

Bài 7 - Pytorch - Buổi 2 - Seq2seq model correct spelling

19 Aug 2019 - phamdinhhkhanh

Menu

- 1. Giới thiệu chung
- 2. Chuẩn bị dữ liệu
 - 2.1. Giải nén dữ liệu
 - 1.2. Load and Trim data
 - 1.2. Chuẩn hóa dữ liệu
 - 1.2.1 Tạo ngram
 - 1.3. Tạo batch huấn luyện
- 2. Mô hình seq2seq sửa lỗi chính tả
 - 2.1. Encoder
 - 2.2. Decoder
 - 2.2.1. Áp dụng Attention layer
 - 2.2.2. Quá trình Decoder
- 3. Huấn luyện model.
 - 3.1 Loss function
 - 3.2. Huấn luyện vòng lặp đơn
 - 3.3. Huấn luyện vòng lặp
- 4. Xác định phương pháp đánh giá
- 5. Huấn luyện model
 - 5.1. Huấn luyện model
 - 5.2 Đánh giá model
- 6. Hạn chế của mô hình
- 7. Tài liệu tham khảo

1. Giới thiệu chung

Cùng với sự phát triển của deep learning nói chung. Ngày nay lớp các mô hình seq2seq càng tỏ ra hiệu quả trong nhiều tác vụ khác nhau như dịch máy, sửa lỗi chính tả, image captioning, recommendation, dự báo chuỗi thời gian,... Nhờ sự phát triển của các kiến trúc mạng RNN hiện đại kèm theo các kĩ thuật learning hiệu quả như sử dụng thêm kiến trúc attention layer, các phương pháp cải thiện accuracy như `teach_forcing, beam search` mà mô hình dịch máy ngày càng đạt độ chính xác cao. Các tài liệu về các phương pháp trên đã có nhiều tuy nhiên đa phần là lý thuyết và khá khó cho người mới bắt đầu tiếp cận. Bài viết này nhằm mục đích tạo ra một bản diễn giải kèm thực hành về các bước xây dựng mô hình seq2seq ứng dụng trong sửa lỗi chính tả. Để hiểu được các nội dung trong bài yêu cầu bạn đọc có kiến thức nền tảng về pytorch (<https://phamdinhhkhanh.github.io/2019/08/10/PytorchTutorial1.html>), nắm vững lý thuyết về mạng LSTM

(https://phamdinhhkhanh.github.io/2019/04/22/L%C3%BD_thuy%E1%BA%BFT_v%E1%BB%81_m%E1%BA%A1ng_LSTM.html).

Ngoài ra bạn đọc cần có sẵn máy tính cài pytorch (<https://pytorch.org/>) hoặc các VM hỗ trợ pytorch. Để tiện cho thực hành tôi khuyến nghị bạn đọc sử dụng google colab (<https://colab.research.google.com>) miễn phí và cài sẵn các deep learning framework cơ bản như tensorflow, pytorch, keras,....

Ngoài ra bài viết được xây dựng và tổng hợp dựa trên nhiều nguồn tài liệu khác nhau. Đặc biệt là từ trang hướng dẫn thực hành pytorch (https://pytorch.org/tutorials/beginner/chatbot_tutorial.html), từ ý tưởng được chia sẻ nằm trong top của cuộc thi thêm dấu cho tiếng việt - aivivn 1st (<https://forum.machinelearningcoban.com/t/aivivn-3-vietnamese-tone-prediction-1st-place-solution/5721>), thêm dấu cho tiếng việt - aivivn 2nd (<https://forum.machinelearningcoban.com/t/aivivn-3-vietnamese-tone-prediction-2nd-place-solution/5759>).

Để huấn luyện mô hình trên google colab chúng ta cần mount folder lưu trữ dữ liệu trên google drive để có thể access dữ liệu dễ dàng. Bên dưới là câu lệnh thực hiện mount dữ liệu.

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

Top

Trong đó `your_data_folder_link` là đường link tới folder chứa dữ liệu. Lưu ý link `gdrive/My Drive` là đường trở mặc định để đi vào thư mục `My Drive` của bạn.

2. Chuẩn bị dữ liệu

2.1. Giải nén dữ liệu

Ở bước này chúng ta sẽ sử dụng đầu vào là dữ liệu của cuộc thi thêm dấu từ tiếng việt tại aivivn (<https://www.aivivn.com/contests/3>). Bạn đọc có thể tải về bộ dữ liệu tại google drive (https://drive.google.com/file/d/1m_5CDQQSavev5zWb8JUq97_zUTnOcVvS/view) dưới dạng zip file. Để giải nén dữ liệu ta cần extract bằng hàm `extract` bên dưới. Nhớ cài đặt package `zipfile` trước khi chạy lệnh.

```
1 # Extract zip file
2 import zipfile
3 import os
4
5 def _extract_zip_file(fn_zip, fn):
6     with zipfile.ZipFile(fn, 'r') as zip_ref:
7         zip_ref.extractall(fn)
8
9 _extract_zip_file(fn_zip = 'vietnamese_tone_prediction.zip', fn = 'vietnamse_tone_predic
10
11 print(os.path.exists('vietnamse_tone_prediction'))
```

Cùng lấy ra 500000 dòng đầu tiên của file `train.txt` trong folder giải nén làm tập huấn luyện.

```
1 def _data_train(fn):
2     with open(fn, 'r') as fn:
3         train = fn.readlines()
4         train = [item[:-1] for item in train[:500000]]
5         return train
6
7 train = _data_train(fn = 'vietnamse_tone_prediction/train.txt')
8 print('length of train: {}'.format(len(train)))
9 train[:5]
```

```
1 length of train: 500000
2
3
4
5
6
7 ['Bộ phim lần đầu được công chiếu tại liên hoan phim Rome 2007 và sau đó được chiếu ở Fa
8 'Những kiểu áo sơ mi may theo chất liệu cotton, KT, hay có chút co giãn năm nay cũng đ
9 'Đương kim tổng thống là Andrés Manuel López Obrador, người nhậm chức vào ngày 1 tháng
10 'Centaurea gloriosa là một loài thực vật có hoa trong họ Cúc.',
11 'Sau này mới thấy người ta nói, đó là con rắn đực đi tìm người ăn thịt để trả thù cho r
```

Bên dưới là một số hàm tiện ích để chuyển các từ có dấu của tiếng việt sang không dấu. Mục đích là để tạo ra các cặp (câu không dấu, câu có dấu) từ câu có dấu.

[illegible]

Để ý thấy hầu hết các dấu câu sẽ dính liền với câu. Chẳng hạn câu tại [khanh blog](#), tôi lưu lại các bài viết như một quyển nhật kí . Thì từ blog và dấu phẩy bị dính liền. Vì vậy ta cần sử dụng hàm `normalizeString()` để tách các dấu ra khỏi từ.

```

1      # Tách dấu ra khỏi từ
2      def normalizeString(s):
3          # Tách dấu câu nếu kí tự liền nhau
4          marks = '!.?,-${}()''
5          r = "([+\"\\\".join(marks)+])"
6          s = re.sub(r, r" \1 ", s)
7          # Thay thế nhiều spaces bằng 1 space.
8          s = re.sub(r"\s+", r" ", s).strip()
9          return s
10
11      normalizeString('vui vẻ, hòa đồng, hoạt bát')

```

Sử dụng 2 hàm số trên để tạo ra các list `train` gồm các câu có dấu và `train_rev_accent` gồm các câu không dấu.

```
1 import itertools
2 train = [normalizeString(item) for item in train]
3 train_rev_accent = [remove_tone_line(item) for item in train]
4
5 print('train top 5:', train[:5])
6 print('train_rev_accent top 5:', train_rev_accent[:5])
```



```
1 train top 5: ['Bộ phim lần đầu được công chiếu tại liên hoan phim Rome 2007 và sau đó đượ
2 train rev accent top 5: ['Bo phim lan dau duoc cong chieu tai lien hoan phim Rome 2007 va
```

1.2. Load and Trim data

Bên dưới ta sẽ xây dựng class Voc để xây dựng từ điển cho tập dữ liệu. Các hàm trong Voc có tác dụng:

- `addword()` : Kiểm tra xem từ đã xuất hiện trong từ điển chưa. Nếu chưa sẽ thêm từ mới đó vào từ điển.
- `addSentence()` : Truyền vào 1 câu và thêm các từ trong câu vào từ điển.
- `trim()` : Loại bỏ các từ hiếm trong từ điển nếu nó có ngưỡng số tần xuất xuất hiện trong toàn bộ tập dữ liệu **Top** ít hơn `min_count`.

```

1      # Default word tokens
2      PAD_token = 0 # Used for padding short sentences
3      SOS_token = 1 # Start-of-sentence token
4      EOS_token = 2 # End-of-sentence token
5
6
7      class Voc:
8          def __init__(self, name):
9              self.name = name
10             self.trimmed = False
11             self.word2index = {}
12             self.word2count = {}
13             self.index2word = {PAD_token: "PAD", SOS_token: "SOS", EOS_token: "EOS"}
14             self.num_words = 3 # Count SOS, EOS, PAD
15
16         def addSentence(self, sentence):
17             for word in sentence.split(' '):
18                 self.addWord(word)
19
20         def addWord(self, word):
21             if word not in self.word2index:
22                 self.word2index[word] = self.num_words
23                 self.word2count[word] = 1
24                 self.index2word[self.num_words] = word
25                 self.num_words += 1
26             else:
27                 self.word2count[word] += 1
28
29         # Remove words below a certain count threshold
30         def trim(self, min_count):
31             if self.trimmed:
32                 return
33             self.trimmed = True
34
35             keep_words = []
36
37             for k, v in self.word2count.items():
38                 if v >= min_count:
39                     keep_words.append(k)
40
41             print('keep_words {} / {} = {:.4f}'.format(
42                 len(keep_words), len(self.word2index), len(keep_words) / len(self.word2index)
43             ))
44
45         # Reinitialize dictionaries
46         self.word2index = {}
47         self.word2count = {}
48         self.index2word = {PAD_token: "PAD", SOS_token: "SOS", EOS_token: "EOS"}
49         self.num_words = 3 # Count default tokens
50
51         for word in keep_words:
52             self.addWord(word)

```

1.2. Chuẩn hóa dữ liệu

1.2.1 Tạo ngram

Trong dữ liệu có những câu rất dài làm giảm hiệu quả dự báo của mô hình. Do đó chúng ta sẽ tìm cách tạo ra các ngram với độ dài đo bằng số lượng từ cố định để dự báo trong một khoảng cách ngắn. Ở bài này ta sẽ lựa chọn ngram = 4.

```

1  def _ngram(text, length = 4):
2      words = text.split()
3      grams = []
4      if len(words) <= length:
5          words = words + ["PAD"]*(length-len(words))
6          return ' '.join(words)
7      else:
8          for i in range(len(words)-length+1):
9              grams.append(' '.join(words[i:i+length]))
10         return grams
11
12     print(_ngram('mùa đông năm nay không còn lạnh nữa. Vì đã có gấu 37 độ ẩm'))
13     print(_ngram('mùa đông'))

```

```

1  ['mùa đông năm nay', 'đông năm nay không', 'năm nay không còn', 'nay không còn lạnh', 'kh
2  ['mùa đông PAD PAD']

```

```

1  import itertools
2
3  train_grams = list(itertools.chain.from_iterable([_ngram(item) for item in train]))
4  train_rev_acc_grams = list(itertools.chain.from_iterable([_ngram(item) for item in train_

```

```

1  corpus = list(zip(train_rev_acc_grams, train_grams))
2  corpus[:5]

```

```

1  [('Bo phim lan dau', 'Bộ phim lần đầu'),
2   ('phim lan dau duoc', 'phim lần đầu được'),
3   ('lan dau duoc cong', 'lần đầu được công'),
4   ('dau duoc cong chieu', 'đầu được công chiếu'),
5   ('duoc cong chieu tai', 'được công chiếu tại')]

```

```

1  import unicodedata
2
3  MAX_LENGTH = 4 # Maximum sentence length to consider
4
5  # Turn a Unicode string to plain ASCII, thanks to
6  # https://stackoverflow.com/a/518232/2809427
7  def unicodeToAscii(s):
8      return ''.join(
9          c for c in unicodedata.normalize('NFD', s)
10         if unicodedata.category(c) != 'Mn'
11     )
12
13 # Lowercase, trim, and remove non-letter characters
14 def normalizeString(s):
15     # Tách dấu câu nếu kí tự liền nhau
16     s = re.sub(r"([.!?,\-\&\(\)\[\]])", r" \1 ", s)
17     # Thay thế nhiều spaces bằng 1 space.
18     s = re.sub(r"\s+", r" ", s).strip()
19     return s
20
21 # Read query/response pairs and return a voc object
22 def readVocs(lines, corpus_name = 'corpus'):
23     # Split every line into pairs and normalize
24     pairs = [[normalizeString(str(s)) for s in l] for l in lines]
25     voc = Voc(corpus_name)
26     return voc, pairs
27
28 voc, pairs = readVocs(corpus)
29
30 # Returns True iff both sentences in a pair 'p' are under the MAX_LENGTH threshold
31 def filterPair(p):
32     # Input sequences need to preserve the last word for EOS token
33     return len(p[0].split(' ')) < MAX_LENGTH and len(p[1].split(' ')) < MAX_LENGTH
34
35 # Filter pairs using filterPair condition
36 def filterPairs(pairs):
37     return [pair for pair in pairs if filterPair(pair)]
38
39 # # Using the functions defined above, return a populated voc object and pairs list
40 def loadPrepareData(voc, pairs):
41     print("Read {!s} sentence pairs".format(len(pairs)))
42     # pairs = filterPairs(pairs)
43     print("Trimmed to {!s} sentence pairs".format(len(pairs)))
44     print("Counting words...")
45     for pair in pairs:
46         voc.addSentence(pair[0])
47         voc.addSentence(pair[1])
48     print("Counted words:", voc.num_words)
49     return voc, pairs
50
51 # Load/Assemble voc and pairs
52 save_dir = os.path.join("data", "save")
53 voc, pairs = loadPrepareData(voc, pairs)
54 # Print some pairs to validate
55 print("\npairs:")
56 for pair in pairs[:10]:
57     print(pair)

```

```

1 Read 9771412 sentence pairs
2 Trimmed to 9771412 sentence pairs
3 Counting words...
4 Counted words: 202153
5
6 pairs:
7 ['Bo phim lan dau', 'Bộ phim lần đầu']
8 ['phim lan dau duoc', 'phim lần đầu được']
9 ['lan dau duoc cong', 'lần đầu được công']
10 ['dau duoc cong chieu', 'đầu được công chiếu']
11 ['duoc cong chieu tai', 'được công chiếu tại']
12 ['cong chieu tai lien', 'công chiếu tại liên']
13 ['chieu tai lien hoan', 'chiếu tại liên hoan']
14 ['tai lien hoan phim', 'tại liên hoan phim']
15 ['lien hoan phim Rome', 'liên hoan phim Rome']
16 ['hoan phim Rome 2007', 'hoan phim Rome 2007']

```

Model sẽ huấn luyện nhanh hơn nếu:

- Giảm bớt số lượng các token không quá phổ biến.
- Loại bỏ bớt các câu không phổ biến.

Bên dưới ta sẽ xây dựng hàm số loại bỏ các token có tần suất nhỏ hơn ngưỡng MIN_COUNT và loại bỏ các câu có chứa token vừa bị loại bỏ.

```

1 MIN_COUNT = 3 # Minimum word count threshold for trimming
2
3 def trimRareWords(voc, pairs, MIN_COUNT):
4     # Trim words used under the MIN_COUNT from the voc
5     voc.trim(MIN_COUNT)
6     # Filter out pairs with trimmed words
7     keep_pairs = []
8     for pair in pairs:
9         input_sentence = pair[0]
10        output_sentence = pair[1]
11        keep_input = True
12        keep_output = True
13        # Check input sentence
14        for word in input_sentence.split(' '):
15            if word not in voc.word2index:
16                keep_input = False
17                break
18        # Check output sentence
19        for word in output_sentence.split(' '):
20            if word not in voc.word2index:
21                keep_output = False
22                break
23
24        # Only keep pairs that do not contain trimmed word(s) in their input or output s
25        if keep_input and keep_output:
26            keep_pairs.append(pair)
27
28        print("Trimmed from {} pairs to {}, {:.4f} of total".format(len(pairs), len(keep_pai
29        return keep_pairs
30
31
32 # Trim voc and pairs
33 pairs = trimRareWords(voc, pairs, MIN_COUNT)

```

```

1 keep_words 171772 / 202150 = 0.8497
2 Trimmed from 9771412 pairs to 9741582, 0.9969 of total

```

Như vậy sau khi loại bỏ các từ hiếm với tần xuất xuất hiện ≤ 3 thì dữ liệu còn lại 87% số lượng các token và 99.81% các câu được giữ lại.

1.3. Tạo batch huấn luyện

Top

```

1  print('EOS_token: ', EOS_token)
2  print('SOS_token: ', SOS_token)
3  print('PAD_token: ', PAD_token)

```

```

1  EOS_token:  2
2  SOS_token:  1
3  PAD_token:  0

```

Bên dưới ta sẽ xây dựng các hàm chức năng, trong đó:

- `indexesFromSentence()` : Mã hóa câu văn thành chuỗi index của các token theo giá trị của cặp `word2index` trong từ điển và đính thêm token EOS ở cuối để đánh dấu kết thúc câu. Chẳng hạn trong từ điển từ các từ tương ứng với index như sau: {'học':5, 'sinh':7, 'đi':9, 'EOS':2}, khi đó giá trị mã hóa index của câu 'học sinh đi học' sẽ là [5, 7, 9, 5, 2].
- `zeroPadding()` : Nhận đầu vào là 1 list các chuỗi index đại diện cho câu. Hàm này sẽ xác định câu có độ dài lớn nhất trong list. Sau đó padding thêm 0 vào cuối mỗi chuỗi index các giá trị 0 về cuối để các câu có độ dài bằng nhau và bằng độ dài lớn nhất.
- `binaryMatrix()` : Một ma trận sẽ biểu diễn một batch của các câu truyền vào. Mỗi dòng của ma trận sẽ đại diện cho 1 câu. Trong ma trận này sẽ tồn tại những index tương ứng với vị trí padding. `binaryMatrix` sẽ đánh dấu các vị trí mà tương ứng với padding bằng 0 và tương ứng với từ bằng 1.
- `inputVar()` : Nhận giá trị truyền vào là 1 list các câu input và từ điển. Hàm sẽ trả về ma trận được padding thêm 0 đại diện cho list các câu input và list độ dài tương ứng thực tế của các câu.
- `outputVar()` : Nhận giá trị truyền vào là 1 list các câu output và từ điển. Về cơ bản cũng giống như `inputVar()` nhưng ngoài trả về ma trận padding và độ dài lớn nhất của các câu còn trả về thêm ma trận mask có kích thước bằng ma trận padding đánh dấu các vị trí là padding (=0) và từ (=1).
- `batch2TrainData()` : Nhận giá trị truyền vào là các cặp câu (input, output) và từ điển. Hàm sẽ khởi tạo batch cho huấn luyện mô hình bao gồm: ma trận batch input, ma trận batch output, ma trận mask đánh dấu padding của output. Ngoài ra còn trả thêm list độ dài thực tế các câu trong input và độ dài lớn nhất của các câu trong output.


```

1  import random
2  import torch
3
4  def indexesFromSentence(voc, sentence):
5      return [voc.word2index[word] for word in sentence.split(' ')] + [EOS_token]
6
7  # Padding thêm 0 vào list nào có độ dài nhỏ hơn về phía bên phải
8  def zeroPadding(l, fillvalue=PAD_token):
9      return list(itertools.zip_longest(*l, fillvalue=fillvalue))
10
11 # Tạo ma trận binary có kích thước như ma trận truyền vào l nhưng giá trị của mỗi phần t
12 def binaryMatrix(l, value=PAD_token):
13     m = []
14     for i, seq in enumerate(l):
15         m.append([])
16         for token in seq:
17             if token == PAD_token:
18                 m[i].append(0)
19             else:
20                 m[i].append(1)
21     return m
22
23 # Returns padded input sequence tensor and lengths
24 def inputVar(l, voc):
25     indexes_batch = [indexesFromSentence(voc, sentence) for sentence in l]
26     lengths = torch.tensor([len(indexes) for indexes in indexes_batch])
27     padList = zeroPadding(indexes_batch)
28     padVar = torch.LongTensor(padList)
29     return padVar, lengths
30
31 # Returns padded target sequence tensor, padding mask, and max target length
32 def outputVar(l, voc):
33     indexes_batch = [indexesFromSentence(voc, sentence) for sentence in l]
34     max_target_len = max([len(indexes) for indexes in indexes_batch])
35     padList = zeroPadding(indexes_batch)
36     mask = binaryMatrix(padList)
37     mask = torch.ByteTensor(mask)
38     padVar = torch.LongTensor(padList)
39     return padVar, mask, max_target_len
40
41 # Returns all items for a given batch of pairs
42 def batch2TrainData(voc, pair_batch):
43     pair_batch.sort(key=lambda x: len(x[0].split(" ")), reverse=True)
44     input_batch, output_batch = [], []
45     for pair in pair_batch:
46         input_batch.append(pair[0])
47         output_batch.append(pair[1])
48     inp, lengths = inputVar(input_batch, voc)
49     output, mask, max_target_len = outputVar(output_batch, voc)
50     return inp, lengths, output, mask, max_target_len
51
52
53 # Example for validation
54 small_batch_size = 4
55 batches = batch2TrainData(voc, [random.choice(pairs) for _ in range(small_batch_size)])
56 input_variable, lengths, target_variable, mask, max_target_len = batches

```

Cuối cùng ta sẽ thử nghiệm giá trị trả ra của hàm `batch2TrainData()` khi lựa chọn ra 4 cặp câu bất kì trong list các cặp (input, output).

```

1  print("input_variable: \n", input_variable)
2  print("lengths: \n", lengths)
3  print("target_variable: \n", target_variable)
4  print("mask: \n", mask)
5  print("max_target_len: \n", max_target_len)

```

```

1  input_variable:
2      tensor([[ 66, 369, 66, 1272],
3              [ 567, 183, 28, 616],
4              [ 392, 1558, 1143, 175],
5              [ 394, 31, 31, 5558],
6              [ 2, 2, 2, 2]])
7  lengths:
8      tensor([5, 5, 5, 5])
9  target_variable:
10     tensor([[ 66, 370, 124, 1275],
11             [ 568, 184, 29, 617],
12             [ 393, 1560, 1144, 6240],
13             [ 395, 31, 31, 5559],
14             [ 2, 2, 2, 2]])
15  mask:
16     tensor([[1, 1, 1, 1],
17             [1, 1, 1, 1],
18             [1, 1, 1, 1],
19             [1, 1, 1, 1],
20             [1, 1, 1, 1]], dtype=torch.uint8)
21  max_target_len:
22     5

```

2. Mô hình seq2seq sửa lỗi chính tả

Model seq2seq sẽ nhận đầu vào là 1 chuỗi và trả ra kết quả output cũng là 1 chuỗi. Chính vì thế tên gọi của mô hình là sequence to sequence (từ câu đến câu).

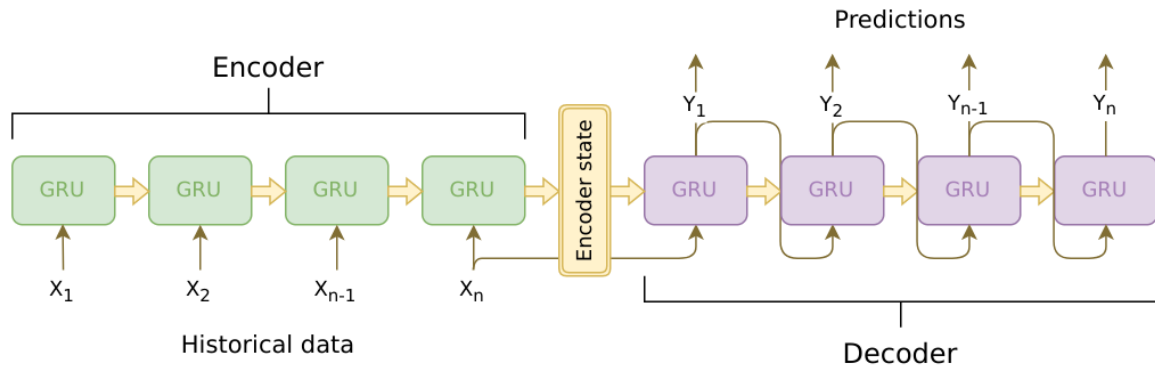
Trong kiến trúc của model seq2seq sẽ gồm 2 phrases: Encoder và decoder.

- **Encoder:** Nhúng các từ thành những véc tơ embedding với kích thước tùy ý. Encoder sẽ xây dựng một chuỗi các xử lý liên hoàn sao cho output của bước liên trước là input của bước liên sau. Khi đó tại mỗi time step sẽ truyền đầu vào là các véc tơ đã được mã hóa ứng với mỗi từ \mathbf{x}_i . Encoder sẽ trả ra 2 kết quả ở đầu ra gồm: encoder outputs đại diện cho toàn bộ câu input trong đó mỗi véc tơ của encoder outputs đại diện cho 1 từ trong câu và hidden state của GRU cuối cùng. hidden state sẽ được sử dụng để làm giá trị hidden khởi tạo cho quá trình Decoder (chi tiết ở hình 1 bên dưới). ma trận encoder outputs được sử dụng để tính attention weight tại mỗi time step trong phrase decoder. Để dễ hình dung output của một mạng RNN chúng ta có thể xem thêm lý thuyết về mạng LSTM (https://phamdinhhkhanh.github.io/2019/04/22/L%C3%BD_thuy%E1%BA%BFt_v%E1%BB%81_m%E1%BA%A1ng_LSTM.htm)
- **Decoder:** Sau phrase encoder ta sẽ thu được một hidden state của GRU cuối cùng và ma trận encoder outputs đại diện cho toàn bộ câu input. Phrase decoder có tác dụng giải mã thông tin đầu ra ở encoder thành các từ. Do đó tại mỗi time step đều trả ra các véc tơ phân phối xác suất của từ tại bước đó. Từ đó ta có thể xác định được từ có khả năng xảy ra nhất tại mỗi time step. Tại time step t , mô hình sẽ kết hợp giữa decoder embedding véc tơ h_t đại diện cho token $word_t$ và ma trận encoder outputs theo cơ chế global attention (được đề xuất bởi anh Lương Mạnh Thắng) để tính ra trọng số attention weight phân bố cho vị trí từ ở câu input lên $word_t$. véc tơ context đại diện cho toàn bộ câu input sẽ được tính bằng tích trọng số của attention weight với từng encoder véc tơ của ma trận encoder outputs (cụ thể hình 3). Tiếp theo để dự báo cho từ kế tiếp $word_{t+1}$ ta cần kết hợp véc tơ decoder hidden state h_{t+1} và véc tơ context như hình 5. Quá trình này sẽ được lặp lại liên tục cho đến khi gặp token cuối cùng là <EOS> đánh dấu vị trí cuối của câu. Như vậy ta sẽ trải qua các bước:
 - Tính decoder input chính là embedding véc tơ h_t của từ đã biết $word_t$ dựa vào embedding layer.
 - Tính attention weight đánh giá mức độ tập trung của các từ input vào từ được dự báo dựa vào ma trận encoder outputs và h_t .
 - Tính context véc tơ c_t chính là tích có trọng số của attention weights với các véc tơ đại diện cho mỗi từ input trong ma trận encoder outputs.
 - Truyền decoder input và last hidden state ở bước t vào mô hình decoder GRU để tính ra được decoder hidden state h_{t+1} .
 - Kết hợp decoder hidden state h_{t+1} và context véc tơ c_t để dự báo từ tại time step $t + 1$.

Quá trình tiếp tục cho đến khi gặp token <EOS> đánh dấu kết thúc câu.

Kết quả đầu ra là chuỗi các indexes (lấy theo vocabulary) đại diện cho từng từ tại mỗi vị trí của câu.

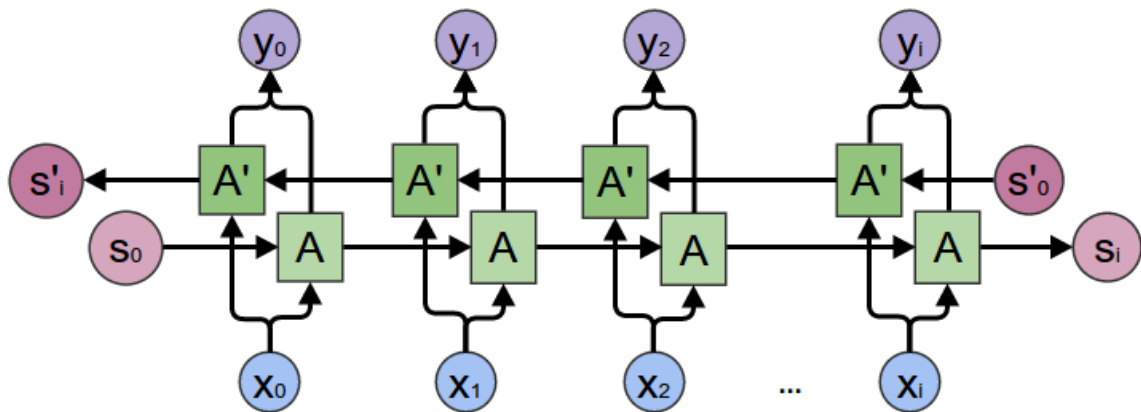
Top



Hình 1: Sơ đồ Encoder và Decoder sử dụng mạng GRU trong mô hình seq2seq.

2.1. Encoder

Tại phrase encoder chúng ta sẽ sử dụng kiến trúc mạng bidirectional GRU (2 chiều) như bên dưới để mã hóa thông tin.



Hình 2: Kiến trúc mạng bidirectional GRU. Khác với mạng unidirectional GRU (1 chiều) chỉ có chiều từ trái qua phải. Mạng GRU 2 chiều sẽ đánh giá thêm chiều từ phải qua trái. Điều này giúp cho việc học được đầy đủ hơn khi sự phụ thuộc của các từ trong câu luôn tuân theo cả 2 chiều.

Lưu ý một embedding layer được áp dụng để mã hóa các từ về một véc tơ với kích thước được khai báo là `hidden_size`. Khi huấn luyện xong model, embedding layer sẽ có kết quả sao cho những từ gần nghĩa sẽ được đại diện bởi những véc tơ sao cho độ tương quan về nghĩa càng lớn. Độ tương quan này được đo bằng `cosin similarity` của các embedding véc tơ của mỗi từ.

Mô hình RNN sẽ nhận đầu vào là một batch. Để padding một batch vào model RNN thì chúng ta phải pack dữ liệu bằng hàm `nn.utils.rnn.pack_padded_sequence` để tự động padding thêm 0 vào các véc tơ từ. Và ở bước decoder chúng ta phải unpack phần zero padding bao quanh đầu ra bằng hàm `nn.utils.rnn.pad_packed_sequence`.

Bên dưới ta sẽ xây dựng Module encoder.

Quá trình tính toán của đồ thị:

1. Mã hóa các từ về index và embedding các index thành các véc tơ.
2. Xác định batch cho mô hình. Padding câu về chung 1 độ dài và đóng gói thành batch bằng hàm `nn.utils.rnn.pack_padded_sequence`.
3. Thực hiện quá trình feed forward.
4. unpack các zero padding bằng hàm `nn.utils.rnn.pad_packed_sequence`
5. Tính tổng của bidirectional GRU outputs theo hai chiều trái và phải.
6. Trả về encoder output của layer GRU cuối cùng và hidden state tại layer cuối cùng.

Các tham số của model:

Inputs:

Top

- `input_seq` : Batch của các câu đầu vào dưới dạng tensor (như tham số input được trả ra của hàm `batch2trainData()`). `shape = max_length x batch_size`. Trong đó `batch_size` là kích thước batch (tương ứng với số câu được đưa vào batch) và `max_length` là kích thước của câu văn đã được pad hoặc trim về độ dài chuẩn.
- `input_lengths` : list của độ dài thực tế (không tính pad hoặc trim) các câu tương ứng trong batch.

Outputs:

- `outputs` : hidden layer cuối cùng của GRU (tổng của bidirectional outputs). Trên hình chính là GRU tại vị trí x_n ; `shape = (max_length, batch_size, hidden_size)`. Trong đó `max_length` là độ dài của câu. `batch_size` là kích thước batch. `hidden_size` là số chiều mã hóa của embedding layer.
- `hidden` : là ma trận concatenate của các hidden state được cập nhật từ mỗi layer của GRU; `shape = (n_layers x num_directions, batch_size, hidden_size)`. Trong đó:
 - `n_layers` là số layers GRU được áp dụng. Số lượng layers này sẽ bằng chính số lượng hidden state của mạng GRU.
 - `hidden_size` chính là kích thước của embedding véc tơ ứng với mỗi từ.
 - `num_bidirections` là số chiều của mạng GRU. Trong trường hợp này là 2 chiều.

Trước tiên ta cần khai báo sử dụng pytorch cuda bằng đoạn code bên dưới. hàm `torch.cuda.is_available()` sẽ tự động kiểm tra xem máy có GPU hỗ trợ CUDA hay không. Nếu có sẽ khai báo `device` là `cuda` và trái lại sẽ sử dụng `cpu` để huấn luyện mô hình. Và tất nhiên là tốc độ huấn luyện trên `cuda` nhanh hơn rất nhiều.

```

1  from __future__ import absolute_import
2  from __future__ import division
3  from __future__ import print_function
4  from __future__ import unicode_literals
5
6  import torch
7  from torch.jit import script, trace
8  import torch.nn as nn
9  from torch import optim
10 import torch.nn.functional as F
11 import csv
12 import random
13 import re
14 import os
15 import unicodedata
16 import codecs
17 from io import open
18 import itertools
19 import math
20
21
22 USE_CUDA = torch.cuda.is_available()
23 device = torch.device("cuda" if USE_CUDA else "cpu")

```

```

1  class EncoderRNN(nn.Module):
2      def __init__(self, hidden_size, embedding, n_layers=1, dropout=0):
3          super(EncoderRNN, self).__init__()
4          self.n_layers = n_layers
5          self.hidden_size = hidden_size
6          self.embedding = embedding
7
8          # Initialize GRU; the input_size and hidden_size params are both set to 'hidden_
9          # set bidirectional = True for bidirectional
10         # https://pytorch.org/docs/stable/nn.html?highlight=gru#torch.nn.GRU to get more
11         self.gru = nn.GRU(input_size = hidden_size, # number of expected feature of input
12                           hidden_size = hidden_size, # number of expected feature of hidden
13                           num_layers = n_layers, # number of GRU layers
14                           dropout=(0 if n_layers == 1 else dropout), # dropout probability
15                           bidirectional=True # one or two directions.
16         )
17
18     def forward(self, input_seq, input_lengths, hidden=None):
19         # Step 1: Convert word indexes to embeddings
20         # shape: (max_length , batch_size , hidden_size)
21         embedded = self.embedding(input_seq)
22         # Pack padded batch of sequences for RNN module. Padding zero when length less than
23         # shape: (max_length , batch_size , hidden_size)
24         packed = nn.utils.rnn.pack_padded_sequence(embedded, input_lengths)
25         # Step 2: Forward packed through GRU
26         # outputs is output of final GRU layer
27         # hidden is concatenate of all hidden states corresponding with each time step.
28         # outputs shape: (max_length , batch_size , hidden_size x num_directions)
29         # hidden shape: (n_layers x num_directions , batch_size , hidden_size)
30         outputs, hidden = self.gru(packed, hidden)
31         # Unpack padding. Revert of pack_padded_sequence
32         # outputs shape: (max_length , batch_size , hidden_size x num_directions)
33         outputs, _ = nn.utils.rnn.pad_packed_sequence(outputs)
34         # Sum bidirectional GRU outputs to reshape shape into (max_length , batch_size ,
35         outputs = outputs[:, :, :self.hidden_size] + outputs[:, :, self.hidden_size:]
36         # outputs shape: (max_length , batch_size , hidden_size)
37         # hidden shape: (n_layers x num_directions , batch_size , hidden_size)
38         return outputs, hidden

```

Để kiểm nghiệm kết quả trả ra của EncoderRNN bên dưới ta sẽ thực hiện 1 ví dụ giả lập encoder với `hidden_size = 3` và `n_layers = 7`.

```

1  # Thử nghiệm phrase Encoder bằng cách giả lập 1 mạng Encoder với:
2  from torch import nn
3
4  hidden_size = 3
5  n_layers = 7
6  embedding = nn.Embedding(voc.num_words, hidden_size)
7  print('input_seq: \n', input_variable)
8  print('input_lengths: \n', lengths)
9  encoder = EncoderRNN(hidden_size = hidden_size, embedding = embedding, n_layers = n_layers)
10
11  print('encoder phrase: \n', encoder)
12
13  output, hidden = encoder.forward(input_seq = input_variable, input_lengths = lengths)

```

```

1  input_seq:
2  tensor([[ 66, 369, 66, 1272],
3          [ 567, 183, 28, 616],
4          [ 392, 1558, 1143, 175],
5          [ 394, 31, 31, 5558],
6          [ 2, 2, 2, 2]])
7  input_lengths:
8  tensor([5, 5, 5, 5])
9  encoder:
10 EncoderRNN(
11   (embedding): Embedding(171775, 3)
12   (gru): GRU(3, 3, num_layers=7, bidirectional=True)
13 )

1  print('output size: ', output.size())
2  print('hidden size: ', hidden.size())

1  output size: torch.Size([5, 4, 3])
2  hidden size: torch.Size([14, 4, 3])

```

Ta nhận thấy đầu ra của output là (max_length, batch_size, hidden_size) và của hidden là (n_layers x num_directions, batch_size, hidden_size). Thông qua ví dụ này bạn đọc đã hình dung được quá trình Encoder rồi chứ.

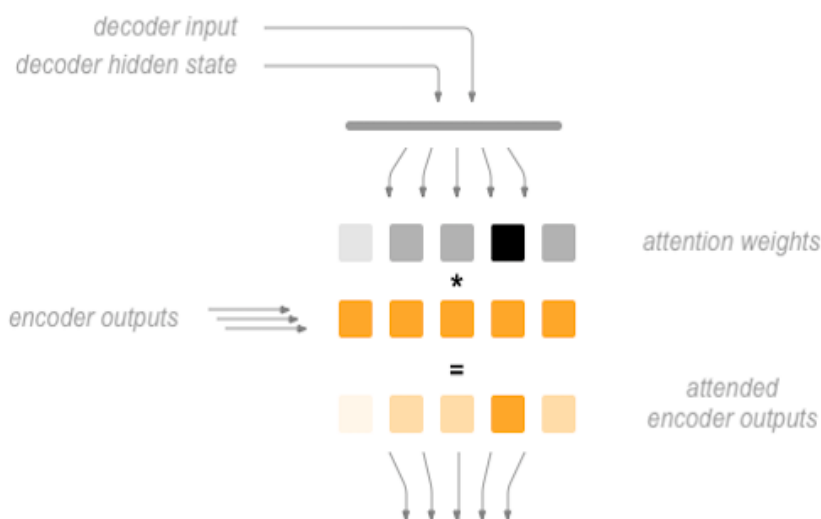
2.2. Decoder

Như vậy sau bước encoder ta thu được output là layer GRU cuối cùng (gọi là encoder outputs) và các hidden state của layer GRU cuối cùng. Bước tiếp theo chúng ta cần giải mã các kết quả thu được từ encoder thành câu hoàn chỉnh.

2.2.1. Áp dụng Attention layer

Ở phrase này chúng ta sẽ sử dụng thêm 1 layer attention để tính phân phối trọng số attention weights cho các véc tơ từ (hidden state) của ma trận encoder outputs. attention weight sẽ có kích thước bằng đúng độ dài của câu. Sau khi tính tổng theo attention weights ta sẽ thu được context véc tơ đại diện cho toàn bộ câu. Quá hình này được thể hiện như hình 5. context véc tơ sẽ kết hợp với các decoder hidden state (các thẻ h_t màu đỏ trong hình 5) tại mỗi time step để dự đoán từ tiếp theo. decoder hidden state chính là output trả ra của model Decoder sau mỗi time step. Quá trình này lặp lại truy hồi (output layer trước làm input cho layer sau) cho đến khi kết quả trả về là <EOS> (end of sequence).

Hình bên dưới sẽ minh họa cho quá trình tính toán attention weights dựa trên sự kết hợp của decoder input (véc tơ embedding của từ được dự báo ở bước trước), decoder hidden state (ma trận của các hidden state sau cùng) và encoder outputs.



Hình 3: Kết hợp giữa decoder input, decoder hidden state và encoder outputs để tạo ra một attended encoder outputs.

Top

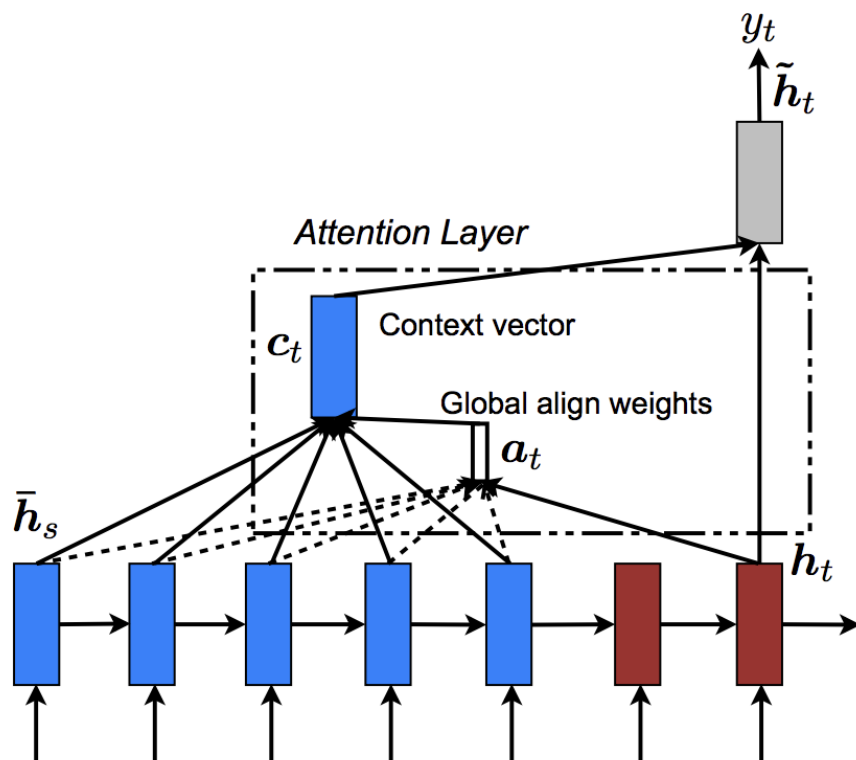
Các bạn đã hình dung được quá trình kết hợp attention vào decoder rồi chứ? Điểm mấu chốt là chúng ta phải tính ra được các trọng số attention tại mỗi time step. Việc tính toán attention weight đã được đề xuất theo rất nhiều cách khác nhau bởi anh Lương Mạnh Thắng (<https://arxiv.org/abs/1508.04025>). Trong đó điểm cải tiến so với cha đẻ của attention layer Bahdanau chính là một global attention được tính toán trên toàn bộ encoder hidden state thay vì local attention được tính toán dựa trên chỉ encoder hidden state của time step hiện tại. Điểm khác biệt thứ 2 là global attention tính attention weight chỉ dựa trên decoder hidden state tại time step hiện tại thay vì local attention được đề xuất bởi Bahdanau yêu cầu thêm các decoder hidden state trước đó. Theo đó điểm của các attention ở time step hiện tại được anh Thắng Lương đề xuất dưới nhiều công thức khác nhau như hình bên dưới.

$$\text{score}(\mathbf{h}_t, \bar{\mathbf{h}}_s) = \begin{cases} \mathbf{h}_t^\top \bar{\mathbf{h}}_s & \text{dot} \\ \mathbf{h}_t^\top \mathbf{W}_a \bar{\mathbf{h}}_s & \text{general} \\ \mathbf{v}_a^\top \tanh(\mathbf{W}_a [\mathbf{h}_t; \bar{\mathbf{h}}_s]) & \text{concat} \end{cases}$$

Hình 4: Tính điểm tại time step t dựa trên decoder hidden state (\mathbf{h}_t) của thời điểm t và toàn bộ các encoder hidden states ($\bar{\mathbf{h}}_s$).

Các bước tiếp theo để tính context vector ta có thể tham khảo ở bài diễn giải về attention is all you need (<https://phamdinhhkhanh.github.io/2019/06/18/AttentionLayer.html>).

Cơ chế tính context vector có thể khái quát hóa như hình bên dưới. Lưu ý chúng ta sẽ triển khai Attention Layer như một nn.Module tách biệt và được gọi là Attn. Kết quả của attention layer là attention weights là một phân phối xác suất hàm softmax dưới dạng tensor có kích thước (batch_size, 1, max_length).



Hình 5: Cơ chế global attention.

Bên dưới là code của class Attention.

```

1      # Luong attention layer
2      class Attn(nn.Module):
3          def __init__(self, method, hidden_size):
4              super(Attn, self).__init__()
5              self.method = method
6              if self.method not in ['dot', 'general', 'concat']:
7                  raise ValueError(self.method, "is not an appropriate attention method.")
8              self.hidden_size = hidden_size
9              if self.method == 'general':
10                 self.attn = nn.Linear(self.hidden_size, hidden_size)
11             elif self.method == 'concat':
12                 self.attn = nn.Linear(self.hidden_size * 2, hidden_size)
13                 self.v = nn.Parameter(torch.FloatTensor(hidden_size))
14
15         def dot_score(self, hidden, encoder_output):
16             # encoder_output shape: (max_length , batch_size , hidden_size)
17             # hidden shape: (1 , batch_size , hidden_size)
18             # return shape: (max_length, batch_size)
19             return torch.sum(hidden * encoder_output, dim=2)
20
21         def general_score(self, hidden, encoder_output):
22             # encoder_output shape: (max_length , batch_size , hidden_size)
23             # hidden shape: (batch_size , hidden_size)
24             # energy shape: (max_length , batch_size , hidden_size)
25             # return shape: (max_length , batch_size)
26             energy = self.attn(encoder_output)
27             return torch.sum(hidden * energy, dim=2)
28
29         def concat_score(self, hidden, encoder_output):
30             # encoder_output shape: (max_length , batch_size , hidden_size)
31             # hidden shape: (batch_size , hidden_size)
32             # energy shape: (max_length , batch_size , 2*hidden_size)
33             # self.v shape: (hidden_size)
34             # return shape: (max_length , batch_size)
35             energy = self.attn(torch.cat((hidden.expand(encoder_output.size(0), -1, -1), enc
36             return torch.sum(self.v * energy, dim=2)
37
38         def forward(self, hidden, encoder_outputs):
39             # Calculate the attention weights (energies) based on the given method
40             # attn_energies.shape: (max_length , batch_size)
41             if self.method == 'general':
42                 attn_energies = self.general_score(hidden, encoder_outputs)
43             elif self.method == 'concat':
44                 attn_energies = self.concat_score(hidden, encoder_outputs)
45             elif self.method == 'dot':
46                 attn_energies = self.dot_score(hidden, encoder_outputs)
47
48             # Transpose max_length and batch_size dimensions
49             # attn_energies.shape: (batch_size , max_length)
50             attn_energies = attn_energies.t()
51             # Return the softmax normalized probability scores (with added dimension)
52             attn_weights = F.softmax(attn_energies, dim=1).unsqueeze(1)
53             # attn_weights shape: (batch_size , 1 , max_length)
54             return attn_weights

```

Để hiểu rõ hơn attention hoạt động như thế nào ta sẽ cùng thử nghiệm truyền vào decoder hidden là các hidden state véc tơ của decoder và encoder outputs là ma trận mã hóa đầu ra của toàn bộ các từ trong câu input. Ví dụ bên dưới ta sẽ xét ở time step 0 nên decoder hidden khi đó sẽ là véc tơ hidden state cuối cùng của encoder (tham số last_hidden_encoder).


```

1     method = 'dot'
2     # Take last hidden encoder vector as a initialize of decoder
3     last_hidden_encoder = hidden[encoder.n_layers, :, :].unsqueeze(0)
4     print('last_hidden_encoder.size: ', last_hidden_encoder.size())
5     print('encoder_outputs.size: ', output.size())
6     print('hidden_size: ', hidden_size)
7     decoder_hidden_size = 7
8     print('decoder_hidden_size: ', decoder_hidden_size)
9
10    attn = Attn(method = method, hidden_size = decoder_hidden_size)
11    attn_weights = attn.forward(hidden = last_hidden_encoder, encoder_outputs = output)
12    print('attn_weights shape: ', attn_weights.size())

1     last_hidden_encoder.size:  torch.Size([1, 4, 3])
2     encoder_outputs.size:  torch.Size([5, 4, 3])
3     hidden_size:  3
4     decoder_hidden_size:  7
5     attn_weights shape:  torch.Size([4, 1, 5])

```

Như vậy tại mỗi time step để tính attention weights chúng ta sẽ truyền vào một decoder hidden véc tơ (chính là output của GRU tại time step đó) và ma trận encoder outputs. Thông qua các phương pháp `dot_score`, `general_score` hoặc `concat_score` ta sẽ tính được attention weights đại diện cho mức độ attention của từng vị trí của encoder outputs lên từ được dự báo.

Qua diễn giải trên chúng ta mới chỉ hiểu cách thức mà attention hoạt động. Vậy thì attention được áp dụng trong toàn bộ quá trình decoder như thế nào tại mỗi time step. Để hiểu rõ chúng ta tìm hiểu qua phrase decoder.

2.2.2. Quá trình Decoder

Trong bước decoder chúng ta sẽ truyền một từ một lần tại mỗi time step. Do đó các véc tơ embedding của từ và GRU output phải có chung shape là $(1, \text{batch_size}, \text{hidden_size})$.

Quá trình tính toán đồ thị

1. Lấy embedding của input word hiện tại.
2. thực hiện lan truyền thuận (feed forward) qua một GRU 1 chiều (unidirectional GRU) để thu được decoder hidden state.
3. Tính attention weights từ decoder hidden state của GRU hiện tại ở bước 2 kết hợp với encoder outputs của encoder.
4. Nhân attention weights với encoder outputs để thu được một tổng có trọng số là context véc tơ.
5. Concatenate context véc tơ và GRU output.
6. Dự báo từ tiếp theo (không sử dụng softmax).
7. Trả về output và hidden state cuối cùng.

Inputs

- `input_step` : một bước thời gian tương ứng với 1 từ của input sequence; shape = $(1, \text{batch_size})$.
- `last_hidden` : layer GRU cuối cùng; shape = $(n_layers \times \text{num_directions}, \text{batch_size}, \text{hidden_size})$.
- `encoder_outputs` : đầu ra của bước encoder; shape = $(\text{max_length}, \text{batch_size}, \text{hidden_size})$.

Outputs

- `output` : softmax normalized tensor trả về xác suất tương ứng với mỗi từ là từ tại vị trí tương ứng của câu; shape = $(\text{batch_size}, \text{voc.num_words})$
- `hidden` : hidden state cuối cùng của GRU; shape = $(n_layers \times \text{num_directions}, \text{batch_size}, \text{hidden_size})$. Do ở phrase decoder ta chỉ áp dụng unidirectional GRU nên shape = $(n_layers, \text{batch_size}, \text{hidden_size})$.

```

1  class LuongAttnDecoderRNN(nn.Module):
2      def __init__(self, attn_model, embedding, hidden_size, output_size, n_layers=1, drop
3          super(LuongAttnDecoderRNN, self).__init__()
4
5          # Keep for reference
6          self.attn_model = attn_model
7          self.hidden_size = hidden_size
8          self.output_size = output_size
9          self.n_layers = n_layers
10         self.dropout = dropout
11
12         # Define layers
13         self.embedding = embedding
14         self.embedding_dropout = nn.Dropout(dropout)
15         self.gru = nn.GRU(hidden_size, hidden_size, n_layers, dropout=(0 if n_layers ==
16         self.concat = nn.Linear(hidden_size * 2, hidden_size)
17         self.out = nn.Linear(hidden_size, output_size)
18
19         self.attn = Attn(attn_model, hidden_size)
20
21     def forward(self, input_step, last_hidden, encoder_outputs):
22         '''
23         input_step: list time step index of batch. shape (1 x batch_size)
24         last_hidden: last hidden output of hidden layer (we can take in right direction
25         encoder_outputs: output of encoder
26         '''
27         #=====
28         # Step 1: Embedding current sequence index
29         # Note: we run this one step (word) at a time
30         # Get embedding of current input word
31         # embedded shape: 1 x batch_size x hidden_size
32         embedded = self.embedding(input_step)
33         embedded = self.embedding_dropout(embedded)
34
35         #=====
36         # Step 2: pass embedded and last hidden into decoder
37         # Forward through unidirectional GRU
38         # rnn_output shape: 1 x batch_size x hidden_size
39         # hidden shape: n_layers x batch_size x hidden_size
40         rnn_output, hidden = self.gru(embedded, last_hidden)
41         # Calculate attention weights from the current GRU output
42         # attn_weights shape: batch_size x 1 x max_length
43         attn_weights = self.attn(rnn_output, encoder_outputs)
44         # Multiply attention weights to encoder outputs to get new "weighted sum" contex
45         # encoder_outputs shape: max_length x batch_size x hidden_size
46         # context shape: batch_size x 1 x hidden_size
47         context = attn_weights.bmm(encoder_outputs.transpose(0, 1))
48         # Concatenate weighted context vector and GRU output using Luong eq. 5
49         # rnn_output shape: batch_size x hidden_size
50         rnn_output = rnn_output.squeeze(0)
51         # context shape: batch_size x hidden_size
52         context = context.squeeze(1)
53
54         #=====
55         # Step 3: calculate output probability distribution
56         # concat_input shape: batch_size x (2*hidden_size)
57         concat_input = torch.cat((rnn_output, context), 1)
58         # concat_output shape: batch_size x hidden_size
59         concat_output = torch.tanh(self.concat(concat_input))
60         # Predict next word using Luong eq. 6
61         # output shape: output_size
62         output = self.out(concat_output)
63         output = F.softmax(output, dim=1)
64         # Return output and final hidden state
65         return output, hidden

```

Bên dưới ta sẽ giải thích quá trình dự đoán các từ tiếp theo tại time step đầu tiên. Lưu ý quá trình được thực hiện trên batch nên các đầu vào sẽ có kích thước scale theo batch_size.

Top

```

1     time_step = 0
2     # Take all index of batch at time step 0. All words are <SOS> mark for start of sentence
3     input_step = torch.tensor([SOS_token] * small_batch_size).unsqueeze(0)
4     n_layers = 7
5     # take last hidden vector of encoder
6     last_hidden = hidden[:n_layers]
7     print('batch_size: ', small_batch_size)
8     print('input_step.size at time_step 0: ', input_step.size())
9     print('last_hidden.size: ', last_hidden.size())
10    attn_model = 'dot'
11    hidden_size = 3
12    # Output size of decoder model is size of vocabulary
13    output_size = len(voc.word2index)
14
15
16    luongAttnDecoderRNN = LuongAttnDecoderRNN(attn_model = attn_model,
17                                              embedding = embedding,
18                                              hidden_size = hidden_size,
19                                              output_size = output_size,
20                                              n_layers = n_layers)
21
22    print('luongAttnDecoderRNN phrase: \n', luongAttnDecoderRNN)
23    dec_output, dec_hidden = luongAttnDecoderRNN.forward(input_step = input_step,
24                                                         last_hidden = last_hidden,
25                                                         encoder_outputs = output)

```

```

1     batch_size: 4
2     input_step.size at time_step 0: torch.Size([1, 4])
3     last_hidden.size: torch.Size([7, 4, 3])
4     luongAttnDecoderRNN phrase:
5     LuongAttnDecoderRNN(
6       (embedding): Embedding(171775, 3)
7       (embedding_dropout): Dropout(p=0.1)
8       (gru): GRU(3, 3, num_layers=7, dropout=0.1)
9       (concat): Linear(in_features=6, out_features=3, bias=True)
10      (out): Linear(in_features=3, out_features=171772, bias=True)
11      (attn): Attn()
12    )

1     print('dec_output.size: ', dec_output.size())
2     print('dec_hidden.size: ', dec_hidden.size())

1     dec_output.size: torch.Size([4, 171772])
2     dec_hidden.size: torch.Size([7, 4, 3])

```

Diễn giải hàm forward:

Như vậy sau khi truyền vào decoder embedding véc tơ đại diện cho từ liền trước kết hợp với last hidden output của chính time step trước (có kích thước = $n_layers \times hidden_size$) ta sẽ thu được decoder hidden state véc tơ (chính là các rnn_output trong hàm `LuongAttnDecoderRNN.forward()`). Đây chính là một quá trình lan truyền thuật thông thường của một mạng RNN.

Tiếp theo kết hợp decoder hidden state véc tơ với encoder outputs thu được ở phrase encoding ta sẽ tính được attention weight và từ đó suy ra context véc tơ. concatenate context véc tơ và decoder hidden state véc tơ tạo thành feature learning. Truyền feature learning véc tơ này qua hàm softmax ta sẽ tìm được véc tơ phân phối của từ tại vị trí time step.

3. Huấn luyện model.

3.1 Loss function

Bởi vì chúng ta huấn luyện mô hình dựa trên các batch đã được padding nên không đơn thuần là đưa toàn bộ các phần tử của output vào để tính toán giá trị của loss function mà cần có thêm một binary mask tensor để xác định các vị trí padding của target tensor. Khi đó loss function được tính toán dựa trên decoder's output tensor, target tensor và binary mask tensor.

mask tensor. Hàm loss function chính là trung bình âm của log likelihood của toàn bộ các phần tử tương ứng với 1 trong 1 trong binary mask tensor.

```

1  def maskNLLLoss(inp, target, mask):
2      nTotal = mask.sum()
3      crossEntropy = -torch.log(torch.gather(inp, 1, target.view(-1, 1)).squeeze(1))
4      loss = crossEntropy.masked_select(mask).mean()
5      loss = loss.to(device)
6      return loss, nTotal.item()

```

3.2. Huấn luyện vòng lặp đơn

Hàm huấn luyện `train()` bên dưới sẽ được xây dựng để huấn luyện cho 1 batch đơn lẻ.

Các kĩ thuật sử dụng

Chúng ta sẽ sử dụng 2 kĩ thuật để hỗ trợ hội tụ:

- **teacher forcing**: Tại một ngưỡng xác suất nào đó được thiết lập thông qua tham số `teacher_forcing_ratio`, chúng ta sẽ sử dụng luôn current target word như là đầu vào cho decoder tại bước tiếp theo hơn là sử dụng dự báo từ mô hình tại bước hiện tại. Kĩ thuật này hỗ trợ cho quá trình huấn luyện hiệu quả hơn. Tuy nhiên điểm hạn chế là có thể khiến cho mô hình thiếu ổn định trong quá trình suy diễn, khi decoder có thể không có cơ hội học đủ các khả năng để tạo ra output sequence trong quá trình huấn luyện. Do đó phải hết sức cẩn thận khi sử dụng `teacher_forcing_ratio`.
- **gradient clipping**: Hay còn gọi là kĩ thuật kẹp gradient để tránh hiện tượng bùng nổ gradient (exploding gradient - gradient lớn quá nhanh) bằng các thiết lập giá trị max cho gradient descent.

Chuỗi các bước triển khai

1. Truyền vào toàn bộ các batch của câu qua encoder.
2. Khởi tạo decoder input bằng token `<SOS>` (start of sequence) và encoder's hidden state cuối cùng.
3. Truyền input batch sequence vào decoder một lần tại một time step.
4. Nếu thực hiện teaching forcing: Thiết lập decoder input tiếp theo chính là giá trị target hiện tại; trái lại: thiết lập decoder input tiếp theo như là đầu ra của decoder hiện tại.
5. Tính loss function.
6. Thực hiện lan truyền ngược.
7. Kẹp gradient tránh hiện tượng bùng nổ tham số.
8. Cập nhật các tham số của encoder và decoder cho mô hình.

```

1      def train(input_variable, lengths, target_variable, mask, max_target_len, encoder, decoder,
2                encoder_optimizer, decoder_optimizer, batch_size, clip, max_length=MAX_LENGTH)
3
4          # Zero gradients
5          encoder_optimizer.zero_grad()
6          decoder_optimizer.zero_grad()
7
8          # Set device options
9          input_variable = input_variable.to(device)
10         lengths = lengths.to(device)
11         target_variable = target_variable.to(device)
12         mask = mask.to(device)
13
14         # Initialize variables
15         loss = 0
16         print_losses = []
17         n_totals = 0
18
19         # Forward pass through encoder
20         encoder_outputs, encoder_hidden = encoder(input_variable, lengths)
21
22         # Create initial decoder input (start with SOS tokens for each sentence)
23         decoder_input = torch.LongTensor([[SOS_token for _ in range(batch_size)]])
24         decoder_input = decoder_input.to(device)
25
26         # Set initial decoder hidden state to the encoder's final hidden state
27         decoder_hidden = encoder_hidden[:decoder.n_layers]
28
29         # Determine if we are using teacher forcing this iteration
30         use_teacher_forcing = True if random.random() < teacher_forcing_ratio else False
31
32         # Forward batch of sequences through decoder one time step at a time
33         if use_teacher_forcing:
34             for t in range(max_target_len):
35                 decoder_output, decoder_hidden = decoder(
36                     decoder_input, decoder_hidden, encoder_outputs
37                 )
38                 # Teacher forcing: next input is current target
39                 decoder_input = target_variable[t].view(1, -1)
40                 # Calculate and accumulate loss
41                 mask_loss, nTotal = maskNLLLoss(decoder_output, target_variable[t], mask[t])
42                 loss += mask_loss
43                 print_losses.append(mask_loss.item() * nTotal)
44                 n_totals += nTotal
45         else:
46             for t in range(max_target_len):
47                 decoder_output, decoder_hidden = decoder(
48                     decoder_input, decoder_hidden, encoder_outputs
49                 )
50                 # No teacher forcing: next input is decoder's own current output
51                 _, topi = decoder_output.topk(1)
52                 decoder_input = torch.LongTensor([[topi[i][0] for i in range(batch_size)]])
53                 decoder_input = decoder_input.to(device)
54                 # Calculate and accumulate loss
55                 mask_loss, nTotal = maskNLLLoss(decoder_output, target_variable[t], mask[t])
56                 loss += mask_loss
57                 print_losses.append(mask_loss.item() * nTotal)
58                 n_totals += nTotal
59
60         # Perform backpropagation
61         loss.backward()
62
63         # Clip gradients: gradients are modified in place
64         _ = nn.utils.clip_grad_norm_(encoder.parameters(), clip)
65         _ = nn.utils.clip_grad_norm_(decoder.parameters(), clip)
66
67         # Adjust model weights
68         encoder_optimizer.step()
69         decoder_optimizer.step()
70
71         return sum(print_losses) / n_totals

```

Top

3.3. Huấn luyện vòng lặp

Cuối cùng đã đến lúc gắn kết toàn bộ quy trình huấn luyện với dữ liệu. Hàm `trainIters` chịu trách nhiệm chạy các `n_interations` quá trình huấn luyện từ các đầu vào gồm các model, hàm tối ưu optimizer, dữ liệu, v.v.

Một lưu ý là khi chúng ta lưu mô hình, chúng ta lưu vào một tarball (một dạng folder) chứa encoder và decoder `state_dicts` (chính là các tham số), `optimizers' state_dicts`, the loss, iteration, Lưu mô hình theo cách này sẽ mang lại sự linh hoạt cho mô hình nhờ checkpoint. Sau khi load một checkpoint, chúng ta sẽ có thể sử dụng các tham số mô hình để chạy suy diễn hoặc có thể tiếp tục huấn luyện ngay tại trạng thái rời đi trước đó.

```

1  def trainIters(model_name, voc, pairs, encoder, decoder,
2                encoder_optimizer, decoder_optimizer,
3                embedding, encoder_n_layers, decoder_n_layers,
4                save_dir, n_iteration, batch_size, print_every,
5                save_every, clip, corpus_name, loadFilename):
6      '''
7      model_name: Tên model
8      voc: bộ từ vựng cho mô hình
9      pairs: list các cặp (input, output)
10     encoder: phrase encoder
11     decoder: phrase decoder
12     encoder_optimizer: phương pháp tối ưu hóa encoder
13     decoder_optimizer: phương pháp tối ưu hóa decoder
14     embedding: layer embedding
15     encoder_n_layers: số lượng layers ở encoder
16     decoder_n_layers: số lượng layers ở decoder
17     save_dir: link save model
18     n_iteration: số iteration tổng cộng được chọn ngẫu nhiên từ pairs. Mỗi iteration = 1
19     batch_size: kích thước của batch.
20     print_every: số iteration của batch để in kết quả
21     save_every: số iteration của batch để save model
22     clip: trimming giá trị của tensor về 1 range
23     corpus_name: tên bộ từ vựng
24     loadFilename: link load file
25     '''
26     # Load batches for each iteration
27     training_batches = [batch2TrainData(voc, [random.choice(pairs) for _ in range(batch_
28                                         for _ in range(n_iteration))])
29
30     # Initializations
31     print('Initializing ...')
32     start_iteration = 1
33     print_loss = 0
34     if loadFilename:
35         start_iteration = checkpoint['iteration'] + 1
36
37     # Training loop
38     print("Training...")
39     for iteration in range(start_iteration, n_iteration + 1):
40         training_batch = training_batches[iteration - 1]
41         # Extract fields from batch
42         input_variable, lengths, target_variable, mask, max_target_len = training_batch
43
44         # Run a training iteration with batch
45         loss = train(input_variable, lengths, target_variable, mask, max_target_len, enc
46                     decoder, embedding, encoder_optimizer, decoder_optimizer, batch_siz
47         print_loss += loss
48
49         # Print progress
50         if iteration % print_every == 0:
51             print_loss_avg = print_loss / print_every
52             print("Iteration: {}; Percent complete: {:.1f}%; Average loss: {:.4f}".forma
53             print_loss = 0
54
55         # Save checkpoint
56         if (iteration % save_every == 0):
57             directory = os.path.join(save_dir, model_name, corpus_name, '{}-{}_{}'.forma
58             if not os.path.exists(directory):
59                 os.makedirs(directory)
60             torch.save({
61                 'iteration': iteration,
62                 'en': encoder.state_dict(),
63                 'de': decoder.state_dict(),
64                 'en_opt': encoder_optimizer.state_dict(),
65                 'de_opt': decoder_optimizer.state_dict(),
66                 'loss': loss,
67                 'voc_dict': voc.__dict__,
68                 'embedding': embedding.state_dict()
69             }, os.path.join(directory, '{}_{}.tar'.format(iteration, 'checkpointTop

```

4. Xác định phương pháp đánh giá

Sau khi huấn luyện mô hình, chúng ta đã có khả năng giao tiếp với bot. Muốn như vậy cần phải xác định xem chúng ta muốn mô hình mã hóa đầu vào như thế nào. Chúng ta có những cách giải mã sau:

Mã hóa tham lam (greedy decoding) Các mã hóa này được sử dụng trong tính huống không sử dụng có teacher forcing. Đầu ra của mô hình tại mỗi time step đơn giản là từ với xác suất cao nhất.

Bên dưới chúng ta tạo ra class GreedySearchDecoder nhận đầu vào là một câu input với shape (input_seq length, 1), một scalar độ dài đầu vào (input_length) tensor, và một max_length để giới hạn độ dài lớn nhất của câu trả về. Câu đầu vào sẽ được đánh giá qua các bước như đồ thị bên dưới.

Đồ thị tính toán:

1. Truyền input qua model encoder.
2. Chuẩn bị hidden layer cuối cùng của encoder làm hidden input đầu tiên của decoder.
3. Khởi tạo input đầu tiên của decoder là <SOS>
4. Khởi tạo tensor để append vào các từ được giải mã vào.
5. Lặp lại quá trình giải mã một word token 1 lần:
6. Truyền vào decoder.
7. Thu được word token dự báo dựa trên xác suất lớn nhất.
8. Lưu lại token và score.
9. Coi token hiện tại như là đầu vào tiếp theo của decoder.
10. Trả về một tập hợp các từ và điểm số.

```

1      class GreedySearchDecoder(nn.Module):
2          def __init__(self, encoder, decoder):
3              super(GreedySearchDecoder, self).__init__()
4              self.encoder = encoder
5              self.decoder = decoder
6
7          def forward(self, input_seq, input_length, max_length):
8              # Forward input through encoder model
9              encoder_outputs, encoder_hidden = self.encoder(input_seq, input_length)
10             # Prepare encoder's final hidden layer to be first hidden input to the decoder
11             decoder_hidden = encoder_hidden[:decoder.n_layers]
12             # Initialize decoder input with SOS_token
13             decoder_input = torch.ones(1, 1, device=device, dtype=torch.long) * SOS_token
14             # Initialize tensors to append decoded words to
15             all_tokens = torch.zeros([0], device=device, dtype=torch.long)
16             all_scores = torch.zeros([0], device=device)
17             # Iteratively decode one word token at a time
18             for _ in range(max_length):
19                 # Forward pass through decoder
20                 decoder_output, decoder_hidden = self.decoder(decoder_input, decoder_hidden,
21                 # Obtain most likely word token and its softmax score
22                 decoder_scores, decoder_input = torch.max(decoder_output, dim=1)
23                 # Record token and score
24                 all_tokens = torch.cat((all_tokens, decoder_input), dim=0)
25                 all_scores = torch.cat((all_scores, decoder_scores), dim=0)
26                 # Prepare current token to be next decoder input (add a dimension)
27                 decoder_input = torch.unsqueeze(decoder_input, 0)
28             # Return collections of word tokens and scores
29             return all_tokens, all_scores

```

Đánh giá câu dự báo Qua bước trên chúng ta đã có hàm decoder xác định câu đầu ra dựa trên câu đầu vào. Hàm đánh giá quản lý quá trình xử lý câu đầu vào mức độ thấp. Trước tiên, chúng ta định dạng câu dưới dạng một batch đầu vào của các word index với batch_size = 1. Chúng ta thực hiện điều này bằng cách chuyển đổi các từ của câu thành các indexes tương ứng và transpose để chuẩn bị tensor cho các mô hình. Chúng ta cũng tạo ra một tensor độ dài chứa độ dài của câu đầu vào. Trong trường hợp này, độ dài là scalar vì chúng ta chỉ đánh giá một câu tại một thời điểm (batch_size = 1). Tiếp theo, chúng ta thu được một tensor các câu trả về được giải mã bằng cách sử dụng object GreedySearchDecoder (trình tìm kiếm). Cuối cùng, chúng tôi chuyển đổi các indexes trả về thành các từ và trả về danh sách các từ được giải mã.

class evaluateInput hoạt động như giao diện người dùng cho chatbot. Khi được gọi, một input text field sẽ sinh ra trong đó chúng ta có thể nhập câu hỏi của mình. Sau khi nhập câu đầu vào và nhấn Enter, text sẽ được chuẩn hóa giống như dữ liệu huấn luyện và cuối cùng được đưa vào hàm đánh giá để thu được câu đầu ra được giải mã. Chúng ta lặp lại quy trình này cho liên tục cho đến khi nhấn q hoặc quit để quit.

Cuối cùng, nếu một câu được nhập có chứa một từ không có trong từ vựng, chúng ta sẽ xử lý việc này một cách tinh tế bằng cách in một thông báo lỗi và nhắc người dùng nhập một câu khác.

```

1  def evaluate(encoder, decoder, searcher, voc, sentence, max_length=MAX_LENGTH):
2      ### Format input sentence as a batch
3      # words -> indexes
4      indexes_batch = [indexesFromSentence(voc, sentence)]
5      # Create lengths tensor
6      lengths = torch.tensor([len(indexes) for indexes in indexes_batch])
7      # Transpose dimensions of batch to match models' expectations
8      input_batch = torch.LongTensor(indexes_batch).transpose(0, 1)
9      # Use appropriate device
10     input_batch = input_batch.to(device)
11     lengths = lengths.to(device)
12     # Decode sentence with searcher
13     tokens, scores = searcher(input_batch, lengths, max_length)
14     # indexes -> words
15     decoded_words = [voc.index2word[token.item()] for token in tokens]
16     return decoded_words
17
18
19 def evaluateInput(encoder, decoder, searcher, voc):
20     input_sentence = ''
21     while(1):
22         try:
23             # Get input sentence
24             input_sentence = input('> ')
25             # Check if it is quit case
26             if input_sentence == 'q' or input_sentence == 'quit': break
27             # Normalize sentence
28             input_sentence = normalizeString(input_sentence)
29             # Evaluate sentence
30             output_words = evaluate(encoder, decoder, searcher, voc, input_sentence)
31             # Format and print response sentence
32             output_words[:] = [x for x in output_words if not (x == 'EOS' or x == 'PAD')]
33             print('Bot:', ' '.join(output_words))
34
35         except KeyError:
36             print("Error: Encountered unknown word.")

```

5. Huấn luyện model

Bên dưới ta sẽ thực hiện huấn luyện mô hình bằng cách thiết lập các tham số. Chúng ta có thể thử nghiệm nhiều tham số cấu hình khác nhau để lựa chọn được 1 mô hình tối ưu. Chúng ta cũng có thể xây dựng mô hình từ đầu hoặc load từ checkpoint một mô hình sẵn có.

Load model

```

1      # Configure models
2      model_name = 'correct_spelling_model'
3      corpus_name = 'corpus_aivivn'
4      attn_model = 'dot'
5      #attn_model = 'general'
6      #attn_model = 'concat'
7      hidden_size = 500
8      encoder_n_layers = 2
9      decoder_n_layers = 2
10     dropout = 0.1
11     batch_size = 100
12
13     # Set checkpoint to load from; set to None if starting from scratch
14     loadFilename = None
15     checkpoint_iter = 5000
16
17     # Load model if a loadFilename is provided
18     if loadFilename:
19         # If loading on same machine the model was trained on
20         checkpoint = torch.load(loadFilename)
21         # If loading a model trained on GPU to CPU
22         #checkpoint = torch.load(loadFilename, map_location=torch.device('cpu'))
23         encoder_sd = checkpoint['en']
24         decoder_sd = checkpoint['de']
25         encoder_optimizer_sd = checkpoint['en_opt']
26         decoder_optimizer_sd = checkpoint['de_opt']
27         embedding_sd = checkpoint['embedding']
28         voc.__dict__ = checkpoint['voc_dict']
29
30
31     print('Building encoder and decoder ...')
32     # Initialize word embeddings
33     embedding = nn.Embedding(voc.num_words, hidden_size)
34     # if loadFilename:
35     #     embedding.load_state_dict(embedding_sd)
36     # Initialize encoder & decoder models
37     encoder = EncoderRNN(hidden_size, embedding, encoder_n_layers, dropout)
38     decoder = LuongAttnDecoderRNN(attn_model, embedding, hidden_size, voc.num_words, decoder_n_layers)
39     if loadFilename:
40         encoder.load_state_dict(encoder_sd)
41         decoder.load_state_dict(decoder_sd)
42     # Use appropriate device
43     encoder = encoder.to(device)
44     decoder = decoder.to(device)
45     print('Models built and ready to go!')
```

```

1      Building encoder and decoder ...
2      Models built and ready to go!
```

Khi muốn load một model từ check point chúng ta chỉ cần thay đổi loadFilename như bên dưới.

```

1      # Khi muốn load model thì enable đoạn dưới để tạo file lưu địa model training.
2      loadFilename = os.path.join(save_dir, model_name, corpus_name,
3                                  '{}-{}_{}'.format(encoder_n_layers, decoder_n_layers, hidden_size),
4                                  '{}_checkpoint.tar'.format(checkpoint_iter))
```

5.1. Huấn luyện model

Để huấn luyện mô hình trước tiên ta cần thiết lập các tham số. Gọi vào hàm `trainIters` để huấn luyện qua các vòng lặp.

```

1      # Configure training/optimization
2      clip = 50.0
3      teacher_forcing_ratio = 1.0
4      learning_rate = 0.0001
5      decoder_learning_ratio = 5.0
6      n_iteration = 5000
7      print_every = 100
8      save_every = 1000
9
10     # Ensure dropout layers are in train mode
11     encoder.train()
12     decoder.train()
13
14     # Initialize optimizers
15     print('Building optimizers ...')
16     encoder_optimizer = optim.Adam(encoder.parameters(), lr=learning_rate)
17     decoder_optimizer = optim.Adam(decoder.parameters(), lr=learning_rate * decoder_learning_ratio)
18     if loadFilename:
19         encoder_optimizer.load_state_dict(encoder_optimizer_sd)
20         decoder_optimizer.load_state_dict(decoder_optimizer_sd)
21
22     # Run training iterations
23     print("Starting Training!")
24     trainIters(model_name, voc, pairs, encoder, decoder, encoder_optimizer, decoder_optimizer,
25               embedding, encoder_n_layers, decoder_n_layers, save_dir, n_iteration, batch_size,
26               print_every, save_every, clip, corpus_name, loadFilename)

```

```

1      Building optimizers ...
2      Starting Training!
3      Initializing ...
4      Training...
5      Iteration: 100; Percent complete: 2.0%; Average loss: 6.7370
6      Iteration: 200; Percent complete: 4.0%; Average loss: 5.0068
7      Iteration: 300; Percent complete: 6.0%; Average loss: 3.5866
8      Iteration: 400; Percent complete: 8.0%; Average loss: 2.3804
9      Iteration: 500; Percent complete: 10.0%; Average loss: 1.8249
10     Iteration: 600; Percent complete: 12.0%; Average loss: 1.5482
11     Iteration: 700; Percent complete: 14.0%; Average loss: 1.3622
12     Iteration: 800; Percent complete: 16.0%; Average loss: 1.2253
13     Iteration: 900; Percent complete: 18.0%; Average loss: 1.1341
14     Iteration: 1000; Percent complete: 20.0%; Average loss: 1.0986
15     Iteration: 1100; Percent complete: 22.0%; Average loss: 1.0536
16     Iteration: 1200; Percent complete: 24.0%; Average loss: 1.0119
17     Iteration: 1300; Percent complete: 26.0%; Average loss: 0.9666
18     Iteration: 1400; Percent complete: 28.0%; Average loss: 0.9291
19     Iteration: 1500; Percent complete: 30.0%; Average loss: 0.9200

```

5.2 Đánh giá model

Bên dưới chúng ta cùng đánh giá mô hình thông qua việc dự đoán một số từ không dấu. Lưu ý rằng do dữ liệu được huấn luyện chỉ là các ngram có 4 từ nên chúng ta sẽ truyền vào các cụm 4 từ. Để thêm dấu cho 1 câu văn với độ dài tùy ý bạn đọc có thể chia câu thành các ngram với kích thước 4 và ghép các kết quả dự báo trên từng ngram đơn lẻ. Hàm này tôi sẽ không viết ở đây và xin dành cho bạn đọc.

```

1      # Set dropout layers to eval mode
2      encoder.eval()
3      decoder.eval()
4
5      # Initialize search module
6      searcher = GreedySearchDecoder(encoder, decoder)
7
8      # Begin chatting (uncomment and run the following line to begin)
9      evaluateInput(encoder, decoder, searcher, voc)

```

Top

```

1 > tai nguyen thien nhien
2 Bot: tài nguyên thiên nhiên
3 > thuong mai dien tu
4 Bot: thương mại điện tử
5 > hoc sinh cap 1
6 Bot: học sinh cấp 1

```

6. Hạn chế của mô hình

Mô hình seq2seq luôn tồn tại một số hạn chế nhất định mà chúng ta sẽ nhận ra trong quá trình huấn luyện đó là:

- Mô hình rất tốn tài nguyên và thời gian huấn luyện lâu. Để huấn luyện mô hình dịch máy, google đã huy động một hệ thống máy chủ mà diện tích có thể trải trên 1 km².
- Mô hình sẽ không thêm dấu chính xác đối với những cụm từ mà chúng chưa từng được học. Do đó để nâng cao mức độ chuẩn xác ngoài cần một kiến trúc mô hình mạnh thì một tập dữ liệu đủ lớn.
- Trong bài tôi sử dụng mô hình theo word level. Chính vì thế kích thước của vocabulary là rất lớn và sẽ ảnh hưởng đến số lượng parameter của model.
- Huấn luyện mô hình theo character level sẽ có lợi thế hơn khi các từ thay đổi chỉ nằm trong tập các chữ cái u eo a i d y . Do đó chúng ta sẽ teacher forcing để chỉ thay đổi các từ nằm trong tập u eo a i d y và không thay đổi các từ còn lại.
- Mô hình mới chỉ xây dựng cho các cụm từ ngram với kích thước bằng 4. Để dự báo cho câu với độ dài tùy ý dựa trên ngram = 4 xin dành cho bạn đọc.
- Lớp model transformer cũng là một trong những mô hình cân nhắc thay thế cho seq2seq trong bài toán này vì tốc độ tính toán nhanh, có thể xử lý trên nhiều GPU. Bạn đọc có thể tham khảo ý tưởng của đội xếp 2nd của cuộc thi thêm dấu tiếng việt.

7. Tài liệu tham khảo

1. hướng dẫn pytorch (<https://phamdinhhkhanh.github.io/2019/08/10/PytorchTutorial1.html>)
2. lý thuyết về mạng LSTM (https://phamdinhhkhanh.github.io/2019/04/22/L%C3%BD_thuy%E1%BA%BFT_v%E1%BB%81_m%E1%BA%A1ng_LSTM.htm)
3. seq2seq pytorch (https://pytorch.org/tutorials/beginner/chatbot_tutorial.html)
4. tensor2tensor project (<https://github.com/tensorflow/tensor2tensor>)
5. attention là tất cả bạn cần (<https://phamdinhhkhanh.github.io/2019/06/18/AttentionLayer.html>)
6. thêm dấu cho tiếng việt - aivivn 1st (<https://forum.machinelearningcoban.com/t/aivivn-3-vietnamese-tone-prediction-1st-place-solution/5721>)
7. thêm dấu cho tiếng việt - aivivn 2nd (<https://forum.machinelearningcoban.com/t/aivivn-3-vietnamese-tone-prediction-2nd-place-solution/5759>)