

变压器注解及PyTorch实现（下）

 mp.weixin.qq.com/s

训练

- 批和蓊
- 训练循环
- 训练数据和批处理
- 硬件和训练进度
- 优化器
- 正则化
 - 标签平滑

第一个例子

- 数据生成
- 损失计算
- 贪心解码

真实示例

- 数据加载
- 迭代器
- 多GPU训练
- 训练系统附加组件：BPE，搜索，平均

结果

- 聚焦可视化

摘要

训练

本节介绍模型的训练方法。

快速穿插介绍训练标准编码器解码器模型需要的一些工具。首先我们定义一个包含源和目标句子的批训练对象进行训练，同时构造构造。

批和阳离子

```
1. class Batch:
2.     "Object for holding a batch of data with mask during training."
3.     def __init__(self, src, trg=None, pad=0):
4.         self.src = src
5.         self.src_mask = (src != pad).unsqueeze(-2)
6.         if trg is not None:
7.             self.trg = trg[:, :-1]
8.             self.trg_y = trg[:, 1:]
9.             self.trg_mask = \
10.                 self.make_std_mask(self.trg, pad)
11.             self.ntokens = (self.trg_y != pad).data.sum()
12.     @staticmethod
13.     def make_std_mask(tgt, pad):
14.         "Create a mask to hide padding and future words."
15.         tgt_mask = (tgt != pad).unsqueeze(-2)
16.         tgt_mask = tgt_mask & Variable(
17.             subsequent_mask(tgt.size(-1)).type_as(tgt_mask.data))
18.         return tgt_mask
```

接下来，我们创建一个通用的训练和叠加函数来跟踪损失。我们放置一个通用的损失计算函数，它也处理参数更新。

训练循环

```

1. def run_epoch(data_iter, model, loss_compute):
2.     "Standard Training and Logging Function"
3.     start = time.time()
4.     total_tokens = 0
5.     total_loss = 0
6.     tokens = 0
7.     for i, batch in enumerate(data_iter):
8.         out = model.forward(batch.src, batch.trg,
9.                             batch.src_mask, batch.trg_mask)
10.        loss = loss_compute(out, batch.trg_y, batch.ntokens)
11.        total_loss += loss
12.        total_tokens += batch.ntokens
13.        tokens += batch.ntokens
14.        if i % 50 == 1:
15.            elapsed = time.time() - start
16.            print("Epoch Step: %d Loss: %f Tokens per Sec: %f" %
17.                  (i, loss / batch.ntokens, tokens / elapsed))
18.            start = time.time()
19.            tokens = 0
20.    return total_loss / total_tokens

```

训练数据和批处理

我们使用标准WMT 2014英语-德语数据集进行了训练，该数据集包含大约450万个句子对。使用字节对的编码方法对句子进行编码，该编码具有大约37000个词的共享源-目标词汇表。对于英语-法语，我们使用了WMT 2014 英语-法语数据集，该数据集由36M个句子组成，并将词分成32000个词片(Word-piece)的词汇表。

句子对按照近似的序列长度进行批处理。每个训练批包含一组句子对，包含大约25000个源词和25000个目标词。

我们将使用torch text来创建批次。下面更详细地讨论实现过程。我们在torchtext的一个函数中创建批次，确保填充到最大批训练长度的大小不超过阈值（如果我们有8个GPU，则阈值为25000）。

```

1. global max_src_in_batch, max_tgt_in_batch

2. def batch_size_fn(new, count, sofar):

3.     "Keep augmenting batch and calculate total number of tokens + padding."

4.     global max_src_in_batch, max_tgt_in_batch

5.     if count == 1:

6.         max_src_in_batch = 0

7.         max_tgt_in_batch = 0

8.     max_src_in_batch = max(max_src_in_batch, len(new.src))

9.     max_tgt_in_batch = max(max_tgt_in_batch, len(new.trg) + 2)

10.    src_elements = count * max_src_in_batch

11.    tgt_elements = count * max_tgt_in_batch

12.    return max(src_elements, tgt_elements)

```

硬件和训练进度

我们在一台配备8个NVIDIA P100 GPU的机器上训练我们的模型。对于使用本文所述的超参数的基本模型，每个训练单步大约需要0.4秒。我们对基础模型进行了总共100,000步或12小时的训练。对于我们的大型模型，每个训练单步时间为1.0秒。大型模型通常需要训练300,000步（3.5天）。

优化器

我们选择Adam[1]作为优化器，其参数为 $\beta_1 = 0.9$ 、 $\beta_2 = 0.98$ 和 $\epsilon = 10^{-9}$ 。根据以下公式，我们在训练过程中改变了学习率：

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

。在预热中随步数线性地增加学习速率，并且此后与步数的反平方根成比例地减小它。我们设置预热步数为4000。

｜ 注意：这部分非常重要，需要这种设置训练模型。

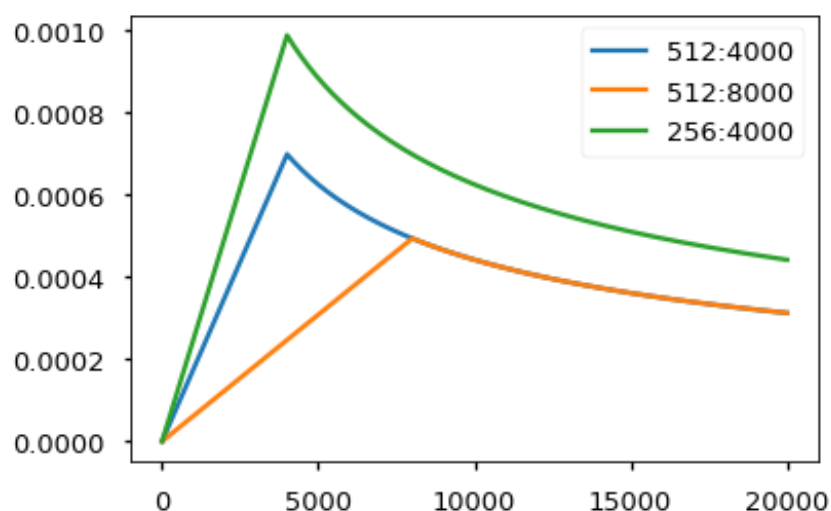
```

1. class NoamOpt:
2.     "Optim wrapper that implements rate."
3.     def __init__(self, model_size, factor, warmup, optimizer):
4.         self.optimizer = optimizer
5.         self._step = 0
6.         self.warmup = warmup
7.         self.factor = factor
8.         self.model_size = model_size
9.         self._rate = 0
10.    def step(self):
11.        "Update parameters and rate"
12.        self._step += 1
13.        rate = self.rate()
14.        for p in self.optimizer.param_groups:
15.            p['lr'] = rate
16.        self._rate = rate
17.        self.optimizer.step()
18.    def rate(self, step = None):
19.        "Implement `lr` above"
20.        if step is None:
21.            step = self._step
22.        return self.factor * \
23.            (self.model_size ** (-0.5) *
24.             min(step ** (-0.5), step * self.warmup ** (-1.5)))
25.    def get_std_opt(model):
26.        return NoamOpt(model.src_embed[0].d_model, 2, 4000,
27.            torch.optim.Adam(model.parameters(), lr=0, betas=(0.9, 0.98), eps=1e-9))

```

当前模型在不同模型大小和超参数的情况下的曲线示例。

1. # Three settings of the lr rate hyperparameters.
2. `opts = [NoamOpt(512, 1, 4000, None),`
3. `NoamOpt(512, 1, 8000, None),`
4. `NoamOpt(256, 1, 4000, None)]`
5. `plt.plot(np.arange(1, 20000), [[opt.rate(i) for opt in opts] for i in range(1, 20000)])`
6. `plt.legend(["512:4000", "512:8000", "256:4000"])`
7. `None`



正则化

标签平滑

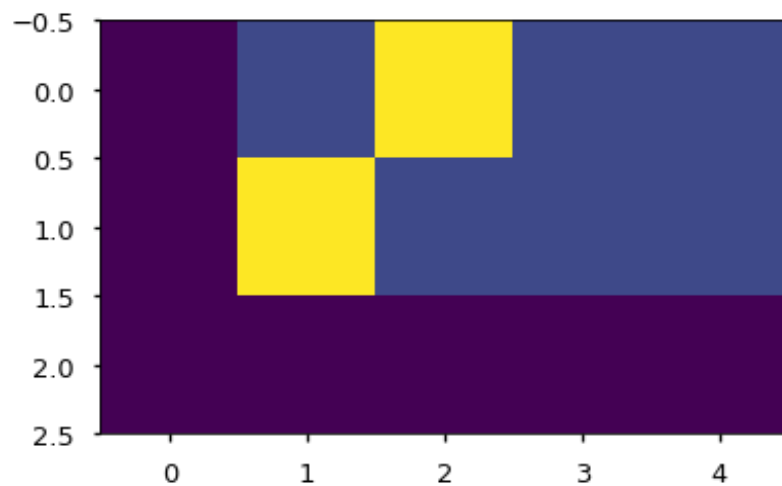
在训练期间，我们采用了值 $\epsilon_{ls} = 0.1$ [2] 的标签平滑。这种做法提高了困惑度，因为模型变得更加不确定，但提高了准确性和BLEU分数。

我们使用KL div loss实现标签平滑。相比使用独热目标分布，我们创建一个分布，其包含正确单词的置信度和整个词汇表中分布的其余平滑项。

```
1. class LabelSmoothing(nn.Module):
2.     "Implement label smoothing."
3.     def __init__(self, size, padding_idx, smoothing=0.0):
4.         super(LabelSmoothing, self).__init__()
5.         self.criterion = nn.KLDivLoss(size_average=False)
6.         self.padding_idx = padding_idx
7.         self.confidence = 1.0 - smoothing
8.         self.smoothing = smoothing
9.         self.size = size
10.        self.true_dist = None
11.        def forward(self, x, target):
12.            assert x.size(1) == self.size
13.            true_dist = x.data.clone()
14.            true_dist.fill_(self.smoothing / (self.size - 2))
15.            true_dist.scatter_(1, target.data.unsqueeze(1), self.confidence)
16.            true_dist[:, self.padding_idx] = 0
17.            mask = torch.nonzero(target.data == self.padding_idx)
18.            if mask.dim() > 0:
19.                true_dist.index_fill_(0, mask.squeeze(), 0.0)
20.            self.true_dist = true_dist
21.            return self.criterion(x, Variable(true_dist, requires_grad=False))
```

在这里，我们可以看到标签平滑的示例。

1. # Example of label smoothing.
2. crit = LabelSmoothing(5, 0, 0.4)
3. predict = torch.FloatTensor([[0, 0.2, 0.7, 0.1, 0],
4. [0, 0.2, 0.7, 0.1, 0],
5. [0, 0.2, 0.7, 0.1, 0]])
6. v = crit(Variable(predict.log()),
7. Variable(torch.LongTensor([2, 1, 0])))
8. # Show the target distributions expected by the system.
9. plt.imshow(crit.true_dist)
10. None



如果对给定的选择非常有信心，标签平滑实际上会开始惩罚模型。


```

1. crit = LabelSmoothing(5, 0, 0.1)

2. def loss(x):

3.     d = x + 3 * 1

4.     predict = torch.FloatTensor([[0, x / d, 1 / d, 1 / d, 1 / d],

5.                                   ])

6.     #print(predict)

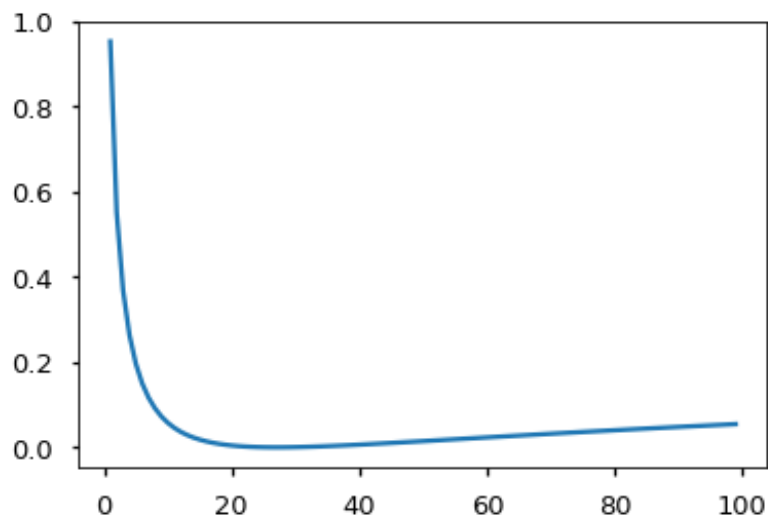
7.     return crit(Variable(predict.log()),

8.                 Variable(torch.LongTensor([1]))).data[0]

9. plt.plot(np.arange(1, 100), [loss(x) for x in range(1, 100)])

10. None

```



第一个例子

我们可以先尝试一个简单的复制任务。给定来自小词汇表的随机输入符号集，目标是生成那些相同的符号。

数据生成

```
1. def data_gen(V, batch, nbatches):
2.     "Generate random data for a src-tgt copy task."
3.     for i in range(nbatches):
4.         data = torch.from_numpy(np.random.randint(1, V, size=(batch, 10)))
5.         data[:, 0] = 1
6.         src = Variable(data, requires_grad=False)
7.         tgt = Variable(data, requires_grad=False)
8.         yield Batch(src, tgt, 0)
```

损失计算

```
1. class SimpleLossCompute:
2.     "A simple loss compute and train function."
3.     def __init__(self, generator, criterion, opt=None):
4.         self.generator = generator
5.         self.criterion = criterion
6.         self.opt = opt
7.     def __call__(self, x, y, norm):
8.         x = self.generator(x)
9.         loss = self.criterion(x.contiguous().view(-1, x.size(-1)),
10.                                y.contiguous().view(-1)) / norm
11.         loss.backward()
12.         if self.opt is not None:
13.             self.opt.step()
14.             self.opt.optimizer.zero_grad()
15.         return loss.data[0] * norm
```

贪心解码

1. # Train the simple copy task.
2. V = 11
3. criterion = LabelSmoothing(size=V, padding_idx=0, smoothing=0.0)
4. model = make_model(V, V, N=2)
5. model_opt = NoamOpt(model.src_embed[0].d_model, 1, 400,
6. torch.optim.Adam(model.parameters(), lr=0, betas=(0.9, 0.98), eps=1e-9))
7. for epoch in range(10):
8. model.train()
9. run_epoch(data_gen(V, 30, 20), model,
10. SimpleLossCompute(model.generator, criterion, model_opt))
11. model.eval()
12. print(run_epoch(data_gen(V, 30, 5), model,
13. SimpleLossCompute(model.generator, criterion, None)))

1. Epoch Step: 1 Loss: 3.023465 Tokens per Sec: 403.074173
2. Epoch Step: 1 Loss: 1.920030 Tokens per Sec: 641.689380
3. 1.9274832487106324
4. Epoch Step: 1 Loss: 1.940011 Tokens per Sec: 432.003378
5. Epoch Step: 1 Loss: 1.699767 Tokens per Sec: 641.979665
6. 1.657595729827881
7. Epoch Step: 1 Loss: 1.860276 Tokens per Sec: 433.320240
8. Epoch Step: 1 Loss: 1.546011 Tokens per Sec: 640.537198
9. 1.4888023376464843
10. Epoch Step: 1 Loss: 1.682198 Tokens per Sec: 432.092305
11. Epoch Step: 1 Loss: 1.313169 Tokens per Sec: 639.441857
12. 1.3485562801361084
13. Epoch Step: 1 Loss: 1.278768 Tokens per Sec: 433.568756
14. Epoch Step: 1 Loss: 1.062384 Tokens per Sec: 642.542067
15. 0.9853351473808288

16. Epoch Step: 1 Loss: 1.269471 Tokens per Sec: 433.388727
17. Epoch Step: 1 Loss: 0.590709 Tokens per Sec: 642.862135
18. 0.5686767101287842
19. Epoch Step: 1 Loss: 0.997076 Tokens per Sec: 433.009746
20. Epoch Step: 1 Loss: 0.343118 Tokens per Sec: 642.288427
21. 0.34273059368133546
22. Epoch Step: 1 Loss: 0.459483 Tokens per Sec: 434.594030
23. Epoch Step: 1 Loss: 0.290385 Tokens per Sec: 642.519464
24. 0.2612409472465515
25. Epoch Step: 1 Loss: 1.031042 Tokens per Sec: 434.557008
26. Epoch Step: 1 Loss: 0.437069 Tokens per Sec: 643.630322
27. 0.4323212027549744
28. Epoch Step: 1 Loss: 0.617165 Tokens per Sec: 436.652626
29. Epoch Step: 1 Loss: 0.258793 Tokens per Sec: 644.372296
30. 0.27331129014492034

为简单起见，此代码使用贪心解码来预测翻译。

```

1. def greedy_decode(model, src, src_mask, max_len, start_symbol):
2.     memory = model.encode(src, src_mask)
3.     ys = torch.ones(1, 1).fill_(start_symbol).type_as(src.data)
4.     for i in range(max_len-1):
5.         out = model.decode(memory, src_mask,
6.                             Variable(ys),
7.                             Variable(subsequent_mask(ys.size(1))
8.                                     .type_as(src.data)))
9.         prob = model.generator(out[:, -1])
10.        _, next_word = torch.max(prob, dim = 1)
11.        next_word = next_word.data[0]
12.        ys = torch.cat([ys,
13.                        torch.ones(1, 1).type_as(src.data).fill_(next_word)], dim=1)
14.    return ys
15. model.eval()
16. src = Variable(torch.LongTensor([[1,2,3,4,5,6,7,8,9,10]]))
17. src_mask = Variable(torch.ones(1, 1, 10))

1. print(greedy_decode(model, src, src_mask, max_len=10, start_symbol=1))
2.  1  2  3  4  5  6  7  8  9  10
3. [torch.LongTensor of size 1x10]

```

真实示例

现在我们通过IWSLT德语-英语翻译任务介绍一个真实示例。该任务比上文提及的WMT任务小得多，但它说明了整个系统。我们还展示了如何使用多个GPU处理加速其训练。

1. `#!pip install torchtext spacy`
2. `#!python -m spacy download en`
3. `#!python -m spacy download de`

数据加载

我们将使用torchtext和spacy加载数据集以进行词语切分。

```
1. # For data loading.
2. from torchtext import data, datasets
3. if True:
4.     import spacy
5.     spacy_de = spacy.load('de')
6.     spacy_en = spacy.load('en')
7.     def tokenize_de(text):
8.         return [tok.text for tok in spacy_de.tokenizer(text)]
9.     def tokenize_en(text):
10.        return [tok.text for tok in spacy_en.tokenizer(text)]
11.    BOS_WORD = '<s>'
12.    EOS_WORD = '</s>'
13.    BLANK_WORD = "<blank>"
14.    SRC = data.Field(tokenize=tokenize_de, pad_token=BLANK_WORD)
15.    TGT = data.Field(tokenize=tokenize_en, init_token = BOS_WORD,
16.                    eos_token = EOS_WORD, pad_token=BLANK_WORD)
17.    MAX_LEN = 100
18.    train, val, test = datasets.IWSLT.splits(
19.        exts=('.de', '.en'), fields=(SRC, TGT),
20.        filter_pred=lambda x: len(vars(x)['src']) <= MAX_LEN and
21.        len(vars(x)['trg']) <= MAX_LEN)
22.    MIN_FREQ = 2
23.    SRC.build_vocab(train.src, min_freq=MIN_FREQ)
24.    TGT.build_vocab(train.trg, min_freq=MIN_FREQ)
```

批训练对于速度来说很重要。我们希望批次分割非常均匀并且填充最少。要做到这一点，我们必须修改torchtext默认的批处理函数。这部分代码修补其默认批处理函数，以确保我们搜索足够多的句子以构建紧密批处理。

迭代器

```
1. class MyIterator(data.Iterator):
2.     def create_batches(self):
3.         if self.train:
4.             def pool(d, random_shuffler):
5.                 for p in data.batch(d, self.batch_size * 100):
6.                     p_batch = data.batch(
7.                         sorted(p, key=self.sort_key),
8.                         self.batch_size, self.batch_size_fn)
9.                     for b in random_shuffler(list(p_batch)):
10.                        yield b
11.             self.batches = pool(self.data(), self.random_shuffler)
12.         else:
13.             self.batches = []
14.             for b in data.batch(self.data(), self.batch_size,
15.                                 self.batch_size_fn):
16.                 self.batches.append(sorted(b, key=self.sort_key))
17.     def rebatch(pad_idx, batch):
18.         "Fix order in torchtext to match ours"
19.         src, trg = batch.src.transpose(0, 1), batch.trg.transpose(0, 1)
20.         return Batch(src, trg, pad_idx)
```

多GPU训练

最后为了真正地快速训练，我们将使用多个GPU。这部分代码实现了多GPU字生成。它不是Transformer特有的，所以我不会详细介绍。其思想是将训练时的单词生成分成块，以便在许多不同的GPU上并行处理。我们使用PyTorch并行原语来做这一点：

- 复制 - 将模块拆分到不同的GPU上
- 分散 - 将批次拆分到不同的GPU上

- 并行应用 - 在不同GPU上将模块应用于批处理
- 聚集 - 将分散的数据聚集到一个GPU上
- `nn.DataParallel` - 一个特殊的模块包装器，在评估之前调用它们。

```

1. # Skip if not interested in multigpu.
2. class MultiGPULossCompute:
3.     "A multi-gpu loss compute and train function."
4.     def __init__(self, generator, criterion, devices, opt=None, chunk_size=5):
5.         # Send out to different gpus.
6.         self.generator = generator
7.         self.criterion = nn.parallel.replicate(criterion,
8.                                                 devices=devices)
9.         self.opt = opt
10.        self.devices = devices
11.        self.chunk_size = chunk_size
12.    def __call__(self, out, targets, normalize):
13.        total = 0.0
14.        generator = nn.parallel.replicate(self.generator,
15.                                           devices=self.devices)
16.        out_scatter = nn.parallel.scatter(out,
17.                                          target_gpus=self.devices)
18.        out_grad = [[] for _ in out_scatter]
19.        targets = nn.parallel.scatter(targets,
20.                                       target_gpus=self.devices)
21.        # Divide generating into chunks.
22.        chunk_size = self.chunk_size
23.        for i in range(0, out_scatter[0].size(1), chunk_size):
24.            # Predict distributions
25.            out_column = [[Variable(o[:, i:i+chunk_size].data,
26.                                   requires_grad=self.opt is not None)]

```



```

27.         for o in out_scatter]
28.     gen = nn.parallel.parallel_apply(generator, out_column)
29.     # Compute loss.
30.     y = [(g.contiguous().view(-1, g.size(-1)),
31.          t[:, i:i+chunk_size].contiguous().view(-1))
32.          for g, t in zip(gen, targets)]
33.     loss = nn.parallel.parallel_apply(self.criterion, y)
34.     # Sum and normalize loss
35.     l = nn.parallel.gather(loss,
36.                            target_device=self.devices[0])
37.     l = l.sum()[0] / normalize
38.     total += l.data[0]
39.     # Backprop loss to output of transformer
40.     if self.opt is not None:
41.         l.backward()
42.         for j, l in enumerate(loss):
43.             out_grad[j].append(out_column[j][0].grad.data.clone())
44.     # Backprop all loss through transformer.
45.     if self.opt is not None:
46.         out_grad = [Variable(torch.cat(og, dim=1)) for og in out_grad]
47.         o1 = out
48.         o2 = nn.parallel.gather(out_grad,
49.                                target_device=self.devices[0])
50.         o1.backward(gradient=o2)
51.         self.opt.step()
52.         self.opt.optimizer.zero_grad()
53.     return total * normalize

```

现在我们创建模型，损失函数，优化器，数据迭代器和并行化。

```
1. # GPUs to use
2. devices = [0, 1, 2, 3]
3. if True:
4.     pad_idx = TGT.vocab.stoi["<blank>"]
5.     model = make_model(len(SRC.vocab), len(TGT.vocab), N=6)
6.     model.cuda()
7.     criterion = LabelSmoothing(size=len(TGT.vocab), padding_idx=pad_idx,
        smoothing=0.1)
8.     criterion.cuda()
9.     BATCH_SIZE = 12000
10.    train_iter = MyIterator(train, batch_size=BATCH_SIZE, device=0,
11.                            repeat=False, sort_key=lambda x: (len(x.src), len(x.trg)),
12.                            batch_size_fn=batch_size_fn, train=True)
13.    valid_iter = MyIterator(val, batch_size=BATCH_SIZE, device=0,
14.                            repeat=False, sort_key=lambda x: (len(x.src), len(x.trg)),
15.                            batch_size_fn=batch_size_fn, train=False)
16.    model_par = nn.DataParallel(model, device_ids=devices)
17. None
```

现在我们训练模型。我将稍微使用预热步骤，但其他一切都使用默认参数。在具有4个Tesla V100 GPU的AWS p3.8xlarge机器上，每秒运行约27,000个词，批训练大小为12,000。

训练系统

```
1. #!wget https://s3.amazonaws.com/opennmt-models/iwslt.pt
```

```

1. if False:
2.     model_opt = NoamOpt(model.src_embed[0].d_model, 1, 2000,
3.         torch.optim.Adam(model.parameters(), lr=0, betas=(0.9, 0.98), eps=1e-9))
4.     for epoch in range(10):
5.         model_par.train()
6.         run_epoch((rebatch(pad_idx, b) for b in train_iter),
7.             model_par,
8.             MultiGPULossCompute(model.generator, criterion,
9.                 devices=devices, opt=model_opt))
10.        model_par.eval()
11.        loss = run_epoch((rebatch(pad_idx, b) for b in valid_iter),
12.            model_par,
13.            MultiGPULossCompute(model.generator, criterion,
14.                devices=devices, opt=None))
15.        print(loss)
16. else:
17.     model = torch.load("iwslt.pt")

```

一旦训练完成，我们可以解码模型以产生一组翻译。在这里，我们只需翻译验证集中的第一个句子。此数据集非常小，因此使用贪婪搜索的翻译相当准确。

```

1. for i, batch in enumerate(valid_iter):
2.     src = batch.src.transpose(0, 1)[:1]
3.     src_mask = (src != SRC.vocab.stoi["<blank>"]).unsqueeze(-2)
4.     out = greedy_decode(model, src, src_mask,
5.                           max_len=60, start_symbol=TGT.vocab.stoi["<s>"])
6.     print("Translation:", end="\t")
7.     for i in range(1, out.size(1)):
8.         sym = TGT.vocab.itos[out[0, i]]
9.         if sym == "</s>": break
10.        print(sym, end=" ")
11.    print()
12.    print("Target:", end="\t")
13.    for i in range(1, batch.trg.size(0)):
14.        sym = TGT.vocab.itos[batch.trg.data[i, 0]]
15.        if sym == "</s>": break
16.        print(sym, end=" ")
17.    print()
18.    break

1. Translation:  <unk> <unk> . In my language , that means , thank you very much .
2. Gold:  <unk> <unk> . It means in my language , thank you very much .

```

附加组件：BPE，搜索，平均

所以这主要涵盖了Transformer模型本身。有四个方面我们没有明确涵盖。我们还实现了所有这些附加功能 OpenNMT-py[3].

1) 字节对编码/ 字片(Word-piece)：我们可以使用库来首先将数据预处理为子字单元。参见Rico Sennrich的subword-nmt实现[4]。这些模型将训练数据转换为如下所示：

`__Die __Protokoll datei __kann __heimlich __per __E - Mail __oder __FTP __an __einen
__bestimmte n __Empfänger __gesendet __werden .`

2) 共享嵌入：当使用具有共享词汇表的BPE时，我们可以在源/目标/生成器之间共享相同的权重向量，详见[5]。要将其添加到模型，只需执行以下操作：

1. if False:
2. `model.src_embed[0].lut.weight = model.tgt_embeddings[0].lut.weight`
3. `model.generator.lut.weight = model.tgt_embed[0].lut.weight`

3) 集束搜索：这里展开说有点太复杂了。PyTorch版本的实现可以参考 OpenNMT-py[6]。

4) 模型平均：这篇文章平均最后k个检查点以创建一个集合效果。如果我们有一堆模型，我们可以在事后这样做：

1. `def average(model, models):`
2. `"Average models into model"`
3. `for ps in zip(*[m.params() for m in [model] + models]):`
4. `p[0].copy_(torch.sum(*ps[1:]) / len(ps[1:]))`

结果

在WMT 2014英语-德语翻译任务中，大型Transformer模型（表2中的Transformer（大））优于先前报告的最佳模型（包括集成的模型）超过2.0 BLEU，建立了一个新的最先进BLEU得分为28.4。该模型的配置列于表3的底部。在8个P100 GPU的机器上，训练需要需要3.5天。甚至我们的基础模型也超过了之前发布的所有模型和集成，而且只占培训成本的一小部分。

在WMT 2014英语-法语翻译任务中，我们的大型模型获得了41.0的BLEU分数，优于以前发布的所有单一模型，不到以前最先进技术培训成本的1/4 模型。使用英语到法语训练的Transformer（大）模型使用dropout概率 $P_{drop} = 0.1$ ，而不是0.3。

1. `Image(filename="images/results.png")`

Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.75			
Deep-Att + PosUnk [39]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [39]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [9]	26.36	41.29	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1	$3.3 \cdot 10^{18}$	
Transformer (big)	28.4	41.8	$2.3 \cdot 10^{19}$	

我们在这里编写的代码是基本模型的一个版本。这里有系统完整训练的版本 (Example Models[7]).

通过上一节中的附加扩展，OpenNMT-py复制在EN-DE WMT上达到26.9。在这里，我已将这些参数加载到我们的重新实现中。

1. `!wget https://s3.amazonaws.com/opennmt-models/en-de-model.pt`
1. `model, SRC, TGT = torch.load("en-de-model.pt")`
1. `model.eval()`
2. `sent = "_The _log _file _can _be _sent _secret ly _with _email _or _FTP _to _a _specified _receiver".split()`
3. `src = torch.LongTensor([[SRC.stoi[w] for w in sent]])`
4. `src = Variable(src)`
5. `src_mask = (src != SRC.stoi["<blank>"]).unsqueeze(-2)`
6. `out = greedy_decode(model, src, src_mask,`
7. `max_len=60, start_symbol=TGT.stoi["<s>"])`
8. `print("Translation:", end="\t")`
9. `trans = "<s> "`
10. `for i in range(1, out.size(1)):`
11. `sym = TGT.itos[out[0, i]]`
12. `if sym == "</s>": break`
13. `trans += sym + " "`
14. `print(trans)`

1. Translation: `<s> _Die _Protokoll datei _kann _ heimlich _per _E - Mail _oder _FTP
_an _einen _bestimmte n _Empfänger _gesendet _werden .`

注意力可视化

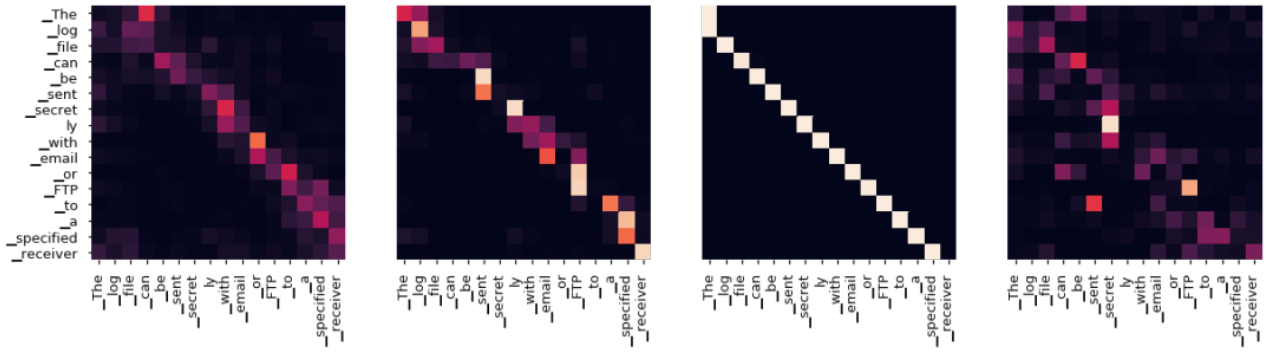
即使使用贪婪的解码器，翻译看起来也不错。我们可以进一步想象它，看看每一层注意力发生了什么。

```

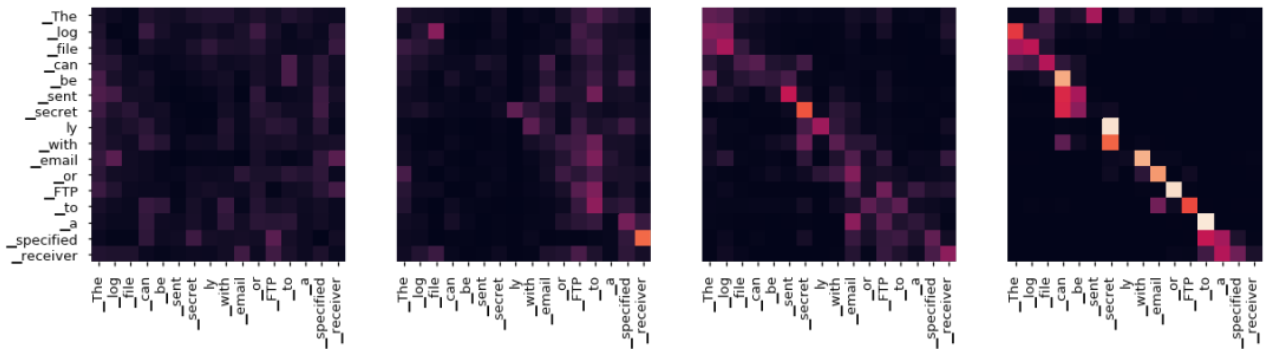
1. tgt_sent = trans.split()
2. def draw(data, x, y, ax):
3.     seaborn.heatmap(data,
4.         xticklabels=x, square=True, yticklabels=y, vmin=0.0, vmax=1.0,
5.         cbar=False, ax=ax)
6. for layer in range(1, 6, 2):
7.     fig, axs = plt.subplots(1,4, figsize=(20, 10))
8.     print("Encoder Layer", layer+1)
9.     for h in range(4):
10.        draw(model.encoder.layers[layer].self_attn.attn[0, h].data,
11.            sent, sent if h ==0 else [], ax=axs[h])
12.    plt.show()
13. for layer in range(1, 6, 2):
14.    fig, axs = plt.subplots(1,4, figsize=(20, 10))
15.    print("Decoder Self Layer", layer+1)
16.    for h in range(4):
17.        draw(model.decoder.layers[layer].self_attn.attn[0, h].data[:len(tgt_sent),
18.            :len(tgt_sent)],
19.            tgt_sent, tgt_sent if h ==0 else [], ax=axs[h])
20.    plt.show()
21.    print("Decoder Src Layer", layer+1)
22.    fig, axs = plt.subplots(1,4, figsize=(20, 10))
23.    for h in range(4):
24.        draw(model.decoder.layers[layer].self_attn.attn[0, h].data[:len(tgt_sent),
25.            :len(sent)],
26.            sent, tgt_sent if h ==0 else [], ax=axs[h])
27.    plt.show()

```

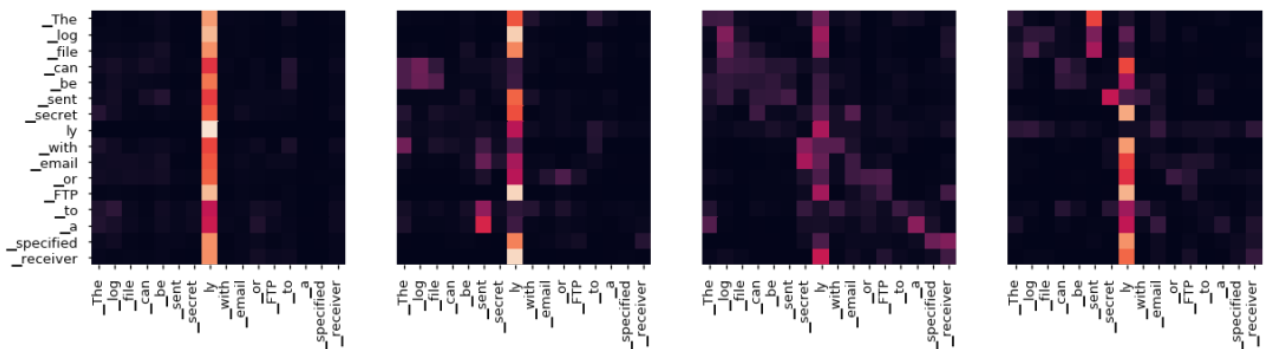
1. Encoder Layer 2



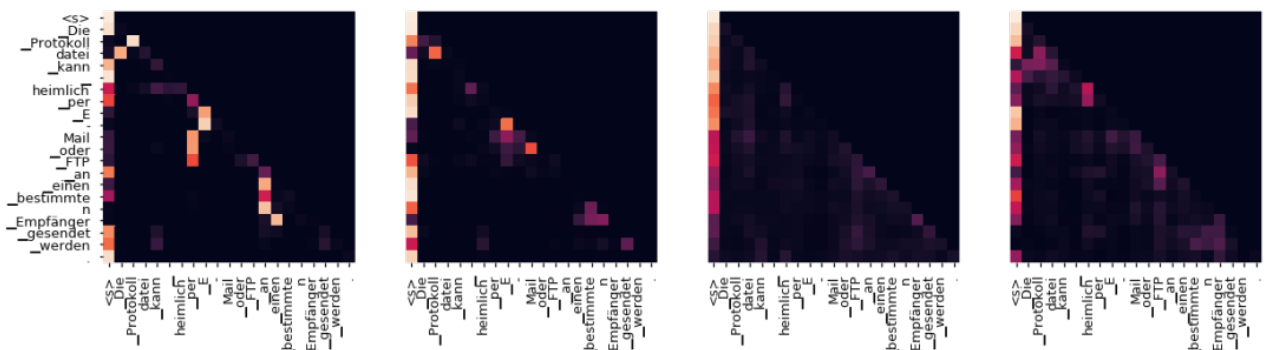
1. Encoder Layer 4



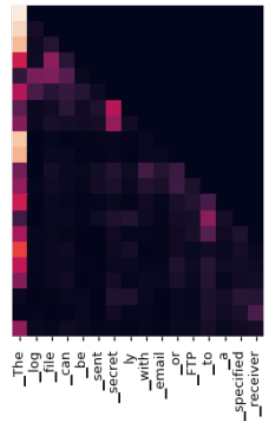
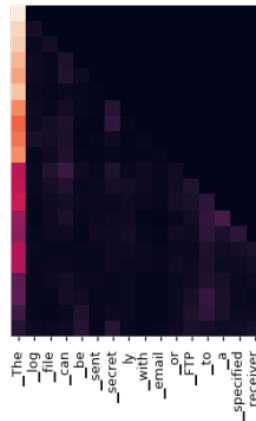
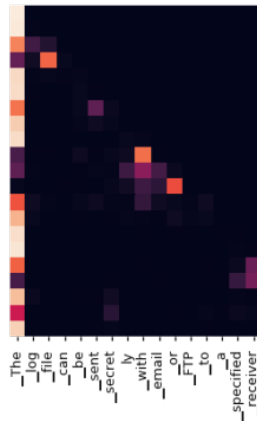
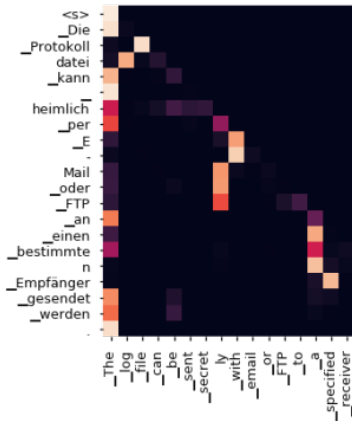
1. Encoder Layer 6



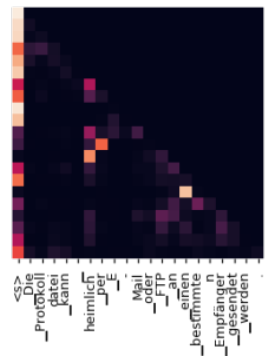
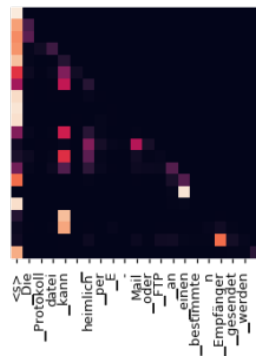
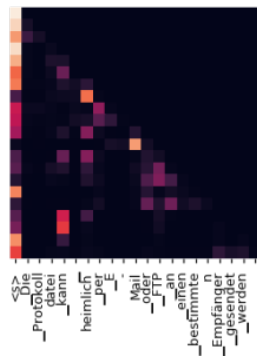
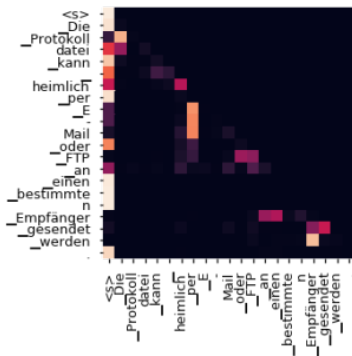
1. Decoder Self Layer 2



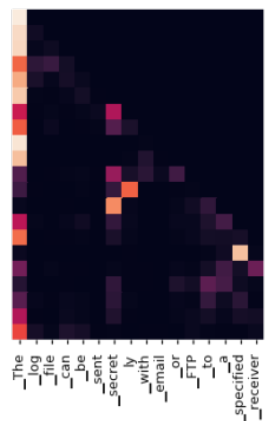
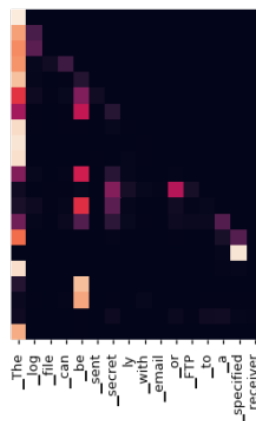
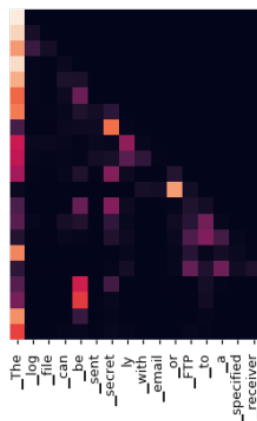
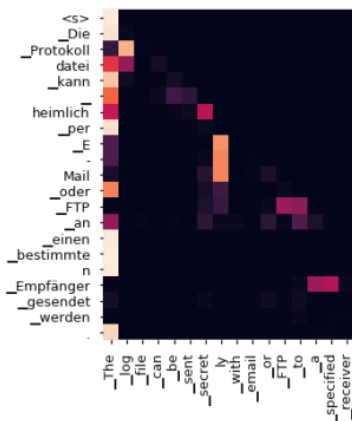
1. Decoder Src Layer 2



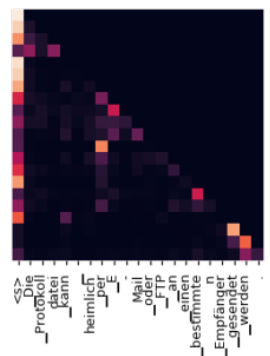
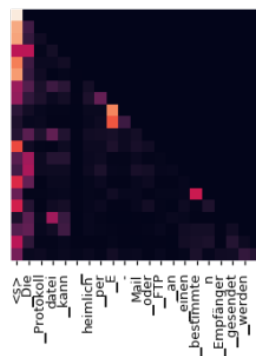
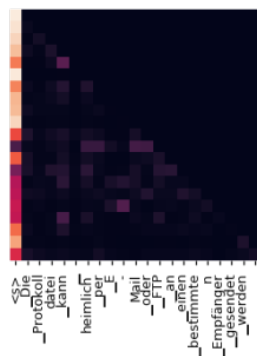
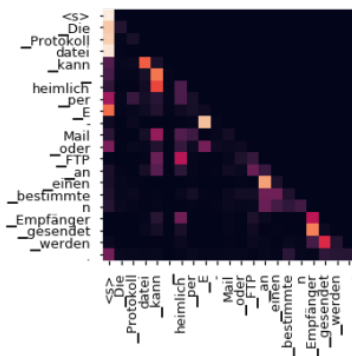
1. Decoder Self Layer 4



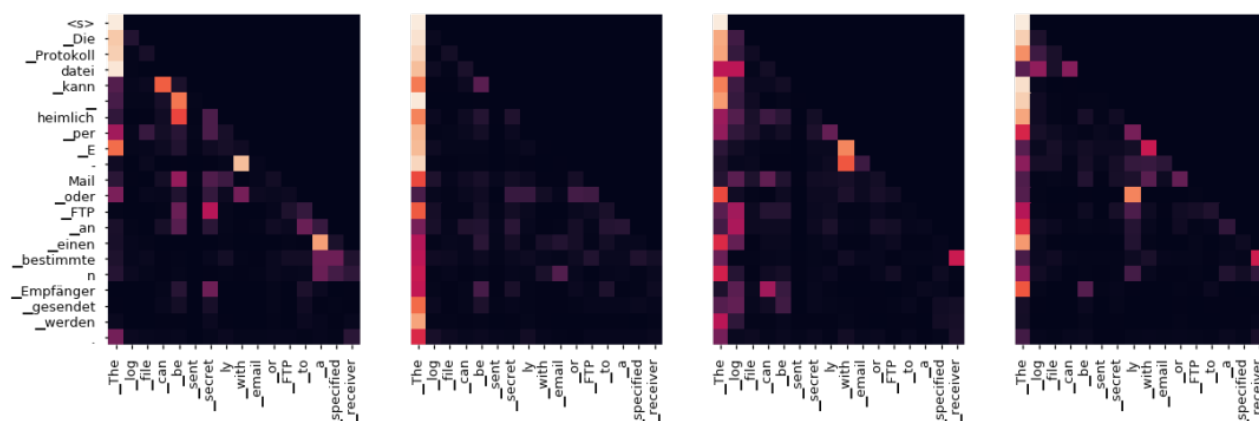
1. Decoder Src Layer 4



1. Decoder Self Layer 6



1. Decoder Src Layer 6



结论

希望这段代码对未来的研究很有用。如果您有任何问题，请与我们联系。如果您发现此代码有用，请查看我们的其他OpenNMT工具。

1. @inproceedings{opennmt,
2. author = {Guillaume Klein and
3. Yoon Kim and
4. Yuntian Deng and
5. Jean Senellart and
6. Alexander M. Rush},
7. title = {OpenNMT: Open-Source Toolkit for Neural Machine Translation},
8. booktitle = {Proc. ACL},
9. year = {2017},
10. url = {https://doi.org/10.18653/v1/P17-4012},
11. doi = {10.18653/v1/P17-4012}
12. }

Cheers , srush

参考链接

- [1] <https://arxiv.org/abs/1412.6980>
- [2] <https://arxiv.org/abs/1512.00567>

- [3] <https://github.com/opennmt/opennmt-py>
- [4] <https://github.com/rsennrich/subword-nmt>
- [5] <https://arxiv.org/abs/1608.05859>
- [6] <https://github.com/OpenNMT/OpenNMT-py/blob/master/onmt/translate/Beam.py>
- [7] <http://opennmt.net/Models-py/>

END