

Forst: A Challenge to the NTCIR-15 QA Lab-PoliInfo-2 Task

Dialog Summarization system A

Hiromu Onogi^{†1}, Kiichi Kondo^{†1},
Younghun Lim^{†1}, Xinnan Shen^{†1},
Madoka Ishioroshi^{†2}, Hideyuki Shibuki^{†2},
Tatsunori Mori^{†1}, Noriko Kando^{†2†3}

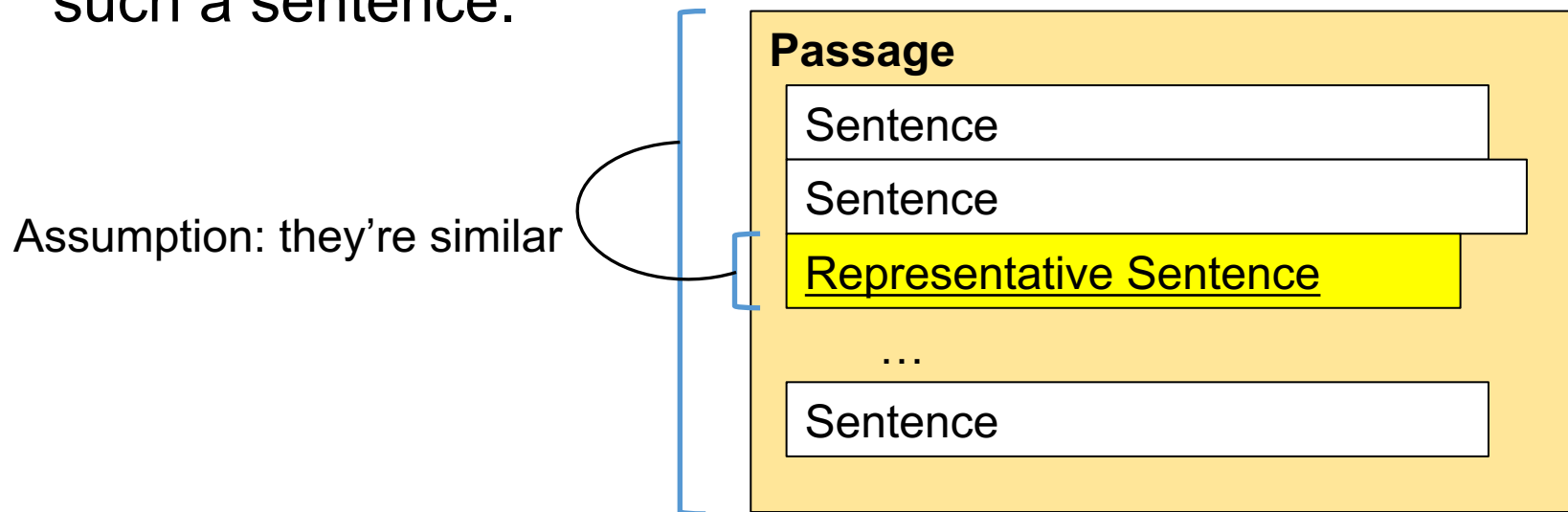
^{†1} Yokohama National University, ^{†2} National Institute of Information, ^{†3} SOKENDAI

December 10, 2020



Approach

We assume that the representative sentence of a passage is similar to the whole passage and so our system extracts such a sentence.



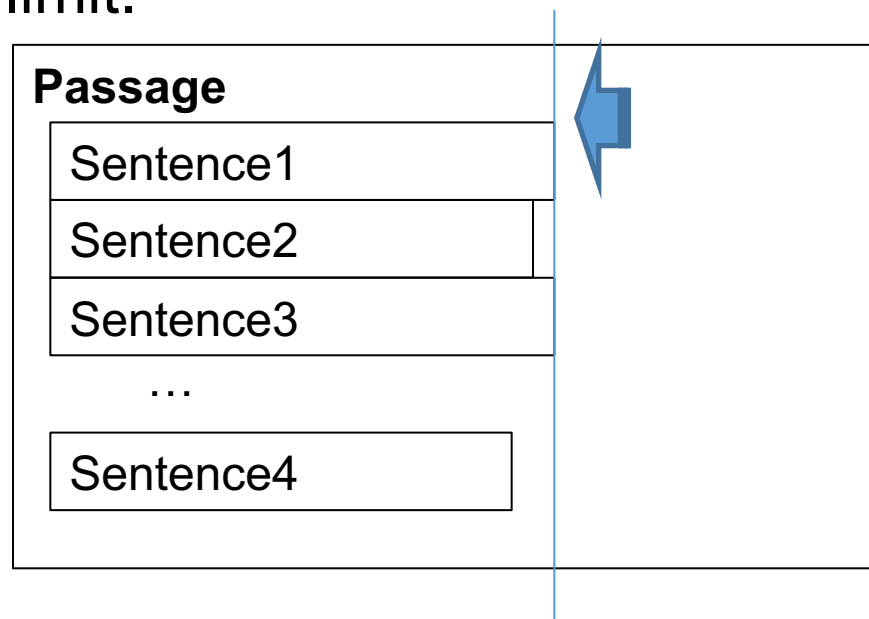
We adopted cosine similarity of distributed representation vectors. Our system uses a trained skip-gram model¹ for the distributed representation.

¹Japanese Wikipedia Entity Vector Model, Inui and Suzuki Lab, Tohoku University (2017)



Approach

The candidate extraction sentences are **pre-compressed** not to exceed the characters limit.



Our pre-compression method regards depth of *bunsetsu*-phrase in the dependency structure as the basic importance and applies MMR to avoid duplicating the content of the selected *bunsetsu*-phrase.



Related Studies

- **Kimura et al.[1]** pointed out that the use of key expressions (e.g., "～を伺います" at the end of a question sentence) is useful for extracting important strings.
- **Noguchi et al.[2]** proposed a method to estimate the importance of sentences using distributed representations of words in summarization for question sentences.

[1]Yasutomo Kimura, Satoshi Sekine, and Kentaro Inui. 2018. Towards Summarizing Local Council Proceedings. Proceedings of the Twenty-fourth Annual Meeting of the Association for Natural Language Processing (NLP2018)P5-3 (3 2018), 596–599. (in Japanese).

[2]Noguchi Masaki, Tanizuka Taichi, and Kobayashi Hayato. 2015. Summarization of Yahoo! Answers with Distributed Representation. Proceedings of the Twenty-first Annual Meeting of the Association for Natural Language Processing (NLP2015)(3 2015). (in Japanese).



Method

1. Split the passage to be summarized into sentences

2. Pre-compress each sentence (the details are described afterwards)

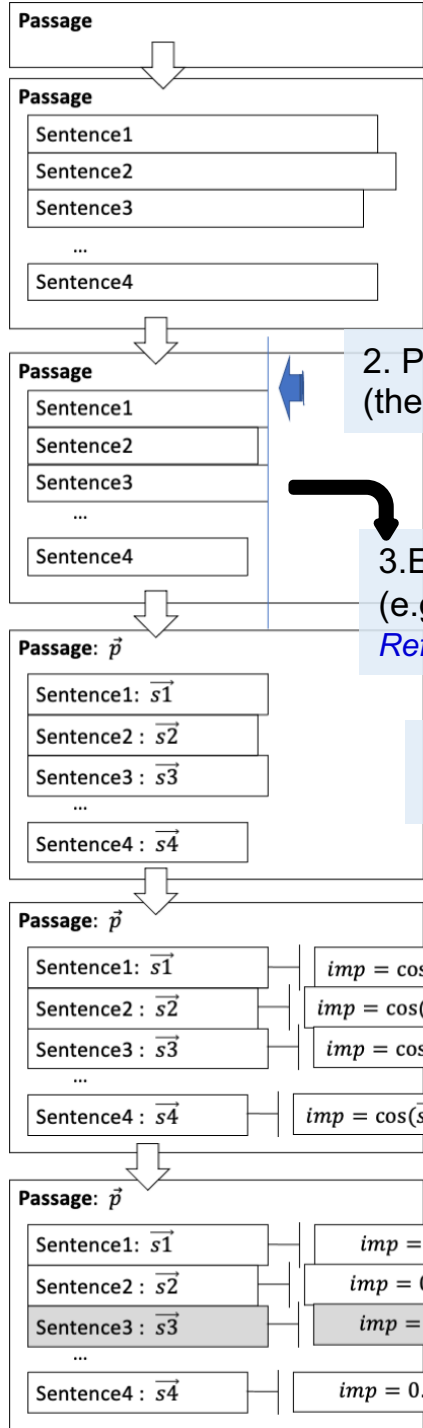
3. Extract a key sentence intending to ask or answer, if any (e.g., sentence containing '～を伺います', '～してまいります')
Ref:Kimura et al.[1]

4. Calculate distributed representation vector of each compressed sentence (obtain the word2vec value of adjectives, nouns, and verbs and take their sum).

5. Compute the cosine similarity of the distributed representation vectors* between each compressed sentence and the passage
**Ref:Noguchi et al.[2]*

7. If there is still room for the characters limit, compress sentences again to fill in the remaining characters (back to step 2).
8. Combine the extracted compressed sentences and output them as a summary.

6. Extract the most important sentence



Method(pre-compression)

このような訓練の成果を検証し、関係機関と綿密に協議しながら、対処要領を策定してまいります。(45 characters)

このような-
訓練の-
成果を-
検証し、
関係機関と-
綿密に-
協議しながら、
対処要領を-
策定して-
まいります。

Analyze the dependency structure using CaboCha

Chunk("このような") -4
Chunk("訓練の") -3
Chunk("成果を") -2
Chunk("検証し、") -1
Chunk("関係機関と") -3
Chunk("綿密に") -2
Chunk("協議しながら、") -1
Chunk("対処要領を") -2
Chunk("策定して") -1
Chunk("まいります。") 0

Assign higher importance to shallow *bunsetsu*-phrases in the dependency structure

$$Imp(C) = -1 \cdot depth(C)$$

C: *bunsetsu*-phrase

Chunk("まいります。") 0
Chunk("検証し、") -0.17420335175590065
Chunk("成果を") -0.5398998343474287
Chunk("対処要領を") -0.5592047125948513
Chunk("綿密に") -0.5806803727871166
Chunk("関係機関と") -0.7022977676902481
Chunk("訓練の") -0.7437282890012684
Chunk("策定して") -0.7881670270757923
Chunk("協議しながら、") -0.8127044042226375
Chunk("このような") -0.877280565273217

apply MMR to *bunsetsu*-phrases by word2vec value to avoid duplicating the content of the selected *bunsetsu*-phrase.

$$MMR = \arg \max_{C_i \in C \setminus S} \left\{ \lambda Imp(C_i) - (1 - \lambda) \max_{C_j \in S} \cos_{sim}(C_i, C_j) \right\}$$

We set the lambda value to 0.15.

Extract important *bunsetsu*-phrases (*)

成果を検証し、対処要領を策定してまいります。(22 characters)

(*)The following rules have been taken to avoid a grammatical breakdown

- When a *bunsetsu*-phrase containing case particles such as “~を” or “~に” is selected, the *bunsetsu*-phrase to which it links must also be extracted.
- Delete grammatically incorrect *bunsetsu*-phrases (e.g., phrases such as ‘、上で、’ and phrases that begin with formal nouns such as ‘ことを望みます’).



Result and Discussion

	Forst215(method described so far)	average of all submissions
ROUGE	0.2410	0.185
Content(X=2)	0.778	0.615
Content(X=0)	0.667	0.533
Well-formed	1.701	1.595
Non-twisted	1.044	0.823
evaluable Non-twisted (C>=1,WF>=1)	1.589	1.552
Sentence goodness	0.780	0.591
Dialog goodness	0.604	0.410

The submitted summary results were approximately 6% to 47% above the average for all evaluation categories.

Control experiments

Modification	ROUGE
Forst215	0.2410
extracting sentences without pre-compression	0.2275
not applying MMR	0.2453
not prioritizing sentences intending to ask or answer	0.1430
(average of all participants)	0.1850



Result and Discussion

Modification	ROUGE
Forst215	0.2410
not pre-compressing each sentence	0.2275
not applying MMR	0.2453
not prioritizing sentences intending to ask or answer	0.1430
(average of all participants)	0.1850

Extracting key expressions (e.g., “～を伺います”, “～してまいります”) are considered to be essential.

However, the method of simply using the cosine similarity as the importance does not consider such expressions important.



Result and Discussion

Modification	ROUGE
Forst215	0.2410
extracting sentences without pre-compression	0.2275
not applying MMR	0.2453
not prioritizing sentences intending to ask or answer	0.1430
(average of all participants)	0.1850

Sentence extraction with pre-compression improved ROUGE score but applying MMR didn't so.

However, applying MMR seemed to have made the meaning of the summary easier to understand.

Pre-compressed sentences often break down the grammar. We will consider methods such as taking into account the strength of the relationship between the *bunsetsu*-phrases.

