

Summarize Medical News using Pretrained Language model



國立臺灣大學
National Taiwan University

Team Members



online collaboration



online collaboration



蔡尚錡 D08922014

Contributions :
method discussion、
assignment slides、
references survey

郭蕙綺 R09922A21

Contributions :
method discussion、
SOTA model building、
main experiments

呂承翰 R09922A01

Contributions :
method discussion、oral
presentation、result
analysis

What's the Problem

Problem description

- Given a Chinese news or post, we want to build a model to learn the meanings of articles and select the most critical words or phrases to generate the summaries.

AI problem Definition

- Input : Chinese articles / news
- Output : Chinese summarized sentence
- Model: sequence-to-sequence with pretrained language model
- Method: supervised learning / reinforcement learning
- Evaluation: Rouge-L

Examples :



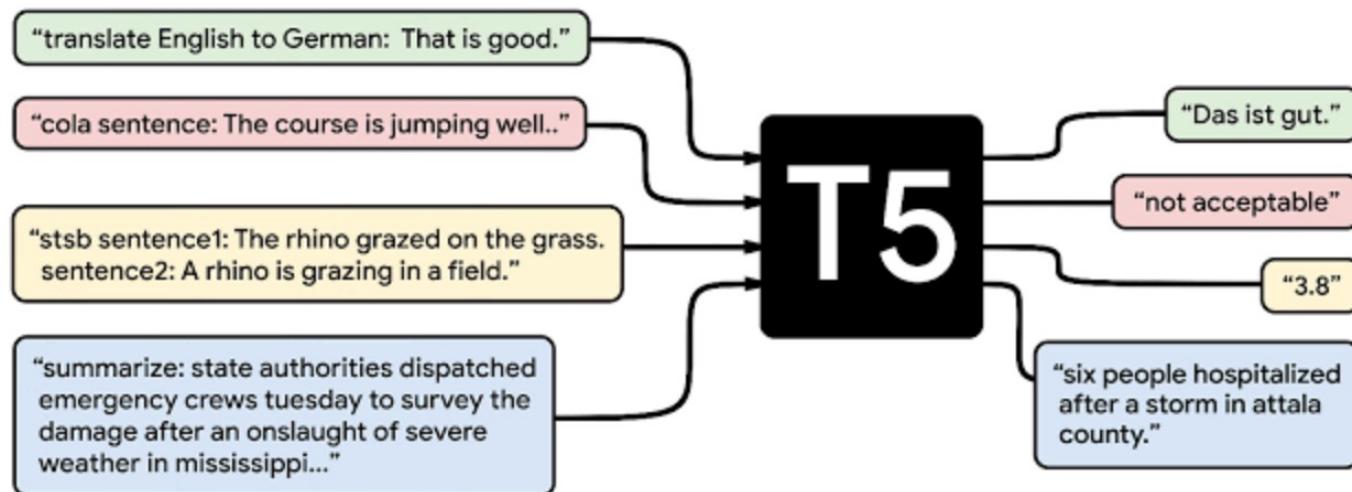
尖端醫攻新冠檢測試劑，快篩國家隊成軍

Why is it important

- Summarization is an important task in NLP field
- Information about coronavirus and disaster dominate the news and social media.
- If we can summarize the news promptly, it is helpful to filter the article and quickly find the related information we need.
- How to learn the semantics and syntax of the articles and generate more natural and informative sentences is a significant challenge.
- Chinese news summarization is a relatively less researched field which contains some different natural language problems

What's the state of the art

- Pretrained language model : **T5 (mT5)**
- T5 is an encoder-decoder model pre-trained on a multi-task mixture of unsupervised and supervised tasks and for which each task is converted into a text-to-text format.
- T5 works well on a variety of tasks, e.g., translation and summarization.
- Achieved 43.52/21.55/40.69 of rouge score on CNN Daily Mail dataset.



What's our big idea ?

- Encoder-Decoder models: they often generate unnatural summaries consisting of repeated phrases.
- teacher forcing algorithm does not always produce the best results on ROUGE
 - exposure bias
 - large number of potentially valid summaries
- ROUGE metrics take some of this flexibility into account, but the maximum-likelihood objective does not.

Apply reinforcement learning to improve the rouge score.

- Learn a policy that maximizes a specific discrete metric with RL
- policy gradient training algorithm
- Define a mixed learning objective function that combines maximum-likelihood training objective and reinforcement learning objective.

How do we plan to solve the problem ?

- Methodology : Combine supervised learning and reinforcement learning
 - Policy gradient : $\nabla \bar{R}_\theta = \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} R(\tau^n) \nabla \log p_\theta(a_t^n | s_t^n)$
 - Reward function : for an output sequence, comparing it with the ground truth sequence with the evaluation ROUGE
 - Learn a policy P_θ to maximize a specific discrete metric
$$\sum_{t=1}^{n'} \log p(y_t | y_1, \dots, y_{t-1}, x)$$
 - $L_{mixed} = \gamma L_{rl} + (1 - \gamma)L_{ml}$: Conditional language model can assist our policy learning algorithm to generate more natural summaries and improve readability.
- Resources
 - Crawl general news on udn.com
 - Crawl medical news on udn.com
- Source codes
 - https://drive.google.com/file/d/1LI2fbUmb3p_rIvBxXroHRwQjEPEmq2q/view?usp=sharing

Methodology - Supervised learning

- Training : Sequence-to-sequence model
 - Use Seq2SeqTrainer from huggingface.co
 - Tokenize: Use pretrained tokenizer from “google/mT5” from huggingface.co
 - Trainer: Sequence to sequence language model from “google/mT5” from huggingface.co
- Predict: predict function from the trainer.
 - Generate algorithm: Beam search.

Methodology - Reinforcement learning

- Training: We build a RLTrainer that inherit the Seq2SeqTrainer from huggingface.co
 - RL algorithm: Policy gradient.
 - RL_loss: $-(\text{Log_probs} * \text{Reward}).\text{sum}().\text{mean}()$
 - Log_probs: We find the maximum probability on each words of a generated sentence.
 - Reward: Use f1-score of rouge score.
 - Reward: $0.2 * \text{RL_loss} + 0.8 * \text{original_loss}$.
- Predict: predict function from the trainer.
 - Generate algorithm: Beam search.

Data Format

○ Training :

- Input : maintext
- Output : title

▼ 3:

title: "【身體不適特輯（下）】「拉肚子」、「想吐」、「嘔吐」英文怎麼說？"

maintext: "嗨，歡迎回到【身體不適特輯】的下集，在這集裡，小編會繼續講自己食物中毒的故事，以及介紹各種症狀的英文喔！一起來看看小編悲慘的遭遇吧！還沒看過上集的朋友點這裡：【身體不適特輯（上）】『食物中毒』、『昏倒』、『急診』英文怎麼說？ the runs / diarrhea 拉肚子小編昏倒就醫之後，後來就出現了「拉肚子」的症狀。大家多少應該有聽過 diarrhea，是「腹瀉」的意思，不過這是比較學術的名詞，口語一點的話可以說 have the runs、get the runs，就是俗稱的「拉肚子」囉。 My wife has had diarrhea for three days. (我老婆已經腹瀉三天了。) My family and I got the runs after dining in that posh restaurant. (我家人跟我在那間高級餐廳吃完晚餐後就開始拉肚子。) 如果要說「某東西讓某人拉肚子」，英文就可以用：something give(s) someone the runs，例如美味的布朗尼竟然讓你拉肚子，你就可以說： I can't believe those delicious brownies gave me the runs.

(不敢相信那些美味的布朗尼竟然讓我拉肚子。) feel nauseous / feel sick 想吐除了「拉肚子」，小編還有「想吐」的症狀，可以用 feel 再加上形容詞 nauseous (感到噁心的、想嘔吐的) 或是 sick (噁心的、想吐的)： I feel nauseous and cold. (我覺得想吐，身體感覺很冷。) She feels sick after having too much wine. (她喝太多酒之後都會想吐。) vomit / puke / throw up 嘔吐「想吐」之外，如果真的吐出來了，英文就可以用 vomit、puke 或是 throw up 表達「嘔吐」： The man kept vomiting; his wife, standing next to him, looked worried and anxious. (那個男人不斷嘔吐；站在他旁邊的妻子看起來非常憂心又焦慮。) The man's funny smell makes me want to puke. (那個男人身上奇怪的味道讓我很想吐。) 那如果要講「把某東西吐出來」可以用 throw up something / throw something up 或是 puke up something / puke something up: My cousin threw up all his breakfast; the doctor diagnosed him with stomach flu. (我表弟把早餐全吐出來了；醫生診斷他得了急性腸胃炎。) The little girl was sick, puking up all the porridge she just ate. (那個小女孩生病了，把剛剛才吃下肚的稀飯全吐出來了。) 來總複習一下，我們學到了「拉肚子」可以用口語常用詞 have the runs，「想吐」則可以說 feel nauseous / feel sick，真的吐出來的話，英文就可以說 vomit、puke、throw up 喔。後來在拉肚子幾天之後，小編終於漸漸痊癒了，希望大家健健康康，永遠不需要用到這些描述自己身體不適的英文 QQ 我們下次見啦！延伸閱讀 1. 【身體不適特輯（上）】『食物中毒』、『昏倒』、『急診』英文怎麼說？ 2. 「和『生病』有關的各種說法」 - Learn Phrases Related to Being Sick 【更多學英文資源，詳見《希平方》】"

id: "3"

Data Format

○ Testing:

- Input: maintext
- Predict: title (as summarization)

title: "快訊／柯文哲發布清零計畫 「2條件」做到才有可能解封"

maintext: "▲柯文哲召開防疫記者會。（圖／北市府提供）記者陳家祥／台北報導新冠肺炎疫情趨緩，各界開始喊話解封。對此，台北市長柯文哲公布清零計畫並訂出條件，第一，確診者不可以住在家裡；第二，以確診者為中心做的疫調，可能的被感染者有義務去做篩檢，「要讓台北要盡速清零，才能思考三級管制可不可以放鬆」。柯文哲表示，當快篩的陽性比率降到1%以下，過去的社區普篩大概是效果不大，所以要有新的思維。居家隔離人數慢慢在走下坡，趨勢看起來是終於明顯下降。至於有1、2個個案藏在社區裡面，到底怎麼解決？柯文哲說，全面打疫苗是一個策略，但從英國、美國的例子，美國昨天感染人數還有1萬，但打滿2劑的超過一半，所以表示說除非打到70%，不然終究還是會有感染。柯文哲說，算一算台北市疫苗才打到6%多，距離以色列的70%還有一段時間，但我們也不能一直封城下去，相信很多產業已經受不了，所以要想辦法。「以前確診者允許你住家裡，現在起通通不准，一但你是確診者，就算沒有症狀，也要去住防疫專責旅館。」柯文哲說，因為從北市的統計，家戶感染占43%，所以常常是家人感染給家人，加護感染又可能去外面傳染給其他人。「所以從現在起，除非有很強的理由，若一但確診，有症狀住醫院，沒症狀的送防疫專責旅館。」柯文哲強調，因為我們必須切斷感染源。你要留在家裡要簽切結書，以後我們不給你任何賠償，「死在家裡怪政府照顧不周，你本來就應該去防疫專責旅館接受醫護人員的監控」，因為我們要做清零計畫，當普篩陽性率不到1%，普篩已經無效，要以確診者為中心做同心圓，找出可能感染的人。柯文哲說，現在要動員健康服務中心所有人力，去做確診者的精準疫調，找出可能被感染的人，我幫你出計程車、快篩費用，但你必須去做檢查，因為不這樣的化沒辦法清零，「要讓台北要盡速清零，才能思考三級管制可不可以放鬆」。柯文哲強調，第一，確診者不可以住在家裡；第二，以確診者為中心做出來的疫調，可能的被感染者，市府會出交通費跟檢查費用，但是你有義務去做篩檢，不然傳染鏈沒有辦法打斷，請市民朋友擴大配合我們要做的篩檢，預防隱形的傳染鏈，這不解決，台北市的清零計畫不會成功。"

Output Comparison

- SL: “台南市長黃偉哲贈3萬份料理調理包 台南市長黃偉哲”
- SL+RL: “台南市「慈揚社會福利基金會」捐贈3萬份料理調理包”

title: "臺南市173處社區照護暫停共餐 黃偉哲感謝燦揚企業捐贈三萬份料理包關懷長輩"

mainText: "▲燦揚企業集團創辦的慈揚社會福利基金會，特別捐贈3萬份料理調理包予社會局，臺南市長黃偉哲17日也頒贈感謝狀感謝善心。（圖／記者林悅翻攝，下同）記者林悅／臺南報導 因應新冠肺炎疫情嚴峻，全國3級警戒延至6月28日，臺南市173處社區照顧關懷據點配合防疫暫停共餐服務，燦揚企業集團創辦的慈揚社會福利基金會，特別捐贈3萬份料理調理包予社會局，臺南市長黃偉哲17日也頒贈感謝狀，感謝慈揚基金會關心據點長輩防疫期間的餐飲需求。臺南市長黃偉哲表示，非常感謝慈揚社會福利基金會王玉飛執行長夫婦以及燦揚企業關心台南市長輩們，體現了企業「取之於社會，用之於社會」的精神，同時，他也肯定社會局與據點志工們兼顧防疫與關懷，讓長輩在據點供餐服務暫停期間，營養餐食不間斷。黃偉哲指出，燦揚企業員工是非常幸福的，很多企業的防疫假是無薪防疫假，但燦揚企業是帶薪的防疫假，當然最好家裡面都不要有人有防疫假的需求，但是一旦碰到了，燦揚企業對他們的同仁提供非常好的照顧，這點也足以作為社會企業的楷模。慈揚社會福利基金會執行長王玉飛表示，看到全國因配合防疫工作而停止據點共餐，因關心據點長輩的用餐需求，特別捐贈3萬包簡易加熱即可食用的調理料理包予臺南市政府社會局，調理料理包葷素皆有，可照顧不同餐食需求的長輩，讓長輩不用擔心防疫期間的餐飲問題。燦揚企業代表黃崇仁副總亦到場表示支持市長黃偉哲打造「高齡友善城市」的施政理念。社會局表示，3萬份料理包已由社會局料理包關懷列車分階段送達新營、柳營、鹽水、左鎮、玉井、楠西、南化、大內、麻豆、善化、北區、安南、白河、後壁、東山、永康、新市、新化等行政區，在料理包送達社區長輩手中的運送過程，為確保人員接觸安全，社會局與據點人員嚴謹落實消毒、戴口罩、頭罩、手套，以及量測體溫，並規劃據點領餐不下車的得來速服務，將人員接觸減到最低，嚴格落實防疫措施。社會局長陳榮枝表示，臺南市截至6月11日止共布建378個社區照顧關懷據點，目前全力配合中央疫情指揮中心的防疫指示，暫停各項具接觸性的活動，並針對高關懷對象持續進行電話問安工作。為了讓長輩在防疫期間中持續學習，社會局協助各據點學習如何運用line直播、視訊來與長輩互動，並提供延緩失智失能、健康促進活動、防疫、防詐騙宣導、才藝教學、衛教課程、體適能課程等多元化課程，讓長輩宅在家也能度過多采多姿的防疫生活。此外，目前社會局為南市37區85歲以上長者提供貼心接送接種疫苗服務，由區公所、里幹事或里鄰長主動通知符合施打資格之長輩就近施打疫苗，提醒長輩不用擔心，若有相關疑問可向戶籍地區公所詢問。防疫面罩、護目鏡269元起"

Output Comparison

- SL: “柯文哲祭出強硬手段 要與確診者有接觸的人”
- SL+RL: “柯文哲祭「硬清零」計畫 不給任何賠償”

title: "「確診有症狀一律住醫院」 柯文哲：死在家裡不能怪政府照顧不周"

mainText: "▲柯文哲祭出強硬手段，要與確診者有接觸的人全部做篩檢、確診者一律住防疫旅館。（圖／記者湯興漢攝）記者陳家祥／台北報導新冠肺炎疫情趨緩，台北市長柯文哲17日公布「硬清零」計畫，要求確診者不可以住在家裡，斷絕家戶感染，一但確診，有症狀住醫院，沒症狀的送防疫專責旅館，如果要留在家裡，要簽切結書，以後不給任何賠償，死在家裡不能怪政府照顧不周。柯文哲表示，當快篩的陽性比率降到1%以下，過去的社區普篩大概是效果不大，所以要有新的思維。居家隔離人數慢慢在走下坡，趨勢看起來是終於明顯下降。對於有1、2個個案藏在社區裡面，到底怎麼解決？柯文哲說，全面打疫苗是一個策略，美國打滿2劑的人已經超過50%，但昨天感染人數還有1萬，所以表示除非像以色列打到70%，不然終究還是會有感染。柯文哲說，台北市算一算疫苗才打到6%多，距離以色列的70%還有一段時間，但也不能一直封城下去，很多產業已經受不了，所以要想辦法。過去允許確診者住在家裡，「現在起通通不准，一但你是確診者，就算沒有症狀，也要去住防疫專責旅館」，經過統計，家戶感染占北市感染的43%，所以常常是家人感染給家人，然後又去外面傳染給其他人。「所以從現在起，除非有很強的理由，若一但確診，有症狀住醫院，沒症狀的送防疫專責旅館。」柯文哲強調，如果要留在家裡，要簽切結書，以後我們不給你任何賠償，死在家裡不能怪政府照顧不周，「你本來就應該去防疫專責旅館接受醫護人員的監控」。會祭出這種強硬手段，柯文哲強調，因為北市要做清零計畫，當普篩陽性率不到1%，代表普篩已經無效，要以確診者為中心做同心圓，做精準疫調找出可能感染的人。我幫你出計程車、快篩費用，但你必須去做檢查，因為不這樣的話沒辦法清零，「要讓台北要盡速清零，才能思考三級管制可不可以放鬆」。柯文哲說，第一，確診者不可以住在家裡；第二，以確診者為中心做出來的疫調，可能的被感染者，市府會出交通費跟檢查費用，但是你有義務去做篩檢，不然傳染鏈沒有辦法打斷。柯文哲說，要請市民朋友擴大配合我們要做的篩檢，預防隱形的傳染鏈，這不解決，台北市的清零計畫不會成功。"

Results

- Train from scratch:
 - SL: {rouge1: 16.3763, rouge2: 5.6533, rougeL: 16.1392}
 - RL: {rouge1: 15.2448, rouge2: 5.451, rougeL: 15.0109}
- Use RL to fine-tune the result of SL:
 - {rouge1: 17.514, rouge2: 6.0042, rougeL: 17.1649}

Reinforcement Learning can help seq2seq pretrained model summarize news or social posts and get better Rouge scores

Analysis

- Ground Truth :
 - 台南市173處社區照護暫停供餐，黃偉哲感謝璨揚企業捐贈三萬份料理包關懷長輩
- Train from scratch:
 - 台南市長黃偉哲贈3萬份料理調理包 台南市長黃偉哲
- Use RL to fine-tune the result of SL:
 - 台南市「慈揚社會福利基金會」捐贈3萬份料理調理包

Reinforcement Learning can help seq2seq pretrained model generate more informative and less repetitive summarized sentences

Reference + AI tool

- Reference Papers

- <https://arxiv.org/pdf/1705.04304.pdf>
- <https://arxiv.org/pdf/1810.06667.pdf>
- <https://arxiv.org/pdf/2009.01325.pdf>
- <https://arxiv.org/pdf/1805.09461.pdf#page=19&zoom=100,65,145>

- AI tools

- *Pretrained mT5 Transformers model from huggingface.co*
- <https://huggingface.co/transformers/>

Feedback questions

- Q : How is the pretrained model chosen in terms of the strength of the model and what other choices are there?
- A : T5 is the state of the art model works well on a variety of tasks, e.g., translation and summarization. And mT5 is multilingual T5 model, also works well on Chinese tasks.

Feedback questions

- Q : Is your summarization abstractive or extractive only?
- A : Abstractive. Because mT5 model decodes the summarization by choosing the most possible word from the vocabularies.

Feedback questions

- Q : Is the model still work when there's new term in the news?
- A : It depends on the meaning can be represented by the characters or not. For example, “新冠” is the abbreviation of “新型冠状病毒”, so it can be recognized correctly ; “莫德納” is a brand of vaccine and its meaning is unrelated to the three characters, so the language model needs lots of data to learn about it.

Feedback questions

- Q : Did you only include the information from TCDC(衛福部疾管署), or all information from the global?
- A : We collect as many news as we can from the udn.com website. As we know, it mainly contains local news and Taiwan-related news.

Feedback questions

- Q : Did you do any postprocessing or use other NLP tool result to get better extraction(Use NER result, etc)?
- A : In this case, we try different search algorithms to achieve better rouge score. For instance, we have tried beam search algorithm and top-k algorithm with different parameters.