

Final Project – Written judgment summarization

第 28 組 書面報告

組長：趙昌昱

組員：吳瀚惟、陳諺樟、鄧新泰

首先，install 與 import 所有會使用到的套件，如下圖所示：

```
import transformers
from transformers import AutoTokenizer, AutoModelForSeq2SeqLM, AutoModelForCausalLM
import torch
from torch.utils.data import Dataset, DataLoader
from torch.optim import AdamW
from tqdm.auto import tqdm
from ignite.metrics import Rouge
import re
import pandas as pd
from sklearn.model_selection import train_test_split
import numpy as np
import jieba
```

1.Data preprocessing (30 pts)

Truncation

我們觀察到原文判決中的資料非常長，有的多達 8000 字以上，因此要讓資料有效率地被放入模型訓練，就必須做 truncation。而我們對 truncation max_length 做了很多的嘗試，最終試到當 max_length=1892 時，模型會有最佳的表現。

```
def get_tensor(sample):
    # 將模型的輸入和ground truth打包成Tensor
    model_inputs = tokenizer([each["origin_context"] for each in sample], max_length=1892, padding=True, truncation=True, return_tensors="pt")
    model_outputs = tokenizer([each["summary"] for each in sample], max_length=1892, padding=True, truncation=True, return_tensors="pt")
    return model_inputs["input_ids"].to(device), model_outputs["input_ids"].to(device)
```

Sliding window

將 input 原文依據 sliding window 大小分割成多段重疊的文字，並透過文字間的關聯性做摘要，最後再將所有結果串接而成，來處理比較長的文本。

```

def __init__(self, window_size=512, overlap=128) -> None:
    super().__init__()
    data_df = pd.read_excel(dataset_path)
    self.data = []
    self.window_size = window_size
    self.overlap = overlap
    for index, row in data_df.iterrows():
        origin_context = str(row['裁判原文'])
        summary = row['摘要']
        windows = self.sliding_window(origin_context)
        for window in windows:
            self.data.append({"origin_context": "summary: " + window, "summary": summary})

def sliding_window(self, text):
    start = 0
    text_length = len(text)
    windows = []
    while start < text_length:
        end = min(start + self.window_size, text_length)
        windows.append(text[start:end])
        start += self.window_size - self.overlap
    return windows

```

2. Model & Training method (20 pts)

Model

我們嘗試過很多 Transformer-based 的模型，盡量找尋有被中文 pretrain 過的模型，最後選用”heack/HeackMT5-ZhSum100k”作為最終模型，原因是它不僅被中文文本 pretrain 過，在 pretrain 的過程中也有學習過中文 summarization 的任務，且實驗的結果，用此模型訓練出的 evaluation score Rouge-2 最高，達到 0.20 左右，此外，訓練時間也不會太久，使用 colab 上的 A100 GPU 進行訓練，一個 epoch 大約花 7 分鐘左右。



Training 過程

我們訓練 10 個 epochs，每個 epoch 的訓練時長約為七分鐘，以下為訓練過程圖。



超參數設定

以下是實驗出最佳的一組超參數，其中，truncation_max_length 如上文所述，用於將資料前處理的 truncation 步驟，而 model_output_max_length 代表 model generate 出的文字最多可達到 2048 個字元。

- ***lr = 1e-4***
- ***epochs = 10***
- ***train_batch_size = 1***
- ***validation_batch_size = 1***
- ***test_batch_size = 1***

- *random_seed = 42*
- *truncation_max_length = 1892*
- *model_output_max_length = 2048*

Evaluation 結果:

以下是透過上述 model 與超參數訓練出來的結果，達到 Rouge-2 分數 0.20 的成績。

```
'Rouge-2-P': 0.20008906707518975,  
'Rouge-2-R': 0.2058318986853808,  
'Rouge-2-F': 0.2058318986853808}
```

3. Analysis (30 pts)

Difficulties

1. 算力不足：在訓練大模型或 sliding window 時需要花大量時間，使用我們的資源 train 不太動。因為沒有能高速運算的本地端 GPU，因此我們購買 colab pro 的雲端 GPU 來補足我們的算力。



2. 嘗試過很多 Transformer-based 的模型，效果大多不佳，最後採用的模型如上文第二點所述，以下列出我們嘗試過的模型及失敗原因。

嘗試過的模型	失敗原因
yihuan/mt5_chinese_small	效果不佳，重複生成一樣的文字 如: 不斷生成判決文本中的名字「葉文勇」
csebuetnlp/mt5_multilingual_XLSum	memory爆掉
IDEA-CCNL/Randeng-Pegasus-238M-Summary-Chinese	tokenizer有bug: “AttributeError: 'PegasusTokenizer' object has no attribute 'vocab' ”

3. 調整 learning rate 做訓練時，比想像中來得困難，不論調大或調小都會導致效果不佳，**仍然是以 sample code 最初提供的 $1e-4$ 為最佳**。以下是將 learning rate 調大為 $1e-3$ 產生的問題：模型對於不同判決會產生相同的輸出，如下圖所示。

```
output: /</s>,於之第。上其人行有法或刑、2及得公本而原部1任已「文月所者明昭年o自(件部前直修正依分裝要犯罪日理法院南未審被告物4元就系起止相3因之事裁判由主土地 此名金非債款可至施行10刑法行使公司表轉target: [「按行保其所有司事招而止之裝,其之裝型是否否該法第266款所載,就案事及整警者,裝製之型特者,依人(係人)(係公司)之性程度之高底何之,即得否自由交付定付之方式(包含工作),自行(例如按所招之係收受之伴保蓋算)model generate examples:  
output: /</s>,於之第。上其人行有法或刑、2及得公本而原部1任已「文月所者明昭年o自(件部前直修正依分裝要犯罪日理法院南未審被告物4元就系起止相3因之事裁判由主土地 此名金非債款可至施行10刑法行使公司表轉target: [「按入基於該第242定代,代人止其第3之人係否償還,而求第3人所有移登入者不,不得代人即(人)列到共犯者,否其人部分之,予回;本件上主土地乃係虛名登於○名中,民法第242定代位他○傳止model generate examples:  
output: /</s>,於之第。上其人行有法或刑、2及得公以本而原部1任已「文月所者明昭年o自(件部前直修正依分裝要犯罪日理法院南未審被告物4元就系起止相3因之事裁判由主土地 此名金非債款可至施行10刑法行使公司表轉target: [「又民事法第408條2數判之行,既明有足既有力,其原因部分主生法律上之效力,而使不利之事人,就部判力,自有上利益,不受原判令文形式上重違背之拘束。原上人以其被上人所為63,205,550元之履保金,其被上人所上返model generate examples:  
output: /</s>,於之第。上其人行有法或刑、2及得公以本而原部1任已「文月所者明昭年o自(件部前直修正依分裝要犯罪日理法院南未審被告物4元就系起止相3因之事裁判由主土地 此名金非債款可至施行10刑法行使公司表轉target: [「按依保罪情之輕重,去核按其在民事事情中,於獻保部下之生意提供,是否盡忠履行犯罪之事。法院於罪之檢察資料,如何正地行使使其裁量,得免欠缺及任何之責任,英美法等○罪量的引○,可否○亦即,在被訴之偵緝後○model generate examples:  
output: /</s>,於之第。上其人行有法或刑、2及得公以本而原部1任已「文月所者明昭年o自(件部前直修正依分裝要犯罪日理法院南未審被告物4元就系起止相3因之事裁判由主土地 此名金非債款可至施行10刑法行使公司表轉target: [「按行改其具有程序法第111合款清之一而致外效,具有勇力,於來類。止或原因因時而失其效,其效乃存在(行政程序法第118節)之,受其他處置之,除非有充分理據的,否該法亦有再效果可能,尚待查○先判而行model generate examples:  
output: /</s>,於之第。上其人行有法或刑、2及得公以本而原部1任已「文月所者明昭年o自(件部前直修正依分裝要犯罪日理法院南未審被告物4元就系起止相3因之事裁判由主土地 此名金非債款可至施行10刑法行使公司表轉target: [「按犯罪罪所,則犯罪行為人,取之、於或部或一不能或,或有不行收,追繳其入止又犯罪抑而得者。其得之物上利益其異,再取得罪者已含法被害人者,不宣當或刑罰權限,歷105年7月1日起施行之刑法第Route-2 rscno, en, enpch B、「Route-1」: 0.5824296451567361, 「Route-1-R」: 0.6566271431739264, 「Route-1-L」: 0.6566271431739264, 「Route-2」: 0.8429991238773275, 「Route-2-R」: 0.84291591675
```

而將 learning rate 調小為 $1e-5$ 也會產生另一種問題：模型會在一筆資料中產生大量重複的文字，如下圖所示。

[illegible]

Testing analysis

將訓練好的模型預測給定的 test 資料集，觀察到以下現象：

testing example1:

以 test 資料集中的第 45 筆資料為例，將原文摘要後的結果大多文句通順且文整。以下為模型生成第 45 筆資料的摘要：

G	H	I	J	K
裁判字號	爭點	裁判日期	裁判案由	裁判原文
最高行政法院106年度判字第305號判決	原住民身分取得要件。	106年6月15日	原住民身分法	最高行政法院判決106年

摘要：

一、原住民身分法第8條第1項之「應具」一謂「取得」之不同用語即知，故「係指「本該具有」，此參原人間身分法採血廠主義精神，應指自然血統而言，其可從二者一稱「發」所謂『取得』之方式即為「兼顧」者，並兼問其子女相關權益，爰增訂第2項如當事人已死亡，及子女準用本法第4條及第7條之規定取得該住民時間，以保障當人直系血親卑親屬之義務。」而該立法理由雖未提及「具有原籍人在80年10月4及過世未能於生前親自申請復當，而未取得新住民變身者。是故，若具本住民血運，為落實原滿住民文書法施行前，因.....或其他原因緝失或未得取得第住民性質者為限。而上訴人之祖母既具山地與泰雅族，依血○○之精神等語，益足證之。二、(二)上開卷(三)第3款、第4項所列情形，包括第6項「原種與非原□結婚所生子女，從具原個人身分之父或母之姓或前住民傳統名字者」，於完成從父 or 母姓及傳運名字後，即能取得。則「並具(本於具有)之母性，自有「即有」或「已」的「未具」者」為不合。

testing example2:

以 test 資料集中的第 24 筆資料為例，裁判原文未出現「行政罰法第5條「從新從輕原則」」等文字，然而摘要中卻能產生出此法條，因此我們推測模型能從之前學習過的資料中預測相關的知識。以下為模型生成第 24 筆資料的摘要：

G	H	I	J	K
裁判字號	爭點	裁判日期	裁判案由	裁判原文
最高行政法院106年度判字第417號	按出賣人究應於何時辦竣戶籍登記	106年8月3日	特種貨物及勞務稅條例	最高行政法院判決106年

摘要：

一、行政法院就行政法上義務之行為，應於出售日次日起3個月內將戶籍地由原房地直接遷移至新房用地」之限制，自無以行政機關內部自行規定之作業標準，作為對外裁罰人民、限制人民權利之依據。又依行政罰法第5條「從新從輕原則」，被上訴人於102年間購置系爭房子地，並於1992年4月10日函釋未發布為人民所周知。

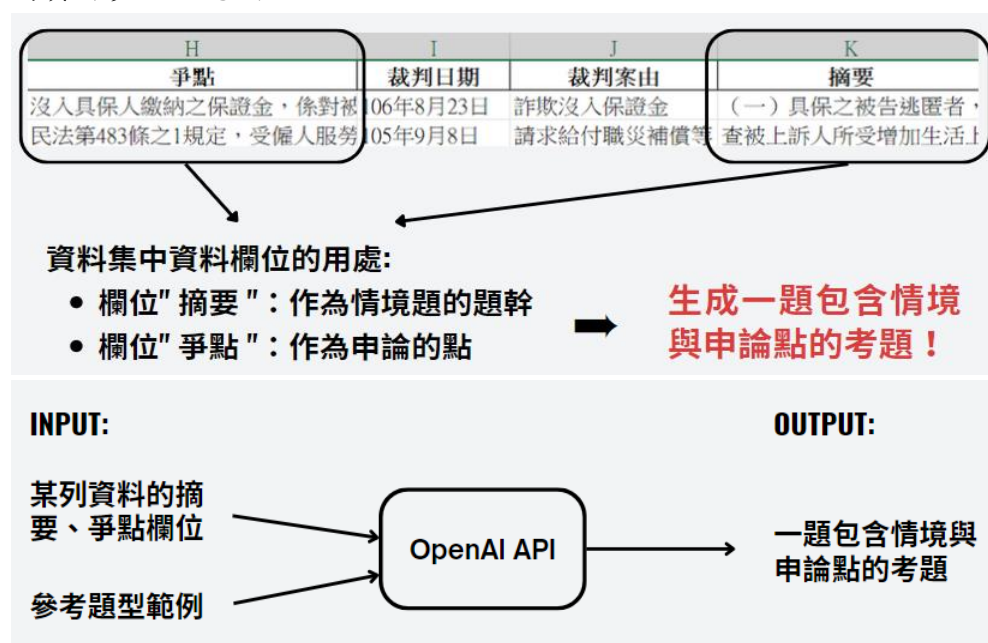
4.Proposal (20 pts)

資料集應用 1：作為法律課程情境與申論題的出題題庫

透過請教身邊的法律系朋友，我們得知法律系學生在課程中的考題很多時候會是以申論題的形式，且題目的來源多為一個判決。因此，我們認為這個資料集可以作為法律課程情境與申論題的出題題庫，來幫助出題教授省去主動摘要判決與想題目的時間。

實作方法：

將資料集的資料欄位"摘要"作為情境題的題幹，欄位"爭點"則作為申論的考點。此外，也提供已存在的範例題型，透過 OpenAI API 的實作，便可以生成一題完整的題目。透過上述做法，一個法律課程的情境與申論題出題題庫就完成了。以下為實作步驟示意圖：



預期結果：

以 train 資料集的第二筆資料為例，是一個關於職業傷害案件的判決。下圖為透過上述步驟生成出的題目，可以看出此題目完整、符合判決內容且也要求學生將爭點作為申論點進行論述與作答。

就上訴人與被上訴人之間的職業傷害事件，評估以下爭點：

- 1.原審對於過失分擔的認定是否合理？請以被上訴人的專業背景和發生事故的情境來分析原審是否應將50%的過失責任歸咎於上訴人。
- 2.探討原審判決在計算賠償金額時先扣抵再進行過失相抵的方法是否適當，並考慮其對案件判決結果的影響。

要求：學生需根據判決摘要及爭點，分析上述問題，並提出合理的法律論點支持其見解。分析時應考慮法律原則、案件事實以及相關法律規定，並需引用適當的案例或法條來支撐其論述。

資料集應用 2：作為種子資料集，生產更大量且有品質的資料

經過對資料集的觀察後發現這樣的法律資料筆數不多，是相當珍貴的資料。因此，若能將這些真實的判例資料作為種子資料集，生成更多高品質的仿真資料集，想必對法律應用面的發展會有幫助。下圖為本次 train 資料集的資料筆數：



實作方式：

- STEP1: 透過 LLM 將種子資料集生成仿真的合成資料集
- STEP2: 兩個 LLM 對於合成出的資料集互相評分
- STEP3: 根據評分、回饋做資料清洗
- STEP4: 得到高品質資料集

預期結果：

經過上述的實作方式，便可以從少量珍貴的真实資料集產生大量高品質的資料集。

以上是我們的報告，謝謝。