| Title | |
|---|---|
| Author(s) | , |
| Citation | |
| Issue Date | 2000-03 |
| Type | Thesis or Dissertation |
| Text version | author |
| URL | http://hdl.handle.net/10119/665 |
| Rights | |
| Description | Supervisor: , , |

# A Study on the Algorithms for Dialogue Segmentation

Mikyoeng Kim

School of Knowledge Science,
Japan Advanced Institute of Science and Technology
March 2000

This paper proposed a new algorithm for segmenting dialogues with high reliability based on the spoken dialogue corpus. The resultant dialogue segments are analyzed in terms of the concept of "initiative" to examine their meaning. The dialogue data used in this study are task-oriented dialogues with a variety of the tasks such as the Map Task and hotel reservation, collected and annotated by the Special Interest Group of Corpus-Based Research for Discourse and Dialogue of Japanese Society for Artificial Intelligence. There are 14 dialogues in the corpus, which amount to approximately 52.7 minutes in total.

Human-to-human dialogues seem to consist of fragmentary utterances, which is partially confirmed by the fact that the length of the utterances (or sentences) is shorter in spoken dialogues than in written texts. However, examining the flow of the utterances participant by participant, in many cases, each participant's utterances have found to be almost full-fledged sentences, segmented by the non-content words or phrases such as "Yes", "No" and "Okay". Based on this observation, a new dialogue segmentation algorithm was proposed (algorithm 1), by which the non-content words or phrases are considered to be segment boundaries. The algorithm 1 was evaluated by the recall and the precision rates: the recall rate was very high (98.3%) but precision rate was lower than expected (68.7%). The mis-recognized segment boundaries were examined for improving the algorithm 1, which can be classified into two types: the repetition of the words or phrases and the adjacency pairs such as question-answer and request-acceptance.

The new algorithm 2 was proposed to alleviate the problems of the algorithm 1. The

1

segment boundaries obtained by the algorithm 1 are excluded if they are of the types of the repetition of the words and phrases and the adjacency pairs. The evaluation result of the algorithm 2 shows that the precision rate was improved by 24.7% (to 93.4%), while maintaining the same level of the recall rate as the algorithm 1. In order to confirm the validity of this algorithm 2, the algorithm using the cue words (fillers, discourse markers and conjunctions) that was proposed in the previous research was evaluated based on the same data, which resulted in 50.2% recall rate and 40.4% precision rate. This shows that the dialogue segmentation with high precision is not an easy task and that the new algorithm 2 performed quite well compared to the previous representative algorithm.

In order to examine the semantic aspect of the dialogue segments, the segments obtained were analyzed in terms of the concept of "initiative". Although the proposed algorithm does not decide on the speaker with the initiative, the dialogue segments are good candidates for examining the initiative, because the algorithm encompassed the substantial utterances segmented by the non-content words and phrases. Currently the definitions of the initiative are different from theory to theory. For example, some takes the utterances of proposing domain goals as the initiative taking utterances while the other takes those of proposing problem solutions.

In this paper, the dialogue segments obtained by the proposed algorithm are examined based on the four definitions of the initiative proposed by Cohen et al. (1998). Especially, the segments consisting of the substantial utterances frequently observed in the question-answering pairs were focused, because the differences of the definitions of the initiative are expected to be salient. Which definition of the initiative is advantageous definitely depends on the purpose of the research. The proposed algorithm does not decide on the initiative. However, it enables us to focus on the segments in which the initiative is difficult to judge, which contributes to, for example, building dialogue corpus annotated with the information of the initiative in terms of the cost and the reliability.