

CSC413PA3

Che Liu, 1002246839

March 19, 2020

Part 1: Gated Recurrent Unit (GRU) [2pt]

1. A screenshot of your full MyGRUCell implementation. [1pt]

```
1 class MyGRUCell(nn.Module):
2     def __init__(self, input_size, hidden_size):
3         super(MyGRUCell, self).__init__()
4
5         self.input_size = input_size
6         self.hidden_size = hidden_size
7
8         # -----
9         # FILL THIS IN
10        # -----
11        ## Input linear layers
12        self.Wiz = nn.Linear(input_size, hidden_size)
13        self.Wir = nn.Linear(input_size, hidden_size)
14        self.Wih = nn.Linear(input_size, hidden_size)
15
16        ## Hidden linear layers
17        self.Whz = nn.Linear(hidden_size, hidden_size)
18        self.Whr = nn.Linear(hidden_size, hidden_size)
19        self.Whh = nn.Linear(hidden_size, hidden_size)
20
21    def forward(self, x, h_prev):
22        """Forward pass of the GRU computation for one time step.
23
24        Arguments
25            x: batch_size x input_size
26            h_prev: batch_size x hidden_size
27
28        Returns:
29            h_new: batch_size x hidden_size
30        """
31
32        # -----
33        # FILL THIS IN
34        # -----
35        sig = nn.Sigmoid()
36        tanh = nn.Tanh()
37
38        z = sig(self.Wiz(x)+self.Whz(h_prev))
39        r = sig(self.Wir(x)+self.Whr(h_prev))
40        g = tanh(self.Wih(x)+r*self.Whh(h_prev))
41        h_new = (1-z)*g+z*h_prev
42        return h_new
```

Figure 1: A screenshot of my full MyGRUCell implementation.

2. The training/validation loss plots. [0pts]

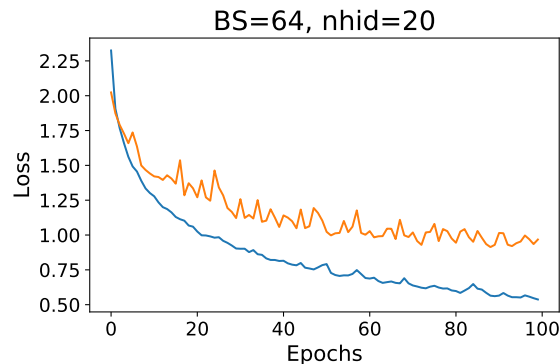


Figure 2: Training and validation loss plot for the GRU.

3. Train the RNN encoder/decoder model. Try a few of your own words. Identify two distinct failure modes and briefly describe them.

Some English to PigLatin outputs of the trained RNN model are as follows:

English: 'the air conditioning is working'

PigLatin (output): 'etay airway ondingingculay isway orkngray'

PigLatin (correct): 'ethay airway onditioningcay isway orkingway'

English: 'i went shopping'

PigLatin (output): 'iway entway opingscay'

PigLatin (correct): 'iway entway oppingshay'

As can be seen in the first example, the model does not seem to learn that the middle portion of the English word should not change. Furthermore, as can be seen in the second example, the model does not seem to learn that consonant pairs such as 'sh' should be moved together.

Part 2: Additive Attention [2pt]

1. Write down the mathematical expression for $\tilde{\alpha}_i^{(t)}$, $\alpha_i^{(t)}$, and c_t as a function of W_1 , W_2 , b_1 , b_2 , Q_t , K_i .

$$\begin{aligned}\tilde{\alpha}_i^{(t)} &= f(Q_t, K_i) = W_2 \text{ReLU}(W_1(\text{concat}(Q_t, K_i)) + b_1) + b_2 \\ \alpha_i^{(t)} &= \text{softmax}(\tilde{\alpha}_i^{(t)})_i \\ c_t &= \sum_{i=1}^T \alpha_i^{(t)} K_i\end{aligned}$$

2. A screenshot of your RNNAttentionDecoder class implementation. [1pt]

```
1 class RNNAttentionDecoder(nn.Module):
2     def __init__(self, vocab_size, hidden_size, attention_type='scaled_dot'):
3         super(RNNAttentionDecoder, self).__init__()
4         self.vocab_size = vocab_size
5         self.hidden_size = hidden_size
6
7         self.embedding = nn.Embedding(vocab_size, hidden_size)
8
9         self.rnn = MyGRUCell(input_size=hidden_size*2, hidden_size=hidden_size)
10        if attention_type == 'additive':
11            self.attention = AdditiveAttention(hidden_size=hidden_size)
12        elif attention_type == 'scaled_dot':
13            self.attention = ScaledDotAttention(hidden_size=hidden_size)
14
15        self.out = nn.Linear(hidden_size, vocab_size)
16
17    def forward(self, inputs, annotations, hidden_init):
18        """Forward pass of the attention-based decoder RNN.
19
20        Arguments:
21            inputs: Input token indexes across a batch for all the time step. (batch_size x decoder_seq_len)
22            annotations: The encoder hidden states for each step of the input.
23                       sequence. (batch_size x seq_len x hidden_size)
24            hidden_init: The final hidden states from the encoder, across a batch. (batch_size x hidden_size)
25
26        Returns:
27            output: Un-normalized scores for each token in the vocabulary, across a batch for all the decoding time steps. (batch_size x decoder_seq_len x vocab_size)
28            attentions: The stacked attention weights applied to the encoder annotations (batch_size x encoder_seq_len x decoder_seq_len)
29
30        """
31
32        batch_size, seq_len = inputs.size()
33        embed = self.embedding(inputs) # batch_size x seq_len x hidden_size
34
35        hiddens = []
36        attentions = []
37        h_prev = hidden_init
38        for i in range(seq_len):
39            # FILL THIS IN - START
40            # -----
41            embed_current = embed[:,i,:] # Get the current time step, across the whole batch
42            context, attention_weights = self.attention(h_prev, annotations, annotations) # batch_size x 1 x hidden_size
43            embed_and_context = torch.cat([np.squeeze(context, axis = 1), embed_current], dim = 1) # batch_size x (2*hidden_size)
44            h_prev = self.rnn(embed_and_context, h_prev) # batch_size x hidden_size
45            # -----
46            # FILL THIS IN - START
47            # -----
48            hiddens.append(h_prev)
49            attentions.append(attention_weights)
50
51        hiddens = torch.stack(hiddens, dim=1) # batch_size x seq_len x hidden_size
52        attentions = torch.cat(attentions, dim=2) # batch_size x seq_len x seq_len
53
54        output = self.out(hiddens) # batch_size x seq_len x vocab_size
55        return output, attentions
```

Figure 3: A screenshot of my RNNAttentionDecoder class implementation.

Part 3: Scaled Dot Product Attention [4pt]

1. Screenshots of your ScaledDotProduct, CausalScaledDotProduct, TransformerEncoder and TransformerDecoder implementations. Highlight the lines you've added. [2pt]

```
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 hidden_size)
18            keys: The encoder hidden states for each step of the input sequence. (batch_size x seq_len x hidden_size)
19            values: The encoder hidden states for each step of the input sequence. (batch_size x seq_len x hidden_size)
20
21        Returns:
22            context: weighted average of the values (batch_size x k x hidden_size)
23            attention_weights: Normalized attention weights for each encoder hidden state. (batch_size x seq_len x k)
24
25        The output must be a softmax weighting over the seq_len annotations.
26        """
27
28        # -----
29        # FILL THIS IN
30        # -----
31        batch_size = keys.size(0)
32        q = self.Q(queries)
33        if len(q.size()) == 2:
34            q = q.view(batch_size, 1, hidden_size)
35        k = self.K(keys)
36        v = self.V(values)
37        unnormalized_attention = torch.bmm(k, q.transpose(1,2))*self.scaling_factor
38        attention_weights = self.softmax(unnormalized_attention)
39        context = torch.bmm(attention_weights.transpose(1,2), v)
40        return context, attention_weights
```

Figure 4: A screenshot of my implementation of the scaled dot-product attention mechanism.

```
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 hidden_size)
19            keys: The encoder hidden states for each step of the input sequence. (batch_size x seq_len x hidden_size)
20            values: The encoder hidden states for each step of the input sequence. (batch_size x seq_len x hidden_size)
21
22        Returns:
23            context: weighted average of the values (batch_size x k x hidden_size)
24            attention_weights: Normalized attention weights for each encoder hidden state. (batch_size x seq_len x k)
25
26        The output must be a softmax weighting over the seq_len annotations.
27        """
28
29        # -----
30        # FILL THIS IN
31        # -----
32        batch_size = keys.size(0)
33        q = self.Q(queries)
34        if len(q.size()) == 2:
35            q = q.view(batch_size, 1, hidden_size)
36        k = self.K(keys)
37        v = self.V(values)
38        unnormalized_attention = torch.bmm(k, q.transpose(1,2))*self.scaling_factor
39        mask = torch.tril(torch.ones(unnormalized_attention.shape), diagonal=-1).cuda()
40        unnormalized_attention = unnormalized_attention+mask*self.neg_inf
41        attention_weights = self.softmax(unnormalized_attention)
42        context = torch.bmm(attention_weights.transpose(1,2), v)
43        return context, attention_weights
```

Figure 5: A screenshot of my implementation of the causal scaled dot-product attention mechanism.

```

1 class TransformerEncoder(nn.Module):
2     def __init__(self, vocab_size, hidden_size, num_layers, opts):
3         super(TransformerEncoder, self).__init__()
4
5         self.vocab_size = vocab_size
6         self.hidden_size = hidden_size
7         self.num_layers = num_layers
8         self.opts = opts
9
10        self.embedding = nn.Embedding(vocab_size, hidden_size)
11
12        # IMPORTANT CORRECTION: NON-CAUSAL ATTENTION SHOULD HAVE BEEN
13        # USED IN THE TRANSFORMER ENCODER.
14        # NEW VERSION:
15        self.self_attentions = nn.ModuleList([ScaledDotAttention(
16            hidden_size=hidden_size,
17            ) for i in range(self.num_layers)])
18
19        # PREVIOUS VERSION:
20        self.self_attentions = nn.ModuleList([CausalScaledDotAttention(
21            hidden_size=hidden_size,
22            ) for i in range(self.num_layers)])
23
24        self.attention_mlp = nn.ModuleList([nn.Sequential(
25            nn.Linear(hidden_size, hidden_size),
26            nn.ReLU(),
27            ) for i in range(self.num_layers)])
28
29        self.positional_encodings = self.create_positional_encodings()
30
31    def forward(self, inputs):
32        """Forward pass of the encoder RNN.
33
34        Arguments:
35            inputs: Input token indexes across a batch for all time steps in the sequence. (batch_size x seq_len)
36
37        Returns:
38            annotations: The hidden states computed at each step of the input sequence. (batch_size x seq_len x hidden_size)
39            hidden: The final hidden state of the encoder, for each sequence in a batch. (batch_size x hidden_size)
40        """
41
42        batch_size, seq_len = inputs.size()
43        # =====
44        # FILL THIS IN - START
45
46        encoded = self.embedding(inputs) # batch_size x seq_len x hidden_size
47
48        # Add positional encodings from self.create_positional_encodings. (a'la https://arxiv.org/pdf/1706.03762.pdf, section 3.5)
49        encoded = encoded + self.positional_encodings[:seq_len]
50
51        annotations = encoded
52
53        for i in range(self.num_layers):
54            new_annotations, self_attention_weights = self.self_attentions[i](annotations, encoded, encoded) # batch_size x seq_len x hidden_size
55            residual_annotations = annotations + new_annotations
56            new_annotations = self.attention_mlp[i](residual_annotations)
57            annotations = residual_annotations + new_annotations
58
59        # =====
60        # FILL THIS IN - END
61
62        # Transformer encoder does not have a last hidden layer.
63        return annotations, None

```

Figure 6: A screenshot of my implementation of TransformerEncoder.

```

1 class TransformerDecoder(nn.Module):
2     def __init__(self, vocab_size, hidden_size, num_layers):
3         super(TransformerDecoder, self).__init__()
4
5         self.vocab_size = vocab_size
6         self.hidden_size = hidden_size
7
8         self.embedding = nn.Embedding(vocab_size, hidden_size)
9         self.num_layers = num_layers
10
11        self.self_attentions = nn.ModuleList([CausalScaledDotAttention(
12            hidden_size=hidden_size,
13            ) for i in range(self.num_layers)])
14
15        self.encoder_attentions = nn.ModuleList([ScaledDotAttention(
16            hidden_size=hidden_size,
17            ) for i in range(self.num_layers)])
18
19        self.attention_mlp = nn.ModuleList([nn.Sequential(
20            nn.Linear(hidden_size, hidden_size),
21            nn.ReLU(),
22            ) for i in range(self.num_layers)])
23
24        self.out = nn.Linear(hidden_size, vocab_size)
25
26        self.positional_encodings = self.create_positional_encodings()
27
28    def forward(self, inputs, annotations, hidden_init):
29        """Forward pass of the attention-based decoder RNN.
30
31        Arguments:
32            inputs: Input token indexes across a batch for all the time step. (batch_size x decoder_seq_len)
33            annotations: The encoder hidden states for each step of the input sequence. (batch_size x seq_len x hidden_size)
34            hidden_init: Not used in the Transformer decoder
35
36        Returns:
37            output: Un-normalized scores for each token in the vocabulary, across a batch for all the decoding time steps. (batch_size x decoder_seq_len x vocab_size)
38            attentions: The stacked attention weights applied to the encoder annotations (batch_size x encoder_seq_len x decoder_seq_len)
39        """
40
41        batch_size, seq_len = inputs.size()
42        embed = self.embedding(inputs) # batch_size x seq_len x hidden_size
43
44        # THIS LINE WAS ADDED AS A CORRECTION.
45        embed = embed + self.positional_encodings[:seq_len]
46
47        encoder_attention_weights_list = []
48        self_attention_weights_list = []
49        contexts = embed
50
51        for i in range(self.num_layers):
52            # FILL THIS IN - START
53
54            new_contexts, self_attention_weights = self.self_attentions[i](contexts, contexts, contexts) # batch_size x seq_len x hidden_size
55            residual_contexts = contexts + new_contexts
56            new_contexts, encoder_attention_weights = self.encoder_attentions[i](residual_contexts, annotations, annotations) # batch_size x seq_len x hidden_size
57            residual_contexts = residual_contexts + new_contexts
58            new_contexts = self.attention_mlp[i](residual_contexts)
59            contexts = residual_contexts + new_contexts
60
61            # FILL THIS IN - END
62
63            encoder_attention_weights_list.append(encoder_attention_weights)
64            self_attention_weights_list.append(self_attention_weights)
65
66        output = self.out(contexts)
67        encoder_attention_weights = torch.stack(encoder_attention_weights_list)
68        self_attention_weights = torch.stack(self_attention_weights_list)
69
70        return output, (encoder_attention_weights, self_attention_weights)

```

Figure 7: A screenshot of my implementation of TransformerDecoder.

2. Training/validation plots you've generated. Your response to question 5. Your analysis should not exceed three sentences (excluding the failure cases you've identified). [1pt]

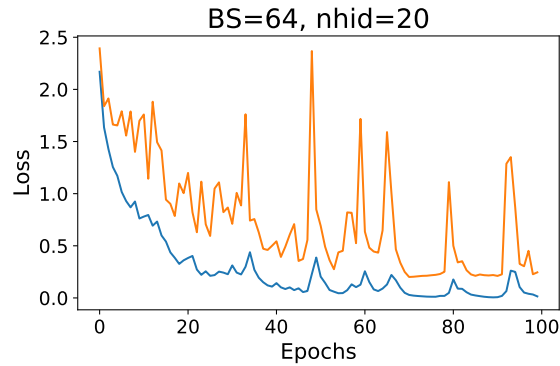


Figure 8: Training and validation loss plot for the Transformer.

Some English to PigLatin outputs of the trained RNN model are as follows:

English: 'the air conditioning is working'

PigLatin (output): 'ethay airway onditioningcay isway orkingway'

PigLatin (correct): 'ethay airway onditioningcay isway orkingway'

English: 'i went shopping'

PigLatin (output): 'iway entway oppingshay'

PigLatin (correct): 'iway entway oppingshay'

The model with Transformer encoder and decoder is significantly better than before. The model translated both sentences perfectly and does not seem to have the failure modes identified before.

3. Your response to question 6. Your response should not exceed three sentences. [1pt]

The training loss went to 0 during training. The trained model always output a single character repeated many times no matter what the input was. This is to be expected because the model would simply learn an identity mapping when the decoder could see the next word (this way the reconstruction error is 0).

Part 4: BERT for arithmetic sentiment analysis [2pt]

1. 10 inference results in question 5 as well as brief comments on why they are interesting or representative results. Your answer should not exceed 3 sentences, you don't need to describe all 10 inference results [1pt]

Ten inference results are shown below:

1.
Input: 'twelve minus fourteen'
Output: 'negative'
2.
Input: 'twelve plus fourteen'
Output: 'positive'
3.
Input: '1 minus 14'
Output: 'negative'
4.
Input: '1 minus twelve'
Output: 'negative'
5.
Input: '1 plus twelve'
Output: 'positive'
6.
Input: '1 minus 1'

Output: 'positive'

7.

Input: 'five plus zero'

Output: 'positive'

8.

Input: 'one plus one plus one'

Output: 'positive'

9.

Input: 'one minus one minus one'

Output: 'negative'

10.

Input: 'one plus one minus one'

Output: 'positive'

Example #4 is interesting because it involves a number and a word that represents a number.

Example #10 is interesting because it involves one operation followed by another.

Example #6 is interesting because it suggests that the trained model did not learn zero.

2. Explanation of what you did for the open question and some preliminary results. Your answer should not exceed 4 sentences. [1pt]

I chose hyperparameter tuning with a grid search for this part. The code is as shown below.

```
1 learning_rates = [1e-5, 2e-5, 3e-5, 4e-5]
2 num_epochs = [5, 6, 7, 8, 9, 10]
3
4 best_hyper = [None, None]
5 best_accuracy = 0
6 for learning_rate in learning_rates:
7     for num_epoch in num_epochs:
8         finttune_bert_loss_vals = train_model(model_finetune_bert, num_epoch, learning_rate)
9         loss_values, val_accuracy = finttune_bert_loss_vals
10        if val_accuracy > best_accuracy:
11            best_hyper = [learning_rate, num_epoch]
12            best_accuracy = val_accuracy
13            print('best hyperparameter:', best_hyper)
14            print('best validation accuracy:', best_accuracy)
```

Figure 9: Grid search for the best learning rate and epoches combination.

The best combination was a learning rate of 1e-5 and number of epoches of 7 which lead to a validation accuracy of 1.