

Bài 37 - Transformer thêm dấu Tiếng Việt

28 May 2020 - phamdinhhkhanh

Menu

- 1. Dữ liệu
- 2. Xây dựng model Transformer
- 3. Position Encoding
- 4. Masking
- 5. Scale dot Product attention
- 6. Multi-head Attention
- 7. Point wise feed forward network
- 8. Encoder and Decoder
 - 8.1. Encoder layer
 - 8.2. Decoder layer
 - 8.3. Encoder
 - 8.4. Decoder
- 9. Create transformer
- 10. Set hyperparameters
- 11. Optimizer
- 12. Loss and Metrics
- 13. Training and checkpointing
- 14. Evaluate
- 15. Pretrain model
- 16. Tổng kết
- 17. Tài liệu tham khảo

1. Dữ liệu

Thêm dấu Tiếng Việt có lẽ đã không còn là một bài toán quá xa lạ đối với cộng đồng AI Việt Nam. Đây là bài toán có tính ứng dụng cao và nhiều doanh nghiệp rất cần. Bản thân công ty mình (Adayroi trước đây) cũng có một dự án như thế. Trước đó thì mình chỉ sử dụng thuật toán LSTM và kết quả trả ra không tốt đối với các câu dài vì hạn chế của các thuật toán có tính tự hồi qui (autoregressive) trong chuỗi thời gian đó là sự phụ thuộc dài hạn kém. Sau đó thì mình chuyển qua kết hợp giữa LSTM là attention và kết quả đã cải thiện hơn. Bạn đọc có thể tham khảo tại Bài 7 - Seq2seq model correct spelling Pytorch

(<https://phamdinhhkhanh.github.io/2019/08/19/CorrectSpellingVietnameseTonePrediction.html>). Trong bài này mình sẽ bỏ qua kiến trúc lõi mền của các thuật toán một chiều (Uni-directional) và thay thế bằng kiến trúc hai chiều (Bi-directional) sử dụng kỹ thuật transformer mà mình đã trình bày ở Bài 4 - Attention is all you need (<https://phamdinhhkhanh.github.io/2019/06/18/AttentionLayer.html>), Bài 36 - BERT model (<https://phamdinhhkhanh.github.io/2020/05/23/BERTModel.html>). Ngoài ra bạn đọc cũng có thể tham khảo các ý tưởng được trình bày tại cuộc thi thêm dấu tiếng việt của 1st (<https://forum.machinelearningcoban.com/t/aivivn-3-vietnamese-tone-prediction-1st-place-solution/5721>), 2st (<https://forum.machinelearningcoban.com/t/aivivn-3-vietnamese-tone-prediction-2nd-place-solution/5759>). Có rất nhiều các ý tưởng hay được trình bày.

Trở lại với bài hướng dẫn này, đầu tiên mình sẽ hướng dẫn các bạn cách lấy dữ liệu. Nếu bạn đọc không quan tâm đến chuẩn bị dữ liệu có thể download `train_tiang_viet.txt` (<https://drive.google.com/file/d/1-7IERkqCoID1691yCXLAOyZoJqYPqhGq/view?usp=sharing>) và chuyển sang mục 2 xây dựng mô hình.

Top

Download dữ liệu viwikipedia

Đầu tiên các bạn cần tải file `viwiki-20200501-pages-articles.xml.bz2` tại wikipedia (<https://dumps.wikimedia.org/viwiki/20200501/>).

Chúng ta cũng có thể download file bằng lệnh `wget` bên dưới:

```
1 from google.colab import drive
2
3 drive.mount('/content/gdrive')
4 path = '/content/gdrive/My Drive/Colab Notebooks/BERT/themdau_tv'
5 %cd {path}
6 !ls

1 !wget https://dumps.wikimedia.org/viwiki/20200501/viwiki-20200501-pages-ar
2 !bzip2 -d viwiki-20200501-pages-articles.xml.bz2
3 !ls
```

Tiếp theo ta sử dụng `wikiextractor` để giải nén dữ liệu từ file `viwiki-20200501-pages-articles.xml.bz2` vừa mới download.

Quá trình giải nén mất khá nhiều thời gian. Bạn đọc có thể download dữ liệu có sẵn tại `viwiki-20200501-pages-articles-output` (<https://drive.google.com/drive/folders/11mkQBCUNuKxyLZEYfGe61INU0WY-SrR0?usp=sharing>). Để giải nén file chúng ta sẽ sử dụng package `wikiextractor` (<https://github.com/attardi/wikiextractor>). Một package chuyên dùng cho khai thác dữ liệu trên wiki.

```
1 !git clone https://github.com/attardi/wikiextractor.git

1 !python wikiextractor/WikiExtractor.py viwiki-20200501-pages-articles.xml.l
```

Các bạn xem thêm cách sử dụng lệnh giải nén bằng file `WikiExtractor.py` (<https://github.com/attardi/wikiextractor>) tại mục **Usage** của README.

Câu lệnh trên sẽ giải nén file và lưu vào folder `/output/` dưới định dạng json. Mỗi file sẽ bao gồm các dòng đại diện cho một văn bản tiếng Việt được giải nén từ wikipedia có định dạng như sau:

```
{"id": "", "revid": "", "url": "", "title": "", "text": "..."}

```

Trong đó:

- id: Mã của bài viết.
- revid: Mã của bài viết.
- url: Link url của bài viết.
- title: Tiêu đề bài viết.
- text: Nội dung của bài viết.

Chọn câu hợp lệ

Trong tập hợp các câu sẽ có một số câu không đạt tiêu chuẩn vì chứa các ký tự Tiếng Trung, Hàn,.... Vì vậy chúng ta sẽ lọc bỏ những câu này bằng cách tạo ra một hàm kiểm tra tính hợp lệ của câu. Một câu được coi là hợp lệ nếu chỉ chứa các ký tự tiếng việt được list trong `accept_strings` bên dưới.

Xây dựng hàm `_check_tiang_viet()` để kiểm tra tính hợp lệ của câu:

Top

[illegible]

1 True

Lưu file huấn luyện

Tiếp theo ta sẽ tạo vòng lặp đi qua toàn bộ các file trong `/output` folder. Kiểm tra câu có thỏa mãn tiêu chuẩn Tiếng Việt không, đánh index cho câu và lưu đồng thời index và câu có dấu vào file `train_tiang_viet.txt` . Thời gian trích xuất sẽ mất vài tiếng trên máy của mình. Vì vậy các bạn có thể download dữ liệu tại `train_tiang_viet.txt` (<https://drive.google.com/file/d/1-7IERkgCoID1691yCXLAOyZoJqYPghGg/view?usp=sharing>) và bỏ qua đoạn code bên dưới.

```

1  import pickle
2  import json
3  from tqdm import tqdm
4  import glob2
5  idx = 0
6
7  for path in tqdm(glob2.glob('output/*/*')):
8      # Đọc nội dung của các văn bản từ folder output. Content sẽ chứa nhiều
9      with open(path, 'r', encoding='utf8') as f:
10         content = f.readlines()
11         for row in content:
12             # Convert row sang json
13             art_json = json.loads(row)
14             # Lấy nội dung văn bản
15             art_cont = art_json['text']
16             art_cont = re.sub("(\\s)+", r"\\1", art_cont)
17             # Chia văn bản thành các câu tại vị trí xuống dòng
18             art_seqs = art_cont.split("\\n")
19             # Lưu các dòng là tiếng việt vào file 'train_tiang_viet.txt'.
20             # Mỗi dòng có định dạng: index{10digits} sequence
21             for seq in art_seqs:
22                 if _check_tiang_viet(seq):
23                     idx_str = str(idx).zfill(10)
24                     with open('train_tiang_viet.txt', 'a') as f:
25                         f.writelines([idx_str+'\\t', seq+'\\n'])
26                     idx += 1

```

Chuẩn bị dữ liệu train/val/test

Sau khi đã lưu file huấn luyện `train_tiang_viet.txt`, chúng ta sẽ load dữ liệu, loại bỏ dấu ở từng câu để tạo dữ liệu input. Dữ liệu output sẽ chính là câu Tiếng Việt có dấu.

```

1  with open('train_tiang_viet.txt', 'r', encoding='utf-8') as f:
2      train_output = f.readlines()
3
4  print('Number of sequences: ', len(train_output))
5  print('First sequence: ', train_output[0])

```



```

1  Number of sequences:  3624432
2  First sequence:  0000000000    Trang Chính

```

Ta thấy tổng số câu của chúng ta lên tới 3.6 triệu câu. Đối với cấu hình máy của google colab thì đây có thể coi là một lượng dữ liệu siêu to khổng lồ. Do đó mình sẽ chỉ lọc ra 500 nghìn câu đầu tiên làm tập huấn luyện (train dataset), 50 nghìn câu tiếp theo làm tập thẩm định (validation dataset) và 50 nghìn câu tiếp theo làm tập kiểm tra (test dataset).

Tạo hàm remove dấu

Hàm `remove_tone_line()` sẽ giúp bạn thực hiện điều này. Đây là hàm tiện ích được lấy từ cuộc thi thêm dấu Tiếng Việt cho từ không dấu aivivn.

Top

[illegible]

1 'Di mot ngay dang hoc 1 sang khon'

Sau đó chúng ta sẽ phân chia tập train/val/test

Top

```

1  from tqdm import tqdm
2  train_idx_500k = []
3  train_opt_500k = []
4  train_ipt_500k = []
5  val_idx_50k = []
6  val_opt_50k = []
7  val_ipt_50k = []
8  test_idx_50k = []
9  test_opt_50k = []
10 test_ipt_50k = []
11
12 for i in tqdm(range(600000)):
13     [idx, origin_seq] = train_output[i].split('\t')
14     try:
15         non_acc_seq = remove_tone_line(origin_seq)
16     except:
17         print('error remove tone line at sequence {}'.format(i))
18         next
19     if i < 500000:
20         train_idx_500k.append(idx)
21         train_opt_500k.append(origin_seq)
22         train_ipt_500k.append(non_acc_seq)
23     elif i < 550000:
24         val_idx_50k.append(idx)
25         val_opt_50k.append(origin_seq)
26         val_ipt_50k.append(non_acc_seq)
27     else:
28         test_idx_50k.append(idx)
29         test_opt_50k.append(origin_seq)
30         test_ipt_50k.append(non_acc_seq)

```

```

1  100%|██████████| 600000/600000 [00:31<00:00, 19184.82it/s]

```

```

1  print(train_ipt_500k[10])
2  print(train_opt_500k[10])

```

```

1  Tiếng Việt là ngôn ngữ có nguồn gốc bản địa, xuất thân từ nền văn minh nông
2
3  Tiếng Việt là ngôn ngữ có nguồn gốc bản địa, xuất thân từ nền văn minh nông

```

Sau khi đã có dữ liệu huấn luyện, thẩm định và kiểm tra. Chúng ta nên lưu lại để tái sử dụng cho những lượt cải thiện mô hình sau. Nếu bạn không lưu lại dữ liệu, bạn sẽ không có căn cứ để đánh giá và so sánh giữa các mô hình. Để không tốn dung lượng thì mình chỉ lưu lại index. Từ index có thể truy xuất ra câu cần lấy. Mình sử dụng google drive nên tài nguyên chỉ có 15GB thôi. Đó là lý do tại sao mình đánh index để tiết kiệm tài nguyên.

Top

```

1  import pickle
2
3  def _save_pickle(filename, obj):
4      with open(filename, 'wb') as f:
5          pickle.dump(obj, f)
6
7  _save_pickle('train_tv_idx_500k.pkl', train_idx_500k)
8  _save_pickle('val_tv_idx_50k.pkl', val_idx_50k)
9  _save_pickle('test_tv_idx_50k.pkl', test_idx_50k)

```

2. Xây dựng model Transformer

Hướng dẫn này sẽ huấn luyện model transformer để tiến hành dịch câu Tiếng Việt không dấu sang có dấu. Mã nguồn tham khảo tại Transformer model for language understanding (<https://www.tensorflow.org/tutorials/text/transformer>) ứng dụng trên thuật toán dịch máy. Có rất nhiều mã nguồn khác nhau về transformer. Nhưng mình chọn mã nguồn này là bởi tác giả giải thích các step rất chi tiết, dễ hiểu ở từng bước xử lý. Bạn đọc sẽ hiểu được transformer sẽ được tiến hành ra sao, các kiến trúc cụ thể như thế nào thông qua tutorial này. Trước khi bắt tay vào xây dựng mô hình, tôi khuyên bạn đọc qua trước Bài 4 - Attention is all you need (<https://phamdinhhkhanh.github.io/2019/06/18/AttentionLayer.html>) để hiểu về kiến trúc transformer và bài Bài 36 - BERT model (<https://phamdinhhkhanh.github.io/2020/05/23/BERTModel.html>) trình bày về mô hình BERT áp dụng kiến trúc transformer trong các tác vụ NLP.

Điểm cải tiến của transformer so với RNN đó là kỹ thuật attention giúp cho học được từ toàn bộ các vị trí khác nhau trong câu input để tính toán ra biểu diễn output của câu. transformer tạo ra một chuỗi stack các block sub-layer và áp dụng các biến đổi Scale dot Product attention và Multi-head attention sẽ được giải thích ở bên dưới.

Một mô hình transformer sẽ kiểm soát kích thước biến của input sử dụng stacks của các self-attention layers thay vì RNNs, CNNs. Kiến trúc tổng quát này có những điểm lợi thế sau:

- Không phụ thuộc vào giả thiết quan hệ thời gian trong toàn bộ dữ liệu.
- Output có thể được tính toán song song thay vì theo chuỗi như RNN.
- Các từ ở xa có thể ảnh hưởng tới những output của những từ khác mà không truyền qua nhiều RNN steps.
- Nó có thể học được sự phụ thuộc dài hạn. Sự phụ thuộc dài hạn là một thách thức của rất nhiều các tác vụ seq2seq.

Hạn chế của transformer:

- Output của tranformer được tính toán từ toàn bộ lịch sử thay vì chỉ input và hidden-state hiện tại. Đây có thể là một nhược điểm vì khoảng cách xa có thể không liên quan.
- Input không có sự phụ thuộc thời gian, do đó position encoding cần được thêm vào để mã hóa sự tương quan về mặt thời gian.

```

1  import tensorflow_datasets as tfds
2  import tensorflow as tf
3
4  import time
5  import numpy as np
6  import matplotlib.pyplot as plt

```

Top

```

1 # train_examples = tf.data.Dataset.from_tensor_slices((train_ipt_100k, tra
2 train_examples = tf.data.Dataset.from_tensor_slices((train_ipt_500k, train
3 val_examples = tf.data.Dataset.from_tensor_slices((val_ipt_50k, val_opt_50k))

```

Khởi tạo tokenize cho tập train input và output

```

1 tokenizer_ipt = tfds.features.text.SubwordTextEncoder.build_from_corpus(
2     (ipt.numpy() for (ipt, opt) in train_examples), target_vocab_size=2**11:
3
4 tokenizer_opt = tfds.features.text.SubwordTextEncoder.build_from_corpus(
5     (opt.numpy() for (ipt, opt) in train_examples), target_vocab_size=2**11:

```

```

1 sample_string = 'Tiếng Việt là ngôn ngữ trong sáng nhất thế giới'
2
3 tokenized_string = tokenizer_opt.encode(sample_string)
4 print ('Tokenized string is {}'.format(tokenized_string))
5
6 original_string = tokenizer_opt.decode(tokenized_string)
7 print ('The original string: {}'.format(original_string))
8
9 assert original_string == sample_string

```

```

1 Tokenized string is [2270, 65, 5, 695, 527, 10, 451, 60, 56, 573]
2 The original string: Tiếng Việt là ngôn ngữ trong sáng nhất thế giới

```

Lưu lại Tokenizer và kích thước vocabulary size của nó:

```

1 import pickle
2
3 def _save_pickle(path, obj):
4     with open(path, 'wb') as f:
5         pickle.dump(obj, f)
6
7 def _load_pickle(path):
8     with open(path, 'rb') as f:
9         obj = pickle.load(f)
10    return obj
11
12 _save_pickle('tokenizer/tokenizer_ipt.pkl', tokenizer_ipt)
13 _save_pickle('tokenizer/tokenizer_opt.pkl', tokenizer_opt)

```

Tokenizer sẽ encoding chuỗi string bằng cách chi nhỏ nó thành những subwords nếu từ không xuất hiện trong từ điển của nó.

Khái niệm subwords: Thuật toán transformer sẽ hiệu quả hơn nếu ta chia các từ theo subwords. subword là một chuỗi các ký tự xuất hiện trong token mà thường được lặp đi lặp lại. Chẳng hạn như từ chính tả thì cụm ký tự như ính sẽ dễ dàng được sử dụng lặp lại ở những từ khác. Do đó nó có thể là một subword. Bạn đọc có thể tham khảo thêm về subwords tại [How subword helps on your nlp model \(https://medium.com/@makcedward/how-subword-helps-on-your-nlp-model-83dd1b836f46\)](https://medium.com/@makcedward/how-subword-helps-on-your-nlp-model-83dd1b836f46)

```

1 for ts in tokenized_string:
2     print ('{} ----> {}'.format(ts, tokenizer_opt.decode([ts])))

```

Top


```

1      2270 ----> Tiếng
2      65 ----> Việt
3      5 ----> là
4      695 ----> ngôn
5      527 ----> ngữ
6      10 ----> trong
7      451 ----> sáng
8      60 ----> nhất
9      56 ----> thể
10     573 ----> giới

```

```

1      BUFFER_SIZE = 20000
2      BATCH_SIZE = 64

```

Tiếp theo chúng ta sẽ thêm token `start` và `end` vào câu input và câu target

```

1      def encode(ipt, opt):
2          ipt = [tokenizer_ipt.vocab_size] + tokenizer_ipt.encode(
3              ipt.numpy()) + [tokenizer_ipt.vocab_size+1]
4
5          opt = [tokenizer_opt.vocab_size] + tokenizer_opt.encode(
6              opt.numpy()) + [tokenizer_opt.vocab_size+1]
7
8          return ipt, opt

```

```

1      # encode(ipt = 'tieng viet la ngon ngu trong sang', opt = 'tiếng việt là n

```

Nếu bạn muốn sử dụng `Dataset.map` để áp dụng hàm số này cho mỗi phần tử của dataset. `Dataset.map` sẽ chạy trên graph mode.

- Graph tensors không có dữ liệu.
- Trong graph model bạn chỉ có thể sử dụng Tensorflow Ops và functions.

Do đó bạn không thể `.map` các hàm số này trực tiếp: Bạn cần wrap nó trong một hàm số gọi là `tf.py_function`. Hàm số này sẽ truyền các tensors thông thường (với một giá trị và một phương thức `.numpy()` để truy cập nó), để wrapped hàm số trong python.

```

1      def tf_encode(ipt, opt):
2          result_ipt, result_opt = tf.py_function(encode, [ipt, opt], [tf.int64, t
3          result_ipt.set_shape([None])
4          result_opt.set_shape([None])
5          return result_ipt, result_opt

```

Note: Để đỡ cho bộ dữ liệu nhỏ và huấn luyện nhanh hơn, chúng ta sẽ loại bỏ những mẫu kích thước lớn hơn 40 tokens.

Khởi tạo tensorflow Dataset cho train và validation.

```

1      next(iter(train_examples))

```

```

1      (<tf.Tensor: shape=(), dtype=string, numpy=b'Trang Chinh\n'>,
2      <tf.Tensor: shape=(), dtype=string, numpy=b'Trang Ch\xcc3\xadnh\n'>),

```

Top

```

1     MAX_LENGTH = 40
2
3     def filter_max_length(x, y, max_length=MAX_LENGTH):
4         return tf.logical_and(tf.size(x) <= max_length,
5                                tf.size(y) <= max_length)
6
7     train_dataset = train_examples.map(tf_encode)
8     train_dataset = train_dataset.filter(filter_max_length)
9     # cache the dataset to memory to get a speedup while reading from it.
10    train_dataset = train_dataset.cache()
11    train_dataset = train_dataset.shuffle(BUFFER_SIZE).padded_batch(BATCH_SIZE)
12    train_dataset = train_dataset.prefetch(tf.data.experimental.AUTOTUNE)
13
14    val_dataset = val_examples.map(tf_encode)
15    val_dataset = val_dataset.filter(filter_max_length).padded_batch(BATCH_SIZE)

```

```

1     print('tokenizer_ipt.vocab_size: ', tokenizer_ipt.vocab_size)
2     x, y = next(iter(train_dataset))
3     print('input shape: ', x.shape)
4     print('output shape: ', y.shape)
5     print('first row of input index: ', x[0, :])

```

```

1     tokenizer_ipt.vocab_size:  8185
2     input shape:  (64, 37)
3     output shape:  (64, 40)
4     first row of input index:  tf.Tensor(
5     [8185 7759 7961 1302 7939 8186    0    0    0    0    0    0    0    0
6      0    0    0    0    0    0    0    0    0    0    0    0    0    0
7      0    0    0    0    0    0    0    0    0], shape=(37,), dtype=int64)

```

3. Position Encoding

Chức năng của position Encoding: Tranformer sẽ không dự báo tuân theo time step như RNN. Toàn bộ các input sẽ được truyền vào mô hình **cùng một thời điểm**. Sẽ rất khó để nhận biết vị trí của các từ input trong câu ở những lớp mô hình không tuân theo time step. Do đó thêm position encoding sẽ cho mô hình thêm các thông tin về vị trí của từ.

Position encoding véc tơ sẽ được cộng trực tiếp vào embedding véc tơ. Embeddings biểu diễn một token trong một không gian d chiều nơi mà các token có cùng ý nghĩa sẽ gần nhau hơn. Nhưng embedding vector không chứa thông tin vị trí của từ trong câu. Do đó sau khi thêm position encoding véc tơ, một từ sẽ gần với những từ khác hơn dựa trên ý nghĩa của chúng và khoảng cách vị trí của chúng trong câu trong không gian d chiều.

Theo dõi note book về position encoding

(https://github.com/tensorflow/examples/blob/master/community/en/position_encoding.ipynb) để hiểu rõ hơn về chức năng và tính chất của nó.

Công thức position Encoding:

$$PE(pos, 2i) = \sin\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right)$$

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

Top

Trong đó pos là vị trí hiện tại của từ, i chỉ số của phần tử nằm trong véc tơ encoding và d_{model} là kích thước các véc tơ positional embedding.

Giả sử $PE(pos)$ là véc tơ encoding tại vị trí pos . Véc tơ này có kích thước phải bằng với các véc tơ embedding từ để phép cộng thực hiện được và kích thước đó bằng d_{model} .

Công thức $PE(pos, 2i)$, $PE(pos, 2i + 1)$ tính giá trị các phần tử véc tơ positional encoding tại lần lượt vị trí $2i$ và $2i + 1$.

```

1 def get_angles(pos, i, d_model):
2     angle_rates = 1 / np.power(10000, (2 * (i//2)) / np.float32(d_model))
3     return pos * angle_rates

1 def positional_encoding(position, d_model):
2     angle_rads = get_angles(np.arange(position)[:, np.newaxis],
3                             np.arange(d_model)[np.newaxis, :],
4                             d_model) # shape (position, d_model)
5
6     # apply sin to even indices in the array; 2i
7     angle_rads[:, 0::2] = np.sin(angle_rads[:, 0::2])
8
9     # apply cos to odd indices in the array; 2i+1
10    angle_rads[:, 1::2] = np.cos(angle_rads[:, 1::2])
11
12    pos_encoding = angle_rads[np.newaxis, ...]
13
14    return tf.cast(pos_encoding, dtype=tf.float32) # shape: (position, d_model)

```

Hàm `positional_encoding()` sẽ tạo ra ma trận mà gồm các véc tơ dòng là những positional encoding véc tơ.

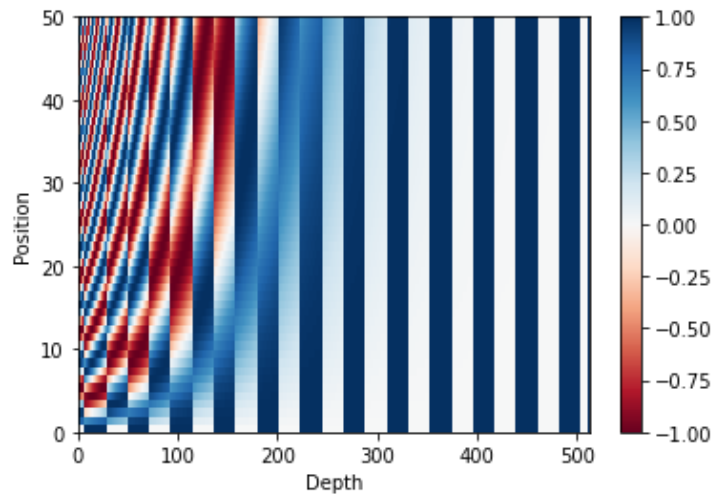
```

1 pos_encoding = positional_encoding(50, 512)
2 print (pos_encoding.shape)
3
4 plt.pcolormesh(pos_encoding[0], cmap='RdBu')
5 plt.xlabel('Depth')
6 plt.xlim((0, 512))
7 plt.ylabel('Position')
8 plt.colorbar()
9 plt.show()

1 (1, 50, 512)

```

Top



4. Masking

Masking nhằm mục đích không đưa các vị trí padding trong câu vào như một input. Mask là ma trận có kích thước bằng với kích thước ma trận input và đánh dấu các vị trí padding (tương ứng với 0) bằng giá trị 1. Các giá trị còn lại bằng 0.

```
1 def create_padding_mask(seq):
2     seq = tf.cast(tf.math.equal(seq, 0), tf.float32)
3
4     # add extra dimensions to add the padding
5     # to the attention logits.
6     return seq[:, tf.newaxis, tf.newaxis, :] # (batch_size, 1, 1, seq_len)
```

```
1 x = tf.constant([[7, 6, 0, 0, 1], [1, 2, 3, 0, 0], [0, 0, 0, 4, 5]])
2 create_padding_mask(x)
```

```
1 <tf.Tensor: shape=(3, 1, 1, 5), dtype=float32, numpy=
2 array([[[[0., 0., 1., 1., 0.]],
3
4
5         [[0., 0., 0., 1., 1.]],
6
7         [[1., 1., 1., 0., 0.]]]], dtype=float32)>
```

Tiếp theo để ngăn cản ảnh hưởng của các từ tương lai vào dự đoán từ hiện tại. Chúng ta sẽ tiếp tục mask các vị trí tương lai bằng 1.

Như vậy tại vị trí cần dự báo từ thứ 3 thì chỉ từ thứ 1 và thứ 2 được đưa vào dự báo.

```
1 def create_look_ahead_mask(size):
2     mask = 1 - tf.linalg.band_part(tf.ones((size, size)), -1, 0)
3     return mask # (seq_len, seq_len)
```

```
1 x = tf.random.uniform((1, 3))
2 temp = create_look_ahead_mask(x.shape[1])
3 temp
```

Top

```

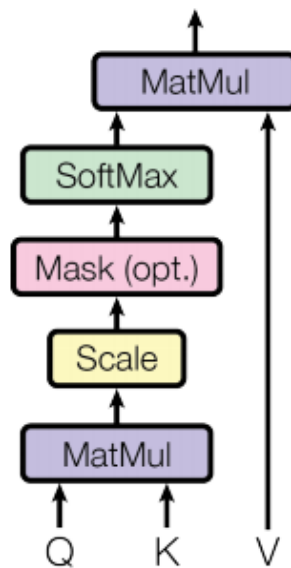
1 <tf.Tensor: shape=(3, 3), dtype=float32, numpy=
2 array([[0., 1., 1.],
3        [0., 0., 1.],
4        [0., 0., 0.]], dtype=float32)>

```

5. Scale dot Product attention

Một attention sẽ sử dụng input là 3 ma trận: **Q** (query), **K** (key) và **V** (value). Theo sơ đồ như sau:

Scaled Dot-Product Attention



Output thu được là ma trận attention của các từ lẫn nhau trong câu.

Diễn giải quá trình tính attention như sau:

- Sau khi thực hiện **MatMul** giữa **Q** và **K** → thu được ma trận attention.
- Scale ma trận attention với nghịch đảo số độ sâu (chính là kích thước véc tơ dòng) của ma trận **K** là $\frac{1}{\sqrt{d_k}}$ để tránh cho giá trị phần tử của ma trận attention quá lớn trong khi gradient thì quá nhỏ → huấn luyện lâu.
- Thực hiện Masking. Ma trận mask là ma trận để ngăn cho những vị trí padding tham gia vào quá trình attention. Ma trận masking được nhân với -1E9 (là một giá trị gần với âm vô cùng). Sở dĩ ta thực hiện như vậy vì sau đó ma trận attention sẽ cộng với ma trận mask. Tại những vị trí padding sẽ có giá trị gần như âm vô cùng và khi tính phân phối softmax theo dòng sẽ thu được output là 0. Một lát nữa chúng ta sẽ thực nghiệm điều này.
- Sau khi tính phân phối softmax cho ma trận attention, chúng ta sẽ nhân ma trận attention với ma trận **V**.

Toàn bộ quá trình phức tạp trên được tổng hợp trong công thức Attention:

$$Attention(Q, K, V) = softmax_k\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad \text{Top}$$

```

1 def scaled_dot_product_attention(q, k, v, mask):
2     """Calculate the attention weights.
3     q, k, v must have matching leading dimensions.
4     k, v must have matching penultimate dimension, i.e.: seq_len_k = seq_len_v
5     The mask has different shapes depending on its type(padding or look ahead)
6     but it must be broadcastable for addition.
7
8     Args:
9         q: query shape == (... , seq_len_q, depth)
10        k: key shape == (... , seq_len_k, depth)
11        v: value shape == (... , seq_len_v, depth_v)
12        mask: Float tensor with shape broadcastable
13              to (... , seq_len_q, seq_len_k). Defaults to None.
14
15    Returns:
16        output, attention_weights
17    """
18
19    matmul_qk = tf.matmul(q, k, transpose_b=True) # (... , seq_len_q, seq_len_k)
20
21    # scale matmul_qk
22    dk = tf.cast(tf.shape(k)[-1], tf.float32) # depth
23    scaled_attention_logits = matmul_qk / tf.math.sqrt(dk)
24
25    # add the mask to the scaled tensor.
26    if mask is not None:
27        scaled_attention_logits += (mask * -1e9)
28
29    # softmax is normalized on the last axis (seq_len_k) so that the scores
30    # add up to 1.
31    attention_weights = tf.nn.softmax(scaled_attention_logits, axis=-1) #
32
33    output = tf.matmul(attention_weights, v) # (... , seq_len_q, depth_v)
34
35    return output, attention_weights

```

Mask sẽ ngăn cản các vị trí được padding. Bên dưới ta sẽ thực nghiệm quá trình masking và xem hệ số attention sau khi đi qua phân phối softmax.

```

1 # with tf.Session() as sess:
2 mask = tf.constant([[0, 1, 1],
3                     [0, 0, 1],
4                     [0, 0, 0]], dtype = tf.float64)
5
6 scaled_attention_logit = tf.constant([[1, 3, 10],
7                                       [1, 2, 5],
8                                       [1, 1, 5]], dtype = tf.float64)
9
10 scaled_attention_logit += (mask * -1e9)
11 attention_weights = tf.nn.softmax(scaled_attention_logit, axis=-1)
12 print('scaled_attention_logit: ', scaled_attention_logit)
13 print('attention_weights: ', attention_weights)

```

Top

```

1 scaled_attention_logits: tf.Tensor(
2 [[ 1.00000000e+00 -9.99999997e+08 -9.99999990e+08]
3  [ 1.00000000e+00  2.00000000e+00 -9.99999995e+08]
4  [ 1.00000000e+00  1.00000000e+00  5.00000000e+00]], shape=(3, 3), dtype=f.
5 attention_weights: tf.Tensor(
6 [[1.         0.         0.         ]
7  [0.26894142 0.73105858 0.         ]
8  [0.01766842 0.01766842 0.96466316]], shape=(3, 3), dtype=float64)

```

Ta có thể thấy tại các vị trí padding được mask với giá trị bằng 1 thì attention_weights bằng 0.

Kết quả sau cùng sẽ thu được là tích giữa ma trận attention_weights với ma trận V . Kết quả này sẽ đảm bảo:

- Phân bố attention giữa các véc tơ trong V
- Loại bỏ những từ không liên quan ra khỏi attention.

Tiếp theo ta sẽ test hàm `scaled_dot_product_attention()`.

```

1 def print_out(q, k, v):
2     temp_out, temp_attn = scaled_dot_product_attention(
3         q, k, v, None)
4     print ('Attention weights are:')
5     print (temp_attn)
6     print ('Output is:')
7     print (temp_out)

1 np.set_printoptions(suppress=True)
2
3 temp_k = tf.constant([[10,0,0],
4                       [0,10,0],
5                       [0,0,10],
6                       [0,0,10]], dtype=tf.float32) # (4, 3)
7
8 temp_v = tf.constant([[ 1,0],
9                       [ 10,0],
10                      [ 100,5],
11                      [1000,6]], dtype=tf.float32) # (4, 2)
12
13 # This `query` aligns with the second `key`,
14 # so the second `value` is returned.
15 temp_q = tf.constant([[0, 10, 0]], dtype=tf.float32) # (1, 3)
16 print_out(temp_q, temp_k, temp_v)

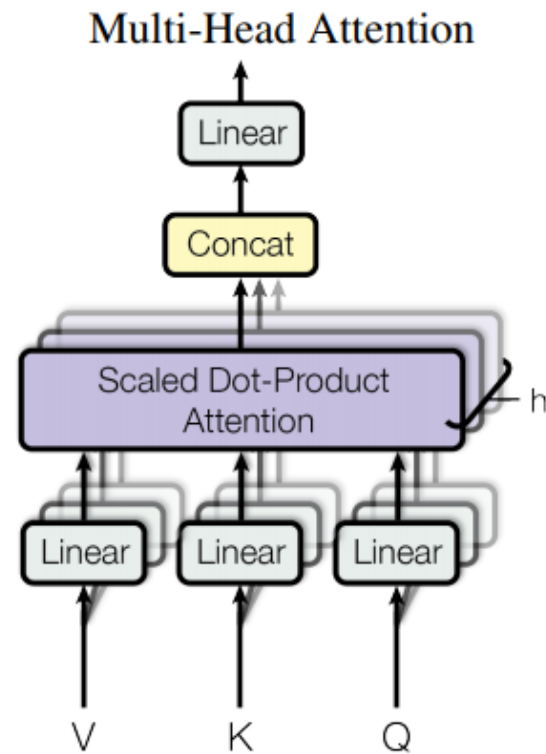
1 Attention weights are:
2 tf.Tensor([[0. 1. 0. 0.]], shape=(1, 4), dtype=float32)
3 Output is:
4 tf.Tensor([[10. 0.]], shape=(1, 2), dtype=float32)

```

Véc tơ attention là $[0, 1, 0, 0]$ chỉ tập trung vào vị trí thứ 2. Do đó kết quả trả về là dòng thứ 2 của ma trận V .

6. Multi-head Attention

Top



Multi-head attention sẽ bao gồm 4 phần:

- Các linear layers và phân chi thành các nhiều heads.
- Scaled dot-product attention
- Concatenate các heads
- Linear layer cuối cùng.

Mỗi một multi-head attention block nhận 3 đầu vào: là các ma trận Q , K , V . Sau đó `scale_dot_product_attention()` sẽ được sử dụng để tính toán trên từng head. Attention output của mỗi layer sau đó được concatenate và truyền qua một Dens layer.

Thay vì sử dụng một single head attention, chúng ta sử dụng nhiều multi-head attention là để mô hình có thể học được attention từ các vị trí trên những biểu diễn không gian khác nhau.

Top


```

1 class MultiHeadAttention(tf.keras.layers.Layer):
2     def __init__(self, d_model, num_heads):
3         super(MultiHeadAttention, self).__init__()
4         self.num_heads = num_heads
5         self.d_model = d_model
6
7         assert d_model % self.num_heads == 0
8
9         self.depth = d_model // self.num_heads
10
11        self.wq = tf.keras.layers.Dense(d_model)
12        self.wk = tf.keras.layers.Dense(d_model)
13        self.wv = tf.keras.layers.Dense(d_model)
14
15        self.dense = tf.keras.layers.Dense(d_model)
16
17    def split_heads(self, x, batch_size):
18        """Split the last dimension into (num_heads, depth).
19        Transpose the result such that the shape is (batch_size, num_heads, seq_len, depth)
20        """
21        x = tf.reshape(x, (batch_size, -1, self.num_heads, self.depth))
22        return tf.transpose(x, perm=[0, 2, 1, 3])
23
24    def call(self, v, k, q, mask):
25        batch_size = tf.shape(q)[0]
26
27        q = self.wq(q) # (batch_size, seq_len, d_model)
28        k = self.wk(k) # (batch_size, seq_len, d_model)
29        v = self.wv(v) # (batch_size, seq_len, d_model)
30
31        q = self.split_heads(q, batch_size) # (batch_size, num_heads, seq_len, depth)
32        k = self.split_heads(k, batch_size) # (batch_size, num_heads, seq_len, depth)
33        v = self.split_heads(v, batch_size) # (batch_size, num_heads, seq_len, depth)
34
35        # scaled_attention.shape == (batch_size, num_heads, seq_len_q, depth)
36        # attention_weights.shape == (batch_size, num_heads, seq_len_q, seq_len_k)
37        scaled_attention, attention_weights = scaled_dot_product_attention(
38            q, k, v, mask)
39
40        scaled_attention = tf.transpose(scaled_attention, perm=[0, 2, 1, 3])
41
42        concat_attention = tf.reshape(scaled_attention,
43                                     (batch_size, -1, self.d_model)) # (batch_size, seq_len_q, d_model)
44
45        output = self.dense(concat_attention) # (batch_size, seq_len_q, d_model)
46
47        return output, attention_weights

```

```

1 temp_mha = MultiHeadAttention(d_model=512, num_heads=8)
2 y = tf.random.uniform((1, 60, 512)) # (batch_size, encoder_sequence, d_model)
3 out, attn = temp_mha(y, k=y, q=y, mask=None)
4 out.shape, attn.shape

```

```

1 (TensorShape([1, 60, 512]), TensorShape([1, 8, 60, 60]))

```

Top

Như vậy shape của attention sẽ là $(batch_size, num_head, seq_len_q, seq_len_k)$ và shape của output sẽ là $(batch_size, seq_len_q, d_model)$.

7. Point wise feed forward network

Point wise feed forward bao gồm 2 fully-connected layers với hàm ReLU activation.

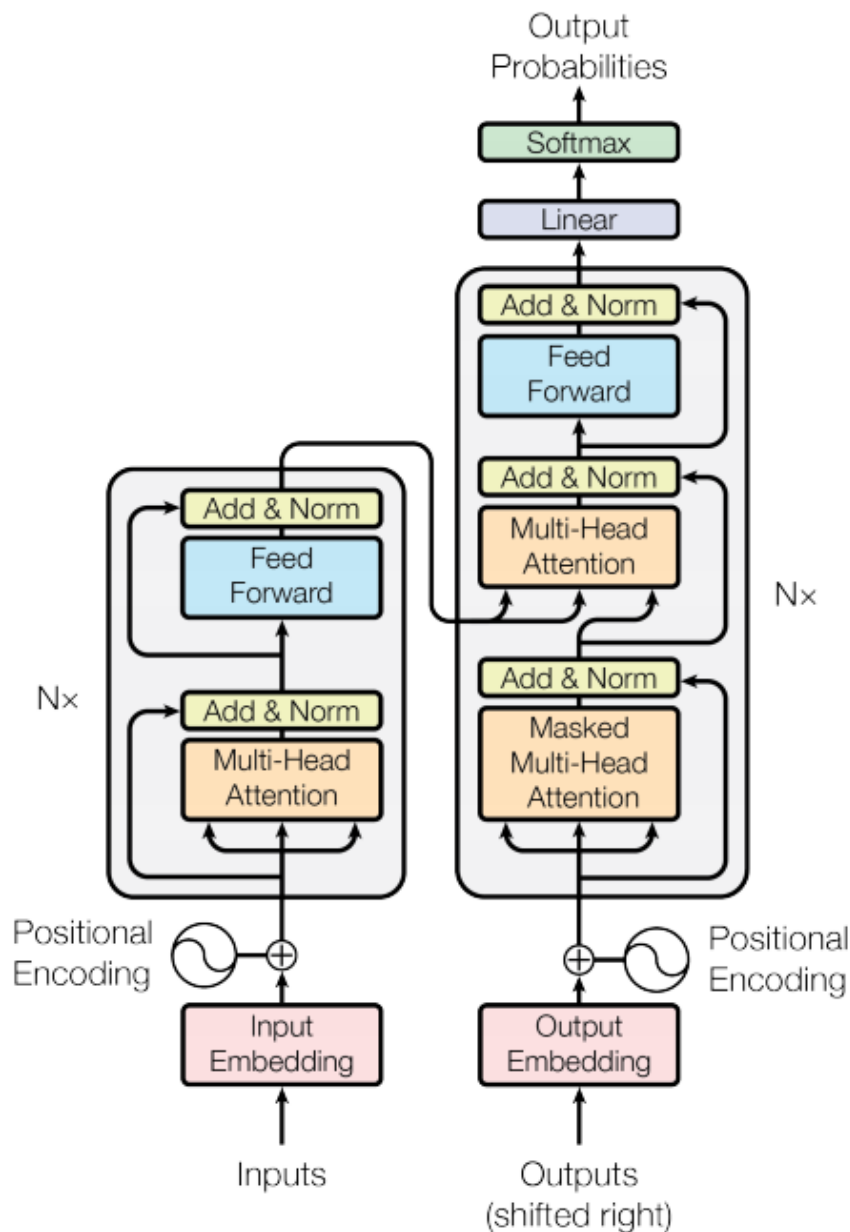
```

1 def point_wise_feed_forward_network(d_model, dff):
2     return tf.keras.Sequential([
3         tf.keras.layers.Dense(dff, activation='relu'), # (batch_size, seq_len, d_model)
4         tf.keras.layers.Dense(d_model) # (batch_size, seq_len, d_model)
5     ])

```

8. Encoder and Decoder

Toàn bộ kiến trúc của transformer sẽ bao gồm 2 nhánh Encoder và Decoder như bên dưới:



Top

- Câu input sẽ được embedding và truyền qua N sub-layers block của encoder để sinh ra output cho mỗi từ trong câu.
- Decoder áp dụng encoder-decoder attention giữa output của encoder và input của decoder để dự báo next word.

8.1. Encoder layer

Mỗi một Encoder layer sẽ gồm 2 sublayers:

- Multi-head Attention (với padding mask)
- Point wise feed forward network.

```

1      class EncoderLayer(tf.keras.layers.Layer):
2          def __init__(self, d_model, num_heads, dff, rate=0.1):
3              super(EncoderLayer, self).__init__()
4
5              self.mha = MultiHeadAttention(d_model, num_heads)
6              self.ffn = point_wise_feed_forward_network(d_model, dff)
7
8              self.layernorm1 = tf.keras.layers.LayerNormalization(epsilon=1e-6)
9              self.layernorm2 = tf.keras.layers.LayerNormalization(epsilon=1e-6)
10
11             self.dropout1 = tf.keras.layers.Dropout(rate)
12             self.dropout2 = tf.keras.layers.Dropout(rate)
13
14         def call(self, x, training, mask):
15
16             attn_output, _ = self.mha(x, x, x, mask) # (batch_size, input_seq_len, d_model)
17             attn_output = self.dropout1(attn_output, training=training)
18             out1 = self.layernorm1(x + attn_output) # (batch_size, input_seq_len, d_model)
19
20             ffn_output = self.ffn(out1) # (batch_size, input_seq_len, d_model)
21             ffn_output = self.dropout2(ffn_output, training=training)
22             out2 = self.layernorm2(out1 + ffn_output) # (batch_size, input_seq_len, d_model)
23
24             return out2

```

```

1      sample_encoder_layer = EncoderLayer(512, 8, 2048)
2
3      sample_encoder_layer_output = sample_encoder_layer(
4          tf.random.uniform((64, 43, 512)), False, None)
5
6      sample_encoder_layer_output.shape # (batch_size, input_seq_len, d_model)

```

```

1      TensorShape([64, 43, 512])

```

Nếu bạn chạy ra shape của output Encoder là (batch_size, input_seq_length, d_model) là đúng.

8.2. Decoder layer

Mỗi một layer của Decoder layer sẽ bao gồm 3 sublayers:

- Masked multi-head attention (với look ahead mask và padding mask).

Top

- Multi-head attention (với padding mask). Ma trận V , K cùng lấy output từ Encoder và ma trận Q nhận output từ *masked multi-head attention*.
- Point wise feed forward network

```

1      class DecoderLayer(tf.keras.layers.Layer):
2          def __init__(self, d_model, num_heads, dff, rate=0.1):
3              super(DecoderLayer, self).__init__()
4
5              self.mha1 = MultiHeadAttention(d_model, num_heads)
6              self.mha2 = MultiHeadAttention(d_model, num_heads)
7
8              self.ffn = point_wise_feed_forward_network(d_model, dff)
9
10             self.layernorm1 = tf.keras.layers.LayerNormalization(epsilon=1e-6)
11             self.layernorm2 = tf.keras.layers.LayerNormalization(epsilon=1e-6)
12             self.layernorm3 = tf.keras.layers.LayerNormalization(epsilon=1e-6)
13
14             self.dropout1 = tf.keras.layers.Dropout(rate)
15             self.dropout2 = tf.keras.layers.Dropout(rate)
16             self.dropout3 = tf.keras.layers.Dropout(rate)
17
18         def call(self, x, enc_output, training,
19                 look_ahead_mask, padding_mask):
20             # enc_output.shape == (batch_size, input_seq_len, d_model)
21
22             attn1, attn_weights_block1 = self.mha1(x, x, x, look_ahead_mask) # (batch_size, target_seq_len, d_model)
23             attn1 = self.dropout1(attn1, training=training)
24             out1 = self.layernorm1(attn1 + x)
25
26             attn2, attn_weights_block2 = self.mha2(
27                 enc_output, enc_output, out1, padding_mask) # (batch_size, target_seq_len, d_model)
28             attn2 = self.dropout2(attn2, training=training)
29             out2 = self.layernorm2(attn2 + out1) # (batch_size, target_seq_len, d_model)
30
31             ffn_output = self.ffn(out2) # (batch_size, target_seq_len, d_model)
32             ffn_output = self.dropout3(ffn_output, training=training)
33             out3 = self.layernorm3(ffn_output + out2) # (batch_size, target_seq_len, d_model)
34
35             return out3, attn_weights_block1, attn_weights_block2

```

```

1      sample_decoder_layer = DecoderLayer(512, 8, 2048)
2
3      sample_decoder_layer_output, _, _ = sample_decoder_layer(
4          tf.random.uniform((64, 50, 512)), sample_encoder_layer_output,
5          False, None, None)
6
7      sample_decoder_layer_output.shape # (batch_size, target_seq_len, d_model)

```

```

1      TensorShape([64, 50, 512])

```

Output shape của Decoder sẽ là (batch_size, target_seq_len, d_model).

8.3. Encoder

Top

Toàn bộ Encoder sẽ bao gồm:

- Input Embedding
- Positional Encoding
- N encoder layers

Input sau khi được Embedding sẽ được cộng với Positional Encoding. Kết quả tổng thu được tiếp tục được truyền qua N encoder layers.

```

1      class Encoder(tf.keras.layers.Layer):
2          def __init__(self, num_layers, d_model, num_heads, dff, input_vocab_size,
3                      maximum_position_encoding, rate=0.1):
4              super(Encoder, self).__init__()
5
6              self.d_model = d_model
7              self.num_layers = num_layers
8
9              self.embedding = tf.keras.layers.Embedding(input_vocab_size, d_model)
10             self.pos_encoding = positional_encoding(maximum_position_encoding,
11                                                     self.d_model)
12
13
14             self.enc_layers = [EncoderLayer(d_model, num_heads, dff, rate)
15                                 for _ in range(num_layers)]
16
17             self.dropout = tf.keras.layers.Dropout(rate)
18
19         def call(self, x, training, mask):
20
21             seq_len = tf.shape(x)[1]
22
23             # adding embedding and position encoding.
24             x = self.embedding(x) # (batch_size, input_seq_len, d_model)
25             x *= tf.math.sqrt(tf.cast(self.d_model, tf.float32))
26             x += self.pos_encoding[:, :seq_len, :]
27
28             x = self.dropout(x, training=training)
29
30             for i in range(self.num_layers):
31                 x = self.enc_layers[i](x, training, mask)
32
33             return x # (batch_size, input_seq_len, d_model)

```

```

1      sample_encoder = Encoder(num_layers=2, d_model=512, num_heads=8,
2                              dff=2048, input_vocab_size=8500,
3                              maximum_position_encoding=10000)
4
5      # Init sample tensorflow with shape 64 x 62 and data type int.
6      temp_input = tf.random.uniform((64, 62), dtype=tf.int64, minval=0, maxval=
7
8      sample_encoder_output = sample_encoder(temp_input, training=False, mask=None)
9
10     print (sample_encoder_output.shape) # (batch_size, input_seq_len, d_model)

```

1 (64, 62, 512)

Top

8.4. Decoder

Decoder sẽ bao gồm:

- Output Embedding
- Positional Embedding
- N encoder layers

Giá trị target được truyền qua Output Embedding và cộng với Positional Encoding. Tổng sau đó tiếp tục được truyền qua N encoder layers (Đã tính attention với output của Encoder). Output của decoder là input của linear layer cuối cùng.

```

1      class Decoder(tf.keras.layers.Layer):
2          def __init__(self, num_layers, d_model, num_heads, dff, target_vocab_size,
3                      maximum_position_encoding, rate=0.1):
4              super(Decoder, self).__init__()
5
6              self.d_model = d_model
7              self.num_layers = num_layers
8
9              self.embedding = tf.keras.layers.Embedding(target_vocab_size, d_model)
10             self.pos_encoding = positional_encoding(maximum_position_encoding, d_model)
11
12             self.dec_layers = [DecoderLayer(d_model, num_heads, dff, rate)
13                               for _ in range(num_layers)]
14             self.dropout = tf.keras.layers.Dropout(rate)
15
16         def call(self, x, enc_output, training,
17                 look_ahead_mask, padding_mask):
18
19             seq_len = tf.shape(x)[1]
20             attention_weights = {}
21
22             x = self.embedding(x) # (batch_size, target_seq_len, d_model)
23             x *= tf.math.sqrt(tf.cast(self.d_model, tf.float32))
24             x += self.pos_encoding[:, :seq_len, :]
25
26             x = self.dropout(x, training=training)
27
28             for i in range(self.num_layers):
29                 x, block1, block2 = self.dec_layers[i](x, enc_output, training,
30                                                         look_ahead_mask, padding_mask)
31
32                 attention_weights['decoder_layer{}_block1'.format(i+1)] = block1
33                 attention_weights['decoder_layer{}_block2'.format(i+1)] = block2
34
35             # x.shape == (batch_size, target_seq_len, d_model)
36             return x, attention_weights

```

Top

```

1 sample_decoder = Decoder(num_layers=2, d_model=512, num_heads=8,
2                           dff=2048, target_vocab_size=8000,
3                           maximum_position_encoding=5000)
4 temp_input = tf.random.uniform((64, 26), dtype=tf.int64, minval=0, maxval=
5
6 output, attn = sample_decoder(temp_input,
7                               enc_output=sample_encoder_output,
8                               training=False,
9                               look_ahead_mask=None,
10                              padding_mask=None)
11
12 output.shape, attn['decoder_layer2_block2'].shape

```

```

1 (TensorShape([64, 26, 512]), TensorShape([64, 8, 26, 62]))

```

9. Create transformer

Mô hình transformer sẽ hình thành từ Encoder, Decoder và linear layer cuối cùng. Kết quả sau cùng là output của linear layer cuối cùng.

```

1 class Transformer(tf.keras.Model):
2     def __init__(self, num_layers, d_model, num_heads, dff, input_vocab_size,
3                 target_vocab_size, pe_input, pe_target, rate=0.1):
4         super(Transformer, self).__init__()
5
6         self.encoder = Encoder(num_layers, d_model, num_heads, dff,
7                               input_vocab_size, pe_input, rate)
8
9         self.decoder = Decoder(num_layers, d_model, num_heads, dff,
10                              target_vocab_size, pe_target, rate)
11
12         self.final_layer = tf.keras.layers.Dense(target_vocab_size)
13
14     def call(self, inp, tar, training, enc_padding_mask,
15             look_ahead_mask, dec_padding_mask):
16         # print('enc_padding_mask: ', enc_padding_mask)
17         enc_output = self.encoder(inp, training, enc_padding_mask) # (batch_size, tar_seq_len, d_model)
18
19         # dec_output.shape == (batch_size, tar_seq_len, d_model)
20         dec_output, attention_weights = self.decoder(
21             tar, enc_output, training, look_ahead_mask, dec_padding_mask)
22
23         final_output = self.final_layer(dec_output) # (batch_size, tar_seq_len, target_vocab_size)
24
25         return final_output, attention_weights

```

Top

```

1     sample_transformer = Transformer(
2         num_layers=2, d_model=512, num_heads=8, dff=2048,
3         input_vocab_size=8500, target_vocab_size=8000,
4         pe_input=10000, pe_target=6000)
5
6     temp_input = tf.random.uniform((64, 38), dtype=tf.int64, minval=0, maxval=
7     temp_target = tf.random.uniform((64, 36), dtype=tf.int64, minval=0, maxval=
8
9     fn_out, _ = sample_transformer(temp_input, temp_target, training=False,
10                                     enc_padding_mask=None,
11                                     look_ahead_mask=None,
12                                     dec_padding_mask=None)
13
14     fn_out.shape # (batch_size, tar_seq_len, target_vocab_size)

```

```

1     TensorShape([64, 36, 8000])

```

10. Set hyperparameters

Chúng ta sẽ cần thiết lập các tham số cho mô hình bao gồm:

1. Các tham số kiến trúc mô hình:

- num_layers: Số lượng Attention sub-layer blocks của encoder và decoder.
- d_model: Kích thước của véc tơ embedding.
- num_heads: Số lượng các head trong một attention layer.
- dff: Số lượng token tối đa cho phép của văn bản đầu vào.

1. Các tham số huấn luyện:

- dropout_rate: Tỷ lệ drop out ở output.
- input_vocab_size: Số lượng từ vựng của input.
- target_vocab_size: Số lượng từ vựng của target.

```

1     num_layers = 4
2     d_model = 128
3     dff = 512
4     num_heads = 8
5
6     input_vocab_size = tokenizer_ipt.vocab_size + 2
7     target_vocab_size = tokenizer_opt.vocab_size + 2
8     dropout_rate = 0.1

```

```

1     print('input_vocab_size: ', input_vocab_size)
2     print('target_vocab_size: ', target_vocab_size)

```

```

1     input_vocab_size: 8197
2     target_vocab_size: 8135

```

11. Optimizer

Sử dụng Adam optimizer với Learning rate đã được scheduler theo công thức tại paper (<https://arxiv.org/abs/1706.03762>).

Top

$$lrate = d_{model}^{-0.5} * \min(step_num^{-0.5}, step_num * warmup_steps^{-1.5})$$

```

1 class CustomSchedule(tf.keras.optimizers.schedules.LearningRateSchedule):
2     def __init__(self, d_model, warmup_steps=4000):
3         super(CustomSchedule, self).__init__()
4
5         self.d_model = d_model
6         self.d_model = tf.cast(self.d_model, tf.float32)
7
8         self.warmup_steps = warmup_steps
9
10    def __call__(self, step):
11        arg1 = tf.math.rsqrt(step)
12        arg2 = step * (self.warmup_steps ** -1.5)
13
14        return tf.math.rsqrt(self.d_model) * tf.math.minimum(arg1, arg2)

```

```

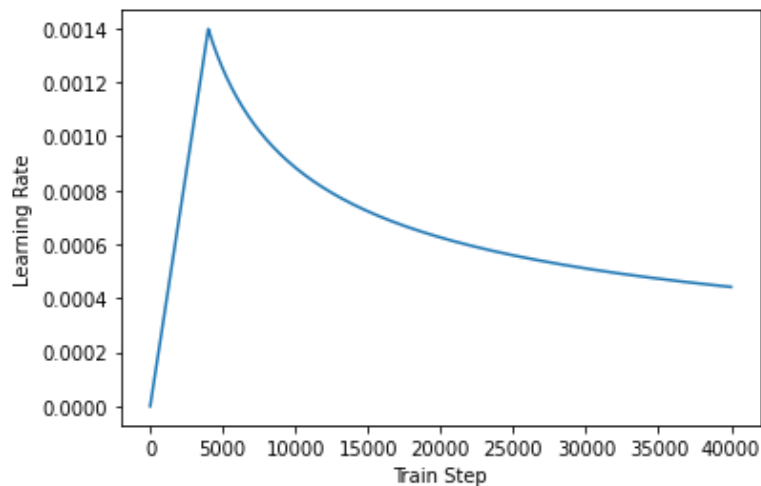
1 learning_rate = CustomSchedule(d_model)
2
3 optimizer = tf.keras.optimizers.Adam(learning_rate, beta_1=0.9, beta_2=0.98,
4                                       epsilon=1e-9)

```

```

1 temp_learning_rate_schedule = CustomSchedule(d_model)
2
3 plt.plot(temp_learning_rate_schedule(tf.range(40000, dtype=tf.float32)))
4 plt.ylabel("Learning Rate")
5 plt.xlabel("Train Step")

```



12. Loss and Metrics

Bởi vì câu mục tiêu và được padded, do đó chúng ta cũng phải áp dụng padding mask khi tính toán loss function.

```

1 loss_object = tf.keras.losses.SparseCategoricalCrossentropy(
2     from_logits=True, reduction='none')

```

Top

```

1 def loss_function(real, pred):
2     mask = tf.math.logical_not(tf.math.equal(real, 0))
3     loss_ = loss_object(real, pred)
4
5     mask = tf.cast(mask, dtype=loss_.dtype)
6     loss_ *= mask
7
8     return tf.reduce_sum(loss_)/tf.reduce_sum(mask)

1 train_loss = tf.keras.metrics.Mean(name='train_loss')
2 train_accuracy = tf.keras.metrics.SparseCategoricalAccuracy(
3     name='train_accuracy')

```

13. Training and checkpointing

```

1 transformer = Transformer(num_layers=num_layers, d_model=d_model, num_head:
2     input_vocab_size=input_vocab_size, target_vocab_size=target_vocab_size,
3     pe_input=input_vocab_size,
4     pe_target=target_vocab_size,
5     rate=dropout_rate)

```

```

1 def create_masks(inp, tar):
2     # Encoder padding mask
3     enc_padding_mask = create_padding_mask(inp)
4
5     # Used in the 2nd attention block in the decoder.
6     # This padding mask is used to mask the encoder outputs.
7     dec_padding_mask = create_padding_mask(inp)
8
9     # Used in the 1st attention block in the decoder.
10    # It is used to pad and mask future tokens in the input received by
11    # the decoder.
12    look_ahead_mask = create_look_ahead_mask(tf.shape(tar)[1])
13    dec_target_padding_mask = create_padding_mask(tar)
14    combined_mask = tf.maximum(dec_target_padding_mask, look_ahead_mask)
15    return enc_padding_mask, combined_mask, dec_padding_mask

```

```

1 x = tf.constant([[7, 6, 0, 0, 0], [1, 2, 3, 0, 0], [6, 0, 0, 0, 0]])
2 y = tf.constant([[1, 4, 5, 0, 0], [1, 4, 3, 0, 0], [1, 2, 0, 0, 0]])
3 enc_padding_mask, combined_mask, dec_padding_mask = create_masks(x, y)

```

```

1 combined_mask

```

Top

```

1      <tf.Tensor: shape=(3, 1, 5, 5), dtype=float32, numpy=
2      array([[[[0., 1., 1., 1., 1.],
3              [0., 0., 1., 1., 1.],
4              [0., 0., 0., 1., 1.],
5              [0., 0., 0., 1., 1.],
6              [0., 0., 0., 1., 1.]]],
7
8
9          [[[0., 1., 1., 1., 1.],
10             [0., 0., 1., 1., 1.],
11             [0., 0., 0., 1., 1.],
12             [0., 0., 0., 1., 1.],
13             [0., 0., 0., 1., 1.]]],
14
15
16          [[[0., 1., 1., 1., 1.],
17             [0., 0., 1., 1., 1.],
18             [0., 0., 1., 1., 1.],
19             [0., 0., 1., 1., 1.],
20             [0., 0., 1., 1., 1.]]]], dtype=float32)>

```

Chúng ta sẽ tạo ra một checkpoint path và checkpoint manager để lưu checkpoints (là trạng thái của model, có thể load lại để huấn luyện lại) sau mỗi $n=5$ epochs.

```

1      checkpoint_path = "./checkpoints/train_500k"
2
3      ckpt = tf.train.Checkpoint(transformer=transformer,
4                                 optimizer=optimizer)
5
6      ckpt_manager = tf.train.CheckpointManager(ckpt, checkpoint_path, max_to_keep=5)
7
8      # if a checkpoint exists, restore the latest checkpoint.
9      if ckpt_manager.latest_checkpoint:
10         ckpt.restore(ckpt_manager.latest_checkpoint)
11         print ('Latest checkpoint restored!!')

```

Note: Layer embedding của Encoder sẽ lưu `input_vocabulary_size` và `target_vocabulary_size` theo bộ dữ liệu. Do đó nếu huấn luyện trên bộ dữ liệu mới thì sẽ không thể load lại model từ checkpoints cũ. Chúng ta cần chuyển sang một bộ checkpoints mới.

Target sẽ được phân thành `tar_inp` và `tar_real`. `tar_inp` được truyền vào như là một input cho decoder. `tar_real` tương tự như `tar_inp` nhưng được dịch chuyển 1 đơn vị: Tại mỗi vị trí trong `tar_inp` ta sẽ có 1 vị trí tương ứng trong `tar_real` là token tiếp theo nên được dự báo.

Chẳng hạn trong câu `sentence = SOS một con sư tử đang ngủ say trong rừng SOS`

`tar_inp` = SOS một con sư tử đang ngủ say trong rừng

`tar_real` = một con sư tử đang ngủ say trong rừng SOS

transformer là một mô hình tự hồi qui (auto-regressive model, một dạng model trong chuỗi thời gian sử dụng chính giá trị quá khứ của chuỗi để dự báo giá trị hiện tại của chuỗi): và tạo ra một dự báo tại mỗi một time step. Nó sử dụng output của nó từ trước đến hiện tại để quyết định xem từ tiếp theo là gì.

Trong quá trình huấn luyện ví dụ này sử dụng kỹ thuật teacher-forcing. Teacher-forcing sẽ bỏ qua giá trị đúng ở output để tới time step tiếp theo mà không quan tâm tới mô hình dự báo ra là gì tại time step hiện tại.

Top

Khi transformer dự báo mỗi từ, self-attention cho phép nó nhìn vào các từ liền trước trong chuỗi input để dự báo tốt hơn từ tiếp theo.

```

1      EPOCHS = 100

1      # The @tf.function trace-compiles train_step into a TF graph for faster
2      # execution. The function specializes to the precise shape of the arguments
3      # tensors. To avoid re-tracing due to the variable sequence lengths or variable
4      # batch sizes (the last batch is smaller), use input_signature to specify
5      # more generic shapes.
6
7      train_step_signature = [
8          tf.TensorSpec(shape=(None, None), dtype=tf.int64),
9          tf.TensorSpec(shape=(None, None), dtype=tf.int64),
10     ]
11
12     @tf.function(input_signature=train_step_signature)
13     def train_step(inp, tar):
14         tar_inp = tar[:, :-1]
15         tar_real = tar[:, 1:]
16         enc_padding_mask, combined_mask, dec_padding_mask = create_masks(inp, tar_inp)
17         with tf.GradientTape() as tape:
18             predictions, _ = transformer(inp, tar_inp,
19                                         True,
20                                         enc_padding_mask,
21                                         combined_mask,
22                                         dec_padding_mask)
23
24             loss = loss_function(tar_real, predictions)
25
26             gradients = tape.gradient(loss, transformer.trainable_variables)
27             optimizer.apply_gradients(zip(gradients, transformer.trainable_variables))
28
29             train_loss(loss)
30             train_accuracy(tar_real, predictions)

```

Top

```

1     for epoch in range(EPOCHS):
2         start = time.time()
3
4         train_loss.reset_states()
5         train_accuracy.reset_states()
6
7         # inp -> non_diacritic, tar -> diacritic
8         for (batch, (inp, tar)) in enumerate(train_dataset):
9             train_step(inp, tar)
10
11        if batch % 50 == 0:
12            print ('Epoch {} Batch {} Loss {:.4f} Accuracy {:.4f}'.format(
13                epoch + 1, batch, train_loss.result(), train_accuracy.result()))
14
15        if (epoch + 1) % 5 == 0:
16            ckpt_save_path = ckpt_manager.save()
17            print ('Saving checkpoint for epoch {} at {}'.format(epoch+1,
18                                                                    ckpt_save_path))
19
20        print ('Epoch {} Loss {:.4f} Accuracy {:.4f}'.format(epoch + 1,
21                                                                train_loss.result(),
22                                                                train_accuracy.result()))
23
24        print ('Time taken for 1 epoch: {} secs\n'.format(time.time() - start))

```

```

1     Epoch 20 Batch 350 Loss 0.2435 Accuracy 0.4345

```

Do quá trình huấn luyện khá lâu nên mình chỉ dừng lại ở epochs 20.

14. Evaluate

Các bước sau đây được sử dụng để đánh giá:

- Mã hóa câu đầu vào bằng cách sử dụng tokenizer input (tokenizer_pt). Ngoài ra , thêm mã bắt đầu và kết thúc để đầu vào tương đương với những gì mô hình đã huấn luyện. Đây sẽ là đầu vào cho encoder.
- Đầu vào cho decoder là start token có index bằng với tokenizer_en.vocab_size.
- Tính toán các padding masks và look ahead masks.
- decoder sau đó đưa ra các dự đoán dựa trên đầu ra của encoder và đầu ra của chính nó (self-attention).
- Chọn từ cuối cùng và tính argmax của từ đó.
- Nối từ dự đoán với đầu vào decoder input khi chuyển từ đó đến decoder.
- Trong phương pháp này, decoder dự đoán từ tiếp theo dựa trên các từ trước đó mà nó dự đoán.

Top

```

1      MAX_LENGTH = 40
2
3      def evaluate(inp_sentence):
4          start_token = [tokenizer_ipt.vocab_size]
5          end_token = [tokenizer_ipt.vocab_size + 1]
6
7          # inp sentence is non_diacritic, hence adding the start and end token
8          inp_sentence = start_token + tokenizer_ipt.encode(inp_sentence) + end_token
9          encoder_input = tf.expand_dims(inp_sentence, 0)
10
11         # as the target is exist diacritic, the first word to the transformer should be
12         # english start token.
13         decoder_input = [tokenizer_opt.vocab_size]
14         output = tf.expand_dims(decoder_input, 0)
15
16         for i in range(MAX_LENGTH):
17             enc_padding_mask, combined_mask, dec_padding_mask = create_masks(
18                 encoder_input, output)
19
20             # predictions.shape == (batch_size, seq_len, vocab_size)
21             predictions, attention_weights = transformer(encoder_input,
22                                                         output,
23                                                         False,
24                                                         enc_padding_mask,
25                                                         combined_mask,
26                                                         dec_padding_mask)
27
28             # select the last word from the seq_len dimension
29             predictions = predictions[:, -1:, :] # (batch_size, 1, vocab_size)
30
31             predicted_id = tf.cast(tf.argmax(predictions, axis=-1), tf.int32)
32
33             # return the result if the predicted_id is equal to the end token
34             if predicted_id == tokenizer_opt.vocab_size+1:
35                 return tf.squeeze(output, axis=0), attention_weights
36
37             # concatenate the predicted_id to the output which is given to the decoder
38             # as its input.
39             output = tf.concat([output, predicted_id], axis=-1)
40
41         return tf.squeeze(output, axis=0), attention_weights

```

Tiếp theo ta sẽ xây dựng biểu đồ heatmap visualize attention weight giữa 2 từ bất kỳ được lấy từ câu input và câu target.

Top

```

1 def plot_attention_weights(attention, sentence, result, layer):
2     fig = plt.figure(figsize=(16, 8))
3
4     sentence = tokenizer_ipt.encode(sentence)
5
6     attention = tf.squeeze(attention[layer], axis=0)
7
8     for head in range(attention.shape[0]):
9         ax = fig.add_subplot(2, 4, head+1)
10
11         # plot the attention weights
12         ax.matshow(attention[head][: -1, :], cmap='viridis')
13
14         fontdict = {'fontsize': 10}
15
16         ax.set_xticks(range(len(sentence)+2))
17         ax.set_yticks(range(len(result)))
18
19         ax.set_ylim(len(result)-1.5, -0.5)
20
21         ax.set_xticklabels(
22             ['<start>']+tokenizer_ipt.decode([i]) for i in sentence)+['<end>']
23             fontdict=fontdict, rotation=90)
24
25         ax.set_yticklabels([tokenizer_opt.decode([i]) for i in result
26                             if i < tokenizer_opt.vocab_size],
27                             fontdict=fontdict)
28
29         ax.set_xlabel('Head {}'.format(head+1))
30
31     plt.tight_layout()
32     plt.show()

```

Hàm translate sẽ có tác dụng chuyển đánh giá câu input và trả ra câu dự báo bằng cách giải mã chuỗi indices token được dự báo từ mô hình.

```

1 def translate(sentence, plot=''):
2     result, attention_weights = evaluate(sentence)
3
4     predicted_sentence = tokenizer_opt.decode([i for i in result
5                                                if i < tokenizer_opt.vocab_size])
6
7     print('Input: {}'.format(sentence))
8     print('Predicted translation: {}'.format(predicted_sentence))
9
10    if plot:
11        plot_attention_weights(attention_weights, sentence, result, plot)

```

```

1 translate("tieng Viet la ngon ngu trong sang nhat the gioi")

```

```

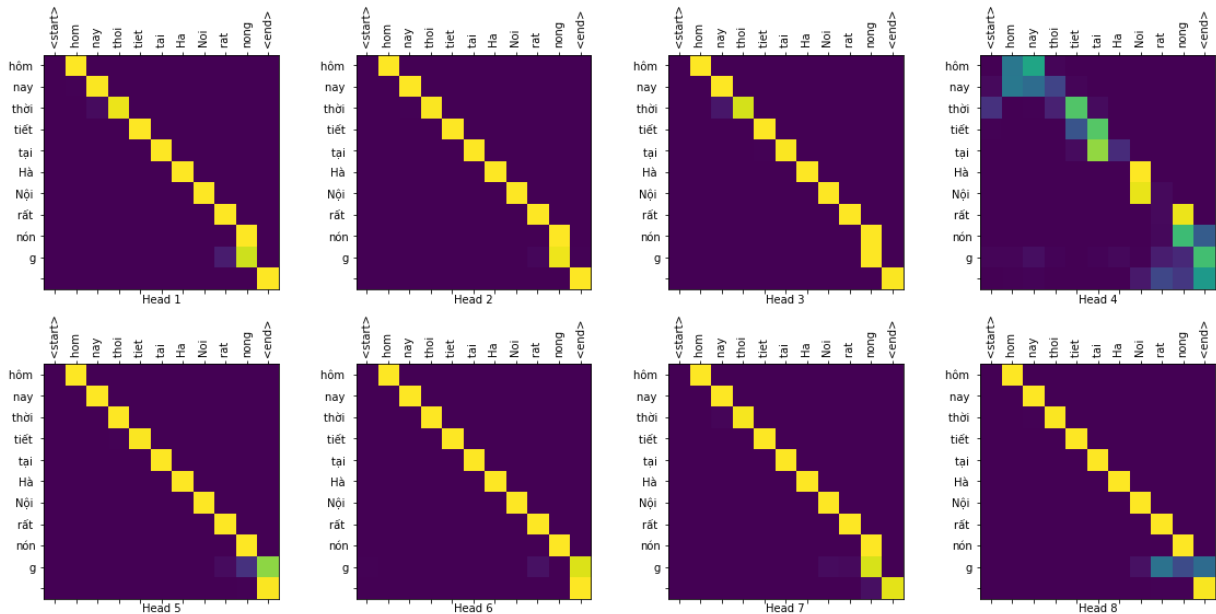
1 Input: tieng Viet la ngon ngu trong sang nhat the gioi
2 Predicted translation: tiếng Việt là ngôn ngữ trong sáng nhất thế giới

```

Top

```
1 translate("hom nay thoi tiet tai Ha Noi rat nong", plot='decoder_layer4_block1')
2 print ("Real translation: hôm nay thời tiết tại Hà Nội rất nóng")
```

```
1 Input: hom nay thoi tiet tai Ha Noi rat nong
2 Predicted translation: hôm nay thời tiết tại Hà Nội rất nóng
```



Ta có thể thấy khi visualize encoder-decoder attention từ một block layer trong decoder cho từng head trong multi-heads thì các vị trí từ ở encoder cùng mà vị trí với từ ở decoder sẽ có attention weight cao hơn.

15. Pretrain model

Nếu bạn không muốn huấn luyện model từ đầu, bạn có thể sử dụng pretrained model đã được tôi huấn luyện bằng cách download các thư mục:

- tokenizer: là tokenizer cho các subwords dùng để chuyển sequence sang indice.
- checkpoint: Là checkpoint, nơi lưu trữ last pretrain model để bạn có thể load lại cho dự báo.

Cả 2 đều được đặt trong folder themdau_tv

(<https://drive.google.com/drive/folders/11KSVIDWINKnVZmoL1ugW8fxHJyrNBCr?usp=sharing>).

```
1 !ls tokenizer
```

```
1 input_vocab.txt    tokenizer_ipt.pkl    tokenizer_opt.pkl
2 output_vocab.txt   tokenizer_ipt.subwords  tokenizer_opt.subwords
```

Top


```

1  import tensorflow_datasets as tfds
2  import tensorflow as tf
3
4  import time
5  import numpy as np
6  import matplotlib.pyplot as plt
7
8  # Load tokenizer
9  import pickle
10
11 def _save_pickle(path, obj):
12     with open(path, 'wb') as f:
13         pickle.dump(obj, f)
14
15 def _load_pickle(path):
16     with open(path, 'rb') as f:
17         obj = pickle.load(f)
18     return obj
19
20 tokenizer_ipt = _load_pickle('tokenizer/tokenizer_ipt.pkl')
21 tokenizer_opt = _load_pickle('tokenizer/tokenizer_opt.pkl')
22
23 # Khai báo tham số
24 num_layers = 4
25 d_model = 128
26 d_ff = 512
27 num_heads = 8
28
29 input_vocab_size = tokenizer_ipt.vocab_size + 2
30 target_vocab_size = tokenizer_opt.vocab_size + 2
31 dropout_rate = 0.1
32 learning_rate = 0.01

```

Nếu các bạn bắt đầu từ mục 15 này ngay từ đầu. Trước khi load model thì các bạn cần chạy lại toàn bộ code cho layers và mô hình ở mục 8 và 9. Đây là hạn chế của notebook so với python file. Trên python file chúng ta có thể import class giữa các files với nhau khá dễ dàng.

```

1  # Load model
2  transformer = Transformer(num_layers=num_layers, d_model=d_model, num_heads=
3                          input_vocab_size=input_vocab_size, target_vocab_size=
4                          pe_input=input_vocab_size,
5                          pe_target=target_vocab_size,
6                          rate=dropout_rate)
7
8  checkpoint_path = "./checkpoints/train_500k"
9
10 ckpt = tf.train.Checkpoint(transformer=transformer)
11
12 ckpt_manager = tf.train.CheckpointManager(ckpt, checkpoint_path, max_to_keep=5)
13
14 # if a checkpoint exists, restore the latest checkpoint.
15 if ckpt_manager.latest_checkpoint:
16     ckpt.restore(ckpt_manager.latest_checkpoint)

```

Top

```

1  def evaluate(inp_sentence):
2      start_token = [ipt_vocab_size]
3      end_token = [ipt_vocab_size+1]
4
5      # inp sentence is non_diacritic, hence adding the start and end token
6      inp_sentence = start_token + tokenizer_ipt.encode(inp_sentence) + end_token
7      encoder_input = tf.expand_dims(inp_sentence, 0)
8
9      # as the target is exist diacritic, the first word to the transformer should be
10     # english start token.
11     decoder_input = [opt_vocab_size]
12     output = tf.expand_dims(decoder_input, 0)
13
14     for i in range(MAX_LENGTH):
15         enc_padding_mask, combined_mask, dec_padding_mask = create_masks(
16             encoder_input, output)
17
18         # predictions.shape == (batch_size, seq_len, vocab_size)
19         predictions, attention_weights = transformer(encoder_input,
20                                                     output,
21                                                     False,
22                                                     enc_padding_mask,
23                                                     combined_mask,
24                                                     dec_padding_mask)
25
26         # select the last word from the seq_len dimension
27         predictions = predictions[:, -1:, :] # (batch_size, 1, vocab_size)
28
29         predicted_id = tf.cast(tf.argmax(predictions, axis=-1), tf.int32)
30
31         # return the result if the predicted_id is equal to the end token
32         if predicted_id == opt_vocab_size+1:
33             return tf.squeeze(output, axis=0), attention_weights
34
35         # concatenate the predicted_id to the output which is given to the decoder
36         # as its input.
37         output = tf.concat([output, predicted_id], axis=-1)
38
39     return tf.squeeze(output, axis=0), attention_weights

```

```

1  def add_diacritic(sentence, plot=''):
2      result, attention_weights = evaluate(sentence)
3      predicted_sentence = tokenizer_opt.decode([i for i in result
4                                              if i < opt_vocab_size])
5      print('Input: {}'.format(sentence))
6      print('Predicted translation: {}'.format(predicted_sentence))
7      if plot:
8          plot_attention_weights(attention_weights, sentence, result, plot)

```

Bên dưới là kết quả của một vài câu thử nghiệm

```
1  add_diacritic("hom nay thoi tiet tai Ha Noi rat nong")
```

```

1  Input: hom nay thoi tiet tai Ha Noi rat nong
2  Predicted translation: hôm nay thời tiết tại Hà Nội rất nóng

```

Top

```
1 add_diacritic("toi la mot nguoi rat yeu thich AI")
```

```
1 Input: toi la mot nguoi rat yeu thich AI
```

```
2 Predicted translation: tôi là một người rất yêu thích AI
```

```
1 add_diacritic("toi muon tro thanh mot AI researcher noi tieng tren the gioi")
```

```
1 Input: toi muon tro thanh mot AI researcher noi tieng tren the gioi
```

```
2 Predicted translation: tôi muốn trở thành một AI researcher nổi tiếng trên
```

16. Tổng kết

Như vậy qua bài viết này chúng ta đã được hướng dẫn từ bước thu thập dữ liệu, xây dựng các layers của kiến trúc transformer, huấn luyện mô hình và dự báo từ pretrained model.

Thông qua việc thực hành chúng ta sẽ hiểu rõ hơn cấu trúc từng layer và các xử lý attention, mask trong tác vụ seq2seq. Bản thân tôi cũng phải bỏ ra rất nhiều thời gian để nghiên cứu bài viết này và viết lại như một tài liệu tham khảo khi cần. Bài viết được lấy từ nhiều nguồn liệt kê bên dưới.

17. Tài liệu tham khảo

1. Transformer model for language understanding
(<https://www.tensorflow.org/tutorials/text/transformer>)
2. Bài 4 - Attention is all you need - Khanh Blog
(<https://phamdinhhkhanh.github.io/2019/06/18/AttentionLayer.html>)
3. Bài 7 - Pytorch Seq2seq model correct spelling - Khanh Blog
(<https://phamdinhhkhanh.github.io/2019/08/19/CorrectSpellingVietnamseTonePrediction.html>)
4. Bài 36 - BERT model - Khanh Blog
(<https://phamdinhhkhanh.github.io/2020/05/23/BERTModel.html>)

Top