

Bài 39 - Thực hành ứng dụng BERT

04 Jun 2020 - phamdinhhkhanh

Menu

- 1. BERT trong Tiếng Việt
- 2. Kiến trúc RoBERTa
- 3. Load model BERT
- 4. Tìm hiểu về mã hóa BPE (Byte Pair Encoding)
- 5. Extract features từ RoBERTa
- 6. Điền từ (Filling mask)
- 7. Trích xuất đặc trưng (Extract feature) cho từ
- 8. Bài toán classification
 - 8.1. Kiến trúc mô hình
 - 8.2. Dữ liệu
 - 8.2.1. Đọc và lưu dữ liệu
 - 8.2.2. Tokenize nội dung
 - 8.3. Tokenize Input và output
 - 8.4. Load model BERT
 - 8.5. Huấn luyện model
- 9. Huấn luyện RoBERTa trên dữ liệu của bạn
- 10. Tổng kết
- 11. Tài liệu

1. BERT trong Tiếng Việt

Ở bài 36 chúng ta đã tìm hiểu về các kiến trúc của model BERT gồm BERT Base, BERT Large và những ứng dụng trong các tác vụ NLP của nó. Sự ra đời của model BERT là một cột mốc rất quan trọng của ngành NLP mà có thể phân chia thành giai đoạn phát triển trước BERT và sau BERT. Các kết quả ứng dụng BERT đã phá vỡ các giới hạn trong NLP với rất nhiều các pretrain model xác lập kết quả SOTA trong các tác vụ. Để chứng minh cho những gì tôi nói là không nhầm nhí, bạn đọc có thể theo dõi tại leader board GLUE benchmark

(<https://gluebenchmark.com/leaderboard>). Bên cạnh đó BERT giúp cho quá trình học chuyển giao trở nên khả thi hơn khi có thể can thiệp và fine tuning mô hình học sâu nhiều tầng ở mức độ sâu thay vì can thiệp nông như các mô hình trước.

Kể từ khi google public mã nguồn mở của BERT, đã có rất nhiều các dự án mã nguồn mở về BERT hỗ trợ huấn luyện và chia sẻ các mô hình pretrain BERT trên ngôn ngữ đơn phương và song ngữ. Đối với Tiếng Việt chúng ta có PhoBERT (<https://github.com/VinAIRResearch/PhoBERT>). Cá nhân mình sử dụng PhoBERT thì thấy các tác vụ NLP trong Tiếng Việt được cải thiện và đạt độ chính xác cao. Bạn đọc cũng có thể tự cảm nhận qua các phần thực hành ở bài hướng dẫn này. Trong Tiếng Việt thì chúng ta có thể ứng dụng BERT trong một số tác vụ như:

- Tìm từ đồng nghĩa, trái nghĩa, cùng nhóm dựa trên khoảng cách của từ trong không gian biểu diễn đa chiều.
- Xây dựng các véc tơ embedding cho các tác vụ NLP như sentiment analysis, phân loại văn bản, NER, POS, huấn luyện chatbot.
- Gợi ý từ khóa tìm kiếm trong các hệ thống search.
- Xây dựng các ứng dụng seq2seq như robot viết báo, tóm tắt văn bản, sinh câu ngẫu nhiên với ý nghĩa tương đồng.

Top

Và nhiều những ứng dụng khác mà mình có thể chưa liệt kê hết, rất mong bạn đọc bổ sung thêm. Mặc dù model BERT có rất nhiều các ứng dụng có thể fine tuning nhưng không thực sự nhiều bạn biết cách áp dụng. Một phần là bởi để fine tuning được BERT đòi hỏi bạn phải có kỹ năng lập trình với các deep learning framework như pytorch, tensorflow và thực sự hiểu sâu về kiến trúc và nguyên lý hoạt động của BERT. Gần đây mình nhận được một vài inbox hỏi về cách áp dụng BERT như thế nào trong các tác vụ NLP. Mình đã dành một thời gian để tìm hiểu và nghiên cứu sâu về mã nguồn và tham khảo các hướng dẫn. Chính vì vậy, bài viết này mình sẽ chia sẻ lại các ứng dụng của model BERT đối với Tiếng Việt mà mình đúc kết được. Nếu bạn đọc có thêm nhiều cách ứng dụng mới của BERT trong Tiếng Việt thì mình rất vui để đón nhận chia sẻ từ các bạn.

Trước khi tìm hiểu bài này mình khuyến nghị các bạn nên đọc qua Bài 36 - BERT model (<https://phamdinhhkhanh.github.io/2020/05/23/BERTModel.html>) để hiểu về model BERT là gì và nguyên lý hoạt động của model BERT.

2. Kiến trúc RoBERTa

RoBERTa là một project của facebook kế thừa lại các kiến trúc và thuật toán của model BERT trên framework pytorch (pytorch cũng là một framework do facebook phát triển, rất được ưa chuộng bởi cộng đồng AI). Đây là một project hỗ trợ việc huấn luyện lại các model BERT trên những bộ dữ liệu mới cho các ngôn ngữ khác ngoài một số ngôn ngữ phổ biến. Kể từ khi ra đời, đã có rất nhiều các mô hình pretrain cho những ngôn ngữ khác nhau được huấn luyện trên RoBERTa.

Ở bài báo gốc cho biết mặc dù RoBERTa lặp lại các thủ tục huấn luyện từ model BERT, nhưng có một thay đổi đó là huấn luyện mô hình lâu hơn, với batch size lớn hơn và trên nhiều dữ liệu hơn. Ngoài ra để nâng cao độ chuẩn xác trong biểu diễn từ thì RoBERTa đã loại bỏ tác vụ dự đoán câu tiếp theo và huấn luyện trên các câu dài hơn. Đồng thời mô hình cũng thay đổi linh hoạt kiểu masking (tức ẩn đi một số từ ở câu output bằng token <mask>) áp dụng cho dữ liệu huấn luyện.

Bạn đọc có thể tìm hiểu thêm về kiến trúc này qua bài báo về RoBERTa (<https://arxiv.org/abs/1907.11692>).

Ở các mục tiếp theo mình sẽ hướng dẫn các bạn triển khai áp dụng model RoBERTa thông qua pretrain model PhoBERT cho Tiếng Việt.

Để bắt đầu bài thực hành, bạn đọc có thể mở file PhoBERT - tutorial Khanh Blog (<https://colab.research.google.com/drive/16a4XFPioXYZQwyTusmzi1liGP8kCHT9t?usp=sharing>) và bắt đầu từ đây.

3. Load model BERT

Để áp dụng được model BERT thì trước tiên chúng ta cần phải load được model. Ví dụ này mình sẽ thực hành trên google colab. Bạn đọc cần mount google drive bằng câu lệnh bên dưới.

```
1 from google.colab import drive
2 import os
3
4 drive.mount('/content/gdrive')
5 path = "/content/gdrive/My Drive/Colab Notebooks/BERT"
6 os.chdir(path)
7 !ls
```

Chúng ta sẽ cần cài đặt các dependency packages sau đây:

Top

- fairseq (<https://github.com/pytorch/fairseq>): Là project của facebook chuyên hỗ trợ các nghiên cứu và dự án liên quan đến model seq2seq.
- fastBPE: Là package hỗ trợ tokenize từ (word) thành các từ phụ (subwords) theo phương pháp mới nhất được áp dụng cho các pretrain model NLP hiện đại như BERT và các biến thể của BERT.
- vncoreNLP: Là một package NLP trong Tiếng Việt, hỗ trợ tokenize và các tác vụ NLP khác.
- transformers (<https://github.com/huggingface/transformers>): Là một project của huggingface hỗ trợ huấn luyện các model dựa trên kiến trúc transformer như BERT, GPT-2, RoBERTa, XLM, DistilBert, XLNet, T5, CTRL,... phục vụ cho các tác vụ NLP trên cả nền tảng pytorch và tensorflow.

```
1 !pip3 install fairseq
2 !pip3 install fastbpe
3 !pip3 install vncoreNLP
4 !pip3 install transformers
```

Tiếp theo là download các model pretrain từ list các pretrain models được liệt kê trong PhoBERT (<https://github.com/VinAIRResearch/PhoBERT>).

Trong hướng dẫn này mình chỉ sử dụng pretrain model BERT base được huấn luyện từ package fairseq. Download và giải nén chúng bằng lần lượt các lệnh wget và tar .

```
1 !wget https://public.vinai.io/PhoBERT_base_fairseq.tar.gz
2 !tar -xzf PhoBERT_base_fairseq.tar.gz
```

Sau khi download và giải nén pretrain model chúng ta sẽ kiểm tra thấy bên trong folder sẽ bao gồm 3 files đó là bpe.codes, dict.txt, model.pt có tác dụng như sau:

- bpe.codes: Là BPE token mà mô hình đã áp dụng để mã hóa văn bản sang index.
- dict.txt: Từ điển subword của bộ dữ liệu huấn luyện.
- model.pt: File lưu trữ của mô hình trên pytorch.

Về BPE và subword là gì mình sẽ lý giải ở chương 4. Tìm hiểu về mã hóa BPE (Byte Pair Encoding) .

```
1 !ls PhoBERT_base_fairseq

1 bpe.codes dict.txt model.pt
```

Load model pretrain PhoBERT

```
1 # Load the model in fairseq
2 from fairseq.models.roberta import RobertaModel
3 phoBERT = RobertaModel.from_pretrained('PhoBERT_base_fairseq', checkpoint_name='dict.txt')
4 phoBERT.eval() # disable dropout (or leave in train mode to finetune)
```



Top

```

1      loading archive file PhoBERT_base_fairseq
2      | dictionary: 64000 types
3
4
5
6      RobertaHubInterface(
7          (model): RobertaModel(
8              (decoder): RobertaEncoder(
9                  (sentence_encoder): TransformerSentenceEncoder(
10                     (embed_tokens): Embedding(64001, 768, padding_idx=1)
11                     (embed_positions): LearnedPositionalEmbedding(258, 768, padding_idx=1)
12                     (layers): ModuleList(
13                         (0): TransformerSentenceEncoderLayer(
14                             (self_attn): MultiheadAttention(
15                                 (k_proj): Linear(in_features=768, out_features=768, bias=True)
16                                 (v_proj): Linear(in_features=768, out_features=768, bias=True)
17                                 (q_proj): Linear(in_features=768, out_features=768, bias=True)
18                                 (out_proj): Linear(in_features=768, out_features=768, bias=True)
19                             )
20                             (self_attn_layer_norm): LayerNorm((768,), eps=1e-05, elementwise_affine=True)
21                             (fc1): Linear(in_features=768, out_features=3072, bias=True)
22                             (fc2): Linear(in_features=3072, out_features=768, bias=True)
23                             (final_layer_norm): LayerNorm((768,), eps=1e-05, elementwise_affine=True)
24                         )
25                     ...
26                     (11): TransformerSentenceEncoderLayer(
27                         (self_attn): MultiheadAttention(
28                             (k_proj): Linear(in_features=768, out_features=768, bias=True)
29                             (v_proj): Linear(in_features=768, out_features=768, bias=True)
30                             (q_proj): Linear(in_features=768, out_features=768, bias=True)
31                             (out_proj): Linear(in_features=768, out_features=768, bias=True)
32                         )
33                         (self_attn_layer_norm): LayerNorm((768,), eps=1e-05, elementwise_affine=True)
34                         (fc1): Linear(in_features=768, out_features=3072, bias=True)
35                         (fc2): Linear(in_features=3072, out_features=768, bias=True)
36                         (final_layer_norm): LayerNorm((768,), eps=1e-05, elementwise_affine=True)
37                     )
38                 )
39                 (emb_layer_norm): LayerNorm((768,), eps=1e-05, elementwise_affine=True)
40             )
41             (lm_head): RobertaLMHead(
42                 (dense): Linear(in_features=768, out_features=768, bias=True)
43                 (layer_norm): LayerNorm((768,), eps=1e-05, elementwise_affine=True)
44             )
45         )
46         (classification_heads): ModuleDict({})
47     )
48 )

```

Ta có thể thấy kiến trúc RoBERTa theo BERT base đã giữ lại 12 block sub-layers là các multi-head attention ở phase Encoder và thêm một linear projection layer ở cuối để tạo ra véc tơ embedding cho từ. Ở mỗi multi-head attention, các output self-attention được khởi tạo từ quá trình nhân kết hợp giữa các ma trận chiếu Key, Value và Query. Các bạn có thể tìm hiểu về quá trình này ở Bài 4 - Attention is all you need (<https://phamdinhhkhanh.github.io/2019/06/18/AttentionLayer.html>) của block.

4. Tìm hiểu về mã hóa BPE (Byte Pair Encoding)

Toknenize là quá trình mã hóa các văn bản thành các index dạng số mang thông tin của văn bản để cho máy tính có thể huấn luyện được. Khi đó mỗi một từ hoặc ký tự sẽ được đại diện bởi một index.

Trong NLP có một số kiểu tokenize như sau:

Tokenize theo word level: Chúng ta phân tách câu thành các token được ngăn cách bởi khoảng trắng hoặc dấu câu. Khi đó mỗi token là một từ đơn âm tiết. Đây là phương pháp token được sử dụng trong các thuật toán nhúng từ truyền thống như GloVe, word2vec.

Tokenize theo multi-word level: Tiếng Việt và một số ngôn ngữ khác tồn tại từ đơn âm tiết (từ đơn) và từ đa âm tiết (từ ghép). Do đó nếu tokenize theo từ đơn âm tiết sẽ làm nghĩa của từ bị sai khác. Ví dụ cụm từ vô xác định nếu được chia thành vô, xác và định sẽ làm cho từ bị mất đi nghĩa phủ định của nó. Do đó để tạo ra được các từ với nghĩa chính xác thì chúng ta sẽ sử dụng thêm từ điển bao gồm cả từ đa âm tiết và đơn âm để tokenize câu. Trong Tiếng Việt có khá nhiều các module hỗ trợ tokenize dựa trên từ điển như VnCoreNLP, pyvivn, underthesea.

Tokenize theo character level: Việc tokenize theo word level thường sinh ra một từ điển với kích thước rất lớn, điều này làm gia chi phí tính toán. Hơn nữa nếu tokenize theo word level thì đòi hỏi từ điển phải rất lớn thì mới hạn chế được những trường hợp từ nằm ngoài từ điển. Tuy nhiên nếu phân tích ta sẽ thấy hầu hết các từ đều có thể biểu thị dưới một nhóm các ký tự là chữ cái, con số, dấu xác định. Như vậy chỉ cần sử dụng một lượng các ký tự rất nhỏ có thể biểu diễn được mọi từ. Từ được token dựa trên level ký tự sẽ có tác dụng giảm kích thước từ điển mà vẫn biểu diễn được các trường hợp từ nằm ngoài từ điển. Đây là phương pháp được áp dụng trong mô hình fasttext.

Phương pháp mới BPE (SOTA): Nhược điểm của phương pháp tokenize theo character level đó là các token không có ý nghĩa nếu đứng độc lập. Do đó đối với các bài toán sentiment analysis, áp dụng tokenize theo character level sẽ mang lại kết quả kém hơn. Token theo word level cũng tồn tại hạn chế đó là không giải quyết được các trường hợp từ nằm ngoài từ điển.

Một phương pháp mới đã được đề xuất trong bài báo Neural Machine Translation of Rare Words with Subword Units (<https://arxiv.org/pdf/1508.07909.pdf>) vào năm 2016, có khả năng tách từ theo level nhỏ hơn từ và lớn hơn ký tự được gọi là subword. Phương pháp đó chính là BPE (byte pair encoding). Theo phương pháp mới này, hầu hết các từ đều có thể biểu diễn bởi subword và chúng ta sẽ hạn chế được một số lượng đáng kể các token <unk> đại diện cho từ chưa từng xuất hiện trước đó. Rất nhanh chóng, Phương pháp mới đã được áp dụng ở hầu hết các phương pháp NLP hiện đại từ các lớp model BERT cho tới các biến thể của nó như OpenAI GPT, RoBERTa, DistilBERT, XLNet. Kết quả áp dụng tokenize theo phương pháp mới đã cải thiện được độ chính xác trên nhiều tác vụ dịch máy, phân loại văn bản, dự báo câu tiếp theo, hỏi đáp, dự báo mối quan hệ văn bản.

Thuật toán BPE:

BPE (Byte Pair Encoding) là một kỹ thuật nén từ cơ bản giúp chúng ta index được toàn bộ các từ kể cả trường hợp từ mới (không xuất hiện trong từ điển) nhờ mã hóa các từ bằng chuỗi các từ phụ (subwords). Nguyên lý hoạt động của BPE dựa trên phân tích trực quan rằng hầu hết các từ đều có thể phân tích thành các thành phần con.

Chẳng hạn như từ: low, lower, lowest đều là hợp thành bởi low và những đuôi phụ er, est. Những đuôi này rất thường xuyên xuất hiện ở các từ. Như vậy khi biểu diễn từ lowTop chúng ta có thể mã hóa chúng thành hai thành phần từ phụ (subwords) tách biệt là low và er.

Theo cách biểu diễn này sẽ không phát sinh thêm một index mới cho từ `lower` và đồng thời tìm được mối liên hệ giữa `lower`, `lowest` và `low` nhờ có chung thành phần từ phụ là `low`.

Phương pháp BPE sẽ thống kê tần suất xuất hiện của các từ phụ cùng nhau và tìm cách gộp chúng lại nếu tần suất xuất hiện của chúng là lớn nhất. Cứ tiếp tục quá trình gộp từ phụ cho tới khi không tồn tại các subword để gộp nữa, ta sẽ thu được tập subwords cho toàn bộ văn bản mà mọi từ đều có thể biểu diễn được thông qua subwords.

Code của thuật toán BPE đã được tác giả chia sẻ tại `subword-nmt` (<https://github.com/rsennrich/subword-nmt>).

Quá trình này gồm các bước như sau:

- Bước 1: Khởi tạo từ điển (vocabulary).
- Bước 2: Biểu diễn mỗi từ trong bộ văn bản bằng kết hợp của các ký tự với token `<\w>` ở cuối cùng đánh dấu kết thúc một từ (lý do thêm token sẽ được giải thích bên dưới).
- Bước 3: Thống kê tần suất xuất hiện theo cặp của toàn bộ token trong từ điển.
- Bước 4: Gộp các cặp có tần suất xuất hiện lớn nhất để tạo thành một n-gram theo level character mới cho từ điển.
- Bước 5: Lặp lại bước 3 và bước 4 cho tới khi số bước triển khai merge đạt đỉnh hoặc kích thước kỳ vọng của từ điển đạt được.

Bạn sẽ dễ hình dung hơn qua ví dụ bên dưới:

Giả sử từ điển của chúng ta gồm các từ với tần suất như sau: `vocab = {'l o w </w>': 5, 'l o w e r </w>': 2, 'n e w e s t </w>': 6, 'w i d e s t </w>': 3}`.

Coi mỗi ký tự là một token. Khi đó thống kê tần suất xuất hiện của các cặp ký tự như sau:

```
1      import collections
2
3      vocab = {'l o w </w>': 5, 'l o w e r </w>': 2, 'n e w e s t </w>': 6,
4
5      def get_stats(vocab):
6          pairs = collections.defaultdict(int)
7          for word, freq in vocab.items():
8              symbols = word.split()
9              for i in range(len(symbols)-1):
10                 pairs[symbols[i],symbols[i+1]] += freq
11          return pairs
12
13      get_stats(vocab)
```

Top

```

1    defaultdict(int,
2        {'d', 'e'): 3,
3        ('e', 'r'): 2,
4        ('e', 's'): 9,
5        ('e', 'w'): 6,
6        ('i', 'd'): 3,
7        ('l', 'o'): 7,
8        ('n', 'e'): 6,
9        ('o', 'w'): 7,
10       ('r', '</w>'): 2,
11       ('s', 't'): 9,
12       ('t', '</w>'): 9,
13       ('w', '</w>'): 5,
14       ('w', 'e'): 8,
15       ('w', 'i'): 3})

```

Lựa chọn cặp từ phụ có tần suất xuất hiện nhỏ nhất và merge chúng thành một từ phụ mới.

```

1    import re, collections
2
3    pairs = get_stats(vocab)
4    best = max(pairs, key=pairs.get)
5    print('max pair frequency: ', best)
6
7    # Hàm merge byte max frequency
8
9    def merge_vocab(pair, v_in):
10       v_out = {}
11       bigram = re.escape(' '.join(pair))
12       # Tìm kiếm các vị trí xuất hiện pair bytes
13       p = re.compile(r'(?!\S)' + bigram + r'(?!\S)')
14       for word in v_in:
15           # Thay thế các cặp pair bytes bằng single byte gộp
16           w_out = p.sub(' '.join(pair), word)
17           v_out[w_out] = v_in[word]
18       return v_out
19
20    merge_vocab(best, vocab)

```

```

1    max pair frequency:  ('e', 's')
2
3
4
5
6
7    {'l o w </w>': 5,
8     'l o w e r </w>': 2,
9     'n e w e s t </w>': 6,
10     'w i d e s t </w>': 3}

```

Lặp lại quá trình thống kê tần suất cặp từ và gộp cặp từ với số lượt gộp là 1000

Top

```

1 vocab = {'l o w </w>': 5, 'l o w e r </w>': 2, 'n e w e s t </w>': 6,
2
3 def get_tokens(vocab):
4     tokens = collections.defaultdict(int)
5     for word, freq in vocab.items():
6         word_tokens = word.split()
7         for token in word_tokens:
8             tokens[token] += freq
9     return tokens
10
11 num_merges = 1000
12
13 for i in range(num_merges):
14     pairs = get_stats(vocab)
15     # max_freq = max(pairs.values())
16     # if max_freq == 1:
17     #     break
18
19     if not pairs:
20         break
21     best = max(pairs, key=pairs.get)
22     # print('best', best)
23     vocab = merge_vocab(best, vocab)
24     print('Iter: {}'.format(i))
25     print('Best pair: {}'.format(best))
26     tokens = get_tokens(vocab)
27     print('Tokens: {}'.format(tokens))
28     print('Number of tokens: {}'.format(len(tokens)))
29     print('=====')

```

```

1 Iter: 10
2 Best pair: ('wi', 'd')
3 Tokens: defaultdict(<class 'int'>, {'low</w>': 5, 'low': 2, 'e': 2, '
4 Number of tokens: 8
5 =====
6 Iter: 14
7 Best pair: ('lower', '</w>')
8 Tokens: defaultdict(<class 'int'>, {'low</w>': 5, 'lower</w>': 2, 'ne
9 Number of tokens: 4
10 =====

```

Ta nhận thấy qua các lượt merge từ phụ, độ dài của các từ phụ trong từ điển tăng dần. Thuật toán hội tụ trước 1000 vòng lặp vì toàn bộ các từ phụ đã được merge và đạt ngưỡng của từng từ đơn.

Khi giới hạn kích thước của từ điển hoặc số lượng lượt merge ta sẽ thu được một từ điển từ phụ là thành phần của các từ trong từ điển. Khi đó mọi từ mới dường như sẽ có thể biểu diễn được theo từ phụ.

Ví dụ: Khi dừng số lượt merge tại bước 10 ta thu được từ điển: {'low</w>': 5, 'low': 2, 'e': 2, 'r': 2, '</w>': 2, 'newest</w>': 6, 'wid': 3, 'est</w>': 3}.

Khi đó ta có thể biểu diễn một token mới chưa từng xuất hiện trong từ điển là wider thành ^{Top}wid e r . Bạn đọc đã hình dung được tác dụng của từ phụ (subword) rồi chứ?

Tác dụng của token </w>

Giả định khi tokenize câu the highest mountain theo từ phụ ta thu được biểu diễn ['the</w>', 'high', 'est</w>', 'moun', 'tain</w>']. Khi đó để khôi phục được thành câu gốc ta chỉ cần nối các token lại theo thứ tự thành the</w>highest</w>mountain</w>. Chỉ cần thay </w> bằng khoảng trắng ta sẽ khôi phục được câu gốc: the highest mountain.

token </w> được thêm vào cuối mỗi từ để phân biệt các từ phụ nằm ở vị trí cuối câu với các vị trí khác để giúp cho việc giải mã token khả thi hơn.

Áp dụng BPE tokenize trong BERT:

Hầu hết các mô hình NLP hiện đại nhất đều đã chuyển sang tokenize theo BPE. Để sử dụng BPE tokenize từ các model pretrain của BERT ta thực hiện như sau:

Load model pretrain RoBERTa

```
1 # Load the model in fairseq
2 from fairseq.models.roberta import RobertaModel
3 phoBERT = RobertaModel.from_pretrained('PhoBERT_base_fairseq', checkpo:
4 phoBERT.eval() # disable dropout (or leave in train mode to finetune
```

Khai báo bpe tokenizer và thực hiện token.

```
1 from fairseq.data.encoders.fastbpe import fastBPE
2
3 # Khởi tạo Byte Pair Encoding cho PhoBERT
4 class BPE():
5     bpe_codes = 'PhoBERT_base_fairseq/bpe.codes'
6
7 args = BPE()
8 phoBERT.bpe = fastBPE(args) #Incorporate the BPE encoder into PhoBERT
9 tokens = phoBERT.encode('Tôn Ngộ Không phò Đường Tăng đi Tây Trúc thi
10 print('tokens list : ', tokens)
11 # Decode ngược lại thành câu từ chuỗi index token
12 phoBERT.decode(tokens) # 'Hello world!'
```

```
1 tokens list : tensor([ 0, 11623, 31433, 453, 44334, 2080, 5922,
2                 31686, 3078, 2])
3
4
5 'Tôn Ngộ Không phò Đường Tăng đi Tây Trúc thỉnh kinh'
```

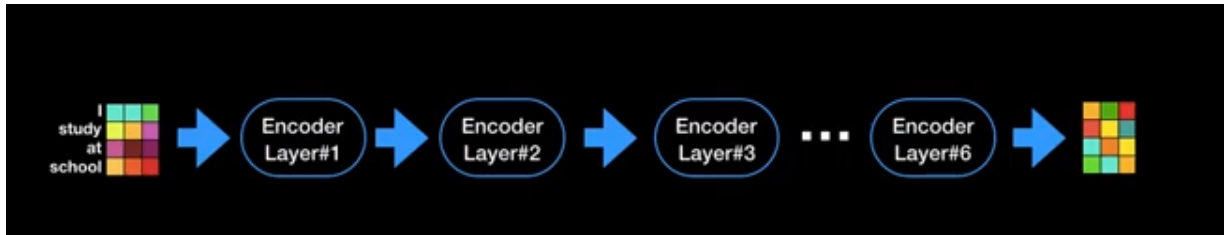
Số lượng các token là 13 lớn hơn số lượng các từ là 11 là vì BERT đã tự thêm các ký tự <s> và </s> đánh dấu từ bắt đầu và kết thúc câu.

5. Extract features từ RoBERTa

Có 2 versions chính trả về 2 kích thước embedding khác nhau khi huấn luyện theo RoBERTa đó là:

- BERT base : 12 sub-layers, kích thước embedding 768, số lượng head attention là 12.
- BERT large : 24 sub-layers, kích thước embedding 1024, số lượng head attention là 16.

Chúng ta có thể trích xuất được các đặc trưng được tạo ra từ BERT của phase Encoder tại layers cuối cùng hoặc toàn bộ các layers. Ngoài phương án trích xuất véc tơ embedding tại layer cuối cùng, một số tác vụ classification trong NLP đã áp dụng trích xuất đặc trưng từ từ những layers trước đó chẳng hạn như trong tác vụ PhoBERT Sentiment Classification (<https://github.com/suicao/PhoBERT-Sentiment-Classification>) tác giả đã trích xuất đặc trưng từ 4 layers cuối cùng thay vì chỉ trích xuất từ một layer cuối.



Hình 1: Kiến trúc gồm nhiều layers tại encoder của model BERT. Mô hình huấn luyện từ RoBERTa cho phép ta trích xuất các đặc trưng từ những layers của encoder. Có thể là layer cuối hoặc toàn bộ các layers.

Kích thước output của mỗi một layer sẽ là `batch_size x seq_len x d_model`. Phương pháp trích xuất như sau:

```
1 # Extract the last layer's features
2 last_layer_features = phoBERT.extract_features(tokens)
3 # assert last_layer_features.size() == torch.Size([1, 5, 1024])
4 print('token size: ', tokens.size())
5 print('size of last layer: ', last_layer_features.size())
6
7 # Extract all layer's features (layer 0 is the embedding layer)
8 all_layers = phoBERT.extract_features(tokens, return_all_hiddens=True)
9 print('number layer in all layers: ', len(all_layers))
10
11 # last_layer_features must equal to last layer in all_layers:
12 print('Last layer features: ', all_layers[-1] == last_layer_features)
```

```
1 token size: torch.Size([13])
2 size of last layer: torch.Size([1, 13, 768])
3 number layer in all layers: 13
4 Last layer features: tensor([[[True, True, True, ..., True, True, T
5     [True, True, True, ..., True, True, True],
6     [True, True, True, ..., True, True, True],
7     ...,
8     [True, True, True, ..., True, True, True],
9     [True, True, True, ..., True, True, True],
10    [True, True, True, ..., True, True, True]]])
```

6. Điền từ (Filling mask)

Trong bài toán này chúng ta sẽ điền các từ hợp lý vào các vị trí còn trống của câu. Trên thực tế có rất nhiều ứng dụng của bài toán filling mask như xây dựng hệ thống suggestion search, gợi ý gõ văn bản, tìm từ đồng nghĩa, tagging.

Mô hình BERT tạo ra các biểu diễn từ từ quá trình ẩn các vị trí token một cách ngẫu nhiên trong câu input và dự báo chính xác từ đó ở output dựa trên bối cảnh là các từ xung quanh.

Top

Như vậy khi đã biết các từ xung quanh, chúng ta hoàn toàn có thể dự báo được từ phù hợp nhất với vị trí đã được masking.

Download package VnCoreNLP để tokenize các câu văn.

```

1 Download VnCoreNLP-1.1.1.jar & its word segmentation component (i.e. R
2 !mkdir -p vncorenlp/models/wordsegmenter
3 !wget https://raw.githubusercontent.com/vncorenlp/VnCoreNLP/master/VnCo
4 !wget https://raw.githubusercontent.com/vncorenlp/VnCoreNLP/master/mod
5 !wget https://raw.githubusercontent.com/vncorenlp/VnCoreNLP/master/mod
6 !mv VnCoreNLP-1.1.1.jar vncorenlp/
7 !mv vi-vocab vncorenlp/models/wordsegmenter/
8 !mv wordsegmenter.rdr vncorenlp/models/wordsegmenter/

```

Giả sử chúng ta có câu gốc là Tôn Ngộ Không phò Đường Tăng đi thỉnh kinh tại Tây Trúc. Từ được ẩn đi trong câu là phò sẽ được thay thế bằng token <mask>.

```

1 from vncorenlp import VnCoreNLP
2 rdrsegmenter = VnCoreNLP("vncorenlp/VnCoreNLP-1.1.1.jar", annotators=
3
4 text = 'Tôn Ngộ Không phò Đường Tăng đi thỉnh kinh tại Tây Trúc'
5 text_masked = 'Học sinh được <mask> do dịch covid-19'
6 # Tokenize câu gốc và thay từ phò bằng <mask>
7 words = rdrsegmenter.tokenize(text)[0]
8 for i, token in enumerate(words):
9     if token == 'phò':
10         words[i] = ' <mask>'
11 text_masked_tok = ' '.join(words)
12 print('text_masked_tok: \n', text_masked_tok)

```

```

1 text_masked_tok:
2 Tôn_Ngộ_Không <mask> Đường_Tăng đi thỉnh_kinh tại Tây_Trúc

```

Tìm ra top 10 từ thích hợp nhất cho vị trí <mask> tại câu trên.

```

1  from fairseq.data.encoders.fastbpe import fastBPE
2  from fairseq import options
3  import numpy as np
4
5  # Khởi tạo Byte Pair Encoding cho PhoBERT
6  class BPE():
7      bpe_codes = 'PhoBERT_base_fairseq/bpe.codes'
8      args = BPE()
9      phoBERT.bpe = fastBPE(args) #Incorporate the BPE encoder into PhoBERT
10
11     # Filling marks
12     topk_filled_outputs = phoBERT.fill_mask(text_masked_tok, topk=10)
13     topk_probs = [item[1] for item in topk_filled_outputs]
14     print('Total probability: ', np.sum(topk_probs))
15     print('Input sequence: ', text_masked_tok)
16     print('Top 10 in mask: ')
17     for i, output in enumerate(topk_filled_outputs):
18         print(output[0])

```

```

1  Total probability:  0.8735223989933729
2  Input sequence:  Tôn_Ngộ_Không <mask> Đường Tăng đi thỉnh_kinh tại T
3  Top 10 in mask:
4  Tôn_Ngộ_Không và Đường Tăng đi thỉnh_kinh tại Tây_Trúc
5  Tôn_Ngộ_Không đưa Đường Tăng đi thỉnh_kinh tại Tây_Trúc
6  Tôn_Ngộ_Không cõng Đường Tăng đi thỉnh_kinh tại Tây_Trúc
7  Tôn_Ngộ_Không hộ_tống Đường Tăng đi thỉnh_kinh tại Tây_Trúc
8  Tôn_Ngộ_Không cùng Đường Tăng đi thỉnh_kinh tại Tây_Trúc
9  Tôn_Ngộ_Không chở Đường Tăng đi thỉnh_kinh tại Tây_Trúc
10  Tôn_Ngộ_Không theo Đường Tăng đi thỉnh_kinh tại Tây_Trúc
11  Tôn_Ngộ_Không dẫn Đường Tăng đi thỉnh_kinh tại Tây_Trúc
12  Tôn_Ngộ_Không , Đường Tăng đi thỉnh_kinh tại Tây_Trúc
13  Tôn_Ngộ_Không tháp_tùng Đường Tăng đi thỉnh_kinh tại Tây_Trúc

```

Ta thấy các từ ở vị trí <mask> tìm được khá tự nhiên và ngữ nghĩa của câu không khác mấy so với con người tạo ra. Sự chuẩn xác và tự nhiên có được dựa trên quá trình huấn luyện mô hình pretrain trên một bộ dữ liệu có kích thước rất lớn (khoảng 20GB).

7. Trích xuất đặc trưng (Extract feature) cho từ

Sau khi load được model BERT, chúng ta hoàn toàn có thể trích xuất đặc trưng cho một từ bất kỳ từ pretrain model. Từ các véc tơ nhúng được trích xuất cho một từ hoặc một câu, chúng ta có thể đo lường similarity để tìm ra các câu tương đồng về nội dung hoặc các từ đồng nghĩa. Một ứng dụng khác đó là chúng ta có thể tận dụng các biểu diễn ngữ nghĩa của từ thông qua véc tơ embedding được huấn luyện từ BERT để chuyển giao sang các tác vụ phân loại văn bản, phân tích cảm xúc bình luận. Phương pháp tiếp cận sẽ tương tự như áp dụng các model GloVe, word2vec, fasttext trong học nông (shallow learning).

Các véc tơ embedding cho từng từ trong câu từ mô hình BERT được trích xuất như sau:

Top

```

1  from fairseq.data.encoders.fastbpe import fastBPE
2
3  # Khởi tạo Byte Pair Encoding cho PhoBERT
4  class BPE():
5      bpe_codes = 'PhoBERT_base_fairseq/bpe.codes'
6
7  args = BPE()
8  phoBERT.bpe = fastBPE(args) #Incorporate the BPE encoder into PhoBERT
9  doc = phoBERT.extract_features_aligned_to_words('học sinh cấp 3 được (
10
11  for tok in doc:
12      print('{:10}{} (...) {}'.format(str(tok), tok.vector[:5], tok.vec

```

```

1  <s>      tensor([ 0.0534,  0.1301, -0.0475, -0.8371,  0.3862], grad_
2  học_sinh tensor([ 0.1764,  0.1603,  0.0792, -0.6043, -0.3138], grad_
3  cấp      tensor([ 0.0679,  0.0194,  0.3450, -0.4951, -0.6394], grad_
4  3        tensor([-0.0465, -0.3846,  0.1337, -1.1276,  0.1910], grad_
5  được     tensor([ 0.1920, -0.0146,  0.2933,  0.0086,  0.0690], grad_
6  đến      tensor([-0.0108, -0.6463, -0.2906, -0.0317,  0.0561], grad_
7  trường   tensor([-0.0270,  0.2676,  0.3856,  0.3514,  0.1169], grad_
8  sau      tensor([-0.1175,  0.4808,  0.0772, -0.2991,  0.0147], grad_
9  nghỉ     tensor([ 0.4385,  0.4162,  0.1529, -0.1419, -0.1928], grad_
10 dịch     tensor([ 0.2958, -0.0976,  0.2024, -0.9278,  0.0270], grad_
11 covid    tensor([ 0.1539,  0.2343,  0.4054, -1.6919, -0.7180], grad_
12 </s>     tensor([ 0.0558,  0.0341, -0.0286, -0.6476,  0.4656], grad_

```

Khi đó mỗi từ sẽ được biểu diễn bằng 768 chiều là số chiều của hidden véc tơ trong mô hình BERT base.

8. Bài toán classification

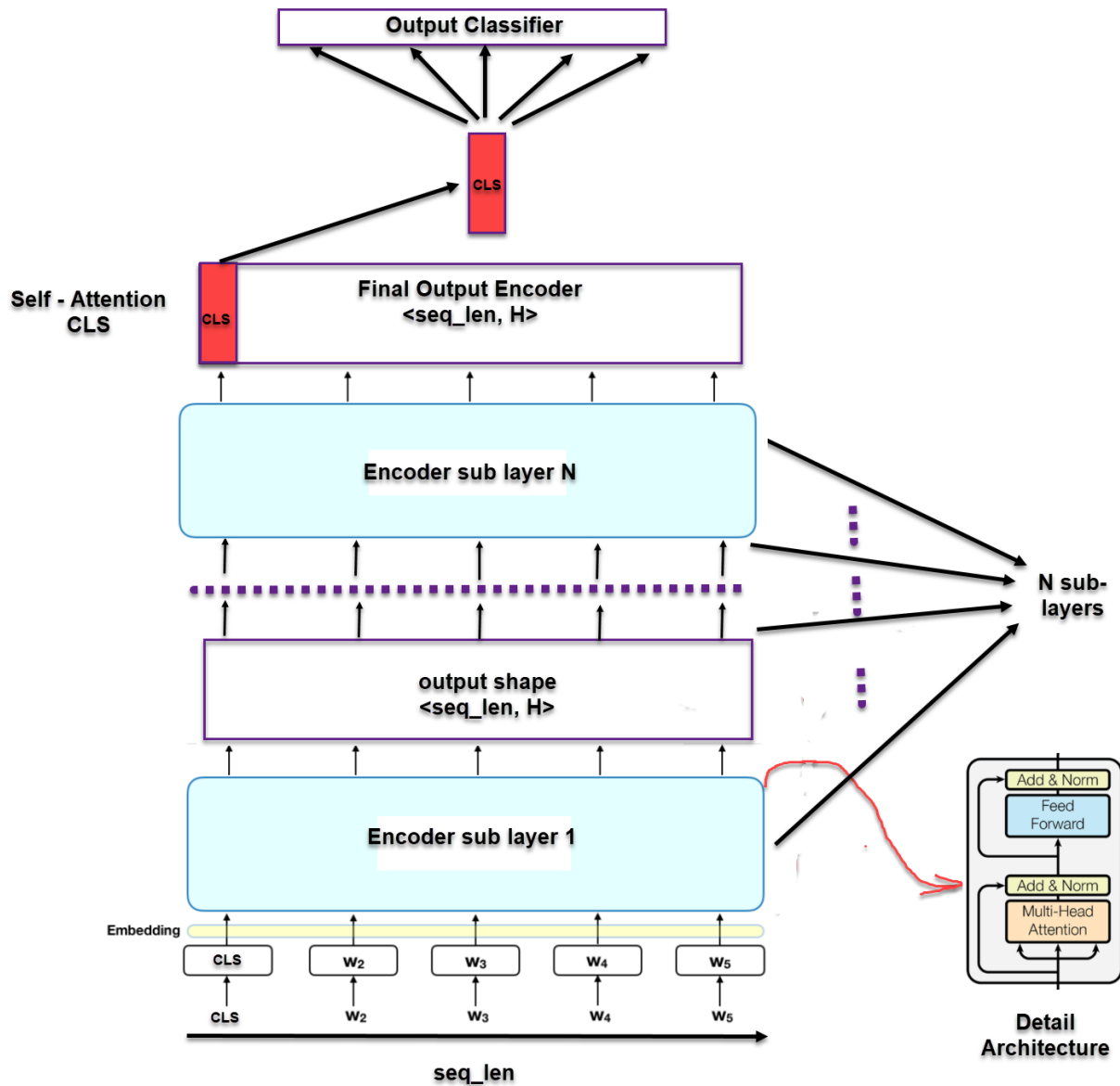
8.1. Kiến trúc mô hình

Ý tưởng fine-tuning được lấy từ bài báo [How to Fine-Tune BERT for Text Classification?](https://arxiv.org/abs/1905.05583) (https://arxiv.org/abs/1905.05583).

Model BERT base sẽ tạo ra một kiến trúc gồm 12 sub-layers ở encoder, 12 heads trong multi-head attention trên mỗi sub-layer. Output là tập hợp các véc tơ self-attention bằng chiều dài của input. Mỗi véc tơ có kích thước là 768.

Để fine-tuning lại kiến trúc của BERT cho tác vụ phân loại văn bản (text classification). Chúng ta truncate decoder của BERT, giữ nguyên kiến trúc encoder của transformer và sau đó trích xuất ra biểu diễn véc tơ của token CLS đánh dấu vị trí đầu tiên. Véc tơ này sẽ được sử dụng làm đầu vào cho thuật toán classifier bằng cách thêm một linear projection layer (cũng chính là fully connected layer) ở cuối có kích thước bằng với số classes cần phân loại. Cụ thể hơn chúng ta cùng xem kiến trúc bên dưới.

Top



Hình 1: Kiến trúc fine-tuning classifier của BERT trong classification. Biểu diễn self-attention của token tại vị trí CLS được sử dụng làm input cho thuật toán phân loại. Chúng ta thêm một linear projection layer ở cuối cùng để tính toán phân phối xác suất.

Để lấy ví dụ cho quá trình fine-tuning lại PhoBERT cho tác vụ phân loại văn bản mình sẽ huấn luyện model phân loại topics báo chí. Chúng ta sẽ tìm hiểu về bộ dữ liệu cho mô hình.

8.2. Dữ liệu

Dữ liệu mà mình sử dụng là VNTC (<https://github.com/duyvuoleo/VNVC.git>) với các bài báo đã được sắp xếp theo 10 topics. Bộ dữ liệu bao gồm 33 nghìn bài báo trên tập train và 50 nghìn bài báo trên tập test có phân bố số lượng theo topics như sau:

```

***Train***
Topic   Topic ID   #files
*****
Chinh tri Xa hoi      XH      5219
Doi song             DS      3159
Khoa hoc             KH      1820
Kinh doanh           KD      2552
Phap luat            PL      3868
Suc khoe             SK      3384
The gioi             TG      2898
The thao             TT      5298
Van hoa              VH      3080
Vi tinh              VT      2481

Total                33759

***Test***
Chinh tri Xa hoi      XH      7567
Doi song             DS      2036
Khoa hoc             KH      2096
Kinh doanh           KD      5276
Phap luat            PL      3788
Suc khoe             SK      5417
The gioi             TG      6716
The thao             TT      6667
Van hoa              VH      6250
Vi tinh              VT      4560

Total                50373

*** Notes:
    * For Train & Test
      + DS_VNE_(...) : VnExpress news agency (http://vnexpress.net/)
      + DS_TT_(...): Youth news agency (http://tuoitre.vn/)
      + DS_TN_(...): Thanh Nien news agency (http://thanhnien.vn/)
      + DS_NLD_(...): Nguoi Lao Dong news agency (http://nld.com.vn/)

```

Dữ liệu sau xử lý được mình chia sẻ. Nếu không muốn tìm hiểu quá trình tạo dữ liệu, bạn đọc có thể chuyển qua mục `Tokenize Input` và `output` và bỏ qua bước này.

8.2.1. Đọc và lưu dữ liệu

- 1 `!git clone https://github.com/duyvuleo/VNTC.git`
- 2 `!ls VNTC/Data/10Topics/Ver1.1`

Sau khi đã download dữ liệu về, chúng ta sẽ đọc và lưu các bài báo vào những list chứa nội dung và nhãn tương ứng theo 2 folders `train` và `test`.

Top

```

1  import glob2
2  from tqdm import tqdm
3
4  train_path = 'Train_Full/**/*.txt'
5  test_path = 'Test_Full/**/*.txt'
6
7  # Hàm đọc file txt
8  def read_txt(path):
9      with open(path, 'r', encoding='utf-16') as f:
10         data = f.read()
11         return data
12
13  # Hàm tạo dữ liệu huấn luyện cho tập train và test
14  def make_data(path):
15      texts = []
16      labels = []
17      for file_path in tqdm(glob2.glob(train_path)):
18          try:
19              content = read_txt(file_path)
20              label = file_path.split('/')[1]
21              texts.append(content)
22              labels.append(label)
23          except:
24              next
25      return texts, labels
26
27  text_train, label_train = make_data(train_path)
28  text_test, label_test = make_data(test_path)

```

Quá trình đọc files sẽ tốn khá nhiều thời gian. Do đó các bạn có thể tạo các hàm lưu trữ lại các list nội dung và nhãn và load lại cho lượt huấn luyện sau.

```

1  import pickle
2
3  def _save_pkl(path, obj):
4      with open(path, 'wb') as f:
5          pickle.dump(obj, f)
6
7  def _load_pkl(path):
8      with open(path, 'rb') as f:
9          obj = pickle.load(f)
10         return obj
11
12  # Lưu lại các files
13  _save_pkl('text_train.pkl', text_train)
14  _save_pkl('label_train.pkl', label_train)
15  _save_pkl('text_test.pkl', text_test)
16  _save_pkl('label_test.pkl', label_test)

```



```

1  print('text content:\n', text_train[0])
2  print('label:\n', label_train[0])

```

Top


```

1 text content:
2   Tấm hít nhỏ xinh
3   Tủ lạnh hay phía tường trước bàn làm việc của bạn sẽ đẹp hơn nếu có nó
4   Chuẩn bị: nam châm dày 3 mm; gỗ mỏng; sơn; keo dán gỗ; cọ, cửa.
5   Thực hiện:
6   Bước 1: Cưa 3 mảnh gỗ vuông làm nền, diện tích 4 cm2. Bạn có thể thay
7   Bước 2: Dùng keo dán những mảnh gỗ nhỏ vào mảnh gỗ nền. Cần chú ý phớt
8   Chú ý: Chọn loại gỗ thật mỏng, nếu không sản phẩm trông rất thô. Không
9
10
11 label:
12   Doi song

```

8.2.2. Tokenize nội dung

Tiếp theo ta sẽ tokenize các câu văn sang chuỗi index và padding câu văn về cũng một độ dài.

```

1 max_sequence_length = 500
2
3 def convert_lines(lines, vocab, bpe):
4     '''
5     lines: list các văn bản input
6     vocab: từ điển dùng để encoding subwords
7     bpe:
8     '''
9     # Khởi tạo ma trận output
10    outputs = np.zeros((len(lines), max_sequence_length)) # --> shape (l
11    # Index của các token cls (đầu câu), eos (cuối câu), padding (padding)
12    cls_id = 0
13    eos_id = 2
14    pad_id = 1
15
16    for idx, row in tqdm(enumerate(lines), total=len(lines)):
17        # Mã hóa subwords theo byte pair encoding(bpe)
18        subwords = bpe.encode('<s> ' + row + ' </s>')
19        input_ids = vocab.encode_line(subwords, append_eos=False, add_if_n
20        # Truncate input nếu độ dài vượt quá max_seq_len
21        if len(input_ids) > max_sequence_length:
22            input_ids = input_ids[:max_sequence_length]
23            input_ids[-1] = eos_id
24        else:
25            # Padding nếu độ dài câu chưa bằng max_seq_len
26            input_ids = input_ids + [pad_id, ]*(max_sequence_length - len(i
27
28        outputs[idx, :] = np.array(input_ids)
29    return outputs

```

8.3. Tokenize Input và output

Các bạn có thể download lại dữ liệu **X, y** mà tôi đã chuẩn bị cho huấn luyện tại Dữ liệu Tokenize (https://drive.google.com/drive/folders/1stRredl0fZ2vE5_SKGggrgDxnV1bxhr1?usp=sharing) và bỏ qua bước này. Thực hiện luôn bước tiếp theo Load model BERT . Top

- Chuẩn bị X input: Tokenize nội dung các văn bản sang chuỗi indices.

- Chuẩn bị y output: Encoding các label output thành indices đánh dấu số thứ tự của văn bản.

```

1      from tqdm import tqdm
2      import torch
3
4      max_sequence_length = 256
5      def convert_lines(lines, vocab, bpe):
6          '''
7          lines: list các văn bản input
8          vocab: từ điển dùng để encoding subwords
9          bpe:
10             '''
11          # Khởi tạo ma trận output
12          outputs = np.zeros((len(lines), max_sequence_length), dtype=np.int32)
13          # Index của các token cls (đầu câu), eos (cuối câu), padding (padding)
14          cls_id = 0
15          eos_id = 2
16          pad_id = 1
17
18          for idx, row in tqdm(enumerate(lines), total=len(lines)):
19              # Mã hóa subwords theo byte pair encoding(bpe)
20              subwords = bpe.encode('<s> ' + row + ' </s>')
21              input_ids = vocab.encode_line(subwords, append_eos=False, add_if_
22              # Truncate input nếu độ dài vượt quá max_seq_len
23              if len(input_ids) > max_sequence_length:
24                  input_ids = input_ids[:max_sequence_length]
25                  input_ids[-1] = eos_id
26              else:
27                  # Padding nếu độ dài câu chưa bằng max_seq_len
28                  input_ids = input_ids + [pad_id, ]*(max_sequence_length - len(i
29
30              outputs[idx,:] = np.array(input_ids)
31          return outputs
32
33          # Load the dictionary
34          vocab = Dictionary()
35          vocab.add_from_file("PhoBERT_base_transformers/dict.txt")
36
37
38          # Test encode lines
39          lines = ['Học_sinh được nghỉ học bắt đầu từ tháng 3 để tránh dịch cov
40          [x1, x2] = convert_lines(lines, vocab, phoBERT.bpe)
41          print('x1 tensor encode: {}, shape: {}'.format(x1[:10], x1.size))
42          print('x1 tensor decode: ', phoBERT_cls.decode(torch.tensor(x1))[:103

```

Top

```

1  from tqdm import tqdm
2  import torch
3
4  max_sequence_length = 256
5  def convert_lines(lines, vocab, bpe):
6      '''
7      lines: list các văn bản input
8      vocab: từ điển dùng để encoding subwords
9      bpe:
10         '''
11      # Khởi tạo ma trận output
12      outputs = np.zeros((len(lines), max_sequence_length), dtype=np.int32)
13      # Index của các token cls (đầu câu), eos (cuối câu), padding (padding)
14      cls_id = 0
15      eos_id = 2
16      pad_id = 1
17
18      for idx, row in tqdm(enumerate(lines), total=len(lines)):
19          # Mã hóa subwords theo byte pair encoding (bpe)
20          subwords = bpe.encode('<s> ' + row + ' </s>')
21          input_ids = vocab.encode_line(subwords, append_eos=False, add_if_
22          # Truncate input nếu độ dài vượt quá max_seq_len
23          if len(input_ids) > max_sequence_length:
24              input_ids = input_ids[:max_sequence_length]
25              input_ids[-1] = eos_id
26          else:
27              # Padding nếu độ dài câu chưa bằng max_seq_len
28              input_ids = input_ids + [pad_id, ]*(max_sequence_length - len(i
29
30          outputs[idx,:] = np.array(input_ids)
31      return outputs
32
33      # Load the dictionary
34      vocab = Dictionary()
35      vocab.add_from_file("PhoBERT_base_transformers/dict.txt")
36
37
38      # Test encode lines
39      lines = ['Học sinh được nghỉ học bắt đầu từ tháng 3 để tránh dịch cov
40      [x1, x2] = convert_lines(lines, vocab, phoBERT_cls.bpe)
41      print('x1 tensor encode: {}, shape: {}'.format(x1[:10], x1.size))
42      print('x1 tensor decode: ', phoBERT_cls.decode(torch.tensor(x1))[:103]

```

Như vậy ta thấy rằng các câu văn đã được encode về token index. Từ token index có thể decode ngược trở lại thành câu input sau khi đã thêm các token đặc biệt đánh dấu vị trí bắt đầu: <s> , kết thúc: </s> câu và các vị trí nằm ngoài câu: <pad> . Ta sẽ token toàn bộ câu input sang index như sau:

```

1  X = convert_lines(text_train, vocab, phoBERT_cls.bpe)
2  print('X shape: ', X.shape)

```

Sau cùng ta thu được các chuỗi index có kích thước là 256, bằng với kích thước của các câu sau khi đã padding. Tiếp theo ta tạo output y bằng index cho các nhãn của câu.

Top

```

1  from sklearn.preprocessing import LabelEncoder
2  lb = LabelEncoder()
3  lb.fit(label_train)
4  y = lb.fit_transform(label_train)
5  print(lb.classes_)
6  print('Top 5 classes indices: ', y[:5])

```

Lưu lại dữ liệu **X** và **y**

```

1  # Save dữ liệu
2  _save_pkl('PhoBERT_pretrain/X1.pkl', X)
3  _save_pkl('PhoBERT_pretrain/y1.pkl', y)
4  _save_pkl('PhoBERT_pretrain/labelEncoder1.pkl', lb)
5
6  # Load lại dữ liệu
7  X = _load_pkl('PhoBERT_pretrain/X1.pkl')
8  y = _load_pkl('PhoBERT_pretrain/y1.pkl')
9
10 print('length of X: ', len(X))
11 print('length of y: ', len(y))

```

8.4. Load model BERT

Tiếp theo ta sẽ load pretrain model BERT từ file weight mà chúng ta đã download trước đó. Để tùy chỉnh mô hình cho phù hợp với tác vụ phân loại văn bản. Chúng ta sẽ thêm vào sau cùng pretrain model một head layer là linear projection có số units ở output bằng với số lượng classes cần phân loại và bằng 10 thông qua hàm

`phoBERT_cls.register_classification_head('new_task', num_classes=10)`. Cụ thể như sau:

```

1  # Load the model in fairseq
2  from fairseq.models.roberta import RobertaModel
3  from fairseq.data.encoders.fastbpe import fastBPE
4  from fairseq.data import Dictionary
5
6  phoBERT_cls = RobertaModel.from_pretrained('PhoBERT_base_fairseq', ch
7  phoBERT_cls.eval() # disable dropout (or leave in train mode to fine
8
9  # Load BPE
10 class BPE():
11     bpe_codes = 'PhoBERT_base_fairseq/bpe.codes'
12
13     args = BPE()
14     phoBERT_cls.bpe = fastBPE(args) #Incorporate the BPE encoder into Pho
15
16     # Add header cho classification với số lượng classes = 10
17     phoBERT_cls.register_classification_head('new_task', num_classes=10)
18     tokens = 'Học sinh được nghỉ học bắt đầu từ tháng 3 do ảnh hưởng của
19     token_idx = phoBERT_cls.encode(tokens)
20     logprobs = phoBERT_cls.predict('new_task', token_idx) # tensor([[ -1
21     logprobs

```

```

1 loading archive file PhoBERT_base_fairseq
2 | dictionary: 64000 types
3
4
5
6
7
8 tensor([[ -2.3722, -2.1128, -2.2945, -2.3484, -2.2294, -2.1600, -2.5104,
9          -2.4144, -2.2299]], grad_fn=<LogSoftmaxBackward>)
```

8.5. Huấn luyện model

Trước khi huấn luyện mô hình chúng ta sẽ xây dựng các hàm đánh giá mô hình theo 2 metric là `accuracy` và `f1_score`.

```

1 import torch
2 import numpy as np
3 from sklearn.metrics import accuracy_score, f1_score
4
5 def evaluate(logits, targets):
6     """
7     Đánh giá model sử dụng accuracy và f1 scores.
8     Args:
9         logits (B,C): torch.LongTensor. giá trị predicted logit cho c
10        targets (B): torch.LongTensor. actual target indices.
11    Returns:
12        acc (float): the accuracy score
13        f1 (float): the f1 score
14    """
15    # Tính accuracy score và f1_score
16    logits = logits.detach().cpu().numpy()
17    y_pred = np.argmax(logits, axis = 1)
18    targets = targets.detach().cpu().numpy()
19    f1 = f1_score(targets, y_pred, average='weighted')
20    acc = accuracy_score(targets, y_pred)
21    return acc, f1
22
23 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
24 logits = torch.tensor([[0.1, 0.2, 0.7],
25                        [0.4, 0.1, 0.5],
26                        [0.1, 0.2, 0.7]]).to(device)
27 targets = torch.tensor([1, 2, 2]).to(device)
28 evaluate(logits, targets)
```

```

1 (0.6666666666666666, 0.5333333333333333)
```

Top

```

1     def validate(valid_loader, model, device):
2         model.eval()
3         accs = []
4         f1s = []
5         with torch.no_grad():
6             for x_batch, y_batch in valid_loader:
7                 x_batch = x_batch.to(device)
8                 y_batch = y_batch.to(device)
9                 outputs = model.predict('new_task', x_batch)
10                logits = torch.exp(outputs)
11                acc, f1 = evaluate(logits, y_batch)
12                accs.append(acc)
13                f1s.append(f1)
14
15        mean_acc = np.mean(accs)
16        mean_f1 = np.mean(f1s)
17        return mean_acc, mean_f1

```

Hàm huấn luyện mô hình trên từng epoch.

```

1     def trainOnEpoch(train_loader, model, optimizer, epoch, num_epochs, c
2         model.train()
3         sum_epoch_loss = 0
4         sum_acc = 0
5         sum_f1 = 0
6         start = time.time()
7         for i, (x_batch, y_batch) in enumerate(train_loader):
8             x_batch = x_batch.to(device)
9             y_batch = y_batch.to(device)
10            optimizer.zero_grad()
11            y_pred = model.predict('new_task', x_batch)
12            logits = torch.exp(y_pred)
13            acc, f1 = evaluate(logits, y_batch)
14            loss = criteria(y_pred, y_batch)
15            loss.backward()
16            optimizer.step()
17
18            loss_val = loss.item()
19            sum_epoch_loss += loss_val
20            sum_acc += acc
21            sum_f1 += f1
22            iter_num = epoch * len(train_loader) + i + 1
23
24            if i % log_aggr == 0:
25                print('[TRAIN] epoch %d/%d  observation %d/%d batch loss:
26                    % (epoch + 1, num_epochs, i, len(train_loader), loss_
27                        len(x_batch) / (time.time() - start)))
28            start = time.time()

```

Quá trình huấn luyện một model classification trên pytorch sẽ bao gồm những bước chính sau đây:

- Khởi tạo DataLoader để quản lý dữ liệu đưa vào huấn luyện và thẩm định.
- Thiết lập kiến trúc mô hình.

Top

- Khai báo hàm loss function.
- Phương pháp optimization giúp tối ưu loss function.
- Huấn luyện mô hình qua các epochs.

Bên dưới chúng ta sẽ lần lượt thực hiện các bước trên.

Top

```

1      import os
2      import time
3      import random
4      import argparse
5      import pickle
6      import numpy as np
7      from tqdm import tqdm
8      from os.path import join
9
10     import torch
11     from torch import nn
12     from torch.utils.data import DataLoader
13     import torch.nn.functional as F
14     import torch.optim as optim
15     from torch.optim.lr_scheduler import StepLR
16     from torch.autograd import Variable
17     from torch.backends import cudnn
18     from sklearn.model_selection import StratifiedKFold
19
20     # Load the model in fairseq
21     from fairseq.models.roberta import RobertaModel
22     from fairseq.data.encoders.fastbpe import fastBPE
23     from fairseq.data import Dictionary
24     from transformers.modeling_utils import *
25     from transformers import *
26
27     # Khởi tạo argument
28     EPOCHS = 20
29     BATCH_SIZE = 6
30     ACCUMULATION_STEPS = 5
31     FOLD = 4
32     LR = 0.0001
33     LR_DC_STEP = 80
34     LR_DC = 0.1
35     CUR_DIR = os.path.dirname(os.getcwd())
36     DEVICE = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
37     FOLD = 4
38     CKPT_PATH2 = 'model_ckpt2'
39
40     if not os.path.exists(CKPT_PATH2):
41         os.mkdir(CKPT_PATH2)
42
43     # Khởi tạo DataLoader
44     splits = list(StratifiedKFold(n_splits=5, shuffle=True, random_state=
45
46     for fold, (train_idx, val_idx) in enumerate(splits):
47         best_score = 0
48         if fold != FOLD:
49             continue
50         print("Training for fold {}".format(fold))
51
52         # Create dataset
53         train_dataset = torch.utils.data.TensorDataset(torch.tensor(X[train_idx, :]))
54         valid_dataset = torch.utils.data.TensorDataset(torch.tensor(X[val_idx, :]))
55
56         # Create DataLoader
57         train_loader = torch.utils.data.DataLoader(train_dataset, batch_size=

```

Top


```

58     valid_loader = torch.utils.data.DataLoader(valid_dataset, batch_
59
60     # Khởi tạo model:
61     MODEL_LAST_CKPT = os.path.join(CKPT_PATH2, 'latest_checkpoint.pt)
62     if os.path.exists(MODEL_LAST_CKPT):
63         print('Load checkpoint model!')
64         phoBERT_cls = torch.load(MODEL_LAST_CKPT)
65     else:
66         print('Load model pretrained!')
67         # Load the model in fairseq
68         from fairseq.models.roberta import RobertaModel
69         from fairseq.data.encoders.fastbpe import fastBPE
70         from fairseq.data import Dictionary
71
72         phoBERT_cls = RobertaModel.from_pretrained('PhoBERT_base_fairseq
73         phoBERT_cls.eval() # disable dropout (or leave in train mode
74
75         ## Load BPE
76         # class BPE():
77         #     bpe_codes = 'PhoBERT_base_fairseq/bpe.codes'
78
79         # args = BPE()
80         # phoBERT_cls.bpe = fastBPE(args) #Incorporate the BPE encoder
81
82         # Add header cho classification với số lượng classes = 10
83         phoBERT_cls.register_classification_head('new_task', num_class
84
85     ## Load BPE
86     print('Load BPE')
87     class BPE():
88         bpe_codes = 'PhoBERT_base_fairseq/bpe.codes'
89
90     args = BPE()
91     phoBERT_cls.bpe = fastBPE(args) #Incorporate the BPE encoder into
92     phoBERT_cls.to(DEVICE)
93
94     # Khởi tạo optimizer và scheduler, criteria
95     print('Init Optimizer, scheduler, criteria')
96     param_optimizer = list(phoBERT_cls.named_parameters())
97     no_decay = ['bias', 'LayerNorm.bias', 'LayerNorm.weight']
98     optimizer_grouped_parameters = [
99         {'params': [p for n, p in param_optimizer if not any(nd in n
100         {'params': [p for n, p in param_optimizer if any(nd in n for
101     ]
102
103     num_train_optimization_steps = int(EPOCHS*len(train_dataset)/BATCH
104     optimizer = AdamW(optimizer_grouped_parameters, lr=LR, correct_b
105     scheduler = get_linear_schedule_with_warmup(optimizer, num_warmup
106     scheduler0 = get_constant_schedule(optimizer) # scheduler với h
107     # optimizer = optim.Adam(phoBERT_cls.parameters(), LR)
108     criteria = nn.NLLLoss()
109     # scheduler = StepLR(optimizer, step_size = LR_DC_STEP, gamma =
110     avg_loss = 0.
111     avg_accuracy = 0.
112     frozen = True
113     for epoch in tqdm(range(EPOCHS)):
114         # warm up tại epoch đầu tiên, sau epoch đầu sẽ phá băng các
115         if epoch > 0 and frozen:

```

Top

```

116         for child in phoBERT_cls.children():
117             for param in child.parameters():
118                 param.requires_grad = True
119         frozen = False
120         del scheduler0
121         torch.cuda.empty_cache()
122         # Train model on EPOCH
123         print('Epoch: ', epoch)
124         trainOnEpoch(train_loader=train_loader, model=phoBERT_cls, optimizer=optimizer, scheduler=scheduler)
125         # scheduler.step(epoch = epoch)
126         # Phá băng layers sau epoch đầu tiên
127         if not frozen:
128             scheduler.step()
129         else:
130             scheduler0.step()
131         optimizer.zero_grad()
132         # Validate on validation set
133         acc, f1 = validate(valid_loader, phoBERT_cls, device=DEVICE)
134         print('Epoch {} validation: acc: {:.4f}, f1: {:.4f} \n'.format(epoch, acc, f1))
135
136         # Store best model checkpoint
137         ckpt_dict = {
138             'epoch': epoch + 1,
139             'state_dict': phoBERT_cls.state_dict(),
140             'optimizer': optimizer.state_dict()
141         }
142         # Save model checkpoint into 'latest_checkpoint.pth.tar'
143         torch.save(ckpt_dict, MODEL_LAST_CKPT)

```

```

1      0%|          | 0/20 [00:00<?, ?it/s] [A
2
3      Init Optimizer, scheduler, criteria
4      ...
5      [TRAIN] epoch 4/20  observation 2700/4502 batch loss: 2.2600 (avg 2.25:

```

Thời gian huấn luyện sẽ khá lâu, các bạn nên kiên nhẫn chờ đợi. Mình chỉ dừng ở epochs số 4 cho mục đích demo. Mặc dù kết quả chỉ đạt 16% accuracy nhưng accuracy luôn tăng. Learning rate của mô hình cũng được thiết lập khá nhỏ để tránh nhảy khỏi điểm tối ưu toàn cục (global optimal value). Ngoài cách fine tuning model từ fairseq như trên, các bạn có thể tham khảo thêm một cách khác của PhoBERT-Sentiment-Classification Khoi Nguyen (<https://github.com/suicao/PhoBert-Sentiment-Classification/>) thực hiện fine tuning dựa trên pretrain model huấn luyện từ transformers.

9. Huấn luyện RoBERTa trên dữ liệu của bạn

Ngoài ra chúng ta có thể tự huấn luyện pretrain model BERT dựa trên kiến trúc RoBERTa theo hướng dẫn tại RoBERTa - README (<https://github.com/pytorch/fairseq/blob/master/examples/roberta/README.pretraining.md>). Việc thực hiện khá đơn giản, bài đã khá dài nên mình gửi link cho bạn đọc tự nghiên cứu.

10. Tổng kết

Top

Như vậy mình đã giới thiệu với các bạn rất nhiều các ứng dụng khác nhau trong việc áp dụng các model pretrain RoBERTa. Trong đó có các tác vụ chính như: điền từ (filling mask), suggest search, phân loại văn bản, phân loại cảm xúc và tìm từ đồng nghĩa, trái nghĩa.

Một số tác vụ khác của NLP có thể fine tuning được từ mô hình RoBERTa như ứng dụng hỏi đáp, sinh văn bản ngẫu nhiên, sinh văn bản có nội dung tương tự nhưng đòi hỏi phải có những bộ dữ liệu chuyên biệt cho các tác vụ này. Nhưng dữ liệu như vậy cho Tiếng Việt đang rất thiếu và hiếm. Bạn đọc cũng có thể thực hành trên các bộ dữ liệu của Tiếng Anh theo hướng dẫn tại mục fine tuning - RoBERTa (<https://github.com/pytorch/fairseq/tree/master/examples/roberta>).

Như vậy qua bài viết này các bạn đã nắm bắt được các phương tiện và công cụ mới trong việc tiếp cận các bài toán của NLP. Đây là những phương pháp rất mạnh và hứa hẹn sẽ mang lại nhiều cải thiện đáng kể cho các tác vụ NLP của bạn. Đừng quên like và share bài viết này nếu bạn cảm thấy kiến thức mình chia sẻ là hữu ích với bạn.

11. Tài liệu

1. fairseq (<https://github.com/pytorch/fairseq>)
2. RoBERTa (<https://github.com/pytorch/fairseq/tree/master/examples/roberta>)
3. transformers (<https://github.com/huggingface/transformers>)
4. PhoBERT-Sentiment-Classification (<https://github.com/suicao/PhoBert-Sentiment-Classification/>)
5. PhoBERT (<https://github.com/VinAIResearch/PhoBERT>)

Top