

Q1 Model

1. Model

T5 為一 encoder-decoder 架構，透過給定當前階段之前產生的文字以及額外的 input X 進行下一個字詞的預測，在作業二 summarization 的任務，X 即為新聞的內文。

2. Preprocessing

優先刪除訓練資料內 date_publish, source_domain, split 等資料，並將 jsonl file 轉為 json file。使用 t5 模型提供之 tokenizer 進行 tokenize。

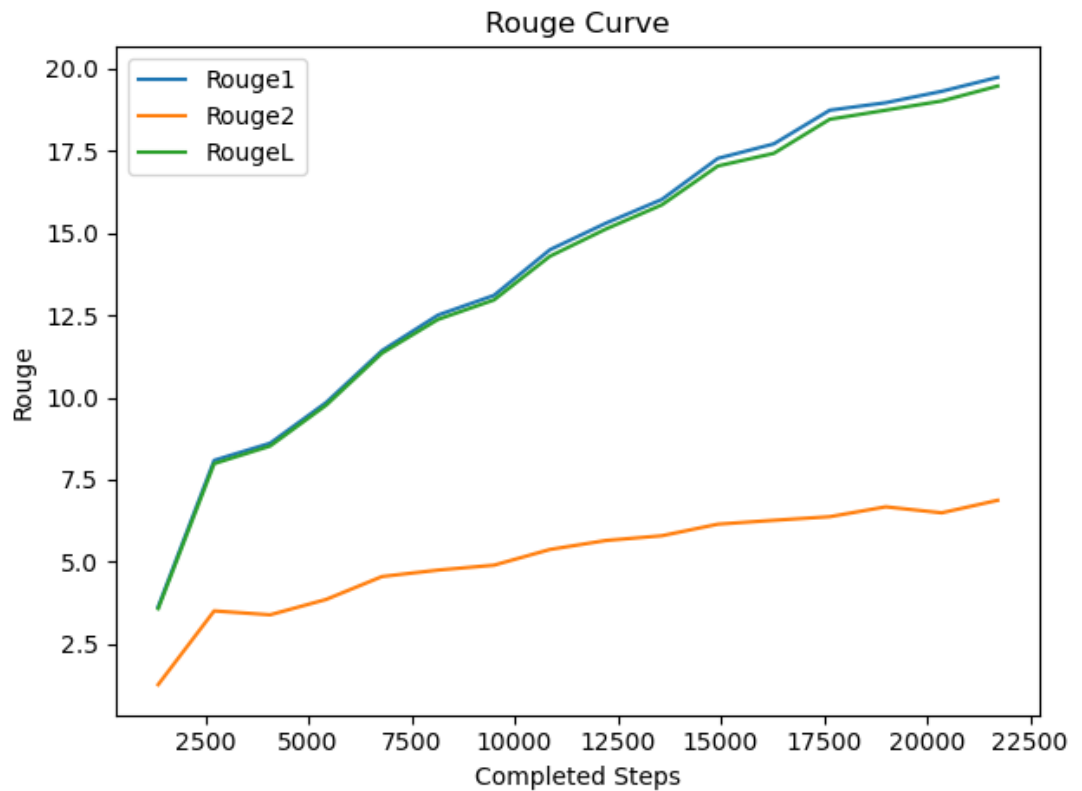
Q2 Training

1. Hyperparameter

```
--train_file data/train_all.json \  
--validation_file data/reference.json \  
--max_source_length 2048 \  
--max_target_length 64 \  
--pad_to_max_length \  
--model_name_or_path google/mt5-small \  
--text_column maintext \  
--summary_column title \  
--per_device_train_batch_size 4 \  
--gradient_accumulation_steps 4 \  
--num_warmup_steps 500 \  
--output_dir HW2_training_result \  
--seed 228 \  
--num_beams 4 \  
--source_prefix summarize: \  
--preprocessing_num_workers 1 \  
--num_train_epochs 16
```

觀察 training 和 testing 資料，內文大多 1000 多字，而標題大多 30-40 字，故選擇 max length 為 2048 及 64。訓練前期容易不穩定，故設定 warmup steps 為 500 輔助模型訓練。

2. Learning Curve



Q3 Generation Strategies

1. Strategies

i. Greedy

每一步都找機率最大的字詞。

ii. Beam Search

持續追蹤機率前幾大的路線，並從中選出最好的。

iii. Top-k Sampling

每一步只保留機率最大的 k 個字可以選擇。較大的 k 使模型有較多元的字詞可以選擇，較小則限制模型的字詞選擇。

iv. Top-p Sampling

每一步只保留機率累積至 p 且機率前幾高的字可以選擇。較大的 p 使模型有較多元的字詞可以選擇，然而當少數字詞的機率較大（較確定）時，模型的選擇也較少。

v. Temperature

$$P(w) = \frac{\exp(s_w / \tau)}{\sum_{w' \in V} \exp(s_{w'} / \tau)}$$

在計算機率的過程中除上 temperature，較高的 temperature 使機率分佈平均，模型產出較多元的字詞。較低的 temperature 使機率分佈極端，模型產出較無變化的字詞。

2. Hyperparameter

i. Greedy

Rouge-1 = 24.33, Rouge-2 = 9.12, Rouge-l = 21.81

ii. Beam Search

num_beams = 4:

Rouge-1 = 25.88, Rouge-2 = 10.39, Rouge-l = 23.19

num_beams = 10:

Rouge-1 = 25.75, Rouge-2 = 10.41, Rouge-l = 23.00

iii. Top-k Sampling

k = 5:

Rouge-1 = 22.94, Rouge-2 = 7.89, Rouge-l = 20.17

k = 20:

Rouge-1 = 20.77, Rouge-2 = 6.71, Rouge-l = 18.33

iv. Top-p Sampling

p = 0.8:

Rouge-1 = 21.28, Rouge-2 = 7.41, Rouge-l = 18.95

p = 0.5:

Rouge-1 = 23.37 Rouge-2 = 8.53, Rouge-l = 20.85

v. Temperature

temperature = 2:

Rouge-1 = 10.11 Rouge-2 = 1.66, Rouge-l = 8.74

temperature = 0.5:

Rouge-1 = 23.37 Rouge-2 = 8.41, Rouge-l = 20.79

最終採用 beam search 策略，beam 數量為 4。