entire sequence, aligning with the specific requirements of the prediction task.
5. **Enhanced Sequence Encoding:** For tasks focused on the complete sequence, a more sophisticated encoding scheme may be necessary. This could involve the incorporation of intricate encoding mechanisms to better represent the holistic aspects of the sequence.

Bibliography

- [1] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.