

1 **Appendix 2**

2

3 Python 3.7.2, NumPy 1.17.2, SciPy 1.2.1, Pandas 0.24.2, Matplotlib 3.0.3, NetworkX 2.3,
4 Apyori 1.1.1, Community 1.0.0, R 4.0.2, and exactRankTests packages 0.8–21 were used
5 for the analysis.

6

7 **1. How to determine the number of tasks and crowd workers of Step 1-b.**

8 The decomposed medical notes of each type were randomly arranged into 130 tasks. The
9 details of 130 tasks are as follows; one text was shown to each worker in each of the 43
10 tasks. Two texts were shown to each worker in each of the 44 tasks. Three texts were
11 shown to each worker in each of the 43 tasks.

12

13 The details of the definitions of the number of tasks or crowd worker are as follows:

14 In all, the tasks add up to 130 ($43 + 44 + 43$) tasks, and the total number of notes,
15 equaled 260 ($43 \times 1 + 44 \times 2 + 43 \times 3$).

16 $(35 \text{ workers}) \times 20 \text{ tasks} / 130 \text{ tasks} = 5.38 \text{ workers}$, therefore, approximately five
17 workers, on average, completed each task.

18

19

20 **2. How to determine the number of tasks and crowd workers of Step 2-a.**

21 The number of notes presented to a crowd worker for each task was randomly determined
22 as 1, 3, 5, 7, or 10, and the same number of decomposed dummy and real progress notes
23 were presented to each worker. For example, in Fig 3, three dummy and three real notes

were shown to the worker. The notes shown were randomly selected. In each job, each worker answered 50 yes/no questions, which were randomly arranged. Therefore, the total number of tasks amounted to 50 tasks, and 260 real progress notes were presented once to each worker

The details of the definitions of the number of tasks or crowd worker are as follows:

The total number of tasks amounted to $5 \text{ tasks} \times 10 = 50 \text{ tasks}$, and $(1 + 3 + 5 + 7 + 10) \times 10 = 260$ real progress notes were presented once to each worker.

We note that 50 tasks consist of five types of 10 tasks in each of the following cases: 1 real progress note was presented, 3 real progress notes were presented, 5 real progress notes were presented, 7 real progress notes were presented, and 10 real progress notes were presented.

3. The detailed protocols to omit the duplicate yes/no questions we obtained of Step 2-a.

The postprocessing rules for candidate features generated by crowd workers in Step 2-a are as follows:

Features that did not satisfy the yes/no questions were excluded.

(e.g., “no difference,” “___ is different.”)

The symbols for bullets and blank characters at the beginning of the yes/no questions were deleted.

47

48 The duplicate descriptions were omitted. We also excluded features that were the same as

49 other features except for punctuation.

50 (e.g., “The patient’s suffering is recorded.” and “The patient’s suffering is recorded” were

51 considered as the same, and the latter was excluded.)

52

53 "Whether or not ____." " Whether ____." " ____?" → " ____"

54 (e.g., "Whether or not the documents use the question form." "Whether the documents

55 use the question form." "Do the documents use the question form?"

56 => "The documents use question form.")

57

58 "There are documents that ____" → "The documents ____"

59 (e.g., "There are documents that consist only of numbers and alphabets."

60 => "The documents consist of numbers and alphabets.")

61

62 "Document group A is ____" "Document group A is more ____" → "Document is ____"

63 (e.g., "Document group A is detailed." "Document group A is more detailed."

64 => "Documents are detailed.")

65

66 “Are ____described?” => “____ are described.”

67 (e.g., "Are the numbers described?" → "The numbers are described.")

68

69

4. How to determine the number of crowd workers of Step 2-b.

$(25 \text{ progress notes} \times 150 \text{ features} \times 860 \text{ workers}) / (100 \text{ progress notes} \times 3,197 \text{ features})$
= 10.08 workers. Therefore, approximately 10 workers checked each progress note and feature pair.

5. The method for determining the lift threshold in Step 2-b.

The candidate features that characterized few progress notes were considered inappropriate for this study, because such candidate features were considered to express the specific contents of a particular patient (e.g., “Are the details of colon cancer, the history of pneumonia, and an egg allergy described?”). Based on these factors, when the lift threshold was 1, 3, 5, 7, or 9 for the pairs of progress notes and candidate features, it represented the maximum value of the lift for which more than half of the candidate features characterized five or more progress notes. The pair of the lift score under the lift threshold would be excluded.

As results, the pairs of real progress notes and candidate features with an appearance count of one or two were excluded, and 28,487 pairs of real progress notes and candidate features were selected. When the lift threshold was set at 1, 3, 5, 7, and 9, the proportions of candidate features with five or more correlated progress notes among all of those selected were 65.62%, 64.85%, 51.47%, 38.91%, and 28.11%, respectively.

6. The method for determining the Jaccard coefficient threshold in Step 2-c.

To define a better network, the Jaccard coefficient was increased from 0.1 to 0.9 in increments of 0.1, and community detection was performed using the Louvain algorithm for each network defined with each Jaccard coefficient. As a result, the determination of the Jaccard coefficient threshold was based on the average number of candidate features in the obtained community. If this number is too large, the obtained communities might include candidate features that are minimally related to each other. This situation would be inappropriate for our purposes, because we wanted to create communities including similar candidate features.

When the Jaccard coefficient was increased from 0.1 to 0.9 (in increments of 0.1) for each community's detection, the number of extracted communities and the average number of candidate features included in each community also showed progressive increments. For Jaccard coefficient values of 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, and 0.9, the corresponding numbers of the extracted communities were 8, 12, 16, 31, 176, 176, 184, 184, and 174, respectively, and the corresponding average numbers of candidate features in each community were 269.1, 171.3, 125.3, 59.8, 8.07, 7.43, 5.77, 4.98, and 4.67. Using this result, we assumed that, when the Jaccard coefficient was 0.5 or more, the possibility of the obtained communities including candidate features that had a minimal relation to them would be reduced compared to when the Jaccard coefficient was 0.4 or less.

7. How to determine the number of crowd workers of Step 2-d.

The details of the definitions of the number of tasks or crowd worker are as follows:

(105 workers) \times 50 communities/176 communities = 29.8 workers; therefore,
approximately 30 workers, on average, worked on each community.

Among 176 communities, seven communities presented with the same number of votes,
and one feature was selected randomly.

8. How to determine the number of crowd workers of Step 3.

The details of the definitions of the number of tasks or crowd worker are as follows:

(500 workers) \times 10 progress notes/100 progress notes = 50 workers, therefore,
approximately 50 workers, on average, worked on each crowd progress note.

9. The reason of the criteria in Step 3.

Long notes inevitably contain a considerable amount of content; that is, they may meet
many crowd features known to be appropriate for medical notes. Therefore, we did not
use the crowd features that were appropriate as a feature of medical notes as an
evaluation metric; instead, we used those that were inappropriate as a feature of medical
notes as an evaluation metric.