

Hemalata Iyer, Mark Giguere
School of Information Science and Policy
SUNY at Albany, NY

Towards Designing an Expert System to Map Mathematics Classificatory Structures

Iyer, H., Giguere, M.: Towards designing an expert system to map mathematics classificatory structures.

Knowl.Org. 22(1995)No.3/4, p.141-147, 9 refs

The convertibility of ordering systems such as thesauri and classification schemes have been investigated for many years. The focus has so far been more on thesauri than on classification schemes. Classification schemes too could differ from one another in several ways: in their structural, semantic, lexical and notational features. These incompatibilities make multiple catalog search difficult for the users. The Dewey Decimal Classification is one of the widely used schemes worldwide that encompasses all of knowledge whereas the Mathematics Subject Classification scheme published by the American Mathematical Society is a special classification scheme that is used in several AMS publications, notably the Mathematics Review. An interface that enables mathematicians to access library collections organized with the Dewey Decimal Classification, using the AMS scheme as an interface will certainly be useful. This paper suggests a prototype expert system interface to map the MSC scheme on to the mathematics (510 schedule of DDC20 and presents the work done so far towards this end. Compares the two schemes and discusses the mapping strategies/rules developed and the features of the prototype expert system design.

(Authors)

1. Introduction

The convertibility of ordering systems, such as classification schemes and thesauri have been investigated for many years. The focus has so far been more on thesauri than on classification schemes. Several methods have been adopted for combining thesauri, of which the most common are mentioned here.

1) Linking domain specific thesauri with a parent umbrella thesaurus. The parent thesaurus contains all terms, and the more specific thesauri draw subsets of terms from it. This is similar to a combining approach used by classification schemes; general classification schemes extend over all of knowledge, and detailed specialized schedules are developed for different disciplines within their structural framework.

2) Translating terms from the form of one thesaurus to another by means of an intermediate neutral language, called a „switching language“ or „intermediate lexicon.“ The terms from the source language are converted to the switching language, and then translated to the target language. Neville (1970) proposed that a unique code number be assigned to every concept in a subject field, as



Hemalata Iyer is a member of the faculty at the School of Information Science and Policy SUNY at Albany. She had previously taught classification and indexing at Atlanta, GA and at top Indian universities. She specializes in information organization and has published one book and several articles.



Mark Giguere is Electronic Records Manager for the City of Philadelphia and is now about to complete his dissertation in the interdisciplinary IS program at SUNY, Albany. He has an MLIS from the Univ. of South Carolina and a Bachelor of Arts degree in the geological sciences with extended research in geophysics at Cornell University.

a means of switching terms in one thesaurus to those in another.

The VSS system developed at the Battelle Memorial Laboratories is an important landmark. It contains over a million terms, from 15 vocabularies, divided into four modules: physical sciences, life sciences, business and social sciences. It enables the user to choose the appropriate database to his query, and translates the search statement into the vocabulary of which ever database is searched (Chamis, 1991). Another example is the UMLS, the Unified Medical Language System (Tuttle, 1989, 1990, 1992). It is a system developed at the National Library of Medicine (NLM), to reconcile the various biomedical vocabularies and classification systems. These systems were integrated into a metathesaurus of concepts, which constitutes the source or foundation of UMLS. UMLS itself is a browsable machine readable reference tool, containing words with their definitions and synonyms, their hierarchical and associative relationships, and their occurrence in databases a number of times.

Like thesauri, classification schemes too may differ from one another in several ways: in their structural, semantic and lexical features. They may cover different subject domains, and even those of the same domain may differ in their scope and coverage. They may have semantic differences that are caused by variations in conceptual structuring. Levels of specificity may vary between schemes. One may find „Trees“ to be specific enough, while another may list individual species. The terminology used may differ too. One may use technical terms and the other layman's vocabulary. Differences may also arise due to syntactic features, such as the word order of terms; lexical differences such as variation in spelling, noun or verb form etc.

In summary, the incompatibilities that occur at structural, conceptual and terminological levels make a multiple catalog search impossible. The users are burdened having to learn unfamiliar classification scheme features in order to efficiently search catalogs and databases. There are a few general schemes that cover all of knowledge of which the Dewey Decimal Classification (DDC), and the Library of Congress Classification (LCC) are widely used in American libraries. These are also used in the MARC records. Besides the general schemes there are also several specialized schemes devoted to a specific discipline or to a group of disciplines. The Mathematics Subject Classification scheme published by the American Mathematical Society is one such scheme. The purpose of this paper is to analyze the MSC classification scheme and the mathematics schedule of the Dewey Decimal Classification scheme and develop mapping strategies for a prototype expert system that would take user input of MSC numbers and output either equivalent DDC numbers where possible or recommend broader, narrower or coordinate DDC classes. It also describes the features of Cxprt, the expert system shell that lends itself to classification mapping.

2. The AMS Mathematics Subject Classification (MSC) Scheme

The MSC scheme is a pragmatic one based upon the current literature of mathematics published in the professional journals. It is a specialized scheme with a primary focus on math while all other disciplines are peripheral. The papers in Mathematical Reviews and other publications of the American Mathematical Society are classified with this scheme. Thus it is used for classification of surrogates as opposed to physical documents.

Bartle observed years ago that „this system has no resemblance to either the Dewey or the Library of Congress system, partly because it is right up-to-date, partly because it was made by the mathematicians, partly because it is designed for papers not for books, and partly because it does not take into consideration many problems that a library classification must consider“ (Bartle, 1960).

The AMS 1991 Mathematics Subject Classification scheme is divided into 94 broad classes. It has three levels of division: the class, subclass and a further specific subclass. However it is not altogether hierarchical. The notation employs a two-digit number for the main class, followed by a letter and a two-digit number for the specific subclass e.g.,

Main Class

Computer Science

Subclass

Theory of computing

Specific subclass

VLSI algorithms

Math bibliographers in academic libraries will agree that mathematicians are more familiar and comfortable

with the MSC scheme than with a general scheme like the DDC. Their familiarity stems from its use in several of the MSC publications, notably the Mathematics Review. With the proliferation of online catalogs and the possibility of access via the Internet, the bibliographic world has come to the scholar's desk. An interface that enables mathematicians to access library collections using the MSC classification scheme as an interface will certainly be very helpful.

This paper suggests an expert system interface that maps the MSC scheme on to the mathematics (510) schedule of DDC20 to facilitate access to mathematics collections in libraries. It is still in the initial phase of development and this paper presents preliminary work done in that direction. Unlike the MSC scheme DDC is used to classify books and documents. DDC is a hierarchically structured scheme with divisions into ten subclasses at each level. The hierarchical notation consists of arabic numerals with decimals.

Both schemes cover the traditionally-accepted divisions of mathematics, familiar to both layman and mathematician alike, such as algebra, analysis, geometry etc. Superficially there seems to be some similarity in the organization of the MSC scheme and DDC in that the traditional divisions of mathematics are covered. MSC however is far more detailed and covers newer areas of mathematics. Figure 1 presents the sections in the MSC scheme and corresponding divisions in the DDC.

MSC Classification	DDC20
(00-01) General (mathematics)	510 Mathematics (General)
(02-04) Logic and Foundations	511 Generalities (covers Logic)
(05-22) Algebra	512 Algebra
(26-49) Analysis	515 Analysis
(51-53) Geometry	516 Geometry
(54-57) Topology	514 Topology
(60-62) Probability and Statistics	519 Probability and applied mathematics
(65-94) Applied mathematics	

Figure 1: Structure of MSC classification and DDC20

There is similarity in the divisions of mathematics covered by both schemes with the exception of arithmetic (513) in DDC which is omitted in the MSC scheme. Besides, the latter lays greater emphasis on applied mathematics. Since they were designed for different purposes they differ in their level of specificity and their emphasis.

3. Mapping Strategies

The two schemes were analyzed and the following mapping strategies were identified. The mapping rules are derived for each of these types.

1) Exact matches:

These are instances where the MSC scheme has an exact corresponding number in DDC at the same level of

6. Suggested Frame Implementation

Typical expert system shells are predicated on the identification of objects termed „frames“, that may possess identifying pieces of characteristic information termed „slots.“ The objects that have been identified as frames in this implementation were done so according to the object management technique [OMT] (Rumbaugh et al 1991). An object can be thought of as an application-domain concept that can be uniquely identified on the basis of its inherent identity. This identity is determined by an object's attributes and its operations. For the remainder of this discussion, object and frame can be considered interchangeable. In the proposed CxPert prototype implement, two frames will be defined, the initial MSC term-frame and the destination DDC term-frame. The slots characterizing each of these frames will represent appropriate pieces of the notation and the accompanying class description terminology. Thus the MSC term-frame 20K12 represents the notation „Ulm sequences“. The following illustrates the correspondence between slots and notation:

slot: superclass 20
 slot: class K
 slot: subclass 12.

Figure 2 displays an OMT methodology representation of the objects MSC term and DDC term and their corresponding operations. Each frame object, represented as boxes according to OMT methodology, consists of an object name, located at the top of the box, object attributes (slots in a frame-based implementation) listed below, and inherent object operations at the bottom. Because of the simplicity of this implementation, both the initial and destination frames can be thought of as generalized subclasses of the superclass Term, which possesses vocabulary/notation pair attributes. The operations associated with this superclass, which also apply to each of the generalized subclasses, allows for the modification of each of the individual classification systems, via the addition, deletion, or editing of various vocabulary/notation pairs. The only way in which these subclasses differ from one another is in their operations. On the basis of shared common vocabulary, an MSC term may map to a corresponding DDC term (i.e., the MSC term-operation „Map forward“), or vice versa (i.e., the DDC term operation „Map backward“).

The cardinality and optionality of interactions between the objects serve to embody the mapping strategies identified in earlier sections. The possibility of zero or many to zero or many forward and backward mapping relationships provides for the realization of strategies: (1) Exact matches, (4) Many to one, (5) Cyclic, (6) No matches. The optional recursive mapping operations within a particular system allow for the realization of cyclic mapping strategies (2) Specific to general and (3) General to specific.

It should be evident that these categories of mapping strategies are not mutually exclusive, but rather, illustrative of the associations indicated in Figure 2. Consider, for instance, the example given for exact matches (see Map-

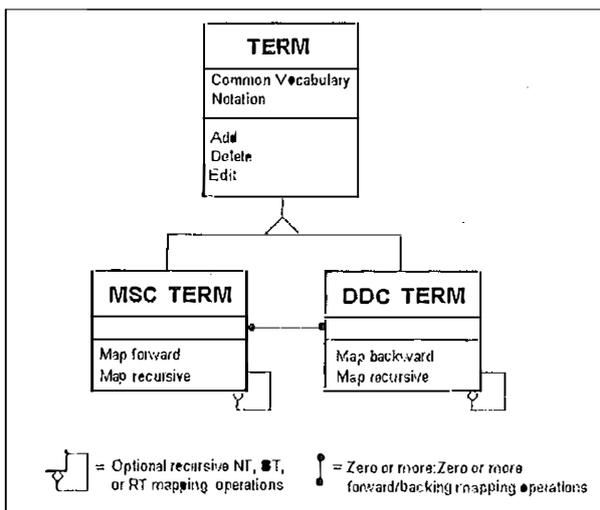


Fig. 2: Object-oriented frame analysis using OMT methodology (Rumbaugh et al. 1991)

ping Strategies - 1), in which mappings occur at differing levels of specificity. This is in part due to the fact that concept linkages occur as explicit pairings in MSC, the vocabulary chosen for matching purposes can come from a variety of places on the DDC term side of the system (e.g., heading, scope notes, instructions).

Another example of the overlap of the cited mapping strategies can be found in the examples (for examples for cyclic mapping see Mapping Strategies - 5). In this example of inter-system, interclass recursion, the MSC notation 53A45 maps to the more general DDC class 515.63 using an exact mapping strategy, but includes an additional general-to-specific recursive mapping strategy iteration.

Similarly, the example in the Many-to-one mapping strategy (see Mapping Strategies - 4) illustrates an instance of inter-system, inter-class recursion, due to the occurrence of the concept „queueing theory“ in three different classes of the MSC. Hypothetically when the end user inputs the MSC number 68M20 that belongs to the class computer science, the corresponding DDC number 519.62 is output, which belongs to the class probability. Since the main classes do not match, inter-system interclass recursion looks for other occurrences of queueing theory back again in the MSC scheme and thus maps two other class occurrences.

7. Basic Rule Structure

The basic structure of the KRL rules in this implementation will be derived from the mapping strategies described in the preceding Mapping Strategies section. Based on appropriate user input, these rules would initiate text retrievals of notation and descriptors from the appropriate database of terms (See hooks under „Special Features“). The reason the rules in this sort of implementation can be so simplistic is that the majority of the intellectual effort is already embodied in the organization of both classification systems.

It should be noted that the shell allows for nesting of rules, hence, if the end user inputs a complete (i.e., full notational specificity mapping at the class level) the MSC-term for which there is no direct mapping, rule-nesting will allow the search engine to reiterate the process, operating on the next highest level of notational complexity. This type of rule cycling would potentially allow the inference engine to recommend broad-class DDC term-notations via an advice function (See „Special Features“ section). Note that complete specificity in the input notation is not required by the system, thus allowing the end user the ability to „browse“ other less-notationally-specific options when confronted with an MSC-term entry that produces no results or advice.

The rule-nesting capability, when combined with the inference engine's backward chaining abilities, also provides additional capabilities for recommending alternate output notations to the enduser. Imagine, for example, the following scenario:

The user inputs a complete MSC term (i.e., specific to the subclass level) that has no direct mapping to a DDC-term. The nesting of rules then broadens the mapping to a class level, at which level a mapping to a DDC-term is found. This DDC-term collocation is also pointed to by another, more general MSC class or superclass which also points to a second DDC point of collocation. The backward chaining abilities of the inference engine would allow movement back from the DDC term that results from the rule's first iteration, to a new MSC term, that might connect to another part of the DDC notation. In this way, it would also be possible for the system to recommend coordinate relations to the DDC term, on the output-side of the application.

8. Special CxPert Features

There are several „canned“ capabilities embodied in the CxPert shell that make it particularly applicable to the development of the prototype described above. The primary special feature of the software that will be used in the design of this prototype is its „HyperWindows“ module development environment. This module of the software allows for the construction of pop-up „point & click“ window boxes that will serve as the primary interface for the prototype end user. These windows will be based on the frame implementation discussed in the previous section.

The second major feature of the software used in this prototype is the ability to implement „hooks“ to external databases. The IN TEXT formatting function allows the inference engine, via the implemented knowledge base, to access external textfiles stored in either a flat file or database record structure.

The implementation of rules in KRL provides several methods for representing the certainty of various decisions made by the inference engine. The differing degrees of specificity that exist between MSC and DDC frequently results in mismatches between points of collocation.

While general-to-specific collocation mismatches might not be problematic, specific-to-general mismatches may introduce errors into the recommendation made by the system. It is for this reason, as well as the previously discussed potential for rules to „cycle“ between MSC and DDC, that KRL's ability to incorporate certainty factors in the recommendation made by the inference engine is particularly attractive, as it allows for a mechanism to communicate this uncertainty to the end user.

One important capability of this shell is its streamlined capability to develop online consultation systems that can assist the end user during the use of the system. In addition to the previously mentioned HyperWindows module to create user consultation windows, there are several hot keys (i.e., F1 - F-12) that can be softwired into functions supported by the shell. These include an EXPLAIN function (to provide an opportunity to explain an inference made by the engine of the shell), a WHY function (to provide reasons why a particular data query was made of the enduser by the system), a HELP function, and an ERROR function.

It should be noted that an end user evaluation of a prototype interface will be collected in the testing of the prototype so as to provide for enhancements in future prototype designs (See the following section).

9. Hypothetical System Operation

The user inputs the attributes of the MSC slots. The system maps forward for a corresponding DDC match. If an exact match is not found it strips the user input MSC number of the digits on the right to the subclass level and class levels and again maps forward for a match with the DDC number. Within this overall operational strategy is built in iterative and recursive mapping as indicated in the preceding sections. It should also be possible to create a transactional log of the expansion of interaction; this would allow for expansion of the rule base by automatically extracting end user expertise in cases where the expert system was not able to make a DDC class recommendation (i.e., what DDC classes did the OPAC retrieve from).

10. Proposed Prototype Testing Experimental Design. Research Hypothesis

When considering a rule-based expert system whose knowledge base is constructed by describing „mapping relationships“ between similar classes in the two classification systems, if exact matches in class description terminology occur. For example, the subject caption „projective geometry“ is the exact description of the MSC facet 53A20 and the DDC facet 516.5., it should be possible to exploit these instances of classification systems' collocation, in conjunction with the hierarchical force embodied in the internal structures of both MSC and DDC to recommend broader or narrower terms/numbers that could be used to modify online retrieval strategies.

Based on this theoretical expectation, the research hypothesis for the testing of the envisioned prototype system is stated as follows: an expert system, of the type described, when used in conjunction with a research library's online public access catalog [OPAC], would allow patrons using the mathematics portion of the collection to construct online search strategies of higher retrieval efficiency as measured by precision. Precision is defined as: (# of relevant retrieved items)/(# of items retrieved) (Cleverdon 1970).

11. Methodology

A two-part, multi-method, experimental methodology is proposed in the evaluation of the aforementioned research hypothesis. The first portion of the methodology attempts to evaluate the accuracy of the prototype's performance, while the second portion addresses the value of the prototype as a tool for enhancing the efficiency of online search strategies. Clearly, accurate functioning of the expert system is a necessary prerequisite for testing of the primary research hypothesis, and this will be discussed first.

Part I: Evaluation of Accuracy of Prototype Performance

This portion of the methodology addresses an experimental design for addressing the accuracy of the expert system's performance. The variable to be operationalized, accuracy of system performance, will be defined as „performance of the system that is in agreement with human experts.“

Towards this end, a sample of bibliographic items that have been classified under both systems will be used to address this issue of accurate expert system performance. Selection of this „faux“ collection (N=100) will proceed according to the following process. Items will be randomly selected, on the basis of sequential accession numbers, from the most recent annual contents of Mathematics Review (i.e., these items have necessarily been classified under MSC). It should then be possible to construct a search strategy that compares these randomly selected items with items found in OCLC's national bibliographic database. Items that are found to possess matches in the OCLC database can then be checked to see if they have been additionally catalogued under DDC via the presence of appropriate cataloguing content in field tag 082. Items found to possess additional DDC cataloguing will be selected for membership in the faux collection.

Having assembled a faux collection for system testing purposes, this portion of the experimental design will then consist of supplying the expert system prototype with the MSC-classification numbers for those items, and comparing the recommended DDC output of the system with the DDC cataloguing contained in the corresponding OCLC record. It is anticipated that differences in the level of precision between the two systems will greatly skew the potential for exact matches in descriptive terminology. For example, the MSC subclass 20K12 (i.e., „Ulm se-

quences“) has no direct matches in DDC. The MSC hierarchical class superstructure related to this subclass is as follows:

Class	20-XX Group theory & generalizations
Subclass	20Kxx Abelian groups
Specific subclass	20K12 Ulm sequences.

While a direct terminology match with DDC does not exist at the specific subclass level, one does exist at the subclass level, with MSC class 20Kxx mapping to DDC class 512.2 (i.e., „Abelian groups“ appears in the scope notes of the DDC-subclass „Group and group theories“). Therefore, if the MSC number 20K12 was input to the expert system, and the prototype output DDC number 512.2, this output would be interpreted as being technically correct for evaluation purposes, albeit to lesser degree of specificity. With this in mind, the accuracy of the expert system's mappings will be evaluated at the class, subclass, and specific subclass levels. The accuracy of system performance at these levels of classification hierarchy will provide baseline data that will serve to modify the outcomes of the second portion of the described methodology.

Part II: Evaluation of Prototype Use on Search Strategy Effectiveness

The second portion of the experimental design involves testing the effect that simultaneous use of the expert system with a research library OPAC has on the retrieval efficiency of end user constructed search strategies. The variable to be operationalized in portion of the methodology include:

relevance - an enduser defined measure of how well a retrieved item fill a particular information need (so as to be used in comparative precision calculations), and participant self-characterization self-description of experiment participant character traits regarding their level of familiarity with mathematics, MSC classification, and the formulation on OPAC search strategies. The sample of experiment participants (N=15,20) will be a purposive sample of DDC organized, research library patrons chosen on the basis of their circulation characteristics. Potential sites for library testing include the New York State Library (i.e., the 20th largest research library in the United States according to ARL) or the Folsom Library at the Rensselaer Polytechnic Institute.

12. Conclusion

The widespread development of microcomputer-based, commercial expert system shells has provided classificationists with the ability to create, via the use of limited knowledge bases (i.e., in this instance, the use of collocated controlled vocabulary), prototype systems by which the inductive reasoning abilities of the software's search engine is able to exploit the hierarchical force inherent in faceted classification systems. This discussion lays the groundwork for the development and testing of

such a rudimentary classification mapping application.

The application of such a system must be viewed in the context of a scenario whereby a user has already identified a relevant document that has been classified under MSC. The MSC class number of that item would then be input into the prototype so as to recommend other DDC classes that might contain relevant materials. The ultimate benefit of such a system would lie in its use as an auxiliary device by patrons constructing search strategies on a research library's OPAC. It is believed that such an application will enhance end user access to specialized materials by providing a relatively transparent, artificially intelligent interface to unfamiliar classification systems.

References

- (1) Babbie, Earl: *The Practice of Social Research*. 6th ed. Belmont, CA: Wadsworth Publ. 1992.
- (2) Bartle, Robert G.: *One Mathematician Looks at the Classification of Mathematics*. In: *The Role of Classification in the Modern American Library*; Institute University of Illinois Graduate Library School of Library Science, Nov. 1-4, 1959. Champaign, IL: Illini Union Bookstore 1960: 93-102.
- (3) Chamis, Alice: *Vocabulary Control and Search Strategies in Online Searching*. Westport, CT: Greenwood Press.
- (4) Cleverdon, C.W.: *Retrieval of Information*. *Works Management* 7((1970)July/August, p.12-15.
- (5) CxPert User's Guide. Crofton, MD: Software Plus, Ltd. 1989.
- (6) Neville, H.H.: *Feasibility Study of a Scheme for Reconciling Thesauri Covering Common Subject*. *J. Doc.*. 26(1970)p. 313-337.
- (7) Rumbaugh, James E. et al.: *Object-Oriented Modelling and Design*. Englewood Cliffs, NJ: Prentice-Hall 1991.
- (8) Tuttle, M.S. et al.: *Using Meta-I: The First Version of the UMLS Metathesaurus*. In: Kingsland, L. C. (Ed.): *Proc. Thirteenth Annual Symposium on Computer Applications in Medical Care*. New York: IEEE Computer Society Press 1989. p.483-487
- (9) Tuttle, M.S. et al.: *Implementing Meta-I: The First Version of the UMLS Metathesaurus*. In: Miller, R.A. (Eds): *Proc. Fourteenth Annual Symposium on Computer-Applications in Medical Care*. New York: IEEE Computer Society Press 1990. p.331-335

Prof. Hemalata Iyer, School of Information Science and Policy, State University of New York at Albany, USA. Phone: (518) 442-5116. H1651@cnsvox.albany.edu

Mark Giguere, Dept. of Records, 166 City Hall, Philadelphia, PA 19107. email: Mg3721@cnsvox.albany-edu