

Q1 : Data processing

使用 spacy 去做斷詞，text 的長度上限為 300 個 tokens，summary 的長度上限為 80 個 tokens，在此將得到的 tokens 以及 special_token 利用 glove.840B.300d.txt (pretrained embedding) 產生相對應的 embedding matrix，而此 matrix 將在模型的 embedding layer 使用。

Q2 : Extractive summarization model

(a) RNNModel : 一個 epoch 處理 batch size 篇文章

1. $Embed_w_t = \text{Dropout}(\text{Embedding Layer}(w_t))$ ， w_t 及 $Embed_w_t$ 是第 t 個時間點的 token 及詞向量
2. $y_t = \text{pack_padded_sequence}(Embed_w_t, len_text_t)$ ， len_text_t 為相對應的文章的長度
3. $output_t, h_t, c_t = \text{LSTM}(y_t, \text{None})$ ， h_t 及 c_t 為第 t 個時間點的 hidden 及 cell，而此 LSTM 為兩層且雙向
4. $output_t, len_text_t = \text{pad_packed_sequence}(output_t)$
5. $pred_t = \text{sigmoid}(\text{Linear}(\text{Dropout}(output_t)))$ ， $pred_t$ 為第 t 個時間點的二分類的結果。

(b) 用 validation data 去做 evaluation，根據 rouge 計算：

rouge-1 : 19%

rouge-2 : 2.76%

rouge-L : 13.7%

(c) 採用 BCEWithLogitsLoss 當作 Loss Function

$\text{Loss}(\text{pred}, \text{target}) = \text{mean}\{l_0, l_1, l_2, \dots, l_{n-1}\}$

$l_m = -(\text{pred}_m * \log(\delta(\text{target}_m)) + (1 - \text{pred}_m) * \log(1 - \delta(\text{target}_m)))$ ，

n 為資料數量， pred_m 為第 m 筆資料的預測， target_m 為第 m 筆資料的標籤， δ 為 sigmoid。

(d) 採用 Adam 當作 optimization algorithm，learning rate 為 0.001 且 batch size 為 128。

(e) $\text{pred}'s \text{ shape} = [\text{batch size}, \text{the number of the tokens of the text}]$ ，各文章的 token 包含數量最多的句子即為 summary。

Q3 : Seq2seq + Attention model

(a)

● Encoder : 一個 epoch 處理 batch size 篇文章

1. $Embed_w_t = \text{Embedding Layer}(w_t)$ ， w_t 及 $Embed_w_t$ 是第 t 個時間點的 token 及詞向量
2. $y_t = \text{pack_padded_sequence}(Embed_w_t, len_text_t)$ ， len_text_t 為相對應文章的長度

3. $output_t, h_t = \text{GRU}(y_t, \text{None})$ ， $output_t$ 為第 t 個時間點的 context vector，而此 GRU 為一層且雙向
4. $output_t, len_text_t = \text{pad_packed_sequence}(output_t)$
5. $final_h_t = \text{Tanh}(\text{Linear}(hidden_t))$ ， $hidden_t$ 為第 t 個時間點最後的 forwards 和 backwards 的 hidden 所組成， $final_h_t$ 為 deocder 的 input hidden。

● Attention：

$attention_t = \text{softmax}(\text{Linear}(h_t \text{ 和 } c_t \text{ 接在一起}))$ ， h_t 為 encoder 的輸出 hidden，也就是上述的 $final_h_t$ ， c_t 為 encoder 的 context vector，也就是上述的 $output_t$ ，其中利用 mask_fill 在 mask 為 1 時用 value 做填充

● Decoder：對於一篇文章的摘要生成是一次一個 token

1. $\text{Embed_target}_t = \text{Embedding Layer}(target_t)$ ， $target_t$ 及 Embed_target_t 是第 t 個時間點的 target 及其詞向量
2. $weight_t = attention_t * output_t$ (Batch Matrix Multiply, bmm)， $output_t$ 為 encoder 的輸出， $attention_t$ 為相對 $output_t$ 應該放多少的注意力， $attention_t$ 的數值介於 0 到 1
3. $decoder_input_t = \text{Embed_target}_t \text{ 和 } weight_t \text{ 接在一起}$
4. $pred_t, h_t = \text{GRU}(decoder_input_t, h_{t-1})$ ， h_t 為第 t 個時間點的 hidden， $pred_t$ 經過 linear 處理後會存入 prediction，而此 GRU 為一層且單向
5. $prediction = [\text{length of target, batch size, vocab size}]$ ，經由 softmax 並挑選第三維機率最大的 token 組成 summary

(b) 用 validation data 去做 evaluation，根據 rouge 計算：

rouge-1：22.5%

rouge-2：5.9%

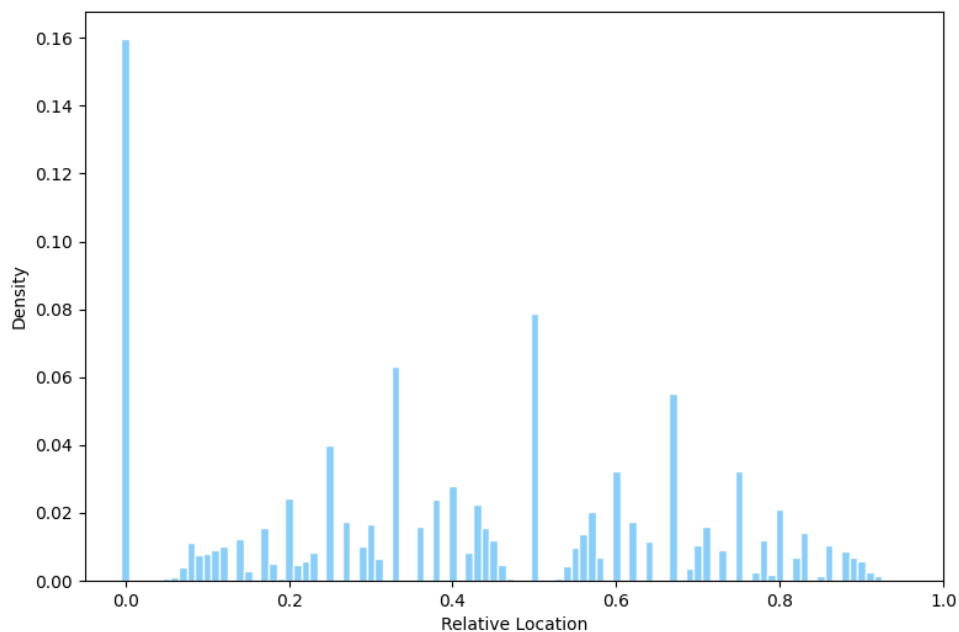
rouge-L：19.4%

(c) 採用 CrossEntropyLoss 當作 loss function

$$\text{Loss}(\text{pred}, \text{target}) = -\log\left(\frac{\exp(\text{pred}[\text{target}])}{\sum_j \exp(\text{pred}[j])}\right)$$

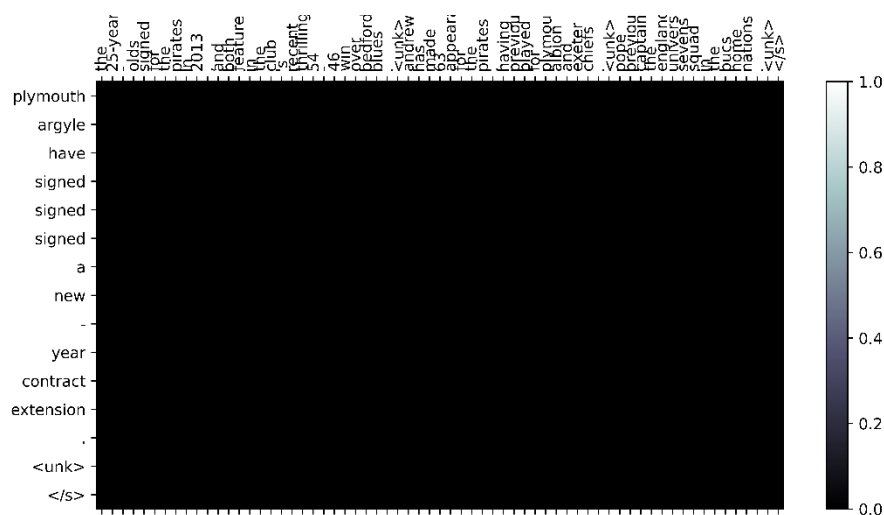
(d) 採用 Adam 當作 optimization alogorithm，learning rate 為 0.001 且 batch size 為 32。

Q4 : The distribution of relative locations



根據上面這張圖可以發現相對位置為 0 的密度最高，這也符合了一般人對於英文文章的認知，topic sentence 通常都是文章的第一句，適合做 summary。

Q5 : Visualize the attention weights



由於 attention + seq2seq 的模型訓練上沒有達到預期效果，所以畫出的 attention weight 呈現 0 的狀態，但在這裡要多加說明，此圖的 x 軸為 input tokens，y 軸為 output tokens，圖中的顏色代表 weight，顏色越白代表預測此字時會越專注於輸入的 input token。

Q6：Rouge-L

Rouge-L 使用了最長共同子序列(LCS)。他的計算方式為下列：

$$R_{lcs} = \frac{LCS(X, Y)}{m} \quad (1)$$

$$P_{lcs} = \frac{LCS(X, Y)}{n} \quad (2)$$

$$F_{lcs} = \frac{(1 + \beta^2)R_{lcs}P_{lcs}}{R_{lcs} + \beta^2 P_{lcs}} \quad (3)$$

$X = [x_1, x_2, \dots, x_m]$ ，是參考摘要(Reference Summaries)的子序列

$Y = [y_1, y_2, \dots, y_n]$ ，是自動摘要 (Predict Summaries)的子序列

m 為參考摘要的長度

n 為自動摘要的長度

數學式(1)可以說是 Recall，數學式(2)可以說是 Precision，數學式(3)就是所謂的 Rouge-L，而式子中的 β 在 DUC 中通常都是一個很大的常數，所以 Rouge-L 通常都只考慮 R_{lcs} 。