



Natural Language Processing CA#3

Student Name:
Pouya Haji Mohammadi Gohari

SID:810102113

Date of deadline
Saturday 11th May, 2024

Dept. of Computer Engineering

University of Tehran

Contents

1	Part 1	3
2	Part 2: LSTM Encoder Model	7
2.1	Implementing	7
2.2	F1-score	10
3	GRU Encoder Model	11
3.1	Implementation	11
3.2	Answering the questions	12
3.2.1	Advantages of LSTM over RNN	12
3.2.2	Differences between LSTM and GRU	13
3.2.3	Why concatenate the verb with words?	13
3.2.4	Vanishing Gradients	14
4	Encoder-Decoder Model	14
4.1	Preprocessing Data for QA	14
4.2	Implementing Model	16
4.3	Validation and F1-score	20
4.4	Answer the Questions	21
4.4.1	Limitation of SRL to QA	21
4.4.2	Why Use BOS and EOS	21
5	Analysis	21
5.1	Quantity	21
5.2	Quality	21

Dataset At first we will investigate through train, valid and test sets. Unfortunately the test set has no "srl_frames". Using "json" library in Python, we load each dataset. The implemented code below shows the procedure:

```
1 with open('data/train.json', 'r') as my_file:
2     train_set = json.load(my_file)
```

At next step we will get a sample from train set in order to be familiar with our dataset.

```
1 keys_list = [key for key in train_set.keys()]
2 for key in keys_list:
3     print(f'The second {key} in traingin set is: {train_set[key][1]}')
```

The second sample for training set is as follow:

- Text: ['The', 'progress', 'of', 'this', 'coordinated', 'offensive', 'was', 'already', 'very', 'entrenched', 'by', 'then', '.']
- SRL¹ Frames: ['O', 'O', 'O', 'O', 'O', 'B-ARG1', 'O', 'O', 'O', 'O', 'O', 'O', 'O']
- Verb Index: 4
- Word Indices: [13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25]

Our datasets is a dictionary consist of a 2D list, where each list represent as a list of words. Another double list "srl_frames" will provide us annotation of each word in corresponding sentence. (e.g., word "offensive" is 'B-ARG1' which is Beginning of Argument 1).

As example above shows, verb of sentence is indexed by 4 which is correspond to 'coordinated' in that sentence. (The word indices is just showing the word index in the whole sentences.)

1 Part 1

As an illustrative example in previous section, dataset is now can be accessed and managed.

We will use a dictionary structure to map the SRL frames to the numerical values. The following code can be used to accomplish this process: (Note that this function will pad first and then mapped them to the numerical values)

```
1 mapper_dict = {
2     'O': 0,
3     'B-ARGO': 1,
4     'I-ARGO': 2,
5     'B-ARG1': 3,
6     'I-ARG1': 4,
7     'B-ARG2': 5,
8     'I-ARG2': 6,
9     'B-ARGM-LOC': 7,
0     'I-ARGM': 8,
1     'B-ARGM-TMP': 9,
2     'I-ARGM-TMP': 10
```

¹Semantic Rule Labels

```

13 }
14 def mapping_dictionary(mapper:dict, srl_labels:list[list[str]]) -> list[list[
    ↳ int]]:
15     """
16         First padding to fix the output length
17         Then simply maps the srl_frames to the tagset
18     """
19     max_length = max((len(labels) for labels in srl_labels), default=0)
20     labels = [[0 for _ in tags] for tags in srl_labels]
21     for index, tag_list in enumerate(srl_labels):
22         if len(tag_list) < max_length:
23             for j in range(max_length-len(tag_list)):
24                 labels[index].append(0)
25             for i, tag in enumerate(tag_list):
26                 labels[index][i] = mapper[tag]
27     return labels

```

Attention: We write an document for each function in order to have more readability and simplify the procedure of using them.

In order to pad the sentences, we will calculate the maximum length and pad all sentences to fix the length of each sentences as follow:

```

1 def Padding(given_set:dict, max_length:int, Pad_token='<PAD>') -> dict:
2     """
3         Adding pad token to the list in order to fix the sizes the sentences
4     """
5     for index, token_list in enumerate(given_set['text']):
6         if len(token_list) < max_length:
7             for j in range(max_length - len(token_list)):
8                 given_set['text'][index].append(Pad_token)
9     return given_set

```

Now it is time to create a class for Vocab:

Adding constructor for vocab class:

```

1 class Vocab(object):
2     def __init__(self, word2id:dict) -> None:
3         """
4             If word2id is None then it would be initialized by a dictionary
5             ↳ otherwise will be initialized by word2id's value
6         """
7         self.word2id = {'<PAD>': 0, '<Start>': 1, '<End>': 2, '<Unknown>': 3} if
            ↳ word2id is None else word2id
8         self.id2word = {value:key for key, value in self.word2id.items()}

```

If word2id is None then vocab will initialize a dictionary with following special tokens:

- '<PAD>': The PAD token will be indexed as 0.

- '<Start>', '<End>': These tokens will be used in the last question of this CA. They stand as Start of sentence and End of sentence.
- '<Unknown>': If any word that is not in vocab then it would be treated as unknown word.

Getting item from Vocab:

```

1  def __getitem__(self, word) -> int:
2      """
3          return the word's key in word2id. If the word does not exist return id
4          ↪ of Unknown word.
5      """
6      if word in self.word2id:
7          return self.word2id[word]
8      return self.word2id['<Unknown>']

```

If word is not in vocab then it would return word's id otherwise it would return '<Unknown>'.

Len method:(magic method)

```

1  def __len__(self) -> int:
2      """
3          Return the vocab size
4      """
5      return len(self.word2id)

```

Add method:

```

1  def add(self, word : str) -> int:
2      """
3          Add a word if it is not already in vocab and return the index of added
4          ↪ word
5      """
6      if word not in self.word2id:
7          length = len(self.word2id)
8          self.word2id[word] = length
9          self.id2word[length] = word
10         return length

```

Word2indices method:

```

1  def word2indices(self, sents:list[list[str]]) -> list[list[int]]:
2      """
3          Get list of sentences and returns the corresponding indices in vocab
4          ↪ dictionary
5      """
6      indices = [[0 for _ in sentence] for sentence in sents]
7      for i,tokens in enumerate(sents):
8          for j,token in enumerate(tokens):
9              indices[i][j] = self[token]
10         return indices

```

Transform into tensor:

```
1 def to_input_tensor(self, sents:list[list[str]]) -> list[list[tensor]]:
2     """
3     First padding the sentences then\\
4     convert sentences to indexes and turn them into tensors
5     """
6     temp_sents = copy.deepcopy(sents)
7     max_length = max((len(sentence) for sentence in sents), default=0)
8     for i,sentence in enumerate(temp_sents):
9         if(len(sentence) < max_length):
10             for _ in range(max_length - len(sentence)):
11                 temp_sents[i].append('<PAD>')
12     return tensor(self.word2indices(temp_sents), dtype=torch.long)
```

From corpus function:

```
1 def from_corpus(corpus: list[list[str]], size:int, remove_frac:float,
2     ↪ freq_cutoff:int) -> Vocab:
3     """
4     Getting an instance of Vocab class and add words by add method
5     """
6     vocab_model = Vocab(None)
7     word_freq = {}
8     for token_list in corpus:
9         for token in token_list:
10             if token in word_freq:
11                 word_freq[token] += 1
12             else:
13                 word_freq[token] = 1
14     sorted_word_freq = sorted(word_freq.items(), key=lambda item: item[1])
15     num_items = len(sorted_word_freq)
16     num_to_remove = int(num_items * remove_frac)
17     remaingin_items = sorted_word_freq[num_to_remove:]
18     remaingin_items = dict(remaingin_items)
19     for word,count in remaingin_items.items():
20         if len(vocab_model) >= size:
21             break
22         if count >= freq_cutoff:
23             vocab_model.add(word)
24     return vocab_model
```

This function first collect all words and their frequencies and then remove a fraction of this words with low frequencies and remaining words will add to the vocab if their frequencies is higher than 'frequency_cutoff'. First we are going to get an instance of this Object and add training words for future purposes.

2 Part 2: LSTM Encoder Model

2.1 Implementing

With the help of 'Pytorch' library, we will create a class of LSTM_Based and inherenct from 'nn.Module' to create our LTSM model as follow:

```
1 class LSTM_Based(nn.Module):
2     def __init__(self, vocab_size, embedding_dim, hidden_dim, output_dim) -> None:
3         super(LSTM_Based, self).__init__()
4         self.word_embeddings = nn.Embedding(vocab_size, embedding_dim)
5         self.lstm = nn.LSTM(embedding_dim, hidden_dim, batch_first=True,
6             ↪ num_layers=1)
7         self.fc = nn.Linear(hidden_dim*2, output_dim)
8         self.output_dim = output_dim
9
10    def forward(self, sentences, predicate_indices):
11        embedded = self.word_embeddings(sentences)
12        lstm_out, (hidden, cell) = self.lstm(embedded)
13        predicate_hidden = lstm_out[torch.arange(len(predicate_indices)),
14            ↪ predicate_indices]
15        predicate_hidden_expanded = predicate_hidden.unsqueeze(1).expand(-1,
16            ↪ lstm_out.size(1), -1)
17        enhanced_hidden = torch.cat((lstm_out, predicate_hidden_expanded), dim=2)
18        output = self.fc(enhanced_hidden)
19        return output
```

After specify the hyper parameters like 'vocab size' and 'embedding dim' and 'output dim' and get an instance of this class. Forward method will simply get sentences and verb indices correspond to that sentence. The shape of sentences would be $[batch_size, seq_len]$ and for verb indices would be $[batch_size]$. After passing the sentences and verb indices to the embedding for each word. After passing sentences through embedding layers the lstm out shape would be $[batch_size, seq_len, embedding_dim]$.

We will get the verb of lstm out and concatenate with all the outputs of lstm. Finally the lstm's outputs will be passed through linear layer to get the label for each word. Thus the shape of the linear output would be $[batch_size, seq_len, 11]$. Before jumping in training phase, we must create a dataset that data loader can interpret with.

Custom Dataset:

```
1 class Custom_Dataset(Dataset):
2     def __init__(self, input_set:dict) -> None:
3         """
4             Specify the dataset as a dictionary
5         """
6         self.txt = input_set['text']
7         self.verb = input_set['verb_index']
8         self.labels = input_set['srl_frames']
9
10    def __len__(self):
11        return len(self.txt)
```

```

2         def __getitem__(self, index):
3             return self.txt[index], self.verb[index], self.labels[index]

```

We will inherit from Dataset class in 'Pytorch' library. For Data loader we will use following implemented code:

```

1 train_dataloader = DataLoader(
2     train_dataset,
3     batch_size=64,
4     shuffle=True,
5     collate_fn=lambda batch : collate_batch(batch, vocab=vocab, mapping_dict=
6         ↪ mapper_dict)

```

The collate function for each batch getting from dataset would be:

```

1 def collate_batch(batch, vocab:Vocab, mapping_dict:dict):
2     text = []
3     verb_index = []
4     labels = []
5     for b in batch:
6         text.append(b[0])
7         verb_index.extend([b[1]])
8         labels.append(b[2])
9     txt_tensor = vocab.to_input_tensor(text)
0     verb_index_tensor = tensor(verb_index, dtype=torch.long)
1     labels_tensor = tensor(mapping_dictionary(mapping_dict, labels), dtype=torch.
2         ↪ long)
3     return txt_tensor, verb_index_tensor, labels_tensor

```

We will get text, verb index, labels from each iteration in batch. From the method '*to_input_tensor*' in Vocab class we will pad them and turn into tensor. Also labels would be padded alongside with text by mapping dictionary implemented in previous section. (Verb indices will be tensor)

Training phase:

```

1 def train(model, optimizer, loss_fn, train_loader, val_loader, epochs, device="
2     ↪ cuda"):
3     accuracy_per_epoch = []
4     loss_per_epoch = []
5     val_loss_per_epoch = []
6     val_accuracy_per_epoch = []
7     reals, all_preds = [], []
8     model.to(device)
9     for epoch in range(epochs):
10         training_loss = 0
11         whole_predictions = []
12         whole_labals = []
13         model.train()
14         for batch in tqdm(train_loader):
15             optimizer.zero_grad()

```



```

5     sentences = batch[0].to(device)
6     verb_indices = batch[1].to(device)
7     srl_labels = batch[2].to(device)
8     preds = model(sentences, verb_indices)
9     loss = loss_fn(preds.view(-1, model.output_dim), srl_labels.view(-1))
10    loss.backward()
11    training_loss += loss.item()
12    optimizer.step()
13    whole_predictions.extend(preds.view(-1, model.output_dim).argmax(axis
    ↪ =1).cpu().numpy().tolist())
14    whole_labals.extend(srl_labels.view(-1).cpu().numpy().tolist())
15    training_loss = training_loss/len(train_loader)
16    train_accuracy = accuracy_score(whole_predictions, whole_labals)
17    loss_per_epoch.append(training_loss)
18    accuracy_per_epoch.append(train_accuracy)
19    print(f'Epoch: {epoch}, training_loss: {training_loss:.2f}, train accuracy
    ↪ {train_accuracy:.2f}')
20    val_loss, val_accuracy, reals, all_preds = validation(model, val_loader,
    ↪ loss_fn, device)
21    val_loss_per_epoch.append(val_loss)
22    val_accuracy_per_epoch.append(val_accuracy)
23    print(f'validation loss: {val_loss:.2f}, validation accuracy: {
    ↪ val_accuracy:.2f}')
24    print(classification_report(reals, all_preds))
25    return loss_per_epoch, accuracy_per_epoch, val_loss_per_epoch,
    ↪ val_accuracy_per_epoch

```

For each epoch, we will append the real outputs to the list of real list. Predicted via model will be appended to predictions list. After seeing all the batches in training data loader, the accuracy will be calculated with accuracy score in 'sklearn' library.

At the end of each epoch, model will validate on the validation data loader. Finally at the end, the classification report will be printed. Validation function will be:

```

1 def validation(model, val_loader, loss_fn, device="cuda"):
2     val_loss = 0
3     model.eval()
4     reals = []
5     all_preds = []
6     for batch in val_loader:
7         sentence = batch[0].to(device)
8         verb_indices = batch[1].to(device)
9         labels = batch[2].to(device)
10        with torch.no_grad():
11            preds = model(sentence, verb_indices)
12            loss = loss_fn(preds.view(-1, model.output_dim), labels.view(-1))
13            val_loss += loss.item()
14            all_preds.extend(preds.view(-1, model.output_dim).argmax(axis=1).cpu().
    ↪ numpy().tolist())

```

```

15     reals.extend(labels.view(-1).cpu().numpy().tolist())
16     val_loss = val_loss / len(val_loader)
17     accuracy = accuracy_score(reals, all_preds)
18     return val_loss, accuracy, reals, all_preds

```

The model's hyper parameters are set as defined in the description but the learning rate is set to $1e - 2$. The loss per epoch for training and validation set alongside with their accuracy is illustrated in figure 1:

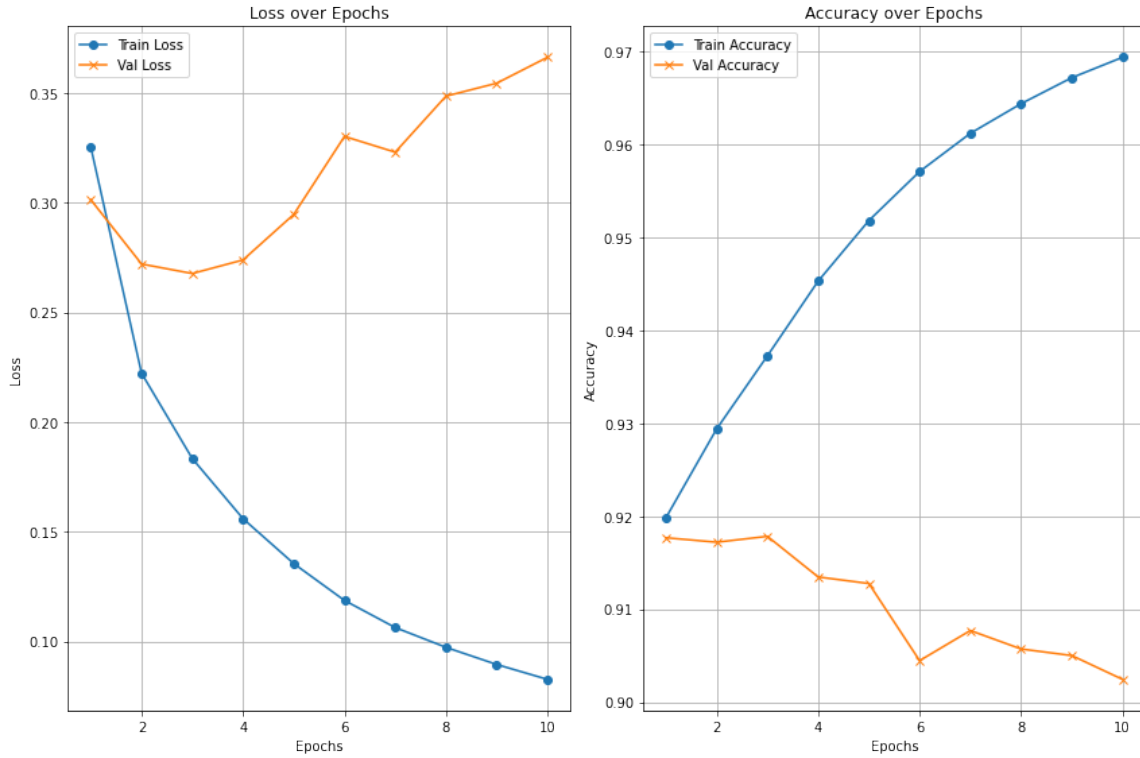


Figure 1: Accuracy and loss for validation and

Interpretation of the loss:

The training loss decreases over the epochs, which indicates that the model is effectively learning from the training data. On the other hand, validation loss shows a different pattern. It initially decreases, indicating that the model is generalizing well from the training data to unseen validation data. But from around epoch 5, the validation loss begins to increase, suggesting the beginning of overfitting.

Accuracy: There is a clear increase in training accuracy, which improves from 0.9 to nearly 0.97 over the epochs. The validation accuracy initially increases along with training accuracy, which is a good indicator of the model's ability to generalize. However, similar to the validation loss pattern, the validation accuracy will drop around epoch 5.

Wrapping this up, Both the training and validation metrics improve, which is typical and expected during the early stages of training. This phase shows that the model is learning useful patterns from the training data and can apply these patterns to validate data. But in later epochs, indicating the overfitting phase.

2.2 F1-score

The f1-score for validation is shown in table 1: Lower macro average for f1-score is indicating that classes are

f1-score	micro	macro
	0.90	0.41

Table 1: F1-score for validation set

imbalance. Classification report shows that the model is likely learning to predict 'o' instead of learning other this will be caused by imbalance classes since the 'o' labels are 34129 out of 37120 labels. That might be due to padding!

In contrast the micro average of 0.9 is average f1-score for all classes.

3 GRU Encoder Model

3.1 Implementation

The procedure for this model is same as previous but instead of using LSTM layer, we will use GRU as follow: (Note that the training and validating is same as LSTM, also dataset and dataloader is same as before).

```

1 class GRU_Baerd(nn.Module):
2     def __init__(self, vocab_size, embedding_dim, hidden_dim, output_dim):
3         super(GRU_Baerd, self).__init__()
4         self.word_embedding = nn.Embedding(vocab_size, embedding_dim)
5         self.gru = nn.GRU(embedding_dim, hidden_dim, num_layers=1, batch_first=
        ↪ True)
6         self.fc = nn.Linear(hidden_dim*2, output_dim)
7         self.output_dim = output_dim
8
9     def forward(self, sentences, predicate_indices):
10        embeded = self.word_embedding(sentences)
11        gru_out, _ = self.gru(embeded)
12        predicate_hidden = gru_out[torch.arange(len(predicate_indices)),
        ↪ predicate_indices]
13        predicate_hidden_expanded = predicate_hidden.unsqueeze(1).expand(-1,
        ↪ gru_out.size(1), -1)
14        enhanced_hidden = torch.cat((gru_out, predicate_hidden_expanded), dim=2)
15        output = self.fc(enhanced_hidden)
16        return output

```

As illustrated in figure 2, the loss and accuracy for epoch is represented in two category training and validation dataset.

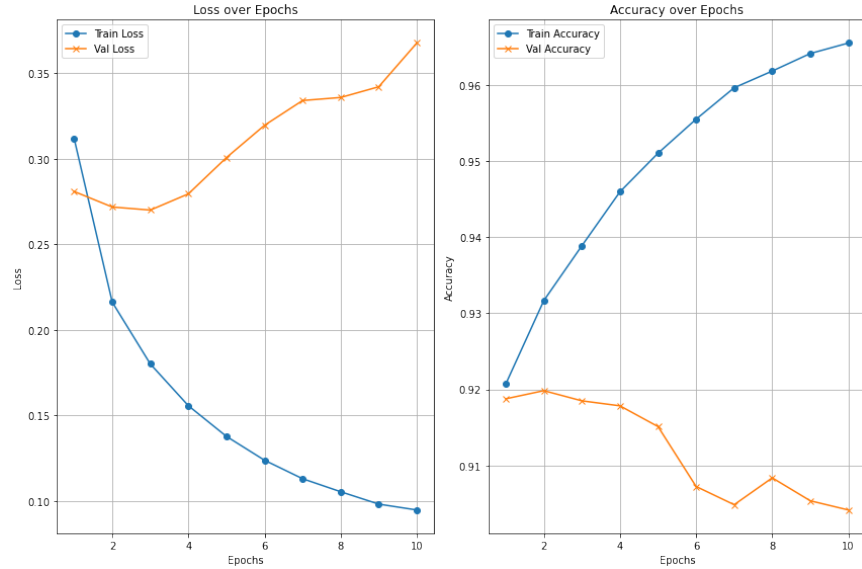


Figure 2: Loss and accuracy for training and validating

Table 2 shows the f1-score metrics for GRU Encoder model. If we compare these two models based on the

f1-score	micro	macro
	0.90	0.40

Table 2: Micro and Macro average of F1-score with Gru

micro and macro average of f1-score metrics interpretation of the results are:

Both models show the same macro F1-score of 0.9. This high score suggests that both models are generally effective across the different label classes when averaged without considering the class imbalance. (Note that macro average treats all classes equally, regardless of how many instances of each class appear in dataset.)

In contrast micro average will treat each class with respect to the classes frequencies. There is slightly difference between their micro F1-score. The slight decrease from LSTM to GRU suggests that the LSTM might be slightly better at correctly identifying true labels across all instances.

Model comparison and interpretation: Both model perform similarly well in terms of their ability to generalize across different SRL labels, as indicated in macro f1-scores. The better micro F1-score for the LSTM model suggests it might be slightly more effective in dealing with class imbalance within dataset. This could be due to the LSTM's more complex gating architecture, which may capture dependencies in the data that the simpler GRU structure slightly misses.

Attention: Personally if i want to chose one, i will prefer the GRU model with respect to computationally efficient that LSTMs due to having fewer parameters.

3.2 Answering the questions

3.2.1 Advantages of LSTM over RNN

The main advantage of LSTM over standard RNN stem from their ability to handle long-term dependencies and avoid problems like vanishing and exploding gradients.

RNN typically struggle to maintain information for long sequences this is due to vanishing gradient problem, where gradients of the loss function become so small that early layers in the sequence stop learning effectively. However LSTM consist gates that regulate the flow of information. These gates (input, output, and forget) can selectively remember or forget information, which allows LSTMs to carry information across many timesteps. This makes them exceptionally good at learning dependencies between events that occur with long gaps in the input data sequences.

3.2.2 Differences between LSTM and GRU

Both are types of RNN architectures used for processing sequential models. Although both are designed to address the problem of vanishing gradients in basic RNNs, they have different structures and operation mechanisms. Aspects of differences are as follow:

1. Architecture: LSTMs have a more complex architecture.They have three gates:

- Input gate decides which values to update.
- Output gate controls what the next hidden state should be.
- Forget gate determines what information should be discarded from the cell state.

Moreover, LSTMs maintain two state vectors;the cell state and the hidden state, which help it to capture long-term dependencies.

2. GRUs simplify the structure seen in LSTMs by using only two gates:

- Update gate acts similarly to the LSTM's input and forget gates combined;it determines both how much the past information(from previous steps) needs to be passed along to the future.
- Reset gate decides how much of the past information to forget.

GRUs combine the cell and hidden state into a single vector, making them less complex often faster to train than LSTMs.

3. Parameters: LSTMs generally have more parameters due to more complex gate mechanisms and state maintenance.However, GRU have fewer parameters due to simplified gate mechanisms and unified state, where can lead to faster training times and require less data to generalize well.

4. Dataset: GRUs need less data to be trained since they have fewer parameters.In contrast LSTMs needs more data to be trained.

3.2.3 Why concatenate the verb with words?

: Intuitively verb has very crucial role in Semantic Role labeling as it is the pivot around which other words in the sentence define their roles.Moreover, each word in the sentence relates to the verb in a way, such as subject(who is performing action), the object(who is receiving action) or other elements like time and location.

Another reason is by concatenating the verb state and output state allows each words representation to be directly informed by the verb's contextual information.This process will ensure verb will influence on each words, leading to more accurate role assignments.

3.2.4 Vanishing Gradients

:

There are several techniques that can be obtained to address this issue. Some of these techniques are:

1. Gradient Clipping: One of the simplest yet effective ways to address both exploding and vanishing gradients is gradient clipping. This method simply involves clipping the gradients during back propagation to ensure they do not exceed a defined threshold.
2. Using Proper Activation Functions: Choosing a good activation function can significantly impact the behaviour of gradients. ReLU or its variants can be helpful for preventing vanishing gradients instead of using sigmoid and tanh.
3. Using dropout can be helpful in RNNs.
4. Instead of padding the sequence we can truncate them in order time steps not getting more.

4 Encoder-Decoder Model

4.1 Preprocessing Data for QA

For this purpose we implement a QA class to make questions and answers in an easy way:

```
1 class QuestionAnswering(object):
2     def __init__(self, dataset:dict, map:dict) -> None:
3         """
4         Pass the Dataset and the vocab model.
5         """
6         self.txt = copy.deepcopy(dataset['text'])
7         self.labels = copy.deepcopy(dataset['srl_frames'])
8         self.verb_indexes = copy.deepcopy(dataset['verb_index'])
9         self.srl_mapping = mapper_dict
10        self.questions = []
11        self.answers = []
12        self.map = map
13        self.label_maps = ['ARGO', 'ARG1', 'ARG2', 'ARGM-LOC', 'ARGM-TMP']
14
15    def mapping_labels(self) -> None:
16        """
17        Map the labels to simplicity
18        """
19        for i, label_list in enumerate(self.labels):
20            for j, label in enumerate(label_list):
21                self.labels[i][j] = self.map[label]
22
23    def make_question_answer(self, index) -> None:
24        sentence = self.txt[index]
25        verb_index = self.verb_indexes[index]
26        labels = self.labels[index]
```

```

27     for tag in self.label_maps:
28         question = [sentence[verb_index], '[SEPT]']
29         question.extend(sentence)
30         answering = []
31         question.extend([tag])
32         if tag in labels:
33             answering.extend(['<Start>'])
34             for i, label in enumerate(labels):
35                 if label == tag:
36                     answering.extend([sentence[i]])
37             answering.extend(['End'])
38         else:
39             answering.extend(['<Start>', '<End>'])
40         self.questions.append(question)
41         self.answering.append(answering)
42
43     def get_dataset(self):
44         """
45         Creating 5 question for each tag with text and verb indexes
46         """
47         for index in tqdm(range(len(self.txt))):
48             self.make_question_answering(index)
49         return self.questions, self.answering

```

Mapping labels method is a simply get the true SRL labels and mapped them with following dictionary:

```

1 mapping = {
2     'O': 'O',
3     'B-ARGO': 'ARGO',
4     'I-ARGO': 'ARGO',
5     'B-ARG1': 'ARG1',
6     'I-ARG1': 'ARG1',
7     'B-ARG2': 'ARG2',
8     'I-ARG2': 'ARG2',
9     'B-ARGM-LOC': 'ARGM-LOC',
10    'I-ARGM': 'ARGM-LOC',
11    'B-ARGM-TMP': 'ARGM-TMP',
12    'I-ARGM-TMP': 'ARGM-TMP'
13 }

```

This method was implemented so we can correctly answer each question with beginning and inside of arguments. Next method called 'make question answering' will get an index of text and for each sentence will make 5 question and 5 answers as description said. Finally with 'get dataset' method we will get all questions and answers. Following example will illustrate that QA class will create questions and answers:

Question = ['inscribed', '[SEPT]', 'A', 'primary', 'stele', ',', 'three', 'secondary', 'steles', ',', 'and', 'two', 'inscribed', 'steles', '.', 'ARG1']

Answer: ['<Start>', 'steles', '<End>']

Now we can treat all questions and answers as list of sentences. Thus we can use vocab class implemented in first

section.

```
1 all_questions_answer = copy.deepcopy(questions)
2 for answer in answerings:
3     all_questions_answer.append(answer)
4 print(len(questions), len(answerings))
5 vocab = from_corpus(all_questions_answer, size=20000, remove_frac=0.3,
    ↪ freq_cutoff=1)
```

Len of vocab will be: 13365 word. Now it is time to get the embedding of each word in vocab if that vocab is in the Glove vectors otherwise we will initialize it with rand. The code below simply doing this: (First loading a 300 dimensional vector as embedding and then check each word in vocab is in the embeddings)

```
1 glove_embedding = {}
2 with open('Glove\glove.6B.300d.txt', 'r') as my_file:
3     for line in my_file:
4         values = line.split()
5         glove_embedding[values[0]] = np.asarray(values[1:], dtype=float)
6 embds = tensor([[0 for _ in range(300)] for __ in range(len(vocab))], dtype=
    ↪ torch.float)
7 for key, value in vocab.word2id.items():
8     if key in glove_embedding:
9         embds[value] = torch.from_numpy(glove_embedding[key]).long()
10    else:
11        embds[value] = torch.rand(embds.shape[1]).reshape((embds.shape[1]))
12 embds[vocab['<PAD>']] = torch.zeros((1, embds.shape[1]))
13 embds[vocab['<Unknown>']] = torch.mean(embds, dim=0, keepdim=True, dtype=float
    ↪ ).squeeze(0)
14 embds.dtype
```

Attention: Embedding for PAD will be assign to zero. Also for Unknown words we will consider the mean of all embedding vectors.

4.2 Implementing Model

We will create 4 classes and inherent each from nn.Module.

Encoder with bidirectional LSTM:

```
1 class Encoder(nn.Module):
2     def __init__(self, embedding_dim, encoder_hidden_dim, pretrained_embeddings):
3         super(Encoder, self).__init__()
4         self.embedding = nn.Embedding.from_pretrained(pretrained_embeddings,
    ↪ freeze=False)
5         self.lstm = nn.LSTM(embedding_dim, encoder_hidden_dim, num_layers=1,
    ↪ batch_first=True, bidirectional=True)
6         self.fc1 = nn.Linear(encoder_hidden_dim*2, encoder_hidden_dim)
7         self.fc2 = nn.Linear(encoder_hidden_dim*2, encoder_hidden_dim)
8
9     def forward(self, src):
```



```

10     embedded = self.embedding(src)
11     encoder_outputs, (hidden, cell) = self.lstm(embedded)
12     hidden = torch.tanh(self.fc1(torch.cat((hidden[0, :, :], hidden[1, :, :]),
13         ↪ dim=1)))
13     cell = torch.tanh(self.fc2(torch.cat((cell[0, :, :], cell[1, :, :]), dim=1))
14         ↪ )
14     return encoder_outputs, hidden.unsqueeze(0), cell.unsqueeze(0)

```

Explain each layers:

- Embedding Layer: we will use the glove embedding vectors as initialize values for embedding vectors. Moreover, freezing argument will be set to 'False' in order that model can update the embedding vectors parameters.
- LSTM: The encoder will get an embedding vector and will return concatenated two vectors obtained by bidirectional lstm.
- Two Linear Layer: For bidirectional LSTMs the hidden and cell shape will be [2, batch, encoder hidden dimension] so we need these layers to transform them into [1, batch, encoder hidden dimension].
- Note that the encoder outputs from LSTM layers would be [batch, sequence length, 2*encoder hidden dimension]

Attention Class:

```

1 class Attention(nn.Module):
2     def __init__(self, encoder_hidden_dim, decoder_hidden_dim):
3         super(Attention, self).__init__()
4         self.attn = nn.Linear(encoder_hidden_dim * 2 + decoder_hidden_dim,
5             ↪ decoder_hidden_dim)
6         self.v = nn.Parameter(torch.rand(decoder_hidden_dim))
7
8     def forward(self, decoder_hidden, encoder_outputs):
9         src_len = encoder_outputs.shape[1]
10        decoder_hidden = decoder_hidden.unsqueeze(1).repeat(1, src_len, 1)
11
12        energy = torch.tanh(self.attn(torch.cat((decoder_hidden, encoder_outputs),
13            ↪ dim=2)))
14        energy = energy.permute(0, 2, 1)
15
16        attn_scores = torch.matmul(self.v, energy)
17        return torch.softmax(attn_scores, dim=1)

```

Explain each layers:

- Linear layers: From previous decoder hidden we will have [batch, decoder hidden dimension] and from encoder outputs we will have [batch, sequence length, 2*encoder hidden dimension]. We will concatenate them and passed through this linear layer so we can transform the [batch, sequence length, 3*decoder hidden dimension] to [batch, sequence length, encoder hidden dimension]. (Note that we will unsqueeze dimension 1 for decoder hidden vector then repeat them due second dimension then concatenate it with encoder outputs)

- Parameters: These parameters will be matmul with energy(explained above) to derive which output of energy we must attend to. Further more, the output of the attention will be [batch, encoder hidden dimension]

Decoder Class:

```

1 class Decoder(nn.Module):
2     def __init__(self, embedding_dim, encoder_hidden_dim, decoder_hidden_dim,
3         ↪ output_dim, pretrained_embedding) -> None:
4         super(Decoder, self).__init__()
5         self.embedding = nn.Embedding.from_pretrained(pretrained_embedding, freeze
6             ↪ =False)
7         self.attention = Attention(encoder_hidden_dim, decoder_hidden_dim)
8         self.lstm = nn.LSTM(encoder_hidden_dim * 2 + embedding_dim,
9             ↪ decoder_hidden_dim, batch_first=True)
10        self.fc_out = nn.Linear(decoder_hidden_dim, output_dim)
11
12    def forward(self, input, hidden, cell, encoder_outputs):
13        input = input.unsqueeze(1)
14        embedded = self.embedding(input)
15
16        attn_weights = self.attention(hidden[-1], encoder_outputs)
17        attn_applied = torch.matmul(attn_weights.unsqueeze(1), encoder_outputs)
18        lstm_inptut = torch.cat((embedded, attn_applied), dim=2)
19        lstm_output, (hidden, cell) = self.lstm(lstm_inptut, (hidden, cell))
20        output = self.fc_out(lstm_output)
21        return output, hidden, cell, attn_weights

```

Explain:

- After getting each word of previous decoder it would be concatenate with attn weights to feed through lstm and finally will passed through the linear function to predict the word in output.

Sequence 2 Sequence Class:

```

1 class Seq2Seq(nn.Module):
2     def __init__(self, encoder, decoder, device) -> None:
3         super(Seq2Seq, self).__init__()
4         self.encoder = encoder
5         self.decoder = decoder
6         self.device = device
7
8     def forward(self, src, trg, teacher_forcing_ratio=0.5):
9         batch_size = src.shape[0]
10        trg_len = trg.shape[1]
11        trg_vocab_size = self.decoder.fc_out.out_features
12        outputs = torch.zeros(batch_size, trg_len, trg_vocab_size).to(self.device)
13        encoder_ouptus, hidden, cell = self.encoder(src)
14        input = trg[:, 0]
15        for t in range(1, trg_len):

```

```

6         output, hidden, cell, _ = self.decoder(input, hidden, cell,
7             ↳ encoder_output)
8         outputs[:, t, :] = output.squeeze(1)
9         teacher_force = np.random.rand() < teacher_forcing_ratio
10        input = trg[:, t] if teacher_force else output.squeeze(1).argmax(1)
11
12    return outputs

```

Explain:

- First questions will be fed to encoder to get the encoder outputs and hidden and cell.
- After that first word of answer will be treated as input of decoder.
- Next step is feed decoder by the encoder outputs and hidden and cell from previous decoder or encoder! (For the first word hidden and cell would be hidden and cell of encoder and for next iterations would be previous decoder's hidden and cell)
- Finally we will teacher force decoder or just simply get the output's decoder.

Before jumping through training and validating we should create Dataset and Data loader for each sets.

```

1 class QAdataset(Dataset):
2     def __init__(self, questions:list[list[str]], answers:list[list[str]]):
3         """
4             Construct with questions and answers
5         """
6         self.questions = questions
7         self.answers = answers
8         print(len(questions))
9         print(len(answers))
10
11    def __len__(self):
12        return len(self.questions)
13
14    def __getitem__(self, index):
15        return self.questions[index], self.answers[index]

```

Collate function and creating dataset and data loader for train and validation sets:

```

1     def collate_batch(batch, vocab):
2         input_ids = []
3         output_ids = []
4         for b in batch:
5             input_ids.append(copy.deepcopy(b[0]))
6             output_ids.append(copy.deepcopy(b[1]))
7         return vocab.to_input_tensor(input_ids), vocab.to_input_tensor(output_ids)
8
9     train_qa_dataset = QAdataset(questions, answerings)
10    train_dataloader = DataLoader(train_qa_dataset, batch_size=16, shuffle=True,
11        ↳ drop_last=True, collate_fn=lambda batch: collate_batch(batch, vocab))

```

```

11 valid_qa_dataset = QAdataset(valid_qustions, valid_answering)
12 valid_dataloader = DataLoader(valid_qa_dataset, batch_size=16, shuffle=True,
    ↪ drop_last=True, collate_fn=lambda batch: collate_batch(batch, vocab))

```

After creating model we should train and validate(Note that the training and validating is just like workshop).After struggling with hyper parameters we choose to set the hyper parameters as in table 3:

batch	encoder dimension	decoder dimension	embedding dimension	epochs
128	128	128	300	20

Table 3: Hyper parameters

The loss and accuracy for each epoch for training and validating has been illustrated in figure 3:

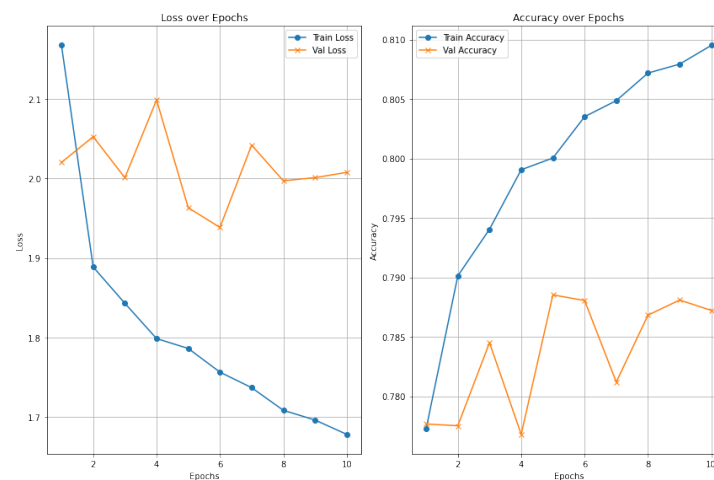


Figure 3: Loss and Accuracy per epoch for training and validation sets

Interpretation of the figure 3: While the model learns the training data effectively, its performance on the validation set raises concerns about its ability to generalize.(That can due to imbalance data especially with pad tokens)

4.3 Validation and F1-score

The F1-score metric is illustrated in table 4:

F1-score	Micro	Macro
	0.00	0.72

Table 4: Macro and Micro f1-score fo sequence 2 sequence

4.4 Answer the Questions

4.4.1 Limitation of SRL to QA

Transforming SRL labels into QA systems using encoder-decoder models has several complexities and limitations. Reason one is for loss of structural information. For instance SRL represent as a structure that understand of sentence semantics, also categorizing different parts of the sentence according to their semantic roles relative to the main verb. When transforming this structural understanding into QA format, there is a risk of losing some of the semantic richness that SRL annotations provide.

This can be challenging to transform SRL labels to QA, since QA systems typically expect input in the form of a natural language question which may not directly map to the structured format used in SRL. Second reason might be effected by generating relevant and meaningful questions for SRL labels involves understanding the context and significant of each role within the broader narrative of the text.

Third reason is might due to higher parameters a encoder-decoder models has. Thus to triaining this models we must have significantly high data while we have 66550 pair of question answering.

4.4.2 Why Use BOS and EOS

In training encoder-decoder models the use of special tokens for the beginning of sequence(BOS in our vocab <Start>) and end of sequence(EOS in our vocab <End>) is crucial.

The BOS token is used to signal the beginning of a sequence. When the decoder part of an encoder-decoder model starts generating an output sequence, the BOS token acts as the initial input. The EOS token, on the other hand, is crucial for indicating when a sequence should end. The EOS token tells the model when to stop generating further tokens, thus defining the length and completeness of the output sequence.

BOS and EOS tokens help standardize sequence lengths by clearly defining starting and stopping points, which are essential for batching and model training where tensors of consistent shapes are required.

Another reason is that the encode-decoder models learns not to just generate content but also where content generation should begin and end. Further, in inference model will continue generation till the EOS token is produced by the decoder!

5 Analysis

5.1 Quantity

For LSTM based we obtain 0.4 micro f1-score. In contrast for sequence to sequence we obtained 0.00! To be certain, we can not compare them with accuracy since the output of encoder-decoder is vocab size that we softmax on them but in LSTM we have 11 labels.

5.2 Quality

We will get two examples from validation since the test set has no srl labels to see the true labels. First two samples are: ["In", "the", "summer", "of", "2005", ",", "a", "picture", "that", "people", "have", "long", "been", "looking", "forward", "to", "started", "emerging", "with", "frequency", "in", "various", "major", "Hong", "Kong", "media", "."] ["The", "most", "important", "thing", "about", "Disney", "is", "that", "it", "is", "a", "global", "brand", "."] Now we are going to predict them with the LSTM therefore the predictions are: ['O', 'O', 'O', 'O', 'O', 'O',

