

# Mining Web Documents for Unintended Information Revelation \*

Rohini K. Srihari  
CEDAR/ Dept. of Computer  
Science and Engineering  
State University of New York at  
Buffalo  
rohini@cedar.buffalo.edu

Munirathnam Srikanth  
CEDAR/ Dept. of Computer  
Science and Engineering  
State University of New York at  
Buffalo  
srikanth@cedar.buffalo.edu

Zhaohui Zheng  
CEDAR/ Dept. of Computer  
Science and Engineering  
State University of New York at  
Buffalo  
zzheng3@cedar.buffalo.edu

## ABSTRACT

This research concerns web site information security. With an increasing number of documents being generated by different individuals and departments in organizations, there is a potential of releasing information which is inconsistent with the overall goals, objectives and operation of the organization. We refer to this as *unintended information revelation* (UIR). This paper focuses on progress made towards detecting UIR scenarios based on text mining of portions of the Federal Aviation Administration website; progress includes the decomposition of this problem into subtasks with the focus on deriving content representation conducive to text mining. Information extraction and information retrieval techniques are employed to first extract concepts and associations from the web documents in a data-driven manner. Various heuristics are used in identifying relevant concepts and associations, including linguistic knowledge, domain knowledge, as well as statistical information. The concepts are automatically mapped into an FAA domain ontology using corpus-based statistical techniques. Associations between concepts reflect both named and unnamed relationships between concepts; these reflect ontological, syntactic, semantic, and general contextual relationships between concepts. Together, the concepts and associations constitute a graph over the knowledge space corresponding to the documents. Graph mining techniques are applied in order to reveal interesting patterns of information dissemination. Graph mining techniques have been employed in situations where precise data is available, typically transactional data in relational databases. This work reflects an attempt to use graph mining on information obtained from processing natural language. Furthermore, this work extends the mining beyond detecting substructures involving named entity categories.

---

\*This work is sponsored by NSF grant IIS-0325404 and a research grant from the FAA 032-G-009

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to publish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.  
Copyright 200X ACM X-XXXXX-XX-X/XX/XX ...\$5.00.

## Categories and Subject Descriptors

H.3.1 [Information storage and retrieval]: Content analysis and indexing—*Linguistic Processing, Indexing methods*

## General Terms

Algorithms, Experimentation, Security

## Keywords

text mining, information extraction, graph mining

## 1. INTRODUCTION

There are potentially valuable nuggets of information hidden in document collections generated by multiple authors, working independently at various times. Such information is not explicitly mentioned, but can be inferred by following chains of concepts and associations connecting them. In the field of scientific discovery, this could reveal interesting connections between genes and the diseases they cause. In the case of the WWW, users generate subsets of documents through purposeful querying or surfing activity; in some cases, users could be seeking information with malicious intent. The problem of *unintended information revelation* (UIR) is a special case of text mining where the goal is to quantify information in a given subset as well as detect any unusual patterns of information dissemination. Figure 1 illustrates the problem that is being tackled; the ultimate goal is to develop a UIR toolkit that can be used to detect such patterns. There are several challenges tasks to be solved before reliable UIR detection is possible. This paper is concerned with the first task, which is deriving a rich content representation such that text/graph mining algorithms can be exploited. Initial experimentation with this graph representation demonstrates its ability to model information patterns beyond what current information retrieval (IR) models can.

A solution to this problem requires a very sophisticated representation of information content, certainly more than what information retrieval (IR) systems use. It is tempting to use information extraction as a substitute for the IR *bag-of-words* representation, but this is not wise. Information extraction systems are fragile in several aspects, particularly with respect to recall. It is a formidable task to extract all the different kinds of entities and relationships that are salient. IR systems on the other hand, are robust.

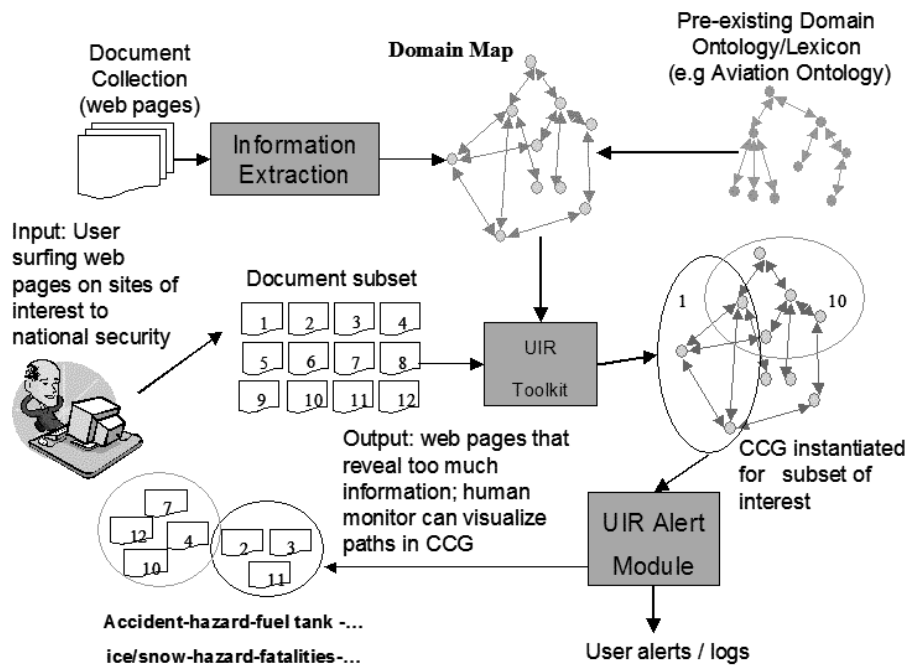


Figure 1: Detecting Unintended Information Revelation

Thus, a hybrid approach combining the best of both is desired: such an approach is presented in this paper. A *concept chain graph* is constructed based on corpus-processing as well as a predefined domain ontology. The concept chain graph reflects statistics pertaining to salient entities (named entities and other concepts) and associations (both named and unnamed) in the data. The CCG can be used either in conjunction with (i) a probabilistic framework in detecting interesting chains or patterns or (ii) a graph mining framework in detecting anomalies. The latter use is the focus here.

This paper focuses on analyzing and mining the Federal Aviation Administration (FAA) website for interesting patterns of information dissemination. Examples of such patterns include an imbalance of information (especially financial data) pertaining to safety measures implemented in different regions of the country. Other scenarios could involve information about individuals and their activities, especially those involved in senior administration or management. At this stage, the website is being analyzed independent of any specific user surfing activity; such analysis could be performed at a later time if required. Thus, this effort takes a different perspective with respect to security, namely the potential hazards associated with providing too much (or imbalanced) information. As an added benefit of this work, a UIR toolkit is being developed that permits information analysts to visualize the manner in which information is presented. Traditionally this has consisted of a site map; in this effort, users are given a concept graph. From this, the most significant chains connecting two or more concepts can be visualized (along with the corresponding document trail).

The solution presented in this paper takes into account both corpus-based statistical techniques as well as graph-based data mining techniques. The former are used in order to process the document collection, extract salient concepts

and associations between concepts, and subsequently map the concepts into a predefined domain ontology. This information is modeled as a graph which is subsequently used in graph mining experiments designed to reveal anomalies corresponding to UIR situations. Figure 1 illustrates these major steps in the approach presented here.

The approach here is novel for several reasons. First, it attempts to use a hybrid representation of information content, by incorporating both an information retrieval index as well as richer content produced by an information extraction system. Whereas most information extraction based text mining approaches rely on named entity categories and named relationships between them, this effort (i) goes beyond traditional named entity categories (e.g. person, organization, location) to include general concept categories, and (ii) considers both named and unnamed associations between these general concept classes. Finally, this represents an attempt to combine statistical text mining approaches with graph-based approaches.

The remainder of the paper is organized as follows. Section 2 discusses related work; Sections 3 and 4 discuss concept extraction and mapping to a domain ontology. Sections 5 and 6 discuss indexing of concepts and associations and the resulting concept chain graph respectively. Section 7 discusses details of the graph modeling and mining framework. Finally, Section 8 discusses experimental results on a large subset of the FAA website.

## 2. RELATED WORK

The work discussed here brings together several research topics in information retrieval and text mining and hence there is a significant amount of related research. Since the focus of this paper is on sophisticated content representations conducive to text/graph mining, it is useful to first examine related work in this area. This is followed by re-

lated work in sub-areas such as ontology mapping as well as graph-based anomaly detection.

Two domains that have seen significant activity in text mining include the biomedical domain as well as the intelligence analysis domain, the latter especially in the context of the DARPA Evidence Extraction and Link Discovery (EELD) program. Srinivasan and her collaborators at The University of Iowa [23, 20] have developed an approach for biomedical text mining that utilizes a predefined medical thesaurus (UMLS) in identifying significant concepts.

Current techniques for text mining are based on shallow information extraction (IE) followed by clustering or link analysis. This leads to more spurious patterns than significant ones. Recently, there have been attempts to guide the mining process by using models, reflecting patterns of interest. Such models require the underlying data to be well structured and noise-free, which is not the case in data obtained from IE systems. Probabilistic Relational Models and more sophisticated variants [24] have been proposed as a technique for combining relational models with Bayesian networks [7]. These are very effective in cases where the user model (types of inferences) is well understood. [11] discusses the need for unifying IE and data mining through conditional-probability, relational models. Their focus however is mostly on the uncertainty associated with IE output; the focus here is on developing models that are designed specifically for determining the best chains/sequence connecting two concepts based on associations between them. Such chains could traverse taxonomical links, or unnamed/named relationships through both named entities and generic concept/topic categories. Thus, the need for extracting and representing general concept categories (beyond named entities). Most important, such chains cut across document boundaries.

A necessary step to finding such chains is the automatic detection of local (sentence, paragraph or document-level) associations. There have been attempts to automatically learn non-taxonomical conceptual relationships from corpora [10]. Some of these methods are based on classifying semantic relationships based on analysis of noun compounds [19].

The work on concept mapping extends current work in ontology acquisition by integrating pre-existing ontologies with automatically derived concepts and associations. Domain ontologies have proved to be an important tool in applications that require semantic knowledge, i.e. semantic web. However, the creation and maintenance of large ontologies require significant investment in time and resources. General ontologies such as Wordnet and Cyc have been used in many applications but require extensive work to tailor them for specific domains. Earlier attempts focused on using natural language processing to extract domain specific terms [15, 16, 17]. Two main approaches have been proposed for trying to solve the ontology creation “bottleneck”. The first approach is to create tools that facilitate the manual labor involved in ontology engineering by integrating ontology editors, consistency checkers, and ontology merging-integration tools. A good example of this approach is Protégé [6]. The second approach uses natural language processing and machine learning to extract important concepts and relationships from existing domain text and databases to build a domain specific ontology. Faure and Poibeau [5] present a system called ASIUM that learns cooperatively seman-

tic knowledge from text. ASIUM learns subcategorization frames of verbs and an ontology of specific terms. Relations between these terms and frames of verbs is learned by using a conceptual clustering algorithm. The system also has a user interface to allow setting clustering parameters and visualize learned concepts and relations. OntoLearn [12] present a system that combines machine learning and natural language processing for ontology construction. OntoLearn uses information extraction to find candidate terms. It uses a domain corpus and a general purpose contrasted-usage corpora. Semantic relations between terms are learned using lexical similarities among terms, identifying possible senses, and mapping them against wordnet synsets. The final step uses inductive machine learning algorithms by manually tagging a small subset of the candidate terms and then building a decision tree using C4.5 that is used to identify semantic relations not included in Wordnet. Cimiano *et al.* [3] have proposed a method for automatic acquisition of taxonomies from text which is based on Formal Concept Analysis (FCA). For their system they first parse the text and then identify verb-object dependencies. They use three different measures for estimating the strength of relation between verb-objects pairs extracted (preference strength). Finally these verb-pairs and their respective preference strength are passed to the FCA module which builds a lattice of all the verb-object pairs in the domain.

Finally, it is useful to examine relevant work pertaining to the ultimate goal here, namely anomaly detection. Anomaly detection has been widely studied in the computer security literature, which is defined as the task of quantifying the usual or acceptable behavior and detecting other irregular behavior as potentially intrusive. [9] Previous research in this area typically used strings or attribute-value data as the medium from which anomalies are to be extracted. Graph mining techniques [4] have also been used effectively for this task. These methods are designed to mine more richly structured datasets; in this paper, an attempt is made to apply them to semi-structured data.

### 3. CONCEPT EXTRACTION

This section discusses the extraction of salient *concepts* and *associations* reflected in a target document collection. The resulting set of extracted concepts are used for two purposes: (i) as the indexing vocabulary for the IR module, and (ii) as the set of candidate concepts that are eventually mapped into an FAA domain ontology; this process significantly prunes this set. Figure 2 illustrates the various steps in processing text, starting with customizing InfoXtract with a domain-specific glossary, using InfoXtract output to generate concepts, indexing the collection using these concepts, and finally mapping concepts into the ontology. The resulting concepts and associations are used to generate a concept chain graph, which in turn is used by the graph mining programs. InfoXtract [8, 21, 22] is a domain-independent and domain-portable, intermediate level IE engine that is used in this effort. The IE outputs from the InfoXtract engine used in concept and association extraction are described:

- **NE:** Named Entity objects represent key items such as names of *person*, *organization*, contact information such as *address*, *email*, *phone number*, *URL*, time and numerical expressions such as *date*, *year*, *etc.*

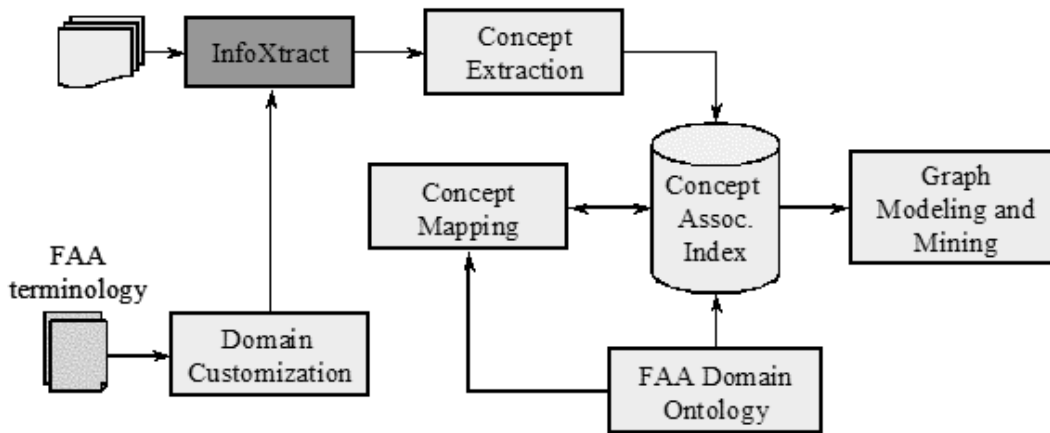


Figure 2: Steps in Mapping Content into Domain-Specific Ontology, Generating Graph

- **CE:** Correlated Entity objects capture relationships between entities such as the *affiliation* associated with a person, for example *Tom Smith, head of Purchasing at Softsoft Inc...*
- **GE:** General Events are verb-centric information objects representing 'who did what to whom when and where' at the logical level. Concept-based GE (CGE) further requires that participants of events be filled by EPs instead of NEs and that other values of the GE slots (the action, time and location) be disambiguated and normalized.

The output from the engine is returned as an XML document reflecting not only the IE output mentioned above, but also a token list dump providing details of the parsing process, including semantic features associated with each word/token. The InfoXtract engine is a series of modules beginning with a tokenizer that reads the text, tokenizes it, and converts it into a data structure known as the token list. Subsequent processing modules, whether they are grammatical (e.g. parsing modules) or machine learning modules (e.g. the NE tagger) all modify the token list.

InfoXtract has been customized for this task by the incorporation of domain knowledge. Domain knowledge is essential for an IE engine to tag documents with named entities in the domain of interest. Various online glossaries of FAA terms were used to construct a single glossary representing the entity class "FAATerm". These were subsequently tagged by InfoXtract which in turn causes nearby concepts to be selected.

The following strategy is adopted in extracting concepts. First, all named entities (and their aliases reflected in the EP) are considered to be concepts. Second, arguments of general events, specifically, the subject, verb and object slots are considered to be key concepts; many of these arguments are noun or verb groups such as *air safety*. InfoXtract performs normalization of all words, rather than the traditional stemming used by IR systems; the root form of the word is returned. Furthermore, phrasal verbs such as *step down* and *take off* are currently recognized and grouped.

With respect to extracting associations, we employ a strategy of using both *named* and *generic* relationships between concepts. This strategy is required in order to increase the

*recall* of the IE system. Association extraction currently uses a simple strategy. All CE relationships extracted (including syntactic relationships) are used. Since general events associate concepts through verbs, these individual relationships are considered. Finally, proximity between concepts (such as sentence or paragraph level) are also taken as evidence of an association. The section on Graph Modeling describes additional associations that are used. The two most important tasks in the IE effort is to increase the precision and recall of concept extraction and association extraction. Future work will include machine learning techniques [13] to automatically learn new NE types and relationships.

## 4. CONCEPT MAPPING

This section describes the method used to map candidate concepts into the FAA domain ontology. The domain ontology (shown in Figure 1) consists of three levels. The concept mapping effort focuses on mapping candidate concepts into the most specific level of the ontology possible. The key to this task is a rich representation of the context of a term which is based on co-occurrence information. The mapping phase takes initial seeds for each target concept and uses these in a bootstrapping technique to map additional concepts.

**Co-occurrence indexing** Indexing concepts are those produced by the concept extraction phase. Co-occurrence frequency between two indexing concepts is the number of times they appear in same sentences. Finally, we have around 600,000 X 600,000 co-occurrence matrix. Only non-zero elements are kept.

**Concept Zoning** For each target concept, a set of initial seeds is generated by (i) using string-matching heuristics to generate matching candidate concepts, and (ii) using these as seeds to generate related terms, and (iii) manual filtering of the results of step (ii). A bootstrapping-like technique was used to generate new concepts given the small number of "seeds".

For each seed:

- calculate the correlation coefficient(CC) between it and any other indexing concept in the collection.

**Table 1: Portion of the FAA Ontology**

Subject Area	Sub-Subject Area	Sample Logical Entities
Safety	Safety Event & Factors	Accident, Incident, Event Factor etc.
	Safety Assessments	Safety Assessment Conclusion, Safety Program
Personnel & Training	Persons	Aviation Person, Aviation Position Type, Job Task etc.
	Positions	Organization Position, Position Type, Federal Job Classification etc.

- sort the indexing concepts according to their CC values with the seed.
- the top few indexing concepts with maximum CC will be selected as most associated concepts of that seed.

Where: Correlation coefficient of two words  $t_1$  and  $t_2$  is defined as

$$CC(t_i, t_j) = \frac{\sqrt{N}[P(t_i, t_j)P(\bar{t}_i, \bar{t}_j) - P(t_i, \bar{t}_j)P(\bar{t}_i, t_j)]}{\sqrt{P(t_i)P(\bar{t}_i)P(t_j)P(\bar{t}_j)}}$$

where:  $N$  is the total number of sentences.

The above procedure can be done iteratively. Two iterations are required. The new concepts generated by the first iteration represent the salient context words of the seed. Since synonyms seldom occurs within the contexts of the seed, the synonyms cannot be obtained after one iteration. However, after the second iteration, we can bring in the synonyms.

**concept mapping** The assumption we made for concept mapping is that the sense of each concept could be decided by its contexts, therefore, two words appearing in the similar contexts should have similar meaning. The semantic similarity between two words can be decided by the similarity between their contexts, e.g. co-occurrence vectors. Our task is to map each new concept to the leaf entity of the predefined hierarchy. Each leaf entity is provided a few keywords. We, therefore

- represent each new concept as its co-occurrence vector
- represent each leaf entity as the centroid of the co-occurrence vectors of its keywords.

The concept mapping consists of the following three steps:

1. each new concept  $t$  is represented as a vector, whose elements are features  $f_i$ , and whose values are the standard tfidf weights:

$$(1 + \log n(t, f_i)) \times \log \frac{N}{sf(f_i)}$$

2. each leaf topic in FAA hierarchy will be represented as the centroid of the co-occurrence vectors of all its elements(keywords)  $e_j$ :

$$(1 + \log \sum_j n(e_j, f_i)) \times \log \frac{N}{sf(f_i)}$$

3. the confidence of the new concept belonging to each leaf topic will be the cosine similarity between the corresponding vectors.

where:  $n(t, t')$  is the co-occurrence frequency between  $t$  and  $t'$ .  $sf(f_i)$  is the number of sentences  $f_i$  appears in.

Initially, approximately 547,000 candidate concepts are generated; this does not include the concepts corresponding to the standard NE categories. Removing those that do not appear at least twice leaves around 275,525 concepts. Of these, 3,467 unique concepts are eventually mapped into the domain ontology (less than 1%). If concept hits (rather than unique concepts) are considered, approximately 12.5% of hits are mapped successfully. A later section discusses the mapping of additional concepts by using verb-subcategorization frames. The next phase of this work will consider the task of mapping additional concepts into the ontology by (i) splitting the set of mapped concepts, thus creating a new bottom level in the ontology, and/or (ii) modifying higher levels of the target ontology.

## 5. CONCEPT-ASSOCIATION INDEX

This section details the integrated index for concepts and association that provides the different views of the document collection as well as FAA domain – the domain of interest. We refer to the index build for a search engine [2] as an IR index. An IR index facilitates efficient search for words and possibly, phrases to retrieve documents relevant to our query. The Concept-Association Index has additional layers built on top of the IR index to facilitate efficient search for concepts and associations. A layered indexing architecture facilitates robust estimation of concept and association weighting based on the corpus statistics available in the different indices.

The three-layered architecture for the Concept-Association Index has the IR index in its bottom most layer. It views of a document as a 'bag-of-words'. The middle layer is the concept index. This is views the document as a 'bag-of-concepts'. The top layer is the association index which views a document as a collection of associations. These layers are linked to facilitate efficient probability estimation. It is noted that the association index can store multiple associations extracted between two concepts. This is possible when the underlying concepts and associations are extracted from a document collection and, in our representation, concept types are also treated as concepts(Sec 6).

## 6. CONCEPT CHAIN GRAPHS

This section describes the graph-based representation of the concepts and associations extracted from documents in a document collection. The concepts mapped into the domain ontology induces a graph that can be used for mining interesting patterns.

A *concept graph* is a labeled graph  $G = (V, E)$  where the nodes are concepts and the edges are the associations be-

tween these concepts. The set of nodes of the graph includes:

- Concepts extracted from the documents using InfoXtract (ref. 3),
- Concepts defined in the domain ontology, and
- Concept type is also viewed as a concept and hence, are represented as nodes in the graph. This includes the types of named entities extracted by the InfoXtract engine.

The edge set includes:

- named and unnamed associations extracted between concepts from the document collection,
- associations described in the domain ontology. This ISA hierarchy of concepts in the FAA domain,
- mapping of named entity instances extracted from the document collection to their types,
- mapping of concepts to the domain ontology(ref. 4, and
- linguistically motivated mapping of concepts to their types.

Linguistically motivated mapping of concepts maps concepts to their concept types based on the surface form of the concept. For example, *en route segment* is *route segment*, *foreign passenger airline* is a *passenger airline*, *department of defense* is a *department*. This mapping is used when a concept cannot be mapped to any of the concepts defined in the domain ontology. While some of the mapping using simple heuristics are incorrect (*vice president* is a *president*), this enables us to bucket concepts into concept types.

It is noted that in a concept graph, concept types are also viewed as nodes in the graph. Hence, 'Person' is a node in the concept graph and an occurrence of a person name in a document results in a node for that person name and an edge between that node and the node 'Person'. While modeling real-world applications in relational settings, one distinguish *properties* or attributes of an entity from its *relations*. *Age\_Of* is modeled as a property of a person, while *Manager\_Of* is modeled as a relation. In concept graph, there is no distinction between these types of relations and all of them are represented as associations. Figure 3 is an portion of concept chains discovered by the system related to the Age-60 rule for pilots. The concept, *Age\_60\_rule*, is typed as a *rule* and both the concept and its type are represented as nodes. Such a representation allows associations be defined between concept types. Over time, machine learning can be employed to learn associations between concept types and include them in the graph.

Since the nodes and edges represent concepts and associations, respectively, discovered in a document collection, a number of weighting schemes can be explored to weight the nodes and edges and use them in the graph mining process. We have used occurrence counts for weighting the nodes and edges in the experiments described in this paper. We plan to explore other weighting schemes that have experimented with and reported in IR literature for text retrieval.

## 7. GRAPH MODELING

The concept graph described in the previous section provides a richer representation for information from the domain ontology and the document collection. Such a representation is transformed to facilitate the search for interesting patterns in the concept graph. The transformed graph retains most of the structure of the graph but takes different labels for the nodes and edges.

The nodes and edges of a concept graph represents the instances of concepts and associations in a document corpus. Here we are interested in mining for patterns between concept types. A node corresponding to an instance of named entity is labeled by the type of the named entity. Nodes corresponding to the extracted concepts are labeled by the target concept names in the domain ontology. This generates what we refer to as *concept type graph*. This graph is used for mining interesting substructures from the data.

Cook and Holder [4] have used graph-based data mining method to discover interesting patterns in databases by representing them as graphs and searching for interesting substructures or subgraph in it. Mining techniques are explored to search for substructures in the *concept type graph*. Subdue [1] is a relational learning system used for finding substructures that appear repetitively in a graph. It is an iterative process by which *best* substructures are identified to be replaced by a node to compress the graph. Best substructures are identified based on their Minimum Description Length (MDL) [18]. We use the SUBDUE software on the *concept type graph* to discover interesting patterns.

Graph-based anomaly detection was recently introduced by Noble and Cook [14], in which anomaly detection was applied to graph-based data. Two methods were described namely anomalous substructure detection and anomalous subgraph detection respectively. The first method looks for specific and unusual substructures within a graph. In the second method, the graph is first partitioned into distinct sets of vertices, (subgraphs), each of which is tested against the others for unusual patterns. We have used the second method in identifying and quantifying information contained in document subsets.

Given a document corpus and a document subset, we use anomalous subgraph detection technique to quantify the information in the document subset. Noble and Cook [14] define a measure,  $A$ , that quantifies the anomaly of the subgraph. It is given by

$$A = 1 - \frac{1}{n} \sum_{i=1}^n (n - i + 1) * c_i \quad (1)$$

where  $n$  is the number of iterations and  $c_i$  is the percentage of the subgraph that is compressed away on the  $i^{th}$  iteration.  $c_i$  is computed using the description length of the graph and is given by

$$c_i(G) = \frac{DL_{i-1}(G) - DL_i(G)}{DL_0(G)} \quad (2)$$

Here  $DL_i(G)$  is the description length of the subgraph after  $i$  iterations.

We generated the concept type graph for the document corpus as well as the document subsets corresponding to the regions. All regional graphs were iteratively compressed using the best substructures found in the concept type graph of the whole document collection. The intuition being the

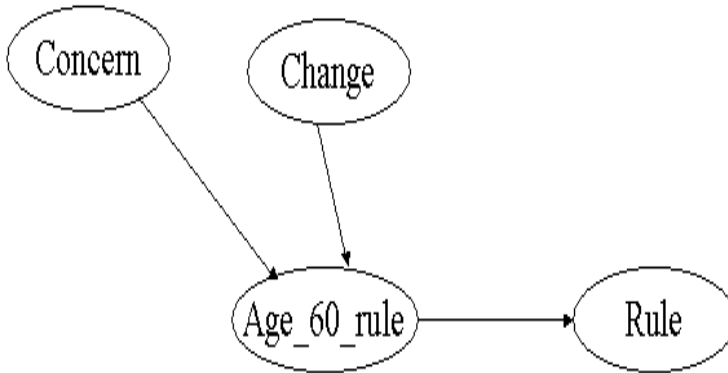


Figure 3: Portion of a concept chain about Age-60-rule

graph which is least compressed at the end of these iterations is a graph representing a subset of document collection that differ from the other subsets. The substructures from the document corpus represent the most frequently occurring concept and association patterns in the whole collection.

## 8. EXPERIMENTS

This sections describes our initial experiments to mine web documents for unintended information revelation. The objectives of our experiments are to motivate and demonstrate that the rich representation of documents using concepts and associations yields well to identifying interesting patterns in document subsets and quantify the information contained in them.

The document corpus for our experiments was crawled from the FAA website. It is noted that the document corpus is a small subset of the complete FAA website whose contents change very frequently. The corpus of 12000 documents includes HTML, Word documents and PDF files. The documents were preprocessed to extract their textual content before they are processed by InfoXtract. Some amount of information that can be extracted based on the structure of the documents is lost due to document preprocessing. The concepts and associations in each document are extracted and mapped to the domain ontology as described in section 4. A concept association index is generated for the document corpus using the extracted concepts and associations. Mapped concepts are also included in the concept association index. If a concept,  $C_e$ , occurs in a document and is mapped to  $C_m$  in the domain ontology, then the occurrence of  $C_m$  in the document is recorded in the concept association index as though it occurred at the same position as concept  $C_e$  in the document.

The UIR problem is defined in the context of document subsets. A user visited a set of documents in the website during a session, we are interested in quantifying the information in the set of visited documents. In the absence of user browsing and querying patterns in FAA website we considered document subsets that are “about” an FAA airport region. FAA has defined 9 airport regions: Alaska, Central, Eastern, Great Lakes, New England, Northwest Mountain, Southern, Southwest, and Western Pacific.

The concept type graph for the whole document corpus

and the document subsets corresponding to airport regions were generated from the concept association index. Table 2 gives information on the size of the document subsets, number of nodes and edges in their corresponding concept type graphs. The procedure described in section 7 is applied to quantify information in the document subsets. These values are also given in Table 2.

Region	Docs.	Concepts	Assocs.	A-Measure
Central	485	68349	50169	0.9784
South	127	10194	6753	0.9893
Eastern	14	1009	370	0.9915
Southwest	510	37433	24809	0.9924
Alaska	392	7633	5783	0.9941
NW Mountain	1443	61980	44404	0.9942
New England	84	14116	11911	1.0031
Great Lakes	129	12384	10044	1.0064
Western Pacific	50	12906	9095	1.0073

Table 2: MDL measure Collection size and Concept type graph size for the different document subsets

The table includes collection size and the size of concept type graph for the document subsets for different regions. The document subsets cover a variety of topics in the aviation domain. In comparison to the concept graph of the whole document corpus, each region seem to concentrate on a few select topics. Safety related concepts and associations are discussed widely in the FAA website. Substructures involving these concepts are generated as best substructures for the document corpus’ concept graph. The ranking of regions in table 2 reflects that such substructures found in the concept graph of the Central region. To estimate the A-measure for the different regions, we considered only the first 5 best substructures of the concept type graph of the corpus. It is noted that the ranking of the subsets is independent of the size of the document subsets, indicating that the rich representation of documents using concepts and associations yields well to quantifying the information contained in them.

It is noted that the graph-based method for quantifying the information in the document subsets has not used the semantics of the concepts and associations as well as the

relative importance of the concepts and associations in the corpus. We intend to incorporate such information in measuring information revealed.

While quantifying information contained in document subsets enables us to quickly identify subsets of the document corpus that reveal unintended information, identifying the text corresponding to this information will enable use to take the required action. Here one is interested in the texts that appear in two or more documents that give information unapparent from only one of those documents. In the concept graph representation, this corresponds to identifying concept chains that involve concepts from more than one document. These concept chains can involve concepts that are from the domain ontology. We have obtained some initial results for identifying these concept chains. These link concepts from different documents. An example of a partial concept chain discovered is given in figure 3.

## 9. SUMMARY

This paper has discussed the steps required in detecting unintended information revelation, a problem faced by many organizations that are obliged to disseminate information publicly. The focus of this paper has been on deriving a rich representation of content conducive to text mining. Initial graph mining experiments based on this representation has shown its ability to model richness of information at a more granular level. This permits comparisons between different document subsets to be made, currently at a very holistic level. Selected partitioning of the dataset permits interesting comparisons to be made.

To go from here to the ultimate task of detecting UIR scenarios requires enhancement of the various steps. The key to detect meaningful scenarios is the proper mapping of concepts into a comprehensive domain ontology. This in turn requires refinement of the concept extraction and ontology mapping/creation stage. The other major area which is currently being worked on, is the improvement of association extraction. Corpus-based statistics are being utilized in order to adjust the weights of existing associations; this will also permit the discovery of associations across documents. Finally, a probabilistic framework is being formulated for determining the best chain connecting two concepts. Such improvements are expected to advance the UIR detection goal significantly.

## 10. ADDITIONAL AUTHORS

Miguel Ruiz, Dept. of Library and Information Studies, State University of New York at Buffalo, meruiz@buffalo.edu, Sudarshan D. Lamkhede, CEDAR/Dept. of Computer Science and Engineering, State University of New York at Buffalo, lamkhede@cedar.buffalo.edu

## 11. REFERENCES

- [1] Subdue knowledge discovery system.  
URL: <http://cygnus.uta.edu/subdue/>.
- [2] R. Baeza-Yates and B. Ribeiro-Neto. *Modern Information Retrieval*. Addison Wesley, New York, New York, 1999.
- [3] P. Cimiano, S. Staab, and J. Tane. Automatic acquisition of taxonomies from text: Fca meets nlp. In F. Ciravenga and N. Kushmerick, editors, *Proceedings of the International Workshop and Tutorial on Adaptive Text Extraction and Mining*, pages 10–17, Cavtat-Dubrovnik, Croatia, 2003.
- [4] D. J. Cook and L. B. Holder. Graph-Based Data Mining. *IEEE Intelligent Systems*, 15(2):32–41, 2000.
- [5] D. Faure and T. Poibeau. First experiments of using semantic knowledge learned by asium for information extraction task using intex. In *Proceedings of the First Workshop on Ontology Learning OL '2000*, 2000.
- [6] J. Gennari, M. A. Musen, R. W. Ferguson, W. E. Grosso, M. Crubézy, H. Eriksson, N. F. Noy, and S. W. Tu. The evolution of protégé: An environment for knowledge-based systems development. Technical Report SMI-2002-0943, Stanford Medical Informatics, Stanford University, 2002.
- [7] D. Heckerman, D. Geiger, and D. M. Chickering. Learning bayesian networks: the combination of knowledge and statistical data. *Machine Learning*, 20(3):197–243, 1995.
- [8] R. k Srihari, C. Niu, and W. Li. A hybrid approach to named entity and sub-type tagging. In *Proceedings of ANLP 2000*, pages 247–254, Seattle, WA, 2000.
- [9] S. Kumar, C. Detection, and C. Intrusions. Department of computer sciences, 1995.
- [10] A. Maedche and S. Staab. Discovering conceptual relations from text. In *Proceedings of ECAI-2000*, pages 321–325. IOS Press, Amsterdam, 2000.
- [11] A. McCallum and D. Jensen. A note on the unification of information extraction and data mining using conditional-probability relational models. In *IJCAI 2003 Workshop on Learning Statistical Models from Relational Data*, 2003.
- [12] R. Navigli, P. Velardi, and A. Gangemi. Ontology learning and its application to automated terminology translation. *IEEE Intelligent Systems*, 18(1):22–31, 2003.
- [13] C. Niu, J. Ding, W. Li, and R. K. Srihari. A bootstrapping approach to named entity classification using successive learners. In *Proceedings of the Association for Computational Linguistics (ACL)*, 2003. to appear.
- [14] C. C. Nobel and D. J. Cook. Graph-based anomaly detection. *Proc. of SIGKDD'03*, pages 631–636, 2002.
- [15] E. Riloff. Automatically constructing a dictionary for information extraction tasks. In *Proceedings of the Eleventh National Conference on Artificial Intelligence AAAI-93*, pages 811–816, 1993.
- [16] E. Riloff. Automatically generating extraction patterns from untagged text. In *Proceedings of the Thirteenth National Conference on Artificial Intelligence AAAI-96*, 1996.
- [17] E. Riloff and J. A. Sheperd. A corpus-based approach for building semantic lexicons. In *Proceedings of the second conference on Empirical methods in natural language processing (EMNLP-2)*, 1997.
- [18] J. Rissanen. *Stochastic Complexity in Statistical Inquiry*. World Scientific Publishing Company, 1989.
- [19] B. Rosario, M. Hearst, and C. Fillmore. The descent of hierarchy and selection in relational semantics. In *Proceedings of ACL 2002*, 2002.
- [20] A. Sehgal, X. Y. Qui, and P. Srinivasan. Mining medline metadata to explore genes and their



- connections. In *Proceedings of the SIGIR 2003 Workshop on Text Analysis and Search for Bioinformatics*, 2003.
- [21] R. K. Srihari and W. Li. Information extraction supported question answering. In T. Stralkowski and S. Harabagiu, editors, *Advances in Open-Domain Question Answering*. Kluwer Academic, 2003.
- [22] R. K. Srihari, W. Li, C. Niu, and T. Cornell. Infoextract: A customizable intermediate level information extraction engine. In *Proceedings of the NAACL 2003 Workshop on Software Engineering and Architecture of Language Technology Systems (SEALTS)*, Edmonton, Canada, 2003.
- [23] P. Srinivasan. Text mining: Generating hypotheses from medline. *JASIST*, 2003.
- [24] B. Taskar, P. Abbeel, and D. Koller. Discriminative probabilistic models for relational data. In *Eighteenth Conference on Uncertainty in Artificial Intelligence (UAI02)*, 2002.