

Application of Data Mining and Text Mining to the Analysis of Medical near Miss Cases

Masaomi Kimura¹, Sho Watabe¹, Toshiharu Hayasaka¹, Kouji Tatsuno¹, Yuta Takahashi¹, Tetsuro Aoto¹, Michiko Ohkura¹ and Fumito Tsuchiya²

¹*Shibaura Institute of Technology,*

²*Tokyo Medical and Dental University*

Japan

1. Introduction

Not only the side effects of medicines themselves, but also their abuse, namely the lack of safety in drug usage, can cause serious medical accidents. The latter applies to the case of the mix-up of medicines, double dose or insufficient dose. Medical equipments can also cause accidents because of wrong treatment, such as wrong input to equipments and wrong power-off. In order to avoid such accidents, it is necessary to investigate past cases to identify their causes and work out counter measures.

Medical near-miss cases caused by wrong treatment with the medicines or the medical equipments are strongly related to medical accidents that occur due to the lack in safety of usage. Medical near-miss cases are incidents, which could be medical accidents avoided owing to certain factors, and happen more frequently than medical accidents. Incorporating Heinrich's law, which shows the tendency of frequency and seriousness of industrial accidents, we estimate that near-miss cases happen three hundred times per one serious medical accident or thirty minor accidents. This can be interpreted as there being many causes of medical accidents, most of which are eliminated by certain suppression factors, which lead to near-miss cases. The rest of the causes lead to medical accidents. From this perspective, we can expect that both medical accidents and near-miss cases originate from the same type of causes, which suggests that the analysis of data on near-miss cases is valid to investigate the cause of medical accidents, since their occurrence frequency is much larger than that of medical accidents.

For the reasons stated above, we analyze the data of medical near-miss cases related to drugs and medical equipments, which have been collected in previous years to determine the root cause of medical accidents caused by the neglect of safety of usage. Though simple aggregation calculations and descriptive statistics have already been applied to them, the analyses are too simple to extract sufficient information such as pairs of medicines that tend to be confused, and the relationships between the contents of incidents and the causes. To realize such analyses, we utilize data mining techniques such as decision-tree and market-basket analysis, and text-mining techniques such as the word linking method.

The related works analyzing medical data by utilizing natural language processing or machine learning were introduced by Hripcsak et al. (Hripcsak et al., 2003), who suggested the framework to detect events such as medical errors or adverse outcome. Recently,

Tsumoto(Tsumoto& Hirano, 2007) collected incident reports independently of a national project and applied decision tree algorithm, whose results show that errors caused by nurses depend on the part of their working hour and that an uncooperative patient tends to diminish nurses' power of attention and causes medication errors.

We have applied data-mining/text-mining approaches to the nation-wide incident reports collected by Japanese government, which is focused on the use of medicines or medical equipments and show the obtained results (Hayasaka et al., 2006; Hayasaka et al., 2007; Kimura et al., 2007; Takahashi et al., 2004; Takahashi et al., 2005; Tatsuno et al., 2005; Tatsuno et al., 2006; Watabe et al., 2007; Watabe et al., 2008). We introduce the results in this paper.

2. Target data and tools

Our target data are the medical near-miss cases related to medication and the use of medical equipments collected by the Japan Council for Quality Health Care, which is an extra-departmental body of the Japanese Ministry of Health, Labor and Welfare.

As for medication, we analyzed 1341 records included in the results from the 1st to the 11th investigations and 858 records in those from the 12th to the 14th. Since some data items in the latter investigations are added to the former ones, as is shown in Table 1 and Table 2, if we use such added items, we restrict the records having the data items, and if not, we use all records in the target data.

As for the use of medical equipments, we analyzed 500 records which are obtained from the investigations conducted at the same time of the investigations about medication. The data items used for investigations for medical equipments are similar to those for medication and increase as the investigations proceed, namely, five items for the first investigation and 26 items for the 14th investigation.

In fact, there were some problems with analyzing the data using a computer program. We introduce some of them as follows:

- Data with many levels of abstractness were contained in one data item. For example, though 'misconception' is a concept which should be included in 'human error', this was adopted as a value in the data item 'major cause of the case'. Because of this, we had to redefine the categories and reclassify the data.
- In the Japanese language environment, there are several ways to express English letters and numbers such as single-byte letters (ASCII), double-byte English letters and Japanese Kana. This causes the diversity of expression. For instance, we can express 0.001g not only as '1mg', but also as '1 m g' or '1 ミリグラム' (which stand for 1 milligram in Japanese Kana).
- The diversity of drug name expression is also caused by the ambiguity when adding the medicine information, such as dosage form. For example, both 'アダラート' (Adalat) and 'アダラート錠' (Adalat tablet) are used to denote the name of the medicine which is administered in near-miss cases.

The diverse expression has to be standardized to ensure correct analysis. Though it is ideal to control the input data by designing the entry user interface appropriately, since it is difficult to realize such control, we standardized the notation of the resultant data before the application of the data/text-mining method.

The standard unit of contents or density consists of a numerical part and a unit part. In the case of standard unit error, we can know how many times the patient (almost) received an overdose of the medicine from the ratio of numerical parts of the wrong standard unit to

that of the correct unit. Since the target data possessed the standard unit of each medicine in one data item, we had to separate it into two parts, respectively.

There also exist many vacancies in the data, which we can fill if we figure out what is referring to here by reading other data items, such as free-description data.

Applying text-mining to the free-description data in the data items such as 'Contents of the incident', 'Background and cause of the incident' and '(Candidates of) counter measures' required the deletion of characters, such as symbols and unnecessary line feed characters, and the standardization of synonyms.

In order to analyze the data, we used Clementine, which is data-mining software released by SPSS Inc., and its text-mining plug-in software, Text Mining for Clementine.

Major cause	Discussed cause	Name of wrong drug	Dosage form of wrong drug
Medical benefit of wrong drug	Name of right drug	Dosage form of right drug	Medical benefit of right drug
Content of incident	Opinion	Remarks	-

Table 1. Data items in 1st to 11th investigations.

Day of the week	Week day or holiday	Time	Place
Department	Content of incident	Psychosomatic state of the patient	Job title
Experience (year)	Experience (month)	Affiliation (year)	Affiliation (month)
Medical benefit class	Nonproprietary name	Name of wrong drug	Dosage form of wrong drug
Effect of wrong drug	Name of right drug	Dosage form of right drug	Medical benefit of right drug
Discussed cause	Concrete contents of the incident	Background/cause of the incidents	Candidates of counter measure
Comment	-	-	-

Table 2. Data items in 12th to 14th investigations.

3. A brief review of data/text-mining techniques used in this study

We mainly utilized the market-basket analysis technique and decision-tree algorithm to extract the relationships between data and the rules to be followed in them. Market-basket analysis originally identifies combinations of goods bought together in a store and generates rules of the purchase tendency, which tell us which goods customers will also buy if some item is put in his/her basket. By applying this method to the values in the multiple data items, we obtain the information on the relations between the values. There are several algorithms supporting the analysis, among which we employ the Apriori algorithm and Web graph.

The decision-tree algorithm identifies the conditions dividing the data into groups by allowing minimization of entropy or the Gini index, and provides a tree that shows the optimal division of data to classify them. As the decision-tree algorithm, we use C5.0, which is suitable for classifying the categorical data.

In order to perform text-mining, we utilized the Word-linking method (Kimura et al., 2005). It is common to analyze textual data based on the frequency of words obtained by morphological analysis applied to the text. Morphological analysis extracts terms (morphemes) in the text, but loses the relationship information between the terms, which forces us to read the original textual data to interpret the result.

To avoid such inconvenience, it is useful to employ dependency structure analysis, which allows us to obtain the relationship between words in a sentence. If there are many sentences whose contents and structures are similar, we can expect that some particular patterns of word-dependencies will emerge and that such sentences can be reconstructed by rearranging these dependencies correctly. This is the basic idea behind the Word-linking method.

The steps of this method are as follows:

- Derive the dependency relationships between words (morphemes) from each sentence using dependency structure analysis. Let W denote the word which depends on another word W' .
- Create a Web graph of the co-occurrence relation between W and W' , where a link is provided between W and W' . We rotate the link to allow W to sit to the left of the link and W' to the right.
- If W' of a link coincides with W of another link, we connect the links by placing W' and W in the same place.
- Read out the sentences from the rearranged links.

Figure 1 illustrates the results obtained with this method. There are links such as the one connecting '経口薬 (oral drug)' to '数 (number)', which indicates that the dependency '経口薬の数 (the number of oral drugs)' occurs many times. Connecting the links shows us the sentences '経口薬の数が減る (the number of oral drugs can be reduced.)', '経口薬が減る (the oral drugs can be reduced.)', and '経口薬が多い時 (when the patients have to take many medicines.)'.

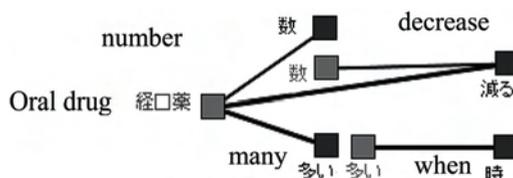


Fig. 1. Word-linking method (example)

4. Results on the analysis on medication

4.1 Relationships between the contents of near-miss cases and their major causes

Incidents occur in various service phases. From the viewpoint of prevention of near-miss cases (and consequently medical accidents), it is important to assess the relationships among the contents of incidents, the service phase in which the incidents take place and their major cause. We therefore applied decision-tree analysis to determine the rules of the reason why the incidents happen, by assigning the data items 'the service phase when the incident occurs' and 'the content of the incident' to explanatory variables, and 'the major cause' to an objective variable.

Figure 2 shows the resultant decision tree. This indicates that the major causes are classified mainly by the contents of incidents, namely medicine error or not.

In the case of medicine error, the major causes mainly consist of resemblance of name and/or external form. This suggests that the source of the problem is the confusing name or shape of the medicines (including their packaging). Moreover, the cases of medicine errors can be classified by the service phases, that is to say, preparation, administration of drug and others (prescription, dispensing and after medication). The result tells us that the major cause is resemblance of external form in the preparation phase, carelessness in the administration phase and name resemblance in the other phase. This suggests that the major cause of medicine errors is different depending on its service phase.

On the other hand, cases other than medicine error mainly stem from carelessness and misconception. The decision tree also shows that they can be classified in more detail, and states that the cases of an error in amount (quantity and/or standard unit) particularly originate in carelessness and misconception and that the errors of route of administration and un-administration of drug are mainly caused by communication errors and unpatency of infusion bags in addition to carelessness. This indicates that, though double checks will have a beneficial effect on the errors in amount, improvement of communication will also be effective to prevent errors related to administration.

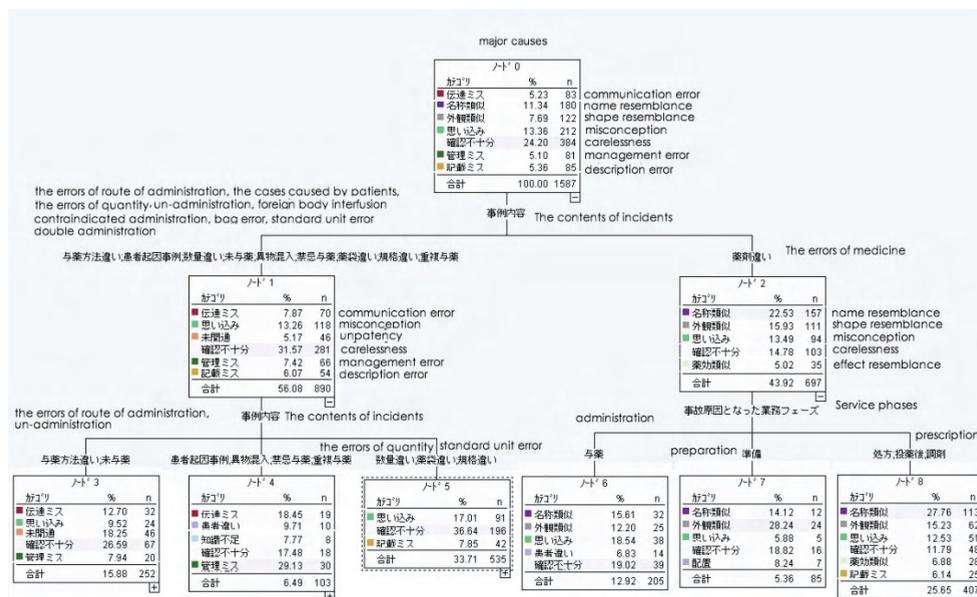


Fig. 2. Decision Tree (major causes versus the contents of incidents and service phases)

4.2 Relationship between Service phases, contents of incidents and Oversight

When a near-miss case occurs, it is crucial to detect the errors by confirmation to prevent it from becoming a medical accident. Moreover, it is also important to determine the rule of occurrences of oversights of errors depending on the circumstance and situation in order to improve counter measures. We, therefore, applied the Apriori algorithm to determine the rules of oversights depending on service phases and contents of incidents. We defined the occurrence of oversight as the difference in the occurrence phase of error and its finding

phase. Table 3 shows the result of the analysis, which indicates that an oversight tends to occur:

- if an error happens in the administration phase,
- if a quantity or standard unit error happens,
- if a medicine error happens in the administration phase,
- if the case of un-administration happens.

Support is the ratio of records for which the rule holds. This suggests to us that medical experts have to pay attention at the administration phase and/or with errors related to amount, whose rules have high support value.

Result	Prerequisite	Support	Confidence	Lift
Oversight='yes'	Occurrence phase='administration'	49.052	98.663	1.055
Oversight='yes'	Content of incident='quantity error'	18.732	98.054	1.049
Oversight='yes'	Content of incident='medicine error' and Occurrence phase= 'administration'	13.732	96.825	1.035
Oversight='yes'	Content of incident='un- administration'	12.682	96.552	1.032
Oversight='yes'	Content of incident='standard unit error'	10.714	95.238	1.018

Table 3. Rules obtained by the Apriori algorithm.

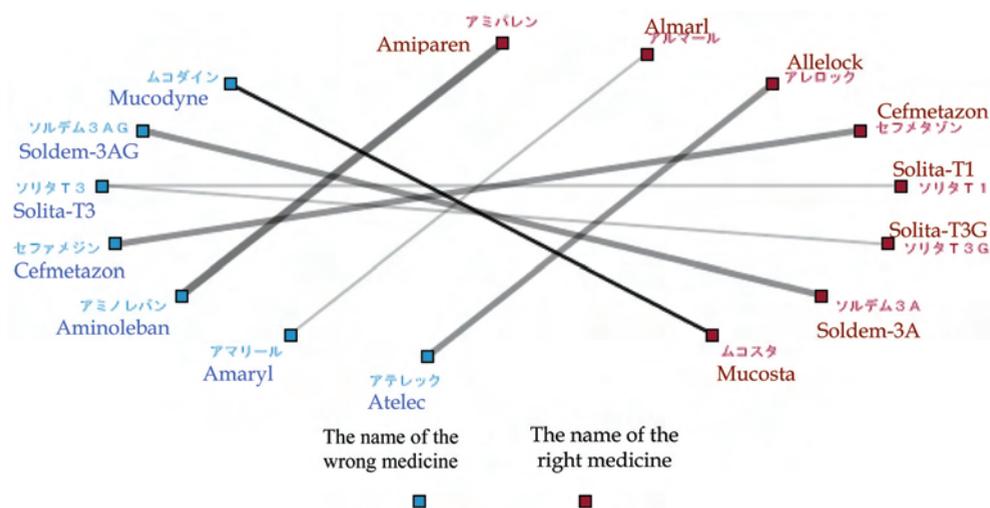


Fig. 3. Combinations of medicines mixed-up by name resemblance.

Shape	Ratio	Degree
similar	83.6%	209
different	16.4%	41

Table 4. Similarity in dosage form of pairs of medicines mixed-up by name resemblance

4.3 The combination of mixed-up medicines

In Section 4.1, we showed that the cause of medicine errors stems from the resemblance of their name. In this section, we identify which medicines are mixed-up because of the name resemblance by means of a Web graph.

Figure 3 shows the pairs of medicines appearing as cases of name resemblance more than five times. We can see that there are two types of name resemblance, one of which is the similar line of letters and the other of which is the same name apart from the last symbol (usually English letter or number). Since most of the paired medicines unfortunately have similar dosage form (both are tablets, etc.) as is listed in Table 2, pharmaceutical companies should keep from naming a medicine with a similar name to one of the existing medicines whose dosage form is also similar, and medical experts should be cautious when they treat people using those medicines.

We propose a similarity index of names of medicines making use of the String Subsequence Kernel, to realize the objective standard of name resemblance which medical experts (especially pharmacists) feels. Using the index is expected to help pharmaceutical companies seeking medicines with similar names.

4.4 Analysis of free description data on the background/cause of the incidents

We applied the Word-linking method to the field 'background/cause of the incidents' in the 12th to 14th investigations, to determine the concrete information on the cause of incidents. Compared with the formatted field 'major cause' in the data, which we deal with in Section 4.1, the field 'The background/cause of the incidents' contains free description data, which gives us more information than the formatted field.

We applied the method by occupation to determine the difference in backgrounds and causes of incidents depending on the job title. Figures 4 and 5 show the result of nurses' and pharmacists' comments, respectively. Though both figures contain the common opinions 'the problem of the checking system of the protocol and the rule' (A) and 'confirmation is insufficient' (B), nurses point out 'the systematic problem of communication' (C) and pharmacists 'the problem of adoption of medicines' (C'). We can see that, though B comes from individual fault, A, C and C' are systematic problems.

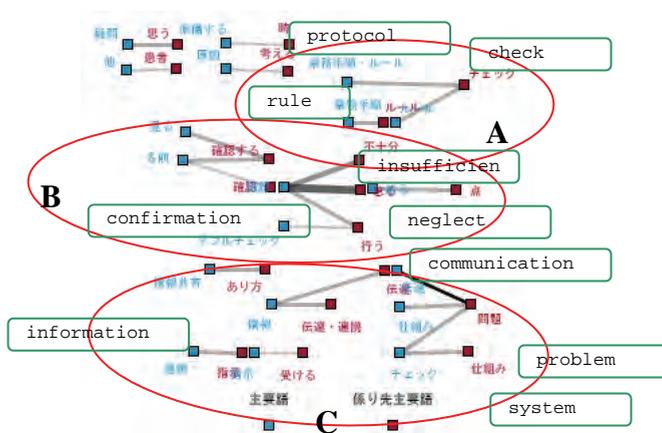


Fig. 4. Backgrounds and causes of incidents (nurse)

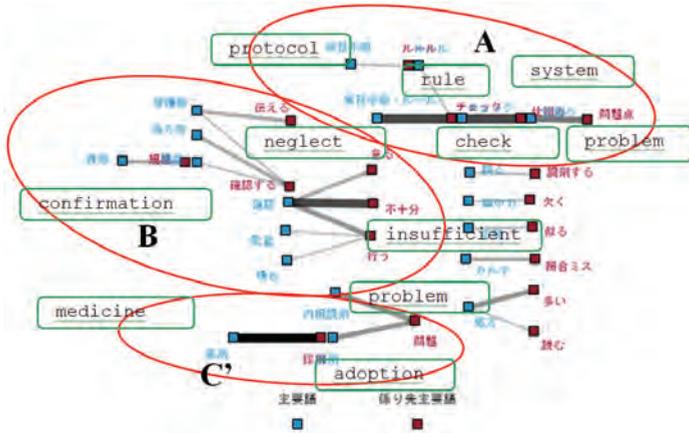


Fig. 5. Backgrounds and causes of incidents (pharmacist)

4.5 Analysis of free description regarding counter measures

We applied the Word-linking method to the field '(Candidates of) counter measures' to summarize the nurses' and the pharmacists' opinions about the counter measures to prevent the incidents. Figure 6 is the summary of the counter measures described by nurses, and suggests that there are many opinions stating '(it is necessary to) instruct to confirm and check', 'make a speech' and 'ensure confirmation'. Figure 7 shows the summary of the counter measures proposed by pharmacists. This says that, besides the confirmation and audit, it is also necessary to invite (pharmacists') attention and to devise ways of displaying medicines such as labels.

Compared with the results in Section 4.4, except for the pharmacists' opinion about the device of labels, there are few opinions on the counter measures related to the system of the medical scenarios pointed out in Section 4.4. This suggests that the medical experts such as nurses and pharmacists tend to try to find the solutions to problems within themselves. To solve the structural problem in medical situations, it is important not only to promote the effort of each medical expert, but also to take measures to improve the organization to which they belong. It is also desirable for them to be aware of the importance of organizational innovation, and to propose measures against the systematic error.

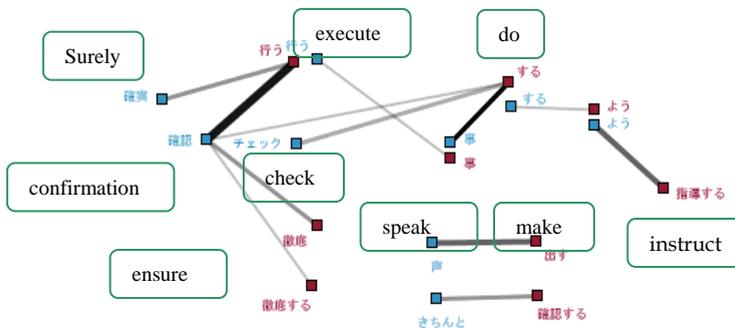


Fig. 6. Counter measures suggested by nurses (with links appearing more than 5 times)

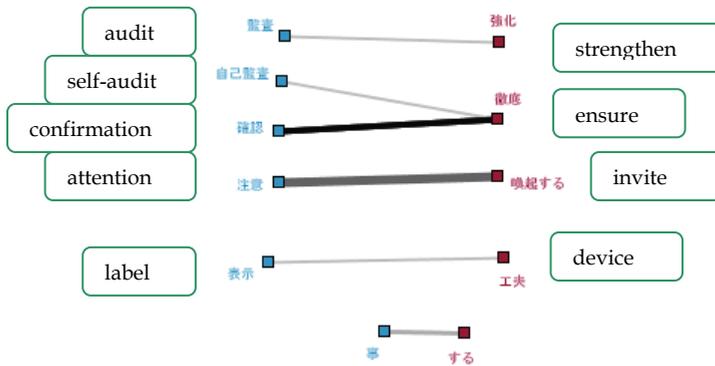


Fig. 7. Counter measure suggested by pharmacists (with links appearing more than 4 times)

5. Results on the analysis on the usage of medical equipments

In order to find the pattern of the relations between the causes of incidents, we visualized the co-occurring relations by use of Web graph(Fig. 8). This shows that the hub node of the graph is 'misuse' and that there is co-occurrence of 'insufficient maintenance', 'failure' and 'malfunction' to no small extent. We can see two groups in Fig. 8, one of which denotes the group of 'misuse', 'inadequate knowledge', 'plural standards', 'hard to handle' and 'an error in connection' related to misuse of operators, and the other of which is the group of 'insufficient maintenance', 'failure' and 'malfunction' related to issues of maintenance of equipments. As to derive these groups from data, we applied TwoStep clustering algorithm to the data and found two clusters (Fig.9), which are clusters corresponding to misuse (Cluster 1) and issues of maintenance(Cluster 2) and are consistent with the groups in Fig.8. Note that the bar charts in Fig.9 denote the ratio of selection (right)/deselection (left) and that some causes, such as 'misuse', characterize the clusters, though their selections do not dominate in the cluster.

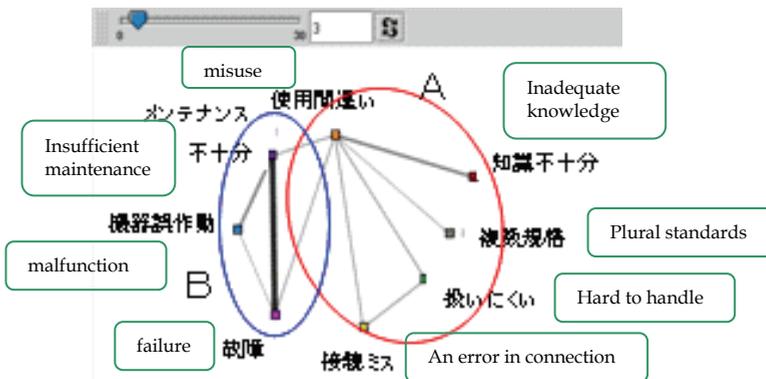


Fig. 8. Co-occurrence relations between the causes of incidents

To understand the condition under which these clusters can be causes of each incident, we applied a decision tree algorithm, where we set the cluster to which the incident belongs as

an objective variable and the types of equipments, the time and the location the incident occurs, and the occupation and the period of job experience of the person concerned as explanatory variables. Fig. 10 shows that the main parameter separating the incident records into the clusters is the type of equipments. This says that the apparatus such as a catheter, a tube, a puncture device is related to Cluster 1 (misuse) and that the incidents associated with the structurally-complex equipments such as a pump set, a mechanical ventilator, a hemodialysis monitor and an X-ray apparatus are mainly caused by the causes in Cluster 2 (issue of maintenance).

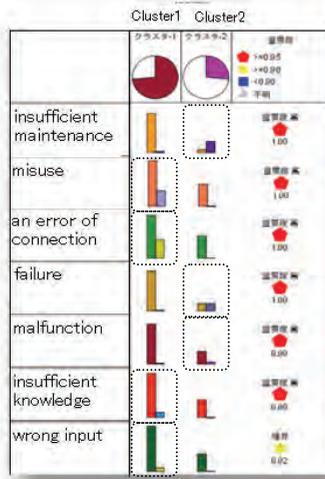


Fig. 9. The clusters of the type of causes based on their co-occurrence relations (main part)

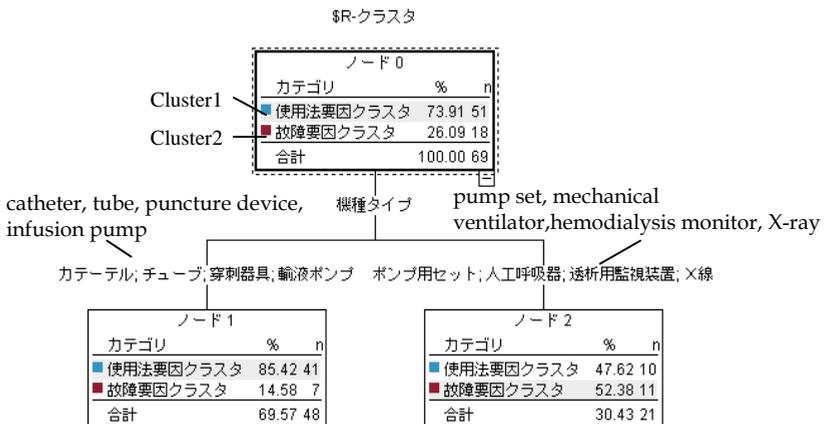


Fig. 10. The decision tree which relates the type of medical equipments and their cause

Notice that Fig. 10 indicates the incidents associated with infusion pump is also caused mainly by misuse (Cluster 1), though it has relatively complex structure. In order to clarify the relationships between the type of equipments and the causes of incidents, we again used

Web graph (Fig. 11) and found that the incidents related to an infusion pump are mainly caused by ‘wrong input’, ‘an error of connection’, ‘misuse’ and ‘oblivescence of holding down a switch’ rather than ‘inadequate knowledge’ or ‘failure’. This suggests that there is a problem of an infusion pump, which itself seduces users into the wrong use because of its human-machine interface, and that it is necessary to improve the interface to prevent from operation mistake.

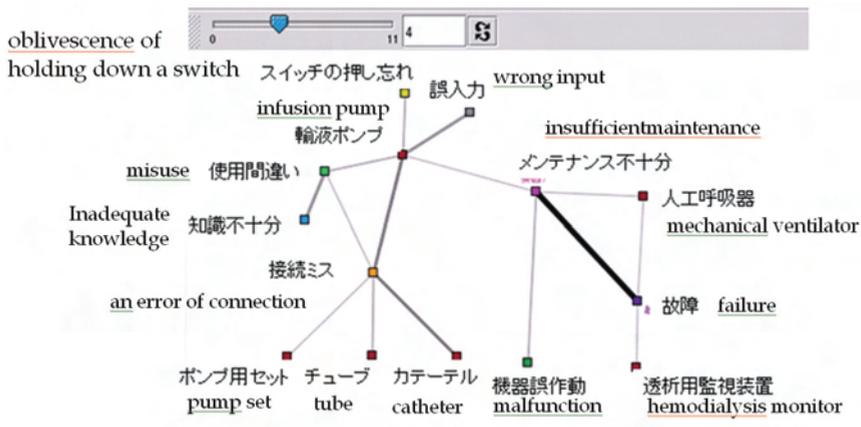


Fig. 11. The relation between the type of equipments and the causes of incidents

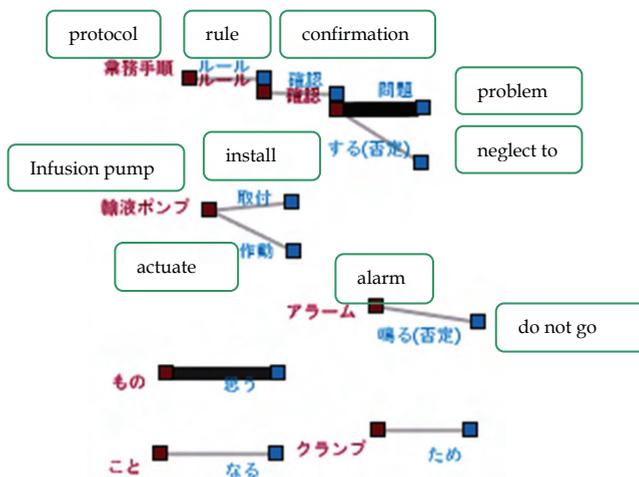


Fig. 12. Backgrounds and causes of incidents related to medical equipments

We applied Word-linking method to the free description data, which describes the background and the causes of incidents (Fig. 12). This indicates that there are major statements such as ‘(The person concerned) neglects to confirm (based on) a rule and a protocol’, ‘alarm did not go off’, ‘actuation/installing of the infusion pump’. Comparing this result to the one on the analysis on medication, it is common that both of them focus on

the necessity of confirmation. Medical equipments, however, leave a room of improvement of a user interface to lead an operator to confirm spontaneously and of mechanism to check his/her configuration under the assumption that people may forget to confirm. This is necessary because it is hard to force medical experts to perfectly confirm under the circumstances where interruption of a task frequently occurs.

5. Summary and conclusion

We applied data-mining and text-mining techniques to medical near-miss case data related to medicines and medical equipments, and analyzed the causal relationship of occurrences of the near-miss cases and the opinions on the counter measures against them.

As for medication, the decision tree obtained by the C5.0 algorithm shows that the major causes are classified mainly by the contents of incidents, namely medicine error or not. In the case of medicine error, the major causes mainly consist of resemblance of name and/or external form. The other cases mainly come from carelessness and misconception. This suggests that, as a counter measure of near-miss cases and medical accidents, it is valid to avoid adopting a confusing name or shape of the medicines themselves or their packages, not just to pay attentions.

To prevent oversights of errors, which may become a medical accident, it is also important to determine the rules of the occurrences of the oversights. We, therefore, applied the Apriori algorithm to determine the rules of oversights depending on service phases and contents of incidents. The result indicates that an oversight tends to occur, if the error happens in the administration phase, or the error is related to amount. Especially, the tendency of oversight in the case of medicine error in the administration phase is consistent with the results of the decision tree.

Since the cause of medicine errors stems from the resemblance of their name, we identified which medicines are mixed-up because of the name resemblance using a Web graph. As a result, we found that there are two types of name resemblance, one of which is the similar line of letters and the other of which is the same name apart from the last symbol. Since most of the paired medicines unfortunately have a similar dosage form, pharmaceutical companies should avoid naming a medicine with a similar name to existing medicines whose dosage form is also similar, and medical experts should pay attention to these.

We applied the Word-linking method to the free description data and found concrete information on the backgrounds and the causes of the incidents depending on job titles. Both describe the common statements on the problems of the checking system of the protocol and the rule, and the unsatisfactory confirmation. Nurses point out the systematic problem of communication and pharmacists indicate the problem of adoption of medicines. Those are systematic problems. In spite of such indications, there are few opinions on the counter measures related to the system of medical situations. This suggests that medical experts such as nurses and pharmacists tend to try to find the solutions to problems within themselves.

As for medical equipments, we utilized Web graph and TwoStep clustering algorithm to find the pattern of the co-occurring relations between the causes of incidents, which consist of two clusters corresponding to misuse and issues of maintenance. To understand the condition under which these clusters can be causes of each incident, we applied a decision tree algorithm, where we set the cluster to which the incident belongs as an objective variable and the types of equipments, the time and the location the incident occurs, and the

occupation and the period of job experience of the person concerned as explanatory variables. This says that the apparatus with relatively simple structure is related to the cluster of misuse and that the incidents associated with the structurally-complex equipment are mainly caused by the causes in the cluster related to maintenance. The exception is an infusion pump, whose incidents are also caused mainly by misuse, though it has relatively complex structure. This suggests that there is a problem of an infusion pump, which itself seduces users into the wrong use because of its human-machine interface.

To prevent near-miss incidents, it is obviously desired to promote the effort of each medical expert, but it is important to take measures against systematic error, such as the adoption policy of medicines with confusing names, the system of communication, better user interface to prevent from misuse and to prompt user to confirm.

As we showed in this paper, data-mining and text-mining are powerful tools to discover information that cannot be found by simple aggregation calculations and descriptive statistics. This is because these methodologies neglect the information contained in the relationship between data items, which can be extracted by data-mining and text-mining approaches.

In order to find countermeasures against near-miss cases and medical accidents related to medicines by means of data/text-mining approaches, it is necessary to collect and disclose the near-miss cases continually to find time series patterns and to confirm the validity of our countermeasures.

8. References

- Hayasaka, T.; Kimura, M.; Ohkura, M.; Tsuchiya, F.(2006). The analysis of medical near-miss cases applying text mining method, The 36th Annual Meeting of the Kanto-branch Japan Ergonomics Society.
- Hayasaka, T.; Kimura, M.; Ohkura, M.; Tsuchiya, F.(2007). The analysis of medical near-miss cases applying text mining method, Proceedings of IPSJ Annual Convention, 3J-3.
- Hripcsak, G.; Bakken, S.; Stetson, D. P.; Patel, L. V.(2003). Mining complex clinical data for patient safety research: a framework for event discovery, Journal of Biomedical Informatics, 36, pp.120-130.
- Kimura, M.; Furukawa, H.; Tsukamoto, H; Tasaki, H.; Kuga, M; Ohkura, M.; Tsuchiya, F.(2005). Analysis of Questionnaires Regarding Safety of Drug Use, Application of Text Mining to Free Description Questionnaires, The Japanese Journal of Ergonomics, Vol.41 No.5 pp.297-3051.
- Kimura, M.; Tatsuno, K.; Hayasaka, T.; Takahashi, Y.; Aoto, T.; Ohkura, M.; Tsuchiya, F.(2007). The Analysis of Near-Miss Cases Using Data-Mining Approach, Proceedings of the 12th International Conference on Human-Computer Interaction, pp.474-483.
- Tatsuno, K.; Kimura, M.; Ohkura, M.; Tsuchiya, F.(2005). Applying the data mining technique to near-miss cases of medicine (III), The 35th Annual Meeting of the Kanto-branch Japan Ergonomics Society.
- Takahashi, Y.; Kimura, M.; Ohkura, M.; Aoto, T.; Tsuchiya, F. (2004). Study on analysis for near-miss cases with medicine (II), The 34th Annual Meeting of the Kanto-branch Japan Ergonomics Society.
- Takahashi, Y.; Kimura, M.; Ohkura, M.; Aoto, T.; Tsuchiya, F. (2005). Study on analysis for near-miss cases of Medication, Proceedings of IPSJ Annual Convention, 6V-6.

- Tsumoto, S.; Hirano, S.(2007). Data Mining as Complex Medical Engineering, Journal of Japanese Society for Artificial Intelligence, Vol.22, No.2, pp.201-207.
- Watabe, S.; Kimura, M.; Ohkura, M.; Tsuchiya, F.(2007).The analysis of the near-miss cases of medical equipments using the text-mining technique, The 37th Annual Meeting of the Kanto-branch Japan Ergonomics Society.
- Watabe, S.; Kimura, M.; Ohkura, M.; Tsuchiya, F.(2008). The analysis of the near-miss cases of medical equipments using the text-mining technique, Proceedings of the IEICE General Conference, A-18-3.