

# Conceptualization Maturity Metrics for Expert Systems

Ovind Hauge, Paola Britos and Ramón García-Martínez

Norsk Teknisk Naturvitenskapelig Universitet. Norway  
Software & Knowledge Engineering Center. Graduate School. Buenos Aires Institute of  
Technology. Argentine  
Intelligent Systems Lab. School of Engineering. University of Buenos Aires. Argentine  
rgm@itba.edu.ar

**Abstract.** Metrics used on development of expert systems is not a well investigated problem area. This article suggests some metrics to be used to measure the maturity of the conceptualization process and the complexity of the decision process in the problem domain. We propose some further work to be done with these metrics. Applying those metrics makes new and interesting problems, concerning the structure of knowledge to surface.

## 1. Metrics

In software development measurement is used to provide some type of quantitative information to a decision making process, in many cases related to a development project [Ford, 2004; SEI, 2004]. The measurement can be on the production process or on the product it self. A metric should have different qualities to be applicable. It should as said be quantitative, but also objective, easy to find and well defined with a defined domain. The process of developing software is not trivial and measurement is done with relatively high uncertainty, but there are several metrics that are widely used today.

## 2. Suggested Metrics

In this section we propose some metrics that will examine the problem domain in expert system development context [García-Martínez & Britos, 2004, Firestone, 2004]. We will give interpretations of the metrics and will describe the expected development of the metrics throughout a development project. A metric should as mentioned have certain qualities as simplicity to be applicable. In the representations of knowledge there are several things that have these qualities. Rules, concepts, attributes and levels of decomposition are easy to count, they are objective and they are easy to find [Menzies,1999; Menzies & Cukic, 1999; 2000; Pasan & Clifford, 1991; Kang & Bahieel, 1990]. These things are therefore good candidates to be included in a metric. Then our suggested metrics are based on rules, concepts, attributes and number of decomposition levels [Nilsson, 1998].

### 2.1. Number of Concepts, Number of Rules or Number of Attributes

These are a very simple metrics. It is just to count the concepts, rules and attributes. But simplicity is good and these could tell us something about the complexity of the domain. we expect these metric to be increasing all the way throughout the project and converge to an unknown number at the end of the project. Since their values will increase all the way throughout the project it is hard to use them as a metric for maturity. But it could be an indication of maturity when their numbers converge. The table 1 shows the interpretations of these metrics.

Result	Cause
Low Few known concepts, rules or attributes	<ul style="list-style-type: none"> <li>• The problem area is simple</li> <li>• We do not know many of the concepts in the domain yet</li> </ul>
High Many known concepts, rules or attributes	<ul style="list-style-type: none"> <li>• The domain is complex with many concepts</li> <li>• We have good knowledge about the domain</li> </ul>

**Table 1.** Interpretations of results from “counting metrics”

These metric could be more useful if the results are compared to history from other projects in the same stages. When comparing to history data it could get an indication of the complexity of the project. These metrics will also be combined to others in the following sections.

### 2.2. Number of Concepts in a Rule / Number of Concepts

The number of concepts in a rule is the concepts that are already included in a rule. If you have 10 concepts and 7 of them are included in one or more rules the ratio will be 0.7. We believe this metric should converge to 1 when the project matures. The value will of course vary when you find new rules and new concepts. The value of this metric will decrease when we discover new concepts and increase when we include a new concept in a rule. If the value of this metric does not converge to 1 we either miss knowledge about relations between concepts in the domain or we have concepts in our knowledge base that are not used and most likely uninteresting. These concepts should therefore be removed. The interpretation of this metric is shown in the table 2.

Result	Cause
Low Many concepts not included in a rule	<ul style="list-style-type: none"> <li>• We miss knowledge about the concepts and the relations between concepts</li> <li>• We have many concepts that are uninteresting in our knowledge base</li> </ul>
High Most concepts included in a rule	<ul style="list-style-type: none"> <li>• We have good knowledge about the concepts</li> <li>• We have few uninteresting concepts in the knowledge base</li> <li>• There are many relations in the domain</li> </ul>

**Table 2.** Interpretation of results from “concepts in rule/concepts”

This metric will give a measure of the maturity of the knowledge base. If the value is close to 1 this it an indication that the knowledge base is mature. But pay attention to those cases where there are many relations in the domain. If there are a plenty of relations this metric can give a high value without a mature knowledge base as well.

This metric is therefore best to use for simple projects or together with a metric for complexity.

**2.3. Number of Attributes in a Rule / Number of Attributes**

This metric is similar to the previous one but we expect it to be easier to discover the concepts that the attributes. Because the attributes may not be discovered before we need them it is a bit difficult to use them as a measure of maturity. But if we have unused attributes we may miss something or we have included attributes that are unnecessary. If this is the case we should look at the reason and especially if the value of this metric is low. This metric could therefore be used as an indicator or alarm.

**2.4. Number of Concepts / Number of Rules**

This metric shows the development of the number of rules compared to the number of concepts. We expect that most concepts contribute to the creation of at least one or most likely several rules. And with good knowledge about the relations in the domain this metric will in most cases decrease below 1.0. In highly related problem domain will the value be much lower than 1.0. This metric can still have a high value at the same time as we have a mature knowledge base. In the cases where the domain only contains a small set of very complex relations the number of rules will be low, but the number of concepts will be high. we recommend combining this metric with some metric for complexity of the domain. Interpretation of the metric is found in Table 3.

Result	Cause
Low Many rules	<ul style="list-style-type: none"> <li>• We know the relations of the domain and have a mature rule-base</li> <li>• The domain is mature</li> <li>• Complex domain with many relations</li> <li>• Redundant rules</li> </ul>
High Few rules	<ul style="list-style-type: none"> <li>• We do not know the rules of the domain well enough</li> <li>• The domain is not very mature. The relations in the domain are not known.</li> <li>• We have too many uninteresting concepts</li> <li>• Many concepts are only included in one or few, very complex rules.</li> </ul>

**Table 3.** Interpretation of “concepts/rules”

**2.5. Average Number of Attributes per Concept**

This metric is an indication of the complexity of the domain. A high value means that each concept has several related attributes and this indicates a more complex domain. It can also be used as a metric for maturity. We expect the value to vary during the project as we discover new concepts and new attributes. In the start of the project it is most likely that we find the most important concepts which have the highest number of related attributes. As the project develops new concepts will be found. We believe that the concepts found in the latter parts of the project will have fewer related attributes than then ones found in the start of the project and the value will therefore decrease. It will converge at the end of the project, when no new

concepts and attributes are found. This indicates that the knowledge base is maturing. The Table 4 shows our interpretations of the metric. As we see a different number of concepts could give this metric different outcome or value.

Result	Cause
Low Few attributes per concept	<ul style="list-style-type: none"> <li>• The problem domain is simple and each concept have few interesting attributes</li> <li>• There are many concepts with few attributes</li> <li>• We do not know the problem domain well, we have not discovered all the necessary attributes</li> </ul>
High Many attributes per concept	<ul style="list-style-type: none"> <li>• The domain is big and complex</li> <li>• There are few concepts with many related attributes</li> <li>• We have good knowledge about the problem domain</li> </ul>

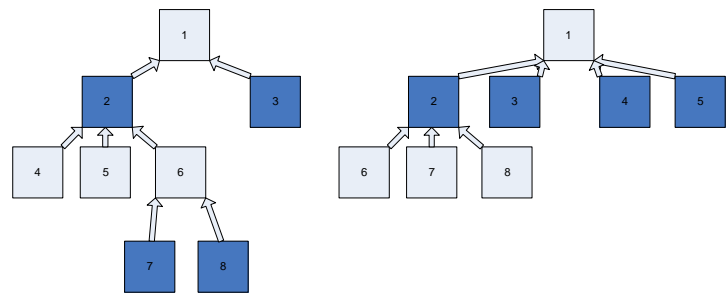
**Table 4.** Interpretations “Average attributes/concept”

**2.6.  $A \cdot (\text{Number of Concepts}) + B \cdot (\text{Average Number of Attributes Per Concept})$**

To get a better indication of the complexity of the project we suggest combining the number of concepts and the average number of attributes per concept. This will remove the different outcomes in average number of attributes per concept that was caused by the number of concepts. To be able to get a reasonable result the two metrics must be weighted by the factors A and B. To be able to find values for these factors we propose using history data. This is not within the scope of this paper and will therefore not be done here.

**2.7. Average Number of Levels in Decision Tree**

For the tasks that are decomposed this average will most likely increase throughout the project and stabilize to the end of the project. The metric is calculated by just counting the levels of the decision trees, add them, and divided the sum on the number of trees. Given the example in Figure 1 we will get the following result:  $(4+3)/2=3.5$ .



**Fig. 1.** Decision trees

The Table 5 shows our interpretations of the metric. A high degree of composition can indicate high complexity but also a high degree of understanding of the decision.

Result	Cause
Low Few levels of decomposition	<ul style="list-style-type: none"> <li>• The domain is simple</li> <li>• We have not decomposed the decisions</li> <li>• We do not have complete knowledge about the domain</li> <li>• We have discovered all decisions but not decomposed they yet</li> </ul>
High Many levels of decomposition	<ul style="list-style-type: none"> <li>• The domain is complex</li> <li>• We have good knowledge about the domain</li> <li>• We totally miss information about some decisions in the domain, which would have decreased the average.</li> </ul>

Table 5. Interpretation of results “Average levels in decisions”

### 2.8. Average Number of Concepts Included in Each Rule

Each rule contains one or more concepts. The number of concepts included in a rule could be a measurement of the complexity of the problem. We expect this number to be increasing as we discover more complex relationships within the problem domain. At the end of the project we suggest that the value converges to a constant. This convergence could be an indication of maturity of the knowledge base. The table 6 shows our interpretations of this metric. We see that the number of rules and the degree of decomposition affects the outcome of this metric, but if the average is high it is likely that we have a complex domain.

Result	Cause
Low Few attributes per concept	<ul style="list-style-type: none"> <li>• The problem domain has low complexity</li> <li>• We do not have completed knowledge about the rules for a concept and interrelations between concepts</li> <li>• Several rules are not complete/mature and they miss one or more concept to be completed</li> <li>• Many simple rules and few complex rules</li> <li>• Rules are decomposed into more rules</li> </ul>
High Many attributes per concept	<ul style="list-style-type: none"> <li>• High complexity</li> <li>• The rules are completed</li> <li>• We have good knowledge about the domain</li> <li>• There are very few but very complex rules</li> <li>• The rules are not decomposed or at least not at a high degree</li> </ul>

Table 6. Interpretations “Average concept in each rule”

### 2.9. Average Number of Attributes Included in Each Rule

This metric will be similar to the last one but it could give a better measure of the complexity of the domain especially in those cases where many rules are dependent of many attributes of few concepts. This metric will then indicate a high complexity where the previous one indicated low complexity. This metric will unfortunately still be dependent of the number of rules and the degree of decomposition.

### 2.10. $A \cdot \text{Average Number of Attributes in Rule} + B \cdot \text{Number of Rules} + C \cdot \text{Average Number of Decomposition Levels}$

To try to remove the dependencies from the previous metric we would suppose to combine attributes, rules and decomposition levels into one metric to better understand the complexity of the domain. The constants A, B and C must be found with use of historical data.

### 2.11. Average Number of Rules Each Concept Is Included in

One concept could be included in one but most likely more than one rule. The average number of rules a concept is included in could give us an indication of complexity. We expect it to increase throughout the project as more rules are made. If there is found a lot of new concepts it may decrease a bit. But in the end of the project we think it is more likely to find more rules than new concepts. If the number of concepts is very high the number of rules could be low and we could still have a very complex domain. At the end of the project we believe this metric should converge and thus it could be used as an indication of maturity. The table 7 shows our interpretations of this metric.

Result	Cause
Low Each concept is included in few rules	<ul style="list-style-type: none"> <li>• The domain is simple</li> <li>• The concepts of the domain is not strongly related</li> <li>• The knowledge about the problem area is sparse</li> <li>• We know all or may of the concepts of the area but we do not know all the relations yet</li> <li>• There are a lot of concepts without many rules</li> </ul>
High Each concept is included in many rules	<ul style="list-style-type: none"> <li>• The domain has many relations and it is complex</li> <li>• We have good knowledge about the domain</li> <li>• We may totally miss some concepts of the domain</li> </ul>

**Table 7.** Interpretations “Average rules each concept is in”

### 2.12. $A \cdot \text{Average Number of Rules Each Concept Is Included in} \cdot B \cdot \text{Number of Concepts}$

To remove the dependency of the number of concepts from the last metric we would propose to combine the previous metric with the number of concepts. The constants must, as mentioned, be found by use of history data.

### 2.13. Average Number of Rules Each Attribute Is Included in

We expect this metric to have a similar development during the project as the previous one with concepts. But we think it is more likely to discover more new attributes throughout the project than new concepts, so the value could vary a bit more than what we saw in Figure 6. We expect this value to converge at the end of the conceptualization phase as well. The table 8 shows our interpretations of this metric.

Result	Cause
Low Each attribute is included in few rules	<ul style="list-style-type: none"> <li>The domain is simple</li> <li>We do not have a mature knowledge base</li> <li>The domain is not strongly related</li> <li>We do not have a lot of knowledge about the domain</li> </ul>
High Each attribute is included in many rules	<ul style="list-style-type: none"> <li>The domain is strongly bound together</li> <li>We have good knowledge about the domain</li> <li>We miss many attributes which would decrease this average.</li> <li>We have good knowledge about just parts of the domain.</li> </ul>

Table 8. Interpretations “Average rules each attribute is in”

**2.14. For all Levels (Number of Decisions at Level i\*i) / Total Number of Decisions**

This metric will give an indication of the tree width of the decision trees. If the main decisions consist of many different decisions of if the decisions and the end of the tree are very detailed. We expect that the value of this metric will be increasing throughout the project and stabilize at some point between 1.0 and the depth of the tree. To better understand the metric please see example 1 in Figure 2 and example 2 in Figure 3.

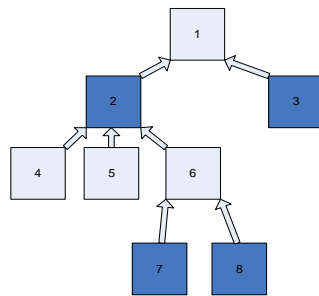


Fig. 2. Example 1: deep tree

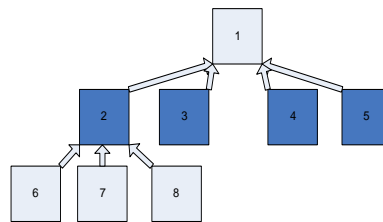


Fig. 3. Example 2: wide tree

With the number of decisions at a level, times the level, for instance 2 decisions at level 4 in Figure 8 will give 2\*4. The two examples in Table would give these results respectively:

Example	Result
1	$\frac{(1*1 + 2*2 + 3*3 + 2*4)}{8} = 2.875$
2	$\frac{(1*1 + 4*2 + 3*3)}{8} = 2.25$

Table 9. Result from examples

We see that the results indicate that the first tree is deeper than second one. We think this can help to show how the decisions in the problem domain are. This metric could give an indication of what kind of decision trees we have on thus what kind of complexity we have. The table 10 shows our interpretations of this metric.

Result	Cause
Low Most decisions at a high level	<ul style="list-style-type: none"> <li>The decisions are based on many decisions at a high level (close to the root of the tree).</li> <li>The decision process is not very complex</li> <li>We have not yet decomposed the tree</li> </ul>
High Many decisions are taken at a low level	<ul style="list-style-type: none"> <li>Few decisions are based on simple decisions. Most decisions contain many decisions at a low level.</li> <li>We have decomposed the tree</li> <li>The decision process is complex</li> </ul>

**Table 10.** Interpretations “Sum of Decision levels/number of decisions”

### 3. Applying the Metrics to Real World

To evaluate our metrics, we have used data from two finished expert systems. They were developed as part of the author’s master thesis at ITBA (see Tables 11 and 12).

#### System 1 Work Accidents

Reference	<i>Help Assistant on Work Risks in Argentinean Law.</i> (in spanish). Master Thesis on Knowledge Engineering. School of Computer Science. Politechnic University of Madrid. 2001.
Author	Paola V. Britos
Description	This system should help the user to search in the Argentinean laws for material regarding occupational accidents. A lot of time is spent by the lawyers to search for the right material and this system is meant to help them in their search.

**Table 11.** Description of system 1

#### System 2 Airport Control

Reference	<i>Expert System for Decision Making Training in an Information &amp; Control Air Traffic Center.</i> (in spanish). Master Thesis on Software Engineering. Graduate School. Buenos Aires Institute of Technology. 2002.
Author	Jorge Salvador Ierache
Description	The system described in this thesis is a decision support system for airport control towers.

**Table 12.** description of system 2

### 4. Some Results

We will here present the results from the expert systems described in the last section.

#### Number of Concepts, Number of Rules or Number of Attributes

System number	1	2
Number of Concepts	17	20
Number of Attributes	81	126
Number of Rules	472	155

These metrics are used as basis for other metrics. But they can also give an indication of the size of the system we have. We see that system 1 has quite many rules. This is because the system contains several simple rules concerning selection of the right document or right law to look up.



**Number of Concepts in a Rule / Number of Concepts**

System number	1	2
Number of Concepts in a Rule	7	19
Number of Concepts	17	20
Result	0.41	0.95

This metric indicates that system 1 has several concepts that are not related to anything and the knowledge engineer should therefore start working with those concepts.

**Number of Attributes in a Rule / Number of Attributes**

System number	1	2
Number of Attributes in a Rule	50	121
Number of Attributes	81	126
Result	0.62	0.96

We see the same indication here as we did with the last metric. System 1 needs to focus on those concepts and attributes not included in any rule or at least find the explanation of the result.

**Number of Concepts / Number of Rules**

System number	1	2
Number of Concepts	17	20
Number of Rules	472	155
Result	0.04	0.13

These resulting numbers are very small and it is hard to give some conclusions based on these numbers. But it could be interesting to follow the development of this figure throughout a project.

**Average Number of Attributes per Concept**

System number	1	2
Number of Attributes	81	126
Number of Concepts	17	20
Average	4.76	6.3

The number of attributes per concept can give us an indication of the complexity of the concepts in the domain. We observe that the result indicates that system 2's domain is more complex.

**Average Number of Levels in Decision Tree**

System number	1	2
Decomposed decisions	NA	NA
Average	NA	NA

Decision trees were not used to represent knowledge in these projects. The structure of the knowledge lead to omitting the application of this and other metrics concerning decomposed decisions.

**Average Number of Concepts Included in Each Rule**

System number	1	2
Average	1.24	1.64

We see that system two has more concepts included in a rule. This is an indication that system 2 may have a more complex domain.

**Average Number of Attributes Included in each Rule**

System number	1	2
Average	2.17	2.81

This metric is very similar to the previous one and it indicates the same. The domain of system 2 is more complex than the one of system 1.

**A\*Average Number of Attributes in Rule + B\*Number of Rules + C\* Average Number of Decomposition Levels**

System number	1	2
Attributes in rule	2.17	2.81
Number of rules	472	155
Average decomposition levels	NA	NA
Sum	NA	NA

We will use all the constants set to 1 since we do not have any historical data from previous projects. Decision trees were as mentioned above not used in any of the projects. Because of that we omitted applying this metric.

**A\*Average Number of Rules each Concept Is Included in\*B\*Number of Concepts**

System number	1	2
Average rules each concept is in	34.5	15
Number of concepts	17	20
Sum	586.5	300

We have also used 1 for the constants in this metric since we do not have any historical data so far. These results indicate that domain 1 is a bigger domain with several relations.

**Average Number of Rules each Attribute Is Included in**

System number	1	2
Result	12.6	3.45

We see the same here as we did in the two last metrics. System 1 has more relations between the attributes then system 2.

**Average Number of Rules Each Concept Is Included in**

System number	1	2
Average	34.5	15.0

We see that system 1 has more discovered rules in average that system 2. This could be an indication of fewer relations in domain 2.

**For all Levels (Number of decisions at level i\*i) / Total Number of Decisions**

System number	1	2
Result	NA	NA

Decision trees were unfortunately not used and applying this metric was omitted.

## 5. Conclusions

The intention of this paper was examining the problem domain and showing the need for metrics in this domain. The metrics were suggested with a theoretical background to create a discussion around use of metrics in the conceptualization phase of an expert system development. We applied most of the proposed metrics to two different expert systems. This is not a large enough data set to draw any statistical conclusions. At this point the metrics serve as indicators and the trend seems to be that system 2 has a more complex domain than system 1. This seems reasonable enough. System 2 is an airport control system and system 1 is a system for finding the right law or text concerning accidents at work. The metrics also show concepts which are not included in any rule. This should alert the knowledge engineer and tell him to focus on these concepts. The application we made can also guide the knowledge engineer in finding unused concepts, attributes or rules where no attributes were found and tell him to review these rules.

## 6. Referentes

- Firestone, J. 2004. *Knowledge Management Metrics Development: A Technical Approach*. Published on-line by Executive Information Systems, Inc. <http://www.dkms.com/> at: <http://www.dkms.com/papers/kmmeasurement.pdf>, downloaded 2004-10-02
- Ford, G. 2004. *Measurement Theory for Software Engineers*, Published on-line by R.S. Pressman & Associates, Inc. <http://www.rsps.com/> at: <http://www2.umassd.edu/SWPI/curriculummodule/em9ps/em9.part3.pdf>, downloaded 2004-09-24
- García Martínez, R. & Britos, P. 2004 *Expert System Engineering*, Nueva Librería Ed. Buenos Aires
- Kang, Y. & Bahieel, T. 1990. *A Tool for Detecting Expert Systems Errors*. *AI Expert*, 5(2): 42-51.
- Menzies, T. & Cukic, B. 1999. *On the Sufficiency of Limited Testing for Knowledge Based Systems*. Proceedings of the 11th IEEE International Conference on Tools with Artificial Intelligence, Pages 431-440.
- Menzies, T. & Cukic, B. 2000. *Adequacy of Limited Testing for Knowledge Based Systems*. *International Journal on Artificial Intelligence Tools* 9(1): 153-172.
- Menzies, T. 1999. *Critical success metrics: evaluation at the business level*. *International Journal of Human-Computer Studies*, 51(4):783-799.
- Nilsson, N. 1998. *Artificial Intelligence: A New Synthesis*, Morgan Kaufmann Publishers
- Pazzani, M. & Clifford, A. 1991. *Detecting and Correcting Errors in Rule-Based Expert Systems: An Integration of Empirical and Explanation-Based Learning*. *Knowledge Acquisition* 3 :157-173.
- SEI. 2004. *Software Metrics, SEI Curriculum Module SEI-CM-12-1.1*, Carnegie Mellon University - Software Engineering Institute, December 1988, <ftp://ftp.sei.cmu.edu/pub/education/cm12.pdf>, downloaded 2004-09-22.

# Toward developing a tele-diagnosis system on fish disease

Daoliang Li<sup>1\*</sup>, Wei Zhu<sup>1</sup>, Yanqing Duan<sup>2</sup>, Zetian Fu<sup>1</sup>

<sup>1</sup> Key Laboratory of Modern Precision Agriculture system Integration, P.O. Box 121, China Agricultural University, Beijing, 100083, P. R. China

<sup>2</sup> Luton Business School, University of Luton, LU1 3JU, UK

**Abstract.** Fish disease diagnosis is a complicated process and requires high level of expertise, an expert system for fish disease diagnosis is considered as an effective tool to help fish farmers. However, many farmers have no computers and are not able to access the Internet. Telephone and mobile uses increase rapidly, so, the provision of call centre service appears as a sound alternative support channel for farmer to acquire counseling and support. This paper presents a research attempt to develop and evaluate a call center oriented Hybrid disease diagnosis & consulting system (H-Vet) in aquaculture in China. This paper looks at why H-Vet is needed and what are the advantages and difficulties in the developing and using such a system. A machine learning approach is adopted, which helps to acquire knowledge when enhancing expert systems with the user information collected through call center. This paper also proposes a fuzzy Group Support Systems (GSS) framework for acquiring knowledge from individual expert and aggregating knowledge into workgroup knowledge by H-Vet in the situation of difficult disease diagnosis. The system's architecture and components are described.

**Keywords** machine learning; Group Decision Support System, expert system, call centre

## 1 Introduction

In China, Aquaculture plays a very important role in agricultural structure adjustment and generating farmers' income (Guo, 2001). However, fish diseases have become one of the most devastating threats to the survival of many Chinese fish farms. Fish disease diagnosis is a complicated process and requires high level of

\* Corresponding author. Tel: +86-10-62336717; Fax:+86-10-62324371.  
Email address: li\_daoliang@yahoo.com or dliangl@cau.edu.cn (D. Li)

domain knowledge, this pose a major challenge for any attempts to provide accurate and timely diagnosis and treatment.

Advances of Internet technologies have offered new opportunities for enhancing traditional decision support systems and expert systems (Power, 2000). With the development of Expert Systems (ES) and multimedia, computers are able to mimic many important roles that normally require human actions. A number of Web-based expert systems are reported in the literature (Grove, 2000; Potter et al., 2000; Riva, Bellazi, & Montani, 1998; Sedbrook, 1998).

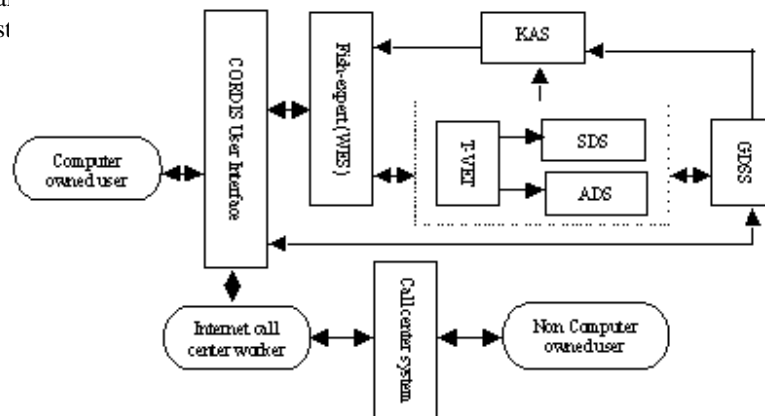
Developments in high performance global communication technologies have also accelerated cooperative image-based medical services to a new frontier. Traditional image-based medical services such as radiology and diagnostic consultation can now be fully enhanced by multimedia technologies to provide novel services, such as tele-medicine. Telemedicine can be defined as medical practices across distance via telecommunications and interactive video technology. It can be used in remote areas or across great distances on the globe (Lim, Pang, & Tan, 2001). It also covers any form of communication between health workers and patients through electronic equipment from remote locations. Similar to tele-medicine, a tele-diagnosis system has been developed to support the diagnosis of various types of problems in China (Duan, et al.,2003), which provide a farmer-to-vet communication in diagnosis.

However, many farmers have no computers and are not able to access the Internet. At the same time, telephone and mobile uses increase rapidly, especially in rural area of China, there are 0.24 billion telephone users, 40% of them are farmers, so, the provision of call centre service appears as a sound alternative support channel for farmer to acquire counseling and support in disease diagnosis in China.

This paper reports a research effort to integrate call center with web-based expert system and tele-diagnosis system in fresh water fish disease diagnosis in China. The system is being developed by Key Laboratory of Modern Precision Agriculture System Integration, China Agricultural University, and was funded by Huo Yingdong foundation in China. The system has emerged as a result of the use of a Web-based expert system called Fish-Expert and tele-diagnosis system called T-Vet. It helps to overcome the limitations and enhance the functionality of traditional ESs. One of the strengths of the system is that it can facilitate computer-not-owned users and knowledge remote acquisition. This research attempts to use this pilot system as a research vehicle to experiment in applying, and to evaluate the usability of, the system with potential users. Feedback collected from the demonstration and evaluation of the Fish-Expert and T-vet has provided valuable insights into the issues related to the development and implementation of H-Vet in China.

## 2 System architecture and components

Based on the user's needs analysis, a call center oriented tele-diagnosis system, called CORDIS, were developed and integrated to Fish-Expert and T-Vet, H-Vet allow fish farmers and technicians get online help in any situations, the system st



**Fig. 1.** Structure and Subsystem in H-Vet

To best meet the different needs of fish farmers and vets, 5 subsystems have been designed and developed in H-vet, they are Web-based Expert System Fish-Expert (WES), Call Center System (CCS), Tele-diagnosis System (T-VET), Group Decision Support System (GDSS), and Knowledge Acquisition System (KAS).

As more details about web-based expert system Fish-Expert and tele-diagnosis system (T-VET) are described in Li, et al. (2002) and Duan, et al. (2003), Fish-Expert and T-vet will be described very shortly. More contents will cover the CCS, GDSS, and KAS.

### **2.1 Web-based expert system Fish-Expert (WES)**

The Fish-Expert system mimics the diagnosing process of human experts and has over 300 rules and 400 images and graphics for different types of diseases and symptoms. It is able to diagnose 126 types of diseases amongst nine species of primary freshwater fish. When using the expert system, various information needs to be provided following different diagnosing steps, such as pond inspection, fish inspection and anatomization, water quality examination, and microscopic examination. A final verdict for the type of the disease and its treatments and prevention will be produced based on the system's knowledge base.

### **2.2 Tele-diagnosis system T-Vet**

T-vet is an add-on subsystem to overcome the limitation of the expert system and was integrated to the Fish-Expert. It includes a synchronous tele-diagnosis subsystem (SDS) and an asynchronous diagnosis subsystem (ADS) (Duan, et al., 2003)

A synchronous tele-diagnosis subsystem has been developed and can be used in situations where an urgent diagnosis is required, but it is impossible for the fish vet to visit the site, and the Web-based expert system is not able to solve the problem. To facilitate the tele-diagnosis, a number of functions have been developed and integrated into the system, such as web-calling, a virtual diagnosis room, Computer Supported Cooperative Work (CSCW) module, video/audio conferencing module, and online help module. An asynchronous diagnosis system (ADS) has been developed, which acts as a practical platform for sending and receiving messages between a fish farmer and a vet. Three support modules—user symptom submission, vet diagnosis and email communication are designed to facilitate asynchronous diagnosis process.

The ADS and the SDS act as good complementary tools to the expert system. The integration of the expert system, synchronous and asynchronous systems complements each other and is able to solve most of the problems fish farmers may encounter.

### 2.3 Call Center System (CCS)

The Call Center System CSS is an add-on subsystem to overcome the limitation of the current low computer owned level. It can provide a bridge between the computer-not-owned fish farmer and Fish-expert, T-Vet through call center agent,. There are 3 main models in the system, such as queuing models to capture the impact of congestion, customer arrival statistics and data collection model, telecommunications resource allocation and telephone-agent staffing model.

As call centers have grown in number and in size, more firms have tried to improve their management by focusing on resource utilization and service levels. This has led to a series of studies dealing with the problem of staffing phone centers, many of which have made use of queuing models to capture the impact of congestion (Aksin & Harker, 2003). So queuing models should be designed to capture the impact of congestion.

Customer arrival statistics and corresponding data collection model is used to evaluate the performance of the service system under study, given customer arrival statistics, servers, buffers, and a shared resource that impacts processing times, to collect all disease case which will be used to acquire fish disease diagnosis knowledge.

Another model take a different approach to telecommunications resource allocation and telephone-agent staffing, a similar approach is taken to determine staffing levels for a multiple class inbound call center.

CCS not only provides a bridge between the computer-not-owned fish farmer and Fish-Expert, but also provides a tool for collecting fish disease diagnosis case which plays a very important role in Fish-expert.

### 2.4 Group Decision support System (GDSS)

Globalization, virtual corporations, telecommuting, empowerment of teams, reduced cycle time and the need to frequently make decisions quickly makes it necessary for groups to work together while the participants may be in different locations (Tung & Turban, 1998). Distributed Group Decision Support System (DGSS) is a technology that can help groups to overcome some of the difficulties associated with being in different places and sometimes in different time zones (Bendoly & Bachrach, 2003).

Fish disease diagnosis is a rather complicated process in aquaculture production activities. The disease commonly resulted from nutritional and environmental problems as well as infections by parasites, viruses, bacteria and fungal agents. Some rare diseases or new diseases normally can't be identified by one fish vet only, and most of them need to be diagnosed by group work, as result, a Group Decision Support System (GDSS) is essential needed to solve the rare disease diagnosis.

Most of the previous research regarding computer support of groups was related to the decision room environment, where a group of participants meet face-to-face, working on a common task (Stohr & Konzynski, 1992). This paper focuses on the technology to support group work in the framework of 3 situations, such as same-time same-place group work, different-place same-time and different-time different-place group work.

Both same-time same-place group work and different-place same-time group work are belong to synchronous GDSSs, the same-time same-place group work is a face-to-face communication work mode for the decision support. This kind of group work has no any limitation of the work environment, the fish vet can discuss together, how ever the shortcoming of this kind of group work all famous vet must be collected together in call center, and answer the user's questions for some rare disease diagnosis.

Different-place same-time group work is synchronous GDSSs allow distributed participants to interact with one another in a 'real time' mode, i.e., they interact with one another at the same time. The participants are distributed across multiple sites linked by various communication technologies. Some of the supporting technologies are screen sharing, whiteboard, audio-conferencing, and various types of video-conferencing. These technologies can be carried on the Intranet, Internet, corporate or public networks, or VANs. Some major issues here include: the loss of face-to-face contact and the ability to manage the group process.

Asynchronous DGSSs allow distributed participants to log into the same meeting but at different times. Participants can log in and catch up with what is going on in the meeting, enter comments if necessary, and log out of the meeting at various times. E-mail, voice-mail, and workflow management systems such as in Lotus Notes are some of the supporting technologies. The issues here are more than just time and distance barriers. For example, the control of the participants and the participation, and the delayed response time could play an important role. Therefore, attention must be given to the coordination of group members to ensure that they stay on task and track and meet the decision deadlines, as well as it is necessary to encourage timely participation by everyone.

As the rare fish disease cases are very short in the fish disease case base of the web-based expert system, all result of group decision will be acquired by the KAS.

## **2.5 Knowledge Acquisition System (KAS)**

Call center system, Fish-expert, T-vet, and GDSS can provide plenty of successful fish disease cases and solutions, these cases and solutions are very important for the Fish-Expert, so how to integrate them together and make them cooperate together, and then get a best effect for the whole system poses a serious challenge.

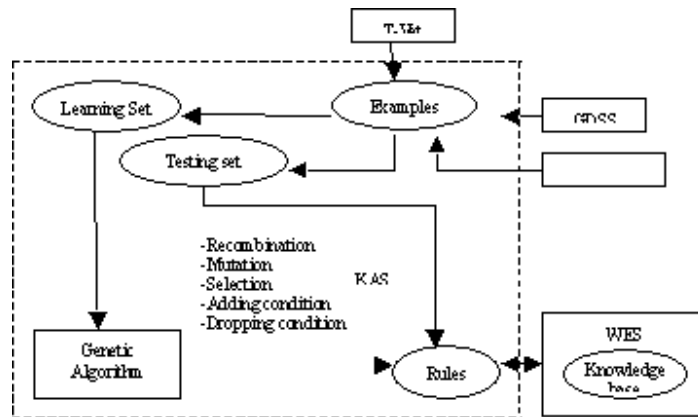


Fig. 2. The Structure of KAS

Case-based reasoning system is a system that solves current problems by adapting to or reusing the used solutions to solve past problems. In case-based reasoning, the study about case adaptation can be divided into two areas. One is about adaptation method and adaptation knowledge that are used at adaptation stage. The other is to reduce the necessity of case adaptation by extracting the most appropriate case for current problems. The former focuses on how to acquire case adaptation knowledge, and the latter focuses on the necessity about case adaptation.

However, the bottleneck phenomenon in acquiring knowledge may be caused in case-based systems. To resolve the bottleneck at the case-based reasoning, a research that automatically acquires knowledge needed for case-based reasoning by using the technologies of machine learning or data mining has been studied (Fig. 2. )

To achieve this intention, we propose a Knowledge Acquisition System (KAS) of adaptation knowledge from derived cases. That is, we construct case base by acquiring cases from Call Center, T-Vet, GDSS automatically acquire adaptation knowledge exploiting data mining concept, and deduce the bottleneck of acquiring adaptation knowledge used at the case adaptation stage of case-based reasoning systems.

### 3 The diagnosing process

There are 4 kinds of users in the H-Vet, such as computer-not-owned fish farmer, computer- owned fish farmer, vet, and call center agent, the work process for all kinds of users can be seen in fig. 3.

For computer-not-owned fish farmer, he or she can describes the symptoms of their fishes' disease to the call center agent by telephone, and the call center agent inputs all these symptoms into the Fish-Expert interface, and tell the fish farmer the diagnosis result based on the web-based expert system if there is a same disease case in the case base, meantime, this case will be input in the knowledge acquisition system as an successful case. the agent also can ask some questions to the fish farmer based on the web-based expert system, which can add some useful information for fish disease diagnosis.



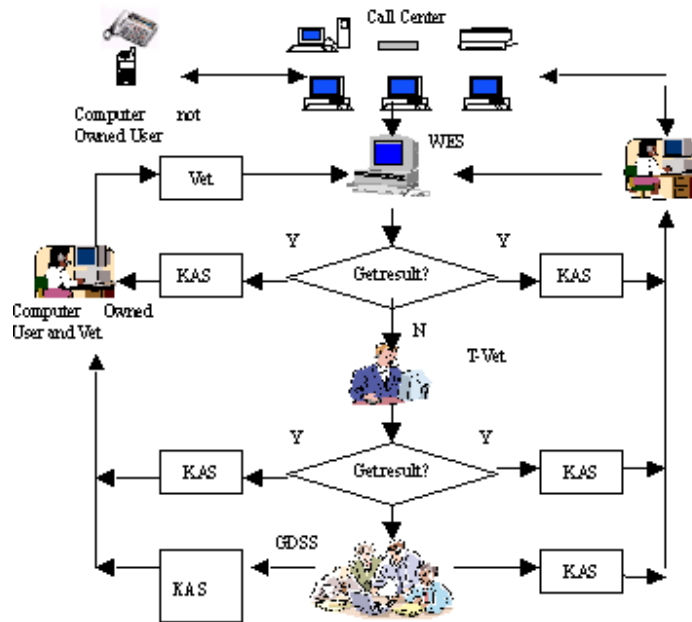


Fig. 3. The work process of H-Vet

The agent can login on the T-Vet (synchronous tele-diagnosis subsystem) interface get as on line help from the vet, and then tell result to the fish farmer, at the same time, this case will be input into the knowledge acquisition system as an successful case. If the agent can't reply the agent's question, the agent will login the GDSS, and submit the question the GDSS, and reply the fish farmer's question base don the GDSS's diagnosis, at the same time the question and the result will be input the knowledge acquisition system as a successful disease diagnosis case.

For computer-owned fish farmer, he/she can login the web-based expert system directly, and use to the Fish-Expert to solve his or her question, if the Fish-Expert can't answer the farmer's question, the farmer can login T-Vet system and get a on-line help from vet or get a synchronous tele-diagnosis help, he can login the GDSS system if T-vet can't help him/her. All successful disease diagnosis will be input the knowledge acquisition system as a successful case during the whole work process.

For call center agent, the main task is to wait computer-not-owned fish farmers' call, and answer their questions base on the Fish-Expert, T-vet and GDSS, the interface and work process for agent same as computer-owned fish farmer.

For Vet, the main task is to answer the questions which the expert system can't solved, there are 2 interfaces for them, one is T-Vet, another is GDSS, the participation of the vet will provide many practical disease diagnosis, which will be added to the case base.

### 4 Implementation

System tests, such as logic tests, debugging, rule checking and sample field tests were carried out by system developers. This was to ensure the system would work correctly before it was distributed to farmers. After the system testing, H-Vet available for pilot implementation in North China, in cities such as Beijing, Tianjin, and Shandong provinces. User feedback was gathered by conducting interviews and collecting information through the system’s built-in visitor feedback form. In general, the system has been an effective aid to fish farmers, fishery experts and a reference system to fish vets. Some interfaces are as follows. (Fig. 4.-7)



Fig. 4 Interface of Call Center



Fig. 5 Interface of Agent



Fig. 6 Interface for Computer Owned User



Fig. 7 Interface of T-Vet

### 5 Conclusions

This paper reports a research attempt in developing a hybrid tele-diagnosis system in aquaculture, the system provides a tool for computer and Internet not owned user using web-based expert system, for knowledge automatically acquisition. The research provides some new ideas to solve the bottleneck of traditional expert system.

The wide spread of Internet, Intranet, call center, telecommunication infrastructures and intelligent software will facilitate the work of hybrid tele-diagnosis system, the frame work proposed in this paper and the specific issues list in section 3 and 4 cover the major topics perceived by the researchers as warrant further

investigation. As our H-Vet knowledge base increase, more research issues will undoubtedly be added.

The next phase in our research could be the refinement of the framework and testing its value by researchers and practitioners. Also, the identification of software tools that will enhance the multi-site and interorganizational meeting process, and the development of procedures and training methods that will help organizations take advantage of the benefits of using a DGSS need to be researched. The data-mining algorithm will be a key work in the future research, which will play a very important role in knowledge refining and acquisition.

The DGSS environment offers many potential areas for investigation,. The research framework and issues raised in this paper are intended for researchers and practitioners who are interested in looking at the impact of the DGSS environment on fish farmers, vet groups, and call center agents.

### **Acknowledgement**

The research was funded by the Huo Yingdong foundation in China (project number: 94032). We would like to thank many domain experts from the Beijing Aquaculture Science Institute, Aquaculture Department of Tianjin Agricultural College, Aquaculture Bureau of Shandong province, for their co-operation and support. Our special thanks should also go to Prof. Kezhi Xing, Mr Yongjun Guo at Tianjin Agricultural College for his valuable suggestions and comments on the system.

### **References**

1. Aksin, O. Z., Harker, P.T.( 2003). Capacity sizing in the presence of a common shared resource: imensioning an inbound call center. *European Journal of Operational Research* 147, 464 – 483
2. Bendoly, E., Bachrach, D. G., (2003). A process-based model for priority convergence in multi-period group decision-making. *European Journal of Operational Research* 148 (2003) 534 – 545
3. Duan, Y., Fu, Z., Li, D. (2003). Toward developing and using Web-based tele-diagnosis in aquaculture. *Expert System with Applications*, 25, 247–254
4. Grove, R. (2000). Internet-based expert systems. *Expert Systems*, 17(2), 129–135.
5. Guo, Z (2001). Analysis and prospects for China's aquatic market in 2000. [www.ifishery.com/jrdd/2001060601.htm](http://www.ifishery.com/jrdd/2001060601.htm)
6. Li, D., Fu, Z., & Duan, Y. (2002). Fish-expert: A web-based expert system for fish disease diagnosis. *Expert System with Applications*, 23(3), 311–320.
7. Lim, S. S., Pang, Y. K., & Tan, H. S. (2001). Telemedicine VSATapplications: Temasek polytechnic experience. *International Telehealth Symposium*, 14 March, Bangkok, United Nations Conference Center.
8. Magrabi, F., Lovell, N. H., & Cellar, B. G. (1999). A Web-based approach for electrocardiogram monitoring in the home. *International Journal of Medical Informatics*, 54, 145–153.

9. Potter, W. D., Deng, X., Li, J., Xu, M., Wei, Y., Lappas, I., Twery, M. J., & Bennett, D. J. (2000). A Web-based expert system for gypsy moth risk assessment. *Computer and Electronics in Agriculture*, 27(1–3), 95–105.
10. Power, D. J. (August 10–13 2000). Web-based and model-driven decision support systems: Concepts and issues. Proceedings of Americas Conference on Information Systems (AMCIS 2000), Long Beach, California.
11. Rajani, R., & Perry, M. (1999). The reality of medical work: The case for a new perspective on telemedicine. *Virtual Reality Journal*, 4, 243–249.
12. Riva, A., Bellazzi, R., & Montani, S (1998). A knowledge-based web server as a development environment for Web-based knowledge servers. IEE Colloquium on Web-based Knowledge Servers (Digest No.1998/307).IEE. 1998, 5/1-5. London, UK
13. Sedbrook, T. A. (1998). A collaborative fuzzy expert system for the Web. *Data Base for Advances in Information Systems*, 29(3), 19–30.
14. Tung, L., Turban, E. (1998). A proposed research framework for distributed group support systems. *Decision Support System*, 23, 175-188
15. Stohr, E.A., Konzynski, B.R. (1992) *Information Systems and Decision Process*, IEEE Computer Society Press, Los Alamitos, CA.

# A new method for fish-disease diagnostic problem solving based on parsimonious covering theory and fuzzy inference model

Jiwen WEN<sup>1</sup>, Daoliang LI<sup>2</sup>, Wei Zhu<sup>2</sup> Zetian FU<sup>2</sup>

1 Economics and Management department, Beijing Forestry University, Beijing, China, 100083

2 China Agricultural University, Key laboratory for modern precision agriculture integration, Ministry of Education, Beijing, China, 100083

**Abstract.** There are three kinds of uncertainty in the process of fish-disease diagnosis, such as randomness, fuzzy and imperfection, which affect the veracity of fish-disease diagnostic conclusion. So, it is important to construct a fish-disease diagnostic model to effectively deal with these uncertainty knowledge's representation and reasoning. In this paper, the well-developed parsimonious covering theory capable of handling randomness knowledge is extended. A fuzzy inference model capable of handling fuzzy knowledge is proposed, and the corresponding algorithms based the sequence of obtaining manifestations are provided to express imperfection knowledge. In the last, the model is proved to be effective and practicality through a set of fish-disease diagnostic cases.

**Key words** Fuzzy set theory, Parsimonious covering theory, Fish-disease diagnosis

## 1 Introduction

In this paper, the well-developed parsimonious covering theory based probability model is extended. A alternative fuzzy inference model capable of handling fuzzy knowledge for fish-disease diagnostic problem solving is proposed, and the corresponding algorithms based the sequence of obtaining diagnostic information are provided. This model can deal with degrees of manifestations simultaneously and sequences of obtaining diagnostic information.

Firstly, PCT and its probability model are introduced into fish-disease diagnosis domain. Secondly, the best probability criterion is determined on through analyzing fish-disease's characteristic, then the plausible function is calculated and the best solution is founded by thoroughly enumerate strategy and the biggest probability criterion. Thirdly, the limitation of the probability model of PCT solving fish-disease diagnosis is pointed out. So that, fuzzy set theory is introduced to represent degrees

of manifestations and the diagnostic problem is newly defined. Finally, responding algorithm is proposed to deal with the imperfection of diagnostic information.

PCT (Parsimonious Covering Theory) based on probability model is extended in such a way that it is integrated with fuzzy set theory and corresponding sequence algorithm in this paper. The new model is applied into fish-disease diagnosis domain to deal with three-uncertainty knowledge's representation and reasoning. In the last, the model is proved to be effective and practical through a set of fish-disease diagnostic cases.

## 2 The fish-disease diagnostic model based on parsimonious covering theory

### 2.1 Parsimonious covering theory

Parsimonious covering theory (PCT) uses two finite sets to define the scope of diagnostic problems. They are the set  $D$ , representing all possible disorders  $d_i$  that can occur, and the set  $M$ , representing all possible manifestations  $m_j$  that may occur when one or more disorders are present. The relation  $C$ , from  $D$  to  $M$ , associates each individual disorder with its manifestations. An association  $\langle d_i, m_j \rangle$  in  $C$  means that  $d_i$  may directly cause  $m_j$ . To complete the problem formulation we need a particular diagnostic case. We use  $M^+$ , a subset of  $M$ , to denote the set of observations, which are manifestations present in a particular diagnostic case.

Definition 1. A diagnostic problem  $P$  is defined as a quadruple  $\langle D, M, C, M^+ \rangle$ , where:

$D$  is a finite, non-empty set of elements, called disorders; and  $M$  is a finite, non-empty set of elements, called manifestations;

$C$  is a binary relation between disorders and manifestations, called causation; it represents with a matrix form. A nonzero element  $c_{ij}$  means the probability that  $d_i$  may directly cause  $m_j$ ;

$M^+$  is a subset of  $M$ , and it identifies all observed manifestations.

Definition 2. The set  $D_i \subseteq D$  is the cover of  $M_j \subseteq M$  if  $M_j \subseteq \text{Effect}(D_i)$

Definition 3. A set  $D_i^* \subseteq D$  is an explanation of  $M^+$  for a diagnostic problem if  $D_i^*$  covers  $M^+$ , and satisfies a given parsimony criterion: simple, minimal, irredundant, and relevant criterion.

Definition 4. The solution of a diagnostic problem  $P$  is the set of all explanations of  $M^+$  and it must satisfy a given parsimony criterion.

The abductive inference model is based upon the notion of parsimoniously covering a set of observed manifestations,  $M^+$ . The premise of the parsimonious covering theory is that a diagnosis hypothesis must be a cover of  $M^+$  in order to account for the presence of all manifestations in  $M^+$ . On the other hand, not all covers of  $M^+$  are equally plausible as the hypotheses for a given problem. It must satisfy a given parsimony criterion: simple, minimal, irredundant and relevant criterion.

## 2.2 Decision about parsimony criterion in fish-disease diagnostic model

The key to solve disease diagnostic problem is to select a preferable parsimony criterion combined with diagnostic practice. In practice, if the diagnosis objective is different, parsimony criterion is too different.

(1) If the restriction of diagnostic object is single disorder, a cover  $D_I$  of  $M^+$  is an explanation if it contains only a single disorder, a single disorder principle should be adopted.

(2) If diagnostic object is a many-disorder system, parsimony criterion is adopted on the characteristic of disorder occurrence: relevancy, a cover  $D_I$  of  $M^+$  is an explanation if it only contains the disorders that causally associate with at least one of the manifestations in  $M^+$ ; irredundancy, a cover  $D_I$  of  $M^+$  is an explanation if it has no proper subsets which also cover  $M^+$  or, in other words, removing any disorder from  $D_I$  results in a noncover of  $M^+$ ; minimality, a cover  $D_I$  of  $M^+$  is an explanation if it has the minimal cardinality among all covers of  $M^+$ , i.e. it contains the smallest possible number of disorders needed to cover  $M^+$ .

So, we must firstly analyze the characteristic of fish-disease occurrence, it is summed as follows:

(1) The relationship between fish-disease and fish-disease is various: Two diseases may repel one another, and two diseases may be interrelated. If the relationship of two fish-diseases is interrelated, one disease may be occur directly, indirectly or along with another disease.

(2) The cause of one fish-disease is different in different environment: Many factors, such as intrinsic body, extrinsic environment and man-made management, will result in fish-disease. In the condition of modern intensive breeding, these factors are more easier to result in fish-diseases occurrence more than before.

(3) The diffusion of fish-diseases is various and interrelated: Any original fish-disease is likely to arise other fish-diseases in many potential ways. Any fish-disease likely arises another disease and is likely aroused by another disease too.

In diagnostic practice, two status, many-disease occurrence and simple disease occurrence, are in existence, so irredundancy criterion is adopted to define the scope of diagnostic solution.

## 2.3 The limitation of PCT's probability model

If the set of observed manifestations,  $M^+$ , is known, plausible function of these solution, disease sets which can explain  $M^+$ , can be found by probability theory, the best solution can be calculated by thoroughly enumerating strategy and the biggest plausible criterion. Plausible function  $L(D_I, M^+)$  represents as follow:

$$L(D_I, M^+) = \prod_{m_j \in M^+} (1 - \prod_{d_i \in D_I} (1 - c_{ij})) \prod_{d_i \in D_I} \frac{p_i}{1 - p_i}$$

The limitation of PCT's probability model in the process of diagnosing fish-disease is summed as follow:

(1) The probability of  $d_i$ ,  $p(d_i)$ , is supposed to be a known numeral in diagnostic model, in fact, the value of  $p(d_i)$  is affected by many factors; it transforms along with factor's variation.

- (2) Degrees of cause-and-effect relation,  $c_{ij}$ , is supposed to be known in diagnostic model, in fact, the value of  $c_{ij}$  is different if the diagnostic expert is different, so it need to statutes.

So, fish-disease diagnostic knowledge must be pretreated before the diagnostic model is applied into fish-disease diagnostic domain.

#### 2.4 Pretreatments with fish-disease knowledge

Because fish-disease diagnostic system is complex, the factors that arise disease are different, but also the knowledge that experts hold is different because their diagnostic experience and the cognition degree of fish disease is different, so the diagnostic knowledge come from different channels constantly differs. It is necessary to foreclose this fish-disease knowledge in some statistic methods before it is introduced into PCT's probability model.

##### 2.4.1 Decision about the scope of fish-disease

In this paper, fresh water fish-diseases occur in Tianjin is illustrated as examples. Some constantly occurring fish-diseases are summed up according to data provided by Tianjin aquatic company, and the occurring frequency is classified into four levels: 1) frequently; 2) often; 3) infrequently; 4) nonoccurrence. We defined the scope of fish-disease by which occurring frequency is bigger than 0.2.

##### 2.4.2 Decision about the foresight probability of fish-disease

We consider that some factors are interrelated with fish-disease by analyzing a lot of investigation data; they are fish-kind, fish-age, occurring season, occurring area and pool-depth. The foresight probability can be decided by firstly confirming these environmental factors.

1. The first factor is fish-kind: If fish-kind is different, kinds of occurring fish-disease are too different. For example, some disease occurs in carp would not occurs in silver carp. We classify fish-kind into six kinds, such as grass carp, silver carp, variegated carp, carp, crucian carp and so on. Variable  $k_i$  represents the corresponding kind; function  $cause(k_i) = \{d_j | j = 1, 2, \dots, n\}$  represents disease sets when fish-kind is  $k_i$ .

2. The second factor is fish-age: If fish-age is different, kinds of occurring fish-disease are too different. For example, some disease occurs in parent fish would not occur in fish-breed. We classify fish-age into four phases, such as fry, fish-breed, grow-up fish, parent fish and so on. Variable  $a_i$  represents the corresponding age, function  $cause(a_i) = \{d_j | j = 1, 2, \dots, n\}$  represents disease sets when fish-age is  $s_i$ .

3. The third factor is season: If season is different, kinds of occurring fish-disease are too different. For example, some disease occurs in spring would not occur in summer. We classify occurring season into: spring, summer, autumn and winter four phrases, variable  $s_i$  represents the corresponding season, function  $cause(s_i) = \{d_j | j = 1, 2, \dots, n\}$  represents disease sets when occurring season is  $s_i$ .



4. The fourth factor is area: If area is different, kinds of occurring fish-disease are too different. For example, some diseases occur in certain area will not occur in other area. Variable  $pl_i$  represents corresponding area name, function  $cause(pl_i) = \{d_j \mid j = 1, 2, \dots, n\}$  represents disease sets when occurring area is of  $pl_i$ .

5. The probability of fish-disease,  $p_i$ , can be calculated by algorithm as follow:

$F_1^+$ : A subset of  $F^+$ , it defines factors, such as fish-kind ( $k_i$ ), fish-age ( $a_i$ ), occurring season ( $s_i$ ), occurring area ( $p_i$ ) and so on;

$D_l$ : Disease sets that can explain the observed manifestations;

The beginning state is that  $F_1^+ = \Phi, D_l = D$ ; the terminating condition is that  $F_1^+ = F^+, F^+ = \{k_i, a_i, s_i, p_i\}$ . Diagnostic algorithm can be described as follow:

- (1) input  $k_i$  from  $F^+$ ;
- (2)  $F_1^+ = F_1^+ \cup k_i$ ;
- (3)  $D_l = D_l \cap cause(k_i)$ ;
- (4)  $F_1^+ = F_1^+ \cup a_i$ ;
- (5)  $D_l = D_l \cap cause(a_i)$ ;
- (6)  $F_1^+ = F_1^+ \cup s_i$ ;
- (7)  $D_l = D_l \cap cause(s_i)$ ;
- (8)  $F_1^+ = F_1^+ \cup p_i$ ;
- (9)  $D_l = D_l \cap cause(p_i)$ ;
- (10) extract  $D_l$  and the number of disease in  $D_l$ ;
- (11) the probability of each disease should be  $p_i = 1/I$ .

#### 2.4.3 Decision about degrees of cause-and-effect relationship $c_{ij}$

In the process of diagnosing fish-disease, degrees of cause-and-effect relation should keep to principles as follow:

1) Because each disease can represent many manifestations, they are classified three levels: necessarily appearing, constantly appearing and occasionally appearing, the level of manifestation should be confirmed at first.

2) Necessarily appearing manifestation plays a decisive role in disease diagnosing, so the value of  $c_{ij}$  is defined as 1.0;

3) Constantly appearing manifestation frequently appears in process of certain fish-disease occurring, which can affect the reliability of diagnostic solution, so the value of  $c_{ij}$  is defined between 0.3 and 0.9;

4) The specialty of occasionally appearing manifestation is probably occurring; the value of  $c_{ij}$  should be between 0 and 0.2.

The equation is  $\bar{c}_{ij} = \frac{1}{n} \sum_1^n c_{ij}$  by analyzing questionnaires of 24 experts in Tianjin.

In the process of diagnosing fish-disease, there are three uncertainties: probability, fuzzy and imperfection. The model describing as up can deal with probability, but degrees of manifestation and the imperfection and sequence of diagnostic information are not taken into account. So we must improve on this model.

### 3 Diagnosis model based on fuzzy theory and diagnosis algorithm

#### 3.1 Diagnosis model based on fuzzy theory

If degrees of observed manifestation are introduced into parsimonious set covering, diagnosis problem is redefined as follow:

Definition 5 “manifestation-disease” diagnosis problem is defined as a quadruple  $(D, M, R, M^+)$ , where:  $M^+ = \{m_1, m_2, \dots, m_k\}$  is a fuzzy subset of  $M$ , and it identifies the observed manifestations. To each  $m_j \in M^+$ ,  $\mu(m_j) \in [0,1]$  is the observed degree of  $m_j$ .

In accordance with the diagnostic problem definition written as up, plausible function of diagnostic solution is defined as follow:

Definition 6  $L(D_I, M^+)$  expresses the degree that  $D_I$  explains the set of manifestations  $M^+$ , and the equation is tenable:

$$L(D_I, M^+) = \prod_{m_j \in M^+} \mu(m_j) (1 - \prod_{d_i \in D_I} (1 - c_{ij})) \prod_{d_i \in D_I} \frac{p_i}{1 - p_i}$$

The plausible function is in keeping with the fact:  $\mu(m_j)$  gradually increases along with degrees of manifestation augments. If the possibility that one manifestation is observed is big, the possibility that one disease can explain the manifestation is comparatively big.

#### 3.2 The selection on fuzzy degree of manifestation

In diagnosing practice, ordinary fishery technicians and fisher folks are accustomed to adopt fuzzy language to describe degrees of manifestation. Each manifestation description is classified five levels, they are “absolutely same, close, relatively close, not same, different”, then corresponding fuzzy value is present to make user easy to describe the observed manifestation.

#### 3.3 The algorithm of diagnosis model

The process of fish-disease diagnosis is to find disease sets that can explain a set of observed manifestations through correlative diagnostic knowledge. It behaves as a repetitious process: preparatory examining- first diagnosis- examining again- diagnose again----make a definite diagnostic solution. In the process, old hypothesis is excluded and new hypothesis is brought up. The circulation repeats again and again until the best solution is brought up. On the basis of analyzing diagnostic process, the corresponding algorithm is present. The basic idea of the algorithm is that the original manifestation set is provided; a trial disease hypothesis can be brought up to explain these manifestations, then current hypothesis can instruct us to find more information for the best solution.

Before introduce the algorithm, we must firstly make some definitions:

$M_1^+$  is a subset of  $M^+$ , and it identifies the set of observed manifestations;

$D_I, cause(M_1^+)$ , it identifies disease sets that can explain observed manifestations;

$S$  is trial solution sets of observed manifestations;

“  $\times$  ” represents Descartes accumulation, and the original state is  $M_1^+ = \Phi, D_1 = \Phi, S = D$ , the terminate condition is  $M_1^+ = M^+$ , in another word, all observed manifestations have been inputted into the model.

The algorithm of diagnosis can be described as follow:

- (1) Input  $m_j$  from  $M^+$ ;
- (2)  $M_1^+ = M_1^+ \cup m_j$ ;
- (3)  $D_1 = D_1 \cup \text{cause}(m_j)$ ;
- (4)  $S'' = S' - S, S' = S, S = S \cap \text{cause}(m_j)$
- (5) IF  $S = \Phi$  THEN;
  - {  $S_1 = S' \times \text{cause}(m_j)$
  - IF  $(\text{effect}(s'') \cup \text{effect}(\{d_i \mid d_i \in \text{cause}(m_j)\}))$  is the smallest covering set of  $M_1^+$
  - THEN  $S_2 = \{d_i\} \times S''$
  - ELSE  $S_2 = \Phi$
  - $S = S_1 \cup S_2$  }
- (6) IF  $M_1^+ = M^+$  THEN end
- ELSE go to (1)

After all probable diagnostic solutions have been found, the plausible function of  $S$  can be calculated and the best solution can be found by thoroughly enumerating strategy and the biggest plausible criterion.

#### 4 Test results

Two sample studies are served for testing the developed fuzzy inference model, and the description of each example and some of the test results are given below.

**Table 1**  $d_i$  and effect ( $d_i$ )

Number	$effect(d_i)$	Number	$effect(d_i)$
$d_1$	$m_{11} m_{46}$	$d_{12}$	$m_{24}$
$d_2$	$m_{10} m_{30} m_{33}$	$d_{13}$	$m_{02} m_{11} m_{18}$
$d_3$	$m_{14} m_{15} m_{43}$	$d_{14}$	$m_{04} m_{15} m_{19} m_{25} m_{32} m_{34}$
$d_4$	$m_{45}$	$d_{15}$	$m_{07} m_{23} m_{29}$
$d_5$	$m_{12}$	$d_{16}$	$m_{09} m_{27}$
$d_6$	$m_{05} m_{33}$	$d_{17}$	$m_{15} m_{16} m_{17} m_{41} m_{42}$
$d_7$	$m_{37} m_{38}$	$d_{18}$	$m_{04} m_{18} m_{19} m_{28} m_{31} m_{41}$
$d_8$	$m_{13}$	$d_{19}$	$m_{41}$
$d_9$	$m_{10} m_{33}$	$d_{20}$	$m_{04} m_{15} m_{16} m_{21} m_{28} m_{30} m_{41} m_{44}$
$d_{10}$	$m_{01} m_{03} m_{20} m_{35}$	$d_{21}$	$m_{08}$
$d_{11}$	$m_{11} m_{15} m_{18} m_{19}$ $m_{26} m_{28} m_{30}$	$d_{22}$	$m_{04} m_{19} m_{22} m_{25} m_{36} m_{39} m_{40}$

Two example is from Jingwu aquatic product breeding factory in Tianjin. After four factors such as  $k_i$  = grass carp,  $a_i$  = fish-seed,  $s_i$  = summer and  $pl_i$  = Jinghai are confirmed, disease sets that can explain these factors can be calculated, it is  $D_i = 22$ , in another word, the probability of each fish-disease is  $p_i = 1/22$ . Which contains 22 disorders and 44 manifestations. The cause-and-effect relationship between disease and manifestations is shown in matrix. More than 20 test cases have been carried out for this example. Only two test cases are illustrated below for saving space.

**4.1 Test case 1**

If  $M^+ = (m_{15}, m_{16}, m_{42})$ , degrees of manifestation separately is:  $\mu(m_{15}) = 0.8$ ,  $\mu(m_{16}) = 0.5$ ,  $\mu(m_{42}) = 0.3$ , the transformation of disease set in the process of diagnosis is described as Table 2.

**Table 2** the diagnosis result of algorithm for test case 1

Manifestations	Original state	$m_{15}$	$m_{16}$	$m_{42}$
$M_1^+$	$\Phi$	$\{m_{15}\}$	$\{m_{15}, m_{16}\}$	$\{m_{15}, m_{16}, m_{42}\}$
$D_1$	$\Phi$	$\{d_{03}, d_{11}, d_{14}, d_{17}, d_{20}\}$	$\{d_{03}, d_{11}, d_{14}, d_{17}, d_{20}\}$	$\{d_{03}, d_{11}, d_{14}, d_{17}, d_{20}\}$
$S'$	$\Phi$	$D$	$\{d_{03}, d_{11}, d_{14}, d_{17}, d_{20}\}$	$\{d_{17}, d_{20}\}$
$S''$	$\Phi$	$\Phi$	$\{d_1, d_{02}, d_{04}, \dots, d_{19}, d_{21}, d_{22}\}$	$\{d_{03}, d_{11}, d_{14}\}$
$S$	$D$	$\{d_{03}, d_{11}, d_{14}, d_{17}, d_{20}\}$	$\{d_{17}, d_{20}\}$	$\{d_{17}\}$

- 1) When manifestations observed is  $\{m_{15}\}$ , diagnostic model will diagnose that disease set is  $\{d_{03}, d_{11}, d_{14}, d_{17}, d_{20}\}$ , user can be instructed to choose other manifestation to obtain the diagnostic solution according to known disease set.
- 2) When manifestations observed is  $\{m_{15}, m_{16}\}$ , diagnosis system will diagnose that disease set is reduced to be  $\{d_{17}, d_{20}\}$ , user can be instructed to choose other manifestation explained by  $\{d_{17}, d_{20}\}$ .
- 3) In the case of manifestation set observed being  $\{m_{15}, m_{16}, m_{42}\}$ , disease set is  $\{d_{17}\}$ .

Finally, the last answer  $S = \{d_{17}\}$ , in another word, when  $M^+ = (m_{15}, m_{16}, m_{42})$ , the result of diagnosis is enteritis. It is obvious that the process of fish-disease diagnosis is a process of seeking more perfection answer 'step by step' based on "hypothesis -test" circulation. In the process of diagnosis, part known manifestation is imported, corresponding disease set will be obtained through diagnosis system. The method can instruct user to collect more manifestation expressed by disease set, the scope of disease will gradually reduced to perfect diagnosis result.

## 5 Conclusion

- (1) Because of three kinds of uncertainty, such as randomness, fuzzy and imperfection, being in existence in the process of fish-disease diagnosis, the degree of manifestations is sufficiently taken into account, a new method for fish-disease diagnostic problem solving based on parsimonious covering theory and fuzzy inference model is constructed in this paper. According to the sequence of diagnosis process, a 'step by step' seeking answer algorithm based on "hypothesis -test" circulation is proposed.
- (2) The diagnosis model proposed in this paper is an efficient attempt to solve the fish-disease diagnosis problem, and realizes the integration probability inference

with fuzzy set theory. Fish-disease can be diagnosed more reliably and practically by means of this model.

(3) The advantage of parsimonious covering model based on fuzzy set theory is analyzed through a great quantity cases. It not only can prompt user to collect more diagnostic information, but also can provide more perfect diagnostic outcome accord with the specialty of fish-disease.

### **Acknowledgement**

The research was funded by the Huo Yingdong foundation in China (project number: 94032). We would like to thank many domain experts from the Beijing Aquaculture Science Institute, Aquaculture Department of Tianjin Agricultural College, Aquaculture Bureau of Shandong province, for their co-operation and support. Our special thanks should also go to Prof. Kezhi Xing, Mr Yongjun Guo at Tianjin Agricultural College for his valuable suggestions and comments on the system.

### **Reference**

1. Chinese academy of aquaculture science, strategies for the aquaculture development in the 21 century of China, Beijing, 2001
2. Wen Ji-wen. A Knowledge-based Fish Diseases Diagnosis Reasoning and Expert System: [D] . Beijing, china agriculture university, 2003
3. Glover F. Tabu search: part 1 ORSA Journal on Computing. 1989,01:190-206
4. Glover F. Tabu search: part 2 ORSA Journal on Computing. 1990,02:4-32
5. Glover F. Future paths foy integer programming and links to artificial in intelligence. Computer and Operations Research.1986.13: 533-549
6. WEN Ji-wen, FU Ze-tian, LI Dao-liang. The Process Imitation and Construction of Reasoning Model of Fish Disease Diagnosis. Aquaculture (In Chinese), 2003(2).

---

# Effective Prover for Minimal Inconsistency Logic

Adolfo Gustavo Serra Seca Neto and Marcelo Finger

Computer Science Department  
Institute of Mathematics and Statistics  
University of São Paulo  
[adolfo,mfinger]@ime.usp.br

**Summary.** In this paper we present an effective prover for **mbC**, a minimal inconsistency logic. The **mbC** logic is a paraconsistent logic of the family of logics of formal inconsistency. Paraconsistent logics have several philosophical motivations as well as many applications in Artificial Intelligence such as in belief revision, inconsistent knowledge reasoning, and logic programming. We have implemented the KEMS prover for **mbC**, a theorem prover based on the **KE** tableau method for **mbC**. We show here that the proof system on which this prover is based is sound, complete and analytic. To evaluate the KEMS prover for **mbC**, we devised four families of **mbC**-valid formulas and we present here the first benchmark results using these families.

## 1 Introduction

In this paper we present new theoretical and practical results concerning paraconsistent logics. On the theoretical side, we have devised a **KE** tableau method for **mbC**, a minimal inconsistency logic, and proved that this proof system is correct, complete and analytic. And on the practical side, we have implemented a theorem prover based on the **mbC KE** proof system and proposed a set of benchmarks for evaluating **mbC** provers.

Paraconsistent logics are tools for reasoning under conditions which do not presuppose consistency [3]. These logics have several philosophical motivations as well as many applications in Artificial Intelligence such as in belief revision [12], inconsistent knowledge reasoning [8], and logic programming [1].

The relevance of reasoning in the presence of inconsistent information can be seen in the following example<sup>1</sup>. Suppose we are working with classical logic and we have a theory (which is a set of formulas)  $\Gamma$  such that  $\Gamma \vdash A$  (i.e. from  $\Gamma$  we can deduce  $A$ ) and also  $\Gamma \vdash \neg A$ . That is, this theory allows us to reach two contradictory conclusions. Suppose also that  $\Gamma \vdash B$ . In classical logic, from  $\Gamma \vdash A$  and  $\Gamma \vdash \neg A$  we can derive  $\Gamma \vdash C$  for any formula  $C$ . In particular,  $\Gamma \vdash \neg B$ .

---

<sup>1</sup> We assume familiarity with the syntax and semantics of propositional classical logic.

In classical logic, a contradictory theory is also trivial, therefore useless. Paraconsistent logics separate these concepts: a contradictory theory needs not to be trivial. Therefore, in a paraconsistent logic such as **mbC**, one can have  $\Gamma \vdash_{\mathbf{mbC}} A$ ,  $\Gamma \vdash_{\mathbf{mbC}} \neg A$  and  $\Gamma \vdash_{\mathbf{mbC}} B$  without necessarily having  $\Gamma \vdash_{\mathbf{mbC}} \neg B$ . Therefore, in paraconsistent logics one can have an inconsistent theory and still draw interesting conclusions from it.

There have been some implementations of paraconsistent formalisms [1, 4], but we do not know of any implementation of a special class of paraconsistency logics: logics of formal inconsistency (**LFIs**) [3]. This class internalizes the notions of consistency and inconsistency at the object-language level. We have extended the KEMS prover [11], originally developed for classical propositional logic, to deal with **LFIs**. The first version of this extension implements a tableau prover for **mbC**, one of the simplest representatives of this class of logics. The KEMS prover for **mbC** is implemented in Java and AspectJ. Java is a well established object-oriented programming language and AspectJ is the major representative of a new programming paradigm: aspect-oriented programming. Its source code available for download in [10].

The KEMS prover is a **KE**-based Multi-Strategy theorem prover. The **KE** system, a tableau method developed by Marco Mondadori and Marcello D'Agostino [7], was presented as an improvement, in the computational efficiency sense, over the Analytic Tableau method [13]. A tableau system for **mbC** had already been presented in [3], but this system is more similar to analytic tableaux than to **KE**: it has five branching rules, which can lead to an inefficient implementation. And although this system is sound and complete it is not analytic. Therefore, to implement the KEMS prover for **mbC** we devised an **mbC KE** system and obtained a sound, complete and analytic tableau proof system with only one branching rule.

To evaluate our prover correctness and performance, we needed some families of **mbC** problems. As we do not know any family of valid formulas elaborated specially for **mbC** or any paraconsistent logic, we devised four families of **mbC**-valid problems for evaluating **mbC** provers. These families are not classically valid, since all of them use the non-classical consistency connective. With these families we obtained the first benchmark results for the KEMS **mbC** implementation.

## 1.1 Outline

In section 2 we present the **mbC** logic. The **mbC KE** system is exhibited in section 3. There we also prove its analyticity, soundness and completeness. In section 4 we show the problem families we devised to evaluate **mbC** provers and in section 5 we present the results obtained with the KEMS prover for **mbC** using these families as benchmarks. Finally, in section 6 we draw some conclusions and point to future work.

## 2 The **mbC** Logic

The **mbC** logic is a member of the family of logics of formal inconsistency [3]. Logics of formal inconsistency are a class of paraconsistent logics that internalize the notions of consistency and inconsistency at the object-language level. Paraconsistent logics



are tools for reasoning under conditions which do not presuppose consistency [3]. Formal characterizations of paraconsistent logics and logics of formal inconsistency can be found, respectively, in [9] and [3].

The logic **mbC** is the weakest<sup>2</sup> **LFI** based on classical logic presented in [3]. It uses the same set of connectives as propositional classical logic (the binary connectives  $\wedge, \vee, \rightarrow$ , and the unary connective  $\neg$ ), plus a new one: the unary *consistency* ( $\circ$ ) connective. The intended reading of  $\circ A$  is ‘ $A$  is consistent’, that is, if  $\circ A$  is true,  $A$  and  $\neg A$  are not both true. In **mbC**,  $\circ A$  is logically independent from  $\neg(A \wedge \neg A)$ , that is,  $\circ$  is a primitive unary connective, not an abbreviation depending on conjunction and negation, as it happens in da Costa’s  $C_n$  hierarchy of paraconsistent logics [5]. Its axiomatization is shown below:

**Axiom schemas**

$$\begin{aligned}
 & A \rightarrow (B \rightarrow A) \\
 & (A \rightarrow B) \rightarrow ((A \rightarrow (B \rightarrow C)) \rightarrow (A \rightarrow C)) \\
 & A \rightarrow (B \rightarrow (A \wedge B)) \\
 & (A \wedge B) \rightarrow A \\
 & (A \wedge B) \rightarrow B \\
 & A \rightarrow (A \vee B) \\
 & B \rightarrow (A \vee B) \\
 & (A \rightarrow C) \rightarrow ((B \rightarrow C) \rightarrow ((A \vee B) \rightarrow C)) \\
 & A \vee (A \rightarrow B) \\
 & A \vee \neg A \\
 & \circ A \rightarrow (A \rightarrow (\neg A \rightarrow B))
 \end{aligned}$$

**Inference rule**

$$\text{(Modus Ponens)} \quad \frac{A, A \rightarrow B}{B}$$

Now we present the formal definition of satisfiable and valid formulas in **mbC** [3]. Let  $\mathbf{2} \stackrel{\text{def}}{=} \{0, 1\}$  be the set of truth-values, where 1 denotes the ‘true’ value and 0 denotes the ‘false’ value. An **mbC-valuation** is any function  $v : \text{For} \rightarrow \mathbf{2}$  subject to the following clauses:

$$\begin{aligned}
 & v(A \wedge B) = 1 \text{ iff } v(A) = 1 \text{ and } v(B) = 1; \\
 & v(A \vee B) = 1 \text{ iff } v(A) = 1 \text{ or } v(B) = 1; \\
 & v(A \rightarrow B) = 1 \text{ iff } v(A) = 0 \text{ or } v(B) = 1; \\
 & v(\neg A) = 0 \text{ implies } v(A) = 1; \\
 & v(\circ A) = 1 \text{ implies } v(A) = 0 \text{ or } v(\neg A) = 0.
 \end{aligned}$$

A formula  $X$  is said to be *satisfiable* if truth-values can be assigned to its propositional variables in a way that makes the formula true, i.e. if there is at least one valuation such that  $v(X) = 1$ . A formula is a *valid* if all possible valuations make the formula true. For instance, the formula  $\neg(A \wedge \neg A \wedge \circ A)$  is a valid in **mbC**, while  $\neg(A \wedge \neg A)$  is satisfiable.

---

<sup>2</sup> It is the weakest because all other **LFI**s presented in [3] prove more theorems.

### 3 A KE System for mbC

The *Analytic Tableau* method is probably the most studied tableau method. It was presented in [13] as “an extremely elegant and efficient proof procedure for propositional logic”. The **KE** System, a tableau method developed by Marco Mondadori and Marcello D’Agostino [7], was presented as an improvement, in the computational efficiency sense, over the Analytic Tableau method. It is a refutation system that, though close to the Analytic Tableau method, is not affected by the anomalies of cut-free systems [6].

In [3], a sound and complete tableau proof system for **mbC** is presented. It was obtained by using a method introduced in [2]. This method is a generic method that automatically generates a set of tableau rules for certain logics. For **mbC**, the rules obtained for its binary connectives are the same as that from classical analytic tableaux. The system also has a branching rule (called  $R_b$ ) similar to **KE** PB rule, as well as rules for negation ( $\neg$ ) and consistency ( $\circ$ ). In total, this tableau system has 5 branching rules.

$$\begin{array}{c}
 \frac{\text{T } A \vee B}{\text{F } A} \quad (\text{T } \vee 1) \qquad \frac{\text{T } A \vee B}{\text{T } A} \quad (\text{T } \vee 2) \qquad \frac{\text{F } A \vee B}{\text{F } A} \quad (\text{F } \vee) \\
 \frac{\text{F } A \wedge B}{\text{T } A} \quad (\text{F } \wedge 1) \qquad \frac{\text{F } A \wedge B}{\text{F } A} \quad (\text{F } \wedge 2) \qquad \frac{\text{T } A \wedge B}{\text{T } A} \quad (\text{T } \wedge) \\
 \frac{\text{T } A \rightarrow B}{\text{T } A} \quad (\text{T } \rightarrow 1) \qquad \frac{\text{T } A \rightarrow B}{\text{F } B} \quad (\text{T } \rightarrow 2) \qquad \frac{\text{F } A \rightarrow B}{\text{T } A} \quad (\text{F } \rightarrow) \\
 \frac{\text{T } \neg A}{\text{T } \circ A} \quad (\text{T } \neg) \qquad \frac{\text{F } \neg A}{\text{T } A} \quad (\text{F } \neg) \\
 \begin{array}{c}
 \diagup \quad \diagdown \\
 \text{T } A \quad \text{F } A
 \end{array} \quad (\text{PB})
 \end{array}$$

**Fig. 1.** mbC KE tableau expansion rules

As explained in [6], branching rules lead to inefficiency. To obtain a more efficient proof system, we devised an original **mbC KE** system using signed formulas (see Figure 1). A *signed formula* is an expression  $SX$  where  $S \in \{\text{T}, \text{F}\}$  is called the *sign* and  $X$  is a propositional *formula*. The symbols T and F, respectively representing

the truth-values true and false, can be used as signs. The *conjugate* of a signed formula  $\mathsf{T}A$  (or  $\mathsf{F}A$ ) is  $\mathsf{F}A$  (or  $\mathsf{T}A$ ). The **mbC** ( $\mathsf{T}\neg$ ) rule is a **LFI** version of classical propositional logic ( $\mathsf{T}\neg$ ) [6]. It states clearly that in **mbC** we need  $\mathsf{T}\neg A$  and  $\mathsf{T}\circ A$  to obtain  $\mathsf{F}A$ . In classical logic, we can obtain  $\mathsf{F}A$  directly from  $\mathsf{T}\neg A$ .

### 3.1 Analyticity, Correctness and Completeness Proof for the **mbC KE** system

An **mbC KE** proof enjoys the *subformula property* if every signed formula in the proof tree is a subformula of some formula in the list of signed formulas to be proved. Let us call *analytic* the applications of PB which preserve the subformula property, and the *analytic restriction of mbC KE* the system obtained by restricting PB to analytic applications. Given a rule  $R$  of an expansion system  $\mathbf{S}$ , we say that an application of  $R$  to a branch  $\theta$  is *analytic* when it has the *subformula property*, i.e. if all the new signed formulas appended to the end of  $\theta$  are subformulas of signed formulas occurring in  $\theta$ . According to [6], a rule  $R$  is *analytic* if every application of it is analytic. It is easy to notice that all **mbC KE** rules except (PB) are analytic.

We prove here that the **mbC KE** system is analytic, sound and complete (some proofs were omitted due to lack of space). It is easy to show a procedure that transforms any proof in the original tableau system for **mbC** ([3]) in an **mbC KE** proof, thus proving that **mbC KE** system is also sound and complete. We will not do this here. Instead, we will demonstrate that even the analytic restriction of **mbC KE** is sound and complete. That is, when performing a proof we can restrict ourselves to analytic applications of PB, applications which do not violate the subformula property, without affecting completeness.

The proof will be as follows. First we will redefine the notion of downward saturatedness for **mbC**. Then we will prove that every downward saturated set is satisfiable. The **mbC KE** proof search procedure for a set of signed formulas  $S$  either provides one or more downward saturated sets that give a valuation satisfying  $S$  or finishes with no downward saturated set. Therefore, if an **mbC KE** tableau for a set of formulas  $S$  closes, then there is no downward saturated set that includes it, so  $S$  is unsatisfiable. However, if the tableau is open and completed, then any of its open branches can be represented as a downward saturated set and be used to provide a valuation that satisfies  $S$ . By construction, downward saturated sets for open branches are analytic, i.e. include only subformulas of  $S$ . Therefore, the **mbC KE** system is analytic. As a corollary, it is also sound and complete.

**Definition 1.** A set of signed formulas  $DS$  is *downward saturated* if

1. whenever a signed formula is in  $DS$ , its conjugate is *not* in  $DS$ ;
2. when all premises of any **mbC KE** rule (except PB) are in  $DS$ , its conclusions are also in  $DS$ ;
3. when the major premise of a **mbC KE** rule is in  $DS$ , either its auxiliary premise or its conjugate is in  $DS$ .

For **mbC KE**, item (3) above is valid for every rule except ( $\mathsf{T}\neg$ ). In this case, if  $\mathsf{T}\neg X$  is in  $DS$ , either  $\mathsf{T}\circ X$  or  $\mathsf{F}\circ X$  is in  $DS$  only if  $\circ X$  is a subformula of some other formula in  $DS$ .

We extend valuations to signed formulas in an obvious way:  $v(\mathbf{T}A) = v(A)$  and  $v(\mathbf{F}A) = 1 - v(A)$ . A set of signed formulas  $L$  is satisfiable if it is not empty and there is a valuation such that for every formula  $SX \in L$ ,  $v(SX) = 1$ . Otherwise, it is unsatisfiable.

**Lemma 1.** (Hintikka's Lemma) Every downward saturated set is satisfiable.

*Proof.* For any downward saturated set  $DS$ , we can easily construct a valuation  $v$  such that for every signed formula  $SX$  in the set,  $v(SX) = 1$ . How can we guarantee this is in fact a valuation? First, we know that there is no pair  $\mathbf{T}X$  and  $\mathbf{F}X$  in  $DS$ . Second, **mbC KE** rules preserve valuations. That is, if  $v(SX_i) = 1$  for every premise  $SX_i$ , then  $v(SC_j) = 1$  for all conclusions  $C_j$ . And if  $v(SX_1) = 1$  and  $v(SX_2) = 0$ , where  $X_1$  and  $X_2$  are, respectively, major and minor premises of an **mbC KE** rule, then  $v(S'X_2) = 1$ , where  $S'X_2$  is the conjugate of  $SX_2$ . For instance, suppose  $\mathbf{T}A \wedge B \in DS$ , then  $v(\mathbf{T}A \wedge B) = 1$ . In accord with the definition of downward saturated sets,  $\{\mathbf{T}A, \mathbf{T}B\} \subseteq DS$ . And by the definition of valuation,  $v(\mathbf{T}A \wedge B) = 1$  implies  $v(\mathbf{T}A) = v(\mathbf{T}B) = 1$ .  $\square$

**Theorem 1.**  $DS'$  is a set of signed formulas.  $DS'$  is satisfiable *if and only if* there exists a downward saturated set  $DS''$  such that  $DS' \subseteq DS''$ .

**Corollary 1.**  $DS'$  is a unsatisfiable set of formulas if and only if there is no downward saturated set  $DS''$  such that  $DS' \subseteq DS''$ .

**Theorem 2.** The **mbC KE** system is analytic.

*Proof.* The **mbC KE** proof search procedure for a set of signed formulas  $S$  either provides one or more downward saturated sets that give a valuation satisfying  $S$  or finishes with no downward saturated set. If an **mbC KE** tableau for a set of formulas  $S$  closes, then there is no downward saturated set that includes it, so  $S$  is unsatisfiable. If the tableau is open and completed, then any of its open branches can be represented as a downward saturated set and be used to provide a valuation that satisfies  $S$ . By construction, downward saturated sets for open branches are analytic, i.e. include only subformulas of  $S$ . Therefore, the **mbC KE** system is analytic.  $\square$

**Corollary 2.** The **mbC KE** system is sound and complete.

## 4 Problem Families

We present below the problem families we devised to evaluate **mbC** theorem provers. We had two objectives in mind. First, to obtain families of **mbC**-valid problems whose **mbC KE** proofs were as complex as possible. And second, to devise problems which required the use of many, if not all, **mbC KE** rules. These families are not classically valid, since all of them have formulas with the non-classical consistency connective.

### 4.1 First family

Here we present the first family ( $\Phi^1$ ) of valid sequents for **mbC**. In this family all **mbC** connectives are used. It is easy to obtain polynomial **mbC KE** proofs for this family of problems. The sequent to be proved for the  $n$ -th instance of this family ( $\Phi_n^1$ ) is:

$$\bigwedge_{i=1}^n (\neg A_i), \bigwedge_{i=1}^n ((\circ A_i) \rightarrow A_i), [\bigvee_{i=1}^n (\circ A_i)] \vee (\neg A_n \rightarrow C) \vdash C \quad (1)$$

The explanation for this family is as follows. Suppose we are working with a database that allows inconsistent information representation.  $A_i$  means that someone expressed an opinion  $A$  about an individual  $i$  and  $\neg A_i$  means that someone expressed an opinion  $\neg A$  about this same individual. For instance, if  $A$  means that a person is nice,  $\neg A_3$  means that at least one person finds 3 is not nice, and  $A_4$  means that at least one person finds 4 nice. Then  $\circ A_i$  means that either all people think  $i$  is nice, or all people think  $i$  is not nice, or there is no opinion  $A$  recorded about  $i$ .  $\circ A_i \rightarrow A_i$  means that if all opinions about a person are the same, then that opinion is  $A$ .

For a subset of individuals numbered from 1 to  $n$ , we have  $\neg A_i$  and  $\circ A_i \rightarrow A_i$  for all of them. From the fact that either  $\neg A_n \rightarrow C$  or for one of them we have  $\circ A_i$ , we can conclude  $C$ .

### 4.2 Second Family

The second family of problems for **mbC** ( $\Phi^2$ ) is a variation over the first family whose proofs are exponential in size. The sequent to be proved for the  $n$ -th instance of this family ( $\Phi_n^2$ ) is:

$$\bigwedge_{i=1}^n (\neg A_i), [\bigwedge_{i=1}^n [(\circ A_i) \rightarrow ((\bigvee_{j=i+1}^n \circ A_j) \vee ((\neg A_n) \rightarrow C))]],$$

$$[\bigvee_{i=1}^m (\circ A_i)] \vee (\neg A_n \rightarrow C) \vdash C$$

This family is a modification of the first family where instead of a conjunction of  $\circ A_i \rightarrow A_i$ , we have a conjunction of  $\circ A_i \rightarrow ((\bigvee_{j=i+1}^n \circ A_j) \vee ((\neg A_n) \rightarrow C))$  meaning that for every person numbered 1 to  $n$ , if all opinions about a person are the same, then either all opinions about some other person with a higher index are the same or  $(\neg A_n) \rightarrow C$  is true.

### 4.3 Third Family

With the third family of problems we intended to develop a family whose instances required the application of all **mbC KE** rules. To devise the third family ( $\Phi^3$ ), we have made some changes to the second family trying to make it more difficult to prove. The  $n$ -th instance of this family ( $\Phi_n^3$ ) is the following sequent:

$$\begin{aligned}
& U_l \wedge U_r, \\
& \bigwedge_{i=1}^n (\neg A_i), \\
& \bigwedge_{i=1}^n [(\circ A_i) \rightarrow (((\neg A_n) \wedge U_i) \rightarrow C) \vee \bigvee_{j=i+1}^n \circ A_j], \\
& (\bigvee_{i=1}^n \circ A_i) \vee ((U_r \wedge (\neg A_n)) \rightarrow C) \\
\vdash & C' \rightarrow (C'' \vee C)
\end{aligned}$$

#### 4.4 Fourth Family

This is the only of these families where negation appears only in the conclusion. The  $n$ -th instance of this family ( $\Phi_n^4$ ) is the following sequent:

$$\bigwedge_{i=1}^n (A_i), \bigwedge_{j=1}^n ((A_j \vee B_j) \rightarrow (\circ A_{j+1})), [\bigwedge_{k=2}^n (\circ A_k)] \rightarrow A_{n+1} \vdash \neg \neg A_{n+1}$$

Note: if  $n \leq 2$ ,  $[\bigwedge_{i=2}^n (\circ A_i)]$  in  $[\bigwedge_{i=2}^n (\circ A_i)] \rightarrow A_{n+1}$  is replaced by the  $\top$  formula.

This family formulas can be explained as follows. We have two formulas to represent two types of opinion:  $A$  and  $B$ . First we assume  $A_i$  for every  $i$  from 1 to  $n$ . Then we suppose for all  $j$  from 1 to  $n$  that  $(A_j \vee B_j)$  implies  $\circ A_{j+1}$ . And finally we assume that for every  $k$  from 2 to  $n$  the conjunction of  $A_k$ 's implies  $A_{n+1}$ . It is easy to see that from these assumptions we can deduce  $A_{n+1}$ . So we can also deduce its double negation:  $\neg \neg A_{n+1}$ .

## 5 Evaluation

Theorem provers are usually compared by using benchmarks. We have extended KEMS prover [11] to prove **mbC** theorems and evaluated it using as benchmarks the problem families presented in section 4. In Table 1 we show some of the results obtained. The tests were run on a personal computer with an Athlon 1100Mhz processor, 384Mb of memory, running a Linux operating system with a 2.26 kernel.

Problem	Time spent (s)	Problem size	Proof size	Tree height
$\Phi_4^1$	0.06	47	197	4
$\Phi_7^1$	0.046	80	491	7
$\Phi_{10}^1$	0.08	113	911	10
$\Phi_4^2$	0.071	77	570	7
$\Phi_7^2$	1.54	164	7350	13
$\Phi_{10}^2$	21.964	278	116037	19
$\Phi_4^3$	0.058	94	706	6
$\Phi_7^3$	1.097	187	5432	9
$\Phi_{10}^3$	17.595	307	52540	12
$\Phi_4^4$	0.007	47	181	3
$\Phi_7^4$	0.013	83	433	3
$\Phi_{10}^4$	0.023	119	793	3

**Table 1.** KEMS results for **mbC**

From these results it is clear that the second and third families are much more difficult to prove than the other two. And interestingly enough it was easier to prove the third than the second family.

## 6 Conclusion

We have presented an effective prover for **mbC**: a minimal inconsistency logic. The **mbC KE** system it implements was proven to be sound, complete and analytic. Besides that, it has only one branching rule. We devised some families of valid problems to evaluate our prover correctness and performance. These families can be used to evaluate any **mbC** theorem prover. The KEMS prover for **mbC** obtained the first benchmark results for these problem families.

In the future we intend to design different KEMS strategies for **mbC**. For instance, we want to implement a strategy that uses some derived rules not presented here. After that, we want to extend the KEMS prover to deal with  $C_1$ , the first logic in da Costa's  $C_n$  hierarchy of paraconsistent logics [5].

This paper has been partially sponsored by FAPESP Thematic Project Grant ConsRel 2004/14107-2.

## References

1. H. A. Blair and V. S. Subrahmanian. Paraconsistent logic programming. *Theor. Comput. Sci.*, 68(2):135–154, 1989.
2. Carlos Caleiro, Walter Carnielli, Marcelo E. Coniglio, and Joao Marcos. Two's company: "The humbug of many logical values". In *Logica Universalis*, pages 169–189. Birkhäuser Verlag, Basel, Switzerland, 2005.
3. Walter Carnielli, Marcelo E. Coniglio, and Joao Marcos. Logics of Formal Inconsistency. In *Handbook of Philosophical Logic*, volume 12. Kluwer Academic Publishers, 2005.
4. Newton C. A. da Costa, Lawrence J. Henschen, James J. Lu, and V. S. Subrahmanian. Automatic theorem proving in paraconsistent logics: theory and implementation. In *CADE-10: Proceedings of the tenth international conference on Automated deduction*, pages 72–86, New York, NY, USA, 1990. Springer-Verlag New York, Inc.
5. Newton C. A. da Costa, Décio Krause, and Otávio Bueno. Paraconsistent logics and paraconsistency: Technical and philosophical developments. *CLE e-prints (Section Logic)*, 4(3), 2004.
6. Marcello D'Agostino. Tableau methods for classical propositional logic. In Marcello D'Agostino et al., editor, *Handbook of Tableau Methods*, chapter 1, pages 45–123. Kluwer Academic Press, 1999.
7. Marcello D'Agostino and Marco Mondadori. The taming of the cut: Classical refutations with analytic cut. *Journal of Logic and Computation*, pages 285–319, 1994.
8. Sandra de Amo, Walter Alexandre Carnielli, and Joao Marcos. A logical framework for integrating inconsistent information in multiple databases. In *FoIKS*

- '02: *Proceedings of the Second International Symposium on Foundations of Information and Knowledge Systems*, pages 67–84, London, UK, 2002. Springer-Verlag.
9. Itala M. Loffredo D'Ottaviano and Milton Augustinis de Castro. Analytical Tableaux for da Costa's Hierarchy of Paraconsistent Logics  $C_n, 1 \leq n \leq \omega$ . *Journal of Applied Non-Classical Logics*, 15(1):69–103, 2005.
  10. Adolfo Gustavo Serra Seca Neto. KEMS - A KE-based Multi-Strategy Theorem Prover, 2006. Retrieved February 01, 2006, from <http://gsd.ime.usp.br/~adolfo/projetos/KEMS.zip>.
  11. Adolfo Gustavo Serra Seca Neto and Marcelo Finger. Implementing a multi-strategy theorem prover. In Ana Cristina Bicharra Garcia and Fernando Santos Osório, editors, *Proceedings of the V ENIA, São Leopoldo-RS, Brazil, July 22-29 2005*, 2005.
  12. G. Priest. Paraconsistent Belief Revision. *Theoria*, 67:214–228, 2001.
  13. Raymond M. Smullyan. *First-Order Logic*. Springer-Verlag, 1968.



# Identification of Important News for Exchange Rate Modeling

Debbie Zhang, Simeon J. Simoff and John Debenham

<sup>1</sup> Faculty of Information Technology,  
University of Technology, Sydney, Australia  
<sup>2</sup> {debbiez, simeon, debenham}@it.uts.edu.au

**Abstract.** Associating the pattern in text data with the pattern with time series data is a novel task. In this paper, an approach that utilizes the features of the time series data and domain knowledge is proposed and used to identify the patterns for exchange rate modeling. A set of rules to identify the patterns are firstly specified using domain knowledge. The text data are then associated with the exchange rate data and pre-classified according to the trend of the time series. The rules are further refined by the characteristics of the pre-classified data. Classification solely based on time series data requires precise and timely data, which are difficult to obtain from financial market reports. On the other hand, domain knowledge is often very expensive to be acquired and often has a modest inter-rater reliability. The proposed method combines both methods, leading to a “grey box” approach that can handle the data with some time delay and overcome these drawbacks.

## 1 Introduction

Until recently, most of exchange rate models are empirical models based on macro economic data. While these models performed reasonably well in predicting long term trends, they have little success in predicting short to middle term movements. Recent research has discovered that irrationality of the market participants, bubbles, and herd behavior are the main driven forces of the short term instability. Many market participants make their decisions before the economic data are announced, trying to beat the market. Their decisions are based on their observations and prediction of market trends. They adjust their investment strategy again after the data announcement according to the differences between their expectation and the actual data. Unarguably, news play an important role in creating such market dynamics. Recent research has confirmed that news has statistically significant effects on daily exchange rate movement. Ehrmann and Fratzscher [1] have evaluated the overall impact of macro news by analyzing the daily exchange rate responses using regression models with news variables. Three key results were found. Firstly, the news about fundamentals can explain relatively well the direction, but only a much smaller extent to the magnitude of exchange rate development. Secondly, news about US economy has a larger impact on exchange rates than news about the euro area. Thirdly,

higher degree of market uncertainty will lead to more significant effects of news releases on exchange rate movements. Prast and De Vor [2] have also studied the reaction of investors in foreign exchange markets to news information about the euro area and the United States on days of large changes in the euro-dollar exchange rates. Unlike the traditional models, daily news about economic variables as well as relevant political events in the United States and Europe were used in the regression model, which is:

$$E_t = \alpha + \sum_{i=1}^8 \beta_i D_i + \varepsilon \quad (1)$$

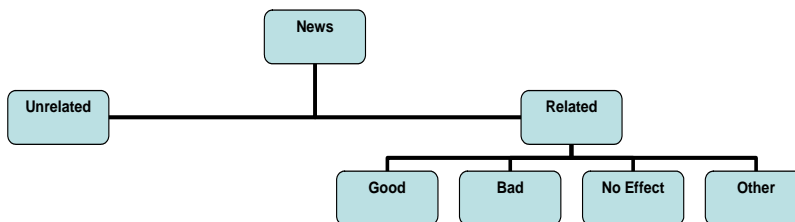
where  $E_t$  is the percentage daily change in the euro-dollar exchange rate;  $D_{1-8}$  represent the following variables: 1 - real economy, euro area; 2 - inflation, euro area; 3 - change in official interest rate, ECB; 4 - statements/political events, euro area; 5 - real economy, United States; 6 - inflation, United States; 7 - change in official interest rate, United States; 8 - statements/political events, United States. It has been found that there is a strong correlation between exchange rate daily movement and the market participants' responses to the daily economy news and political events.

Motivated by their findings, Zhang, Simoff and Debenham [3] proposed a method to automate the exchange rate prediction process using text mining techniques. Nowadays, news retrieval using a computer program is effective and efficient as news data are widely available from the internet. However, to develop a program to understand the news articles and make the correct decision according to the news content is still an extreme difficult task. In this approach, the problem was broken down to multiple text classification problems and a regression modeling problem using numerical data. Relevant news are firstly identified and then being classified into the "good", "bad" or "no effect" news categories, which would have dramatically different impact on the market behaviors. "good" or "bad news" in this application are defined as the news that indicate the US dollar is destined to appreciate or depreciate versus Euro respectively. A regression prediction model is built based on these effects. Since the classification approaches are used, the training data are required to be pre-classified before the training process. The pre-classified training data contains the underlying knowledge of how data is being classified. Therefore, classification of training data is critical in any classification projects. However, many papers ignore this issue and provide no details of the methodology used in this step. This paper proposed a semi-automatic approach to assist the manual classification process.

## 2 News classification using a rule-based system

The training news data are required to be classified into one of the categories as shown in Figure 1 to incorporate their effects into the exchange rate model. This step can be done automatically or manually.

The automatic approach can be done by linking the news events with the changes in the exchange rate data if precise data are available. The exchange



**Fig. 1.** News categories

rate data are recorded regularly over a time period, often referred to as “time series”. At the same time, news about the events that may cause the data fluctuation are also recorded. The correlation between these two observations plays an important role in developing the prediction models. Methodologies like: Apriori-like Discovery of Association Rules, Template-Based Mining for Sequences, and Classification of Temporal Data can be used to discover their correlations. Recently, Ting et al [4] proposed an approach to classify the company announcements into three categories according to changes of the share prices. Lavrenko et al [5] also presented a method for identifying news stories that influence the behavior of financial markets using language models. However, these approaches require precise data which are recorded with high sampling frequencies [6]. Testing the influence of the information requires precise data on both the information available prior to a public announcement and the actual released information. Recording the information flow accurately can be very difficult, particularly for the exchange rate data. Although the times at which scheduled macroeconomic releases are available, identification of the pre and post announcement market behavior is often very difficult. This is due to the decentralized nature of currency exchange market and leaking of inside information to some of the market participants. Delays in publication of the announcements on the internet news sources further contribute to the difficulty. These are the typical issues in the data retrieved from online news web sites including [www.bloomberg.com](http://www.bloomberg.com), which are used in this study. At the time the news is published on the web site, some market operators already have the information and have reacted correspondingly.

More common approaches of data sorting relies on analysts’ knowledge of the domain. News data are manually classified into several categories according to a set of rules that are specified by the domain experts. Rule based systems such as expert systems can be employed to help collecting and maintenance of the rules. However, the knowledge acquisition and rule generation process from domain experts is time consuming, expensive, and sometimes has only a modest inter-rater reliability. The rules derived from the domain experts can only a good starting point. It would be a sensible approach to use the data patterns to improve the rules. Therefore, a “grey box” method using domain knowledge in

conjunction with time series features of exchange rate data to refine the rules is proposed. The domain knowledge to form the initial rules is the white part of the model. However, these rules are not complete and can only be used to describe the data partially. Therefore, the rules need to be fine tuned using the data.

The intention of rule-based classification is to apply a validated, fully disclosed and understandable set of rules to perform the classification. The classification algorithm includes following steps:

- Specify a set of rules that reasonable to experts.
- Associate the text data with the exchange rate trends and classify the news according to the trend.
- Refine the rules if they are contradictive to the classified data in the previous step.
- Reclassify the data according to the new rules.

## 2.1 Obtaining Information for Developing Rules

Without an on-site domain expert, rules can be developed using published models that represent economists' understanding of the system. The rules used in this paper were developed using regression models provided by [1] and [2]. According to the regression result tables in the literature, the independent variables represent the factors that affect the exchange rate. The sign on the coefficient (positive or negative) indicates the direction of the effect. Although the size of the coefficient for each independent variable gives the size of the effect that variable is having on the exchange rate, they are ignored in the rules as they are highly sensitive to the data sets. It is difficult to quantify how the information is assessed by market participants since they have different responses for the same set of data when they have different expectations.

## 2.2 Association the News Data with the Exchange Rate Data

Each news item can also be classified according to its association with a particular trend in the exchange rate data. Considering most of the news web sites have time delay in publishing latest news and their refreshing frequency is about twice a day, the time window is set to 12 hours before and 12 hours after the exchange rate data. Figure 2 illustrates the alignment.

The trends are then aligned with news stories for classification. The results are used to examine the rules derived in Section 2.1. Adya etc [7] developed and automated heuristics to detect six features of the time series based on statistics: outliers, level shifts, change in basic trends, unstable recent trends, unusual last observation and functional form. However, to simplify the analysis process, only the upward trend and downward trend are identified. The exchange rate changes of each day are calculated. The days with large changes are selected, such as t1, t2 and t3 shown in Figure 2. The news articles that associated with these days are also selected and put into either "good" news or "bad" news categories

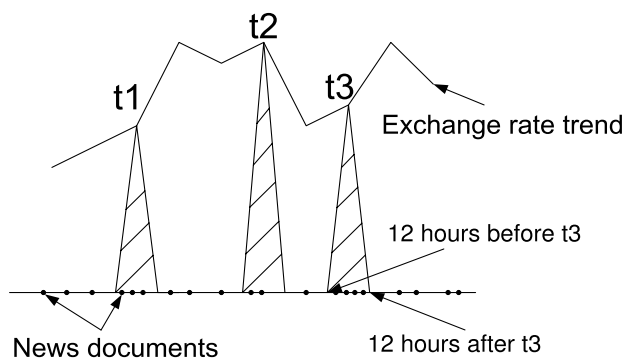


Fig. 2. News Alignment

according to their associated trend. At this stage, both categories contain many unrelated news since many news items of that day may have no impact on the exchange rate.

### 2.3 Refining the Rules and Classification Results

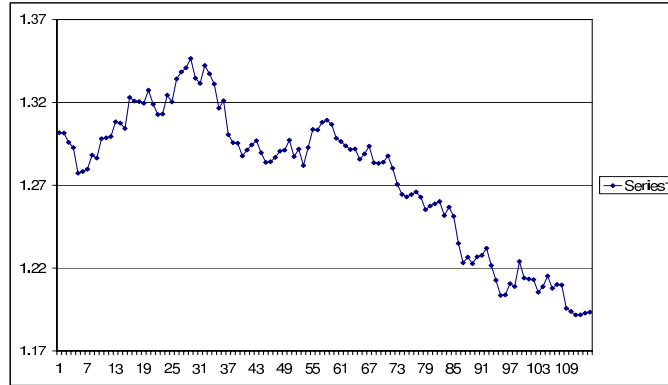
The rules derived in the section 2.1 can be validated by the classification results obtained in the section 2.2. The news in the “good” or “bad” categories are examined and summarized manually. Some simple heuristics are used to refine the rules. For example, if a rule describes an event has positive impact on the exchange rate, but in the data, both “good” and “bad” news contains multiple times of this event, it suggests that this event has no impact in the exchange rate at least within the period the data were collected. The rule should be removed. If multiple news items in the same category indicate that an event has impact in the exchange rate data, a new rule is added to link the event to the exchange rate movement. After the rules are being refined, the whole data set is manually classified according to the rules. In the final classification results, both “good” and “bad” news items may occur in the same day. This is reasonable as the changes of exchange rate are modeled as the the combined effect of multiple news items.

## 3 Experiment results

The experiment was conducted using the news articles and the US dollar versus Euro exchange rate data, both collected from [www.bloomberg.com](http://www.bloomberg.com) between 7/2/2005 to 4/7/2005.

Rules are derived from Ehrmann and Prast’s models (see [1] and [2]), which are summarized in Table 1.

Their models also include “statements/political events”, the area in which the rules are mainly expanded.



**Fig. 3.** News Alignment

**Table 1.** Summary of Ehrmann and Prast’s models of positive and negative factors

	Positive factors	Negative factors
US announcement surprised	Inflation, Retail sales Nonfarm payrolles, Industrial production, N.A.P.M, Trade balance, Advance GDP, Average workweek Consumer confidence	Unemployment rate, Housing starts, CPI, PPI
German announcement surprises	CPI, Unemployment rate	Ifo Business Climate Industry production Manufacturing orders Retail sales Trade balance Business confidence PPI

To align trends with news data, a predefined threshold is used to define the trend. If  $\Delta g_t$ , which is defined in equation , is within the positive and negative range of the predefined threshold, it is considered that there is no changes in the price.

$$\Delta g_t = \frac{p_t - p_{t-1}}{p_{t-1}} \times 100\% \quad (2)$$

where,  $\Delta g_t$  denotes the change of the exchange rate of day t,  $p_t$  and  $p_{t-1}$  denotes the exchange rate of day t and t-1.

In the experiment, the threshold of 0.1% is used. Within the experiment period, there were 40 days of upward trend and 49 days of downward trend. The news were compared on a daily base manually. Data have shown that Greenspan has profound inferences in the exchange rate. Although this is a well known fact, the original rules do not include this due to the models' limitation. Therefore, Greenspan's indication of positive or negative economic trend becomes one of the rule. The other rule worth to mention is the effect of U.S. 10 year treasury. The news of U.S. 10 year treasury appear very often. The first thought was that it could be an indicator of the rate trend. After carefully examining the data, no direct association was found. The final classification results using the refined rules are listed in Table 2.

**Table 2.** Classification results

Total number of News: 2589
Number of "Unrelated" News: 1885
Number of "Related" News: 704
Number of "Good" News: 200
Number of "Bad" News: 113
Number of "No_Effect" News: 230
Number of "Other" News: 161

As the experiment is still in an early stage, having data collected in a longer period would be better to evaluate the rules. More data are being collected using a program that can retrieve news data regardless the news web page format [8]. However, the preliminary results show that this approach provides better classification results than the classification solely based on the domain knowledge.

## 4 Conclusions

Exchange rate prediction using news article requires high quality training data. To achieve this, the news articles that have impact on exchange rate should be correctly associated with one of its trends, namely the upward and downward trends. Classification solely based on time series data requires precise and timely data, which are difficult to obtain for this application. On the other hand, domain

knowledge is often very expensive to be acquired and often has a modest inter-rater reliability. The proposed method combines both methods, leading to a “grey box” approach, which well suits this particular application.

## References

1. Ehrmann, M., Fratzscher, M.: Exchange rates and fundamentals: new evidence from real-time data. *Journal of International Money and Finance* **24** (2005) 317–341
2. Prast, H.M., de Vor, M.P.H.: Investor reactions to news: a cognitive dissonance analysis of the euro-dollar exchange rate. *European Journal of Political Economy* **21** (2005) 115–141
3. Zhang, D., Simoff, S., Debenham, J.: Exchange rate modelling using news articles and economic data. In: 18th Australian Joint Conference on Artificial Intelligence. (2005)
4. Yu, T., Jan, T., Debenham, J., Simoff, S.J.: Incorporate domain knowledge into support vector machine to classify price impacts of unexpected news. In: Australasian Data Mining Conference. (2005)
5. Lavrenko, V., Schmill, M., Lawrie, D., Ogilvie, P., Jensen, D., Allan, J.: Language models for financial news recommendation. In: Ninth International Conference on Information and Knowledge Management (CIKM), Washington (2000) 389–396
6. Hautsch, N., Hess, D.: Bayesian learning in financial markets - testing for the relevance of information precision in price discovery. *Journal of Financial and Quantitative Analysis* (2005)
7. Adya, M., Collopy, F., Armstrong, J.S., Kennedy, M.: Automatic identification of time series features for rule-based forecasting. *International Journal of Forecasting* **17** (2001) 143–157
8. Zhang, D., Simoff, S.: Informing the curious negotiator: Automatic news extraction from the internet. In: Australasian Data Mining Conference, Cairns, Australia (2004)