

Transformer Practice

컴퓨터소프트웨어학부 심승현

Contents

<https://www.tensorflow.org/text/tutorials/transformer>

Tensorflow.org, transformer tutorial

translating Portuguese to English using transformer model

Review : Positional Encoding

```
def get_angles(pos, i, d_model):  
    angle_rates = 1 / np.power(10000. (2 * (i//2)) / np.float32(d_model))  
    return pos * angle_rates  
  
def positional_encoding(position, d_model):  
    angle_rads = get_angles(np.arange(position)[:, np.newaxis],  
                             np.arange(d_model)[np.newaxis, :],  
                             d_model)  
  
    # apply sin to even indices in the array; 2i  
    angle_rads[:, 0::2] = np.sin(angle_rads[:, 0::2])  
  
    # apply cos to odd indices in the array; 2i+1  
    angle_rads[:, 1::2] = np.cos(angle_rads[:, 1::2])  
  
    pos_encoding = angle_rads[np.newaxis, ...]  
  
    return tf.cast(pos_encoding, dtype=tf.float32)
```

input parameters are both int type

$$PE_{(pos, 2i)} = \sin(pos / 10000^{2i/d_{model}})$$

$$PE_{(pos, 2i+1)} = \cos(pos / 10000^{2i/d_{model}})$$

Review : Positional Encoding

```
def get_angles(pos, i, d_model):  
    angle_rates = 1 / np.power(10000, (2 * (i//2)) / np.float32(d_model))  
    return pos * angle_rates  
  
def positional_encoding(position, d_model):  
    angle_rads = get_angles(np.arange(position)[:, np.newaxis],  
                             np.arange(d_model)[np.newaxis, :],  
                             d_model)  
  
    # apply sin to even indices in the array; 2i  
    angle_rads[:, 0::2] = np.sin(angle_rads[:, 0::2])  
  
    # apply cos to odd indices in the array; 2i+1  
    angle_rads[:, 1::2] = np.cos(angle_rads[:, 1::2])  
  
    pos_encoding = angle_rads[np.newaxis, ...]  
  
    return tf.cast(pos_encoding, dtype=tf.float32)
```

$$PE_{(pos, 2i)} = \sin\left(pos / 10000^{2i/d_{model}}\right)$$
$$PE_{(pos, 2i+1)} = \cos\left(pos / 10000^{2i/d_{model}}\right)$$

Review : Positional Encoding

```
def get_angles(pos, i, d_model):  
    angle_rates = 1 / np.power(10000, (2 * (i//2)) / np.float32(d_model))  
    return pos * angle_rates  
  
def positional_encoding(position, d_model):  
    angle_rads = get_angles(np.arange(position)[:, np.newaxis],  
                             np.arange(d_model)[np.newaxis, :],  
                             d_model)  
  
    # apply sin to even indices in the array; 2i  
    angle_rads[:, 0::2] = np.sin(angle_rads[:, 0::2])  
  
    # apply cos to odd indices in the array; 2i+1  
    angle_rads[:, 1::2] = np.cos(angle_rads[:, 1::2])  
  
    pos_encoding = angle_rads[np.newaxis, ...]  
  
    return tf.cast(pos_encoding, dtype=tf.float32)
```

$$PE_{(pos, 2i)} = \sin(pos / 10000^{2i/d_{model}})$$
$$PE_{(pos, 2i+1)} = \cos(pos / 10000^{2i/d_{model}})$$

dimension of return matrix is
(1, 10000, 512)

Review : Multi Head Attention

- First, get Query, Key, Value matrix from input matrix X

$$\begin{aligned} Q &= X @ W_Q \\ K &= X @ W_K \\ V &= X @ W_V \end{aligned}$$

- Then, derive attention matrix Z through the formula behind

$$Z = \text{softmax} \left(\frac{Q \cdot K^T}{\sqrt{d_K}} \right) V$$

※ Z is attention matrix, d_K is dimension of matrix K.

Review : Multi Head Attention

```
class EncoderLayer(tf.keras.layers.Layer):  
    def __init__(self, d_model, num_heads, dff, rate=0.1):  
        super(EncoderLayer, self).__init__()  
        self.mha = MultiHeadAttention(d_model, num_heads)  
    def call(self, x, training, mask):  
        attn_output, _ = self.mha(x, x, x, mask)
```

```
class MultiHeadAttention(tf.keras.layers.Layer):  
    def call(self, v, k, q, mask):  
        batch_size = tf.shape(q)[0]  
        q = self.wq(q) # q.shape (batch_size, seq_len, d_model)  
        k = self.wk(k)  
        v = self.wv(v)  
        q = self.split_heads(q, batch_size)  
        k = self.split_heads(k, batch_size)  
        v = self.split_heads(v, batch_size)  
        scaled_attention, attention_weights = scaled_dot_product_attention(  
            q, k, v, mask)
```

- First, get Query, Key, Value matrix from input matrix X

$$\begin{aligned}Q &= X @ W_Q \\K &= X @ W_K \\V &= X @ W_V\end{aligned}$$

- Then, derive attention matrix through the formula behind

$$Z = \text{softmax}\left(\frac{Q \cdot K^T}{\sqrt{d_K}}\right) V$$

※ Z is attention matrix, d_K is dimension of matrix K.

Review : Multi Head Attention

```
class EncoderLayer(tf.keras.layers.Layer):  
    def __init__(self, d_model, num_heads, dff, rate=0.1):  
        super(EncoderLayer, self).__init__()  
        self.mha = MultiHeadAttention(d_model, num_heads)  
    def call(self, x, training, mask):  
        attn_output, _ = self.mha(x, x, x, mask)
```

```
class MultiHeadAttention(tf.keras.layers.Layer):
```

```
    def call(self, v, k, q, mask):  
        batch_size = tf.shape(q)[0]
```

```
        q = self.wq(q) # (batch_size,  
        k = self.wk(k)  
        v = self.wv(v)
```

```
        q = self.split_heads(q, batch_size)  
        k = self.split_heads(k, batch_size)  
        v = self.split_heads(v, batch_size)
```

```
        scaled_attention, attention_weights = scaled_dot_product_attention(  
            q, k, v, mask)
```

```
        self.wq = tf.keras.layers.Dense(d_model)  
        self.wk = tf.keras.layers.Dense(d_model)  
        self.wv = tf.keras.layers.Dense(d_model)
```

v, k, q passes different dense layers,
which have different weight matrix and bias.
In this step, Q, K, V are defined.

- First, get Query, Key, Value matrix from input matrix X

$$\begin{aligned}Q &= X @ W_Q \\ K &= X @ W_K \\ V &= X @ W_V\end{aligned}$$

- Then, derive attention matrix through the formula behind

$$Z = \text{softmax}\left(\frac{Q \cdot K^T}{\sqrt{d_k}}\right) V$$

d_k is dimension of matrix K.

Review : Multi Head Attention

```
class EncoderLayer(tf.keras.layers.Layer):  
    def __init__(self, d_model, num_heads, dff, rate=0.1):  
        super(EncoderLayer, self).__init__()  
  
        self.mha = MultiHeadAttention(d_model, num_heads)  
    def call(self, x, training, mask):  
  
        attn_output, _ = self.mha(x, x, x, mask)
```

```
class MultiHeadAttention(tf.keras.layers.Layer):
```

```
    def call(self, v, k, q, mask):  
        batch_size = tf.shape(q)[0]
```

```
        q = self.wq(q) # (batch_size, seq_len, d_model)  
        k = self.wk(k)  
        v = self.wv(v)
```

```
    def split_heads(self, x, batch_size):  
        x = tf.reshape(x, (batch_size, -1, self.num_heads, self.depth))  
        return tf.transpose(x, perm=[0, 2, 1, 3])
```

```
        q = self.split_heads(q, batch_size)  
        k = self.split_heads(k, batch_size)  
        v = self.split_heads(v, batch_size)
```

```
        scaled_attention, attention_weights =  
            self._scaled_dot_product_attention(  
                q, k, v, mask)
```

- First, get Query, Key, Value matrix from input matrix X

$$Q = X @ W_Q$$

$$K = X @ W_K$$

$$V = X @ W_V$$

perm of transpose
tf.transpose(x, perm=[0, 2, 1, 3]) moves
x[a][b][c][d] to x[a][c][b][d]

q.shape
(batch_size, seq_len, d_model)

->

(batch_size, self.num_heads, seq_len, self.depth)

$$\frac{QK^T}{\sqrt{d_K}} V$$

d_K is dimension of matrix K.

Review : Multi Head Attention

```
class EncoderLayer(tf.keras.layers.Layer):  
    def __init__(self, d_model, num_heads, dff, rate=0.1):  
        super(EncoderLayer, self).__init__()  
  
        self.mha = MultiHeadAttention(d_model, num_heads)  
    def call(self, x, training, mask):  
  
        attn_output, _ = self.mha(x, x, x, mask)
```

```
class MultiHeadAttention(tf.keras.layers.Layer):  
    def call(self, v, k, q, mask):  
        batch_size = tf.shape(q)[0]  
  
        q = self.wq(q) # (batch_size, seq_len, d_model)  
        k = self.wk(k)  
        v = self.wv(v)  
  
        q = self.split_heads(q, batch_size)  
        k = self.split_heads(k, batch_size)  
        v = self.split_heads(v, batch_size)  
  
        scaled_attention, attention_weights = scaled_dot_product_attention(  
            q, k, v, mask)
```

let's go to scaled dot product attention!

- First, get Query, Key, Value matrix from input matrix X

$$\begin{aligned}Q &= X @ W_Q \\K &= X @ W_K \\V &= X @ W_V\end{aligned}$$

- Then, derive attention matrix through the formula behind

$$Z = \text{softmax}\left(\frac{Q \cdot K^T}{\sqrt{d_K}}\right) V$$

※ Z is attention matrix, d_K is dimension of matrix K.

Review : Multi Head Attention

q.shape
(batch_size, self.num_heads, seq_len, self.depth)

```
def scaled_dot_product_attention(q, k, v, mask):  
  
    matmul_qk = tf.matmul(q, k, transpose_b=True)  
  
    dk = tf.cast(tf.shape(k)[-1], tf.float32)  
    scaled_attention_logits = matmul_qk / tf.math.sqrt(dk)  
  
    if mask is not None:  
        scaled_attention_logits += (mask * -1e9)  
    attention_weights = tf.nn.softmax(scaled_attention_logits, axis=-1)  
    output = tf.matmul(attention_weights, v)  
    return output, attention_weights
```

► First, get Query, Key, Value matrix
matrix X

$$\begin{aligned} Q &= X @ W_Q \\ K &= X @ W_K \\ V &= X @ W_V \end{aligned}$$

► Then, derive attention matrix
through the formula behind

$$Z = \text{softmax}\left(\frac{Q \cdot K^T}{\sqrt{d_K}}\right) V$$

※ Z is attention matrix, d_K is dimension of matrix K.

Review : Multi Head Attention

- First, get Query, Key, Value matrix
matrix X

q.shape
(batch_size, self.num_heads, seq_len, self.depth)

```
def scaled_dot_product_attention(q, k, v, mask):
```

```
    matmul_qk = tf.matmul(q, k, transpose_b=True)
```

```
    dk = tf.cast(tf.shape(k), tf.float32)
    scaled_attention_logits =
```

matmul_qk.shape
(batch_size, self.num_heads, seq_len, seq_len)

```
    if mask is not None:
```

```
        scaled_attention_logits += (mask * -1e9)
```

```
    attention_weights = tf.nn.softmax(scaled_attention_logits, axis=-1)
```

```
    output = tf.matmul(attention_weights, v)
```

```
    return output, attention_weights
```

$$Q = X @ W_Q$$

$$K = X @ W_K$$

$$V = X @ W_V$$

- Then, derive attention matrix
through the formula behind

$$Z = \text{softmax}\left(\frac{Q \cdot K^T}{\sqrt{d_K}}\right) V$$

※ Z is attention matrix, d_K is dimension of matrix K.

Review : Multi Head Attention

- First, get Query, Key, Value matrix
matrix X

q.shape
(batch_size, self.num_heads, seq_len, self.depth)

```
def scaled_dot_product_attention(q, k, v, mask):  
    matmul_qk = tf.matmul(q, k, transpose_b=True)  
  
    dk = tf.cast(tf.shape(k)[-1], tf.float32)  
    scaled_attention_logits = matmul_qk / dk  
  
    if mask is not None:  
        scaled_attention_logits += (mask * -1e9)  
    attention_weights = tf.nn.softmax(scaled_attention_logits, axis=-1)  
    output = tf.matmul(attention_weights, v)  
    return output, attention_weights
```

dk is a tf.float32 number

$$\begin{aligned} Q &= X @ W_Q \\ K &= X @ W_K \\ V &= X @ W_V \end{aligned}$$

- Then, derive attention matrix through the formula behind

$$Z = \text{softmax} \left(\frac{Q \cdot K^T}{\sqrt{d_K}} \right) V$$

Note: In the original image, the term $Q \cdot K^T$ is highlighted with a red box and labeled 'matmul_qk' in red, and d_K is highlighted with a blue box and labeled 'dk' in blue.

※ Z is attention matrix, d_K is dimension of matrix K.

Review : Multi Head Attention

- First, get Query, Key, Value matrix from input matrix X

$$\begin{aligned} Q &= X @ W_Q \\ K &= X @ W_K \\ V &= X @ W_V \end{aligned}$$

```
def scaled_dot_product_attn
```

```
    matmul_qk = tf.matmul(q, k, transpose_b=True)
```

```
    dk = tf.cast(tf.shape(k)[-1], tf.float32)
```

```
    scaled_attention_logits = matmul_qk / tf.math.sqrt(dk)
```

```
    if mask is not None:
```

```
        scaled_attention_logits = tf.nn.softmax(scaled_attention_logits, axis=-1)
```

```
    output = tf.matmul(scaled_attention_logits, v)
```

```
    return output, scaled_attention_logits
```

matmul_qk.shape
(batch_size, self.num_heads, seq_len, seq_len)

scaled_attention_logits.shape
(batch_size, self.num_heads, seq_len, seq_len)

- then, derive attention matrix through the formula behind

$$Z = \text{softmax}\left(\frac{Q \cdot K^T}{\sqrt{d_K}}\right) V$$

※ Z is attention matrix, d_K is dimension of matrix K.

Review : Multi Head Attention

```
def scaled_dot_product_attention(q, k, v, mask):  
    matmul_qk = tf.matmul(q, k, transpose_b=True)  
  
    dk = tf.cast(tf.shape(k)[-1], tf.float32)  
    scaled_attention_logits = matmul_qk / tf.math.sqrt(dk)  
  
    if mask is not None:  
        scaled_attention_logits += (mask * -1e9)  
    attention_weights = tf.nn.softmax(scaled_attention_logits, axis=-1)  
    output = tf.matmul(attention_weights, v)  
    return output, attention_weights
```

- Padding mask matrix consists of 1 and 0
1 means mask

seq is token ID matrix, and its size is
(batch_size, seq_len)

```
def create_padding_mask(seq):  
    seq = tf.cast(tf.math.equal(seq, 0), tf.float32)  
  
    # add extra dimensions to add the padding  
    # to the attention logits.  
    return seq[:, tf.newaxis, tf.newaxis, :] # (batch_size, 1, 1, seq_len)
```

- First, get Query, Key, Value matrix from input matrix X

$$\begin{aligned} Q &= X @ W_Q \\ K &= X @ W_K \\ V &= X @ W_V \end{aligned}$$

- Then, derive attention matrix through the formula behind

$$Z = \text{softmax}\left(\frac{Q \cdot K^T}{\sqrt{d_K}}\right) V$$

※ Z is attention matrix, d_K is dimension of matrix K.

Review : Multi Head Attention

```
def scaled_dot_product_attention(q, k, v, mask):  
    matmul_qk = tf.matmul(q, k, transpose_b=True)  
  
    dk = tf.cast(tf.shape(k)[-1], tf.float32)  
    scaled_attention_logits = matmul_qk / tf.math.sqrt(dk)  
  
    if mask is not None:  
        scaled_attention_logits += (mask * -1e9)  
    attention_weights = tf.nn.softmax(scaled_attention_logits, axis=-1)  
    output = tf.matmul(attention_weights, v)  
    return output, attention_weights
```

- Padding mask matrix consists of 1 and 0
1 means mask

```
def create_padding_mask(seq):  
    seq = tf.cast(tf.math.equal(seq, 0), tf.float32)  
  
    # add extra dimensions to add the padding  
    # to the attention logits.  
    return seq[:, tf.newaxis, tf.newaxis, :]
```

seq is token ID matrix, and its size is
(batch_size, seq_len)

seq is now mask matrix, and its size is
(batch_size, 1, 1, seq_len)

- First, get Query, Key, Value matrix
from input matrix X

$$\begin{aligned}Q &= X @ W_Q \\ K &= X @ W_K \\ V &= X @ W_V\end{aligned}$$

- Then, derive attention matrix
through the formula behind

$$Z = \text{softmax}\left(\frac{Q \cdot K^T}{\sqrt{d_K}}\right) V$$

※ Z is attention matrix, d_K is dimension of matrix K.

Review : Multi Head Attention

- First, get Query, Key, Value matrix from input matrix X

$$\begin{aligned}Q &= X @ W_Q \\K &= X @ W_K \\V &= X @ W_V\end{aligned}$$

```
def scaled_dot_product_attention(q, k, v, mask):
```

```
    matmul_qk = tf.matmul(q, k, transpose_b=True)
```

```
    dk = tf.cast(tf.shape(k)[-1], tf.float32)
```

```
    scaled_attention_logits = matmul_qk / dk
```

```
    if mask is not None:
```

```
        scaled_attention_logits += (mask * -1e9)
```

```
    attention_weights = tf.nn.softmax(scaled_attention_logits, axis=-1)
```

```
    output = tf.matmul(attention_weights, v)
```

```
    return output, attention_weights
```

scaled_attention_logits.shape
(batch_size, self.num_heads, seq_len, seq_len)

mask.shape
(batch_size, 1, 1, seq_len)

- Then, derive attention matrix through the formula behind

$$Z = \text{softmax}\left(\frac{Q \cdot K^T}{\sqrt{d_K}}\right) V$$

※ Z is attention matrix, d_K is dimension of matrix K.

- Padding mask matrix consists of 1 and 0
1 means mask
- In mask matrix, mask elements are -10^9 and otherwise 0
- Add mask matrix to scaled_attention_logits

Review : Multi Head Attention

```
def scaled_dot_product_attention(q, k, v, mask):  
    matmul_qk = tf.matmul(q, k, transpose_b=True)  
  
    dk = tf.cast(tf.shape(k)[-1], tf.float32)  
    scaled_attention_logits = matmul_qk / tf.math.sqrt(dk)  
  
    if mask is not None:  
        scaled_attention_logits += (mask * -1e9)  
  
    attention_weights =  
    output = tf.matmul(scaled_attention_logits, v)  
    return output, attention_weights
```

for the result of masking
scaled_attention_logits will have
 -10^9 for the pad token.

- ▶ Padding mask matrix consists of 1 and 0
1 means mask
- ▶ In mask matrix, mask elements are -10^9
and otherwise 0
- ▶ Add mask matrix to scaled_attention_logits

- ▶ First, get Query, Key, Value matrix
from input matrix X

$$\begin{aligned} Q &= X @ W_Q \\ K &= X @ W_K \\ V &= X @ W_V \end{aligned}$$

- ▶ Then, derive attention matrix
through the formula behind

$$Z = \text{softmax}\left(\frac{Q \cdot K^T}{\sqrt{d_K}}\right) V$$

※ Z is attention matrix, d_K is dimension of matrix K.

Review : Multi Head Attention

```
def scaled_dot_product_attention(q, k, v, mask):  
    matmul_qk = tf.matmul(q, k, transpose_b=True)  
  
    dk = tf.cast(tf.shape(k)[-1], tf.float32)  
    scaled_attention_logits = matmul_qk / tf.math.sqrt(dk)  
  
    if mask is not None:  
        scaled_attention_logits += (mask * -1e9)  
    attention_weights = tf.nn.softmax(scaled_attention_logits, axis=-1)  
    output = tf.matmul(attention_weights, v)  
    return output, attention_weights
```

- First, get Query, Key, Value matrix from input matrix X

$$\begin{aligned}Q &= X @ W_Q \\ K &= X @ W_K \\ V &= X @ W_V\end{aligned}$$

- Then, derive attention matrix through the formula behind

$$\boxed{Z} = \text{softmax}\left(\frac{Q \cdot K^T}{\sqrt{d_K}}\right) V$$

output attention_weights

※ Z is attention matrix, d_K is dimension of matrix K.