

仿 PTT 鄉民問答 Bot

108321030 劉承熙
108321032 王廷郡



動機

- 因為平常喜歡看 PTT 的鄉民梗，雖然有時候PTT鄉民發言時常政治不正確，但是有些用語很有趣而且很特別，算是一種次文化
- 剛好 Kaggle 上有人整理 PTT 中文語料庫 PTT-Gossiping-Corpus

【板主:ubcs/ffooxx/moth...】				[八卦] Covid-19exp開新板徵文中!		看板《Gossiping》	
[←]離開 [→]閱讀 [Ctrl-P]發表文章 [d]刪除 [z]精華區 [i]看板資訊/設定 [h]說明							
編號	日期	作者	文章標題				人氣:14418
780883	+ 3	6/07 zakijudelo	[問卦] 大家吃過最貴的一餐是吃什麼				
780884	+ 8	6/07 Daniel0712	[問卦] 最希望哪家美國速食店來台灣?				
780885		6/07 yoyoruru	R: [新聞] 陳時中自打臉! 去年看別國每日上千確診				
780886	+ 2	6/07 shrinkage	[問卦] 英國開始實驗週休三日				
780887		6/07 mdm	R: [問卦] 四叉貓噁心成這樣, 以後公務員都不要				
780888	+13	6/07 doro0202	[問卦] 吃到飽明明很便宜, 為何還要用wifi				
780889		6/07 mk203125	R: [問卦] 眼鏡一副配到4000是不是太貴了?				
780890	+ 7	6/07 bornwinner	[問卦] 滴妹男友才24歲, 不急結婚正常吧?				
780891	+ 3	6/07 KennethC	[問卦] 台式點心怎麼輸這麼慘?				
80892	+ 4	6/07 w1230319	R: [問卦] 公務員不可以上班玩PTT可以看影片嗎				
80893	+ 1	6/07 icrol11	R: [新聞] 快訊/退休消防員曝: 恩恩爸聽偽造錄音				
80894	+14	6/07 Busufu	[問卦] 幹雨傘的人是不是從爸爸屁眼出生的?				
80895	+ 1	6/07 app325	[問卦] 牙醫是不是太多了, 預約都要一個月?				
80896	+ 1	6/07 poeta	[新聞] 8個月後返台! 三原震驚台灣5大變化嘆:				
80897	+ 2	6/07 Tenc	[問卦] 三菱是不是很不會做生意				
80898	+ 2	6/07 ilv1181023	[問卦] 31歲一事無成, 焦慮中				
80899	+12	6/07 t21	R: [新聞] 高雄+12555創新高 陳其邁: 廣發快篩、減				
80900	+ 6	6/07 sted0101	R: [新聞] 快訊/退休消防員曝: 恩恩爸聽偽造錄音				
80901	+ 1	6/07 ccyaztfe	R: [新聞] 新竹房市好瘋! 房子淪法拍、房價反漲逾5				
780902		6/07 NONOTV	[新聞] 海科館新船名方舟 遭蘇貞昌嗆「像逃命那				
文章選讀 (v)回應(X)推文(*)轉錄 (=[]<>)相關主題(/?a)找標題/作者 (b)進板畫面							

Dataset - 資料來源

- PTT-Gossiping-Corpus 的資料太舊，最新的資料還停留在2019年
 - 網路流行語變化很快，可能不太適用
 - GitHub 爬蟲程式無法運作
- 自己寫爬蟲程式
 - 爬取八卦板中分類為 “[問卦]” 的「標題」及「推文」
 - 每 2.5 秒爬一個頁面，爬了大約五天
 - 取 2022-02-01 ~ 2022-05-28 共計 108900 篇文章

2018年度10大PTT鄉民流行用語揭曉			
排名	用語	網友怎麼說	網路聲量
1	韓導	「被政治耽誤的導演」	72,189
2	1124滅東廠	「年底用選票給民進黨教訓」	34,245
3	發大財	「OO進的來，OO出的去，OO發大財！」	31,849
4	又老又窮	「台灣有哪裡不是又老又窮？」	25,601
5	越想越不對勁	「台男要好好保護自己」	20,043
6	主流民意	「大家好我們是對同志不友善的中國台北人」	17,817
7	喜韓兒	「韓粉出征，寸草不生」	16,914
8	貪婪老人	「1500=台北價值」	16,765
9	擊潰丁守中	「擊潰丁守中！」	13,304
10	關三天	「政客、名嘴要不要先關一關？」	9,973

資料分析：DailyView網路溫度計 透過 KEYPO大數據關鍵引擎(keypo.tw) · 以國際級的語意分析架構、先進的機器學習技術與人工智慧推論引擎 · 感知網友語意脈絡與情緒，分析時事網路大數據。
分析期間: 2017/12/17~2017/12/16



Dataset - QA對應

- 每篇問卦的文章有不定數量的推文，要找出適當的回應當作回應
 - 「標題」: 問題(Q)
 - 「推文」: 回應(A)
- 使用詞出現的數量當作 weight, 每個推文計算加權分數, 取最高者

{‘貓’: 4, ‘可愛’: 1, ‘這’: 1, ‘隻’: 1,}

推 TPDC: 貓貓可愛
推 horseorange: 這隻貓太肥了吧
噓 iamshiba: 拿貓騙推

$$4 + 4 + 1$$

$$\frac{\quad}{3}$$

Dataset - QA對應

- 無法作到完美, 但是大多數都看得出來問答對應關係

平板該買哪一台
除了蘋果 都是電子垃圾

有沒有外國月亮比較圓的八卦
一堆人移民去美國啊 你問他們

急！搭雲霄飛車時遇到一高一胖的黑衣人？
你要被灌藥了

猴痘是怎麼傳播到多國的啊
悟空的如意金菇棒打的

美國人發明不出抖音？
看YT的演算法有多廢就知道了

接到博客來電話是詐騙嗎
接到直接回我是買金石堂

請問宵夜要吃什麼好？
樓下要請我ㄅ什麼

今天違規停車在公車站牌會被檢舉嗎？
公車站牌記得可以檢舉

彥均這篇文章是不是太早發了
等死幾十個也不用說了

急！女兒一邊讀書一邊給我打瞌睡怎麼辦？
比大雄好多了~要知足~

美國野豬氾濫成災？
沒有從小閹割的豬騷味很重的

住外面的旅館或飯店或民宿最在意什麼？
獨立衛浴 有浴缸

Dataset - PTT 常見詞

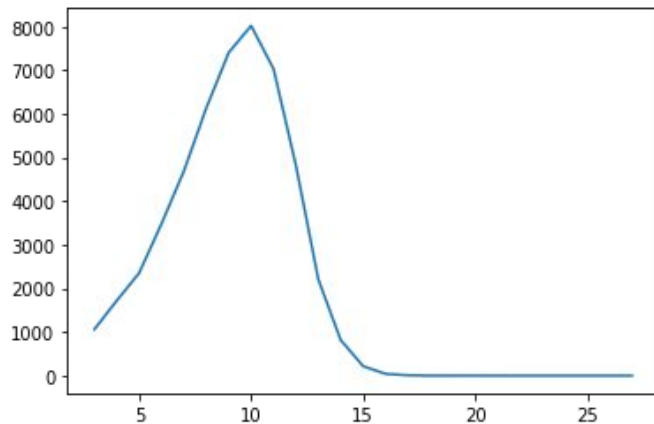
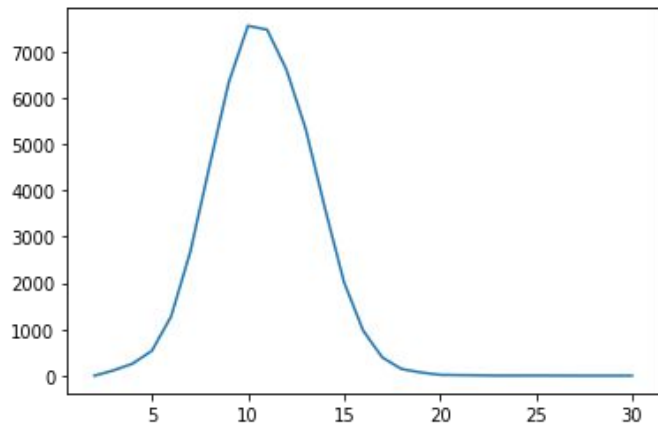
- Python WordCloud 套件



Preprocess - Word-based

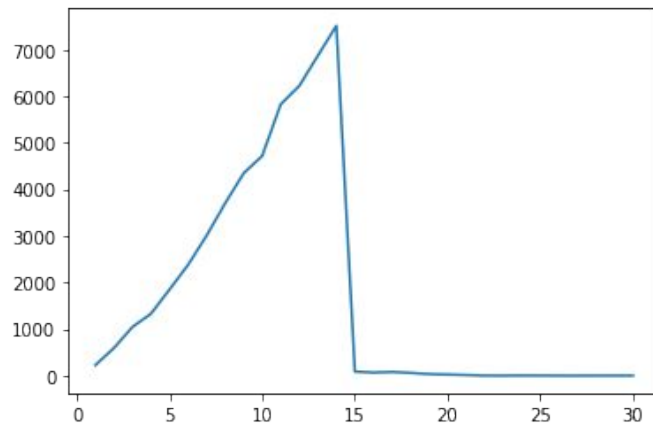
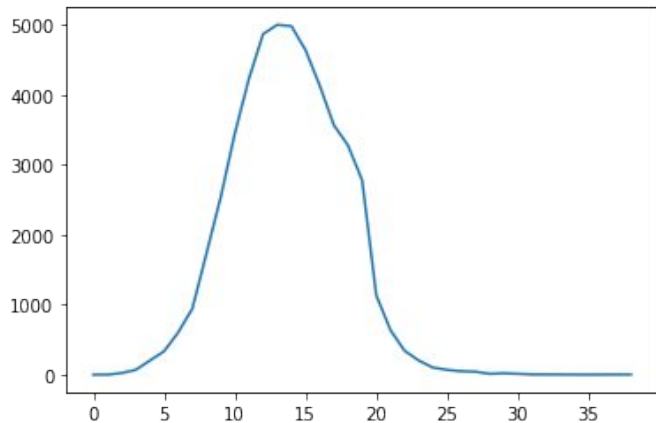
- 為所有詞建立 Tokenizer, 總共 56130 個詞
- 分析 Q 和 A 的 token 數量分佈, 大都落在 7~13 左右
 - QA 最長都可能到 30, 因此補 padding token 到長度 30

['<sos>', '平板', '該', '買', '哪', '一', '台', '<end>', '<pad>', '<pad>', '<pad>', '<pad>']



Preprocess - Char-based

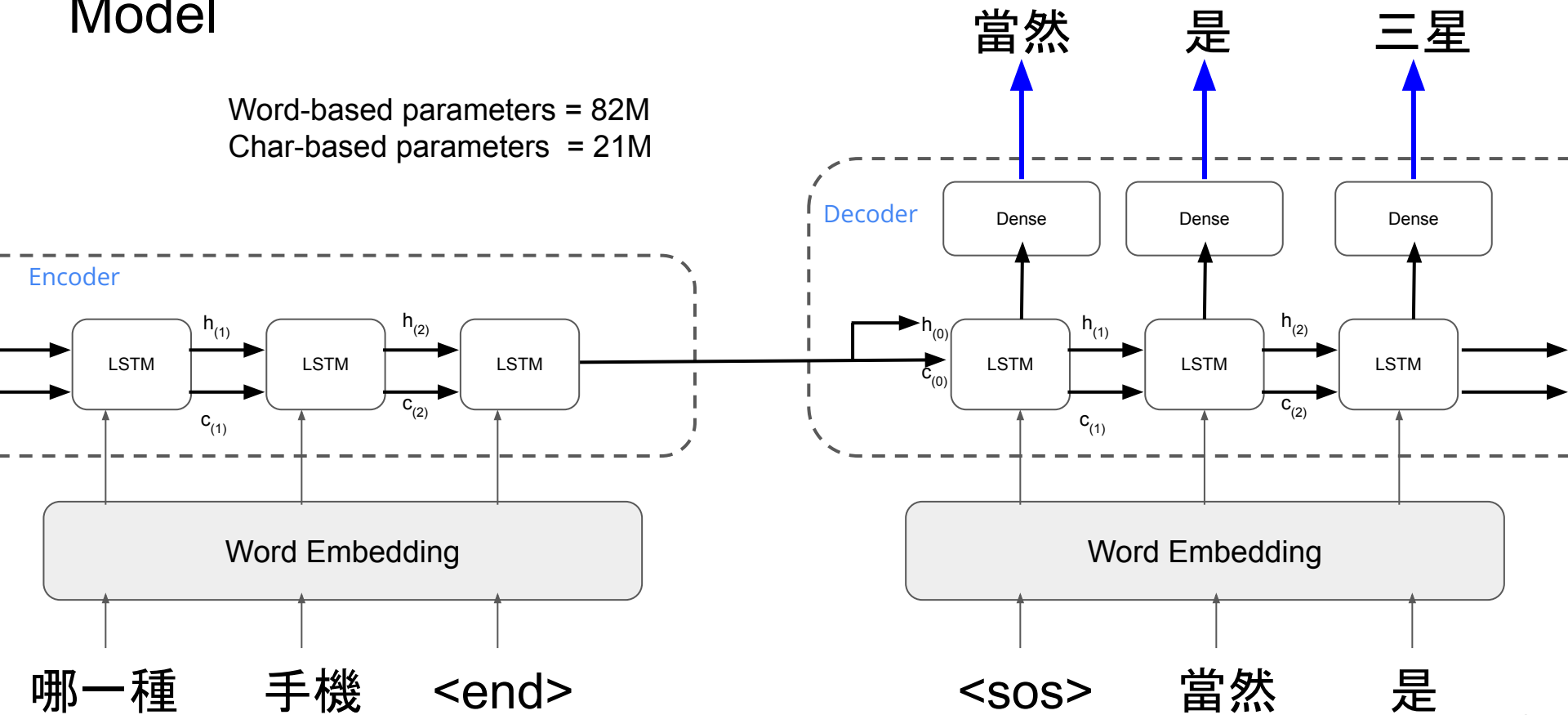
- 為所有字元建立 Tokenizer, 總共 5434 個詞
- 分析 Q 和 A 的 token 數量分佈
 - Q 最長都可能到 37, 因此補 padding token 到長度 40
 - A 最長都可能到 30, 因此補 padding token 到長度 30



Model

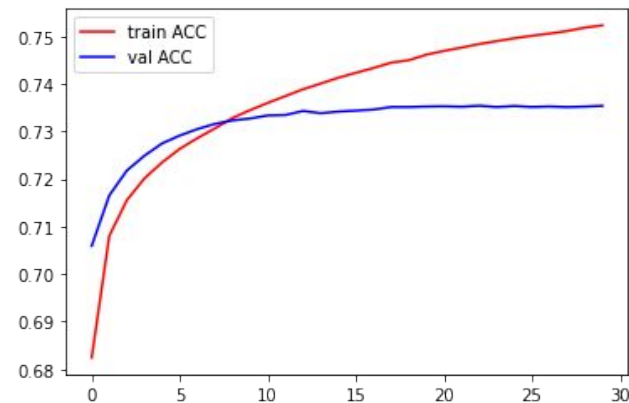
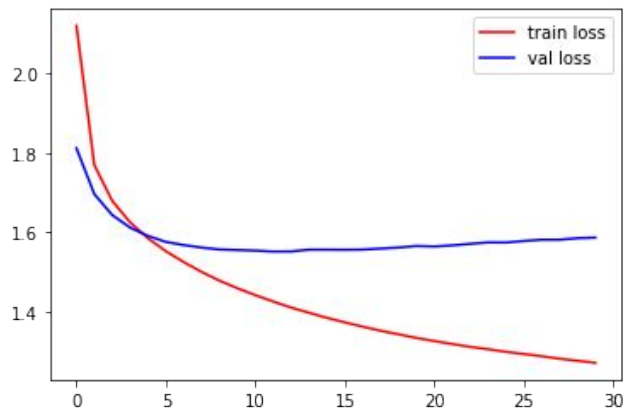
Word-based parameters = 82M

Char-based parameters = 21M



Char-based Training

- Epoch: 30
- Loss function: Cross Entropy
- Optimizer: Adam
- Train:Val = 10000:8900



Word-based Predict

- 確診了體溫只有35度多這樣正常嗎？
 - 現在都低加盟金了
- 死很多的定義是幾個啊？
 - 台灣的人情味是世界前段班
- 清明節大家都有吃潤餅捲嗎？
 - 我都用傳送回家的卷軸

Char-based Predict

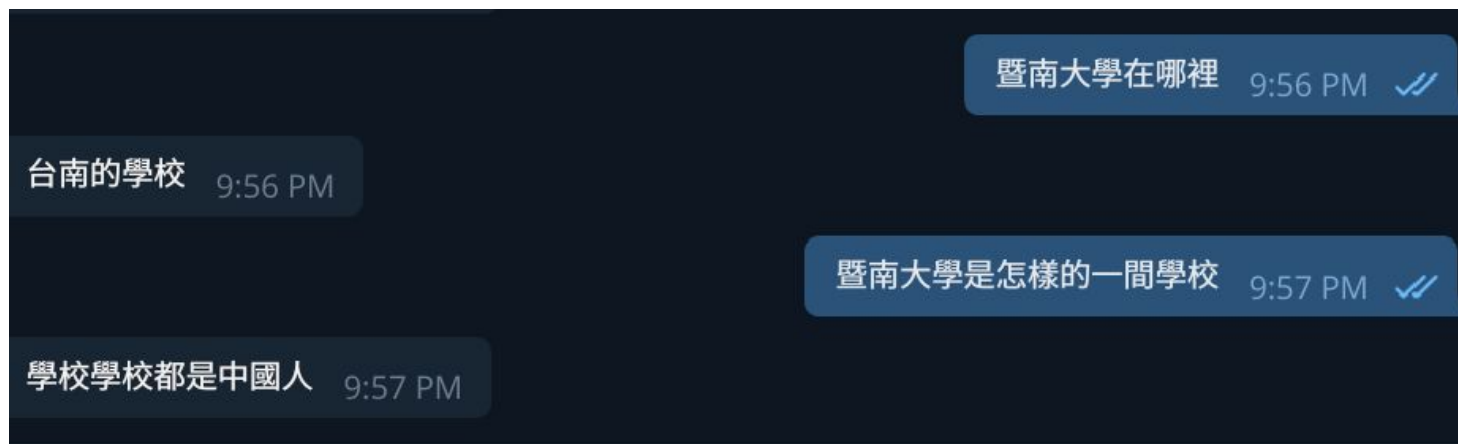
- 清心的飲料點全冰會出現多少冰塊？
 - 我都喝冰冰冰冰塊冰塊冰塊
- 排隊群聚搶打疫苗？
 - 打疫苗的人都是打疫苗的
- 現在還敢生小孩的人都是怎樣的人
 - 你的人生都是老人的人

Telegram Bot

- 使用 python-telegram-bot 開發，開發前要先跟 Telegram 的 BotFather 申請 Token
- 執行在 i5-8265U、MX150 顯示卡的筆電上，需要 1~2 秒的預測時間



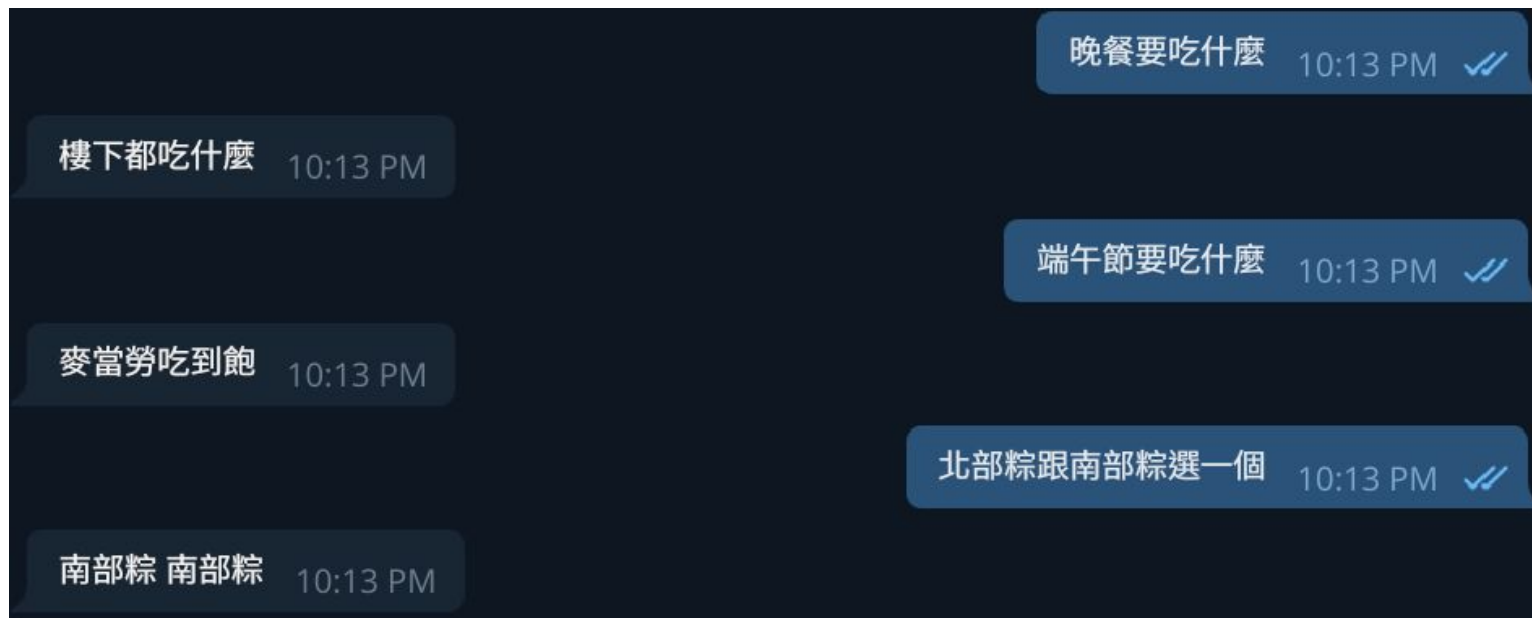
Example 1 - 暨大



Example 2 - 投資



Example 3 - 食物



Example 4 - 食物



Example 5 - 防疫



Example 6 - 防疫



Conclusion

- Char-based 成果比 Word-based 好推測是斷詞不夠好的關係
 - 有許多詞會被錯誤的切割開
 - 尤其在新詞彙、特殊詞彙比較多的 PTT
- Improvement
 - 使用 CKIP 的中文 Pretrain Word Embedding
 - Attention
 - BERT