

# Github Issue Summarization

Gianluca Rea  
m. 278722

January 13, 2023

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>The Dataset</b>	<b>2</b>
<b>3</b>	<b>Preprocessing</b>	<b>3</b>
<b>4</b>	<b>Tokenization</b>	<b>4</b>
4.1	Body Tokenization . . . . .	4
4.2	Title Tokenization . . . . .	5
<b>5</b>	<b>Model</b>	<b>5</b>
<b>6</b>	<b>Diagnostic plots</b>	<b>7</b>
<b>7</b>	<b>Inference</b>	<b>8</b>
<b>8</b>	<b>Conclusions</b>	<b>9</b>

## 1 Introduction

In natural language processing, there are two categories of summarization:

- Extractive Summarization
- Abstractive Summarization

Abstractive Summarization includes heuristic approaches to train the system in making an attempt to understand the whole context and generate a summary based on that understanding. This is a more human-like way of generating summaries, which are more effective than the extractive approaches.

Extractive Summarization essentially involves extracting particular pieces of text (usual sentences) based on predefined weights assigned to the important words where the selection of the text depends on the weights of the words in it. Usually, the default weights are assigned according to the frequency of occurrence of a word. Here, the length of the summary can be manipulated by defining the maximum and minimum number of sentences to be included in the summary.

In our case, we are going to implement an abstractive summarization. The problem concerns the auto-generation of GitHub titles from the summarization of the GitHub issue body.

## 2 The Dataset

For this task, I used the open-source project gharchive.org. The Dataset is called github-issues.



Figure 1: github issue

As we can see in the image above there are two main components of a GitHub issue. The first one is the title highlighted in orange meanwhile in green, we can see the body of the issue.

	issue_title	body
0	can't load the addon. issue to: <a href="https://github.com/zhangyuanwei/node-images/issues/error:/lib64/libc.so.6: version glibc_2.14' not found required by /usr/local/app/taf/fileserver.fileserver/bin/s...">https://github.com/zhangyuanwei/node-images/issues/error:/lib64/libc.so.6: version glibc_2.14' not found required by /usr/local/app/taf/fileserver.fileserver/bin/s...</a>	can't load the addon. issue to: <a href="https://github.com/zhangyuanwei/node-images/issues/error:/lib64/libc.so.6: version glibc_2.14' not found required by /usr/local/app/taf/fileserver.fileserver/bin/s...">https://github.com/zhangyuanwei/node-images/issues/error:/lib64/libc.so.6: version glibc_2.14' not found required by /usr/local/app/taf/fileserver.fileserver/bin/s...</a>
1	hcl accessibility a11yblocking a11ymas mas4.2.10 hcl-makecode win10-edge -title screen reader-help-javascript-call a function narrator focus does not moving to expand side a documentation button a...	user experience: user who depends on screen reader will get confused if narrator focus does not retain on expand side a documentation button after pressing enter on collapse side a documentation b...
2	issue 1265: issue 1264: issue 1261: issue 1260: issue 1257: issue 1256: issue 1253: issue 1252: issue 1250: issue 1247: issue 1246: issue 1243: issue 1242: issue 1239: issue 1237: issue 1236: issu...	! attachments: <a href= <a href="https://github.com&amp;x2f;matiseipl&amp;x2f;czekolada&amp;x2f;issues&amp;x2f;1265">https:&amp;x2f;&amp;x2f;github.com&amp;x2f;matiseipl&amp;x2f;czekolada&amp;x2f;issues&amp;x2f;1265</a> >https:&x2f;&x2f;github.com&x2f;matiseipl&x2f;czekolada&x2f;issues&x2f;1265</a>
df.shape		
(7000, 2)		

Figure 2: github issue dataset head

From figure 2 we can see what part of the dataset we are going to use. After the "issue URL" column drops we are left with two columns "issue title" and "body". Our task will be to summarize the body to have a new text title similar to the issue title. From the entire dataset, we also took only 14000 entries. This choice was forced by the low computational capacity of the computer.

### 3 Preprocessing

For the preprocessing part, a few changes were made to the dataset for cleaning purposes. The first thing done was to eliminate the contraction by substituting them with the formal form of the phrase. We also eliminated special characters and word issues from the various text with the text cleaner function.

```
1 stop_words = set(nltk.corpus.stopwords.words('english'))
2
3 def text_cleaner(text,num):
4     newString = text.lower()
5     newString = BeautifulSoup(newString, "lxml").text
6     newString = re.sub(r'\([^\)]*\)', '', newString)
7     newString = re.sub('"', '', newString)
8     newString = ' '.join([contraction_mapping[t] if t in contraction_mapping else t
9                             for t in newString.split(" ")])
10    newString = re.sub(r"'s\b", "", newString)
11    newString = re.sub("[^a-zA-Z]", " ", newString)
12    newString = re.sub('[m]{2,}', 'mm', newString)
13    newString = re.sub('issue', "", newString)
14    newString = re.sub('issu', "", newString)
15
16    if(num==0):
17        tokens = [w for w in newString.split() if not w in stop_words]
18    else:
19        tokens=newString.split()
20        long_words=[]
21        for i in tokens:
22            if len(i)>1:
23                long_words.append(i)
24    return (" ".join(long_words)).strip()
25    return M
```

and we also deleted the empty rows.

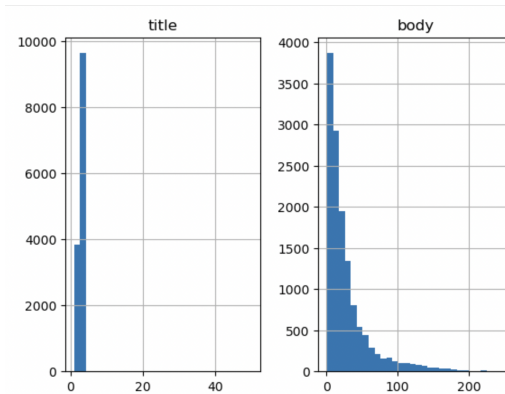


Figure 3: length of text Histogram

From this histogram, we tried to calculate a good number of words for the title and body to use as the limit to choose the rows to use.

We found out that in our case a good compromise was:

- 8 Words for title taking 0.997 percentage of the title group
- 120 Words for body taking 0.967 percentage of the body group

After this, we removed the rows that were not inside the predefined range of words and we applied special tokens to the start and the end of the title for future usage to check if it was not empty.

Finally, we split the title and body into train and test collections with a test size of 0.1 of the entire collection.

## 4 Tokenization

We then proceeded to the tokenization of the text for both title and the body.

### 4.1 Body Tokenization

The code below was used to tokenize the body

```
1 x_tokenizer = Tokenizer()
2 x_tokenizer.fit_on_texts(list(x_tr))
3
4 thresh=3
5 cnt=0 ### Number of rare words -> words appearing less than the thresh
6 tot_cnt=0 ### Vocabulary size
7 freq=0
8 tot_freq=0
9
10 for key,value in x_tokenizer.word_counts.items():
11     tot_cnt=tot_cnt+1
12     tot_freq=tot_freq+value
13     if(value<thresh):
14         cnt=cnt+1
15         freq=freq+value
16
17 #prepare a tokenizer for reviews on training data
18 x_tokenizer = Tokenizer(num_words=tot_cnt-cnt)
19 x_tokenizer.fit_on_texts(list(x_tr))
20
21 #convert text sequences into integer sequences
22 x_tr_seq = x_tokenizer.texts_to_sequences(x_tr)
23 x_val_seq = x_tokenizer.texts_to_sequences(x_val)
24
25 #padding zero upto maximum length
26 x_tr = pad_sequences(x_tr_seq, maxlen=max_body_len, padding='post')
27 x_val = pad_sequences(x_val_seq, maxlen=max_body_len, padding='post')
28
29 #size of vocabulary ( +1 for padding token)
30 x_voc = x_tokenizer.num_words + 1
```

As a result of the code, we found out that the percentage of rare words, defined as the word which has appeared less than the thresh in this case 3, was 69.342. The total coverage of rare words was 9.503. In the end, we have collected a vocabulary with a size equal to 9927

## 4.2 Title Tokenization

The code below was used to tokenize the title

```
1 #prepare a tokenizer for reviews on training data
2 y_tokenizer = Tokenizer()
3 y_tokenizer.fit_on_texts(list(y_tr))
4
5 thresh=6
6 cnt=0
7 tot_cnt=0
8 freq=0
9 tot_freq=0
10
11 for key,value in y_tokenizer.word_counts.items():
12     tot_cnt=tot_cnt+1
13     tot_freq=tot_freq+value
14     if(value<thresh):
15         cnt=cnt+1
16         freq=freq+value
17
18 #prepare a tokenizer for reviews on training data
19 y_tokenizer = Tokenizer(num_words=tot_cnt-cnt)
20 y_tokenizer.fit_on_texts(list(y_tr))
21
22 #convert text sequences into integer sequences
23 y_tr_seq = y_tokenizer.texts_to_sequences(y_tr)
24 y_val_seq = y_tokenizer.texts_to_sequences(y_val)
25
26 #padding zero upto maximum length
27 y_tr = pad_sequences(y_tr_seq, maxlen=max_title_len, padding='post')
28 y_val = pad_sequences(y_val_seq, maxlen=max_title_len, padding='post')
29
30 #size of vocabulary
31 y_voc = y_tokenizer.num_words +1
32
33 y_tokenizer.word_counts['sostok'],len(y_tr)
```

As a result of the code, we found out that the percentage of rare words, defined as the word which has appeared less than the thresh in this case 6, was 86.09. The total coverage of rare words was 14.012. In the end, we collected a vocabulary with a size equal to .

After this, we deleted all the rows in which the body or title was formed only by the start and end tokens.

## 5 Model

We are finally at the model-building part. But before we do that, we need to familiarize ourselves with a few terms which are required prior to building the model.

- Return Sequences = True: When the return sequences parameter is set to True, LSTM produces the hidden state and cell state for every timestep
- Return State = True: When return state = True, LSTM produces the hidden state and cell state of the last timestep only
- Initial State: This is used to initialize the internal states of the LSTM for the first timestep
- Stacked LSTM: Stacked LSTM has multiple layers of LSTM stacked on top of each other. This leads to a better representation of the sequence.

Here, we are building a 3-stacked LSTM for the encoder:

```
1 K.clear_session()
2
3 latent_dim = 300
4 embedding_dim=100
5
6 # Encoder
7 encoder_inputs = Input(shape=(max_body_len,))
8
9 #embedding layer
10 enc_emb = Embedding(x_voc, embedding_dim,trainable=True)(encoder_inputs)
11
12 #encoder lstm 1
13 encoder_lstm1 = LSTM(latent_dim,return_sequences=True,return_state=True,dropout=0.4,
14 recurrent_dropout=0.4)
15 encoder_output1, state_h1, state_c1 = encoder_lstm1(enc_emb)
16
17 #encoder lstm 2
18 encoder_lstm2 = LSTM(latent_dim,return_sequences=True,return_state=True,
19 dropout=0.4,recurrent_dropout=0.4)
20 encoder_output2, state_h2, state_c2 = encoder_lstm2(encoder_output1)
21
22 #encoder lstm 3
23 encoder_lstm3=LSTM(latent_dim, return_state=True,
24 return_sequences=True,dropout=0.4,recurrent_dropout=0.4)
25 encoder_outputs, state_h, state_c= encoder_lstm3(encoder_output2)
26
27 # Set up the decoder, using 'encoder_states' as initial state.
28 decoder_inputs = Input(shape=(None,))
29
30 #embedding layer
31 dec_emb_layer = Embedding(y_voc, embedding_dim,trainable=True)
32 dec_emb = dec_emb_layer(decoder_inputs)
33
34 decoder_lstm = LSTM(latent_dim, return_sequences=True,
35 return_state=True,dropout=0.4,recurrent_dropout=0.2)
36 decoder_outputs,decoder_fwd_state, decoder_back_state =
37 decoder_lstm(dec_emb,initial_state=[state_h, state_c])
38
39 # Attention layer
40 attn_layer = AttentionLayer(name='attention_layer')
41 attn_out, attn_states = attn_layer([encoder_outputs, decoder_outputs])
42
43 # Concat attention input and decoder LSTM output
44 decoder_concat_input = Concatenate(axis=-1, name='concat_layer')
45 ([decoder_outputs, attn_out])
46
47 #dense layer
48 decoder_dense = TimeDistributed(Dense(y_voc, activation='softmax'))
49 decoder_outputs = decoder_dense(decoder_concat_input)
50
51 # Define the model
52 model = Model([encoder_inputs, decoder_inputs], decoder_outputs)
53
54 model.summary()
```

After this, we set up the batch size to 128 and we set the number of epoch to 8.

```
1 history=model.fit([x_tr,y_tr[:, :-1]],  
2 y_tr.reshape(y_tr.shape[0],y_tr.shape[1], 1)[: ,1:] ,  
3 epochs=8,callbacks=[es],batch_size=128,  
4 validation_data=([x_val,y_val[:, :-1]],  
5 y_val.reshape(y_val.shape[0],y_val.shape[1], 1)[: ,1:] ))
```

with the following results:

Epoch	Loss	Val Loss
1	2.3261	2.0039
2	1.9497	1.9345
3	1.8714	1.8975
4	1.8162	1.9302
5	1.7769	1.8365
6	1.7367	1.8292
7	1.6969	1.7967
8	88	788

## 6 Diagnostic plots

We analyzed the diagnostic plots to understand the behavior of the model over time and we can infer that validation loss has increased after epoch 2 for 6 successive epochs. Although, it could be useful to train the model for other epochs to be sure about the increase would be maintained over time.

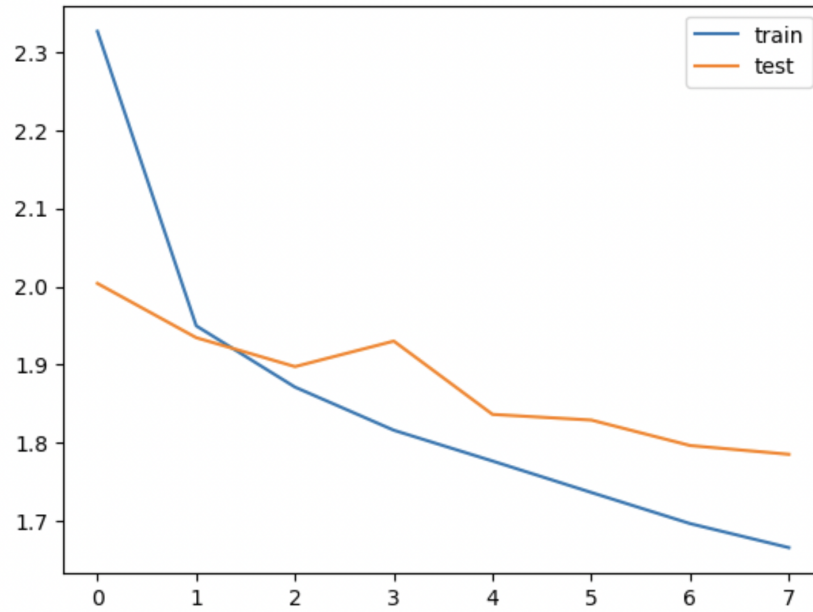


Figure 4: Diagnostic plots

## 7 Inference

We set up the inference of the encoder and decoder:

```
1 # Encode the input sequence to get the feature vector
2 encoder_model = Model(inputs=encoder_inputs, outputs=[encoder_outputs, state_h, state_c
3 ])
4 # Decoder setup
5 # Below tensors will hold the states of the previous time step
6 decoder_state_input_h = Input(shape=(latent_dim,))
7 decoder_state_input_c = Input(shape=(latent_dim,))
8 decoder_hidden_state_input = Input(shape=(max_body_len, latent_dim))
9
10 # Get the embeddings of the decoder sequence
11 dec_emb2= dec_emb_layer(decoder_inputs)
12 # To predict the next word in the sequence, set the initial states to the states from
13 # the previous time step
14 decoder_outputs2, state_h2, state_c2 = decoder_lstm(dec_emb2, initial_state=[
15     decoder_state_input_h, decoder_state_input_c])
16
17 #attention inference
18 attn_out_inf, attn_states_inf = attn_layer([decoder_hidden_state_input,
19     decoder_outputs2])
20 decoder_inf_concat = Concatenate(axis=-1, name='concat')([decoder_outputs2,
21     attn_out_inf])
22
23 # A dense softmax layer to generate prob dist. over the target vocabulary
24 decoder_outputs2 = decoder_dense(decoder_inf_concat)
25
26 # Final decoder model
27 decoder_model = Model(
28     [decoder_inputs] + [decoder_hidden_state_input, decoder_state_input_h,
29     decoder_state_input_c],
30     [decoder_outputs2] + [state_h2, state_c2])
```

We are defining a function below which is the implementation of the inference process:

```
1 def decode_sequence(input_seq):
2     # Encode the input as state vectors.
3     e_out, e_h, e_c = encoder_model.predict(input_seq)
4     # Generate an empty target sequence of length 1.
5     target_seq = np.zeros((1,1))
6     # Populate the first word of the target sequence with the start word.
7     target_seq[0, 0] = target_word_index['sostok']
8     stop_condition = False
9     decoded_sentence = ''
10    while not stop_condition:
11        output_tokens, h, c = decoder_model.predict([target_seq] + [e_out, e_h, e_c])
12        # Sample a token
13        sampled_token_index = np.argmax(output_tokens[0, -1, :])
14        sampled_token = reverse_target_word_index[sampled_token_index]
15        if(sampled_token!='eostok'):
16            decoded_sentence += ' '+sampled_token
17        # Exit condition: either hit max length or find stop word.
18        if (sampled_token == 'eostok' or len(decoded_sentence.split()) >= (
19            max_summary_len-1)):
20            stop_condition = True
21        # Update the target sequence (of length 1).
22        target_seq = np.zeros((1,1))
23        target_seq[0, 0] = sampled_token_index
24        # Update internal states
25        e_h, e_c = h, c
26    return decoded_sentence
```



## 8 Conclusions

Below few good predictions are displayed to show the correct work of the tool.

```
Body: https github com prj rev bwfs dasmoto blob master dasmoto art resource css style css e
xcellent implementation dry principles setting font family elements one selector use set ent
ire page
Original title: great dry css
1/1 [=====] - 2s 2s/step
1/1 [=====] - 0s 17ms/step
1/1 [=====] - 0s 21ms/step
1/1 [=====] - 0s 15ms/step
1/1 [=====] - 0s 15ms/step
Predicted title: great dry css
```

Figure 5: Good Result 1

```
Body: create organization test https try gitea io supertest new migration organization test
https try gitea io supertest installer framework see watchers test https try gitea io supert
est installer framework watchers look picture please https user images githubusercontent com
aaf dbd png
Original title: repo bug
1/1 [=====] - 2s 2s/step
1/1 [=====] - 0s 16ms/step
1/1 [=====] - 0s 17ms/step
1/1 [=====] - 0s 16ms/step
Predicted title: found bug
```

Figure 6: Good Result 2

Following is an example of a bad prediction:

```
Body: make sure site deploys appropriate information db
Original title: migrate the db
1/1 [=====] - 2s 2s/step
1/1 [=====] - 0s 15ms/step
1/1 [=====] - 0s 15ms/step
1/1 [=====] - 0s 15ms/step
Predicted title: add page
```

Figure 7: A bad result

We can say that deep learning could be improved by taking a bigger size of the dataset to work with. Also increasing the number of epochs will also result in better predictions.