```
13  import math
14
15  import torch
16  import torch.nn as nn
17
18  from labml_nn.utils import clone_module_list
19  from .feed_forward import FeedForward
20  from .mha import MultiHeadAttention
21  from .positional_encoding import get_positional_encoding
```

# Transformer Encoder and Decoder Models

CO Open in Colab

## Embed tokens and add fixed positional encoding

```
24  class EmbeddingsWithPositionalEncoding(nn.Module):
```

```
31      def __init__(self, d_model: int, n_vocab: int, max_len: int = 5000):
32          super().__init__()
33          self.linear = nn.Embedding(n_vocab, d_model)
34          self.d_model = d_model
35          self.register_buffer('positional_encodings', get_positional_encoding(d_model, max_len))
```

```
37      def forward(self, x: torch.Tensor):
38          pe = self.positional_encodings[:x.shape[0]].requires_grad_(False)
39          return self.linear(x) * math.sqrt(self.d_model) + pe
```

## Embed tokens and add parameterized positional encodings

```
42  class EmbeddingsWithLearnedPositionalEncoding(nn.Module):
```

```
49      def __init__(self, d_model: int, n_vocab: int, max_len: int = 5000):
50          super().__init__()
51          self.linear = nn.Embedding(n_vocab, d_model)
52          self.d_model = d_model
53          self.positional_encodings = nn.Parameter(torch.zeros(max_len, 1, d_model), requires_grad=
True)
```

```
55      def forward(self, x: torch.Tensor):
56          pe = self.positional_encodings[:x.shape[0]]
57          return self.linear(x) * math.sqrt(self.d_model) + pe
```

```
60  class TransformerLayer(nn.Module):
```

## Transformer Layer

This can act as an encoder layer or a decoder layer.

☐ Some implementations, including the paper seem to have differences in where the layer-normalization is done. Here we do a layer normalization before attention and feed-forward networks, and add the original residual vectors. Alternative is to do a layer normalization after adding the residuals. But we found this to be less stable when training. We found a detailed discussion about this in the paper On Layer Normalization in the Transformer Architecture.

- `d_model` is the token embedding size
- `self_attn` is the self attention module
- `src_attn` is the source attention module (when this is used in a decoder)
- `feed_forward` is the feed forward module
- `dropout_prob` is the probability of dropping out after self attention and FFN

```
78      def __init__(self, *,
79                   d_model: int,
80                   self_attn: MultiHeadAttention,
81                   src_attn: MultiHeadAttention = None,
82                   feed_forward: FeedForward,
83                   dropout_prob: float):
```

```python
        super().__init__()
        self.size = d_model
        self.self_attn = self_attn
        self.src_attn = src_attn
        self.feed_forward = feed_forward
        self.dropout = nn.Dropout(dropout_prob)
        self.norm_self_attn = nn.LayerNorm([d_model])
        if self.src_attn is not None:
            self.norm_src_attn = nn.LayerNorm([d_model])
        self.norm_ff = nn.LayerNorm([d_model])
```

Whether to save input to the feed forward layer

```python
        self.is_save_ff_input = False
```

```python
    def forward(self, *,
                x: torch.Tensor,
                mask: torch.Tensor,
                src: torch.Tensor = None,
                src_mask: torch.Tensor = None):
```

Normalize the vectors before doing self attention

```python
        z = self.norm_self_attn(x)
```

Run through self attention, i.e. keys and values are from self

```python
        self_attn = self.self_attn(query=z, key=z, value=z, mask=mask)
```

Add the self attention results

```python
        x = x + self.dropout(self_attn)
```

If a source is provided, get results from attention to source. This is when you have a decoder layer that pays attention to encoder outputs

```python
        if src is not None:
```

Normalize vectors

```python
            z = self.norm_src_attn(x)
```

Attention to source. i.e. keys and values are from source

```python
            attn_src = self.src_attn(query=z, key=src, value=src, mask=src_mask)
```

Add the source attention results

```python
            x = x + self.dropout(attn_src)
```

Normalize for feed-forward

```python
        z = self.norm_ff(x)
```

Save the input to the feed forward layer if specified

```python
        if self.is_save_ff_input:
            self.ff_input = z.clone()
```

Pass through the feed-forward network

```python
        ff = self.feed_forward(z)
```

Add the feed-forward results back

```python
        x = x + self.dropout(ff)

        return x
```

## Transformer Encoder

```python
class Encoder(nn.Module):
```

```python
    def __init__(self, layer: TransformerLayer, n_layers: int):
        super().__init__()
```

Make copies of the transformer layer

```python
        self.layers = clone_module_list(layer, n_layers)
```

Final normalization layer

```python
        self.norm = nn.LayerNorm([layer.size])
```

```python
    def forward(self, x: torch.Tensor, mask: torch.Tensor):
```

Run through each transformer layer

```python
        for layer in self.layers:
            x = layer(x=x, mask=mask)
```

Finally, normalize the vectors

```python
        return self.norm(x)
```

## Transformer Decoder

```python
class Decoder(nn.Module):
```

```python
    def forward(self, *,
```

```python
169    def __init__(self, layer: TransformerLayer, n_layers: int):
170        super().__init__()
```

Make copies of the transformer layer

```python
172        self.layers = clone_module_list(layer, n_layers)
```

Final normalization layer

```python
174        self.norm = nn.LayerNorm([layer.size])
```

```python
176    def forward(self, x: torch.Tensor, memory: torch.Tensor, src_mask: torch.Tensor, tgt_mask: torch.Tensor):
```

Run through each transformer layer

```python
178        for layer in self.layers:
179            x = layer(x=x, mask=tgt_mask, src=memory, src_mask=src_mask)
```

Finally, normalize the vectors

```python
181        return self.norm(x)
```

# Generator

This predicts the tokens and gives the lof softmax of those. You don't need this if you are using `nn.CrossEntropyLoss`.

```python
184 class Generator(nn.Module):
```

```python
194    def __init__(self, n_vocab: int, d_model: int):
195        super().__init__()
196        self.projection = nn.Linear(d_model, n_vocab)
```

```python
198    def forward(self, x):
199        return self.projection(x)
```

# Combined Encoder-Decoder

```python
202 class EncoderDecoder(nn.Module):
```

```python
209    def __init__(self, encoder: Encoder, decoder: Decoder, src_embed: nn.Module, tgt_embed: nn.Module, generator: nn.Module):
210        super().__init__()
```

```python
211        self.encoder = encoder
212        self.decoder = decoder
213        self.src_embed = src_embed
214        self.tgt_embed = tgt_embed
215        self.generator = generator
```

This was important from their code. Initialize parameters with Glorot / fan_avg.

```python
219        for p in self.parameters():
220            if p.dim() > 1:
221                nn.init.xavier_uniform_(p)
```

```python
223    def forward(self, src: torch.Tensor, tgt: torch.Tensor, src_mask: torch.Tensor, tgt_mask: torch.Tensor):
```

Run the source through encoder

```python
225        enc = self.encode(src, src_mask)
```

Run encodings and targets through decoder

```python
227        return self.decode(enc, src_mask, tgt, tgt_mask)
```

```python
229    def encode(self, src: torch.Tensor, src_mask: torch.Tensor):
230        return self.encoder(self.src_embed(src), src_mask)
```

```python
232    def decode(self, memory: torch.Tensor, src_mask: torch.Tensor, tgt: torch.Tensor, tgt_mask: torch.Tensor):
233        return self.decoder(self.tgt_embed(tgt), memory, src_mask, tgt_mask)
```