

Using Dependence Relations in MeSH as a Framework for the Analysis of Disease Information in Medline

Lowell Vizenor

Medical Ontology Research
U. S. National Library of Medicine
Bethesda, MD 20894 - USA
vizenorl@mail.nih.gov

Olivier Bodenreider

Medical Ontology Research
U. S. National Library of Medicine
Bethesda, MD 20894 - USA
olivier@nlm.nih.gov

Abstract

Motivation: Many terminologies such as MeSH have a hierarchical structure, but no trans-ontological relations. Such relations, however, are useful for characterizing the relations identified in data sets. **Methods:** We identified trans-ontological relations in MeSH (between diseases and other categories) based on formal ontological principles and compared them to co-occurrence data in MEDLINE. **Results:** Dependence relations identified between a disease and other categories generally correspond to the highest proportion of relations between this disease and any other category under investigation. The other relations observed between this disease and other categories correspond mostly to contingent relations.

1 Introduction

There exist many kinds of relationships between biomedical entities. For example, a viral meningitis is *located in* the meninges, it is *caused by* some virus, and it can be *treated by* some antiviral drugs. Such relations are recorded explicitly as symbolic relations in biomedical knowledge bases and, to a lesser degree, in terminological resources such as SNOMED CT and MeSH. Moreover, the association between indexing terms (i.e., term co-occurrence) in the citations from a bibliographic database such as MEDLINE also represents statistical relations among these indexing terms. For example, the MeSH term *Viral meningitis* co-occurs frequently with MeSH terms for various virus species, including *Enterovirus B*, *Human* and *Herpesvirus 2*, *Human*. One major difference be-

tween symbolic and statistical relations is that, whereas the nature of the symbolic relations is explicit (e.g., *location of*), the nature of the statistical relations is implicit. However, the frequency of co-occurrence can be analyzed to assess the salience of the association.

Formal ontology provides another perspective on relations, distinguishing between two major kinds of symbolic relations: dependence relations (inherent to the nature of related entities) and contingent relations. These distinctions are presented in detail in section 2.

The objective of this study is to analyze dependence relations in MeSH and to compare them to statistical relations obtained from co-occurrence data. We restrict our analysis to the relations between disease categories and other categories of biomedical interest.

Our hypothesis is that systematic associations will be found between diseases and the types of entities on which they are dependent, namely between diseases classified by location and their corresponding anatomical sites and between diseases classified by etiology and their corresponding causes or agents. In practice, for a given disease, the largest proportion of relations to another category should be to a category on which this disease is dependent, and this systematically for each disease. In contrast, we expect to find a smaller proportion of relations between diseases and other categories of biomedical interest, corresponding to contingent relations.

Besides clarifying the link between dependence relations and co-occurrence, this paper seeks to identify associative relations in MeSH (i.e., relations across trees), which, we expect, will support information retrieval and semantic mining applications.

2 Background

2.1 Relations in Biomedicine

Broadly speaking, there are two types of meaningful relations that hold between biomedical classes. First, there are those relations in which every instance of a class is related to some instance of another class. For example, *every* instance of protein synthesis is (necessarily) related to *some* strand of mRNA.

Second, there are those relations in which only some instances of a class are related to some instances of another class. For example, *some* aspirin tablets are related to *some* headaches. Certainly, the relation between aspirin and headaches is important (the former alleviates the latter), but merely contingent.

Here, though, we will focus exclusively on one type of necessity relation, namely, the relation of ontological dependence. Briefly, an entity *A* is ontologically dependent on *B* if and only if *A* exists then *B* exists. Moreover, we are interested in dependence relations between what we refer to as biomedical continuants and biomedical processes. But before we can state these relations more precisely, we need to introduce some ontological distinctions.

2.2 Formal-Ontological Distinctions

All real world entities in the biomedical domain fall into one of two exclusive categories of continuant and process. Think of the difference between a human being and the event of losing weight. Informally, what changes (the human) is the continuant and the change itself (the weight-loss event) is the process. More precisely, continuants are entities which continue to exist through time; they preserve their identity from one moment to the next even while undergoing a variety of different sorts of changes (Smith and Grenon 2004). Examples include molecules, cells, anatomical structures, and organisms.

Processes differ from continuants in several important respects. Most importantly, though, whereas continuants exist fully at a given time—i.e. all their parts are present at a time—processes never exist in full at a time; instead, they unfold through successive phases. Processes have a beginning, middle and end (Smith and Grenon 2004). Examples include cell division, ion transport, protein synthesis and respiration.

In addition to the distinction between continuants and processes, there is a distinction between independent and dependent continuants. *Independent continuants* (also sometimes called

substances) are entities such as organisms, organs, cells and genes which do not require the existence of any other entity in order to exist. *Dependent continuants* (sometimes called accidents) are entities such as dispositions, functions, properties, qualities, roles and states. A dependent continuant is such as to be fully present at a given time, but nevertheless requires the existence of some independent continuant in order to exist.

Likewise, every process depends for its existence on some (independent) continuant. So, for example, *every* instance of photosynthesis (process) depends on *some* instance of chlorophyll (continuant).

The ontological categories of continuant and process provide a basis for the distinction between *intra-ontological relations* and *trans-ontological relations* (Grenon, Smith and Goldberg 2004; Smith and Grenon 2004; Rosse et al. 2005). Examples of intra-ontological relations are *subsumption* (the relation of one class being wholly included in another), *instantiation* (the relation between an individual [or instance] and a universal [or class]) and *meronymy* (the part-whole relation). These relations are always restricted to a given ontological category. A sound ontological principle for the construction of hierarchies is that all the classes employed should belong to one and only one of these categories. So, no continuant has a process as a part and vice versa. Similarly, no continuant is a subclass of a process and conversely. Trans-ontological relations are those sorts of relations that transcend the ontological divide between the formal ontological categories of continuant and process and are just those ontological dependence relations that hold between continuants and processes.

2.3 Dependence Relations between Biomedical Continuants and Processes

Think of the relation between the process of cell transport and the cell that participates in this process. Without the cell there is no cell transport. What is more, every instance of the processual class *cell transport* is ontologically dependent on some instance of the continuant class *cell*. *Participation* is the relation that holds between any continuant that is involved in some form or another in a process. Examples of participation include (Smith et al. 2005):

death has_participant organism

breathing has_participant thorax.

In what follows, we focus on the (implicit) participation relations in MeSH between processual diseases classified in terms of the anatomical structure or bodily system in which they are located.

(Smith and Grenon 2004) note that special types of participation relations can be distinguished according to whether a continuant is agent or patient in a process. In the case of *disease D has_participant anatomical structure S*, the continuant *S* is a patient in the process *D*. In other words, *S* *passively participates in D*.

But there are also cases where the continuant is an agent in the process. That is, *C actively participates in P*. For example, a pathogen actively participates in (i.e. is the cause of) a pathological process. We refer to relations of active participation as *has_agent*.

These relations are defined more formally in (Smith et al. 2005) as follows:

p has_participant c at t – a primitive relation between a process (*p*), continuant (*c*), and a time (*t*)

p has_agent c at t – a primitive relation between a process (*p*), a continuant (*c*) and a time (*t*) at which the continuant is causally active in the process.

Note that these relations are defined at the level of *instances*. Because terminologies such as MeSH contain essentially names for biomedical classes, we are primarily interested in relations between *classes*. (Smith et al. 2005) define the class-level counterparts of these relations as follows:

P has_participant C = [definition] for all *p*, if *Pp* then there is some *c*, *t* such that *Cct* and *p has_participant c at t*

P has_agent C = [definition] for all *p*, if *Pp* then there is some *c*, *t* such that *Cct* and *p has_agent c at t*

(where *Pp* reads *p* is an instance of processual class *P* and *Cct* reads *c* is an instance of continuant class *C* at time *t*).

2.4 Statistical Relations

In probability theory, two events E_1 and E_2 are independent when the probability of occurrence of the two events simultaneously, $P(E_1 \cap E_2)$, is not greater than the product of the probabilities of occurrence for each event, $P(E_1) \cdot P(E_2)$. Conversely, when $P(E_1 \cap E_2) > P(E_1) \cdot P(E_2)$, E_1 and E_2 are not independent. What we are interested in identifying here are pairs of “non-independent” MeSH terms, whose frequency of co-occurrence

(i.e., simultaneous presence as indexing terms in a MEDLINE document) is higher than would be expected if the two terms had been used independently by the indexers.

3 Materials

3.1 MeSH

The Medical Subject Headings (MeSH[®]) is the controlled vocabulary used by the National Library of Medicine (NLM) for indexing articles from 4,800 of the world's leading biomedical journals for the MEDLINE/PubMED[®] database. The MeSH thesaurus is organized in 16 tree-like hierarchies. Examples of top-level categories in MeSH include *Anatomy*, *Organisms* and *Diseases*. Hierarchical relations among MeSH descriptors are indicated by tree numbers assigned to the descriptors. For example, the tree number **C05.550** indicates that *Joint Diseases* is a descendant of *Musculoskeletal Diseases* [C05] in the C tree (*Diseases*). The 2004 version of MeSH is used in this study and comprises 22,658 descriptors.

3.2 MEDLINE

MEDLINE[®] is NLM's premier bibliographic database that contains approximately 13 million references to journal articles in life sciences with a concentration on biomedicine. MEDLINE citations are indexed manually using MeSH as a controlled vocabulary. Descriptors denoting the central topic of an article are indicated with a star. In this study, the co-occurrence of MeSH descriptors is restricted to co-occurrence between starred descriptors. In order to focus on salient associations, only those pairs of MeSH descriptors co-occurring at least 10 times in our collection are considered. The set of MEDLINE citations used in this study contains 385,642 citations and consists of the citations entered in MEDLINE and completed (indexed) between December 2003 and November 2004. Of these, we discarded 151 citations with no descriptor indicating a central topic, retaining a total of 385,491 citations. 20,085 distinct MeSH descriptors were used to index these citations, resulting in 1,378,027 citation-descriptor associations.

4 Methods

4.1 Identifying Dependence Relations

In this part of our study, we identify dependence relations across hierarchies manually, by examining the top-level categories in MeSH in the light

of both the formal ontological principles presented earlier and domain knowledge. We focus on relations between diseases and other categories, distinguishing between dependence and contingent relations. Moreover, we restrict our study to those diseases that are best understood as pathological processes. Note that not all diseases are processes. Think, for example, of velocardiofacial syndrome.

The first step, then, is to identify those MeSH disease descriptors that can be interpreted as a pathological process. For this reason we omit *Neoplasms* [C04], which is a term that refers to pathological continuants (i.e. tumors). The second step involves identifying the biological continuants that participate (actively or passively) in these pathological processes. Once we identify the dependence relations that exist between diseases and other categories we contrast these relations with contingent relations.

4.2 Identifying Statistical Relations

For a given pair of MeSH terms (A, B), information about their association in documents can be summarized in a two-way contingency table and analyzed statistically (Agresti 1996):

- n_{AB} , the number of documents indexed with both term A and term B
- $n_{A\bar{B}}$, the number of documents indexed with term A but not term B
- $n_{\bar{A}B}$, the number of documents indexed with term B but not term A
- $n_{\bar{A}\bar{B}}$, the number of documents indexed with neither term A or term B

In order to evaluate the statistical significance of the association between two MeSH terms, we use the likelihood ratio test (also called G-test or G-square test). The G^2 statistic compares the maximum of the likelihood function under two circumstances: 1) under the hypothesis of independence and 2) under the general, observed conditions. The G^2 statistic does not have the minimum expected frequency requirements imposed by the chi-square test. However, for the G^2 statistic to be computed, all observed frequencies must be greater than 0.

Because this study focuses on trans-ontological relations, statistically significant associations between MeSH terms are restricted to those pairs (A, B) where A and B belong to different MeSH hierarchies (e.g., between a disease in the C tree and an organism in the B tree).

Finally, in order to amplify the relations between disease categories and other categories in MeSH, the frequencies of co-occurrence are ag-

gregated upwards, more precisely up to the second level of MeSH hierarchies. For example, the frequency of co-occurrence between *Cholera* [C01.252.400.959.347] and *Vibrio cholerae* [B03.440.450.900.859.225] is recorded as co-occurrence between *Bacterial Infections and Mycoses* [C01] and *Bacteria* [B03]. For those terms linked to more than one high-level category (e.g., *Lung Neoplasms*, linked to both *Neoplasms* [C04] and *Respiratory Tract Diseases* [C08]), the frequency of co-occurrence is divided among the corresponding high-level categories during the aggregation process, so as to avoid double-counting pairs.

5 Results

5.1 Expected Dependence Relations

Examples of dependence relations identified based on formal ontological principles among disease categories and other high-level categories in MeSH are presented below. Table 1 presents *has_participant* relations between pathological processes and anatomical entities, while Table 2 shows *has_agent* relations between pathological processes and pathogens.

Pathological process	Anatomical entity
Musculoskeletal Diseases	Musculoskeletal System
Digestive System Diseases	Digestive System
Stomatognathic Diseases	Stomatognathic System
Respiratory Tract Diseases	Respiratory System
Nervous System Diseases	Nervous System
Eye Diseases	Sense Organs (+)
Urological and Male Genital Diseases	Urogenital System
Female Genital Diseases and Pregnancy Complications	Urogenital System Embryonic Structures
Cardiovascular Diseases	Cardiovascular System
Hemic and Lymphatic Diseases	Hemic and Immune Systems
Skin Diseases	Integumentary System
Endocrine Diseases	Endocrine System

Table 1. Dependence relations of the type *has_participant* (between a pathological process and an anatomical entity).

Pathological Process	Pathogen
Bacterial Infection and Mycoses	Bacteria Fungi
Virus Diseases	Viruses
Parasitic Diseases	Animals (+)

Table 2. Dependence relations of the type *has_agent* (between a pathological process and a pathogen).

The MeSH descriptors marked with plus sign (+) represent cases where the corresponding term is too far down on the tree to be recognized here. For example, *Parasitic Diseases has_agent Animals* is true, but it would be more precise to say *has_agent Parasites*, *Parasites* being a subcategory of *Animals*.

5.2 Statistical Relations

25,376 pairs of co-occurring “starred” MeSH terms were extracted from the 385,491 citations in our 2004 MEDLINE set. The associations were statistically significant in all but 68 cases. Of these, 7,896 pairs of terms had a frequency of co-occurrence of at least 10 and were used in our analysis. In what follows, we report frequencies of co-occurrence after aggregation to the first subdivision of the MeSH trees, with emphasis on the associations between disease categories and other categories (6525 pairs).

Table 3 shows the number of associations between a given disease category and all other categories. Remarkably, for each disease category, there is generally one subcategory of the *Anatomy* and *Organisms* trees accounting for highest number of associations between this disease and other categories (e.g., between *Cardiovascular Diseases* and *Cardiovascular System*). Exceptions include *Neoplasms* [C04], *Congenital, Hereditary, and Neonatal Diseases and Abnormalities* [C16], *Endocrine Diseases* [C19] and *Immunologic Diseases* [C20].

Similarly, most *Anatomy* and *Organisms* categories are preferentially associated with one disease category. In contrast, other categories are frequently associated with most disease categories, but not one in particular, including *Pathological Conditions, Signs and Symptoms* [C23], *Amino Acids, Peptides, and Proteins* [D12], *Diagnosis* [E01], *Therapeutics* [E02] and *Surgical Procedures, Operative* [E04].

6 Discussion

6.1 Findings

Our hypothesis, i.e., that dependence relations are systematically corroborated by co-occurrence relations, has been for the most part verified in this experiment. In fact, among the 15 dependence relations identified between diseases and other categories, 11 (73%) correspond to the highest proportions of relations between this disease and any other category under investigation. The most remarkable exception is probably *Neoplasms* [C04]. The top category related to *Neo-*

plasms is *Amino Acids, Peptides, and Proteins* [D12], followed by *Surgical Procedures, Operative* [E04], which both correspond to contingent rather than dependence relations. As already mentioned, tumors are independent continuants, which helps to explain why there is not one particular category on which they depend. For *Endocrine Diseases* [C19], other dependence relations not studied here (namely to hormones) seem stronger than the dependence relation to *Endocrine System* [A06].

In some cases (e.g., *Otorhinolaryngologic Diseases* [C09]), the highest proportion of relations does not correspond to anatomical structures as expected, but to *Surgical procedures, Operative* [E04]. In fact, there is no such thing as an “otorhinolaryngologic system” and the corresponding anatomical structures (i.e., ear, nose and throat) are categorized in MeSH under *Respiratory System* [A04] and *Sense Organs* [A09]. These two categories, when taken together, correspond indeed to the highest proportion of relations for *Otorhinolaryngologic Diseases*. The same observation applies to groupings such as *Female Genital Diseases and Pregnancy Complications*, associated with *Urogenital System* and *Embryonic Structures*.

Of note, in most cases, each independent continuant such as anatomical entities and organisms is also preferentially associated with one disease category: the process which depends on this continuant. Exceptions include *Tissues* [A10], *Cells* [A11] and *Fluids and Secretions* [A12], all associated with several diseases. Finally, as expected, categories other than that predominantly associated with a disease generally stand in a contingent rather than dependence relation with this disease.

6.2 Significance

At first glance, the relations identified between diseases on the one hand and anatomical entities and organisms on the other appear essentially trivial. In fact, the curators of MeSH seem to purposely maintain a parallel in the names of these high-level descriptors (e.g., *Virus Diseases / Viruses*). However, none of the approaches to identifying relations presented in this study relies on the lexical properties of MeSH terms. Identifying dependence relations entirely relies on formal ontological principles and domain knowledge. The statistical relations were computed from MeSH identifiers rather than names. Moreover, relations between *Mycoses* and *Fungi* and between *Skin Diseases* and *Integumentary Sys-*

tem could probably not have been discovered by lexical analysis of the terms.

6.3 Limitations

For practical reasons, this statistical analysis is limited to co-occurrence data extracted from a subset of MEDLINE selected somewhat arbitrarily. The results would need to be confirmed on a larger subset.

Analogously, the statistical analysis of co-occurrence data is the only technique used for identifying associations between MeSH terms in this study. Other techniques such as a vector space model and association rule mining could also be used and may yield different results (Bodenreider, Aubry and Burgun 2005). Our hypothesis should be tested against these different techniques as well.

Finally, we chose to amplify the results by aggregating frequencies of co-occurrence to the top-level categories in MeSH. While useful for this purpose, the aggregation prevented us from identifying fine associations. For example, as mentioned earlier, *Parasitic Diseases* [C03] is associated with *Animals* [B01], rather than *Parasites* [B01.500.714].

6.4 Applications

Semantic mining. The purpose of semantic mining is to identify and characterize the relations among entities of interest in a given domain. The statistical associations identified by most data mining techniques have to be characterized semantically in order to fully interpret the dataset. While this elucidation process is often part of the interpretation of the data, we argue that identifying dependence relations *a priori* may provide a useful framework for interpreting the relations discovered in a dataset.

Terminology creation and maintenance. Terminologies such as MeSH are built manually by experts and the relations explicitly represented in MeSH are limited to organizing descriptors in hierarchies. We showed that the approaches proposed in this paper can help identify and partially characterize trans-ontological relations. Representing dependence relations explicitly would undoubtedly facilitate terminology maintenance. Such relations could, for example, inform the MeSH editing environment that the modification of a descriptor must trigger the review of all descriptors standing in a dependence relation with the modified descriptor.

6.5 Generalizability

To other subdomains. This study was purposely limited to the relations between diseases and other categories. As mentioned earlier, many other categories are involved in dependence relations. Namely, all processes, not only pathological processes such as diseases, are dependent on some continuant. Physiology, for example, is a subdomain where many dependence relations can be identified.

To other terminologies. Similar to MeSH in many respects is the Gene Ontology™ (GO), a controlled vocabulary used to annotate (i.e., index) gene products across various model organisms. No trans-ontological relations are currently represented explicitly in GO. We showed in previous work that the approaches proposed in this paper would be applicable to GO as well (Bodenreider et al. 2005). While the statistical analysis of co-occurrence data requires a dataset, the formal ontological analysis can be applied to any terminology.

Acknowledgments

This research was supported in part by the Intramural Research Program of the National Institutes of Health (NIH), National Library of Medicine (NLM). This work was done while Lowell Vizenor was a visiting fellow at the Lister Hill National Center for Biomedical Communications, NLM, NIH. Our thanks go to Jim Mork and Kelly Zeng who helped extract the set of MEDLINE citations.

References

- Agresti, A. 1996. *An introduction to categorical data analysis*. New York: Wiley.
- Bodenreider, O., Aubry, M. and Burgun, A. 2005. Non-lexical approaches to identifying associative relations in the gene ontology. In R. B. Altman, A. K. Dunker, L. Hunter, T. A. Jung and T. E. Klein (eds.). *Pacific symposium on biocomputing 2005*, pp. 91-102. World Scientific.
- Grenon, P., Smith, B. and Goldberg, L. 2004. Biodynamic ontology: Applying BFO in the biomedical domain. In D. M. Pisanelli (ed.). *Ontologies in medicine*, pp. 20-38. Amsterdam: IOS Press.
- Rosse, C., Kumar, A., Mejino, J. L. V., Cook, D. L., Detwiler, L. T. and Smith, B. 2005. A strategy for improving and integrating biomedical ontologies. *Proc AMIA Symp*: 639-643.
- Smith, B., Ceusters, W., Klagges, B., Kohler, J., Kumar, A., Lomax, J., Mungall, C., Neuhaus, F., Rector, A. L. and Rosse, C. 2005. Relations in biomedical ontologies. *Genome Biol* 6(5): R46.
- Smith, B. and Grenon, P. 2004. The cornucopia of formal ontological relations. *Dialectica* 58(3): 279-296.

