Analysis of Similarity Measures between Short Text for the NTCIR-12 Short Text Conversation Task

Kozo Chikai

Graduate school of information science and technology, Osaka University

Yuki Arase

Graduate school of information science and technology, Osaka University





1 <u>Introduction</u>

≻Overview of STC task and goal of our study

② Methods & System Design

- WTMF model
- Word2vec
- Random Forest and its learning method

③ <u>Experiment</u>

④ Result & Discussion

INTRODUCTION

• <u>Short Text Conversation (STC) Task</u>

✓Retrieve suitable replies from a pool of tweets

- \rightarrow System's output = Ranking of replies
- ✓System design : Input → Text (similar to input) → Output (reply)





① <u>Introduction</u>

➢Overview of STC task and Goal of our study



2 Methods & System Design

- WTMF model
- Word2vec
- Random Forest and its learning method

③ <u>Experiment</u>

④ Result & Discussion

• <u>Weighted Text Matrix Factorization (WTMF)</u>

- ✓Simply models the text-to-word information without leveraging the correlation between short text.
- ✓ X: term-document matrix (row: document vector, cell: TF-IDF value) → Factorized into two matrix such that $X \approx P^T Q$ (where $P \in \mathbb{R}^{K \times M}$, $Q \in \mathbb{R}^{K \times N}$)

$$\sum_{i,j} W_{i,j}(P_{\cdot,i} \cdot Q_{\cdot,j} - X_{i,j}) + \lambda ||P||_2^2 + \lambda ||Q||_2^2$$

$$W_{i,j} = \begin{cases} 1 & \text{(if } X_{i,j} \neq 0) \\ w_m & \text{(if } X_{i,j} = 0) \end{cases}$$



Source : Modeling Sentences in the Latent Space

• Word Embedding

✓ Generate word vectors using a neural network

✓Liner operation between words is possible

(ex: "king"-"man"+"woman" = "queen")

- ✓Adopted the implementation known as Word2vec tool
 - → Output: low-dimensional vector representing words



Source : Efficient estimation of word representations in vector space

• Learning Post-Comment Pairs using Random Forest



• Combination of Random Forest and TF-IDF



DESIGN OF CONVERSATION SYSTEM



AGENDA

① Introduction

➢Overview of STC task and Goal of our study

② Methods & System Design

- WTMF model
- Word2vec
- Random Forest and its learning method



④ Result & Discussion

EXPERIMENT

<u>Setting</u>

✓Goal: Evaluate a ranked list of potential replies to an input tweet

✓10 annotators assigned a score to each reply (tweet)

- \rightarrow score: +2, +1, and 0
- \rightarrow The larger score the tweet has, the better reply it is!
- ✓ Evaluation criteria: nDCG@1, nERR@5, and Accuracy
 - \rightarrow Each method has a value between 0.0 and 1.0 using these criteria
 - \rightarrow The larger score is better

Name of method	Technique
Oni-J-R1	④Random Forest → TF-IDF
Oni-J-R2	⑤Random Forest + TF-IDF
Oni-J-R3	①TF-IDF
Oni-J-R4	③Word2Vec → TF-IDF
Oni-J-R5	②Weighted Text Matrix Factorization

OVERVIEW OF THIS PRESENTATION

① <u>Introduction</u>

≻Overview of STC task and Goal of our study

② Methods & System Design

- WTMF model
- Word2vec
- Random Forest and its learning method

③ <u>Experiment</u>

Result & Discussion

RESULT



DISCUSSION



Figure 8: Mean of labels for each rank

DISCUSSION

input	$0 \le \#$ of cha	aracters ≤ 20	$20 < \#$ of characters ≤ 40		$40 < \# \text{ of characters} \le 60$		60 < # of characters	
system	R1	R3	R1	R3	R1	R3	R1	R3
rank1	0.810	0.587	0.424	0.408	0.376	0.279	0.350	0.350
rank2	0.860	0.793	0.722	0.709	0.703	0.685	0.691	0.644
rank3	0.503	0.550	0.401	0.328	0.371	0.365	0.375	0.341
rank4	0.583	0.670	0.731	0.749	0.694	0.732	0.775	0.806
rank5	0.227	0.600	0.359	0.417	0.471	0.356	0.397	0.294

Table 2: Effect of R1 and R3 due to the number of characters

(Oni-J-R1: 'Random Forest \rightarrow TF-IDF', Oni-J-R3's method: 'TF-IDF')

DISCUSSION(2)

	フォローありがとうございます!
Input post	Thank you for following!
	こちらこそ、懐かしい画像をありがとうございました♪ よろしくお願いします!
rank1	Thank you too for sharing the nostalgic photos! Keep in touch.
	RT ありがとです♪ フォローさせていただきました!
rank2	Thank you for retweeting. I had followed you!
	お休みなさい♪ 今年も宜しくお願いします!
rank3	Good night. Wish we'll have fun together!
	寝ます。おやすみなさい。 さっそく今日から仕事だ!
rank4	Going to bed. Good night. Gotta work from today!
	んう いくう
rank5	hmmm

(Oni-J-R1's method is 'Random Forest \rightarrow TF-IDF')

DISCUSSION(3)

Table 4: Example of output from Oni-J-R3 and Oni-J-R4

Input post	LINE 使えんので用のある人はこっちでお願いします!					
	Cannot use LINE now. Please contact me via Twitter!					
Method	R3	R4				
rank1	新しいアカウント作ったから教えてー!	PC で扱うようなでかい解像度で色々したりするんでなければいよ				
		いよ十分な感じなんです?				
	Let me know your account! I've just created one.	It should be satisfactory unless you process high-resolution im-				
		ages that require PC-like performance.				
rank2	LINE 使えるようになった?	んー、パソコンで使えた機能の大半はもう使えてるよね。 単純に使				
		える、だけだけど。 次の 8 コアのスマートフォン向けの cpu が				
		出るからそれで単純に使える機能から卒業して、もうスマートフォ				
		ンで十分になる機能は出るだろう				
	Has LINE recovered?	Most functions in PCs are already covered by smart phones,				
		although they are just functioning They should be actually				
		usable with the next 8-core CPU for smart phones.				
rank3	パソコンの LINE ツンデレすぎるよ! もう! 早く携帯よとどけ!	単発でやる環境に長くいたせいで、最近キャンペ参加できるように				
		なったのですがなかなか上手い感じにキャンペ用キャラが組めない				
		んですよね				
	LINE app for PC is hard to use Get me a cell phone!	I cannot figure out how to organize my party for online campaign				
		as I've been playing alone				

(Oni-J-R3's method: 'TF-IDF', Oni-J-R4: 'Word2vec \rightarrow TF-IDF')

SUMMARY

- In this study, we compare the conventional methods to handle short text.
- We use unsupervised methods to generate vectors
- We also use supervised methods to learn if a pair of tweets can be a postcomment pair.
- As the result of the formal run in STC, 1 Random Forest \rightarrow TF-IDF outperformed other methods.
- The method using Word2vec shows interesting results in some context, however there is room for improvement in this method.

REFERENCES

- <u>Modeling Sentences in the Latent Space</u> W. Guo, and M. Diab In Proceeding of ACL (2012)
- <u>Distributed Representations of Sentences and Documents</u> T. Mikolov, Q. V. Le In Proceedings of ICML(2014)
- Overview of the NTCIR-12 short text conversation task
 L. Shang, T. Sakai, Z. Lu, H. Li, R. Higashinaka, and Y. Miyao
 In Proceedings of the NTCIR Workshop Meeting on Evaluation of
 Information Access Technologies(2016)

$\mathsf{APPENDIX}(1)$

• The way of Word2Vec \rightarrow TF-IDF

- 1. Extract nouns and verbs in an input post
- 2. Extract 20 most similar words for each noun and verb using word2vec
- 3. Search tweets containing these words in the corpus
- 4. Generate the reply list from these tweets in the same manner with 1TF-IDF



APPENDIX (2)

- <u>Method of Generating the reply list</u>
 - ✓We use the relationship of a pair of post-comment in Twitter.
 - ✓Why is the similar tweet itself included in the reply list?
 - \rightarrow Because the similar tweet can be useful as a reply in some context.

For example…

Post : It is pretty cool in Hokkaido today. Comment : Summer is the best season in Hokkaido.

✓In this case, the direction of post-comment can be reversed. Thus each tweet can be a reply for an input post.