

به نام خدا
درس یادگیری عمیق
تمرین پنجم

غزل زمانی نژاد

401722244

1. رویکردهای معمولی این تسک بر اساس شبکه عصبی بازگشتی هستند و در آن کپشن به صورت کلمه

به کلمه تولید می شود و پیش بینی هر کلمه بر اساس محتوای تصویری و کلمات تولید شده قبلی

صورت می گیرد. اما بر اساس تحقیقات، در ویدئو محتوای اضافی و نامربوط قابل توجهی وجود دارد

که ممکن است باعث تداخل در تولید کپشن درست شود. ایده اصلی این مقاله ارائه یک مدل برای

تولید کپشن برای ویدئو بر اساس مکانیزم توجه است. ۴ بخش کلی این معماری عبارتند از:

Visual attention module: این ماژول برای استخراج ویژگی های تصویری به کار می رود. شامل دو لایه

است و می تواند بر روی مرتبط ترین ویژگی های تصویری تمرکز کند: لایه اول یاد می گیرد که روی

برجسته ترین مناطق در هر فریم تمرکز کند در حالیکه لایه دوم تلاش می کند به مرتبط ترین فریم ها

توجه داشته باشد. این ویژگی ها در قالب یک بردار با طول ثابت انکود می شوند.

Text attention module: در طول generation این ماژول بر روی کلمات تولید شده قبلی عمل می کند

و به طور انتخابی بر روی مرتبط ترین عبارت برای تولید کلمه فعلی تمرکز می کند.

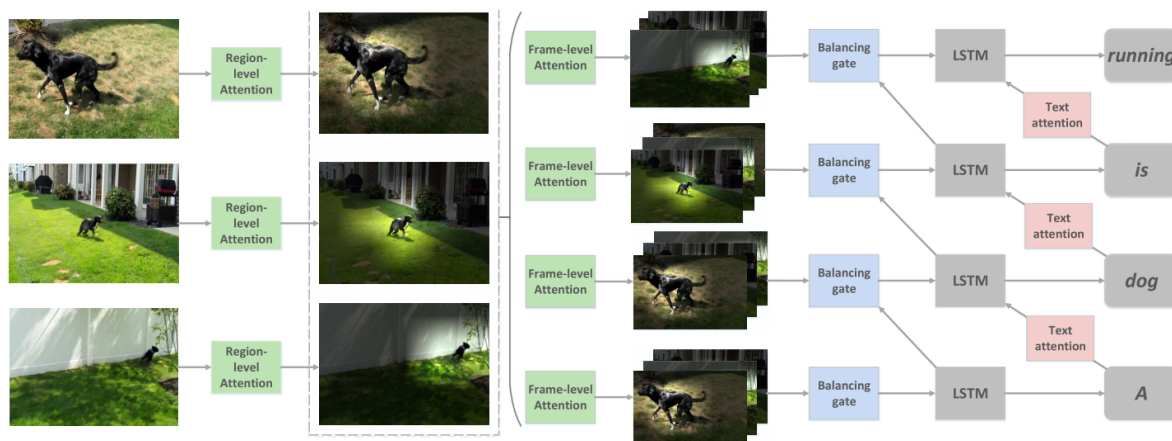
Caption generator: از بردار انکود شده ویژگی های تصویر و ویژگی های متن به همراه اطلاعاتی که

توسط LSTM ذخیره شده برای تولید کپشن به صورت کلمه به کلمه استفاده می شود.

Balancing gate: از این گیت برای تنظیم تاثیر ویژگی های تصویری و ویژگی های متن در فرآیند تولید

کپشن استفاده شده است.

در تصویر زیر معماری این مدل به خوبی مشاهده می شود:



از ایده مطرح شده در این مقاله می توان برای تولید کپشن برای تصاویر نیز استفاده کرد. با این تفاوت که در تصاویر دیگر بحث چندین فریم مطرح نیست و در ماژول visual attention می توانیم از لایه اول (که به نقاط برجسته هر فریم اشاره دارد) برای آموزش آن استفاده کنیم.

2. الف) گیت آپدیت به مدل کمک می کند تا تعیین کند چه مقدار از اطلاعات گذشته باید به آینده منتقل شود. و گیت ریست به مدل کند که چه اطلاعاتی را فراموش کند. اگر گیت ریست وجود نداشته باشد فرمول ها بدین شکل تغییر می کنند:

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{hx}\mathbf{x}_t + \mathbf{b}_h)$$

$$\mathbf{h}_t = \mathbf{u}_t \cdot \tilde{\mathbf{h}}_t + (1 - \mathbf{u}_t) \cdot \mathbf{h}_{t-1}$$

$$\mathbf{u}_t = \sigma(\mathbf{W}_{uh}\mathbf{h}_{t-1} + \mathbf{W}_{ux}\mathbf{x}_t + \mathbf{b}_u)$$

در صورت حذف این گیت شبکه تنها می تواند تصمیم بگیرد چه مقدار از اطلاعات گذشته را در استیت فعلی استفاده کند و در فراموش کردن اطلاعات قدیمی موفق نخواهد بود. مثلا در تسک تشخیص احساسات، فرض کنید فردی در ابتدای کامنت های خود در مورد یک کتاب ذکر کرده باشد: این یک کتاب فانتزی است که ... را به تصویر می کشد. و پس از چند پاراگراف، در جملات آخر گفته باشد: اما من از این کتاب به دلیل ذکر جزئیات لذت نبردم. در اینجا جمله آخر برای تشخیص حس فرد کافی است و گیت ریست می تواند در پاک کردن اطلاعات گذشته که مفید نبودند کمک کند و در صورت نبود آن کار شبکه برای تعیین کلاس نهایی سخت تر خواهد شد.

ب) مطابق فرمول های GRU برای اینکه hidden state تنها به x مرحله قبل بستگی داشته باشد باید مقدار گیت ریست ۰ و آپدیت ۱ باشد. و برای اینکه تنها به h های مرحله قبل بستگی داشته باشد، برخلاف حالت قبل باید مقدار گیت ریست ۱ و آپدیت ۰ شود.

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_{hh}(\mathbf{r}_t \cdot \mathbf{h}_{t-1}) + \mathbf{W}_{hx}\mathbf{x}_t + \mathbf{b}_h)$$

$$\mathbf{h}_t = \mathbf{u}_t \cdot \tilde{\mathbf{h}}_t + (1 - \mathbf{u}_t) \cdot \mathbf{h}_{t-1}$$

$$\mathbf{u}_t = \sigma(\mathbf{W}_{uh}\mathbf{h}_{t-1} + \mathbf{W}_{ux}\mathbf{x}_t + \mathbf{b}_u)$$

$$\mathbf{r}_t = \sigma(\mathbf{W}_{rh}\mathbf{h}_{t-1} + \mathbf{W}_{rx}\mathbf{x}_t + \mathbf{b}_r)$$

3. الف) تعداد عملیات ها به نحوه پیاده سازی دو تابع مربوط است. در اینجا فرض میکنیم در پیاده سازی softmax، مقدار e^{**x} برای هر یک از نورون ها یک بار محاسبه و ذخیره می شود. همینطور مقدار مخرج (مجموع e^{**x} ها) تنها یک بار حساب می شود.

softmax $\rightarrow s(x) = \frac{e^x}{\sum_i e^x}$

محاسبه e^x 20 تبار
 محاسبه \sum 19 جمع
 محاسبه softmax برای هر نورون 20 تقسیم
 $20 + 19 + 20 = 59$

sigmoid $\rightarrow \sigma(x) = \frac{1}{1 + e^{-x}}$

محاسبه $-x$ 20 ضرب
 محاسبه e^{-x} 20 تبار
 محاسبه جمع 20 جمع
 محاسبه sigmoid برای هر نورون 20 تقسیم
 $20 + 20 + 20 + 20 = 80$

تعداد ۲۱ عملیات اضافه می شود.

ب) با اضافه کردن لایه dense، ظرفیت یادگیری دو مدل حدوداً شبیه به هم خواهد شد. با این تفاوت که در لایه dense می توانیم بایاس و همچنین تابع فعال سازی داشته باشیم.

مزایای لایه embedding نسبت به لایه dense: این لایه می تواند به خوبی ارتباط بین کلمات را آموزش ببیند و بر این اساس یک ماتریس از امبدینگ کلمات بسازد اما به عنوان ورودی لایه dense باید بردارهای one-hot کلمات را بدهیم که در آنها اطلاعات خاصی از کلمات وجود ندارد و نمی تواند نمایانگر فاصله معنایی باشد (فاصله میان هر دو کلمه ۱ است). همچنین در انجام محاسبات لایه dense، نیاز است که ضرب ماتریسی انجام شود اما در لایه embedding، بردار وزن ها را به عنوان یک lookup

table در نظر می گیرد که سطر n ام آن بردار امبدینگ کلمه n در مجموعه کلمات است و همین کار باعث سرعت بخشیدن در فرایند آموزش می شود. بنابراین استفاده از لایه embedding به کاهش اندازه ورودی و کاهش پیچیدگی محاسبات کمک خواهد کرد.

4. الف) ابتدا یک کلاس از نوع Dataset می سازیم و در constructor از مسیر عکس ها و کپشن ها آنها را می خوانیم. داده ها را split کرده و پیش پردازش های مربوطه را با تابع preprocess_caption انجام می دهیم. " را از ابتدا و انتهای کپشن حذف می کنیم و تمامی کلمات را به lower case تبدیل می کنیم. از دیکشنری img_caption_dict برای نگهداری هر ۵ کپشن یک تصویر استفاده می کنیم. همچنین مجموعه کلمات را با تابع build_vocab می سازیم و در دیکشنری مربوطه نگه می داریم. در دیکشنری vocab کلید، توکن و مقدار، ایندکس است. دیکشنری idx_to_vocab برعکس دیکشنری قبلی است یعنی کلید، ایندکس و مقدار، توکن است.

```
2
3 PAD_VALUE = 0
4
5 class MyDataset(Dataset):
6     def __init__(self, img_dir, annotation_dir, transform=None, tokenizer="spacy"):
7         super(MyDataset, self).__init__()
8         #####
9         # your code here, you can add new parameter in constructor
10        #####
11        self.img_dir = img_dir
12        self.transform = transform
13        self.tokenizer = get_tokenizer(tokenizer)
14
15        with open(annotation_dir) as f:
16            self.img_caption_pair = f.readlines()
17            self.img_caption_pair = self.img_caption_pair[1:]
18            # each element of list is a tuple: (img_name, caption)
19            self.img_caption_pair = [self.split_img_caption(line) for line in self.img_caption_pair]
20            # dict[img] = [cap1, cap2, cap3, cap4, cap5]
21            self.img_caption_dict = {}
22            self.build_img_cap_dict()
23
24            self.vocab = {"<PAD>": PAD_VALUE, "<SOS>": 1, "<EOS>": 2, "<UNK>": 3}
25
26            self.build_vocab()
27            self.idx_to_vocab = {self.vocab[k]:k for k in self.vocab.keys()}
28
```

```

29
30 def preprocess_caption(self, caption):
31     # remove " from the beginning and the end of a caption
32     if caption.startswith('"'):
33         caption = caption[1:]
34     if caption.endswith('"'):
35         caption = caption[:-1]
36
37     # convert characters to lower case
38     caption = caption.lower()
39     return caption
40
41
42 def split_img_caption(self, l):
43     img, *cap = l.strip().split(',')
44     # if the caption itself has comma, join the other parts of it
45     cap = ','.join(cap)
46     cap = self.preprocess_caption(cap)
47
48     return (img, cap)
49
50
51 def build_img_cap_dict(self):
52     for img, cap in self.img_caption_pair:
53         if img in self.img_caption_dict.keys():
54             self.img_caption_dict[img].append(cap)
55         else: self.img_caption_dict[img] = [cap]
56
57
58 def build_vocab(self):
59     idx = 4
60     for _, c in self.img_caption_pair:
61         tokens = self.tokenizer(c)
62         for t in tokens:
63             if t not in self.vocab.keys():
64                 self.vocab[t] = idx
65                 idx += 1
66
67
68 def build_vocab_from_words(self, words):
69     self.vocab = {}
70     c = 0
71     for w in words:
72         self.vocab[w] = c
73         c += 1
74
75     # index_to_vocab dictionary should be updated
76     self.idx_to_vocab = {self.vocab[k]:k for k in self.vocab.keys()}
77

```

در تابع `get_item` ایندکس را به عنوان ورودی دریافت کرده و عکس مربوطه را می خوانیم و `transform` را روی آن اعمال می کنیم. سپس کپشن را توکنایز می کنیم. توکن های `<EOS>` و `<SOS>` را به ابتدا و انتهای کپشن اضافه می کنیم و اگر کلمه ای داخل `vocab` قرار نداشته باشد به جای آن توکن `<UNK>` می گذاریم. در نهایت بردار تصویر و بردار ایندکس کلمات کپشن را برمی گردانیم.

```

79 def __len__(self):
80     #####
81     #your code here, you should be return size of vocabulary here
82     #####
83
84     # return len(self.vocab)
85     return len(self.img_caption_pair)
86
87 def __getitem__(self, index):
88     #####
89     #your code here, you should be retrun image and caption of th
90     #####
91
92     # read image
93     img_name, caption = self.img_caption_pair[index]
94     img = cv2.imread(os.path.join(self.img_dir, img_name))
95     if self.transform:
96         img = self.transform(Image.fromarray(img))
97
98     # tokenize caption and get its index
99     tokenized_caption = self.tokenizer(caption)
100
101     indices = [self.vocab["<SOS>"]]
102     for t in tokenized_caption:
103         if t in self.vocab.keys():
104             idx = self.vocab[t]
105         else: idx = self.vocab["<UNK>"]
106         indices.append(idx)
107
108     indices.append(self.vocab["<EOS>"])
109     return img, indices
110

```

برای اینکه داده ها را در قالب batch به عنوان ورودی مدل بدهیم، باید طول کپشن ها یکسان باشد.

برای این کار از تابع padify برای پد کردن استفاده می کنیم. همچنین با استفاده از transform تصاویر را

به ۲۲۴ * ۲۲۴ (اندازه ورودی رزنت) resize و سپس نرمالیزه می کنیم.

```

1 def padify(batch):
2     images = []
3     captions = []
4
5     for i, c in batch:
6         images.append(i)
7         captions.append(torch.tensor(c))
8
9     padded_caps = pad_sequence(captions, batch_first=False, padding_value=PAD_VALUE)
10    return (torch.stack(images),
11            padded_caps)
12
13
14 transform = transforms.Compose(
15     [
16         transforms.Resize((224, 224)),
17         transforms.ToTensor(),
18         transforms.Normalize((0.5, 0.5, 0.5), (0.5, 0.5, 0.5)),
19     ]
20 )

```

دیتاست را می سازیم و از آن برای ساختن data loader آموزش و ارزیابی استفاده می کنیم. از batch

size = 64 و نسبت ۰.۸ برای داده آموزشی و ۰.۲ برای داده ارزیابی استفاده می کنیم.

```
1 dataset = MyDataset("/content/flickr8k/images", "/content/flickr8k/captions.txt", transform=transform)
2
3 #####
4 #your code here
5 #####
6 batch_size = 64
7
8 train_split = 0.8
9 train_size = int(train_split * len(dataset))
10 val_size = len(dataset) - train_size
11
12 train_set, val_set = random_split(dataset, [train_size, val_size])
13
14 train_loader = torch.utils.data.DataLoader(train_set, batch_size=batch_size, shuffle=True, collate_fn=padify)
15 val_loader = torch.utils.data.DataLoader(val_set, batch_size=batch_size, shuffle=False, collate_fn=padify)

/usr/local/lib/python3.10/dist-packages/torchtext/data/utils.py:105: UserWarning: Spacy model "en" could not be l
warnings.warn(

1 print(len(dataset))
2 print(len(dataset.vocab))
3 print(len(train_loader.dataset))
4 print(len(val_loader.dataset))

40455
8508
32364
8091
```

ب) در مدل دو تغییر اعمال می کنیم: افزودن padding_idx به لایه امبدینگ و افزودن تابع

predict_caption برای پیش بینی داده جدید (ورودی این تابع تنها تصویر خواهد بود از برای پیش بینی

کپشن، از کلمات قبلی که خود پیش بینی کرده استفاده می کند. یعنی X پیش بینی شده در این مرحله

به عنوان ورودی lstm در مرحله بعد مورد استفاده قرار می گیرد.


```

1 class ImageCaptioning(nn.Module):
2     def __init__(self, vocab_size, embed_size, hidden_size, num_layers):
3         super(ImageCaptioning, self).__init__()
4         self.vocab_size = vocab_size
5         self.embed_size = embed_size
6         self.hidden_size = hidden_size
7         self.num_layers = num_layers
8         # backbone for feature extract
9         self.featuresCNN = models.resnet50(pretrained=True)
10        # convert features to feature vector
11        for param in self.featuresCNN.parameters():
12            param.requires_grad = False
13        self.featuresCNN.fc = nn.Linear(self.featuresCNN.fc.in_features, embed_size)
14        self.fc = nn.Linear(hidden_size, vocab_size)
15        #activation function
16        self.relu = nn.ReLU()
17        #RNN
18        # I added padding idx
19        self.embed = nn.Embedding(self.vocab_size, self.embed_size, padding_idx=PAD_VALUE)
20        self.lstm = nn.LSTM(self.embed_size, self.hidden_size, self.num_layers)
21        self.linear = nn.Linear(self.hidden_size, self.vocab_size)
22
23    def forward(self, images, captions):
24        features = self.featuresCNN(images)
25        features = self.relu(features)
26        embeddings = self.embed(captions)
27        embeddings = torch.cat((features.unsqueeze(0), embeddings), dim=0)
28        hiddens, _ = self.lstm(embeddings)
29        outputs = self.fc(hiddens)
30        return outputs
31
32    def predict_caption(self, image, idx_to_vocab, maxlength=40):
33        result_caption = []
34
35        with torch.no_grad():
36            x = self.featuresCNN(image)
37            x = self.relu(x).unsqueeze(0)
38            states = None
39
40            for _ in range(maxlength):
41                hiddens, states = self.lstm(x, states)
42                output = self.fc(hiddens.squeeze(0))
43                predicted = output.argmax(1)
44
45                result_caption.append(predicted.item())
46                x = self.embed(predicted).unsqueeze(0)
47
48                if idx_to_vocab[predicted.item()] == "<EOS>":
49                    break
50        return [idx_to_vocab[i] for i in result_caption]
51

```

هایپرپارامترها را به صورت زیر تنظیم می کنیم:

```

1 device = 'cuda' if torch.cuda.is_available() else 'cpu'
2 print("device:", device)
3
4 num_epochs = 10
5 vocab_size = len(dataset.vocab)
6 embed_size = 256
7 hidden_size = 256
8 num_layers = 3
9 learning_rate = 1e-3

```

device: cuda

سپس مدل را در ۱۰ اپاک با بهینه ساز Adam و تابع ضرر cross entropy آموزش می دهیم.

```
5 model = ImageCaptioning(vocab_size, embed_size, hidden_size, num_layers)
6 model.to(device)
7 optimizer = optim.Adam(model.parameters(), lr=learning_rate)
8 criterion = nn.CrossEntropyLoss()
9
10 for epoch in range(num_epochs):
11     print(f"Epoch {epoch+1}")
12
13     # train loop
14     running_loss = 0
15     model.train()
16     # CHANGE to train_loader
17     for idx, (imgs, captions) in enumerate(tqdm(train_loader)):
18         #####
19         #your code here, complete train loop
20         #####
21         imgs = imgs.to(device)
22         captions = captions.to(device)
23
24         pred = model(imgs, captions[:-1])
25         optimizer.zero_grad()
26         loss = criterion(pred.reshape(-1, pred.shape[2]), captions.reshape(-1))
27
28         loss.backward()
29         optimizer.step()
30         running_loss += loss.item()
31
32     running_loss /= len(train_loader)
33     print("Train Loss:", running_loss)
34
35     # evaluation loop
36     eval_loss = 0
37     model.eval()
38     with torch.no_grad():
39         for idx, (imgs, captions) in enumerate(tqdm(val_loader)):
40             #####
41             # your code here, complete validation loop
42             #####
43             imgs = imgs.to(device)
44             captions = captions.to(device)
45
46             pred = model(imgs, captions[:-1])
47             loss = criterion(pred.reshape(-1, pred.shape[2]), captions.reshape(-1))
48             eval_loss += loss.item()
49     eval_loss /= len(val_loader)
50     print("Eval Loss:", eval_loss)
```

خروجی به صورت زیر است:

```

100%|██████████| 506/506 [06:08<00:00, 1.37it/s]
Train Loss: 2.7491757448011707
100%|██████████| 127/127 [01:25<00:00, 1.49it/s]
Eval Loss: 2.261204177000391
*****
Epoch 2
100%|██████████| 506/506 [05:56<00:00, 1.42it/s]
Train Loss: 2.062711037194776
100%|██████████| 127/127 [01:24<00:00, 1.51it/s]
Eval Loss: 1.9737005646773211
*****
Epoch 3
100%|██████████| 506/506 [05:51<00:00, 1.44it/s]
Train Loss: 1.8527477378901758
100%|██████████| 127/127 [01:24<00:00, 1.50it/s]
Eval Loss: 1.8453665051873274
*****
Epoch 4
100%|██████████| 506/506 [05:54<00:00, 1.43it/s]
Train Loss: 1.7473391402851453
100%|██████████| 127/127 [01:23<00:00, 1.51it/s]
Eval Loss: 1.7636045930892463
*****
Epoch 5
100%|██████████| 506/506 [05:50<00:00, 1.44it/s]
Train Loss: 1.6457291823837596
100%|██████████| 127/127 [01:24<00:00, 1.51it/s]
Eval Loss: 1.6854843483196469
*****
Epoch 6
100%|██████████| 506/506 [05:49<00:00, 1.45it/s]
Train Loss: 1.554649254785696
100%|██████████| 127/127 [01:24<00:00, 1.50it/s]
Eval Loss: 1.6305737195052499
*****
Epoch 7
100%|██████████| 506/506 [05:54<00:00, 1.43it/s]
Train Loss: 1.4803795881657733
100%|██████████| 127/127 [01:24<00:00, 1.50it/s]
Eval Loss: 1.587299199554864
*****
Epoch 8
100%|██████████| 506/506 [05:51<00:00, 1.44it/s]
Train Loss: 1.412818374252131
100%|██████████| 127/127 [01:23<00:00, 1.52it/s]
Eval Loss: 1.559513425263833
*****
Epoch 9
100%|██████████| 506/506 [05:49<00:00, 1.45it/s]
Train Loss: 1.3574818484632394
100%|██████████| 127/127 [01:23<00:00, 1.53it/s]
Eval Loss: 1.5363937991810597
*****
Epoch 10
100%|██████████| 506/506 [05:48<00:00, 1.45it/s]
Train Loss: 1.30267662481357
100%|██████████| 127/127 [01:24<00:00, 1.50it/s] Eval Loss: 1.523096230555707

```

مقدار loss در حین آموزش کاهش یافته و به مقدار ۱.۳۰ برای داده آموزشی و ۱.۵۲ برای داده ارزیابی رسیده است.

پس از پایان آموزش از تابع `corpus_blue` برای محاسبه امتیاز BLEU استفاده می کنیم. وزن های آن را به گونه ای تنظیم می کنیم که برای unigram و bigram محاسبات را انجام دهد.

تابع `predict_whole_data` را پیاده سازی می کنیم تا برای همه تصاویر موجود در دیتاست پیش بینی انجام دهد و سپس امتیاز BLEU را با یک پیش بینی و ۵ رفرنس حساب می کنیم. توکن های `<EOS>` و `<SOS>` را از کپشن حذف می کنیم.

```
1 def compute_bleu(list_of_references, hypotheses):
2     weights = [
3         (1., 0, 0, 0),
4         (1./2., 1./2., 0, 0)
5     ]
6
7     return corpus_bleu(list_of_references, hypotheses, weights=weights)


1 def predict_whole_data(dataset, model):
2     references = []
3     hypotheses = []
4
5     # compute bleu on whole data
6     for img_name, caps in dataset.img_caption_dict.items():
7
8         # read image
9         img = cv2.imread(os.path.join(dataset.img_dir, img_name))
10        if transform:
11            img = transform(Image.fromarray(img))
12
13        img = img.to(device)
14        predicted_caption = model.predict_caption(img.unsqueeze(0), dataset.idx_to_vocab)
15        if '<SOS>' in predicted_caption:
16            predicted_caption.remove('<SOS>')
17        if '<EOS>' in predicted_caption:
18            predicted_caption.remove('<EOS>')
19
20        # print(predicted_caption)
21        hypotheses.append(predicted_caption)
22
23        tokenized_captions = [dataset.tokenizer(cap) for cap in caps]
24        # print(tokenized_captions)
25        references.append(tokenized_captions)
26
27    return references, hypotheses
```

امتیاز BLEU برای این بخش به صورت زیر است:

BLEU-1: 0.5575184961055694

BLEU-2: 0.3455166215514869

ج) ابتدا Glove را دانلود می کنیم. از وزن ها با ابعاد ۲۰۰ استفاده می کنیم. فایل مربوطه را می خوانیم و کلمات آن را داخل لیست vocab و وزن ها را داخل لیست embeddings ذخیره می کنیم. سپس این دو لیست را به آرایه numpy تبدیل می کنیم. توکن های pad, sos, eos, unk را به مجموعه کلمات و امبدینگ اضافه می کنیم.

```
1 vocab, embeddings = [], []
2 # use embedding_size 200 from glove
3 with open('glove.6B.200d.txt', 'rt') as f:
4     full_content = f.read().strip().split('\n')
5
6 for i in range(len(full_content)):
7     i_word = full_content[i].split(' ')[0]
8     i_embeddings = [float(val) for val in full_content[i].split(' ')[1:]]
9     vocab.append(i_word)
10    embeddings.append(i_embeddings)

1 vocab_np = np.array(vocab)
2 embeds_np = np.array(embeddings)
3
4 #insert '<PAD>' and '<UNK>' tokens at start of vocab_np.
5 vocab_np = np.insert(vocab_np, 0, '<PAD>')
6 vocab_np = np.insert(vocab_np, 1, '<SOS>')
7 vocab_np = np.insert(vocab_np, 2, '<EOS>')
8 vocab_np = np.insert(vocab_np, 3, '<UNK>')
9 # print(vocab_np[:10])
10
11 pad_emb_np = np.zeros((1, embeds_np.shape[1])) #embedding for '<PAD>' token.
12 sos_emb_np = np.random.rand(1, embeds_np.shape[1]) #embedding for '<SOS>' token.
13 eos_emb_np = np.random.rand(1, embeds_np.shape[1]) #embedding for '<EOS>' token.
14 unk_emb_np = np.mean(embeds_np, axis=0, keepdims=True) #embedding for '<UNK>' token
15
16 #insert embeddings for pad, sos, eos, unk tokens at top of embeds_np.
17 embeds_np = np.vstack((pad_emb_np, sos_emb_np, eos_emb_np, unk_emb_np, embeds_np))
18 print(embeds_np.shape)

(400004, 200)
```

دیتاست جدید و loaderهای جدید را تشکیل می دهیم.

```

1 transform = transforms.Compose(
2     [
3         transforms.Resize((224, 224)),
4         transforms.ToTensor(),
5         transforms.Normalize((0.5, 0.5, 0.5), (0.5, 0.5, 0.5)),
6     ]
7 )
8
9 dataset_glove = MyDataset("/content/flickr8k/images", "/content/flickr8k/captions.txt", transform=transform)
10
11 batch_size = 64
12
13 train_split = 0.8
14 train_size = int(train_split * len(dataset_glove))
15 val_size = len(dataset_glove) - train_size
16
17 train_set2, val_set2 = random_split(dataset_glove, [train_size, val_size])
18
19 train_loader2 = torch.utils.data.DataLoader(train_set2, batch_size=batch_size, shuffle=True, collate_fn=padify)
20 val_loader2 = torch.utils.data.DataLoader(val_set2, batch_size=batch_size, shuffle=False, collate_fn=padify)

```

برای اینکه به مشکل پر شدن RAM کولب بر نخوریم و بتوانیم آموزش را انجام دهیم، از کل ۴۰۰ هزار کلمه GLOVE استفاده نمی کنیم، بلکه از مجموعه کلماتی که در دیتای ما و GLOVE مشترک است استفاده می کنیم. برای این کار باید اشتراکات را پیدا کنیم و reduced_vocab و reduced_embeds را بسازیم. سپس مجموعه وکب را آپدیت کنیم.

```

1 # a boolean array of same shape as vocab_np
2 mask = np.isin(vocab_np, np.array(list(dataset_glove.vocab.keys())))
3 print("mask:", mask)
4 # indices of words which are common in both (glove_vocab and dataset_vocab)
5 intersection = np.argwhere(mask)
6 intersection = np.squeeze(intersection, 1)
7 print("intersection", intersection)
8
9 reduced_vocab = vocab_np[intersection]
10 reduced_embeds = embeds_np[intersection]
11 print(reduced_vocab[:20])
12
13 # update dataset vocabulary
14 dataset_glove.build_vocab_from_words(reduced_vocab)

```

mask: [True True True ... False False False]

intersection [0 1 2 ... 396597 398567 399875]

['<PAD>' '<SOS>' '<EOS>' '<UNK>' 'the' ',' '.' 'of' 'to' 'and' 'in' 'a' ' ' 's' 'for' '-' 'that' 'on' 'is' 'was']

```

1 print(len(dataset_glove.vocab))
2 print(reduced_vocab.shape)

```

7846
(7846,)

در تعریف مدل جدید، از همان مدل قبل ارث بری می کنیم با این تفاوت که برای لایه امبدینگ از وزن های pretrain استفاده می کنیم.

```
1 class ImageCaptioningPretrained(ImageCaptioning):
2     def __init__(self, vocab_size, embed_size, hidden_size, num_layers, embeddings):
3         super(ImageCaptioningPretrained, self).__init__(vocab_size, embed_size, hidden_size, num_layers)
4
5         # different from the previous part
6         self.embed = nn.Embedding.from_pretrained(torch.from_numpy(embeddings).float(), padding_idx=PAD_VALUE)
```

سپس یک مدل instantiate می کنیم و آن را در ۱۰ اپیاک آموزش می دهیم. خروجی آموزش به صورت زیر است:

```
100%|██████████| 506/506 [06:02<00:00, 1.39it/s]
Train Loss: 2.8903727321756687
100%|██████████| 127/127 [01:24<00:00, 1.49it/s]
Eval Loss: 2.3872270415148398
*****
Epoch 2
100%|██████████| 506/506 [05:52<00:00, 1.44it/s]
Train Loss: 2.160399492314682
100%|██████████| 127/127 [01:24<00:00, 1.51it/s]
Eval Loss: 2.0730672567848147
*****
Epoch 3
100%|██████████| 506/506 [05:48<00:00, 1.45it/s]
Train Loss: 1.9235348213802685
100%|██████████| 127/127 [01:24<00:00, 1.50it/s]
Eval Loss: 1.8981473286320845
*****
Epoch 4
100%|██████████| 506/506 [05:50<00:00, 1.44it/s]
Train Loss: 1.7973678368824744
100%|██████████| 127/127 [01:33<00:00, 1.36it/s]
Eval Loss: 1.8083393113819632
*****
Epoch 5
100%|██████████| 506/506 [06:00<00:00, 1.40it/s]
Train Loss: 1.7051508106261846
100%|██████████| 127/127 [01:25<00:00, 1.49it/s]
Eval Loss: 1.7323606661924227
*****
Epoch 6
100%|██████████| 506/506 [05:49<00:00, 1.45it/s]
Train Loss: 1.6367170751801592
100%|██████████| 127/127 [01:23<00:00, 1.51it/s]
Eval Loss: 1.6731525922384787
*****
Epoch 7
100%|██████████| 506/506 [05:48<00:00, 1.45it/s]
Train Loss: 1.5663855751983733
100%|██████████| 127/127 [01:24<00:00, 1.50it/s]
Eval Loss: 1.6271703308961523
*****
Epoch 8
100%|██████████| 506/506 [05:47<00:00, 1.45it/s]
Train Loss: 1.5141459580937864
100%|██████████| 127/127 [01:24<00:00, 1.50it/s]
Eval Loss: 1.5910822904016089
*****
Epoch 9
100%|██████████| 506/506 [05:51<00:00, 1.44it/s]
Train Loss: 1.4750279853702062
100%|██████████| 127/127 [01:24<00:00, 1.51it/s]
Eval Loss: 1.5633375776095653
*****
Epoch 10
100%|██████████| 506/506 [05:50<00:00, 1.44it/s]
Train Loss: 1.4285901988683483
100%|██████████| 127/127 [01:23<00:00, 1.51it/s] Eval Loss: 1.5406883438741128
```

مقدار loss در حین آموزش کاهش یافته و به مقدار ۱.۴۲ برای داده آموزشی و ۱.۵۴ برای داده ارزیابی رسیده است.

امتیاز BLEU برای این بخش به صورت زیر است:

BLEU-1: 0.5240019886572881

BLEU-2: 0.31634241107845434

د) در بخش دوم از وزن های pretrained برای لایه امبدینگ استفاده کردیم اما امتیاز BLUE تغییر چندانی نداشت و حتی کمتر شد. البته در بخش اول از $\text{embed_size} = ۲۵۶$ ولی در بخش دوم از اندازه ۲۰۰ استفاده کردیم و این ممکن است مقایسه را کمی دشوار کند. زمانی که از وزن های pretrained استفاده می کنیم:

۱. نقطه شروع آموزش بهتر خواهد بود. در حالتی که از وزن های رندوم برای آموزش استفاده می کنیم ممکن است در بعضی مینیمم های محلی گیر کنیم و وزن ها به درستی آموزش نیینند.

۲. در صورت استفاده از وزن های pretrained، می توانیم با تعداد ایپاک کمتری آن را برای دیتاست خودمان fine-tune کنیم اما در حالت رندوم به تعداد زیادی ایپاک برای آموزش بردارهای امبدینگ نیاز داریم و این کار فرآیند آموزش را کند می کند.

ه) یک کلاس جدید تعریف می کنیم و از کلاس ImageCaptioning ارث بری می کنیم. مقدار require_grad تمامی پارامترها به جز تعدادی از لایه ها را false می کنیم. و دوباره لایه امبدینگ را با وزن های از پیش آموخته GLOVE لود می کنیم.

```
1 class ImageCaptioningFinetuneResnet(ImageCaptioning):
2     def __init__(self, vocab_size, embed_size, hidden_size, num_layers, embeddings):
3         super(ImageCaptioningFinetuneResnet, self).__init__(vocab_size, embed_size, hidden_size, num_layers)
4
5         for name, param in self.featuresCNN.named_parameters():
6             if 'layer4.2' in name or 'layer4.1' in name:
7                 continue
8             param.requires_grad = False
9             # print(name)
10
11         self.featuresCNN.fc = nn.Linear(self.featuresCNN.fc.in_features, embed_size)
12
13         self.embed = nn.Embedding.from_pretrained(torch.from_numpy(embeddings).float(), padding_idx=PAD_VALUE)
14
```


یک مدل جدید می سازیم و آن را در ۱۰ اپیک آموزش می دهیم. نتیجه آموزش به صورت زیر است:

```
100%|██████████| 506/506 [06:07<00:00, 1.38it/s]
Train Loss: 2.8924215056679468
100%|██████████| 127/127 [01:25<00:00, 1.48it/s]
Eval Loss: 2.2913984229245523
*****
Epoch 2
100%|██████████| 506/506 [05:57<00:00, 1.42it/s]
Train Loss: 2.1378590271406965
100%|██████████| 127/127 [01:25<00:00, 1.49it/s]
Eval Loss: 1.9961898655403318
*****
Epoch 3
100%|██████████| 506/506 [05:55<00:00, 1.42it/s]
Train Loss: 1.9337210963837244
100%|██████████| 127/127 [01:25<00:00, 1.48it/s]
Eval Loss: 1.8522482162385474
*****
Epoch 4
100%|██████████| 506/506 [05:55<00:00, 1.42it/s]
Train Loss: 1.805910297062086
100%|██████████| 127/127 [01:25<00:00, 1.49it/s]
Eval Loss: 1.7621322363380372
*****
Epoch 5
100%|██████████| 506/506 [05:53<00:00, 1.43it/s]
Train Loss: 1.7172745810195862
100%|██████████| 127/127 [01:24<00:00, 1.50it/s]
Eval Loss: 1.6988021379380713
*****
Epoch 6
100%|██████████| 506/506 [05:54<00:00, 1.43it/s]
Train Loss: 1.6491905919647971
100%|██████████| 127/127 [01:26<00:00, 1.48it/s]
Eval Loss: 1.6562527446296271
*****
Epoch 7
100%|██████████| 506/506 [05:52<00:00, 1.44it/s]
Train Loss: 1.5866548351855145
100%|██████████| 127/127 [01:26<00:00, 1.47it/s]
Eval Loss: 1.604901394506139
*****
Epoch 8
100%|██████████| 506/506 [05:50<00:00, 1.44it/s]
Train Loss: 1.5438391149279629
100%|██████████| 127/127 [01:23<00:00, 1.52it/s]
Eval Loss: 1.5661863041674997
*****
Epoch 9
100%|██████████| 506/506 [05:52<00:00, 1.44it/s]
Train Loss: 1.482220013504443
100%|██████████| 127/127 [01:23<00:00, 1.52it/s]
Eval Loss: 1.536145264708151
*****
Epoch 10
100%|██████████| 506/506 [05:49<00:00, 1.45it/s]
Train Loss: 1.4399805086639088
100%|██████████| 127/127 [01:24<00:00, 1.51it/s] Eval Loss: 1.5082164923037131
```

مقدار loss در حین آموزش کاهش یافته و به مقدار ۱.۴۴ برای داده آموزشی و ۱.۵۰ برای داده ارزیابی رسیده است.

امتیاز BLEU برای این بخش به صورت زیر است:

BLEU-1: 0.5225900808585162

BLEU-2: 0.31742193678211106

مقدار امتیاز BLEU در بخش ه و ج تقریباً برابر است. Fine-tune کردن این لایه تأثیر چندانی در استخراج ویژگی های تصویر نداشته. البته مقدار loss روی داده ارزیابی کمتر شده اما در BLEU اختلافی مشاهده نمی شود. باید ذکر کنیم که BLEU در محاسبات تنها به n-gram ها توجه دارد و تعداد شباهت ها را محاسبه می کند. بهتر است از semantic metrics (مثل SPICE) برای بررسی کپشن پیش بینی شده و کپشن های رفرنس استفاده کنیم تا بتواند از لحاظ معنایی مقایسه را انجام دهد.

و) در صورتی که gradient vanishing داشته باشیم، یعنی مقدار گرادیان ها بسیار کوچک شده و در نتیجه پارامترها نمی توانند به خوبی آپدیت شوند و احتمالاً loss به خوبی کاهش نمی یابد اما در هر ۳ قسمت کاهش loss مشاهده می شود و احتمالاً دچار gradient vanishing نشده اند. البته برای بررسی بیشتر می توانیم از kernel weights distribution استفاده کنیم. اگر وزن ها کم کم به صفر میل کنند، gradient vanishing رخ داده است.

ز) مزایای استفاده از امبدینگ در سطح کاراکتر:

اندازه ماتریس امبدینگ به طور چشمگیری کاهش می یابد. در هر زبان تنها تعدادی کاراکتر داریم در صورتی که تعداد کلمات آن ممکن است هزار-میلیون باشد. در نتیجه به حافظه بسیار کمتری برای ذخیره ماتریس امبدینگ در سطح کاراکتر نیاز داریم.

تعداد کاراکترهای یک زبان ثابت است اما ممکن است با گذر زمان کلمات جدیدی به vocabulary یک زبان اضافه شود و در این صورت ماتریس امبدینگ مرتباً آپدیت و سایز آن بیشتر می شود. اما در سطح کاراکتر می توانیم کلمات جدید تولید کنیم.

معایب استفاده از امبدینگ در سطح کاراکتر:

اگر به خوبی آموزش نبیند، ممکن است دنباله ای کاراکترها را پشت هم بچیند و کلمات بی معنی تولید کند.

میان بردارهای کاراکترها فاصله معنایی وجود ندارد در صورتی که در سطح کلمه، فاصله بین بردارها می تواند اطلاعات خوبی به ما بدهد.

طول دنباله ورودی می تواند بسیار زیاد شود و این باعث کاهش سرعت محاسبات می شود.

منابع

<https://towardsdatascience.com/understanding-gru-networks-2ef37df6c9be>

<https://medium.com/logivan/neural-network-embedding-and-dense-layers-whats-the-difference-fa177c6d0304>

<https://plainenglish.io/blog/understanding-collate-fn-in-pytorch-f9d1742647d3>

<https://medium.com/mllearning-ai/load-pre-trained-glove-embeddings-in-torch-nn-embedding-layer-in-under-2-minutes-f5af8f57416a>

<https://ai.stackexchange.com/questions/18132/how-to-detect-vanishing-gradients>

<https://www.lighttag.io/blog/character-level-NLP/>

<https://stats.stackexchange.com/questions/216000/advantage-of-character-based-language-models-over-word-based>