

Transformers PseudoCode for Donut

March 6, 2024

1 Donut Encoder Pseudocode

Algorithm 1: Pseudocode for Donut's SwinTransformer Encoder

Input: $x \in \mathbb{R}^{C \times H \times W}$, a tensor representing the input image
Output: A set of embeddings $\{z_i | z_i \in \mathbb{R}^d, 1 \leq i \leq n\}$, where n is the number of patches and d is the dimension of the embeddings
Hyperparameters: P , patch size,
 L , list of number of layers in each stage,
 H , list of number of heads in each layer,
 W , window size.
Parameters: W_{embed} , parameters for patch embedding,
 W_{norm} , parameters for normalization layer (omitted in Donut),
 W_{trans} , parameters of Swin Transformer layers.

```
1  $Z \leftarrow \text{ImageToTensor}(x)$  ; // Converts image to tensor
2  $Z \leftarrow \text{PatchEmbedding}(Z, W_{\text{embed}}, P)$  ; // Maps image patches to
   embeddings
3  $Z \leftarrow \text{DropPositional}(Z)$  ; // Applies dropout after adding
   positional encodings
4 foreach stage in SwinTransformer do
5   for  $l = 1$  to  $L[\text{stage}]$  do
6      $Z \leftarrow \text{LayerNorm}(Z)$  ; // Applies layer normalization
7      $Z \leftarrow \text{WindowBasedMSA}(Z, W, H[\text{stage}])$  ; // Window based
   multi-head self-attention
8      $Z \leftarrow \text{LayerNorm}(Z)$  ; // Applies layer normalization
   again
9      $Z \leftarrow \text{MLP}(Z)$  ; // Applies MLP block
10    if  $\text{stage} < \text{NumberOfStages}$  then
11       $Z \leftarrow \text{PatchMerging}(Z)$  ; // Downsamples feature map for
   next stage
12    end
13  end
14 end
```

2 Donut Decoder Pseudocode

Algorithm 2: Pseudocode for Donut’s BART-based Decoder

Input: $z \in \mathbb{R}^{n \times d}$, a set of embeddings from the encoder
Output: A sequence of token probabilities $P \in \mathbb{R}^{m \times v}$, where each row represents a distribution over the vocabulary for the corresponding token
Hyperparameters: m , maximum sequence length,
 v , size of token vocabulary,
 L , number of decoder layers,
 d , dimension of latent vectors.
Parameters: W_e , embedding matrix for token embeddings,
 W_p , positional embedding matrix,
 $W_{\text{layer_norm}}$, parameters for layer normalization,
 W_{MHAtt} , multi-head attention parameters,
 W_{MLP} , parameters for the feed-forward network,
 $W_{\text{lm_head}}$, parameters for the language model head.

```

1 Initialize token sequence  $Y = []$ 
2 for  $i = 1$  to  $m$  do
3    $Y_{\text{embed}} \leftarrow W_e Y + W_p$ ;    // Embed tokens and add positional
   encoding
4    $Y \leftarrow \text{LayerNorm}(Y_{\text{embed}}, W_{\text{layer\_norm}})$ ;    // Apply layer
   normalization
5   for  $l = 1$  to  $L$  do
6      $Y \leftarrow \text{MultiHeadAttention}(Y, z, W_{\text{MHAtt}})$ ;    // Self-attention
     over the sequence
7      $Y \leftarrow \text{LayerNorm}(Y, W_{\text{layer\_norm}})$ ;    // Apply layer
     normalization
8      $Y \leftarrow \text{MLP}(Y, W_{\text{MLP}})$ ;    // Feed-forward network
9   end
10   $P_i \leftarrow \text{Softmax}(W_{\text{lm\_head}} Y)$ ;    // Project to vocabulary and
   apply softmax
11   $Y \leftarrow \text{Append}(Y, \text{ArgMax}(P_i))$ ;    // Append the most probable
   next token
12 end

```
