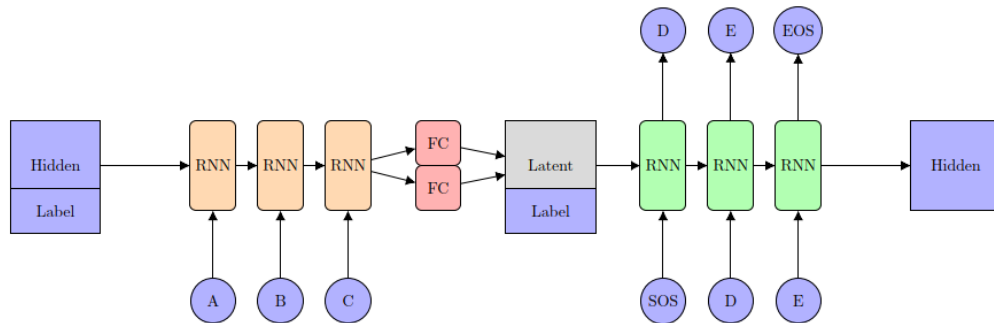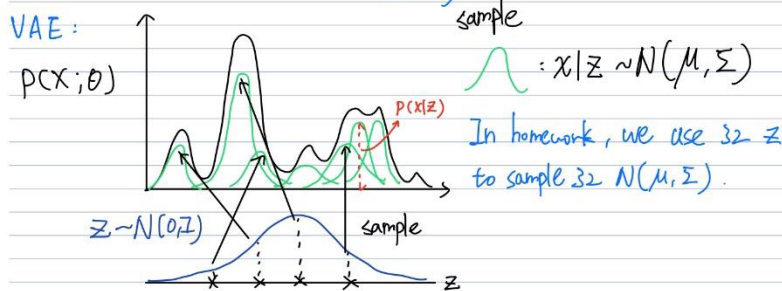# Lab 5 Conditional Sequence to Sequense VAE

1. Introduction

   Using CVAE to change the tense of the verbs , ex. abandon -> abandoned , input "abandon" to encoder with its tense(sp) , and use output and the target tense(p) as decoder input , and expect the verb "abandoned" as output .



2. Derivation of CVAE

$$(q(z|x;\theta) = q(z))$$

decoder output

$$\Rightarrow \log P(x;\theta) \ge L(x,q,\theta) = \int_z q(z) \log\left(\frac{P(x|z;\theta)\,P(z)}{q(z)}\right) dz$$

decoder    encoder

★ Find $P(x|z)$ and $q(z)$ to maximum $\log P(x;\theta)$

Object function: $L(x,q,\theta)$:

$$L(x,q,\theta) = \int_z q(z) \log\left(\frac{P(x|z;\theta)\,P(z)}{q(z)}\right) dz$$

$$= \int_z q(z)\, \log(P(x|z;\theta))\, dz + \int_z q(z) \log\left(\frac{P(z)}{q(z)}\right) dz$$

$$= \int_z q(z|x;\theta) \log(P(x|z;\theta))\, dz - KL\big(q(z|x;\theta) \,\|\, P(z)\big)$$

$$= E_{z \sim q(z|x;\theta)} \log P(x|z;\theta) - KL\big(q(z|x;\theta) \,\|\, P(z)\big)$$

maximum      minimum

Change to CVAE, the object function is:

$$E_{z \sim q(z|x,c;\theta)} \log P(x|z,c;\theta) - KL\big(q(z|x,c;\theta) \,\|\, P(z|c)\big)$$

3. Implementation details

   Dataloader (dataloader.py):



```python
def Dataloader(path,mode='train'):
    file = open(path, "rt")
    line = file.readlines()

    tense2num = {'sp': 0, 'tp': 1, 'pg': 2, 'p': 3}
    condition_list = list(tense2num.values())

    alphabet2num = {'SOS': 0, 'EOS': 1}
    alphabet2num.update([(chr(i+97),i+2) for i in range(0,26)])

    input_batch = []
    cond_batch = []

    if mode == 'test':
        for i in range(4):
            cond_batch.append(condition_list[i])

    for s in line:
        # read word each line
        word = s.split('\n')[0].split(' ')
        if mode == 'train':
            tmp = [word[0], word[1], word[2], word[3]]
        else:
            tmp = [word[0]]
```

```python
        # encode word to numbers
        for n in range(len(tmp)):
            encode = []
            for id in tmp[n]:
                encode.append(alphabet2num[id])

            encode.append(1)
            input_batch.append(encode)

            if mode == 'train':
                cond_batch.append(condition_list[n])

    return input_batch, cond_batch
```

I encode the "SOS" , "EOS" and 26 alphabet to 0~27 (input size is 28 ) and tense encoded "sp" , "tp" , "pg" , "p" to 0~3 . My dataloader will return two object , "input_batch" , which contains 4908 verbs for training or 10 for testing ; "cond_batch" , 4908 tense for training or 10 for testing .

CVAE (model.py):

**Encoder** :

```
16    class CVAE(nn.Module):
17        class EncoderRNN(nn.Module):
18            def __init__(self, input_size, input_cond_size, hidden_size):
19                super(CVAE.EncoderRNN, self).__init__()
20                self.hidden_size = hidden_size
21                self.embedding_N = nn.Embedding(input_size, hidden_size)
22                self.rnn_N = nn.LSTM(hidden_size, hidden_size)
23                for weight in self.rnn_N.parameters():
24                    if len(weight.size()) > 1:
25                        init.orthogonal_(weight.data)
26
27            def initial(self):
28                h = torch.zeros(1,1,self.hidden_size-8, device=device)
29                c = torch.zeros(1,1,self.hidden_size, device=device)
30
31                return h ,c
32
33            def forward(self, input, hidden, cell):
34                embedded = self.embedding_N(input).view(1, 1, -1)
35                embedded = embedded.permute(1,0,2)
36
37                output, (h, c) = self.rnn_N(embedded, (hidden, cell))
38
39                return output, h, c
```

Includes one embedding layer and one LSTM layer , and I found that initial an orthogonal weight will get the better result .

**Decoder** :

```
41        class DecoderRNN(nn.Module):
42            def __init__(self, input_size, input_cond_size, hidden_size):
43                super(CVAE.DecoderRNN, self).__init__()
44                self.hidden_size = hidden_size
45                self.embedding_D = nn.Embedding(input_size, hidden_size)
46                self.rnn_D = nn.LSTM(hidden_size, hidden_size)
47                for weight in self.rnn_D.parameters():
48                    if len(weight.size()) > 1:
49                        init.orthogonal_(weight.data)
50
51            def initial(self):
52                h = torch.zeros(1,1,self.hidden_size, device=device)
53                c = torch.zeros(1,1,self.hidden_size, device=device)
54
55                return h, c
56
57            def forward(self,input, hidden, cell):
58                embedded = self.embedding_D(input).view(1, 1, -1)
59                embedded = embedded.permute(1,0,2)
60
61                output, (h, c) = self.rnn_D(embedded, (hidden, cell))
62
63                return output, (h, c)
```

Includes one embedding layer and one LSTM layer , and I found that initial an orthogonal weight will get the better result .

**Other layers** :

```
65        def __init__(self, input_size, input_cond_size, hidden_size):
66            super(CVAE,self).__init__()
67            self.input_size = input_size
68            self.encoder = self.EncoderRNN(input_size, input_cond_size, hidden_size).cuda()
69            self.decoder = self.DecoderRNN(input_size, input_cond_size, hidden_size).cuda()
70            self.embedding_cond = nn.Embedding(input_cond_size, 8)
71            self.fc_meam_h = nn.Linear(hidden_size, 32)
72            self.fc_logvar_h = nn.Linear(hidden_size, 32)
73            self.fc_st_D = nn.Linear(40, hidden_size)
74            self.fc_out = nn.Linear(hidden_size, input_size)
75
```

**self.embedding_cond** : Embeds the conditional to size 8 .

**self.mean_h** : A fully connection layer to get size 32 mean from final hidden unit of encoder .

**self.logvar_h** : A fully connection layer to get size 32 logvar from final hidden unit of encoder .

**self.fc_st_D** : A fully connection layer to change the size of hidden unit put into the decoder at initial state from 32+8 to 256 .

**self.fc_out** : A fully connection layer to change the size of the output of decoder from 256 to 28(input size) .

**Reparameterize** :

```
77      def reparameterize(self, mean, var):
78          std = torch.exp(0.5*var)
79          eps = torch.randn_like(std)
80          z = eps.mul(std).add_(mean)
81
82          return z
```

**Forward** :

```
85  def forward(self, input_word, cond, teach_ratio, KLD_ratio, criterion, optimizer):
86      loss = 0
87      optimizer.zero_grad()
88
89      # concatenate conditional to hidden unit 0
90      Cond = self.embedding_cond(cond)
91      hidden_n, cell_n = self.encoder.initial()
92      hidden_n = torch.cat((hidden_n, Cond.view(1, 1, -1)),2)
93
94      # Encoder
95      for alphabet in input_word:
96          _, hidden_n, cell_n = self.encoder(alphabet, hidden_n, cell_n)
97
98      # get mean and variance
99      h_t = torch.squeeze(hidden_n)
100     mean_h = self.fc_mean_h(h_t)
101     logvar_h = self.fc_logvar_h(h_t)
102
103     KLD_loss = -0.5 * torch.sum(1 + logvar_h - mean_h.pow(2) - logvar_h.exp())
104
105     # sample
106     latent_h = self.reparameterize(mean_h, logvar_h)
107
108     pred = None
109     out_fc = None
110
111     _, cell_d = self.decoder.initial()
112     decoder_input = torch.tensor([[[0]]]).cuda()
```

```
113     # concatenate conditional to latent vector
114     hidden_d = latent_h.view(1, -1)
115     hidden_d = torch.cat((hidden_d, Cond),1)
116     hidden_d = self.fc_st_D(hidden_d).view(1, 1, -1)
117
118     # Decoder
119     for idx in range(input_word.shape[0]):
120         use_teacher_forcing = True if random.random() < teach_ratio else False
121         if (use_teacher_forcing) and (idx !=0):
122             output_d, (hidden_d, cell_d) = self.decoder(input_word[idx - 1], hidden_d, cell_d)
123         else:
124             output_d, (hidden_d, cell_d) = self.decoder(decoder_input, hidden_d, cell_d)
125
126
127         out_fc = self.fc_out(output_d) # classification
128         out_alphabet = torch.argmax(out_fc).item() # return to alphabet code
129         decoder_input = torch.tensor([[[out_alphabet]]]).cuda()
130
131         if idx == 0:
132             pred = out_fc.view(1,self.input_size)
133         else:
134             pred = torch.cat((pred, out_fc.view(1,self.input_size)),0)
135
136
137     loss = criterion(pred, input_word.long())
138     loss += KLD_loss*KLD_ratio
139     loss.backward()
140     optimizer.step()
141
142     return loss.item(), KLD_loss.item()
```

**Evaluation** (for bleu testing ):

```
144  def Eval(self, input_word, cond1, cond2):
145      with torch.no_grad():
146          hidden_n, cell_n = self.encoder.initial()
147          Cond1 = self.embedding_cond(cond1).view(1, 1, -1)
148          hidden_n = torch.cat((hidden_n, Cond1),2)
149
150          for alphabet in input_word:
151              _, hidden_n, cell_n = self.encoder(alphabet, hidden_n, cell_n)
152
153          h_t = torch.squeeze(hidden_n)
154          mean_h = self.fc_mean_h(h_t)
155          logvar_h = self.fc_logvar_h(h_t)
156
157          latent_h = self.reparameterize(mean_h, logvar_h)
158
159          pred_word = []
160          out_fc = None
```

```
162          _, cell_d = self.decoder.initial()
163          decoder_input = torch.tensor([[[0]]]).cuda()
164          Cond2 = self.embedding_cond(cond2).view(1, 1, -1)
165          hidden_d = latent_h.view(1, 1, -1)
166          hidden_d = torch.cat((hidden_d, Cond2),2)
167          hidden_d = self.fc_st_D(hidden_d).view(1, 1, -1)
168
169          for i in range(input_word.shape[0]):
170              output_d, (hidden_d, cell_d) = self.decoder(decoder_input, hidden_d, cell_d)
171              out_fc = self.fc_out(output_d)
172              out_alphabet = torch.argmax(out_fc).item()
173              if out_alphabet == 1:
174                  break
175              decoder_input = torch.tensor([[[out_alphabet]]]).cuda()
176              pred_word.append(out_alphabet)
177
178          return pred_word
```

**Gaussian generation** (In line 185 : torch.randn_like(torch.zeros(1, 1, 32)) is noise):

```
180  def gaussion(self, Cond):
181      Result = []
182      with torch.no_grad():
183          for i in range(100):
184              print(i,end='\r')
185              latent = torch.randn_like(torch.zeros(1, 1, 32)).cuda()
186              pred_batch = []
187
188              for cond in Cond:
189                  word = []
190                  _, cell_d = self.decoder.initial()
191                  out_fc = None
192                  decoder_input = torch.tensor([[[0]]]).cuda()
193                  Cond_ = self.embedding_cond(cond.long()).view(1, 1, -1)
194
195                  hidden_d = latent.view(1, 1, -1)
196                  hidden_d = torch.cat((hidden_d, Cond_),2)
197                  hidden_d = self.fc_st_D(hidden_d).view(1, 1, -1)
```

```
199                  for i in range(100):
200                      output_d, (hidden_d, cell_d) = self.decoder(decoder_input, hidden_d, cell_d)
201                      out_fc = self.fc_out(output_d)
202                      out_alphabet = torch.argmax(out_fc).item()
203                      if out_alphabet == 1:
204                          break
205                      decoder_input = torch.tensor([[[out_alphabet]]]).cuda()
206                      word.append(out_alphabet)
207                  pred_batch.append(word)
208
209              Result.append(pred_batch)
210
211          return Result
```

**Test** (test.py) : Test tense switching and gaussian score .

**blue.py** : compute bleu score .

**main.py ( for training and testing ):**

```python
parser = argparse.ArgumentParser(description='Set up')
parser.add_argument('--lr', type=float, default = 0.005)
parser.add_argument('--epochs', type=int, default = 500)
parser.add_argument('--hidden_size', type=int, default = 256)
parser.add_argument('--feq', type=float, default = 2.0)
parser.add_argument('--Mode', type=str, default = 'train')
args = parser.parse_args()
print(args)

if args.Mode == 'train':
    wandb.init(project='CVAE')
    wandb.save('/home/austin/nctu_hw/DL/DL_hw5/model.py')
    config = wandb.config
    config.hidden_size = args.hidden_size
    config.epochs = args.epochs
    config.learning_rate = args.lr

    model = CVAE(28, 4, args.hidden_size)
    model = model.cuda()
    wandb.watch(model)
    input_batch, cond_batch = Dataloader('/home/austin/nctu_hw/DL/DL_hw5/lab5_dataset/train.txt','train')
    Batch = []

    for i in range(len(input_batch)):
        input_batch[i] = torch.LongTensor(input_batch[i]).cuda()
        cond_batch[i]= torch.LongTensor([cond_batch[i]]).cuda()
        Batch.append([input_batch[i], cond_batch[i]]) # size : 4908

    criterion = torch.nn.CrossEntropyLoss().cuda()
    optimizer = optim.SGD(model.parameters(), lr=args.lr, momentum=0.9)
    BLeU = Bleu()
    test = Test()
```

```python
54    def gen_teach_ratio(Epoch):
55        return 1.0 - (Epoch/args.epochs)
56    def sigmoid(x):
57        return 1.0 / (1.0 + np.exp(-x))
58    def gen_KLD_ratio(mode, Epoch):
59        if mode == 'mon':
60            return (Epoch/args.epochs)*0.25
61        else:
62            period = args.epochs//args.feq
63            Epoch %= period
64            ratio = sigmoid((Epoch - period // 2.0) / (period // 10)) / 2.0
65            return ratio*0.5
66
```

```python
68        for epoch in range(args.epochs):
69            print('Epoch : ', epoch+1 ,'...',end='\r')
70            random.shuffle(Batch)
71            Loss = []
72            KLD_avg = []
73
74            teach_ratio = gen_teach_ratio(epoch)
75            KL_ratio = gen_KLD_ratio('cyc', epoch)
76
77            for sample in Batch:
78                loss, KLD = model(sample[0], sample[1], teach_ratio, KL_ratio, criterion, optimizer)
79                Loss.append(loss)
80                KLD_avg.append(KLD)
81
82            CE_loss = sum(Loss)/len(Loss) - (sum(KLD_avg)/len(KLD_avg))*KL_ratio
83            model.eval()
84            test.bleU_test(model)
85            model.train()
86            score = BLeU.get_score()
87
88            wandb.log({"teach_ratio": teach_ratio})
89            wandb.log({"BLeU": score})
90            wandb.log({"KL_ratio": KL_ratio})
91            wandb.log({"KLD": sum(KLD_avg)/len(KLD_avg)})
92            wandb.log({"loss": sum(Loss)/len(Loss)})
93            wandb.log({"CE_loss": CE_loss})
94            file = open('/home/austin/nctu_hw/DL/DL_hw5/Record.txt', "a+")
95            file.write(str(teach_ratio)+'/'+str(KL_ratio)+'/'+str(sum(Loss)/len(Loss))+'/'+str(CE_loss)+'/'+str(sum(KLD_avg)/len(KLD_avg))+'/'+str(score)+'\n')
96
97            print('Epoch : ',epoch+1,' Loss : ',sum(Loss)/len(Loss),' CE_Loss :', CE_loss, 'KLD_loss : ',sum(KLD_avg)/len(KLD_avg), ' Bleu : ', score)
98            if score >= 0.7:
99                name = '/home/austin/nctu_hw/DL/DL_hw5/weight/CVAE_'+str(epoch+1)+'_'+str(sum(Loss)/len(Loss))+'_'+str(score)+'.pth'
100               torch.save(model.state_dict(), name)
101               wandb.save(name)
```

```python
102    else:
103        model = CVAE(28, 4, args.hidden_size)
104        model.load_state_dict(torch.load('/home/austin/nctu_hw/DL/DL_hw5/weight/CVAE_45_0.38194182919856_0.8323583241361134.pth'))
105        model = model.cuda()
106        model.eval()
107        score = 0
108        BLeU = Bleu()
109        test = Test()
110        test.gaussian_test(model)
111        for i in range(1000):
112            print(i,end='\r')
113            test.bleU_test(model)
114            score = BLeU.get_score()
115            if score >= 0.8:
116                print('Bleu score : ', score)
117                break
```
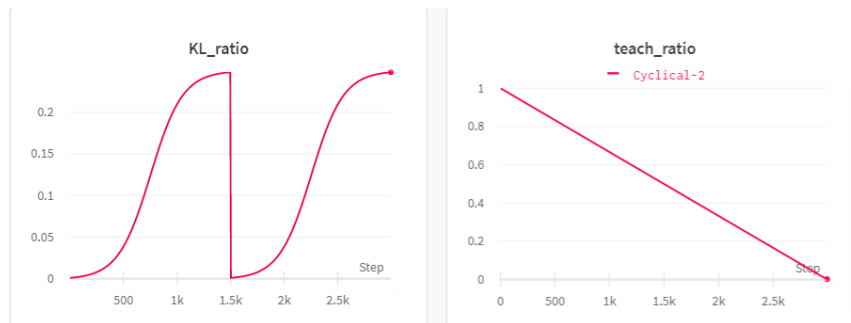
Hyperparameters :

Learning rate : 0.005 , Epochs : 500

I shuffle the training data with each epoch .

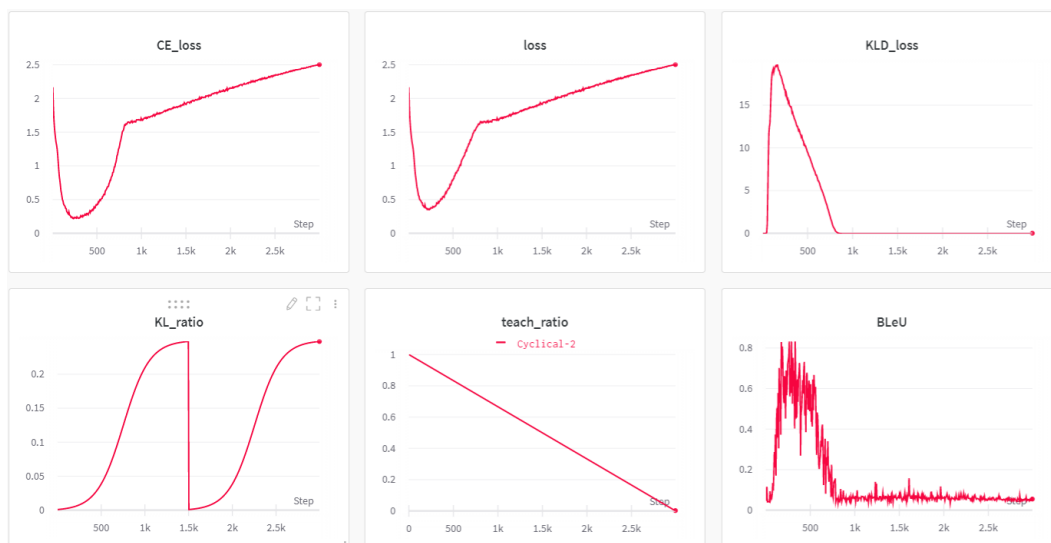KL weight : sigmoid with two cycles (max : 0.25)

Teach : from 1 down to 0 .



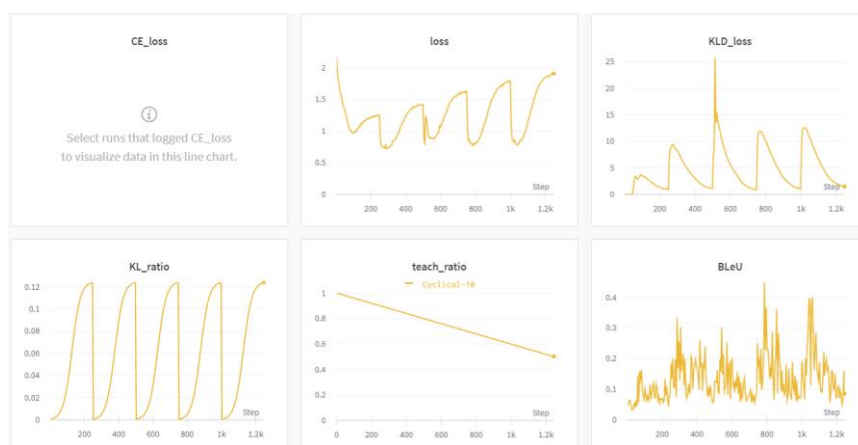4. Result and discussion

( I use the model weight saved at 45 epochs . ):





Bleu / Gaussian example :

I found that the bleu score is positively related to KLD loss , the two curves are very similar . In the early stages of training , KL weight was very low , and reconstruction loss was keep going down . When the KL weight up to 0.002 , reconstruction loss started to go up , and made the KLD loss and bleu down . I assume that when the KLD loss low enough , the reconstruction loss will keep down , so I set the max value of kl weight only 0.25 to ensure that the reconstruction loss drops enough to get high bleu score . For learning rate , because I update the model for every words in one epoch , It is better to use lower rate , and I find that 0.005 can make it almost drop vertically . For teach ratio , the higher ratio at the early stages can make loss drop steadily , so simply using monotonic mode . I found the issue that in the second cycle of kl weight , when it went down again , it doesn't make reconstruction loss down , and KLD loss doesn't go up again too . So I redo the experiment , this time I set 10 cycles of kl weight :



Although my computer was shutdown when epoch closed to 250 , only half of data , but it still can prove that my assumption is correct : When the KLD loss low enough , the reconstruction loss will keep down . And take a look at higher cycles , it makes the bleu score drops before it has time to rise enough , but the highest score in each cycle has an upward trend .

The curve of using monotonic KL weight (from 0 down to 0.25) :