

DLP Lab4: Conditional Sequence-to-Sequence VAE

Zhi-Yi Chin

joycenerd.cs09@nycu.edu.tw

August 16, 2021

1 Introduction

In this lab, we are implementing conditional seq2seq VAE (Figure 1) for English tense conversion and generation. The following are the detailed tasks we need to do:

- Implementing encoder and decoder of CVAE using LSTM.
- Using reparameterization trick to generate the latent vector z .
- Adopt teacher forcing ratio and KL loss throughout the training process.
- Apply to KL annealing method (monotonic and cyclical).
- Visualize the results.

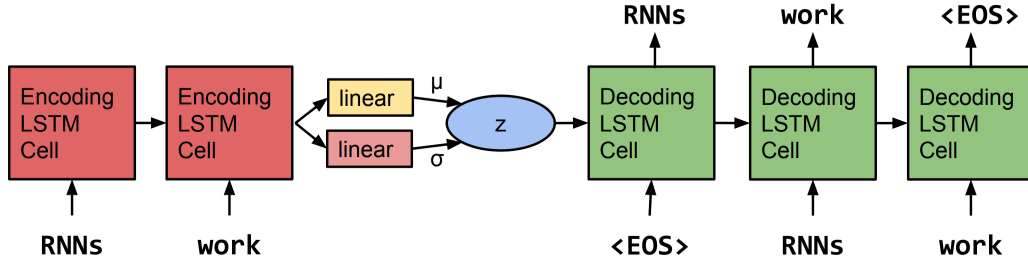


Figure 1: Conditional seq2seq VAE.

2 Derivation of CVAE

To determine θ , we would hope to maximize the marginal distribution $p(x; \theta)$ shown in Eq. 1

$$p(x; \theta) = \int p(x|z; \theta) p(z) dz \quad (1)$$

However, this becomes difficult as the integration over z is generally intractable when $p(x|z; \theta)$ is modeled by a neural network. To circumvent this difficulty, we recall Eq.3.

$$\log p(X; \theta) = \mathcal{L}(X, q, \theta) + \text{KL}(q(z) || p(Z|X, \theta)) \quad (2)$$

where

$$\begin{aligned}\mathcal{L}(X, q, \theta) &= \int q(Z) \log p(X, Z; \theta) dZ - \int q(Z) \log q(Z) dZ \\ \text{KL}(q(Z) || p(Z|X; \theta)) &= \int q(Z) \log \frac{q(Z)}{p(Z|X; \theta)} dZ\end{aligned}\tag{3}$$

As the equality holds for any choice of $q(Z)$, we introduce a distribution $q(Z|X; \theta')$ modeled by another neural network with parameter θ' . We can replace every $q(Z)$ to $q(Z|X)$ in Eq.3 and organize into Eq.4.

$$\begin{aligned}\mathcal{L}(X, q, \theta) &= \int q(Z|X) \log p(X, Z) dZ - \int q(Z|X) \log q(Z|X) dZ \\ &= \int q(Z|X) \log p(X|Z)p(Z) dZ - \int q(Z|X) \log q(Z|X) dZ \\ &= \mathbb{E}_{Z \sim q(Z|X; \theta')} \log p(X|Z; \theta) + \mathbb{E}_{Z \sim q(Z|X; \theta')} \log p(Z) - \mathbb{E}_{Z \sim q(Z|X; \theta')} \log q(Z|X; \theta') \\ &= \mathbb{E}_{Z \sim q(Z|X; \theta')} \log p(X|Z; \theta) - \text{KL}(q(Z|X; \theta') || p(Z))\end{aligned}\tag{4}$$

From Eq.4, we can easily see that $\mathbb{E}_{Z \sim q(Z|X; \theta')} \log p(X|Z; \theta)$ is the reconstruction term and $\text{KL}(q(Z|X; \theta') || p(Z))$ is the regularization term.

3 Derivation of KL Divergence Loss

First of all, why do we need KL loss? If the loss function only cares about output accuracy, it cannot prove the latent vector z is Gaussian normal distribution $\mathcal{N}(0, 1)$. So we derive KL loss as Eq.5.

$$\begin{aligned}\text{KL}(\mathcal{N}(\mu|\sigma^2) || \mathcal{N}(0, 1)) &= \int \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{(x-\mu)^2}{2\sigma^2}) \left(\log \frac{\exp(-\frac{(x-\mu)^2}{2\sigma^2})/\sqrt{2\pi\sigma^2}}{\exp(-\frac{x^2}{2})/\sqrt{2\pi}} \right) dx \\ &= \int \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{(x-\mu)^2}{2\sigma^2}) \log \left(\frac{1}{\sqrt{\sigma^2}} \exp(\frac{1}{2}(x^2 - \frac{(x-\mu)^2}{\sigma^2})) \right) dx \\ &= \int \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{(x-\mu)^2}{2\sigma^2}) \frac{1}{2} \left(-\log \sigma^2 + x^2 - \frac{(x-\mu)^2}{\sigma^2} \right) dx \\ &= \frac{1}{2} \left(\mathbb{E}_{x \sim \exp(-\frac{(x-\mu)^2}{\sigma^2})} (-\log \sigma^2) + \mathbb{E}_{x \sim \exp(-\frac{(x-\mu)^2}{\sigma^2})} (x^2) - \mathbb{E}_{x \sim \exp(-\frac{(x-\mu)^2}{\sigma^2})} \left(\frac{(x-\mu)^2}{\sigma^2} \right) \right) \\ &= \frac{1}{2} \left(-\log \sigma^2 + (\mu^2 + \sigma^2) - \frac{1}{\sigma^2} \int \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{(x-\mu)^2}{\sigma^2}) (x-\mu)^2 dx \right) \\ &= \frac{1}{2} (-\log \sigma^2 + \mu^2 + \sigma^2 - 1)\end{aligned}\tag{5}$$

4 Implementation Details

4.1 Data pre-processing

For **training** data, we flatten the data to 1 vector since for training data, there are four tenses in one line. The condition of the training data is the index of the flatten data modulo 4. For **testing**

data, we have the input word, and target word and the tenses conversion have been mentioned in `readme.txt`, we code it into integer numbers.

```
1 class TextDataset(Dataset):
2     def __init__(self, root_dir, mode):
3         self.root_dir = root_dir
4         self.mode = mode
5
6         file = Path(self.root_dir).joinpath(f"{self.mode}.txt")
7         self.data = np.loadtxt(file, dtype=str)
8
9         if self.mode == 'train':
10             self.data = self.data.reshape(-1)
11         else:
12             self.tense = [[0, 3], [0, 2], [0, 1], [0, 1], [3, 1], [0, 2], [3, 0], [2, 0], [2, 3],
13                           [2, 1]]
14             self.tense = np.asarray(self.tense, dtype=int)
15
16     def __len__(self):
17         return self.data.shape[0]
18
19     def __getitem__(self, index):
20         if self.mode == 'train':
21             data = self.data[index]
22             c = index % 4
23             return data, c
24         else:
25             data = self.data[index, 0]
26             target = self.data[index, 1]
27             c1 = self.tense[index, 0]
28             c2 = self.tense[index, 1]
29             return data, target, c1, c2
```

Listing 1: Dataset class.

Deep learning models only accept numbers, so we transform all the words into one hot encoder according to the alphabet. Since we only have lower case alphabets, so a...z mapped to 2...27. Moreover, we also add SOS token at the beginning of each word which is of value 0, and EOS token at the end of each word which is of value 1.

```

1 class OneHotEncoder():
2     def __init__(self):
3         self.sos=0
4         self.eos=1
5
6     def tokenize(self,word):
7         chars=['SOS']+list(word)+['EOS']
8         token=[]
9         for ch in chars:
10             if ch=='SOS':
11                 token.append(self.sos)
12             elif ch=='EOS':
13                 token.append(self.eos)
14             else:
15                 token.append(ord(ch)-ord('a')+2)
16         token=torch.from_numpy(np.asarray(token))
17         return token
18
19     def inv_tokenize(self,token):
20         word=''
21         for val in token:
22             if val==1:
23                 break
24             word+=chr(val-2+ord('a'))
25         return word

```

Listing 2: Onehot encoder.

4.2 Encoder

We initialize the hidden state and cell state to 0. Then we embed condition c and concatenate c with hidden state and cell state, respectively. After that, we perform a fully connected layer to make the hidden state and cell state size to hidden layer size. Then we embed the input sequence x and passed the input sequence, hidden state, and cell state into LSTM and output the new hidden state and the new cell state. Then we perform a full connection to get mean and log variance for both hidden state and cell state. We get the latent z by $\exp(\logvar/2) * \epsilon + \mu$, ϵ is sample from $\mathcal{N}(0,1)$, this is known as the **reparameterization trick**.

```

1 class EncoderRNN(nn.Module):
2     def __init__(self, input_size, hidden_size,c_size,c_hidden_size,z_size,
3         device):
4         super(EncoderRNN, self).__init__()
5         self.hidden_size = hidden_size
6         self.input_size=input_size
7         self.c_size=c_size
8         self.c_hidden_size=c_hidden_size
9         self.z_size=z_size

```

```

10     self.device=device
11
12     self.embedding = nn.Embedding(input_size, hidden_size)
13     self.c_embedding=nn.Embedding(c_size,c_hidden_size)
14     self.lstm=nn.LSTM(self.hidden_size,self.hidden_size)
15     self.hidden_fc=nn.Linear(self.hidden_size+self.c_hidden_size,
16                               self.hidden_size)
17     self.cell_fc=nn.Linear(self.hidden_size+self.c_hidden_size,
18                             self.hidden_size)
19
20     self.hidden_mean=nn.Linear(hidden_size,z_size)
21     self.hidden_log_var=nn.Linear(hidden_size,z_size)
22     self.cell_mean=nn.Linear(hidden_size,z_size)
23     self.cell_log_var=nn.Linear(hidden_size,z_size)
24
25     def forward(self,x,hidden_state,cell_state,c):
26         # hidden state
27         c=self.c_embedding(c).reshape(1,1,-1)
28         hidden_state=torch.cat((hidden_state,c),dim=2)
29         hidden_state=self.hidden_fc(hidden_state)
30
31         # cell state
32         cell_state=torch.cat((cell_state,c),dim=2)
33         cell_state=self.cell_fc(cell_state)
34
35         x=self.embedding(x).reshape(-1,1,self.hidden_size)
36         out,(hidden_state,cell_state)=self.lstm(x,(hidden_state,cell_state))
37
38         # mean and variance
39         hidden_mu=self.hidden_mean(hidden_state)
40         hidden_log_sigma=self.hidden_log_var(hidden_state)
41         hidden_eps=torch.normal(torch.zeros(self.z_size),
42                                  torch.ones(self.z_size)).to(self.device)
43         # reparameterization trick
44         hidden_z=torch.exp(hidden_log_sigma/2)*hidden_eps+hidden_mu
45
46         cell_mu=self.cell_mean(cell_state)
47         cell_log_sigma=self.cell_log_var(cell_state)
48         cell_eps=torch.normal(torch.zeros(self.z_size),
49                                torch.ones(self.z_size)).to(self.device)
50         # reparameterization trick
51         cell_z=torch.exp(cell_log_sigma/2)*self.std_normal+cell_mu
52
53         return hidden_mu,hidden_log_sigma,hidden_z,cell_mu,cell_log_sigma,cell_z
54
55     def init_hidden_and_cell(self):
56         hidden_state=torch.zeros(1,1,self.hidden_size,device=self.device)
57         cell_state=torch.zeros(1,1,self.hidden_size,device=self.device)

```

```
return hidden_state, cell_state
```

Listing 3: Encoder of CVAE

4.3 Decoder

To initialize the decoder, we first embed the condition. Then we concatenate hidden state latent code z and c and perform a full connection to make it of size hidden layer size, which is the decoder’s initial hidden state. For the initial cell state, we do the same thing. For the tense conversion task, the latent code z is from the encoder, and in the word generation task, the latent code z is a sample from $\mathcal{N}(0, 1)$.

In the decoder, we take one token at a time and perform embedding and activate by ReLU. Then we use LSTM with the token, hidden state, cell state as inputs and output the output token, new hidden state, and new cell state.

```

1 class DecoderRNN(nn.Module):
2     def __init__(self, hidden_size, output_size, c_size, c_hidden_size, z_size,
3         device):
4         super(DecoderRNN, self).__init__()
5         self.hidden_size = hidden_size
6         self.output_size=output_size
7         self.c_size=c_size
8         self.c_hidden_size=c_hidden_size
9         self.z_size=z_size
10        self.device=device
11
12        self.c_embedding=nn.Embedding(c_size,c_hidden_size)
13        self.hidden_fc=nn.Linear(z_size+c_hidden_size,hidden_size)
14        self.cell_fc=nn.Linear(z_size+c_hidden_size,hidden_size)
15        self.embedding = nn.Embedding(output_size, hidden_size)
16        self.lstm = nn.LSTM(hidden_size,hidden_size)
17        self.fc = nn.Linear(hidden_size, output_size)
18
19    def forward(self,x,hidden_state,cell_state):
20        x=self.embedding(x).reshape(-1,1,self.hidden_size)
21        x=F.relu(x)
22        out,(hidden_state,cell_state)=self.lstm(x,(hidden_state,cell_state))
23        out=self.fc(out).reshape(-1,self.output_size)
24        return out,hidden_state,cell_state
25
26    def init_hidden_and_cell(self,c,hid_z,cell_z):
27        c=self.c_embedding(c).reshape(1,1,-1)
28        hid_state=torch.cat((hid_z,c),dim=2)
29        hid_state=self.hidden_fc(hid_state)
30
31        cell_state=torch.cat((cell_z,c),dim=2)
32        cell_state=self.cell_fc(cell_state)
33
34        return hid_state,cell_state

```

Listing 4: CVAE decoder

```

1 def decode(decoder,tokenizer,device,token,hid_z,cell_z,c,is_tf):
2     in_token=torch.from_numpy(np.asarray(tokenizer.sos))
3     in_token=in_token.to(device, dtype=torch.long)
4     out_distribution=[]
5     if token==None:
6         max_len=30
7     else:
8         max_len=token.shape[0]-1
9
10    hidden_state,cell_state=decoder.init_hidden_and_cell(c, hid_z, cell_z)
11    for i in range(max_len):
12        output,hidden_state,cell_state=decoder(in_token,hidden_state,cell_state)
13        out_distribution.append(output)
14        out_token=torch.max(torch.softmax(output,dim=1),1)[1]
15        if out_token.item()==tokenizer.eos:
16            break
17        if is_tf==True:
18            in_token=token[i+1]
19        else:
20            in_token=out_token
21    out_distribution=torch.cat(out_distribution,dim=0).to(device)
22    return out_distribution

```

Listing 5: The decoding process.

4.4 KL weight annealing

For monotonic annealing, in the first 20% of total epochs, the KL weight is set to 0, and after that, we increase linearly to the final KL weight, which we set as **0.2**. For cyclical annealing, we do the same thing as the monotonic method for each cycle.


```

1 def klw_sched(anneal_method,cur_epoch,epochs,final_klw,anneal_cyc):
2     if anneal_method=="monotonic":
3         thres=0.2*epochs
4         if cur_epoch<=thres:
5             kl_w=0
6         else:
7             kl_w=final_klw*(cur_epoch-thres)/(epochs-thres)
8     elif anneal_method=="cyclic":
9         T=int(epochs/anneal_cyc)
10        thres=int(T*0.2)
11        cur_epoch%=T
12
13        if cur_epoch<=thres:
14            kl_w=0
15        else:
16            kl_w=final_klw*(cur_epoch-thres)/(T-thres)
17
18    return kl_w

```

Listing 6: KL weight scheduler.

4.5 Teacher forcing ratio

In the first 20% of the total epochs, we set the teacher forcing ratio to 1. Furthermore, after that, we linearly decrease to final teacher forcing ratio, which we set to **0.8**.

```

1 def tf_sched(cur_epoch,epochs,final_tf_ratio):
2     thres=int(0.2*epochs)
3     if cur_epoch<thres:
4         tf_ratio=1
5     else:
6         tf_ratio=final_tf_ratio+
7             (epochs-cur_epoch)*(1-final_tf_ratio)/(epochs-thres)
8     if tf_ratio>1:
9         tf_ratio=1
10    elif tf_ratio<final_tf_ratio:
11        tf_ratio=final_tf_ratio
12    return tf_ratio

```

Listing 7: Teacher forcing ratio scheduler.

4.6 Word generator

In word generation task, we random sample the latent z from $\mathcal{N}(0,1)$. And for the four tenses, the condition c is range from 0...3. For each condition c , we give the decoder the latent z and the condition, then the decoder will generate a word with proper tense. We randomly generate 100 words, each with 4 tenses.

```

1 def gen_word(decoder,z_size,device,tokenizer):
2     words_list=[]
3     for i in range(100):
4         hid_z,cell_z=gen_gauss_noise(z_size)
5         hid_z=hid_z.to(device,dtype=torch.float)
6         cell_z=cell_z.to(device,dtype=torch.float)
7         words=[]
8         for c in range(4):
9             c=torch.from_numpy(np.asarray(c))
10            c=c.to(device,dtype=torch.long)
11            outputs=decode(decoder,tokenizer,device,None,hid_z,cell_z,c,
12                           is_tf=False)
13            out_token=torch.max(torch.softmax(outputs,dim=1),1)[1]
14            out_word=tokenizer.inv_tokenize(out_token)
15            words.append(out_word)
16        words_list.append(words)
17    return words_list

```

Listing 8: Word generation.

4.7 Hyperparameters

- Final KL weight: 0.1.
- Learning rate: 0.05.
- Final teacher forcing ratio: 0.8.
- Epochs: 150.

5 Experimental Results

Some of the hyperparameters have been mentioned in Section 4. We provide the full lists here:

- Input size: 28.
- Hidden layer size: 256.
- Input condition size: 4.
- Condition hidden layer size: 8.
- Latent size: 32.
- Final KL weight: 0.1.
- Learning rate: 0.05.
- Final teacher forcing ratio: 0.8.
- Epochs: 150.

KL annealing	BLEU score	Gaussian score
monotonic	0.8312	0.472
cyclical	0.9527	0.452

Table 1: Best results of the two KL annealing method

5.1 Tense conversion

<input: functioning<br=""></input:> target: functioned prediction: functioned <input: abet<br=""></input:> target: abetting prediction: abeting <input: functioning<br=""></input:> target: function prediction: function <input: split<br=""></input:> target: splitting prediction: splitting <input: abandon<br=""></input:> target: abandoned prediction: abandon <input: flared<br=""></input:> target: flare prediction: flan <input: begin<br=""></input:> target: begins prediction: begins <input: sent<br=""></input:> target: sends prediction: sents Average BLEU-4 score: 0.8036	<input: sent<br=""></input:> target: sends prediction: sends <input: functioning<br=""></input:> target: functioned prediction: functioned <input: flared<br=""></input:> target: flare prediction: flare <input: abet<br=""></input:> target: abetting prediction: abetting <input: expend<br=""></input:> target: expends prediction: endeats <input: functioning<br=""></input:> target: function prediction: function <input: healing<br=""></input:> target: heals prediction: heals <input: split<br=""></input:> target: splitting prediction: splitting Average BLEU-4 score: 0.9186
(a) Monotonic	(b) Cyclical

Figure 2: Results of tense conversion.

5.2 Word generation

```
[ 'fidget', 'fidgets', 'fidgeting', 'firmulated' ]
[ 'overthrow', 'overthrows', 'overthrowing', 'overthrew' ]
[ 'lurn', 'lives', 'livering', 'lived' ]
[ 'assure', 'assures', 'assuming', 'assured' ]
[ 'barr', 'barges', 'barring', 'barred' ]
[ 'befall', 'befalls', 'befalling', 'befalled' ]
[ 'appear', 'appears', 'appearing', 'appeared' ]
[ 'pick', 'picks', 'picking', 'picked' ]
[ 'scramble', 'scrambles', 'scrambling', 'scrambled' ]
[ 'backslide', 'backslides', 'backsliding', 'backslided' ]
[ 'spend', 'spends', 'preserving', 'preserved' ]
[ 'intermingle', 'intermingles', 'intermingling', 'intermingled' ]
[ 'grab', 'grasps', 'grabbing', 'grabbed' ]
[ 'abound', 'abounds', 'abounding', 'abounded' ]
[ 'bear', 'bears', 'bearing', 'bearded' ]
[ 'speak', 'speaks', 'speaking', 'speaked' ]
[ 'reply', 'returns', 'retouching', 'retouched' ]
[ 'assure', 'assures', 'assuring', 'assured' ]
[ 'convert', 'resettles', 'resettling', 'resettled' ]
[ 'indulge', 'indulges', 'indulging', 'indulged' ]
[ 'bristle', 'brightens', 'bristling', 'bristled' ]
[ 'provoke', 'provokes', 'provoking', 'provoked' ]
[ 'hurries', 'hurries', 'houseing', 'hurried' ]
[ 'dive', 'dives', 'diving', 'dived' ]
[ 'glower', 'glows', 'glowing', 'glowered' ]
[ 'oblige', 'obliges', 'obliging', 'obliged' ]
[ 'swill', 'swills', 'swilling', 'swilled' ]
[ 'brow', 'brows', 'browing', 'browled' ]
[ 'dangle', 'dangles', 'disagreeing', 'dangled' ]
[ 'behold', 'beholds', 'beholding', 'behold' ]
[ 'assay', 'assays', 'assaying', 'assayed' ]
[ 'illumine', 'illumines', 'illumining', 'illuminated' ]
[ 'undrie', 'undress', 'undrising', 'undred' ]
[ 'enact', 'enacts', 'enacting', 'enacted' ]
Gaussian score: 0.4200
```

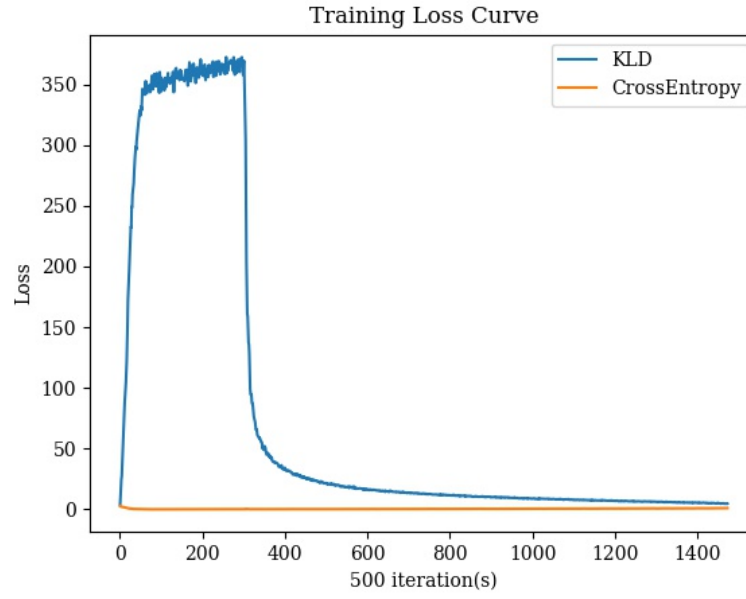
(a) Monotonic

```
[ 'comprise', 'comprises', 'comprising', 'comprised' ]
[ 'cestrate', 'cestrates', 'cestining', 'cestrated' ]
[ 'inflict', 'inlays', 'involving', 'inflicted' ]
[ 'supervise', 'supervises', 'supervising', 'supervised' ]
[ 'sirclease', 'sircles', 'sircling', 'signalled' ]
[ 'peer', 'peers', 'peering', 'peered' ]
[ 'compile', 'combines', 'compiling', 'compiled' ]
[ 'bicycle', 'bids', 'bidding', 'bided' ]
[ 'unwrap', 'unwraps', 'unwrapping', 'unwrapped' ]
[ 'receive', 'receives', 'receiving', 'received' ]
[ 'holler', 'holds', 'holding', 'holded' ]
[ 'correspond', 'remembers', 'remembering', 'corresded' ]
[ 'lunge', 'lunges', 'lunging', 'lunged' ]
[ 'soirl', 'soirs', 'soiring', 'sought' ]
[ 'father', 'fathers', 'averting', 'fathered' ]
[ 'disunite', 'disunites', 'disuniting', 'disunited' ]
[ 'shove', 'shoves', 'shoving', 'shoved' ]
[ 'cough', 'coughs', 'coughing', 'coughed' ]
[ 'destroy', 'destroys', 'destroying', 'destroyed' ]
[ 'purchase', 'purses', 'putting', 'putted' ]
[ 'transfer', 'transfers', 'transferring', 'transferred' ]
[ 'calculate', 'calculates', 'calculating', 'calculated' ]
[ 'blish', 'blishes', 'blishing', 'blished' ]
[ 'gulp', 'gulps', 'gulping', 'gulled' ]
[ 'continue', 'continues', 'contributing', 'continued' ]
[ 'inquire', 'inquires', 'inquiring', 'inquired' ]
[ 'transform', 'transferrs', 'transforming', 'transferred' ]
[ 'avow', 'avows', 'avowing', 'avowed' ]
[ 'advocate', 'advocates', 'advocating', 'advocated' ]
[ 'snarl', 'rears', 'snarling', 'rared' ]
[ 'applie', 'applies', 'appling', 'applied' ]
[ 'seize', 'seizes', 'glessing', 'seized' ]
[ 'overthrow', 'overthrows', 'overthrowing', 'overthrew' ]
[ 'rrash', 'rrashes', 'regaining', 'rrashed' ]
Gaussian score: 0.4600
```

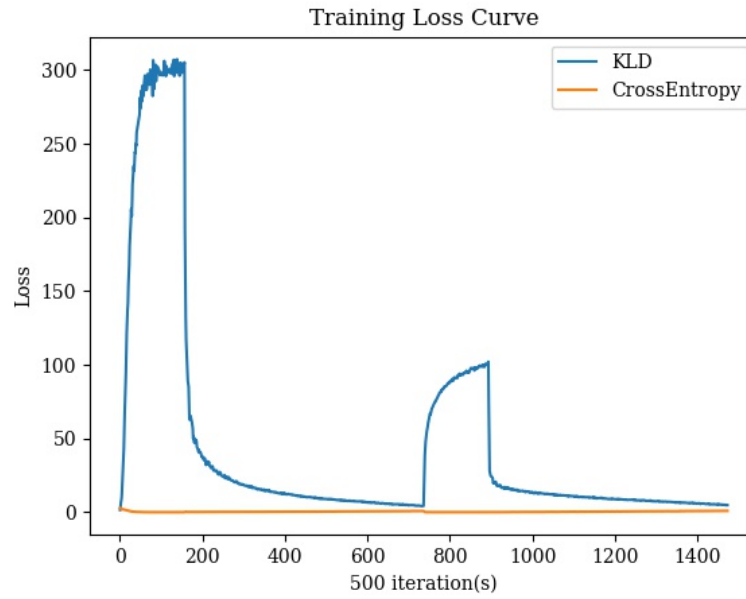
(b) Cyclical

Figure 3: Results of word generation.

5.3 Loss curve



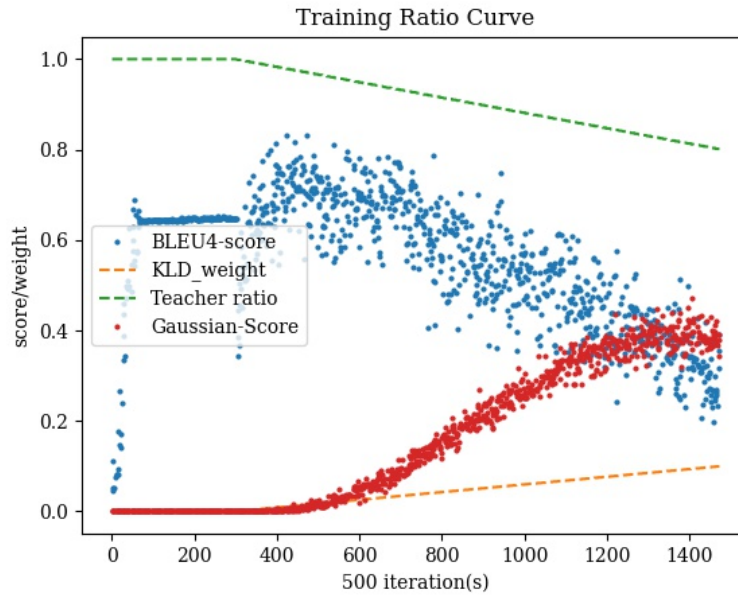
(a) Monotonic



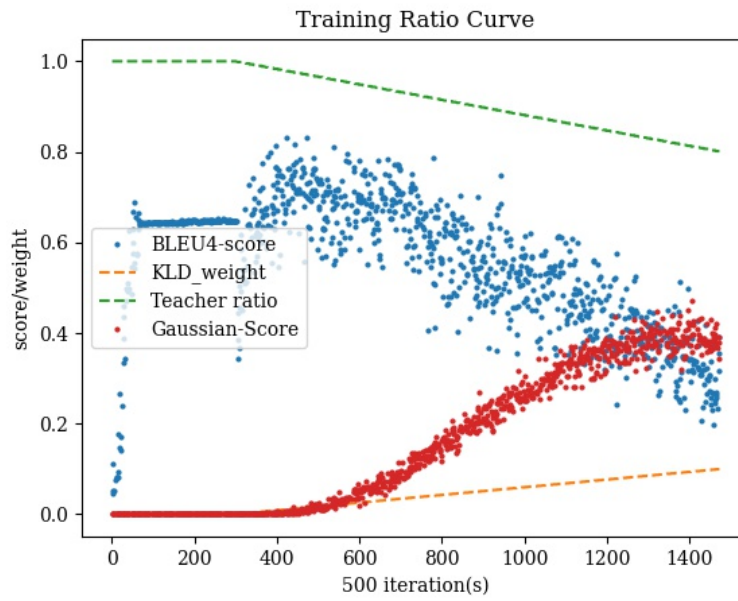
(b) Cyclical

Figure 4: Loss curve.

5.4 Ratio/score



(a) Monotonic



(b) Cyclical

Figure 5: Ratio and score curve.

6 Discussion

How does the teacher forcing ratio affect the results? We have tried out numerous teacher forcing ratios. In extreme case 0, the model performs poorly. When we set the teacher forcing ratio to 0, the model only learns from the output in the previous step, so if the output is incorrect, the results will always be wrong. In another extreme case 1, the bleu score rise to 78% then stopped ever since. Then we have tried 0.4, 0.6, and 0.8, and the bleu score increases when the teacher forcing ratio increases. The best bleu score we get is when the teacher forcing ratio is set to 0.8.

How does KLD weight affect the results? We found out that KLD weights have a significant effect on the Gaussian score. At the beginning of the training, we set the KLD weight to 0. In this phase, the Gaussian score never exceeds 0.02. Then after we start to raise the KLD weights, the Gaussian score rises as well. In KL cyclical annealing, KLD weight will drop back to 0 in every cycle. In this phase, the Gaussian score drops as well! KLD weight and bleu score work the opposite; when the KLD weight rises, the bleu score drops. We have tried 0.1,0.2,0.4 as our KLD weight, and we found out that 0.1 yields the best results. Since this is the regularization, we think 0.2 and 0.4 regularize too much, so the accuracy suffers.

Which KLD annealing method is better? No matter which method we choose, the Gaussian score can be around 0.4. However, the BLEU score is higher when using the cyclical annealing method, so we think the cyclical annealing method is better than the monotonic annealing method.