# Q1 Data processing :

1. Tokenizer

助教給的 train.json 檔如下圖:



我將所有的 paragraphs 接在一起變成一整大篇的文章，並且重新計算 answer
start 的位置，也因為這樣做的關係我避免了做分類的任務(判斷 relevant)， 當
初會想這樣做是因為覺得如果把功課分成兩個任務去做，這樣整體準確率不就
會變成 accuracy of stage 1 ＊ accuracy of stage 2 ，這樣可能導致我的準確率
較低(但結果並不是我想的那樣)， 而雖然我們文章長度變為原本的 6~7 倍，但
是因為 bert model 有使用 doc_stride 因此應該不會降低模型準確率太多，於是
就這樣做下去了。

而在丟資料進 bert model 的方式如下圖顯示



也就是 [CLS] question [SEP] paragraph context [SEP]，每個字會去 bert 的 vocab.txt 檔(如下圖) 找尋對應的 id ，[CLS] 起頭符號、[SEP] 分割符號以及 [UNK] 未知詞。



由於使用 doc_stride 的關係，且 bert input size 最大只有 512 (我最大設 200 個字，因為要在 8G 的空間上 train RoBERTa large)故需要把文章拆成許多小文章，這樣就必須要做(token to original map)的動作，input id 為把 input 每個字對應到 vocab.txt 的索引，segment ids 就是把 question 用 0 表示 ，paragraph context 用 1 表示，下圖可見:



而標準的 bert input 有三層:

Token Embeddings、Segment Embeddings、Position Embeddings



2. Answer Span

(a)How did you convert the answer span start/end position on characters to position on tokens after BERT tokenization?

```
unique_id: 1000000000
example_index: 0
doc_span_index: 16
tokens: [CLS] 舍 本 和 諾 的 數 據 能 推 算 出 蓮 星 的 恆 星 的 質 量 ？ [SEP] 千 億 個 。 2010 年 對 恆 星 數 量 的 估 計 是 在 可 觀 測 宇 宙 中 有 3000 [UNK] 顆 。 儘 管 人 們 往 往 認 為 恆 星 僅 存 在 於
token_to_orig_map: 22:800 23:801 24:802 25:803 26:804 27:805 28:806 29:807 30:808 31:809 32:810 33:811 34:812 35:813 36:814 37:815 38:816 39:817 40:818 41:819 42:820 43:821 44:822 45:823
token_is_max_context: 22:False 23:False 24:False 25:False 26:False 27:False 28:False 29:False 30:False 31:False 32:False 33:False 34:False 35:False 36:False 37:False 38:False 39:False 40
input_ids: 101 5650 3315 1469 6306 4638 3149 3087 5543 2972 5050 1139 6865 3215 4638 2604 3215 4638 6549 7030 8043 102 1283 1023 943 511 8166 2399 2205 2604 3215 3149 7030 4638 844 6243
input_mask: 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
segment_ids: 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
start_position: 181
end_position: 184
answer: 斯 特 魯 維
*** Example ***
unique_id: 1000000001
example_index: 0
doc_span_index: 17
tokens: [CLS] 舍 本 和 諾 的 數 據 能 推 算 出 蓮 星 的 恆 星 的 質 量 ？ [SEP] 恆 星 已 經 被 發 現 了 。 在 19 世 紀 雙 星 觀 測 所 獲 得 的 成 就 使 重 要 性 也 加 入 了 。 在 183 ##4 年 ， 白 塞 爾 觀 測
token_to_orig_map: 22:850 23:851 24:852 25:853 26:854 27:855 28:856 29:857 30:858 31:859 32:860 33:861 34:862 35:863 36:864 37:865 38:866 39:867 40:868 41:869 42:870 43:871 44:872 45:873
token_is_max_context: 22:False 23:False 24:False 25:False 26:False 27:False 28:False 29:False 30:False 31:False 32:False 33:False 34:False 35:False 36:False 37:False 38:False 39:False 40
input_ids: 101 5650 3315 1469 6306 4638 3149 3087 5543 2972 5050 1139 6865 3215 4638 2604 3215 4638 6549 7030 8043 102 2604 3215 2347 5195 6158 4634 4412 749 511 1762 8131 686 5145 7427
input_mask: 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
segment_ids: 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
start_position: 131
end_position: 134
```

Answer 並不會出現在所有文章中(因為我將所有文章合成一個大文章)，所以在做 training data 的時候，我只會將滑窗有節錄到答案的文章段落做成 training data 如上圖所示 start_position 為當前滑窗 答案:斯特魯維 用 str.find 去找並返回當前此答案的位置。而 end_position 就會是 start_position+len(answer)。

(b)After your model predicts the probability of answer span start/end position, what rules did you apply to determine the final start/end position?

"hidden size":1024 出來外面會接上一層 nn.linear(1024, 2)，2 代表一個是 start position 另一個是 end position，最後會接上 softmax 來做 normalize 機率，最後再選出機率最大的 start position 和 end position。

```
Model= BertForQuestionAnswering.from_pretrained()
sequence_output=model(input_ids, segment_ids, input_mask, start_positions, end_positions)
sequence_output .shape = [batch_size, max seq len, 1024]

self.qa_outputs = nn.Sequential(nn.Linear(config.hidden_size, 2))
logits = self.qa_outputs(sequence_output)
logits.size=[batch_size, max seq len, 2]
```

## Q2: Modeling with BERTs and their variants

1. Describe (2%)
    a. your model
       Chinese-roberta-wwm-ext-large
    b. performance of your model.
       EM:76.6
       F1:83

c.  the loss function you used.

Loss 分成 start logit 以及 end logit 兩個部分:

loss_fct = CrossEntropyLoss()

start_loss = loss_fct(start_logits, start_positions)

end_loss = loss_fct(end_logits, end_positions)

total_loss = (start_loss + end_loss) / 2

d.  The optimization algorithm (e.g. Adam), learning rate and batch size.

Implements BERT version of Adam algorithm with weight decay fix

BertAdam: lr=3e-5

warmup_proportion=0.1(前面 10 趴的 data 會使用較小的 LR 訓練)

b1=0.9, b2=0.98, e=1e-6, weight_decay=0.01

2.  Try another type of pretrained model and describe (2%)

a. your model

Chinese-bert-base

b. performance of your model.

EM: 68

F1: 71

c.  the difference between pretrained model

RoBERTa 採用的是 dynamic masking 方式，且在 optimizer adam 超參數 β2 改為 0.98 這樣一來使的 training 在大的 batch size 比較穩定，在訓練數據上讓 model 看了更多的資料，包含了 BookCorpus, CC-News, OpenWebText, Stories，這也是 RoBERTa 比 Bert 還要強大的主要原因。

在 model config 上得差異如下:

Chinese-bert-base:

```
{
  "attention_probs_dropout_prob": 0.1,
  "directionality": "bidi",
  "hidden_act": "gelu",
  "hidden_dropout_prob": 0.1,
  "hidden_size": 768,
  "initializer_range": 0.02,
  "intermediate_size": 3072,
  "max_position_embeddings": 512,
  "num_attention_heads": 12,
  "num_hidden_layers": 12,
  "pooler_fc_size": 768,
  "pooler_num_attention_heads": 12,
  "pooler_num_fc_layers": 3,
  "pooler_size_per_head": 128,
  "pooler_type": "first_token_transform",
  "type_vocab_size": 2,
  "vocab_size": 21128
}
```

Chinese-roberta-wwm-ext-large:

```
{
  "architectures": [
    "BertForMaskedLM"
  ],
  "attention_probs_dropout_prob": 0.1,
  "bos_token_id": 0,
  "directionality": "bidi",
  "eos_token_id": 2,
  "hidden_act": "gelu",
  "hidden_dropout_prob": 0.1,
  "hidden_size": 1024,
  "initializer_range": 0.02,
  "intermediate_size": 4096,
  "layer_norm_eps": 1e-12,
  "max_position_embeddings": 512,
  "model_type": "bert",
  "num_attention_heads": 16,
  "num_hidden_layers": 24,
  "output_past": true,
  "pad_token_id": 1,
  "pooler_fc_size": 768,
  "pooler_num_attention_heads": 12,
  "pooler_num_fc_layers": 3,
  "pooler_size_per_head": 128,
  "pooler_type": "first_token_transform",
  "type_vocab_size": 2,
  "vocab_size": 21128
}
```

d. The optimization algorithm (e.g. Adam), learning rate and batch size.

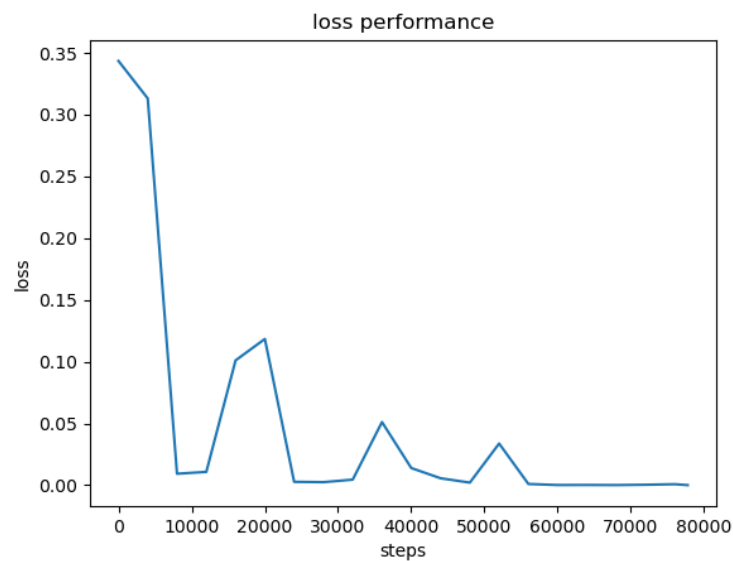Implements BERT version of Adam algorithm with weight decay fix

BertAdam: lr=3e-5

warmup_proportion=0.1(前面 10 趴的 data 會使用較小的 LR 訓練)

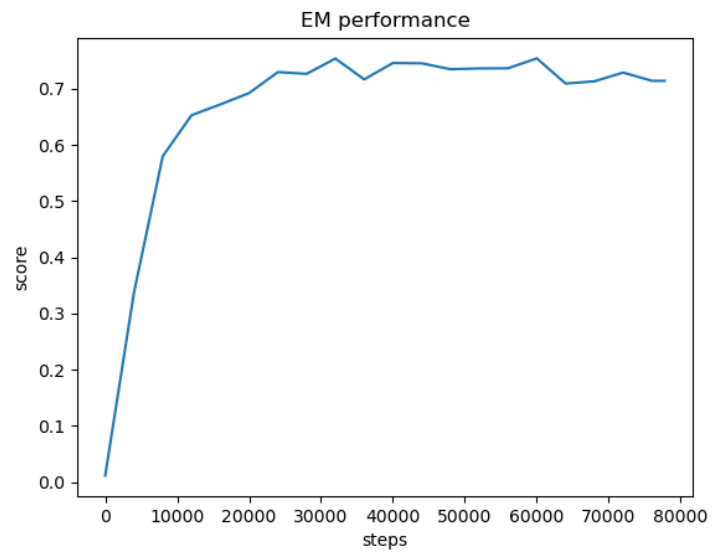b1=0.9, b2=0.999, e=1e-6, weight_decay=0.01

# Q3: Curves

使用 matplotlib 套件繪出，每訓練 4000 步得到一組資料，且為第一個 epoch。

## 1. Plot



loss performance

## a. learning curve of EM



## b. learning curve of F1