# CSC413: Programming Assignment 3

Tianyu Du (1003801647)
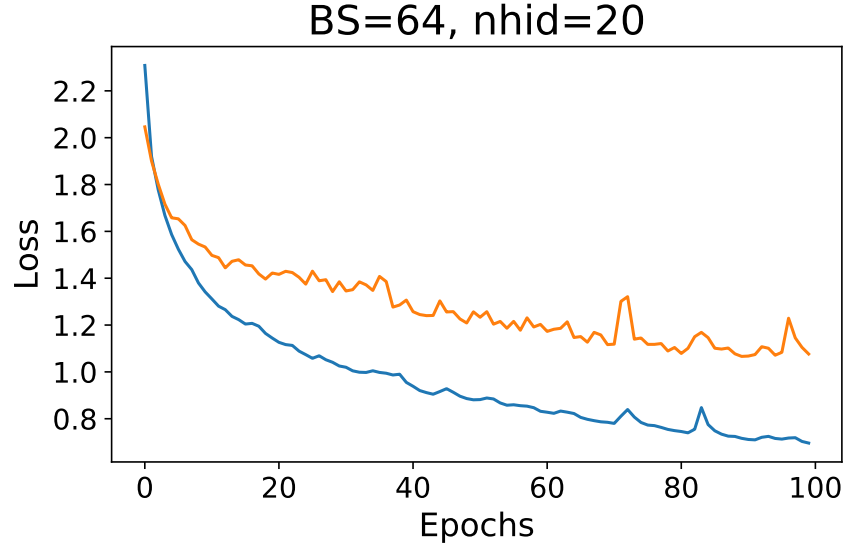
2020/03/16 at 22:51:38

# 1 Part 1: Gated Recurrent Units

## 1.1

```
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, bias=False)
13         self.Wir = nn.Linear(input_size, hidden_size, bias=False)
14         self.Win = nn.Linear(input_size, hidden_size, bias=False)
15
16         ## Hidden linear layers
17         self.Whz = nn.Linear(hidden_size, hidden_size, bias=True)
18         self.Whr = nn.Linear(hidden_size, hidden_size, bias=True)
19         self.Whn = nn.Linear(hidden_size, hidden_size, bias=True)
20
21
22     def forward(self, x, h_prev):
23         # ------------
24         # FILL THIS IN
25         # ------------
26         r = torch.sigmoid(self.Wir(x) + self.Whr(h_prev))
27         z = torch.sigmoid(self.Wiz(x) + self.Whz(h_prev))
28         g = torch.tanh(self.Win(x) + r * self.Whn(h_prev))
29         h_new = (1 - z) * g + z * h_prev
30         return h_new
```

**1.2**



**1.3**

**Failure Type 1** Long vocabularies. The model fails to translate long vocabularies, the middle part of vocabularies get messed up.

```
1 source:      computer science
2 translated: opcorchyway ipenceway
```

**Failure Type 2** Vocabularies containing dashes ("-"). The model fails to distinguish parts of word before and after the dash. Sometime, the dash is missing after translation.

```
1 source:      electric-powered
2 translated: elcercorstentway
3 ===================================
4 source:      to-buy
5 translated: otay-othay
```

# 2 Additive Attention

**2.1**

$$\tilde{\alpha}_i^{(t)} = f(Q_i, K_i) = W_2 \, \text{ReLU}(W_1[Q_i, K_i] + b_1) + b_2 \tag{2.1}$$

$$\alpha_i^{(t)} = \text{softmax}(\tilde{\alpha}^{(t)})_i = \frac{\exp(\tilde{\alpha}_i^{(t)})}{\sum_{t=1}^{\texttt{seq\_len}} \exp(\tilde{\alpha}_i^{(t)})} \tag{2.2}$$

$$c_t = \sum_{t=1}^{\texttt{seq\_len}} \alpha_i^{(t)} K_i \tag{2.3}$$

## 2.2

```python
class RNNAttentionDecoder(nn.Module):
    def __init__(self, vocab_size, hidden_size, attention_type='scaled_dot'):
        super(RNNAttentionDecoder, self).__init__()
        self.vocab_size = vocab_size
        self.hidden_size = hidden_size

        self.embedding = nn.Embedding(vocab_size, hidden_size)

        self.rnn = MyGRUCell(input_size=hidden_size*2, hidden_size=hidden_size)
        if attention_type == 'additive':
          self.attention = AdditiveAttention(hidden_size=hidden_size)
        elif attention_type == 'scaled_dot':
          self.attention = ScaledDotAttention(hidden_size=hidden_size)

        self.out = nn.Linear(hidden_size, vocab_size)


    def forward(self, inputs, annotations, hidden_init):
        """Forward pass of the attention-based decoder RNN.

        Arguments:
            inputs: Input token indexes across a batch for all the time step. (batch_size x
    decoder_seq_len)
            annotations: The encoder hidden states for each step of the input.
                         sequence. (batch_size x seq_len x hidden_size)
            hidden_init: The final hidden states from the encoder, across a batch. (
    batch_size x hidden_size)

        Returns:
            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)
            attentions: The stacked attention weights applied to the encoder annotations (
    batch_size x encoder_seq_len x decoder_seq_len)
        """

        batch_size, seq_len = inputs.size()
        embed = self.embedding(inputs)  # batch_size x seq_len x hidden_size

        hiddens = []
        attentions = []
        h_prev = hidden_init
        for i in range(seq_len):
            # ------------
            # FILL THIS IN - START
            # ------------
            embed_current = embed[:,i,:]  # Get the current time step, across the whole
    batch
            context, attention_weights = self.attention(
                h_prev, # queries @ (bs, hidden_size)
                annotations, # keys @ (bs, sl, hs)
                annotations # values @ (bs, sl, hs)
            )  # @ (batch_size, 1,  hidden_size) and (batch_size, seq_len, 1)
            embed_and_context = torch.cat((
                embed_current.view(batch_size, -1),
                context.view(batch_size, -1)),
                dim=1
            )  # batch_size x (2*hidden_size)
```

```
53        h_prev = self.rnn(embed_and_context, h_prev)  # batch_size x hidden_size
54        # ------------
55        # FILL THIS IN - END
56        # ------------
57
58        hiddens.append(h_prev)
59        attentions.append(attention_weights)
60
61     hiddens = torch.stack(hiddens, dim=1) # batch_size x seq_len x hidden_size
62     attentions = torch.cat(attentions, dim=2) # batch_size x seq_len x seq_len
63
64     output = self.out(hiddens)  # batch_size x seq_len x vocab_size
65     return output, attentions
```

# 3 Scaled Dot Product Attention

## 3.1 Implementations

### 3.1.1 ScaledDotAttention

```python
1 class ScaledDotAttention(nn.Module):
2     def __init__(self, hidden_size):
3         ...
4
5     def forward(self, queries, keys, values):
6         """..."""
7
8         # ------------
9         # FILL THIS IN
10        # ------------
11        hidden_size = self.hidden_size
12        batch_size = queries.shape[0]
13        d = hidden_size
14        # Convert tensor to 3D.
15        # k is the number of queries.
16        queries = queries.view(batch_size, -1, hidden_size)
17        num_queries = queries.shape[1]
18        seq_len = keys.shape[1]
19        # Expand.
20        # keys = keys.expand(batch_size, seq_len, hidden_size)
21        # keys = torch.transpose(keys, dim0=0, dim1=1)
22
23        q = self.Q(queries) # @ (batch_size, k, hidden_size)
24        k = self.K(keys) # @ (batch_size, seq_len, hidden_size)
25        v = self.V(values) # @ (batch_size, seq_len, hidden_size)
26        q = torch.transpose(q, 1, 2) # @ (batch_size, hidden_size, k)
27        # print("q @", q.shape)
28        # print("k @", k.shape)
29        unnormalized_attention = torch.bmm(k, q) * self.scaling_factor
30        # unnormalized_attention @ (batch_size, seq_len, k)
31        # print(unnormalized_attention.shape)
32
33        attention_weights = self.softmax(unnormalized_attention).transpose(1, 2) # @ (
    batch_size, k, seq_len)
34        context = torch.bmm(attention_weights, v) # @ (batch_size, k, hidden_size)
35        attention_weights = attention_weights.transpose(1, 2) # @ (batch_size, seq_len, k)
36        return context, attention_weights
```

### 3.1.2 CausalScaledDotAttention

```python
class CausalScaledDotAttention(nn.Module):
    def __init__(self, hidden_size):
        ...

    def forward(self, queries, keys, values):
        """..."""

        # ------------
        # FILL THIS IN
        # ------------
        hidden_size = self.hidden_size
        batch_size = queries.shape[0]
        d = hidden_size
        # Convert tensor to 3D.
        # k is the number of queries.
        queries = queries.view(batch_size, -1, hidden_size)
        num_queries = queries.shape[1]
        seq_len = keys.shape[1]
        # keys = keys.expand(batch_size, seq_len, hidden_size)
        # keys = torch.transpose(keys, dim0=0, dim1=1)

        q = self.Q(queries) # @ (batch_size, k, hidden_size)
        k = self.K(keys) # @ (batch_size, seq_len, hidden_size)
        v = self.V(values) # @ (batch_size, seq_len, hidden_size)
        q = torch.transpose(q, 2, 1) # @ (batch_size, hidden_size, k)
        # print("q @", q.shape)
        # print("k @", k.shape)
        unnormalized_attention = torch.bmm(k, q) * self.scaling_factor
        # unnormalized_attention @ (batch_size, seq_len, k)
        # print(unnormalized_attention.shape)
        # ==== Enforce Casual ====
        mask = torch.tril(torch.ones_like(unnormalized_attention)) * self.neg_inf
        unnormalized_attention += mask
        # ==== End ====
        attention_weights = self.softmax(unnormalized_attention).transpose(1, 2) # @ (
    batch_size, k, seq_len)
        context = torch.bmm(attention_weights, v) # @ (batch_size, k, hidden_size)
        attention_weights = attention_weights.transpose(1, 2) # @ (batch_size, seq_len, k)
        return context, attention_weights
```

### 3.1.3 TransformerEncoder

```python
class TransformerEncoder(nn.Module):
    def __init__(self, vocab_size, hidden_size, num_layers, opts):
        ...

    def forward(self, inputs):
        """..."""

        batch_size, seq_len = inputs.size()
        # ------------
        # FILL THIS IN - START
        # ------------
        encoded = self.embedding(inputs)  # @ (batch_size, seq_len, hidden_size)
        for i in range(self.num_layers):
            new_annotations, self_attention_weights = self.self_attentions[i](
                annotations, annotations, annotations
```

```
16              )  # batch_size x seq_len x hidden_size
17              # annotation with residual added.
18              residual_annotations = annotations + new_annotations
19              new_annotations = self.attention_mlps[i](residual_annotations)
20              # Update annotations, the output of this layer.
21              annotations = residual_annotations + new_annotations
22          # ------------
23          # FILL THIS IN - END
24          # ------------
25
26          # Transformer encoder does not have a last hidden layer.
27          return annotations, None
28
29      def create_positional_encodings(self, max_seq_len=1000):
30          ...
```

### 3.1.4 TransformerDecoder

```
1  class TransformerDecoder(nn.Module):
2      def __init__(self, vocab_size, hidden_size, num_layers):
3          ...
4
5      def forward(self, inputs, annotations, hidden_init):
6          """..."""
7
8          batch_size, seq_len = inputs.size()
9          embed = self.embedding(inputs)  # batch_size x seq_len x hidden_size
10
11          # THIS LINE WAS ADDED AS A CORRECTION.
12          embed = embed + self.positional_encodings[:seq_len]
13
14          encoder_attention_weights_list = []
15          self_attention_weights_list = []
16
17          # Decoder: the input fed to the first layer.
18          contexts = embed # batch_size x seq_len x hidden_size
19          for i in range(self.num_layers):
20              # ------------
21              # FILL THIS IN - START
22              # ------------
23              new_contexts, self_attention_weights = self.self_attentions[i](
24                  contexts, contexts, contexts
25              ) # batch_size x seq_len x hidden_size
26              residual_contexts = contexts + new_contexts
27              new_contexts, encoder_attention_weights = self.encoder_attentions[i](
28                  residual_contexts, annotations, annotations
29              ) # batch_size x seq_len x hidden_size
30              residual_contexts = residual_contexts + new_contexts
31              new_contexts = self.attention_mlps[i](residual_contexts)
32              contexts = residual_contexts + new_contexts
33
34              # ------------
35              # FILL THIS IN - END
36              # ------------
37
38              encoder_attention_weights_list.append(encoder_attention_weights)
39              self_attention_weights_list.append(self_attention_weights)
40
41          output = self.out(contexts)
42          encoder_attention_weights = torch.stack(encoder_attention_weights_list)
```
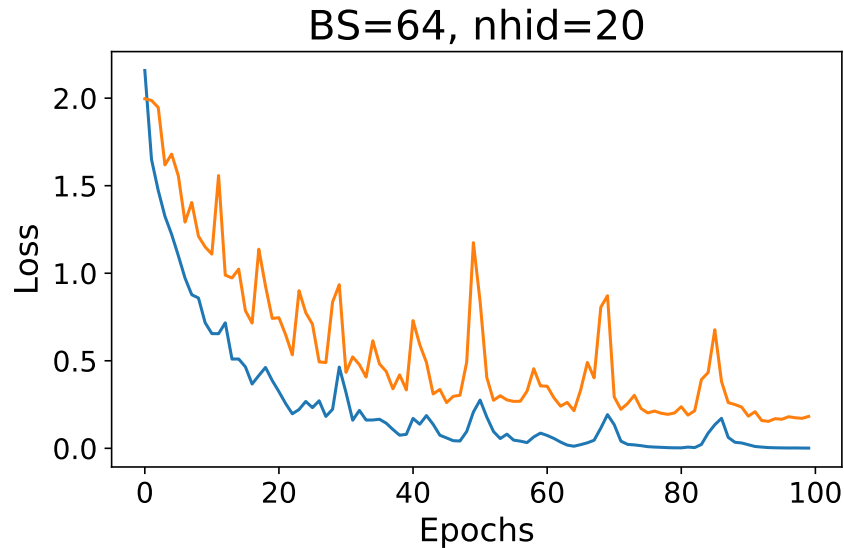
```
43            self_attention_weights = torch.stack(self_attention_weights_list)
44
45            return output, (encoder_attention_weights, self_attention_weights)
46
47        def create_positional_encodings(self, max_seq_len=1000):
48            ...
```

## 3.2 Question 5: Training and Validation Plots



Training logs of the last few steps from three models are reported below, the validation loss of transformer is significantly lower than the previous two decoders. However, the translation results from additive attention is overall better. Additive attention model failed to translate one word (conditioning), but the transformer only translated conditioning correctly, it could be that additive attention works better for short words (because it's recurrent and might suffers from gradient vanishing/exploding problems) but transformer works better for long vocabularies (because it reads the entire sequence the same time).

```
1  =============== GRU ===============
2  Epoch:  95 | Train loss: 0.650 | Val loss: 1.069 | Gen: ethay airway onintoidingsday isway
       oulgefray
3  Epoch:  96 | Train loss: 0.647 | Val loss: 1.048 | Gen: ethay ariway onsidtoingray isway
       oulfrway
4  Epoch:  97 | Train loss: 0.647 | Val loss: 1.120 | Gen: ethay aringpay ondintingshingbay
       isway orkingway
5  Epoch:  98 | Train loss: 0.670 | Val loss: 1.123 | Gen: ethay aringpay onsidtenfay-onsay
       isway orkgingway
6  Epoch:  99 | Train loss: 0.673 | Val loss: 1.053 | Gen: ethay aisray onsiditiongray issway
       oulfreday
7  =============== Additive Attention ===============
8  Epoch:  95 | Train loss: 0.006 | Val loss: 0.138 | Gen: ethay airway onditioningcay isway
       orkingway
9  Epoch:  96 | Train loss: 0.006 | Val loss: 0.133 | Gen: ethay airway onditioningcay isway
       orkingway
10 Epoch:  97 | Train loss: 0.006 | Val loss: 0.137 | Gen: ethay airway onditioningcay isway
       orkingway
11 Epoch:  98 | Train loss: 0.056 | Val loss: 1.192 | Gen: ethay airway ondicecgcay isway
       orkiwway
```

```
12 Epoch:  99 | Train loss: 0.123 | Val loss: 0.332 | Gen: ethay airway onditionwway isway
      orkingway
13 ================ Transformer (Enforcing Causal) ================
14 Epoch:  95 | Train loss: 0.002 | Val loss: 0.166 | Gen: ethhay iarway onditioningcay iseway
      orkingway
15 Epoch:  96 | Train loss: 0.002 | Val loss: 0.180 | Gen: ethhay iirway onditioningcay
      isiiiiiiiiiiisssssacy orkingwaay
16 Epoch:  97 | Train loss: 0.002 | Val loss: 0.175 | Gen: ethhay iirway onditioningcay iswway
      orkingway
17 Epoch:  98 | Train loss: 0.002 | Val loss: 0.171 | Gen: ethhay iirway onditioningcay iswway
      orkingway
18 Epoch:  99 | Train loss: 0.001 | Val loss: 0.182 | Gen: ethhay iirway onditioningcay iswway
      orkingway
```

## 3.3   Question 6: Non-causal Decoder

The outputs from the last few training iterations suggested the modified transformer achieves both lower training and validation loss compared with the original transformer. However, the generated translation is non-sense compared with the transformer with causal decoder. The can be resulted from the fact that, without enforcing causal mask, we allow the model to peak into the future, which discourages the decoder from learning the sequential structure of sentences (i.e., the model failed to learn the importance of character orders).

```
1 ================ Output From Transformer with Causal Decoder ================
2 Epoch:  95 | Train loss: 0.002 | Val loss: 0.166 | Gen: ethhay iarway onditioningcay iseway
      orkingway
3 Epoch:  96 | Train loss: 0.002 | Val loss: 0.180 | Gen: ethhay iirway onditioningcay
      isiiiiiiiiiiisssssacy orkingwaay
4 Epoch:  97 | Train loss: 0.002 | Val loss: 0.175 | Gen: ethhay iirway onditioningcay iswway
      orkingway
5 Epoch:  98 | Train loss: 0.002 | Val loss: 0.171 | Gen: ethhay iirway onditioningcay iswway
      orkingway
6 Epoch:  99 | Train loss: 0.001 | Val loss: 0.182 | Gen: ethhay iirway onditioningcay iswway
      orkingway
7 ================ Output From Transformer with Normal Decoder ================
8 Epoch:  95 | Train loss: 0.000 | Val loss: 0.001 | Gen: - - - - -
9 Epoch:  96 | Train loss: 0.000 | Val loss: 0.001 | Gen: - - - - -
10 Epoch:  97 | Train loss: 0.000 | Val loss: 0.001 | Gen: - - - - -
11 Epoch:  98 | Train loss: 0.000 | Val loss: 0.001 | Gen: - - - - -
12 Epoch:  99 | Train loss: 0.000 | Val loss: 0.001 | Gen: - - - - -
```

## 3.4   Question 7: Advantages and Disadvantages of Additive Attentions and Scaled Dot Product Attention

It seems that the scaled dot attention is better at translating long vocabularies, since a transformer takes the entire sequence of characters once, and can better exploit the correlation between characters distant apart from each other. Additive attention models is based on recurrent neural networks, and RNNs may suffer from vanishing and exploding gradient problems, depends on the specific types of RNN cell used. Therefore, RNN together with additive attentions works better for short vocabularies.

# 4  BERT

## 4.1  Question 1

The `BertCSC413_MLP` class uses 512 hidden neurones (in contrast to the 768 hidden neurones in the original implementation), and a sigmoid activation function (in contrast to the ReLU activation).

## 4.2 Question 2

## 4.3 Question 3

```
[72]  1 what_is("twelve minus fourteen")
```
negative

```
[73]  1 what_is("twelve plus fourteen")
```
positive

```
[74]  1 what_is("eight plus thousand")
```
positive

```
[75]  1 what_is("eight minus thousand")
```
negative

```
[76]  1 what_is("thousand minus eight")
```
positive

```
[77]  1 what_is("eight minus thousand")
```
negative

```
[78]  1 what_is("1 minus 14")
```
negative

```
[79]  1 what_is("1 minus two") # interesting.
```
negative

```
[80]  1 what_is("one minus two")
```
negative

```
[81]  1 what_is("three minus two minus eight")
```
negative

```
[82]  1 what_is("three minus two")
```
positive

```
[83]  1 what_is("one minus one minus one")
```
positive

```
[84]  1 what_is("one minus one minus one plus ten")
```
positive

```
[85]  1 what_is("one minus one plus ten minus one")
```
positive

```
[86]  1 what_is("minus three plus eight")
```
positive

These inference tasks involves both standard usages of binary operator (i.e., number + operator + number) and longer compound usages (i.e., using multiple binary operators consecutively). Moreover, three types of representations of numbers are used: plain English(e.g., three), plain numerical (e.g., 3), and English multipliers (e.g., thousand). Interestingly, the model processes ambiguous compound operations differently, for example "three minus two minus eight" is interpreted as $3 - 2 - 8 < 0$, but "one minus one minus one" as $1 - (1 - 1) > 0$.

## 4.4 Question 4

I changed some hyper-parameters to the training of `model_finetune_bert` by reducing the learning rate while increasing the number of training epochs. Specifically, learning rate is changed to `5e-6` (originally `2e-5`) and the model is now trained for 6 epochs (originally 4 epochs). The (overall) validation accuracy improved from 97% to 98%. The number of incorrect predictions in validation set reduced from two instances to one instance. The detailed training logs and validation reports are attached below:

```
1  ============================================================
2                    Old Hyperparameters
3  ============================================================
4  ======== Epoch 1 / 4 ========
5  Training...
6
7    Average training loss: 1.08
8    Training epcoh took: 0:01:20
9  Running Validation...
10   Accuracy: 0.88
11   Validation took: 0:00:01
12
13 ======== Epoch 2 / 4 ========
14 Training...
15
16   Average training loss: 0.79
17   Training epcoh took: 0:01:19
18 Running Validation...
19   Accuracy: 0.98
20   Validation took: 0:00:01
21
22 ======== Epoch 3 / 4 ========
23 Training...
24
25   Average training loss: 0.59
26   Training epcoh took: 0:01:19
27 Running Validation...
28   Accuracy: 0.97
29   Validation took: 0:00:01
30
31 ======== Epoch 4 / 4 ========
32 Training...
33
34   Average training loss: 0.53
35   Training epcoh took: 0:01:19
36 Running Validation...
37   Accuracy: 0.97
38   Validation took: 0:00:01
39
40 Training complete!
41
42 Predicting labels for 160 test sentences...
```

```
43  Number of expressions with negative result 47
44   45   predicted correctly , accuracy   0.9574468085106383
45
46  Number of expressions with 0 result 2
47   0   predicted correctly , accuracy   0.0
48
49  Number of expressions with positive result 111
50   111   predicted correctly , accuracy   1.0
51
52  =============================================================
53                     New Hyperparameters
54  =============================================================
55  ======== Epoch 1 / 6 ========
56  Training...
57
58    Average training loss: 0.49
59    Training epcoh took: 0:01:27
60  Running Validation...
61    Accuracy: 0.98
62    Validation took: 0:00:01
63
64  ======== Epoch 2 / 6 ========
65  Training...
66
67    Average training loss: 0.45
68    Training epcoh took: 0:01:26
69  Running Validation...
70    Accuracy: 0.98
71    Validation took: 0:00:01
72
73  ======== Epoch 3 / 6 ========
74  Training...
75
76    Average training loss: 0.41
77    Training epcoh took: 0:01:26
78  Running Validation...
79    Accuracy: 0.98
80    Validation took: 0:00:01
81
82  ======== Epoch 4 / 6 ========
83  Training...
84
85    Average training loss: 0.38
86    Training epcoh took: 0:01:28
87  Running Validation...
88    Accuracy: 0.98
89    Validation took: 0:00:01
90
91  ======== Epoch 5 / 6 ========
92  Training...
93
94    Average training loss: 0.37
95    Training epcoh took: 0:01:29
96  Running Validation...
97    Accuracy: 0.98
98    Validation took: 0:00:01
99
100 ======== Epoch 6 / 6 ========
101 Training...
```

```
102
103    Average training loss: 0.36
104    Training epcoh took: 0:01:28
105 Running Validation...
106    Accuracy: 0.98
107    Validation took: 0:00:01
108
109 Training complete!
110
111 Predicting labels for 160 test sentences...
112 Number of expressions with negative result 47
113  47  predicted correctly , accuracy  1.0
114
115 Number of expressions with 0 result 2
116  0   predicted correctly , accuracy  0.0
117
118 Number of expressions with positive result 111
119  110  predicted correctly , accuracy  0.990990990990991
```