



# Data Mining Challenges in the Management of Aviation Safety

Olli Sjöblom

## ► To cite this version:

Olli Sjöblom. Data Mining Challenges in the Management of Aviation Safety. 13th Conference on e-Business, e-Services and e-Society (I3E), Nov 2014, Sanya, China. pp.213-223, 10.1007/978-3-662-45526-5\_21 . hal-01342151

**HAL Id: hal-01342151**

**<https://inria.hal.science/hal-01342151>**

Submitted on 5 Jul 2016

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

# Data Mining Challenges in the Management of Aviation Safety

Olli Sjöblom

Turku University School of Economics, Turku, Finland

oljusj@utu.fi

**Abstract.** This paper introduces aviation safety data analysis as an important application area for data mining. Safety is a key strategic management concern for safety-critical industries and management needs new, more efficient tools and methods for more effective management routines. The aviation field is confronted with increasing challenges to provide safe and fluent services. Air travel has grown steadily during the last decades with a direct impact on the air traffic control. At the same time, the competition has become tougher because of increasing fuel prices and growing demand for air travel.

**Keywords:** Management, Flight Safety, Strategic Management, Data Mining, Text Mining, Analysis Method

## 1 Introduction

Organisational decision making, especially in safety-critical systems, such as nuclear power and air traffic, is a complicated task. For successful operations, an acceptable air safety record has been required from the airline [1]. Air traffic has generally been forecasted to grow 5 – 6 % annually over the next two decades [2], or even over the next 10 – 15 years, the global air travel will probably double [3]. Consequently, the number of accidents will respectively increase if nothing were done to improve it, which development would, clearly, be unacceptable. This is why new and efficient ways for improving air safety need to be explored [4]. The conventional safety tools and methods based on data collection have reached their peak performance because of their inability to create new knowledge. Usually, data accumulates faster than it can be processed [5]. For further improvements new methods and tools are urgently needed [6].

## 2 Management in safety-related context

Any system can be recognised to consist of elements, or factors, or parts that make up the whole [7]. Managing the organisation is exercised largely through management processes, in which the means of managerial communication inter-links with the environment. Johnsen (2002) defines the management process as “*the interaction between*

*people who want to attain mutual ends through mutual means.”* [8]. The strategy of the corporation is according to Johnson et al. [9] to concern the organisation’s mission, vision and objectives, developing plans and policies to use resources for enhancing the performance of the organisation.

Kettunen et al. [10] emphasise the managerial challenges in the safety-critical industries, which are typically related to finding a balance between diverging demands and expectations, like economy- and safety-related objects without forgetting the priorities-setting and maintaining focus on these components. The key action is a continuous balancing between taking risks and allocating resources for risk management. A scale with theoretical ends can be displayed, where at one end there is a situation where risks do not exist because the resources allocated are infinite; at the other end no resources are allocated because the risks are ignored and thus they are (practically) infinite. The reality is found somewhere in between, but no fixed location can be defined because all environments are somewhat unique and are also changing all the time. In daily operations perhaps existing hidden threats produce the need to maintain extra safety level naturally causing additional costs.

In studying risk management, the concept of tension cannot be ignored. It refers to the challenges of balancing conflicting objectives or expectations, like safety and other goals. These might exist for various reasons, even in the situation in which the executives of the organisation have set a high safety level as the priority official goal [11]. In case warning signals appear, responding to those should happen without delay allocating safety resources to the critical area.

The safety decisions in an air traffic company follow the same pattern as other strategic decisions. Risk management should be carried out in parallel with safety management, referring to measures seeking to identify, assess and control risks on the organisational level having the goal to ensure the organisational and environmental safety. The executive management is responsible for recognising the safety significance of the ways the organisation is operated and maintained [12]. Managing risk and safety has been problematic in air transport: very high levels of safety are too costly – high levels of risk are unacceptable. Therefore, safety reports have been collected through decades to investigate and assess risks and to define risk standards, which are consistent with the value systems of the society [13, 14].

The value of safety cannot be estimated in any traditional way, because it has no determined price. Theoretically, limitless resources should be allocated to it, because one single failure may lead to significant losses in the form of missed business possibilities and claims for covering the damage caused to a third party. Kaplanski and Haim [15] have presented some estimates for the accident costs. A very large disaster with hundreds of casualties will cause a loss of about \$1 billion for an airline company. However, the observed market effect has been found to be about 60 times larger; Kaplanski and Haim (2010) have found the evidence of a significant negative effect with an average market loss of more than \$60 billion per aviation disaster. However, budget constraints set limits in practise and therefore a certain risk has to be accepted by achieving a sufficient safety level. There is never a 0-level risk. In case sufficient resources could not be allocated to achieve the required level of safety, the whole air traffic business would be critical. When confronting such a situation, the

operations are to be adjusted by diminishing or changing them to correspond with the allocable safety resources so that a sufficient safety level is maintained.

Estimating the significance and importance of different alternatives in managing risks also needs tools, the exact definition of which is important for making strategic decisions. After the executive management has set goals as the thresholds of achievement, there must be methods and models to measure to what degree the achievements have been realised. In the decision process, there is always question about evaluating different alternatives. Any matter having significance enough to be taken into account in the evaluation process should be considered for evaluation [16]. Rumsfeld [17] has defined (simply expressed) three categories for knowledge: first, we know what we do know; second, we know what we do not know; and, finally, we do not know what we do not know. The hidden dangers belong to the last group, so in case we know what we are searching for, we obviously have means to reach it, but otherwise we need tools for finding something we do not know we are looking for. Thus, a deeper understanding is required for developing better methods and refining rules and practices that will contribute to higher levels of safety.

The unknown lethal factors brought into daylight could be eliminated; at least a significant part of them and a sufficient safety level could be reached with reduced investment allocation. For air traffic, there is theoretically no upper limit to allocate resources to safety in different forms. The relation between safety and cost efficiency could be illustrated explicitly comparing the costs between comprehensive maintenance programs and maintenance-induced accidents, the benefits that outweigh the accident costs [18]. The process for allocating extra resources to special projects might become even more troublesome in case there are interdependencies among the projects [16].

### **3 Flight Safety**

According to the ICAO Safety Management manual [19], safety is defined as “*a state in which the risk of harm to persons or property damage is reduced to, and maintained at or below, an acceptable level through a continuing process of hazard identification and risk management*”. Safety is not a matter-of-course, but the result of a rather complicated, carefully structured and comprehensive management process approaching to all airline safety aspects, particularly those of flight operations.

Air traffic is full of incidents and deviations that do not contain any hazard as such, but need to be reported and investigated to find out potential lethal trends. These undesirable, but very minor events are valuable investigation subjects for risk and safety specialists to build an understanding about their causes and to detect unsafe trends. Investigation also reveals whether countermeasures are warranted and how to reduce or eliminate potential accidents [20]. The appearance of similar recurring cases (a cluster, cf. Chapter 6) may indicate a hazardous trend that should be analysed very carefully to find out whether a real danger exists or not. The possibly existing lethal trends are trying to penetrate through the layers of defences, barriers and safeguards (cf. Figure 1) that, fortunately, usually stop them from proceeding. Because serious

incidents and even accidents do happen, it can be presumed that after a certain amount of time they pass all the layers but the last one; then they will pass the last layer as well, which leads to accidents.

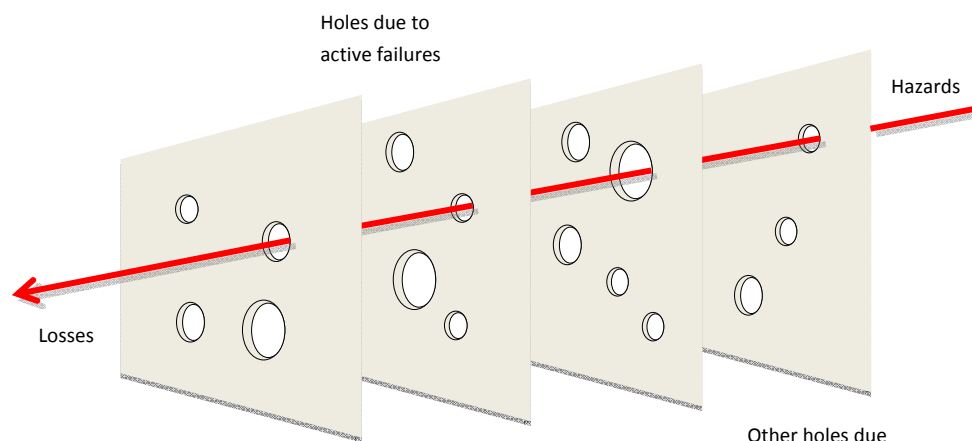
Finding trends from flight safety data, especially from narrative data has required significant human involvement. Thus, the analysis process and its possible results rely on the skill, memory and experience of the safety officers [21]. Watson [22] found that with conventional techniques it might take years to find meaningful relationships. Before text mining systems (one sub-class of data mining) were developed, there were no tools for analysing textual data with computers. Data mining provides a worthy analysis method in order to illustrate the safety indicators and to reveal undesired trends.

## 4 Safety Tools and Systems

Accident analysis as well as flight and operations modelling and simulation enhance the understanding of risk, but this is usually reactive and produces knowledge about causal factors potentially at the human and/or financial cost. Risk modelling typically collects knowledge resulting from flight safety analysis, human experience and theoretical and empirical studies. The goal of aviation risk assessment is to be comprehensive, timely and proactive, and this is why the analysis methods should be enhanced [23].

In aviation, the quantitative assessment of risk is particularly challenging, because the deviation events are extremely rare and the causal factors are non-linearly related to the events which makes them difficult to quantify [23]. The eventuality for the incident or accident occurring may be markedly reduced in case the risks can be efficiently diagnosed [24]. Then the question is: how to find and identify deviations leading to incidents and those leading to accidents? Reason [25] has modelled the process for the occurrence of accidents in his Swiss Cheese model, which is presented in Figure 1. The hazards appear from the right-hand side. Normally, their progress is stopped by successive layers of defences, barriers and lifeguards. If the process goes through all of these 'holes in the cheese slices', formally called the limited windows of accident opportunities, an accident will happen.

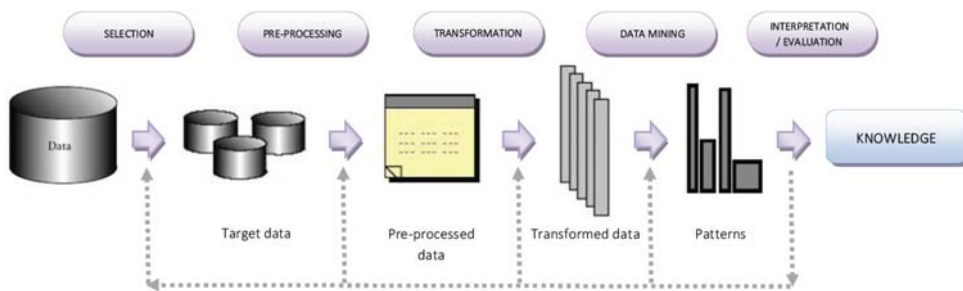
Kettunen et al. [10] regard redundancy as a method in improving safety by the duplication and overlap of critical factors like systems, functions and/or personnel. In general, redundancy can augment safety as such, but may also have counter-productive or unexpected effects, especially in case it is not managed properly. These unwanted effects can increase the complexity of the systems, which may hide individual failures and make them latent, so that they remain unnoticed and uncorrected and may even accumulate over time. Under these circumstances, a rather rare event might act as a trigger for an avalanche of unexpected events, which may be difficult to handle [10]. For situations of this kind, the Reason's Swiss Cheese model would work out excellently.



**Figure. 1.** The Swiss Cheese model (adapted from Reason 1997, 2000 [25, 26])

## 5 Data Mining in Flight safety

Several different methods are recognised as data mining methods and a mining system can use the combinations of several of these methods. Parsaye [27] describes data mining as searching in the data for the patterns of information to guide a decision support process. These, often called “the nuggets of knowledge”, are hidden in vast amounts of data and are practically undiscoverable with conventional techniques [22]. Using mining software, knowledge of data is combined by an analyst with advanced machine learning technologies to discover the relationships. In the discovery process to find hidden patterns, there are neither hypotheses nor any other predetermined model of the characteristics of the patterns. Obviously large databases, like those of aviation incidents and other deviations, contain a large number of patterns, so that the user of the discovery system can practically never ask the right question. The mining process acts as a decision support system that will not give straight answers to the questions; that is why skilled analytical and technical specialists are still required to interpret the created output [28]. The process contains several steps or phases (cf. Figure 2) that must be gone through to form knowledge from raw data. To be understandable the information must be presented with reports, graphs or in other suitable forms once found.



**Figure. 2.** The Knowledge Discovery in a Database Process (adapted from Fayyad et al. 1996 [29])

With structured data, the explanation of a case usually tells the truth to a certain extent, but completed with narrative data it can be close to 100 %, at least theoretically. Mining combined with other methods will give significant contributions to the decision processes. The idea to use text mining in the analysis of flight safety reports occurred along the need to analyse large amounts of narrative reports and when reports

about successful text mining projects in the flight safety data analysis of English narratives were published [21, 30].

## **6 Testing Three Tools - Data Mining in Finnish**

The basic idea of cluster analysis is that all the texts within each cluster have a high similarity in content [31]. This method was chosen for this study because it is an essential mining function in searching for similar documents, able to reveal a recurring hazard that might lead to an accident. It explores the data set and determines the structure of natural groupings without any preliminary assumptions. Another reason for its choice was the direct applicability to Reason's Swiss Cheese model presented in Figure 1. A third reason was that English literature gives several examples about using clustering in mining flight safety reports. These results have proved its better performance compared with more traditional statistical methods [32].

The beginning was finding text mining tools for processing Finnish. Three different systems seemed to be appropriate for benchmarking. The author was aware of one prototype (GILTA), one commercial product (TEMIS) with a Finnish module prototype, and one commercial system (PolyVista) with encouraging results mining Spanish, which seemed worth testing in Finnish. The Finnish Civil Aviation Authority granted the test data of 1240 cases (Target data on Figure 2), which created "a critical mass" for study.

The pre-processing produced filtered data containing 10572 word tokens, numbers and special characters, call signs, headings, the temperature, etc. The amount could be reduced to 8294 when parentheses and other similar characters without relevance were removed. The next procedure was preparing the lists of stop words (those to be ignored because of having no information) and synonyms. No transformation was needed because the data was extracted from one database.

The first round produced already promising results. Due to the Finnish module of TEMIS, no pre-processing was necessary. It created 26 clusters, their size varying between 108 and 21 reports. As the biggest cluster contained more than 100 clusters, the operator allowed the tool divide it into two sub-clusters with 58 and 50 documents. After the division, the biggest cluster included 78 reports. The similarity (range 5-1) of the five closest clusters varied from 3.41 to 2.07 %, which supports the assumption that the clusters are different from each other and thus this method in this data selection is reliable. Because the maximum degrees of explanation of the clusters, about 18 %, are relatively high, they prove that the clusters are composed of relevant reports and the most explaining reports alone might well reveal a trend that should be examined more thoroughly.

As Klopitchenko [33] says, interpreting the mining results is more art and common sense than science. The one single mining round of TEMIS made the direct comparison of the results challenging. Despite it, due to the high efficiency of the system with its in-built module for Finnish and because the mining results did not seem to require major changes, missing the second mining round was not considered a cause for losing significant information.

The smallest clusters began to produce some directly applicable information indicating that the sizes of the clusters play a significant role in the applicability of the results. This must, however, be scaled with the amount of production data. Additionally, a couple of similar cases found do not automatically create a dangerous trend; the way they occur and the reasons causing them can only be estimated by a thorough examination and investigation by human analysts. The results of TEMIS ought to be examined differently from the two other systems due to its interface and way of producing results which differ remarkably from the others. This, however, does not mean that these mining results would not be coherent with those of the other ones.

GILTA (manaGIng Large Text mAsses) divided the data on both rounds into 100 clusters (named classes) on the basis of the nine most significant words. Hence, on the first round 63 clusters contained less than 10 reports. These were easily analysable by a human analyst and could already be considered good mining results, proving clustering to be a useful method for this type of data. Some of the bigger classes could be interpreted as being real clusters, but according to experience the sizes should be reduced to less than 20. The results that were produced in Excel-form made it possible to carry out a comprehensible analysis and comparison of them with the results found with other tools. The system left out four reports beyond defined clusters.

PolyVista was originally built for using in English, but due to encouraging results with Spanish, its applicability for Finnish was tested, too. The system set score 100 for the most content describing word of the cluster and correspondent values to the others. The scores of the ten most important words of each cluster were only available, not the reports. The reports of the clusters could be ‘guessed’ by comparing the scores with the most important words in GILTA changing their relative weights for comparison. The data was processed determining the number of clusters first to be 6 and then raising it up to 20 in a second step. When there were 20 clusters, the smallest of them contained 10 reports and the biggest 232. In the case of 20 clusters, in eleven of them the scores of the three most important words were more than 50. In the last cluster containing 10 reports, the scores of the 10 most important words were 50 or more, which can be considered a good mining result.

As one result of the first mining round, the need for tuning, especially the definition of stop words and synonyms was discovered. Some pure mistakes, like some common stop words and synonyms forgotten from the list, were noticed. A more significant problem was the appearance of some frequently used “common” words (like ‘plane’ with its synonyms ‘airplane’ and ‘aircraft’) skewing the results. Their role in the data was carefully analysed [34], using an application called NVivo to get a deeper analysis. NVivo itself has no mining characteristics, but is used in analysing qualitative information, especially meeting the requirements of deep levels of analyses on different quantities of data, varying between a couple of sentences and thousands of text rows. In this context, the most important feature was cross-examining the mining results applying its search engine and query functions. Almost one hundred checking procedures were made with synonyms and stop words to prepare the data for the second round. After the careful estimation of the impact of possible changes, no major ones were made to keep the process unchanged but making the results more accurate.



After the second mining round with GILTA and PolyVista was performed, the results were studied carefully using the professional skills of a flight safety inspector. The coherent clusters were taken into more detailed inspection. The progress as the change of distribution can be recognised through the increased percentage of ‘sense making’<sup>1</sup> clusters, as for GILTA, illustrated in Table 1 displaying the minor, but perceptible change. First, the number of the relevant clusters increased from 9 to 11, and their average size diminished from 11.9 to 10.5 reports per cluster, shown in columns two and three. Further, the average weight of the nine most important words increased from 5.88 to 6.44 and the correspondent standard deviation diminished from 5.588 to 5.065, as shown in the two next columns. All these changes indicate the movement towards the aimed more homogenous clusters.

**Table 1.** Results illustration in GILTA rounds I and II

Round	Clusters	Average size	Average weight	Correspondent Standard Deviation
I	9	11.9	5.88	5.588
II	11	10.5	6.44	5.065

As already mentioned, the mining results of PolyVista must be analysed differently. Although comparing the weights of the most significant words is a cursory method, it was noticed to be relevant in this context. The results are illustrated in Table 2 showing an obvious progress between the two rounds. On the first round, 40.0 % of the clusters seemed to belong to the ‘sense making’ clusters, on the second 52.3 %. The size of the clusters did not seem to have any linear impact, but on both rounds those were found among the smallest ones. The average sizes changed from the first round being 37.9 compared with 62.5 of all clusters to the second, being then 20.6 compared with 28.2. These numbers illustrate that more information is achieved from the results of round II.

**Table 2.** Cluster distribution change between rounds 1 and 2 in PolyVista

Criteria / Round #	1	2
‘Sense making’ clusters of all content	40.0 %	52.3 %
Average size of all clusters (reports)	62.5	28.2
Average size of ‘sense making’ clusters	37.2	20.6

Proceeding with the same test and putting the results in a graphic presentation in Excel, the increased homogeneity was seen also from the ‘centre of gravity’ moving from the beginning of the rows rightwards as well as from the top downwards, when the clusters were sorted by the weight of the most significant words. It means that the number of clusters having more significant words increased. This occurred with both systems, indicating a slight improvement using this method, too.

Based on the professional skills and experience of the author, in case the safety personnel know what they are looking for, business intelligence (BI) methods could

<sup>1</sup>Clusters, from which information can be seen clearly as such

be applicable, allowing database queries using numerous keywords to search for known cases of a certain type or their combinations. BI could also be applied as a complementary method when mining is used to find something worth examining.

## **7 Results and Discussion**

As already expressed before, the mining process does not give straight answers to the questions, but it acts as a support system for producing information for decision making. That is why experienced analytical and technical specialists are needed to interpret the created output. The testing process proved that data mining is neither an easy nor a fast method, but might be the only one for uncovering hidden information. All the results support the premise that it could reveal important safety information from fast accumulating, vast amounts of data, not accessible with other methods, to be used as an essential factor for strategic safety management. It is worth noticing that the test data was that contained no lethal trends, but in other case they could have been discovered and revealed using the method and tools as done in this study. An additional detail is worth noting - all the used tools left out almost the same reports as outliers.

The research process confirmed that text mining is a challenging task, especially in small language groups, where tools for text mining are scarcer than for big languages such as English which is an “easy” language for search technologies. Narrative text mining is generally demanding due to the multiplicity of languages spoken in the world. Especially languages with small user groups, such as Finnish, have to wait for efficient tools being developed much longer than the major languages. The search technologies are challenged by inflected forms and compounds. In Finnish, for example, the words may have thousands of inflected forms and in addition to that, they can be parts of compounds in almost countless combinations [35]. On average, every seventh word can be found in its basic form in fluent Finnish texts [36]. From the point of view of language processing, two significant results were achieved: first, Finnish texts were successfully mined with a tool originally to be used in an English environment. Secondly, the Finnish module for TEMIS was successfully production tested with real Finnish production data.

The number of clusters proved to be significant in the process: the more clusters, the better results. Mining is an iterative process although it makes no sense to increase the amount of rounds too much. Although this study has offered data mining as one solution to growing challenges, it is to be noticed that it is only one among several methods. Its special characteristic simply expressed is the ability to find something that is not known but expected to exist. Data mining has been used successfully for several years by a couple of airlines and other actors in the aviation industry. The process chain, beginning from the collection of safety data and ending in revised regulations for improving flight safety, going through several mining rounds and analyses to produce issued aviation rules and instructions, is rather long and demanding. Despite its complexity, it is worth going through, even for avoiding one single accident.

## References

1. Liou, J.J.H., L. Yen, and G.-H. Tzeng, Building an effective safety management system for airlines. *Journal of Air Transport Management*, 2008. 14(1): p. 20–26.
2. Netjasov, F. and M. Janic, A review of research on risk and safety modelling in civil aviation. *Journal of Air Transport Management*, 2008. 14(4): p. 213–220.
3. Global Airline Industry Program. Analysis: The Airline Industry. Global Airline Industry Program [WWW-page] 2008 [cited 2011 9.5.]; Available from: [http://web.mit.edu/airlines/analysis/analysis\\_airline\\_industry.html](http://web.mit.edu/airlines/analysis/analysis_airline_industry.html).
4. European Commission, Proposal for a DIRECTIVE OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL on occurrence reporting in civil aviation, Commission of the European Communities, Editor 2000: Brussels.
5. Wang, X., et al., LSSVM with Fuzzy Pre-processing Model Based Aero Engine Data Mining Technology, in *Advanced Data Mining and Applications 2007*, Springer Berlin / Heidelberg: Heidelberg. p. 100-109.
6. Evans, B., A.I. Glendon, and P.A. Creed, Development and initial validation of an Aviation Safety Climate Scale. *Journal of Safety Research*, 2007. 38(6): p. 675–682.
7. Barnard, C.I., *The Functions of the Executive*. Thirtieth Anniversary ed 1938, Cambridge, MA and London, UK: Harvard University Press. 334.
8. Johnsen, E., *Managing the Managerial Process. A Participative Process* 2002, Copenhagen: DJØF Publishing. 606.
9. Johnson, G., K. Scholes, and R. Whittington, *Exploring Corporate Strategy* 2005: Pearson Education Limited.
10. Kettunen, J., T. Reiman, and B. Wahlström, Safety management challenges and tensions in the European nuclear power industry. *Scandinavian Journal of Management*, 2007. 23(4): p. 424-444.
11. Sagan, S.D., *The limits of safety. Organizations, accidents, and nuclear weapons*. 1993, Princeton, NJ: Princeton University Press.
12. OECD/NEA, State-of-the-art report on systematic approaches to safety management, O.N.E. Agency, Editor 2006, OECD Nuclear Energy Agency: Issy-les-Moulineaux.
13. Janic, M., An assessment of risk and safety in civil aviation. *Journal of Air Transport Management*, 2000. 6(1): p. 43-50.
14. Sage, A.P. and E.B. White, Methodologies for Risk and Hazard Assessment: A Survey and Status Report. *IEEE Transaction on Systems, Man, and Cybernetics*, 1980. SMC-10(8): p. 425-446.
15. Kaplanski, G. and L. Haim, Sentiment and stock prices: The case of aviation disasters. *Journal of Financial Economics*, 2010. 95(2): p. 174-201.
16. Kirkwood, C.W., *Strategic Decision Making* 1997, Belmont, CA: Wadsworth Publishing Company.
17. Rumsfeld, D.H., News Transcript, 2002, U.S. Department of Defense Office of the Assistant Secretary of Defense (Public Affairs).
18. Castro, R., A Holistic Approach to Aviation Safety, in *Flight Safety Digest* 1988. p. 1-12.
19. ICAO, *Safety Management Manual*, 2009, International Civil Aviation Organization: Montreal, Canada. p. 264.
20. Kirwan, B., Incident reduction and risk migration. *Safety Science*, 2011. 49(1): p. 11–20.
21. Nazeri, Z., Application of Aviation Safety Data Mining Workbench at American Airlines. Proof-of-Concept Demonstration of Data and Text Mining., 2003, Center for Advanced Aviation Systems Development, MITRE Corporation Inc.: McLean, Virginia, US.

22. Watson, R.T., Data Management: Databases and Organizations. 2nd Edition ed, ed. J.W.S. 2nd Edition, 1999.1999: John Wiley & Sons.
23. Hadjimichael, M., A fuzzy expert system for aviation risk assessment. *Expert Systems with Applications*, 2009. 36(3): p. 6512–6519.
24. Lee, W.-K., Risk assessment modeling in aviation safety management. *Journal of Air Transport Management*, 2006. 12(5): p. 267–273.
25. Reason, J.T., *Managing the Risks of Organizational Accidents*1997, Aldershot: Ashgate Publishing Limited. 252.
26. Reason, J.T., Human error: models and management. *British Medical Journal*, 2000. 320(7237): p. 768-770.
27. Parsaye, K., A Characterization of Data Mining Technologies and Processes. *Journal of Data Warehousing*, 1997. 2(3): p. 2-15.
28. Kutais, B.G., ed. *Focus on the Internet*. 2006, Nova Science Publishers, Inc. 225.
29. Fayyad, U., G. Piatetsky-Shapiro, and P. Smyth, From Data Mining to Knowledge Discovery in Databases. *AI Magazine*, 1996. 17(3): p. 18.
30. Megaputer Intelligence. Flight safety data analysis for Southwest Airlines. 2004 [cited 2004 17 December]; Available from: <http://www.megaputer.com/company/cases/southwest.php3>
31. Rosell, M., Text Clustering Exploration. Swedish Text Representation and Clustering Results Unraveled, in *School of Computer Science and Communication*2009, Kungliga Tekniska Högskolan: Stockholm. p. 71.
32. Saracoglu, R., K. Tütüncü, and N. Allahverdi, A new approach on search for similar documents with multiple categories using fuzzy clustering. *Expert Systems with Applications: An International Journal*, 2008. 34(4): p. 2545-2554.
33. Kloptchenko, A., Text Mining Based on the Prototype Matching Method, in *Turku Centre for Computer Science*2003, Åbo Akademi University: Turku. p. 117 plus additional pages including original papers.
34. Lindén, K., Word Sense Discovery and Disambiguation, in *General Linguistics*2005, University of Helsinki: Helsinki. p. 191.
35. Karlsson, F., *Yleinen kielitiede*1994, Helsinki: Yliopistopaino.
36. Karlsson, F., *Finnish grammar*1987, Porvoo: WSOY.