

Part 1: Long Short-Term Memory

1. The code for this question is shown below.

```

1  class MyLSTMCell(nn.Module):
2      def __init__(self, input_size, hidden_size):
3          super(MyLSTMCell, self).__init__()
4
5          self.input_size = input_size
6          self.hidden_size = hidden_size
7
8          # -----
9          # FILL THIS IN
10         # -----
11         self.Wii = nn.Linear(input_size, hidden_size)
12         self.Whi = nn.Linear(hidden_size, hidden_size)
13
14         self.Wif = nn.Linear(input_size, hidden_size)
15         self.Whf = nn.Linear(hidden_size, hidden_size)
16
17         self.Wig = nn.Linear(input_size, hidden_size)
18         self.Whg = nn.Linear(hidden_size, hidden_size)
19
20         self.Wio = nn.Linear(input_size, hidden_size)
21         self.Who = nn.Linear(hidden_size, hidden_size)
22
23
24     def forward(self, x, h_prev, c_prev):
25         """Forward pass of the LSTM computation for one time step.
26
27         Arguments
28             x: batch_size x input_size
29             h_prev: batch_size x hidden_size
30             c_prev: batch_size x hidden_size
31
32         Returns:
33             h_new: batch_size x hidden_size
34             c_new: batch_size x hidden_size
35         """
36
37         # -----
38         # FILL THIS IN
39         # -----
40         i = torch.sigmoid(self.Wii(x) + self.Whi(h_prev))
41         f = torch.sigmoid(self.Wif(x) + self.Whf(h_prev))
42         g = torch.tanh(self.Wig(x) + self.Whg(h_prev))
43         o = torch.sigmoid(self.Wio(x) + self.Who(h_prev))
44         c_new = torch.mul(f, c_prev) + torch.mul(i, g)
45         h_new = torch.mul(o, torch.tanh(c_new))
46         return h_new, c_new
47

```

According to Figure 1, there is no large difference between models using the small dataset and large dataset. The LSTM model using small dataset performs better during the training, while the validation loss of the large dataset model is lower, which is reasonable. Since more training examples could reduce the impact of overfitting, it might be the reason for the higher training loss and lower validation loss for the smaller dataset model.

2. I tested five sentences and results are shown below.

```

1  source:   his son is an eight-year-old student
2  translated: ishay onsay isway anyway ewtigay-ouncehsay udentday
3  correct:  ishay onsay isway anyway eightway-earyay-oldway udentstay
4

```

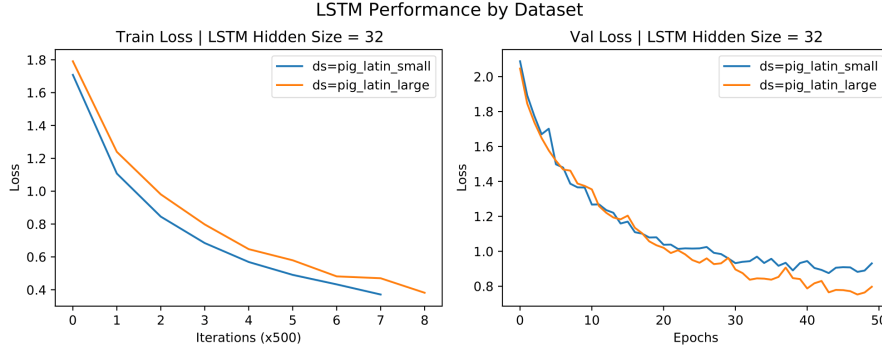


Figure 1: Train and validation loss comparison by dataset.

```

5 source: the model with smaller dataset performs significantly better
6 translated: ethay odedmay ithway alerlsmay atesstay erstoursway intigslaingsay
   ettetyway
7 correct: ethay odelmay ithway allersmay atasetday erformsay ignificantlysay
   etterbay
8
9 source: people are facing a severe situation of resources shortage and
   environmental pollution and degradation of ecosystems
10 translated: eoplephay areway acingfay away eerkenay ituationsday ofway esurechtray
   orpaleday andway indenverelngway olluptionway andway egdgatinesway ofway
   eccospestlay
11 correct: eoplepay areway acingfay away everesay ituationsay ofay esourcesray
   ortageshay andway environmentalway ollutionpay andway egradationday ofway
   ecosystemsway
12
13 source: turning off the lights when leaving is not common
14 translated: untringway offay ethay ylustlay engay ealundway isway otnay ommouncay
15 correct: urningtay offway ethay ightslay enwhay eavinglay isway otnay ommoncay
16
17 source: there are many other ways to achieve a low-carbon lifestyle
18 translated: eterhay areway anmay otherway ackay otay achemhway away owlungay-ag-
   odhay ifulestleway
19 correct: erethay areway anmay otherway aysway otay achieveay away owl-arboncay
   ifestylelay
20

```

One significant failure is that the model performs badly when dealing with the word including the dash symbol. The dash symbol is either missing or too many, which may count two words as one, resulting in the incorrect combination of transformed word. For example, “eight-year-old” becomes “ewtigay-ouncehsay”. “low-carbon” is translated to “owlungay-ag-odhay”.

- For each character, we trained the input, output, forget, and tanh gate separately. Since the output includes both next hidden units and next cell units, there are $8H$ training units in total. For a K length word, the number of LSTM units should still be $8H$.

For the number of connections, since the dimension of each gate is H , the number of connections is $4(H + D) \times H$ for training four gates. Then, to generate the new cell state, it requires the input gate, the forget gate, the tanh gate and the previous cell state. Thus, for this process, the number of connections should be $(2H + 2H) \times H$. Similarly, for processing the new hidden state, we just take the product of output gate and new cell state, which gives the number of connections should be $2H$. Thus, the total number of connections is

$$\begin{aligned}
 & (4H^2 + 4DH + 4H * H + 2H)K \\
 & = (8H^2 + 4DH + 2H)K \\
 & = 8H^2K + 4DHK + 2HK
 \end{aligned}$$

1. Based on the given information, we can compute $\tilde{\alpha}_i^{(t)}$, α_i , c_t as follow.

$$\begin{aligned}
\tilde{\alpha}_i^{(t)} &= f(Q_t, K_i) \\
&= W_2 \text{ReLU}(W_1 \text{concat}(Q_t, K_i) + b_1) + b_2 \\
&= W_2 \max(0, W_1 [Q_t, K_i] + b_1) + b_2 \\
\alpha_i^{(t)} &= \text{softmax}(\tilde{\alpha}_i^{(t)}) \\
&= \frac{\exp(\tilde{\alpha}_i^{(t)})}{\sum_t \exp(\tilde{\alpha}_i^{(t)})} \\
c_t &= \sum_i \alpha_i^{(t)} K_i
\end{aligned}$$

2. (Optional)

3. (Optional)

4. For each character, we trained the hidden state and attention weight separately. For the learning parameters from the concatenation, since it concatenates each query and corresponding keys by sequences, there are $2H$ training units for one sequence. After activation, the model would train them again to get the attention weight, which also requires $2H$ training units for one sequence. In addition, for RNN step, since we used the model in Part one, there are $8H$ LSTM units for training. Thus, for a K length word, the number of LSTM units should be $12H$.

For the number of connections, let's calculate it step by step. Since it calculates each attention by sequences, there are K steps. For each step, it firstly get the attention based on the query and the key. After the concatenation, it trained the attention by two fully connected layer with one activation function. It firstly translates from $2H$ to H , and then H to 1. After we get the normalized attention, the model also get the context by multiplication of the attention weight and values. Thus, the number of connections should be $2H^2K^2 + HK^2 + HK^2$. Then, it goes to RNN layer to generation the prediction based on the hidden unit and the cell unit. It requires $8H^2 + 4DH + 2H$ connections from part one. Thus, the total number of connections is

$$\begin{aligned}
&(2H^2K^2 + HK^2 + HK^2 + 8H^2 + 4DH + 2H)K \\
&= (2H^2K^2 + 2HK^2 + 8H^2 + 4DH + 2H)K \\
&= 2H^2K^3 + 2HK^3 + 8H^2K + 4DHK + 2HK
\end{aligned}$$

At the end, it connects to the fully connected layer to output the probabilities for each vocabulary. Since the size of the output layer is KV and the previous concatenation has a shape of HK , the number of connections should be HK^2V . Thus, the total number of connections is $2H^2K^3 + 2HK^3 + 8H^2K + 4DHK + 2HK + HK^2V$.

Part 3: Scaled Dot Product Attention

1. The code for this question is shown below.

```

1 class ScaledDotAttention(nn.Module):
2     def __init__(self, hidden_size):
3         super(ScaledDotAttention, self).__init__()
4
5         self.hidden_size = hidden_size
6
7         self.Q = nn.Linear(hidden_size, hidden_size)
8         self.K = nn.Linear(hidden_size, hidden_size)
9         self.V = nn.Linear(hidden_size, hidden_size)
10        self.softmax = nn.Softmax(dim=1)
11        self.scaling_factor = torch.rsqrt(torch.tensor(self.hidden_size, dtype=
torch.float))
12
13    def forward(self, queries, keys, values):

```

```

14         """The forward pass of the scaled dot attention mechanism.
15
16         Arguments:
17             queries: The current decoder hidden state, 2D or 3D tensor. (
18                 batch_size x (k) x hidden_size)
19             keys: The encoder hidden states for each step of the input sequence. (
20                 batch_size x seq_len x hidden_size)
21             values: The encoder hidden states for each step of the input sequence.
22                 (batch_size x seq_len x hidden_size)
23
24         Returns:
25             context: weighted average of the values (batch_size x k x hidden_size)
26             attention_weights: Normalized attention weights for each encoder
27             hidden state. (batch_size x seq_len x 1)
28
29         The output must be a softmax weighting over the seq_len annotations.
30         """
31
32         # -----
33         # FILL THIS IN
34         # -----
35         batch_size = queries.size()[0]
36         # 2D to 3D
37         if len(queries.size()) == 2:
38             queries = queries.view(batch_size, 1, self.hidden_size)
39         q = self.Q(queries)
40         k = self.K(keys)
41         v = self.V(values)
42         unnormalized_attention = torch.bmm(k, torch.transpose(q, 1, 2)) * self.
scaling_factor
43         attention_weights = self.softmax(unnormalized_attention)
44         context = torch.bmm(torch.transpose(attention_weights, 1, 2), v)
45         return context, attention_weights

```

2. The code for this question is shown below.

```

1 class CausalScaledDotAttention(nn.Module):
2     def __init__(self, hidden_size):
3         super(CausalScaledDotAttention, self).__init__()
4
5         self.hidden_size = hidden_size
6         self.neg_inf = torch.tensor(-1e7)
7
8         self.Q = nn.Linear(hidden_size, hidden_size)
9         self.K = nn.Linear(hidden_size, hidden_size)
10        self.V = nn.Linear(hidden_size, hidden_size)
11        self.softmax = nn.Softmax(dim=1)
12        self.scaling_factor = torch.rsqrt(torch.tensor(self.hidden_size, dtype=
torch.float))
13
14    def forward(self, queries, keys, values):
15        """The forward pass of the scaled dot attention mechanism.
16
17        Arguments:
18            queries: The current decoder hidden state, 2D or 3D tensor. (
19                batch_size x (k) x hidden_size)
20            keys: The encoder hidden states for each step of the input sequence. (
21                batch_size x seq_len x hidden_size)
22            values: The encoder hidden states for each step of the input sequence.
23                (batch_size x seq_len x hidden_size)
24
25        Returns:
26            context: weighted average of the values (batch_size x k x hidden_size)
27            attention_weights: Normalized attention weights for each encoder
28            hidden state. (batch_size x seq_len x 1)
29
30        The output must be a softmax weighting over the seq_len annotations.

```

```

27     """
28
29     # -----
30     # FILL THIS IN
31     # -----
32     batch_size = queries.size()[0]
33     # 2D to 3D
34     if len(queries.size()) == 2:
35         queries = queries.view(batch_size, 1, self.hidden_size)
36     q = self.Q(queries)
37     k = self.K(keys)
38     v = self.V(values)
39     unnormalized_attention = torch.bmm(k, torch.transpose(q, 1, 2)) * self.
scaling_factor
40     mask = torch.tril(torch.ones_like(unnormalized_attention) * self.neg_inf,
diagonal=-1)
41     # print(mask)
42     attention_weights = self.softmax(torch.triu(unnormalized_attention) + mask
)
43     context = torch.bmm(torch.transpose(attention_weights, 1, 2), v)
44     return context, attention_weights
45

```

3. The positional encoding is required to process the information about the order and the position of each sequence in a context since the previous steps do not give any positional information to the model while the word is context-related. It is better since it not only gives the information about the absolute position in a context, but also tells the relative position within sequences.
4. The validation loss of this model is 0.7783684799447655, which is the highest loss compared with the previous two, where the additive attention model gets 0.2094228501940961 validation loss and the best validation loss of LSTMs model is 0.7528372138308791. The transformers failed to translate “air” correctly, and the LSTMs model performed badly when facing the word “conditioning”, while additive attention model translate the text sentence correctly.
5. The plots and the loss results are shown below.

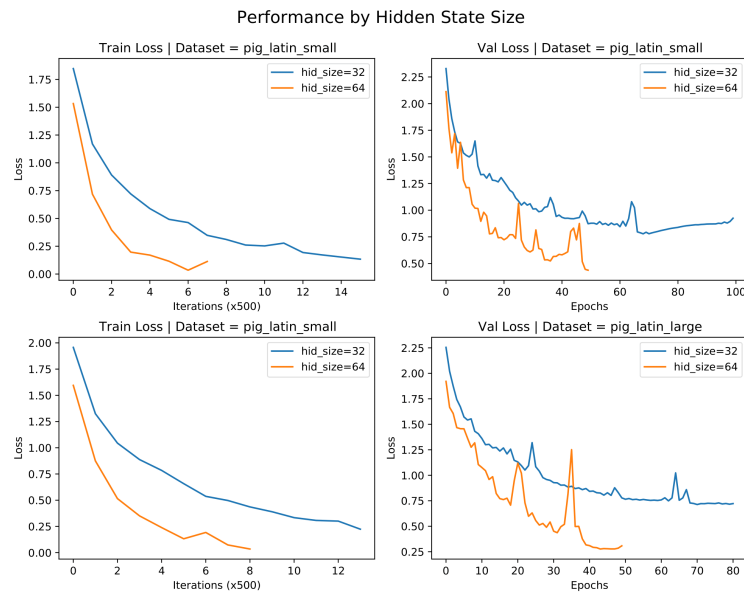


Figure 2: Train and validation loss comparison by hidden size.

According to Figure 2, we can see that the hidden state size affect the model performance a lot, where larger hidden state size would get lower train and validation loss. According to Figure 3,



Figure 3: Train and validation loss comparison by dataset size.

Hidden size	Dataset size	small	large
32		0.7783684799447655	0.7129870925206118
64		0.4370648428797722	0.276622651008298

Table 1: The best validation loss of four models with different hidden size and dataset size.

we found that the model using larger dataset would perform relatively better, but not significant. As we increased the number of iterations, the training loss would generally decrease, while the validation loss follows a decreasing trend with some oscillation. These results follow what I expected, since more hidden state size gain more trainable parameters and larger dataset provide more examples.

6. (Optional)
7. (Optional)

Part 4: Fine-tuning for arithmetic sentiment analysis

1. Similar to the BERT model, I simply add a fully connected layer as the classifier where it takes previous predictions and outputs the probabilities of each label. The code for this question is shown below.

```

1 from transformers import OpenAIGPTForSequenceClassification
2 class GPTCSC413(OpenAIGPTForSequenceClassification):
3     def __init__(self, config):
4         super(GPTCSC413, self).__init__(config)
5         # Your own classifier goes here
6         self.score = nn.Linear(config.hidden_size, self.config.num_labels)
7 
```

2. (Optional)

3. (Optional)
4. GPT architecture is more preferred than BERT in the language generation task such as left-to-right language generation. GPT model is the autoregressive language model, which generate predictions based on the previous several inputs, while BERT model is a self-encoding language model, which does not consider the correlation between the prediction. Since BERT model has a different pre-training process and the generation process, it leads a poorer performance on the language generation task comparing with the GPT model.
5. (Optional)