

On the Use of the Dempster Shafer Model
in Information Indexing and Retrieval Applications

*Shimon Schocken*¹

Stern School of Business, New York University

*Robert A. Hummel*²

Courant Institute of Mathematical Sciences, New York University

October 12, 1992

Working Paper Series

STERN IS-92-27

¹ Department of Information Systems, Management Education Center, Stern School of Business, New York University, 44 W. 4th Street, New York, NY 10003.

² Courant Institute of Mathematical Sciences, New York University, 251 Mercer Street, New York, NY 10012.

On the Use of the Dempster Shafer Model in Information Indexing and Retrieval Applications

The Dempster Shafer theory of evidence concerns the elicitation and manipulation of degrees of belief rendered by multiple sources of evidence to a common set of propositions. Information indexing and retrieval applications use a variety of quantitative means – both probabilistic and quasi-probabilistic – to represent and manipulate relevance numbers and index vectors. Recently, several proposals were made to use the Dempster Shafer model as a relevance calculus in such applications. The paper provides a critical review of these proposals, pointing at several theoretical caveats and suggesting ways to resolve them. The methodology is based on expounding a canonical indexing model whose relevance measures and combination mechanisms are shown to be isomorphic to Shafer's belief functions and to Dempster's rule, respectively. Hence, the paper has two objectives: (i) to describe and resolve some caveats in the way the Dempster Shafer theory is applied to information indexing and retrieval, and (ii) to provide an intuitive interpretation of the Dempster Shafer theory, as it unfolds in the simple context of a canonical indexing model.

Keywords: Theory of evidence, Dempster Shafer model, relevance measures, information indexing and retrieval.

1 Introduction

Consider a finite and exhaustive set of mutually-exclusive propositions and a body of evidence that supports some subsets of propositions and discounts others. Many theories were put forward to describe how one should represent and update one's degrees of belief in such propositions when new or additional evidence is brought to bear. The classical approach is to cast degrees of belief as probabilities – a set of numbers between 0 and 1 that obeys the axioms of subjective probability – and use Bayesian inference rules to revise them in light of new evidence. One problem with this approach is that it doesn't offer a clear way to model the various degrees of 'uncommitted beliefs,' or 'second order uncertainties,' that characterize most realistic inference problems. For example, consider the extreme case of 'insufficient reason,' in which one knows absolutely nothing about a given set of n propositions. The common solution, which goes back to LaPlace, is to assign a degree of belief of $1/n$ to each of the propositions under consideration. Incidentally, this is also the solution that emerges from maximizing the unconstrained entropy function associated with the n unknown probabilities.

Over the years, many students of belief revision theories have objected to this crude quantification of insufficient reason. Why, the argument goes, should ignorance be translated to the strong statement that every proposition (or state of nature) is equally likely? This criticism has led to several alternative models that attempt to capture the elusive notion of uncommitted belief by modifying the axiomatic framework of probability theory. Perhaps the best known model in this category is the 'theory of evidence,' originated by Dempster's (1967, 1967a) , work on upper and lower probabilities. Dempster's ideas, which were

based on a frequentist view of inference, were refined and extended by Shafer (1976), who also gave them a subjective interpretation. This led to the Dempster Shafer model – an elaborate formalism for representing and revising degrees of support rendered by multiple sources of evidence to a common set of propositions¹.

When the work of Dempster and Shafer was ‘discovered’ by the artificial intelligence community, it immediately stirred a considerable interest in two application areas in which normative models of belief formation play a key role: expert systems, and information indexing and retrieval systems. For expert systems, the Dempster Shafer (DS) model provides a mathematically-sound model for representing and manipulating rule-based degrees of belief, an area that was traditionally dominated by ad-hoc belief revision calculi whose relationship to probability theory was at best murky. For information indexing and retrieval systems, the DS model can be used as a *relevance calculus*, designed to quantify and revise the degrees of relevance between documents, keywords, and user-supplied queries.

This line of thought has led to the development of several DS-based information indexing and retrieval applications. For example, Biswas, Bezdek, Marques, and Subramanian (1987) built a document retrieval system in which the relevance of documents to taxonomical classes was measured and manipulated, respectively, by belief functions and Dempster’s rule: “*We choose to define similarity functions based on the Dempster Shafer theory of evidence ... one of the advantages of this approach is that it reflects the process of belief revision and updating just as in human reasoning processes.*” (Biswas et al, 1987). Coming from a different direction, Turtle and Croft (1991) describe a canonical representation in which relevance is handled through inference networks that are structured as directed

¹In this paper, the terms the *Dempster Shafer theory of evidence* and the *Dempster Shafer model* are used interchangeably.

acyclic graphs. The nodes in the networks correspond to keywords, documents, and queries, and “the arcs joining the nodes are interpreted as assertions that the parent node provides support for the child node.” Turtle and Croft proposed to operationalize these degrees of support through either subjective probabilities, or DS belief functions. A similar approach was undertaken in RUBRIC, a full-text information retrieval system described by Tong and Shapiro (1985) . RUBRIC can be instantiated to operate with several alternative relevance calculi, the DS model being a prime example.

The importance of such applications is obvious, as they attempt to take the DS model ‘out of the lab’ and implement it in realistic settings. In doing so, however, many adopters of the DS model have taken the model’s validity for granted, either explicitly or implicitly. With that in mind, it is important to point out that both the cognitive and the normative roots of the DS model are still a matter of intense controversy: whereas Shafer (1987) argues that the theory of evidence is a natural extension of probability theory , many critics, e.g. Lindley (1987) , view it as a reformulated version of a specialized, albeit interesting, case of classical probabilistic inference. The debate is not helped by the somewhat forbidding notation of the DS model, which prevents an intuitive understanding of its underlying structure and philosophy.

In fact, the gap between the theory and practice of the DS model seems to be two-directional. On the one hand, many practitioners believe that the normative correctness of the DS model is a ‘closed case,’ proceeding to implement it without questioning its underlying rationale. On the other hand, many researchers try to defend the DS model on complex philosophical and mathematical grounds, without realizing that simpler justifications can be found *in the field*, i.e. in the way the model is actually used in certain canonical settings.

The latter point is worth emphasizing: a close examination of certain applications of the DS model can provide not only a better understanding of the model, but, furthermore, a compelling normative justification.

The plan of the paper is depicted in figure 1 and described as follows. §2 presents the notion of index vectors and the challenge of eliciting and measuring relevance in a normative, rather than ad-hoc, fashion. §3 gives an overview of the DS model, as it unfolds in the context of a typical information indexing and retrieval (IR) application. This sets the stage to four critical questions regarding the theoretical fit between the general features of the DS model and the specific requirements of IR applications. In order to answer these questions, §4 presents a canonical indexing model in which the notions of lexicons, taxonomies, and relevance, are treated formally and unambiguously. It is then shown that the canonical model is completely isomorphic to the DS model, leading to a new intuitive understanding of the latter. §5 offers concluding remarks about the implications of the research on the DS model and on IR applications.

Put figure 1 around here

2 The Problem

Models of bibliographical indexing concern the construction of data structures that enable rapid content-based access to collections of documents. Given a document, on the one hand, and a *keyword lexicon*, on the other, the goal of the indexing model is to select a subset of keywords that ‘best’ describes the document to its prospective users. Since some

keywords are more relevant to the document than others, a numeric scale is often used to express the strength of association between the document and the selected keywords. The result is an *index vector*, consisting of pairs of keywords and their respective relevance weights. Several models exist for representing and manipulating such relevance vectors, and the reader is referred to Salton and McGill (1983) and to Salton and Buckley (1988) for comprehensive treatments of the general approach to the subject.

Formally, let D be a set of documents about a certain domain of interest, and let $\mathcal{K} = \{k_1, \dots, k_m\}$ be a lexicon, or a set of domain keywords. The index of each document $d \in D$ is a set of pairs of the form:

$$S_d = \{(K_1, r_1), \dots, (K_n, r_n)\} \quad (1)$$

where $K_i \subseteq \mathcal{K}$ and $0 \leq r_i \leq 1, i = 1, \dots, n$. The K_i 's are *lexical subsets*, representing different groupings of keywords, and the r_i 's are called *relevance numbers*. Taken together, the pair $(K_i, r_i) \in S_d$ says that the degree of relevance between the document d and the lexical subset K_i is r_i . Had we restricted the K_i 's to be singletons only, (1) would become the familiar 'term-weight vectors' that are normally used in information indexing and retrieval applications. Further, had we required that all the r_i be 0 or 1 only, (1) would be reduced to the familiar keyword list (also called 'subject headings') that is normally used to classify articles in professional journals. Given the obvious simplicity of a *Boolean* indexing scheme, why bother about developing formalisms for *weighted* indexing?

The answer is that relevance is a subjective and composite relation, based on an aggregation

of several indexing opinions. Specifically, each document has many *classifiers*, or discerning characteristics, that determine its relevance. For example, the *title* of a document can suggest one index, whereas the *abstract* can suggest another. Other aspects of the document, obtained through lexical, linguistic, and citations, analyses will yield additional indexing opinions that must be taken into consideration. Hence, even if the individual opinions were forced to be binary, their aggregation would probably induce a continuous index. In addition, the indexing opinions are not cast automatically; rather, they are elicited from human catalogers who inject yet another level of uncertainty and subjectivity to the indexing process. That is, when two catalogers are given access to the same classifier as background information, they may well supply two different (but hopefully similar) indexing opinions.

Different IR applications use different models to handle this pluralism in a formal way. From a theoretical perspective, the credibility of these models hinges on their capacity to elicit, represent, and synthesize, relevance opinions in a *normative*, rather than *ad-hoc*, fashion. In order to do so, the relevance numbers and the rules that combine them must be given a compelling interpretation. So far, the leading interpretation in the study of relevance has been probabilistic, beginning with the seminal work of Maron and Kuhns (1960). Recently, however, several attempts were made to handle relevance in IR applications using the Dempster Shafer model, which is widely considered to be a less restricted extension of probability theory. The strengths and weaknesses of the latter approach are discussed in the next section.

3 A Dempster-Shafer Indexing Model

The DS theory of evidence concerns the representation and manipulation of degrees of support rendered by different sources of evidence to a common set of propositions, denoted θ and called the *frame of discernment*. In contrast to a standard Bayesian design, in which degrees of belief are normally assigned to elements of θ directly, the DS model assigns degrees of belief to *subsets* of propositions, i.e. to members of the power set 2^θ , also called ‘possibilities.’ The DS model offers several complementary ways to express evidential support in possibilities. In particular, the model defines three mappings from 2^θ to $[0, 1]$ termed *mass*, *belief*, and *plausibility*, functions. The three mappings are mathematically equivalent in the sense that knowledge of any one of them (for every possibility) can be used to compute the other two. Therefore, we view them as alternative means to keep score of the same primitive set of degrees of support. In the standard model, when several sources of evidence support a common set of possibilities (the support can be cast in terms of either mass, belief, or plausibility functions), the overall support in the possibilities is computed through Dempster’s rule of combination.

What is the nexus of the DS model and information indexing and retrieval applications? In one way or another, all DS-based IR applications are based on the following premises: (i) The DS notion of *degrees of support* can be used to operationalize the IR notion of *relevance numbers*; and (ii) When two or more classification criteria supply different sets of relevance numbers concerning the same document, Dempster’s rule provides a plausible mechanism to combine them into a composite index (said otherwise: revise the relevance of the document to certain keywords in light of new evidence). The goal of this section is

to motivate a critical analysis of these premises. Specifically, we intend to:

- First, provide a rigorous but accessible overview of the DS model, as it unfolds in the familiar context of an IR application;
- Second, present a series of questions regarding the *theoretical fit* between the general features of the DS model and the specific characteristics of IR applications.

The Frame of Discernment: The frame of discernment θ is an exhaustive set of mutually exclusive elements that can be interpreted as hypotheses, propositions, or simply ‘labels.’ The power-set that contains all the subsets of θ (including θ itself and the empty set) is denoted 2^θ . In general, the semantics of the labels depends on the context in which the DS model is applied. In information indexing and retrieval applications, the frame of discernment is normally taken to be a keywords lexicon $\mathcal{K} = \{k_1, \dots, k_n\}$. To illustrate, a lexicon that supports a collection of documents about modern art might be $\mathcal{K} = \{\text{Arp}, \text{Braque}, \text{Cezanne}, \dots, \text{Zorn}\}$, enumerating all the major artists of the Twentieth Century. The power set in this case is $2^\mathcal{K} = \{\{\text{Arp}\}, \{\text{Braque}\}, \{\text{Cezanne}\}, \dots, \{\text{Arp}, \text{Braque}\}, \{\text{Arp}, \text{Cezanne}\}, \{\text{Braque}, \text{Cezanne}\}, \dots, \{\text{Arp}, \text{Braque}, \text{Cezanne}\}, \dots, \emptyset, \mathcal{K}\}$, the last two elements being the empty set and \mathcal{K} itself. Each element in $2^\mathcal{K}$ represents a disjunction of keywords, denoted hereafter a *lexical subset*. The act of indexing a document using \mathcal{K} amounts to choosing, among all the indexing possibilities in $2^\mathcal{K}$, the one or more lexical subsets that best describe the document to its potential users.

For example, suppose that an art scholar is asked to index the document “*The Influence of Cezanne on early Cubism*” using \mathcal{K} , based on partial information such as the document’s title or abstract. Without loss of generality, assume that (i) the main focus of

the document is Cezanne; and (ii) the only Cubist artists in the lexicon are Braque and Picasso. Under these assumptions, the scholar will probably supply an index of the form $S = \{(\{\text{Cezanne}\}, r_1), (\{\text{Braque, Picasso}\}, r_2)\}$, with $r_1 > r_2$. This would entail the following information: (i) the document is relevant to Cezanne; (ii) it is also relevant, to a lesser extent, to either Braque or to Picasso. This is quite different from the indexing opinion $S' = \{(\{\text{Cezanne}\}, r_1), (\{\text{Braque}\}, r_2), (\{\text{Picasso}\}, r_2)\}$, which would be more appropriate if the document's title were, say, "*The Influence of Cezanne on the early work of Braque and Picasso*".

We arrive at our first question:

Question Q1: When the DS model is applied to information indexing and retrieval applications, the keyword lexicon \mathcal{K} is taken to be the *frame of discernment*, and indexing possibilities are taken to be elements of the lexical *power set* $2^{\mathcal{K}}$. What are the taxonomical implications and limitations of this representation?

To motivate this question, consider again the document "*The Influence of Cezanne on early Cubism*". Note that the most reasonable index of this document would be $S'' = \{(\{\text{Cezanne}\}, r_1), (\{\text{Cubism}\}, r_2)\}$, especially if the document's abstract makes no references to specific artists other than Cezanne. However, Cubism is not an element of the original lexicon \mathcal{K} , so it doesn't entail an indexing possibility. To solve the problem, we may want to extend the original frame of discernment, creating a lexicon of the form $\mathcal{K}' = \mathcal{K} \cup \{\text{Cubism}\}$. However, the keywords Braque, Picasso and Cubism, have a great deal in common from a bibliographical standpoint. Therefore, \mathcal{K}' is not a valid frame of discernment, because

some of its elements are no longer mutually exclusive. Before we present a solution to this problem, we have to be very specific about the proper relationship among *frames of discernment*, *keyword lexicons*, and *taxonomies of classes*. We'll return to this issue in section 3, where an answer to Q1 is given.

Mass Functions: A mapping $m : 2^\theta \rightarrow [0, 1]$ with the properties:

$$m(\emptyset) = 0 \tag{2}$$

$$\sum_{X \in 2^\theta} m(X) = 1 \tag{3}$$

is called a *mass function*². In the DS model, the mass $m(X)$ represents the degree to which a certain source of evidence supports the possibility X , where $X \subseteq \theta$. As a convention, the mass which is 'left over' after all the *proper* subsets of θ have been assigned masses is allocated to θ itself and denoted the *uncommitted belief* displayed by m , or $m(\theta)$.

In DS-based IR applications, where θ is taken to be a keywords lexicon \mathcal{K} , the mass $m(X)$ is taken to represent (to a first approximation that will be discussed shortly) a degree of relevance, or, more accurately, the degree of belief that the document is relevant to the lexical subset $X \subseteq \mathcal{K}$, according to a certain classifier. Hence, if a classifier (say, classifier number 1) supplies the relevance opinion $S_1 = \{(\{\text{Cezanne}\}, 0.6), (\{\text{Braque}, \text{Picasso}\}, 0.3)\}$, then the mass function that is induced by this opinion is defined as follows:

²Throughout the paper, upper case variables refer to sets and lower case variables refer to scalars.

$$\begin{aligned}
m_1(\{\text{Cezanne}\}) &= 0.6 \\
m_1(\{\text{Braque, Picasso}\}) &= 0.3 \\
m_1(\mathcal{K}) &= 0.1 \\
m_1(X) &= 0 \text{ for all other proper subsets of } \mathcal{K}
\end{aligned}
\tag{4}$$

Note that the uncommitted belief induced by the opinion is assigned by default to the frame of discernment by means of $m_1(\mathcal{K}) = 1 - 0.6 - 0.3 = 0.1$. The rationale for this assignment is as follows. If a certain classifier provides no information whatsoever about indexing possibilities, the classifier's 'ignorance' can be represented by the index $S = \{(\mathcal{K}, 1)\}$. This implies the mass function $m(\{\text{Arp, Braque, Cezanne, } \dots, \text{Zorn}\}) = 1$ and $m(X) = 0$ elsewhere, reflecting the (not very useful) opinion that the document is relevant to Arp, or to Braque, or to Cezanne, or to any other artist in the lexicon. Other classifiers can provide more focused relevance opinions, resulting with lower levels of $m(\mathcal{K})$. Hence, unlike a standard probabilistic design, where the notion of uncommitted belief is not well-defined, the DS model provides explicit means to quantify and manipulate it via $m(\mathcal{K})$. Although uncommitted beliefs, or 'second-order uncertainties,' can and have been treated in the standard framework of subjective probability, (e.g. Baron, 1987), there is no *simple* way to do it. The theory of evidence is unique in that it treats the notion of uncommitted belief explicitly, at the axiomatic level.

It's important to observe that mass functions represent indivisible, or atomic, degrees of belief. For example, the magnitudes of $m(\{\text{Braque, Picasso}\})$, $m(\{\text{Braque}\})$, and $m(\{\text{Picasso}\})$ are unrelated, and a mass function like $m(\{\text{Braque, Picasso}\}) = 0.9$, $m(\{\text{Braque}\}) = 0$, and $m(\{\text{Picasso}\}) = 0$ is not inconsistent with the theory. This particular function represents a cataloger who strongly believes that the document is relevant to either Braque or to Picasso, although he is not willing to say anything more specific

beyond this assessment.

But what does this notion of relevance *mean*? We arrive at our next question:

Question Q2: A mass function is a formal, domain-independent, component of the DS model. Relevance is an informal, but highly intuitive, concept that plays a key role in information indexing and retrieval applications. If a mass function is taken to represent relevance, then what is the *exact* semantics of this representation? Said otherwise, what *type* of relevance do mass functions represent?

Question Q2 suggests the premise that mass functions are not necessarily a natural representation of the intuitive notion of relevance, as it is typically construed in information indexing and retrieval applications. For example, if mass functions are used to represent relevance, then the relevance numbers in each index must sum up to 1. That is, the set of allowable indexing opinions $\{(K_1, r_1), \dots, (K_n, r_n)\}$ is constrained by $\sum_1^n r_i = 1$. Many would argue that this constraint doesn't make sense, and that an indexing opinion like, say, $\{(\{Albers\}, 0.8), (\{Kandisnki\}, 0.4), (\{Klee\}, 0.4)\}$ is perfectly reasonable. The only 'wrong' thing about this opinion is that it is inconsistent with the DS notion of a mass function, but this seems to be a limitation of the model's application, not of the opinion.

One pragmatic solution is to treat the relevance numbers not as *absolute*, but rather as *relative*, measures of subjective relevance. According to this position, the two indexes $S = \{(A, 0.8), (B, 0.4), (C, 0.4)\}$ and $S' = \{(A, 0.4), (B, 0.2), (C, 0.2)\}$ are equally informative, as both imply exactly the same relative information: the document is twice as relevant

to A as it is to B, and it is as relevant to B as it is to C. However, this immediately leads to another snag: according to the same principle, the index is also equivalent to $S'' = \{(A, 0.2), (B, 0.1), (C, 0.1)\}$. Yet S' and S'' reflect two different states of uncommitted belief (0.2 and 0.6, respectively), and thus they don't induce the same mass function.

To get around the problem, we can elicit uncommitted beliefs directly from the catalogers³. For example, having specified a relevance opinion, say $\{A, 0.8, B, 0.4, C, 0.4\}$, the cataloger can be asked to rate his confidence in the opinion on a scale of 0 to 1. If the confidence level is 1, the index is normalized to $\{A, 0.5, B, 0.25, C, 0.25\}$, reflecting an uncommitted belief of 0. If the confidence level is 0.8, the index is normalized to $\{A, 0.4, B, 0.2, C, 0.2\}$, reflecting an uncommitted belief of 0.2. In general, for any unconstrained indexing opinion $\{(K_1, r_1), \dots, (K_n, r_n)\}$ and a confidence level c , we can find a unique mass function $\{m(K_1), \dots, m(K_n), m(\mathcal{K})\}$ such that (i) the $m(K_i)$'s preserve the relative properties of the unconstrained r_i 's; and (ii) $m(\mathcal{K}) = 1 - c$.

The shift from an absolute to a relative scale of relevance has several justifications. First, a significant body of psychological and cognitive evidence indicates that relevance is indeed a relative property (Saracevic, 1975). Second, we are motivated by the observation that ultimately, an IR application must satisfy the information needs of library patrons, and that relevance numbers should be used pragmatically to that end. For example, according to Maron (1982)'s 'Ranking Principle', the chief objective of relevance numbers is to present to the patron a set of documents, sorted in decreasing order of perceived relevance to his or her query. A similar principle is used in diagnostic expert systems, where ordinal, rather than cardinal, degrees of beliefs are often used to guide the inference engine to

³In this section, the terms *classifier* and *cataloger* are used interchangeably. The distinction between the two terms is made explicit in the next section.

promising directions and to explain the system's reasoning to the people who consult it. If we accept Maron's Ranking Principle as a working assumption, then normalization is not an issue, since rankings are invariant under normalization. However, when multiple indexing opinions are aggregated into a pooled index (something that we haven't done yet), normalization becomes a tricky manipulation. Specifically, let S_1 and S_2 be two indexing opinions, \oplus an aggregation operator, and N a normalization operator. In many cases (depending on the specific definitions of \oplus and N), it can be shown that $N(S_1 \oplus S_2) \neq N(S_1) \oplus N(S_2)$, i.e. that N is not homomorphic.

In conclusion, we see that even though relevance numbers can be represented by mass functions, the representation has some theoretical caveats. Clearly, these limitations are related to the fact that we are still lacking explicit domain semantics. That is, we don't know yet what is the exact *meaning* of relevance numbers. This analysis is taken up in section 3, where a complete answer to question Q2 is given.

The Core: The *core* of a mass function $m : 2^\theta \rightarrow [0, 1]$ is the set of possibilities $X \in 2^\theta$ for which $m(X) > 0$. When the frame θ is taken to be a keyword lexicon \mathcal{K} , the core becomes a list of indexing possibilities, in the view of one particular classifier. For example, the core of the mass function induced by classifier 1 (Eqn. 4) is $C_1 = \{\{\text{Cezanne}\}, \{\text{Braque}, \text{Picasso}\}, \mathcal{K}\}$. Suppose now that the *same* document is indexed by another classifier (classifier no. 2), whose indexing opinion is captured by the following mass function:

$$\begin{aligned}
 m_2(\{\text{Picasso}\}) &= 0.8 \\
 m_2(\mathcal{K}) &= 0.2 \\
 m_2(X) &= 0 \text{ for all other proper subsets of } \mathcal{K}
 \end{aligned}
 \tag{5}$$

The core of this mass function is $C_2 = \{\{\text{Picasso}\}, \mathcal{K}\}$. Is there a credible way to combine the two indexing opinions (4),(5) into an aggregate index? As a first approximation, one can focus on all the lexical subsets that both classifiers agree are relevant to the document. In particular, if classifier 1 thinks that X is relevant and classifier 2 thinks that Y is relevant, then *both* classifiers agree that $X \cap Y$ is relevant (recall that both X and Y are interpreted as disjunctions of keywords). This leads to the following definition of a *pooled core*: Let $m_1, m_2 : 2^\theta \rightarrow [0, 1]$ be two mass functions with cores C_1 and C_2 . The pooled core $C = C_1 \oplus C_2$ will be:

$$C_1 \oplus C_2 = \{X \cap Y \mid X \in C_1, Y \in C_2, X \cap Y \neq \emptyset\}. \quad (6)$$

For example, the pooled core of $C_1 = \{\{\text{Cezanne}\}, \{\text{Braque, Picasso}\}, \mathcal{K}\}$ and $C_2 = \{\{\text{Picasso}\}, \mathcal{K}\}$ is $C_1 \oplus C_2 = \{\{\text{Cezanne}\}, \{\text{Picasso}\}, \{\text{Braque, Picasso}\}, \mathcal{K}\}$ ⁴. In general, then, the pooled core can be viewed as a first approximation of the degree of consensus or disagreement displayed by two independent indexing opinions. If $C_1 \oplus C_2 = C_1 = C_2$, we have a consensus regarding which possibilities are likely. If $C_1 \oplus C_2 = \emptyset$, the classifiers agree on nothing. If $C_1 \oplus C_2$ is not empty, we have an overlap of some opinions. Of course the problem of (6) is that it merely *identifies* areas of mutual agreement (or lack thereof) between two classifiers. In order to compute the *intensity* of such agreements, a more sensitive pooling mechanism is required. Dempster's rule provides one such mechanism.

Dempster's Rule: The most fundamental (and debateable) pillar of the DS model is the convention that once degrees of support are cast in terms of mass functions, Demp-

⁴Note that \mathcal{K} acts as an attractor, in that $A \cap \mathcal{K} = A$ for all $A \subseteq \mathcal{K}$.

ster's rule provides a proper mechanism to combine them. Let m_1 and m_2 be two mass functions defined over the same frame of discernment: $m_1, m_2 : 2^\theta \rightarrow [0, 1]$, with cores $C_1 = \{A_1, \dots, A_{n_1}\}$ and $C_2 = \{B_1, \dots, B_{n_2}\}$, respectively. Dempster's rule computes the pooled mass function $m = m_1 \oplus m_2 : 2^\theta \rightarrow [0, 1]$ as follows:

$$m'(X) = \sum_{A_i \cap B_j = X} m_1(A_i) \cdot m_2(B_j), \quad (7)$$

$$m(X) = \begin{cases} \frac{1}{1-m'(\emptyset)} \cdot m'(X) & X \neq \emptyset \\ 0 & X = \emptyset \end{cases} \quad (8)$$

The rationale behind (7-8) can be explicated through an 'intersection table.' In our two-classifiers scenario (4-5), the table has the following form:

	$m_1(\text{Cezanne}) = 0.6$	$m_1(\text{Picasso, Braque}) = 0.3$	$m_1(\mathcal{K}) = 0.1$
$m_2(\text{Picasso}) = 0.8$	$m'(\emptyset) = 0.48$	$m'(\text{Picasso}) = 0.24$	$m'(\text{Picasso}) = 0.08$
$m_2(\mathcal{K}) = 0.2$	$m'(\text{Cezanne}) = 0.12$	$m'(\text{Picasso, Braque}) = 0.06$	$m'(\mathcal{K}) = 0.02$

The top row of the table records the mass function of the first classifier excluding its zero elements, i.e. the set of values $m_1(A_1), \dots, m_1(A_{n_1})$ for elements A_i in the core C_1 . The left column of the table records the mass values of the second classifier for its core elements,

i.e. the set of values $m_2(B_1), \dots, m_2(B_{n_2})$ (The curly brackets are dropped for the sake of brevity, e.g. $m(\text{Picasso, Braque})$ stands for $m(\{\text{Picasso, Braque}\})$, etc.). Inside the table, the (i, j) 'th cell records the pooled mass contributed to $A_i \cap B_j$ by the pair A_i and B_j , which is taken to be the product $m(A_i) \cdot m(B_j)$. Using these entries and combining cells with equivalent intersections following (7-8), one obtains:

$$\begin{aligned}
 m'(\text{Cezanne}) &= 0.12, \\
 m'(\text{Picasso}) &= 0.24 + 0.08 = 0.32, \\
 m'(\text{Picasso, Braque}) &= 0.06, \\
 m'(\mathcal{K}) &= 0.02, \\
 m'(\emptyset) &= 0.48,
 \end{aligned} \tag{9}$$

After multiplying by $\frac{1}{1-m'(\emptyset)} = 1.02$ one obtains:

$$\begin{aligned}
 m(\text{Cezanne}) &= 0.23 \\
 m(\text{Picasso}) &= 0.62 \\
 m(\text{Picasso, Braque}) &= 0.11 \\
 m(\mathcal{K}) &= 0.04 \\
 m(\emptyset) &\stackrel{\text{def}}{=} 0
 \end{aligned} \tag{10}$$

Since the $m(\cdot)$'s sum up to 1 and $m(\emptyset) = 0$, the mapping $m = m_1 \oplus m_2$ that emerges from Dempster's rule is also a mass function, consistent with (3).

In words, Dempster's rule computes a measure of agreement between two sources of evidence concerning various possibilities drawn from a common frame of discernment. The rule is conservative in the sense that it focuses only on those possibilities that *both* sources support. The magnitude of the pooled support that a possibility X collects is computed by summing the products of the two masses $m_1(X)$ and $m_2(X)$, which explains the product

operator in (7). Because the sources of evidence express their opinions over 2^θ rather than over θ , a joint agreement on a possibility can occur in more than one way, i.e. whenever the two sources support supersets of X . This explains the summation operator in (7), which runs over all the possible supersets of X . Finally, when a pairing of two opinions results in a null possibility (the empty set), the multiplication of their masses may still be positive. This is an anomaly, since the definition of a mass function (3) requires that the mass of the null possibility be zero. This explains the role of (8), in which $m'(\emptyset)$ is deducted from the total mass and the remaining mass is divided by $(1 - m'(\emptyset))$ to ensure that the pooled mass will sum up to 1.

Dempster's rule is often compared to and contrasted with Bayes rule, because both rules concern the combination of probabilistic opinions into an aggregate (posterior) opinion. It is crucial to observe however that unlike Bayes rule, which is a trivial consequence of the axioms of probability theory, Dempster's rule is a *prescriptive* pooling mechanism which is neither right nor wrong, and thus it is less of a 'rule,' and more of a 'recipe.' Therefore, we take the position that the ultimate justification of Dempster's rule should be sought in the field, i.e., in the various applications in which the rule is supposed to have a certain sense of domain validity. This leads to the following question:

Question Q3: What is the intuitive justification of Dempster's rule in the context of information indexing and retrieval applications? If one wishes to aggregate indexing opinions via a certain pooling mechanism, then why use (7-8) and not another set of formulae?

A typical way to avoid this question is to invoke the argument: "If one uses mass functions

to represent relevance numbers, then one should combine them using Dempster's rule, because that's how mass functions are combined in the DS model." This argument could have been valid if Dempster's rule had a normative, domain-independent, and non-controversial justification. But this is not the case. In fact, many researchers have struggled to make sense of Dempster's rule, and the debate is still going strong: *"Shafer's theory has been strongly criticized for its failure to give a meaning to the measures of belief and plausibility, or to show how someone might arrive at a particular numerical assessment. In the absence of a definite interpretation, it is difficult to see how the rules of the theory, and in particular Dempster's rule, can be justified"* (Buxton, 1989). Given this controversy, the importance of question Q3 is obvious. Hence, our goal is to clarify, and to a certain extent defend, the *meaning* of Dempster's rule in the specific bibliographical context of an information indexing and retrieval application. This analysis is carried out in section 3, where a complete answer to Q3 is presented.

Belief Functions: Building on the elementary notion of a mass function $m : 2^\theta \rightarrow [0, 1]$, the function $\text{Bel} : 2^\theta \rightarrow [0, 1]$, denoted a *belief* function, can be defined as follows:

$$\text{Bel}(X) = \sum_{A \subseteq X} m(A) \quad (11)$$

Whereas $m(X)$ measures the support rendered to X (a subset of propositions) directly, $\text{Bel}(X)$ measures the total support rendered to X and to all its subsets (each being a more specific proposition). This relationship is depicted in figure 3, which illustrates how a $\text{Bel}(\cdot)$ function can be derived from the $m(\cdot)$ function given by (10). Note that (3) and (11) imply that $\text{Bel}(\emptyset) = 0$ and $\text{Bel}(\theta) = 1$ always. In fact, (11) implies that the Bel function

is completely determined by the mass function m , and, likewise, that m can be recovered from Bel's definition (Shafer, 1976, , p. 39).

Plausibility Functions: Whereas $\text{Bel}(X)$ measures the total support rendered to a possibility X , the *plausibility* of X , denoted $\text{Pl}(X)$, measures the maximal support that X can possibly attain under a given mass function m . Specifically:

$$\text{Pl}(X) = \sum_{A \cap X \neq \emptyset} m(A) \quad (12)$$

In words, $\text{Pl}(X)$ records the total mass allocated to all the possibilities with which X intersects. For a pictorial description of this relationship, refer again to figure 3.

Put figure 3 around here

The intuitive relationship between the three functions $m(\cdot)$, $\text{Bel}(\cdot)$, and $\text{Pl}(\cdot)$ can be described as follows. Beginning with the definition of Bel, consider the two possibilities $X, A \subseteq \theta$. Since both X and A are disjunctions of propositions, the set-theoretic statement $A \subseteq X$ is equivalent to the logical rule $A \rightarrow X$, which we will interpret as: 'If the truth lies in A , it must also lie in X .' Therefore, the sum of all the masses associated with premises A that imply X can be viewed as a measure of the total support rendered to X . As regards Pl's definition, suppose now that $A \cap X \neq \emptyset$ (but A is not necessarily a subset of X). Since the possibility A is a disjunction of propositions, the mass $m(A)$ rendered to it can 'float' freely to any one of its subsets, including those that intersect X . In the extreme case, the

intersection $A \cap X$ may inherit the *entire* mass of A . It follows that $Pl(X)$ is the upper bound of $Bel(X)$.

To do justice to the theory of evidence, it should be noted that the construction of Bel and Pl using m is only one way to define these functions. Shafer provided direct definitions of mass, belief and plausibility functions in terms of each other. He has also emphasized the key role that *subadditivity* plays in the theory of evidence, a point which we now turn to discuss in the specific context of information indexing and retrieval.

Sub Additivity: The *complement* of a set $X \subseteq \theta$, i.e. the set of all propositions that are in θ and not in X , is denoted hereafter \bar{X} . Definitions (11) and (12) imply the following important relationships:

$$Pl(X) = 1 - Bel(\bar{X}) \quad (13)$$

$$0 \leq Bel(X) \leq Pl(X) \leq 1 \quad (14)$$

If a certain Bel_b were a *Bayesian* representation of degrees of belief, the additivity axiom of probability theory ($X \cap Y = \emptyset$ implies $Bel_b(X \cup Y) = Bel_b(X) + Bel_b(Y)$) would mean that

$$Bel_b(X) = 1 - Bel_b(\bar{X}), \quad (15)$$

yet (13) and (14) imply that in the general case $\text{Bel}(X) \leq 1 - \text{Bel}(\overline{X})$, leading to the famous subadditivity property of the theory of evidence:

$$\text{Bel}(X) + \text{Bel}(\overline{X}) \leq 1 \tag{16}$$

In other words, the belief that one holds in a possibility does not automatically imply one's disbelief in the negation of that possibility. In information indexing and retrieval applications, where θ is taken to be a keyword lexicon \mathcal{K} , this tenant has important implications. For example, if the admittance of new evidence causes a cataloger to increase his belief in the document's relevance to a lexical subset X , the same evidence should not necessarily decrease his belief in the document's relevance to lexical subsets in \overline{X} , especially if the cataloger is not confident in his relevance opinion. In particular, the difference $1 - \text{Bel}(X) - \text{Bel}(\overline{X})$ is called the uncommitted belief with respect to X . If Bel were a Bayesian representation of degrees of belief, the uncommitted belief would be zero by definition. This is best illustrated in the 'state of insufficient reason,' in which one knows absolutely nothing about a set of propositions $\theta = \{q_1, \dots, q_n\}$. Whereas the common solution is to set $\text{Bel}(q) = 1/n$ for all $q_i \in \theta$, the theory of evidence would set $\text{Bel}(\theta) = 1$ and $\text{Bel}(X) = 0$ for all the other proper subsets of θ . This is the case when the uncommitted belief is at maximum.

The interpretation of $\text{Bel}(\cdot)$ and $\text{Pl}(\cdot)$ as lower and upper-probabilities has led many to view the theory of evidence as a novel calculus for eliciting and manipulating interval-valued, rather than point-valued, degrees of beliefs. Indeed, the theory allows one to express the belief in every hypothesis X by means of the interval $[\text{Bel}(X), \text{Pl}(X)]$, which

may be updated as new evidence about X is admitted. Further, the width of the interval, $Pl(X) - Bel(X)$, is by definition $1 - Bel(X) - Bel(\bar{X})$, or the uncommitted belief with respect to X . If the uncommitted beliefs induced by a certain mass function m were zero for all the hypotheses under consideration, the intervals would degenerate to zero widths and Bel would be a standard probability function. Yet in the more general case in which the mass reflects some 'second-order uncertainty,' or 'ambiguity,' the degree of belief in possibilities X drawn from θ is allowed to 'float' between $Bel(X)$ and $Pl(X)$. One benefit of such a model is that it is more robust and less prone to human errors in assessing subjective degrees of support.

We arrive at our last question:

Question Q4: The designer of a DS-based IR application can choose to elicit and represent relevance through three alternative languages: mass functions, belief functions, and belief intervals. What is the relationship among these three representation in the specific context of information indexing and retrieval applications?

Recall that the three functions m , Bel , and Pl , are mathematically equivalent, in the sense that knowledge of any one of them (for every possibility) can be used to compute the other two. This equivalence might lead one to concur that the question of whether to use m , Bel , or $[Bel, Pl]$ to elicit and manipulate degrees of support depends on cognitive and efficiency considerations. As it turns out, this conclusion is quite naïve. For example, belief intervals are not as flexible a representation as we would like them to be. That is, when one elicits $[Bel, Pl]$ intervals from a source of evidence, it is not true that the only restriction is that

$0 \leq \text{Bel} \leq \text{Pl} \leq 1$. Again, a full understanding of these constraints requires a semantic interpretation, which we now proceed to present.

4 A Canonical Indexing Model

As figure 1 illustrates, the key theme of this paper is the interplay of the theory and practice of the Dempster Shafer model, as viewed through the ‘lens’ of a particular application. The previous section was structured around the key constructs of the *theory*: the frame of discernment, mass and belief functions, and Dempster’s rule. Coming from the other extreme, this section is structured around the key constructs of the *application*: taxonomies, relevance functions, and index aggregation operators. This leads to the development of a *canonical indexing model*, around which the remainder of the paper evolves. In building this model, our intention is to articulate an indexing mechanism which is simple, intuitive, and, most importantly, probabilistic.

The main result that we are aiming at is this: notwithstanding its domain-specific origin and its strict probabilistic nature, the canonical model that is expounded here is completely isomorphic to the DS model. This has several important implications. First, the canonical model provides concrete answers to all the questions that were raised about the theoretical fit between the DS theory and information indexing and retrieval applications. Second, because the limitations of the former will be explicit, implicit limitations of the latter will become apparent. Third, because the canonical model makes no use of extra probabilistic arguments, it also provides a simple probabilistic interpretation to the DS theory, which is often claimed to be an extension of probability theory.

4.1 The Taxonomy

In most IR applications, documents are indexed and sought within a data structure that is called a *taxonomy*. The taxonomy is a finite set of *classes*, or categories, designed to organize documents in a particular subject of interest. For example, consider the following set of classes, taken from an art-related taxonomy: $C = \{\text{Art, Braque, Cubism, Dada, Impressionist, Janco, Modern, Picasso}\}$. Taxonomies are constructed by domain experts — in this case art scholars — who provide two types of information: (i) a set of classes; and (ii) a taxonomical data structure, expressed as ordered pairs of classes. Specifically, if we let (x, y) code the assertion ‘ y is a direct sub-class of x ’, then the expert might specify a relation like $H = \{(\text{Art, Modern}), (\text{Art, Impressionists}), (\text{Modern, Cubism}), (\text{Cubism, Braque}), (\text{Cubism, Picasso}), (\text{Dada, Picasso}), (\text{Dada, Janco})\}$, resulting with the taxonomy depicted in figure 2.

Put figure 2 around here

Formally, a taxonomy is a rooted directed acyclic graph $\langle C, H \rangle$. The nodes set C represents taxonomical classes, and the edges set H represents a relation on C (i.e., a subset of $C \times C$, giving directed pairs) with two restrictions: (i) no cycles exist in the digraph, and (ii) the digraph contains exactly one *root*, i.e. a class $r \in C$ such that no edge (x, r) exists in H . The descendants of a class x are the subclasses of x , and the predecessors are the generalization of x . The root of the taxonomy is the only class that (i) has no predecessors, and (ii) generalizes all other classes, e.g. the class art in figure 2. Since the taxonomy is a finite and acyclical graph, it contains a ‘boundary,’ or a set of

terminal classes. A class $k \in C$ is said to be terminal if has no descendants, i.e. if no edges of the form (k, x) exist in H . In the example of figure 2, the terminal classes are {Braque, Picasso, Janco}. We will denote the set of terminal classes by \mathcal{K} . This notation is not coincidental, as \mathcal{K} is precisely the keyword lexicon discussed in the previous section.

In a taxonomy, each class $c \in C$ is characterized by two sets that we denote $\text{LIB}(c)$ and $\text{VOL}(c)$ and call the *library rooted in c* and the *libraries that intersect c*, respectively. $\text{LIB}(c)$ contains all the classes that can be reached by paths beginning at c and following edges all the way ‘down’ to the terminal classes. This set of classes, which includes c itself, represents the entire set of classes into which c may be decomposed. Conversely, $\text{VOL}(c)$ contains all the classes that can be reached by following ‘upward’ paths beginning at c and ending at the taxonomy’s root. This set, which also includes c itself, represents all the classes to which c can be generalized. For example, $\text{LIB}(\text{Cubism}) = \{\text{Cubism}, \text{Braque}, \text{Picasso}\}$ and $\text{VOL}(\text{Cubism}) = \{\text{Cubism}, \text{Modern}, \text{Art}\}$. These sets can be given a recursive definition, as follows:

$$\text{LIB}(c) = \{c\} \cup \{x \in C \mid \exists y \in \text{LIB}(y) \text{ with } (y, x) \in H\} \quad (17)$$

$$\text{VOL}(c) = \{c\} \cup \{x \in C \mid \exists y \in \text{VOL}(y) \text{ with } (x, y) \in H\} \quad (18)$$

A taxonomy is similar to a hierarchical tree structure, with a difference. Unlike a common tree, each class in the taxonomy can have as many parents as we desire. For example, in figure 2 {Picasso} is a subclass of both {Cubism} and {Dada}. Also note in passing

that definitions (18-17) imply that (i) $\text{LIB}(\text{root}) = C$, so that the root library contains all the classes in the taxonomy; (ii) $\text{VOL}(\text{root}) = \{\text{root}\}$, the root library is not contained by any other library, and (iii) $\text{LIB}(k) = \{k\}$ if and only if $k \in \mathcal{K}$, so that terminal classes are characterized by libraries that contain singleton classes only.

The indexing process: The act of indexing a document within a taxonomy can be described as a top-down, depth-first search process. To illustrate, suppose that an art-related document is to be indexed within the art taxonomy from figure 2. Without loss of generality, assume that the document is relevant to modern art. Beginning at the first level under Modern and proceeding left to right, we test if the document is relevant to Cubism. If the answer is 'yes,' we step down one level and test if it's relevant to Braque. If the answer is 'yes,' we index the document in Braque. If the answer is either 'no' or 'unsure,' we test if it's relevant to Picasso. If the answer is either 'no' or 'unsure,' and assuming that Picasso is the last class below Cubism, we backtrack one level, index the document in Cubism, and proceed to explore Dada. If the document is deemed irrelevant to any one of the classes thus visited, we backtrack one level and index the document under Modern. This would reflect the notion that even though the document is related to modern art, the existing taxonomy fails to discern the exact category to which it belongs. Thus the indexing process involves a depth first search which is cut off at any class that is deemed to be irrelevant to the indexed document.

We see that the notion of relevance that is consistent with this process is defined over *subsets* of classes, not over *individual* classes in \mathcal{K} . That is, if a document is indexed under, say, Cubism, it implies that the document belongs to the library $\text{LIB}(\text{Cubism})$, i.e., to the collection of documents about Cubism, Braque, or Picasso. This definition of relevance is

convenient because it allows us to be as specific as we wish in our relevance statements. If we're sure that a document is relevant to a certain class, we index it under that class. If we're not sure, we can step back and index the document in a library that contains the class. We can do this all the way up to the root of the taxonomy, at which point the indexing decision root would express the opinion that the document belongs somewhere in the library, without specifying exactly where.

Relationship to the theory of evidence: The relationship between a taxonomy $\langle C, H \rangle$ and a lexical frame of discernment \mathcal{K} is simple, but not trivial. From a mathematical perspective, the taxonomy can be viewed as a subset of a graph G whose vertices are indexed by $2^{\mathcal{K}}$. In the graph G , there is an edge from the vertex indexed by the subset $A \in 2^{\mathcal{K}}$ to the vertex indexed by the subset $B \in 2^{\mathcal{K}}$ if and only if $A \supseteq B$ such that no other subset C satisfies $A \supseteq C \supseteq B$. In the taxonomy, which is a subset of G , each vertex c corresponds to the vertex of G that is indexed by the subset of \mathcal{K} obtained from the terminal elements in $\text{LIB}(c)$. Thus, each class in the taxonomy can be associated with an element in $2^{\mathcal{K}}$, namely the subset obtained from the keywords of the terminal classes that can be reached by looking 'downward' from the class in the taxonomy.

While there is this mathematical association, there are important differences between the notion of a taxonomy and the power set of \mathcal{K} as used in the DS model. First, we may distinguish between two types of taxonomies: static and adaptive. A *static taxonomy* consists of a fixed and unmodifiable set of classes, like the Dewey decimal system or the Library of Congress index. An *adaptive taxonomy* is a dynamic data structure that evolves from the indexing process itself. Such a taxonomy consists of a fixed set of keywords,

denoted \mathcal{K} , and an ‘open-ended’ set of *classes*, each class being a different grouping of keywords from \mathcal{K} . That is, when a new document is deemed relevant to a subset of keywords that don’t make up an existing category, we simply announce this subset a new class and add it to the taxonomy. Hence, a document titled “*A letter from Braque to Janco*” may well be indexed in the class {Braque, Janco}, something that would have been impossible in a static taxonomy that doesn’t contain such a predefined category. The only restriction that is placed on an adaptive taxonomy is that it must contain at least all the elements in \mathcal{K} (as singletons, or classes that are made up of single keywords), as well as \mathcal{K} itself. Hence, we begin with the initial set of classes $C = \{\{k_1\}, \dots, \{k_n\}\}, \mathcal{K}$, and add more classes to it as we go along.

Thus, the precise relationship between the IR notion of a taxonomy and the theoretical DS notion of a lexical frame of discernment \mathcal{K} can be described in two steps. First, any static taxonomy is conceptually a ‘frozen’ and ‘named’ version of some adaptive taxonomy. Second, any adaptive taxonomy, in turn, is a subset of the lexical power set $2^{\mathcal{K}}$. An example is illustrated in figure 4, using the simple lexicon $\mathcal{K} = \{\text{Braque, Picasso, Janco}\}$. Figure 4-a depicts the lexical power set $2^{\mathcal{K}}$ (excluding \emptyset). In practice, dealing with the power set of keywords is unrealistic, since the set of all possible classes becomes prohibitively large even with only a few dozen keywords. However, once the *semantics* of the lexicon is taken into consideration, many if not most of the classes in $2^{\mathcal{K}}$ become irrelevant, since they represent arbitrary grouping of keywords that can be excluded from the taxonomy for all practical purposes. If we choose to focus on *tree* taxonomies only, the power set can be restricted further by disregarding all its non-hierarchical subsets.⁵ Figure 4-b depicts a

⁵Using the notation $|X|$ to represent the cardinality of a set X , characterize each class $X \in C$ by the set $L(X) = \{Y \in C \mid |X| = |Y|\}$. A taxonomy $\langle C, H \rangle$ will be a *tree taxonomy* if and only if for every class $X \in C$, $L(X)$ contains only disjoint sets.

specific adaptive taxonomy that might have emerged from a hypothetical indexing process. By definition, this taxonomy is a subset of the exhaustive 4-a taxonomy. Finally, figure 4-c depicts a 'frozen' and 'named' version of the 4-b taxonomy. The naming procedure is domain-dependent: if certain classes 'make sense' on semantic grounds, they can be given descriptive names that reflect their contents. For example, the class {Braque, Picasso} can be named Cubism, the class {Braque, Picasso, Janco} Modern, etc.

Put figure 4 around here

We now turn to question Q1, which asked whether the DS concept of a lexical power set provides an adequate 'skeleton' for indexing documents in IR applications. The answer to this question is 'yes,' but there is a caveat. Note that there is a subtle difference between a bibliographical taxonomy and a subset of the DS power set: in the former, the classes have *names*; in the latter, the classes correspond to *anonymous* lexical subsets. That is, in the logical context of the DS model, to say that a document is relevant to $\{k_1, k_2\}$ is tantamount to saying that the document is relevant to either k_1 , or to k_2 . Yet in the context of a bibliographical taxonomy, most lexical subsets have meaningful names, like Cubism and Dada, just like the elementary keywords that make up their contents. Therefore, indexing a document in a named class might mean something quite different than the implication that the document should be indexed in one or more of the class's constituent keywords.

For example, suppose that a cataloger decided to index the title "*Cubist Landscapes*" directly in the class Cubism. In the standard DS model, this indexing opinion would imply that "the document is relevant either to Picasso, or to Braque, or to another Cubist artist." Although this interpretation is logically correct, it clearly entails a loss of concrete

information about the document's *direct* relevance to Cubism at large. Also, it leads to a situation in which the set of documents relevant to a class is *larger than or equal to* the union of the sets of documents relevant to all of its children classes, which is inconsistent with the disjunctive interpretation of a standard DS power set.

How can we augment the power set representation of the DS model so as not to force a cataloger to disregard information about a documents's direct relevance to non-singleton classes? By viewing the power set (or the portion of the power set that is in use for indexing) as a taxonomy $\langle C, H \rangle$, the problem may be solved by adding to the taxonomy a new set of *net classes*, as follows. For each non-terminal class $c \in C$, add (i) a new class named *net.c* to C , and (ii) a new link $(c, \text{net.c})$ to H . The new class *net.c*, which is a direct terminal descendant of c , can now serve as the index of the documents that are relevant specifically and directly to c . With this modification, each class c becomes a mere tag, or a pointer, and the proposition 'the document is relevant to the class c ' is once again equivalent to the proposition 'the document is relevant to the library rooted at c .'

Since the net classes are terminal classes, they become elements of the lexicon. Therefore, in a taxonomy which is augmented with a set of net classes, every indexing decision *can be interpreted* as selecting subsets of relevant keywords (which may include net classes) from the lexicon, so we are back in the familiar disjunctive stance of the DS frame of discernment. Purists may find this solution crude, but the adjustment is necessary if one wants to apply the DS model to information indexing and retrieval applications without violating, or misinterpreting, the set theoretic premise of the model.

4.2 Relevance Functions

The fundamental rule of indexing is that a document should be indexed using certain keywords if prospective users of the document would find it *relevant* to these keywords. In its most primitive form, then, relevance is a Boolean and subjective relation, indicating categorically that a document $d \in D$ is relevant to a lexical subset $X = \{k_1, \dots, k_m\}$ in the view of a particular library patron. However, due to the fact that bibliographical classes don't have crisp boundaries, and due to the multitude of relevance opinions expressed by different catalogers and library patrons, a more reasonable question is not whether d is relevant to X , but rather what is the *intensity* of this relation. In other words, we seek to represent relevance in terms of a mapping $r : 2^K \times D \rightarrow [0, 1]$, rather than in terms of a characteristic function $r : 2^K \times D \rightarrow \{0, 1\}$.

There have been many efforts to interpret relevance on probabilistic grounds, Maron and Kuhns (1960) being the defining article. One of the fundamental problems in this area has been the proper definition of the *sample space* from which relevance propositions are drawn. This point was alluded to by Maron, as follows:

“The notion of probability of relevance can be interpreted in two different perspectives: of the *document*, as the proportion of patrons of a given type who would judge that document relevant, and of the *patron* himself, as the proportion of documents of a given type which he would judge relevant. The first model leads to a theory of probabilistic indexing; The second model leads to a theory of probabilistic query formulation (Maron, 1982).”

In what follows we will focus on Maron's first perspective, in which multiple patrons form relevance opinions about a fixed document. Consistent with Maron's observation, this perspective yields a model of inexact indexing. Unlike Maron, though, the uncertainty associated with the indexes in our model will lead not to probability functions, but rather to Dempster Shafer mass functions, i.e. functions that conform to definition (3).

Let $U = \{u_1, \dots, u_n\}$ be a set of catalogers, and let \mathcal{K} be a keyword lexicon. Suppose that each cataloger in U is asked to index the same document using \mathcal{K} , i.e. to specify one or more keywords from \mathcal{K} that are relevant to the document. Suppose that cataloger u_i supplies the opinion that the document is relevant to the lexical subset $X \subseteq \mathcal{K}$; we then record this opinion by means of the following Boolean function:

$$v_i(X) = \begin{cases} 1 & \text{if } u_i \text{ indexed the document using } X \\ 0 & \text{otherwise} \end{cases} \quad (19)$$

$$i = 1, \dots, n$$

Since each cataloger u_i supplies one set of relevant keywords, there will be exactly one subset $X \in 2^{\mathcal{K}}$ such that $v_i(X) = 1$. Also, the empty set is not allowed to be a valid relevance opinion. If a cataloger is unwilling to give an opinion or is unsure about the proper classification of the document, the document is indexed by default in the root class \mathcal{K} , which is also an element of $2^{\mathcal{K}}$. This convention makes sense because the root class represents the *entire library*, and is therefore the natural place to store documents whose specific class membership is indiscernible.

After all n catalogers have cast their indexing opinions regarding *the same document* d , we

compute for each lexical subset $X \in 2^{\mathcal{K}}$ three ‘relevance counters,’ as follows:

$$r(X) = \sum_{i=1}^n v_i(X, d) \quad (20)$$

$$r_{\text{LIB}}(X) = \sum_{Y \in \text{LIB}(X)} r(Y, d) \quad (21)$$

$$r_{\text{VOL}}(X) = \sum_{Y \in \text{VOL}(X)} r(Y, d) \quad (22)$$

In words, $r(X)$, $r_{\text{LIB}}(X)$, and $r_{\text{VOL}}(X)$ count the number of catalogers who classified the document in X , in the library rooted in X , and in libraries that intersect (or in a hierarchical taxonomy, *contain*) X , respectively. (When d is fixed in our analysis, we will suppress the explicit dependence, and write $r(X)$ instead of $r(X, d)$.)

Relationship to the theory of evidence: Suppose now that the Boolean relevance opinions of the catalogers are averaged over the space of catalogers U through the following computation:

$$m(X) = \frac{1}{n} \cdot r(X) = \frac{1}{n} \sum_{i=1}^n v_i(X) \quad (23)$$

The resulting function $m(X)$ is a DS mass function over the lexical space \mathcal{K} . Formally, we have the following proposition (the proofs are given in a separate appendix):

Proposition 1: Let $U = \{u_1, \dots, u_n\}$ be a set of catalogers with their Boolean relevance opinions $v_1, \dots, v_n : 2^K \rightarrow \{0, 1\}$. The real function $m : 2^K \rightarrow [0, 1]$ defined by $m(X) = \frac{1}{n} \cdot r(X) = \frac{1}{n} \sum_{i=1}^n v_i(X)$ is a mass function, satisfying definition (3).

The consequence of the proposition is that DS mass functions arise *naturally* when we view the relevance functions as derived from averages of multiple Boolean indexing opinions. We begin with a space U of n catalogers who are asked to index the same document using the same lexicon \mathcal{K} . Each cataloger supplies an individual opinion that specifies which keywords are relevant to the document. Note that the cataloger's indexes are not restricted, and that they are free to choose any keyword or combination of keywords that, in their opinion, are relevant to the document. Next, shifting our attention from the catalogers space U to the lexical space \mathcal{K} , we compute for each lexical subset $X \subseteq \mathcal{K}$ a measure of 'average relevance,' $\frac{1}{n} \cdot r(X)$, which represents the fraction of catalogers who thought that the document was relevant to X . Disregarding the lexical subsets that no cataloger has chosen, we obtain a set of pairs of the form $\{(K_1, r_1), \dots, (K_n, r_n)\}$ in which $K_i \in 2^K$ and $0 \leq r_i = \frac{1}{n} \cdot r(K_i) \leq 1$.

We are now in a position to answer question Q2, regarding the 'type' of relevance that DS mass functions represent, given the context of multiple relevance opinions. First, the canonical model has yielded the type of *relevance numbers* that are at the center of any probabilistic indexing model. Second, according to Proposition 1, these numbers form a mass function, consistent with the standard DS model. Finally, the meaning of the mass $m(X)$ is simply the fraction of catalogers who thought that the document was relevant to the set of keywords X .

Following the same line of reasoning, we can also provide an answer to question Q4, that

sought an IR interpretation of the meaning and relationship of mass functions, belief functions, and belief intervals. Given the IR context in which $\theta \equiv \mathcal{K}$, it is easily seen that the relevance counters (20-22) are proportional to the mappings that represent degrees of belief in the DS model. Specifically, dividing each counter by n — the number of catalogers — yields the mass, belief, and plausibility, functions defined in (3), (11), and (12), respectively:

$$m(X) = \frac{1}{n} \cdot r(X) \quad (24)$$

$$\text{Bel}(X) = \frac{1}{n} \cdot r_{\text{LIB}}(X) \quad (25)$$

$$\text{Pl}(X) = \frac{1}{n} \cdot r_{\text{VOL}}(X) \quad (26)$$

If we combine these observations with the interpretation of the power set of the lexicon as a taxonomy, we see that the mass on a lexical subset X is given by the fraction of catalogers who indexed the document using X directly. Similarly, the belief in X is the fraction of catalogers who indexed the document in libraries within X , and the plausibility of X is the fraction of catalogers who indexed the document in libraries that intersect (in a hierarchical taxonomy, *contain*) X .

The key component of the canonical model that enables this interpretation of the DS functions is the assumption of multiple patrons and the $v_i(\cdot)$ functions that keep track of their individual indexing opinions. In the canonical model, the assumption of multiple patrons is explicit and is the foundation on which the entire analysis rests. In the DS model, the assumption of multiple Boolean opinions and their respective $v_i(\cdot)$ functions are implicit.

4.3 Aggregating Relevance

So far, we have assumed that (i) relevance is a two-place function $r(X, d)$ between a document d and a lexical subset X , and that (ii) all the catalogers from whom $r(X, d)$ was elicited were of the same ‘type,’ using Maron’s terminology (see quote in Section 4.2). In this section we retract both assumptions. Specifically, we argue that relevance, in its most elementary form, is a three-place relation $r(X, d, q)$ in which q is the *classifier* dimension, or *context*, in which d is judged to be relevant to X . With that in mind, $r(X, d)$ can be viewed as a measure of *aggregate relevance* that runs over all the possible contexts in which d ’s relevance to X is judged. We now turn to describe a pooling mechanism that implements such an aggregation.

Let $U_1 = \{u_1, \dots, u_{n_1}\}$ be a group of n_1 catalogers who are asked to index a document d using a keyword lexicon \mathcal{K} based on a certain classifier, denoted q_1 . Similarly, let $U_2 = \{u'_1, \dots, u'_{n_2}\}$ be a group of n_2 catalogers who are asked to index the same document, based on another classifier, denoted q_2 . The semantics of the classifiers depends on the indexing scenario. For example, q_1 might be the document’s title, whereas q_2 might be the document’s abstract. Alternatively, in a dynamic model in which the relevance indexes of documents are continuously revised to reflect actual use, q_1 and q_2 can represent two different information needs, or queries, in the context of which the relevance of d to X was judged, either explicitly or through an automatic keywords extraction algorithm.

For example, let $\mathcal{K} = \{A, B, C\}$ and let U_1 and U_2 consist of 4 and 3 catalogers, respectively. Assume that within the U_1 group, two catalogers index the document in $\{A, B\}$, one in $\{A\}$, and one in $\{B\}$. Within the U_2 group, one cataloger indexes the document in $\{B, C\}$,

one in $\{A, B\}$, and one in $\{B\}$. These indexing opinions are tabulated in the two tables on the left side of figure 5. The columns of each table represent the common lexicon $\mathcal{K} = \{A, B, C\}$. The i th tuple in each table represents the relevance opinion elicited from the i th cataloger in the respective group as a binary vector. To be precise, 1 in the (i, j) th table entry indicates that cataloger i has included the j th keyword in his indexing opinion and 0 indicates that he didn't.

Put figure 5 around here

In what follows, we denote the binary vector that represents the relevance opinion of cataloger u_i by w_i . Similarly, the set of all relevance opinions of a group of catalogers will be denoted $W = \{w_i | u_i \in U\}$. Finally, the group of catalogers U together with their relevance opinions W will be denoted $T = \langle U, W \rangle$ and referred to as a *model*. With this notation, consider two groups of catalogers U_1 and U_2 together with their relevance opinions W_1 and W_2 . If all the catalogers in both groups are considered equally qualified to cast relevance opinions, then a variety of different pooling mechanisms may be used to compute the aggregate index induced by *all* the catalogers. Symbolically, we seek an operator \otimes to compute the model $\langle U, W \rangle = \langle U_1, W_1 \rangle \otimes \langle U_2, W_2 \rangle$.

The pooling mechanism depicted in figure 5, denoted hereafter by \otimes , implements an operator that was described by Hummel and Landy (1988) as "a consensus opinion formed by the committees of two." Here, the set of all possible committees is $U = U_1 \times U_2$, consisting of all the $n_1 \cdot n_2$ unique pairs of catalogers that can be drawn from U_1 and from U_2 . The combined relevance opinion associated with the pair $(u_i, u'_j) \in U$ is defined to be the *binary conjunction* of the individual opinions of $u_i \in U_1$ and $u'_j \in U_2$, which we denote

$w_{i,j} = w_i \cdot w'_j$. For example, consider the first tuple in the U table in figure 5. This tuple gives the opinion of the committee (u_1, u'_1) , i.e. $w_{1,1'} = (0, 1, 0)$. This opinion is the binary conjunction of the individual opinion $w_1 = (1, 1, 0)$ and $w'_1 = (0, 1, 1)$ as given by catalogers u_1 and u'_1 respectively.

The pooling operation \otimes is completed by treating U as a new group of catalogers and using (23) to compute the mass function that it induces:

$$\begin{aligned}
 m'(A) &= m'(1, 0, 0) = 1/12 \\
 m'(B) &= m'(0, 1, 0) = 7/12 \\
 m'(A, B) &= m'(1, 1, 0) = 2/12 \\
 m'(\emptyset) &= m'(0, 0, 0) = 2/12
 \end{aligned} \tag{27}$$

Note that m' is not necessarily a mass function, since \otimes can yield a result like $m'(\emptyset) > 0$. This happens when there is a pair of opinions (e.g. u_2 and u'_1 in our example), such that the conjunction of the opinions gives the empty set even though neither opinion gives the empty set individually. To resolve the problem, we normalize $m'(\cdot)$ as follows:

$$\begin{aligned}
 m(X) &= \frac{1}{1-m'(\emptyset)} \cdot m'(X) = 1/10 \\
 m(B) &= \frac{1}{1-m'(\emptyset)} \cdot m'(B) = 7/10 \\
 m(A, B) &= \frac{1}{1-m'(\emptyset)} \cdot m'(AB) = 2/10 \\
 m(\emptyset) &\stackrel{\text{def}}{=} 0
 \end{aligned} \tag{28}$$

In words, for each lexical subset $X \in \mathcal{K}$, $m'(X)$ is the fraction of the (paired) catalogers who classified the document in that subset. Next, the fraction of the catalogers who agreed on *nothing* – $m'(0, 0, 0)$ – is distributed evenly among the fractions of catalogers who agreed on *something*, yielding a new mass that sums up to unity. This function is now taken to

be the ‘aggregate index’ of the document d , implying the taxonomy depicted at the top right of the figure. We may also view $m(X)$ as the fraction of (paired) catalogers who index the document in X among those paired catalogers who do not index the document in the empty set \emptyset . That is, if we discard pairs that agree on no relevant keywords, then the remaining pairs can compute their pooled relevance and then yield a mass function m .

Relationship to the theory of evidence: In order to explore the relationship of the multiple catalogers/multiple classifiers scenario to the DS model, we first have to step back and say a few words about the role of ‘sources of evidence’ in the latter. Basically, the DS theory models a situation in which a finite set of ‘pieces’ or ‘sources’ of evidence $E = \{e_1, \dots, e_n\}$ is used to discern the likelihoods of various possibilities X drawn from a common frame of discernment. Yet the *identity* of the sources of evidence is rather implicit in the model’s language. That is, the common notation $m_i(X)$ and $\text{Bel}_i(X)$ is meant to be shorthand of the mass and belief functions $m(X|e_i)$ and $\text{Bel}(X|e_i)$, where e_i is the source of evidence whose ‘support’ of the possibility X we are trying to capture. The total support that the body of evidence E lends to X is computed through Dempster’s rule (7-8), which yields a new function of the form $m(X|e_1, \dots, e_n) = m(X|e_1) \oplus, \dots, \oplus m(X|e_n)$.⁶ For simplicity’s sake, we denote the latter function $m(X)$, which reads ‘the mass that the possibility X attains after all the available evidence has been taken into consideration.’

With that, the relationship between the canonical model and the DS model is as follows: possibilities correspond to lexical subsets, and sources of evidence correspond to classifiers, i.e. to different aspects of the document (title, abstract, author, etc.) that help discern the

⁶Dempster’s rule (7-8) is commutative and associative, so its extension from 2 to n operands is straightforward.

document's proper classification. The missing piece in the analogy is the set of catalogers who inspect each classifier individually and cast Boolean relevance opinions based on that information. In the DS model, these catalogers are implicit. In the canonical model, they are the driving force of the entire analysis. Another way to interpret the group of catalogers is to view them as a group of library patrons who approach the same document with different information needs (or queries) in mind, each corresponding to a piece of evidence that highlights one facet of the composite relation that we call 'relevance.'

How should we combine this multitude of relevance opinions into an aggregate index? In the canonical indexing model, the opinions are combined at the catalogers level, through the cartesian consensus operator \otimes . In the DS model, where the catalogers space is implicit, the opinions are combined at the classifiers level, via Dempster's rule \oplus . The key point, as illustrated in figure 5, is that both combination methods lead to precisely the same result. Formally, we have the following proposition:

Proposition 2: Let $T_1 = \langle U_1, W_1 \rangle$ and $T_2 = \langle U_2, W_2 \rangle$ be two sets of catalogers together with their Boolean relevance opinions, and let $T = \langle U, W \rangle$ be the outcome of $T = T_1 \otimes T_2$, as follows: (i) $U = U_1 \times U_2$; and (ii) $W = \{w_{i,j} = w_i \cdot w'_j | w_i \in W_1 \text{ and } w'_j \in W_2\}$. Let \oplus be Dempster's rule as it is applied to mass functions. Let m_{T_1} , m_{T_2} , and $m_{T_1 \otimes T_2}$, be the mass functions induced by the models T_1 , T_2 , and $T_1 \otimes T_2$. Then we have the following: $m_{T_1 \otimes T_2} = m_{T_1} \oplus m_{T_2}$.

Proposition 2 serves to shed light on the prescriptive nature of Dempster's rule. That is, once we accept the fact that Dempster's rule \oplus is isomorphic to the cartesian product operator \otimes , a whole set of questions emerges: (1) why are the individual catalogers forced

to specify only *Boolean*, and not probabilistic, relevance opinions? (2) why are the groups of catalogers joined using a *set product* operator, as opposed to other set combination operators, e.g. union? (3) why committees of *two*, and not, say, comatous of three? (4) why are the individual relevance opinions combined using a binary *conjunction* rule? (5) why are all cataloger opinions given the *same weight*, where in practice some opinions may be more informed or worthy than others?

A proper answer to these questions requires an elaborate research program, involving both theoretical and empirical work. Also, the exact nature of the combination rule can vary from one situation to another. In the specific context of information indexing and retrieval, one can think of a family of indexing models, designed to operate under different sets of assumptions. For example, if the catalogers prefer to express binary relevance opinions, we can use Dempster's rule (or the equivalent \otimes) to combine them. If they wish to express relevance by selecting a number between 0 and 1, we can modify the combination rule to accommodate this language as well (this will be similar to the way Yen (1989) extended Dempster's rule in the GERTIS system). If the catalogers wish to use a discrete language such as 'remotely relevant,' 'somewhat relevant,' etc., we can develop a fuzzy version of the rule. The key point here is that the precise definition of \otimes , along with Proposition 2, provide clear guidelines as to (i) which aspect of the combination rule has to be modified, and (ii) what will be the normative relationship between the modified rule, Dempster's rule, and probability theory.

5 Conclusion

All the major implications of the research were already discussed in the body of the paper. We conclude with several comments regarding (i) efforts to apply the DS model to information indexing and retrieval applications; and (ii) efforts to interpret the theory of evidence on logical or probabilistic grounds.

Information Indexing and Retrieval: One objective of the paper was to articulate a concrete relationship between the Dempster Shafer model and information indexing and retrieval applications. The relationship that we expounded can be summarized as follows:

IR application	Dempster-Shafer model
keyword lexicon (\mathcal{K})	frame of discernment (θ)
taxonomy ($\langle C, H \rangle$)	named subset of 2^θ
classification criteria (q_i)	sources of evidence (e_i)
groups of catalogers (U_j)	implicit
individual indexing opinions (W_i)	implicit
relevance measure to class (r)	mass function (m)
relevance measure to library ($r_{\text{LIB}}(X)$)	belief function (Bel)
relevance measure to volume ($r_{\text{VOL}}(X)$)	plausibility function (Pl)
relevance aggregation operator (\otimes)	Dempster's rule (\oplus)

We hope that the details of this 'mapping,' as discussed in the paper, will promote a better understanding of the proper way to apply the DS model to IR applications. In addition,

the mapping provides a practical foundation for building a variety of different indexing algorithms. These algorithms can use the \otimes combination rule, or versions thereof, as called by the application. Ultimately, the success of one relevance calculus or another will depend on face validity and on field performance considerations.

The Dempster Shafer theory of evidence: Several authors provided canonical examples that explain the rationale of the DS model in the way of analogy. Zadeh (1986) illustrated how mass functions and Dempster's rule can be mapped on fuzzy queries about *interval-valued*, rather than point-valued, attributes, in a relational database. Gordon and Shortliffe (1985) gave a compelling interpretation of how a DS calculus can be used to represent and combine the degrees of belief that clinical symptoms (pieces of evidence) render to classes of bacterial organisms (disjunctions of hypotheses), whose set relationships forms a hierarchy. Coming from a different, domain-independent, direction, Hummel and Landy (1988) analyzed the probabilistic foundation of the theory of evidence *in general*, without making any assumptions on the underlying domain or the logical structure of the hypotheses. In contrast to other researchers who attempted to interpret high-level constructs of the DS model *directly* (e.g. Baron, 1987, Kyburg, 1987, and Schocken and Kleindorfer, 1989), Hummel and Landy took a more fundamental viewpoint that showed how the theory's constructs were implicitly linked to statistics of the opinions of hypothetical experts. However, their abstract mathematical analysis made no use of canonical examples, and it is difficult to interpret its implications on practical domains of application.

With that in mind, one objective of this paper was to illustrate how constructs of the DS model that up until now defied simplistic interpretations yield to a plausible interpretation in the practical context of a multi-classifier/multi-cataloger model. We have seen, in

propositions 1 and 2, that the canonical model leads to exactly the same set of functions and formulae of the DS model. Hence, from a mathematical perspective, the canonical model is isomorphic to the DS model. Yet from a semantic perspective, it invokes the notion of multiple catalogers. Although the notion of multiple patrons appears in several major interpretations of bibliographical relevance (Maron and Kuhns, 1960 , Maron, 1982), it may or may not exist in other applications.

To what extent, then, are we forced to accept the canonical interpretation of multiple catalogers *in principle*? One can simply reject the notion, avoiding the isomorphism by denying the possibility of multiple opinions, and relying simply on the DS theory as presented in Section 2. In that case, however, one is left with philosophical questions like Q1 through Q4. There could, of course, be other interpretations. However, in a real sense, *all* valid interpretations must be accepted or explained. That is, either the interpretation is accepted as is, or one must show how another set of semantic constructs provides a plausible interpretation of the theory. One advantage of our approach is that new calculi can be developed, different from the DS combination rule, that might better suit particular applications, based on modifications of the canonical model. It is precisely the unsatisfactory elements of this canonical model that permit us to systematically seek improved methods for managing uncertainty.

Since our analysis was strictly probabilistic, it seems to be consistent with Lindley's observation that "*Anything that can be done with belief functions can better be done with probability theory*" (Lindley, 1987, , p. 38). However, we believe that this argument misses an important point. To use a crude but useful analogy, it will be unreasonable to write off a programming language like Pascal simply because every Pascal program can be rewritten

in machine language. Just like high-level languages provide complex structures for dealing with specialized problems, the DS model provides non-elementary functions and operators that lend themselves nicely to certain domains of application, information indexing and retrieval being one such example.

We conclude that the Dempster Shafer theory of evidence provides an attractive framework for supporting information indexing and retrieval applications, and that these applications, in turn, serve to highlight the internal validity and limitations of the theory. Dempster's rule remains a controversial operator for combining degrees of beliefs, but this paper has illustrated that it is just one member in a family of parametric combination rules, and that the question of whether to use this rule or another is more a matter of reasoned choice than a matter of adhering to a fixed set of formulae.

Appendix: Proofs

Proposition 1: Let $U = \{u_1, \dots, u_n\}$ be a set of catalogers with their Boolean relevance opinions $v_1, \dots, v_n : 2^K \rightarrow \{0, 1\}$. The real function $m : 2^K \rightarrow [0, 1]$ defined by $m(X) = \frac{1}{n} \cdot r(X) = \frac{1}{n} \sum_{i=1}^n v_i(X)$ is a mass function, satisfying definition (3).

Proof: For each class $X \in 2^K$, either *all*, *some*, or *none* of the catalogers indexed the document in X . Hence, $r(X) = n$, or $r(X) < n$, or $r(X) = 0$, respectively, implying that $0 \leq m(X) \leq 1$. Hence, $m(\cdot)$ is a mapping from 2^K to $[0, 1]$, satisfying the first requirement of being a mass function. The second requirement is that the function will sum up to 1 over all the subsets of \mathcal{K} . This is proved as follows. For each cataloger u_i , exactly one

of the subsets $X \subseteq \mathcal{K}$ is such that $v_i(X) = 1$. For all other subsets Y , $v_i(Y) = 0$. Thus $\sum_{X \in 2^{\mathcal{K}}} v_i(X) = 1$. We thus have the following:

$$\begin{aligned} \sum_{X \in 2^{\mathcal{K}}} m(X) &= \sum_{X \in 2^{\mathcal{K}}} \frac{1}{n} \sum_{i=1}^n v_i(X) \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{X \in 2^{\mathcal{K}}} v_i(X) \\ &= \frac{1}{n} \sum_{i=1}^n 1 = 1 \end{aligned}$$

Further, since no cataloger gives \emptyset as his opinion, it is always true that $v_i(\emptyset) = 0$. Therefore, the third requirement of definition (3) is satisfied. Thus m is a mass function.

Definition of the \otimes combination rule: Let $U = \{u_1, \dots, u_n\}$ be a set of catalogers with their Boolean relevance opinions $v_1, \dots, v_n : 2^{\mathcal{K}} \rightarrow \{0, 1\}$. To denote the fact that the keyword $k \in \mathcal{K}$ was included in the indexing opinion of the i th cataloger, we use the following notation:

$$w_i(k) = \begin{cases} 1 & \text{if } v_i(X) = 1 \text{ and } k \in X \\ 0 & \text{otherwise} \end{cases} \quad (29)$$

If $\mathcal{K} = \{k_1, \dots, k_n\}$, the binary vector obtained by $w_i(k_1), \dots, w_i(k_n)$ is denoted w_i and called the Boolean relevance opinion of u_i . The collection of all such opinions of members of U is denoted $W = \{w_i | u_i \in U\}$. To combine the relevance opinions of two sets of catalogers $\langle U_1, W_1 \rangle$ and $\langle U_2, W_2 \rangle$, we use the following formulae (\otimes):

$$U = U_1 \times U_2, \quad (30)$$

$$W = \{w_{i,j}(\cdot) | u_i \in U_1, u_j \in U_2\}, \quad (31)$$

$$w_{i,j}(k) = w_i(k) \cdot w'_j(k). \quad (32)$$

Where $w_i(k)$ and $w'_j(k)$ are as defined in (29) for $u_i \in U_1$ and of $u'_j \in U_2$.

In section 4.2 we have shown how a mass function can be constructed from a set of catalogers (Proposition 1). Specifically, recall that the mass function induced by the model $T = \langle U, W \rangle$, denoted hereafter $m_T(X)$, gives the fraction of catalogers in U , among those catalogers who express an opinion (i.e. $w_i \neq \vec{0}$), whose relevance opinion exactly matched X . This is the same as those catalogers for whom $w_i(k_j) = 1$ if and only if $k_j \in X$. For $T = \langle U, W \rangle$, This fraction can be written down exactly:

$$m_T(X) = \frac{\#\{u_i \in U | w_i(k_j) = 1 \text{ if } k_j \in X \text{ and } w_i(k_m) = 0 \text{ if } k_m \notin X\}}{\#\{u_i \in U | w_i \neq \vec{0}\}}, \quad (33)$$

for $X \neq \emptyset$. Of course, $m_T(\emptyset) = 0$. We are now in a position to prove the following.

Proposition 2: Let $T_1 = \langle U_1, W_1 \rangle$ and $T_2 = \langle U_2, W_2 \rangle$ be two sets of catalogers together with their Boolean relevance opinions, and let $T = \langle U, W \rangle$ be the outcome of $T = T_1 \otimes T_2$, as follows: (i) $U = U_1 \times U_2$; and (ii) $W = \{w_{i,j} = w_i \cdot w'_j | w_i \in W_1 \text{ and } w'_j \in W_2\}$. Let \oplus be Dempster's rule as it is applied to mass functions. Let m_{T_1} , m_{T_2} , and $m_{T_1 \otimes T_2}$, be the mass functions induced by the models T_1 , T_2 , and $T_1 \otimes T_2$. Then we have the following: $m_{T_1 \otimes T_2} = m_{T_1} \oplus m_{T_2}$.

Proof: This proposition asserts a relationship between the general Dempster Shafer model and the canonical indexing model presented in section 4. The fact that the mapping from one model to the other is homomorphic follows from Hummel and Landy (1988) , but we will supply an independent argument here in the context of the indexing model.

Let us assume that there are n_1 catalogers in U_1 and n_2 catalogers in U_2 , and let us fix a particular nonempty lexical subset X of the lexicon \mathcal{K} . We wish to show that

$$m_{T_1 \oplus T_2}(X) = (m_{T_1} \oplus m_{T_2})(X) \quad (34)$$

Beginning with the right hand side of (34) and using the definition of Dempster's rule \oplus , $(m_{T_1} \oplus m_{T_2})(X)$ is equivalent to

$$\frac{\sum_{A \cap B = X} m_{T_1}(A) \cdot m_{T_2}(B)}{\sum_{A \cap B \neq \emptyset} m_{T_1}(A) \cdot m_{T_2}(B)}. \quad (35)$$

Multiplying top and bottom by $n_1 \cdot n_2$ and distributing, we obtain

$$\frac{\sum_{A \cap B = X} n_1 m_{T_1}(A) \cdot n_2 m_{T_2}(B)}{\sum_{A \cap B \neq \emptyset} n_1 m_{T_1}(A) \cdot n_2 m_{T_2}(B)}. \quad (36)$$

Recalling how mass functions are induced from the opinions of groups of catalogers (Eqn. 23 in Section 4.2), we may interpret this expression as follows. The value $n_1 m_{T_1}(A)$ counts the number of catalogers in U_1 who have indexed the document in the lexical subset A . Likewise, $n_2 m_{T_2}(B)$ counts the number of catalogers in T_2 who have indexed the document

in the lexical subset B . Hence, the product $n_1 m_{T_1}(A) \cdot n_2 m_{T_2}(B)$ counts the number of distinct pairs of catalogers $(u_i, u_{j'})$ in $U_1 \times U_2$ where $u_i \in U_1$ has indexed the document in A and $u_{j'} \in U_2$ has indexed the document in B . Now, according to the way \otimes is defined, if u_i has indexed in A and $u_{j'}$ has indexed in B , then the pair of catalogers $(u_i, u_{j'})$ end up indexing the document in $A \cap B = X$. Thus, the numerator of expression (36) counts *all* the cataloger pairs that end up indexing the document in X .

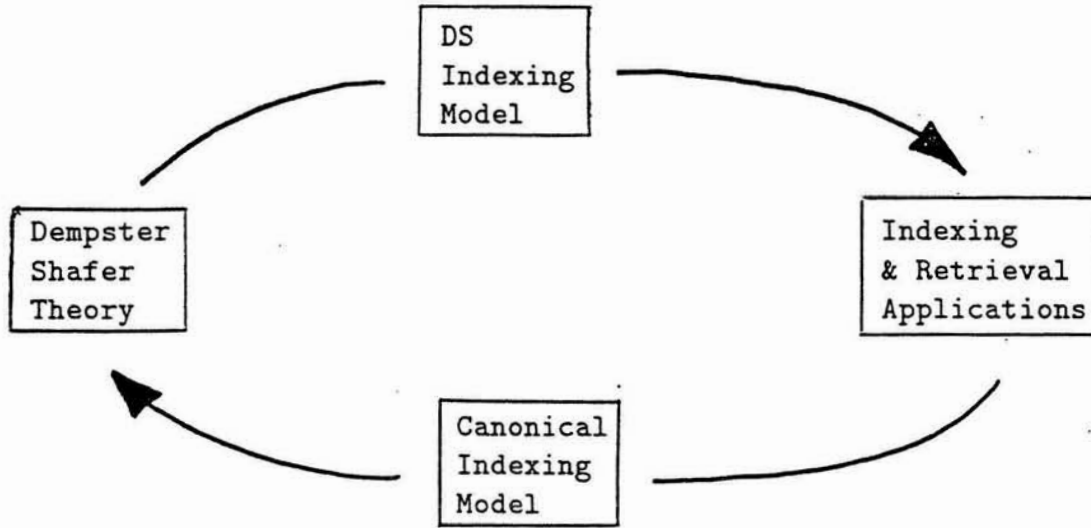
Precisely the same argument can be used to show that the denominator of (36) counts all the pairs of catalogers who don't index the document on \emptyset . Thus (36) gives the fraction of cataloger pairs in $U_1 \times U_2$ that have indexed the document in X out of the pairs of catalogers in $U_1 \times U_2$ who have indexed the document in some non-empty lexical subset, which is exactly the definition of $m_{T_1 \otimes T_2}$, the left hand side of (34).

References

- [1] J. Baron. (1987). Second-order probabilities and belief functions. *Theory and Decision*, 22.
- [2] G. Biswas, J.C. Bezdek, M. Marques, and V. Subramanian. (1987). Knowledge assisted document retrieval (I and II). *Journal of the American Society for Information Science*, 38(2):83-110.
- [3] R. Buxton. (1989). Modeling uncertainty in expert systems. *International Journal Man-Machine Studies*, 31:415-476.
- [4] A.P. Dempster. (1967). Upper and lower probabilities induced by a multivalued mapping. *Annals Mathematics Statistics*, 38:325-339.
- [5] A.P. Dempster. (1967). Upper and lower probability inferences based on a sample from a finite univariate population. *Biometrika*, 54:515-528.
- [6] J. Gordon and E.H. Shortliffe. (1985). A method for managing evidential reasoning in a hierarchical hypothesis space. *Artificial Intelligence*, 26:323-357.
- [7] R.A. Hummel and M.S. Landy. (1988). A statistical viewpoint on the theory of evidence. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 10(2):235-247.
- [8] H.E. Kyburg. (1987). Bayesian and non-Bayesian evidential updating. *Artificial Intelligence*, 31:271-293.
- [9] D.V. Lindley. (1987). The probability approach to the treatment of uncertainty in artificial intelligence and expert systems. *Statistical Science*, 2(1):17-24.
- [10] M.E. Maron. (1982). Associative search techniques versus probabilistic retrieval models. *Journal of the American Society for Information Science*, 308-310.
- [11] M.E. Maron and J.L. Kuhns. (1960). On relevance, probabilistic indexing and information retrieval. *Journal of the Association for Computing Machinery*, 7(3):216-244.
- [12] G. Salton and C. Buckley. (1988). Term weighing approaches in automatic text retrieval. *Information Processing Management*, 24(5):513-523.

- [13] G. Salton and M.J. McGill. (1983). *Introduction to Modern Information Retrieval*. McGraw-Hill, New York.
- [14] T. Saracevic. (1975). Relevance: a review of and a framework for thinking on the notion in information science. *Journal of the American Society for Information Science*, Nov-Dec:321-342.
- [15] S. Schocken and P.R. Kleindorfer. (1989). Artificial intelligence dialects of the Bayesian belief revision language. *IEEE Transactions on Systems, Man, and Cybernetics*, 19:1106-1121.
- [16] G. Shafer. (1976). *A Mathematical Theory of Evidence*. Princeton University Press.
- [17] G. Shafer. (1987). Probability judgement in artificial intelligence and expert systems. *Statistical Science*, 2(1):3-44.
- [18] R.M. Tong and D.G. Shapiro. (1985). Experimental investigations of uncertainty in a rule-based system for information retrieval. *International Journal of Man-Machine Studies*, 22:265-282.
- [19] H. Turtle and W.B. Croft. (1991). Evaluation of an inference network-based retrieval model. *ACM Transactions on Information Systems*, 9(3):187-222.
- [20] J. Yen. (1989). Gertis: a Dempster-Shafer approach to diagnosing hierarchical hypotheses. *Communications of the ACM*, 32(5):573-585.
- [21] L.A. Zadeh. (1986). A simple view of the dempster-shafer theory of evidence and its implication on the rule of combination. *The AI Magazine*, 85-90.

EXTERNAL VALIDITY



INTERNAL VALIDITY

Figure 1: A pictorial description of the paper's methodology. Section 3 uses the terminology and rationale of the Dempster Shafer theory to derive a DS indexing model for IR applications (top arrow). Taking the opposite direction, Section 4 builds a canonical indexing model that is based on the domain specific requirements of IR applications. As it turns out, the canonical model provides a probabilistic and domain-independent interpretation of the Dempster Shafer theory of evidence.

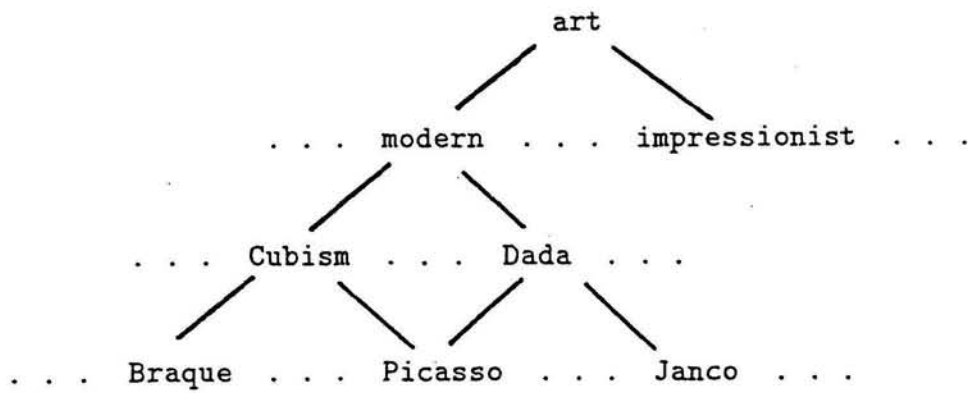


Figure 2: An excerpt from an art-related taxonomy designed to classify documents on major artists and artistic movements.

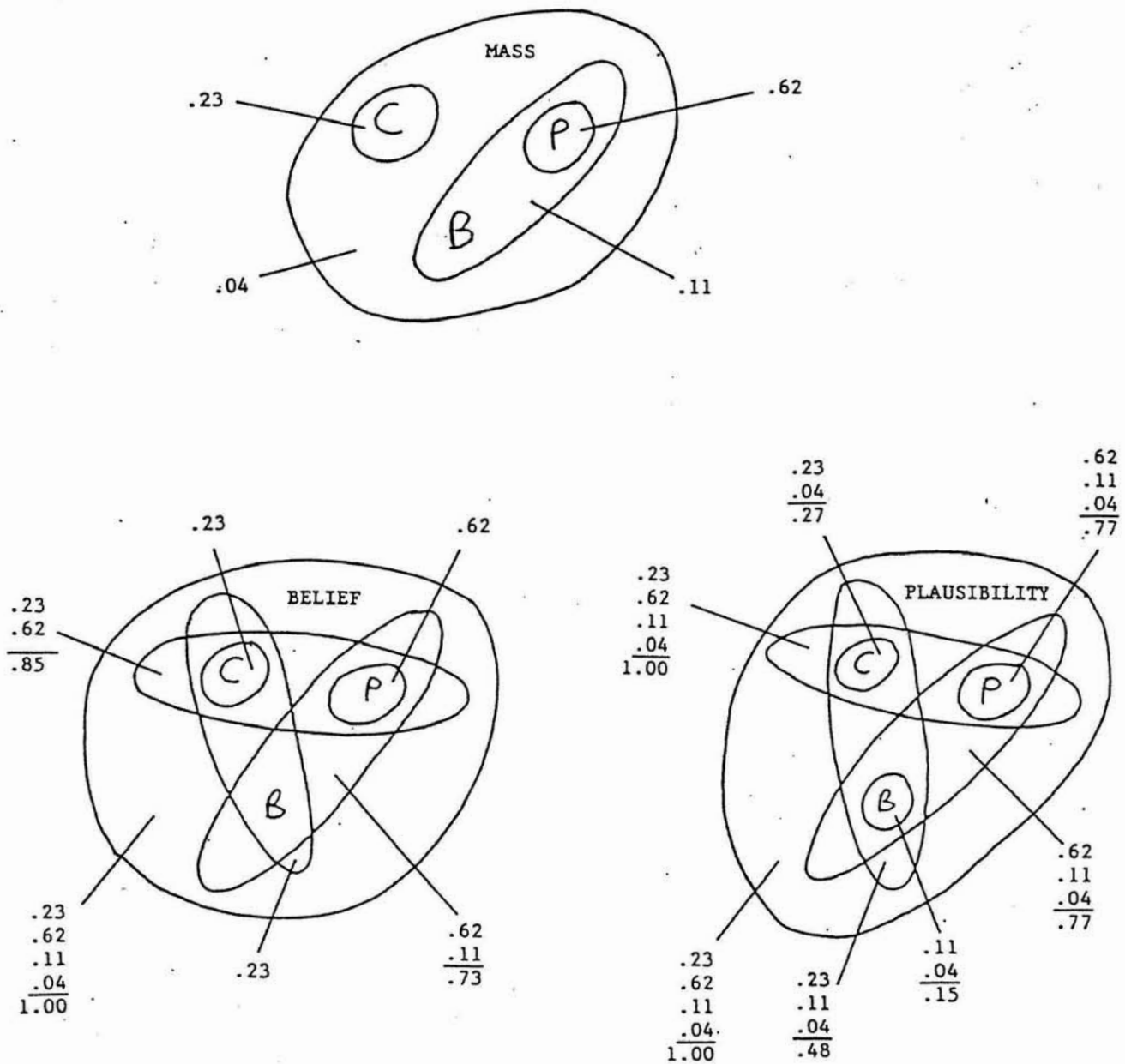
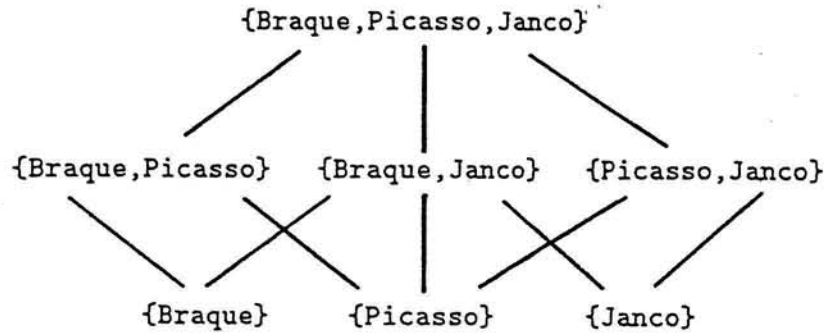
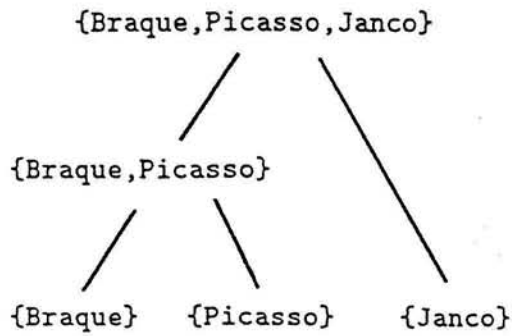


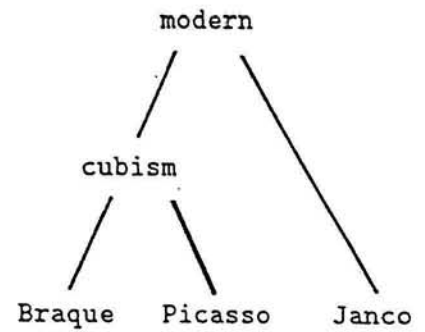
Figure 3: An illustration of the relationship that exists among the a mass (top), belief (left), and plausibility (right) functions that represent the same set of primitive degrees of support.



3-a



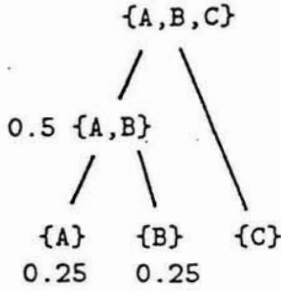
3-b



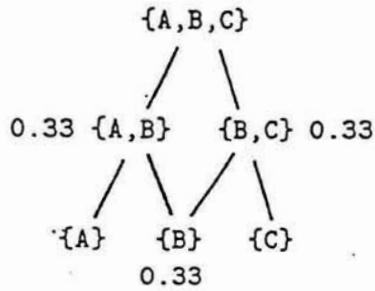
3-c

Figure 4: The evolution of a taxonomy from a lexical power set

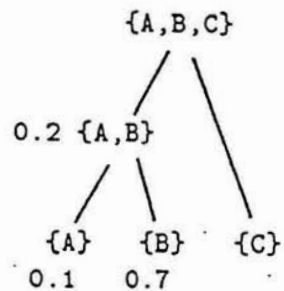
The U1 taxonomy
with $m(X,d,q_1)$ values:



The U2 taxonomy
with $m(X,d,q_2)$ values:



The U taxonomy
with $m(X,d)$ values:



\oplus

=



U1	A	B	C	\otimes	U2	A	B	C	=	U1xU2	A	B	C
u1:	1	1	0		u1':	0	1	1		u1,u1':	0	1	0
u2:	1	0	0		u2':	1	1	0		u1,u2':	1	1	0
u3:	0	1	0		u3':	0	1	0		u1,u3':	0	1	0
u4:	1	1	0							u2,u1':	0	0	0
										u2,u2':	1	0	0
										u2,u3':	0	0	0
										u3,u1':	0	1	0
										u3,u2':	0	1	0
										u3,u3':	0	1	0
										u4,u1':	0	1	0
										u4,u2':	1	1	0
										u4,u3':	0	1	0

Figure 5: The individual indexing opinions of two groups of catalogers (U_1 and U_2) are recorded at the bottom of the figure. These opinions induce two different taxonomies and two different relevance functions, $m(X,d,q_1)$ and $m(X,d,q_2)$, depicted at the top of the figure. The combination of the relevance taxonomies via Dempster's rule \oplus at the classifiers level and the combination of the opinions via the cartesian consensus rule \otimes at the catalogers level leads to the same pooled index depicted at the top right of the figure.