

CSC413: Programming Assignment 3: Attention-Based Neural Machine Translation

Anzhe Dong, 1002608163

March 2021

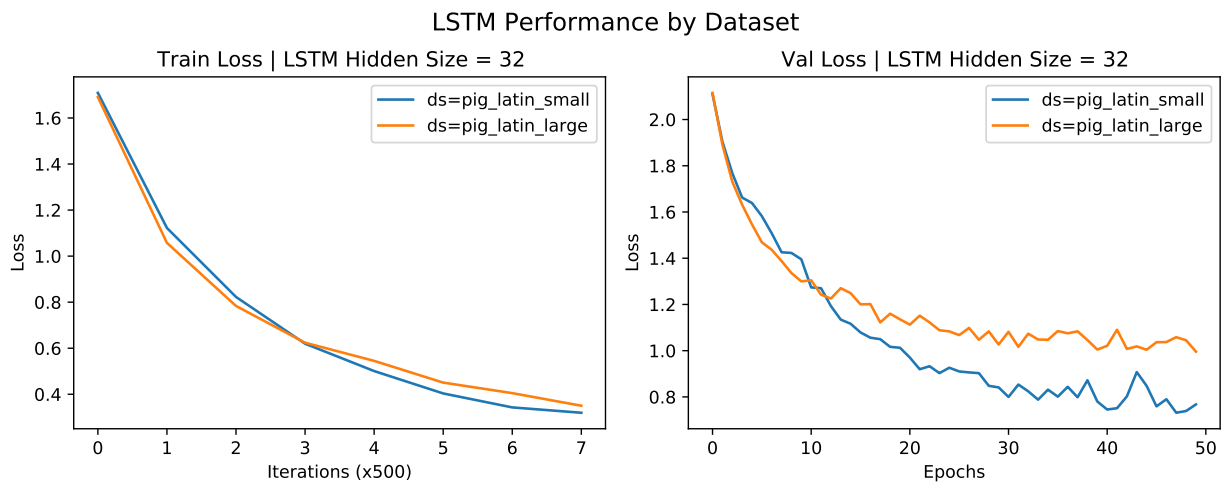
Part 1: LSTMs

1. LSTM Training

A screenshot of your fullMyLSTMCell implementation

```
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 = f * c_prev + i * g
45        h_new = o * torch.tanh(c_new)
46        return h_new, c_new
```

... the loss plots output by save_loss.comparison_lstm



... and your analysis

Does either model perform significantly better? Why might this be the case?

The `pig_latin_small` model performed better with lower validation loss. The small model was trained with smaller batch sizes (`'batch_size':64`), which tend to have better generalization ability than ones trained with larger batch sizes (`'batch_size':512`). [Reference: Keskar NS, Mudigere D, Nocedal J, Smelyanskiy M, Tang PT. On large-batch training for deep learning: Generalization gap and sharp minima. arXiv preprint arXiv:1609.04836. 2016 Sep 15.]

Model Failure

Identify a distinct failure mode and briefly describe it.

```

1 source:  the air conditioning is working
2 translated: ethay airway oniningbay-intway isway orkingway
3
4 source:  this has a lot of errors
5 translated: isthay ashay away otlay ofay errorsway
6
7 source:  the princess listened peacefully to what the frogs had to sing
8 translated: ethay insecspray istenedlay eaffulecay otay athay ethay ogsgay adhay otay ingsay

```

The model fails apparently at (1) when the leading character is an vowel, (2) when the leading consonant pairs such as "sh" and "wh", (3) long words.

Model size

Write down the number of neurons and connections of this encoder model as a function of H, K, and D. For simplicity, you may ignore the bias units.

Number of neurons:
 Number of connections:

Part 2: Additive Attention

1

$$\tilde{\alpha}_i^{(t)} = W_2^T \text{ReLU}(W_1^T \begin{bmatrix} Q_t \\ K_i \end{bmatrix})$$

$$\alpha_i^{(t)} = \text{softmax}(\tilde{\alpha}^{(t)})_i = \frac{\exp(\tilde{\alpha}_i^{(t)})}{\sum_{j=1}^{\text{seq_len}} \exp(\tilde{\alpha}_j^{(t)})}$$

$$c_t = \alpha^{(t)\top} K = \sum_{i=1}^{\text{seq_len}} \alpha_i^{(t)} K_i$$

2

3

The training speed is faster with the attention model which reached less than 1.0 validation loss within 7 epochs, whereas previously it took the “small” non-attention model 20 epochs to reach this level of validation loss.

Part 3: Scaled Dot Product Attention

1. Implement the scaled dot-product attention mechanism.

A screenshot of your ScaledDotAttention implementation

```

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. (batch_size x (k) x
18            hidden_size)
19            keys: The encoder hidden states for each step of the input sequence. (batch_size
20            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 hidden state. (
27            batch_size x seq_len x 1)
28
29        The output must be a softmax weighting over the seq_len annotations.
30        """
31        # -----
32        # FILL THIS IN
33        # -----
34        batch_size, _i, _j = queries.size()
35        q = self.Q(queries).view(batch_size, -1, self.hidden_size)
36        k = self.K(keys).view(batch_size, -1, self.hidden_size)
37        v = self.V(values).view(batch_size, -1, self.hidden_size)
38        unnormalized_attention = torch.bmm(k, q.transpose(1,2)) * self.scaling_factor
39        attention_weights = self.softmax(unnormalized_attention)
40        context = torch.bmm(attention_weights.transpose(1,2),v)
41        return context, attention_weights

```

2. Implement the causal scaled dot-product attention mechanism.

A screenshot of your CausalScaledDotAttention implementation

```

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. (batch_size x (k) x
19            hidden_size)
20            keys: The encoder hidden states for each step of the input sequence. (batch_size
21            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 hidden state. (
28            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, _i, _j = queries.size()
33         q = self.Q(queries).view(batch_size, -1, self.hidden_size)
34         k = self.K(keys).view(batch_size, -1, self.hidden_size)
35         v = self.V(values).view(batch_size, -1, self.hidden_size)
36         unnormalized_attention = torch.bmm(k, q.transpose(1,2)) * self.scaling_factor
37         mask = self.neg_inf * torch.tril(torch.ones_like(unnormalized_attention), diagonal
=-1)
38         attention_weights = self.softmax(unnormalized_attention + mask)
39         context = torch.bmm(attention_weights.transpose(1,2),v)
40         return context, attention_weights

```

3.

4.

5.