

# Attention in PyTorch

In this tutorial, you will implement the attention mechanism of transformers as described in Attention is all you need on your own.

Source: [https://d2l.ai/chapter\\_attention-mechanisms-and-transformers](https://d2l.ai/chapter_attention-mechanisms-and-transformers)

---

First, implement the single-head QKV attention mechanism. This does not include learnable parameters so far. You can use the `masked_softmax` utility below to account for inputs of different lengths that have been padded. Use batch matrix multiplication for efficiency and apply random dropout to the attention scores.

Also note that we need to keep the order of magnitude of the arguments in the exponential function under control. Assume that all the elements of the query  $\mathbf{q} \in \mathbb{R}^d$  and the key  $\mathbf{k}_i \in \mathbb{R}^d$  are independent and identically drawn random variables with zero mean and unit variance. The dot product between both vectors has zero mean and a variance of  $d$ . To ensure that the variance of the dot product still remains 1 regardless of vector length, we use the *scaled dot product attention* scoring function. That is, we rescale the dot product by  $1/\sqrt{d}$ . We thus arrive at the first commonly used attention function that is used, e.g., in Transformers:

$$a(\mathbf{q}, \mathbf{k}_i) = \mathbf{q}^\top \mathbf{k}_i / \sqrt{d}.$$

Note that attention weights  $\alpha$  still need normalizing. We can simplify this further by using the softmax operation:

$$\alpha(\mathbf{q}, \mathbf{k}_i) = \text{softmax}(a(\mathbf{q}, \mathbf{k}_i)) = \frac{\exp(\mathbf{q}^\top \mathbf{k}_i / \sqrt{d})}{\sum_{j=1} \exp(\mathbf{q}^\top \mathbf{k}_j / \sqrt{d})}.$$

```
In [1]: import math
import torch
from torch import nn
```

```
In [ ]: def masked_softmax(X, valid_lens): #@save
        """Perform softmax operation by masking elements on the last axis."""
        # X: 3D tensor, valid_lens: 1D or 2D tensor
        def _sequence_mask(X, valid_len, value=0):
            maxlen = X.size(1)
            mask = torch.arange((maxlen), dtype=torch.float32,
                                device=X.device)[None, :] < valid_len[:, None]
            # In Task 2: output mask = [[True, True, True, False],
            #                             [True, True, False, False]]

            X[~mask] = value # in-place update: assigning value to elements where the m
            return X

        if valid_lens is None:
            return nn.functional.softmax(X, dim=-1)
        else:
            shape = X.shape
```

```

    if valid_lens.dim() == 1:
        valid_lens = torch.repeat_interleave(valid_lens, shape[1])
    else:
        valid_lens = valid_lens.reshape(-1)
    # On the last axis, replace masked elements with a very large negative
    # value, whose exponentiation outputs 0
    X = _sequence_mask(X.reshape(-1, shape[-1]), valid_lens, value=-1e6)
    return nn.functional.softmax(X.reshape(shape), dim=-1)

```

```

In [ ]: class DotProductAttention(nn.Module):
        """Scaled dot product attention."""
        def __init__(self, dropout):
            super().__init__()
            self.dropout = nn.Dropout(dropout)

        # Shape of queries: (batch_size, no. of queries, d)
        # Shape of keys: (batch_size, no. of key-value pairs, d)
        # Shape of values: (batch_size, no. of key-value pairs, value dimension)
        # Shape of valid_lens: (batch_size,) or (batch_size, no. of queries)
        def forward(self, queries, keys, values, valid_lens=None):
            d = queries.shape[-1]
            # Swap the last two dimensions of keys with keys.transpose(1, 2)
            scores = torch.bmm(queries, keys.transpose(1, 2)) / math.sqrt(d)
            self.attention_weights = masked_softmax(scores, valid_lens)
            return torch.bmm(self.dropout(self.attention_weights), values)

```

Dropout is a regularization technique used to prevent overfitting in neural networks. It works by randomly setting a fraction of the elements in a tensor to zero during training. This forces the model to not overly rely on any specific connection or weight.

dropout=0.5 means that 50% of the attention weights will be set to zero randomly.

Test your implementation with the following code. Is the shape of the output correct?

```

In [4]: queries = torch.normal(0, 1, (2, 3, 2))
        keys = torch.normal(0, 1, (2, 10, 2))
        values = torch.normal(0, 1, (2, 10, 4))

        attention = DotProductAttention(dropout=0.5)

        output = attention(queries, keys, values)
        output.shape

```

```

Out[4]: torch.Size([2, 3, 4])

```

Now, implement the multihead attention. Here, we use learnable, linear mappings to project the input queries, keys, and values to queries, keys, and values for different heads.

```

In [5]: class MultiHeadAttention(nn.Module):
        """Multi-head attention."""
        def __init__(self, num_hiddens, num_heads, dropout, bias=False, **kwargs):
            super().__init__()
            self.num_heads = num_heads
            self.attention = DotProductAttention(dropout)
            self.W_q = nn.Linear(num_hiddens, num_hiddens, bias=bias)
            self.W_k = nn.Linear(num_hiddens, num_hiddens, bias=bias)
            self.W_v = nn.Linear(num_hiddens, num_hiddens, bias=bias)
            self.W_o = nn.Linear(num_hiddens, num_hiddens, bias=bias)

        def transpose_qkv(self, X):

```

```

        """Transposition for parallel computation of multiple attention heads."""
        # Shape of input X: (batch_size, no. of queries or key-value pairs,
        # num_hiddens). Shape of output X: (batch_size, no. of queries or
        # key-value pairs, num_heads, num_hiddens / num_heads)
        X = X.reshape(X.shape[0], X.shape[1], self.num_heads, -1)
        # Shape of output X: (batch_size, num_heads, no. of queries or key-value
        # pairs, num_hiddens / num_heads)
        X = X.permute(0, 2, 1, 3)
        # Shape of output: (batch_size * num_heads, no. of queries or key-value
        # pairs, num_hiddens / num_heads)
        return X.reshape(-1, X.shape[2], X.shape[3])

    def transpose_output(self, X):
        """Reverse the operation of transpose_qkv."""
        X = X.reshape(-1, self.num_heads, X.shape[1], X.shape[2])
        X = X.permute(0, 2, 1, 3)
        return X.reshape(X.shape[0], X.shape[1], -1)

    def forward(self, queries, keys, values, valid_lens):
        # Shape of queries, keys, or values:
        # (batch_size, no. of queries or key-value pairs, num_hiddens)
        # Shape of valid_lens: (batch_size,) or (batch_size, no. of queries)
        # After transposing, shape of output queries, keys, or values:
        # (batch_size * num_heads, no. of queries or key-value pairs,
        # num_hiddens / num_heads)
        queries = self.transpose_qkv(self.W_q(queries))
        keys = self.transpose_qkv(self.W_k(keys))
        values = self.transpose_qkv(self.W_v(values))

        if valid_lens is not None:
            # On axis 0, copy the first item (scalar or vector) for num_heads
            # times, then copy the next item, and so on
            valid_lens = torch.repeat_interleave(
                valid_lens, repeats=self.num_heads, dim=0)

        # Shape of output: (batch_size * num_heads, no. of queries,
        # num_hiddens / num_heads)
        output = self.attention(queries, keys, values, valid_lens)
        # Shape of output_concat: (batch_size, no. of queries, num_hiddens)
        output_concat = self.transpose_output(output)
        return self.W_o(output_concat)

```

Test your implementation with the following code. Is the shape of the output correct?

```

In [6]: num_hiddens, num_heads = 100, 5
        attention = MultiHeadAttention(num_hiddens, num_heads, 0.5)
        batch_size, num_queries, num_kvpairs = 2, 4, 6
        valid_lens = torch.tensor([3, 2])
        X = torch.normal(0, 1, (batch_size, num_queries, num_hiddens))
        Y = torch.normal(0, 1, (batch_size, num_kvpairs, num_hiddens))
        output = attention(X, Y, Y, valid_lens)
        output.shape

```

```

/opt/anaconda3/envs/codefusion/lib/python3.9/site-packages/torch/nn/modules/lazy.py:
180: UserWarning: Lazy modules are a new feature under heavy development so changes
to the API or functionality can happen at any moment.

```

```

.. warnings.warn('Lazy modules are a new feature under heavy development '
Out[6]: torch.Size([2, 4, 100])

```