

Automatic Evaluation of Search Ontologies in the Entertainment Domain using Text Classification

Michael Elhadad, David Gabay, and Yael Netzer
(elhadad—gabayd—yaeln)@cs.bgu.ac.il
Department of Computer Science

Ben Gurion University, Israel

Abstract. Information Retrieval (IR) research has recently started addressing the information need of *exploratory search*, where the searcher may be unfamiliar with the domain or not have decided what is the goal of his query. A popular tool to support exploratory search is the use of *faceted search*. The implementation of faceted search requires that documents be annotated by metadata in the form of attributes and hierarchical categories. In many applications, the metadata is maintained manually, in the form of a search ontology. Recent work has also investigated methods to automatically acquire such metadata from sample documents [1, 2]. In this work, we propose a new method to automatically evaluate the quality of such a search ontology.

Our method relies on mapping ontology instances to textual documents. On the basis of this mapping, we evaluate the adequacy of ontology relations by measuring their classification potential over the textual documents. This data-driven method provides concrete feedback to ontology maintainers and a quantitative estimation of the functional adequacy of the ontology relations towards search experience improvement. We specifically evaluate whether an ontology relation can help the search engine support exploratory search in the form of effective facets.

We test this ontology evaluation method on an ontology in the Movies domain, that has been acquired automatically from the integration of multiple semi-structured and textual data sources (*e.g.*, IMDb and Wikipedia). We automatically construct a domain corpus from a set of movie instances by crawling the Web for movie reviews (both professional and user reviews). The 1-1 relation between textual documents (reviews) and movie instances in the ontology enables us to translate ontology relations into text classes. We verify that the text classifiers induced by key ontology relations (genre, keywords, actors) achieve high performance and exploit the properties of the learned text classifiers to provide concrete feedback on the ontology.

The proposed ontology evaluation method is general: it only relies on the possibility to automatically align textual documents to ontology instances.

1 Introduction

Faceted search is a popular interaction method that allows users to navigate through complex unfamiliar data [3, 4]. Faceted search has been widely used in recent e-commerce sites and spread by software providers such as Endeca. Faceted search supports the information need of *exploratory search* [5]. Exploratory search corresponds to a shift in IR from focus on navigational queries and document ranking to the higher level goals of content extraction, user goal recognition and content aggregation [6].

When using a faceted search interface, a user first submits a general query, and then navigates through several facet hierarchies that describe the result set. By navigating through the facet hierarchies, the user interactively refines his query and discovers further facets that may guide him through the document repository and help him discover relations he did not know when first issuing his query. Typically, faceted search applications show possible refinements to the current query together with the number of search results that the refinement would bring in focus. These counts provide helpful quantitative feedback to the user and guide him when exploring the data.

The effectiveness of this interaction strategy relies heavily on the quality of the metadata associated with the documents: which attributes (facets) describe the documents, and the hierarchy of values associated to each facet. In many applications such as e-commerce Web sites, information architects carefully prepare this metadata in the form of a search ontology [4, 7, 8]. Some researchers have also investigated automatic and semi-automatic methods to construct such a search ontology by mining information in the document repository and related resources [1, 2].

Both when manually maintaining the search ontology or acquiring it automatically, there is a pressing need to evaluate its quality. In this work, we present a new method to evaluate a *search ontology* [9]. We tested our method specifically in the entertainment domain: a semantic search engine enables users to search for movies and songs recommendations. The search engine relies on an explicit internal ontology of the domain, which captures a structured representation of objects (movies, actors, directors...). The ontology is acquired and maintained semi-automatically from semi-structured resources (such as IMDb and Wikipedia). The ontology supports improved search experience at different stages: content indexing, query interpretation, search result ranking and presentation (faceted search, aggregated search result presentation and search result summarization).

We focus in this paper specifically on evaluating the quality of the ontology as it impacts the search process. As noted by [10], one can distinguish ontology evaluation methods at three levels: structural (measure properties of the ontology viewed as a formal graph), usability (how is the ontology accessed - through API or search tools, versioned, annotated and licensed) and functional (which services does the ontology deliver to applications). The method we present addresses *functional evaluation*, that is, we investigate how one can measure the adequacy of an ontology to support a semantic search engine.

As part of this functional evaluation, we distinguish two forms of information needs expressed by users: fact finding (the user expects to retrieve a precise set of results or to navigate to a specific movie), and exploratory search (the user seeks recommendations for several movies according to non-specific requirements). The ontology provides services to the application for both types of information needs, but in this paper, we focus on support for exploratory search.

The key idea of our evaluation method is that one can evaluate the functional adequacy of an ontology by investigating a corpus of textual documents anchored to the ontology. The textual documents are collected automatically and associated to ontology instances (outside of the document repository covered by the search engine). Hypotheses about the ontology can then be transformed into classification tests on the corpus.

The rest of the paper is organized as follows: we first review previous work in ontology evaluation and ontology-based information retrieval (ObIR). We then present our ontology evaluation method and a set of experiments we ran to evaluate the functional adequacy of our ontology in the entertainment domain. The experiments validate the adequacy of the specific ontology acquired as part of our semantic engine for exploratory search, and provide specific, concrete indications on how to improve the ontology.

2 Dimensions of Ontology Evaluation

Evaluation of ontologies is designed and performed according to two main scenarios: assessing the quality of an ontology by its developers and ranking ontologies in order to choose the most suitable one for a particular task.

As a general task, evaluation of ontologies is complicated, since ontologies vary in their domain, size, purpose, language and more. Therefore, it is not possible to define a general ontology evaluation paradigm. In addition, the ontology evaluation process depends on the way the ontology was constructed: ontologies may be constructed manually, by scholars or domain experts, or may be the product of an automatic or semi-automatic process. In most applications, ontology quality is best measured in terms of cost/profit effectiveness.

Ontology evaluation can focus on one or more of the following dimensions:

- *Functionality (task-based)*: measure how well an ontology serves its purpose as part of a larger application;
- *Usability*: assess the pragmatic aspects of the ontology, *i.e.*, metadata and annotation [11];
- *Structural evaluation*: identify structural properties of the ontology viewed as a graph-like artefact [12].

It is also useful to distinguish between *extrinsic* and *intrinsic* evaluation methods. Extrinsic evaluation requires either external information in order to evaluate qualities of the ontology, such as a corpus that represents the domain knowledge (data-driven evaluation), expert opinion, or it requires a particular task which defines the context of the evaluation. Intrinsic evaluation reflects the

quality of the ontology as a standalone body of knowledge. Naturally, intrinsic evaluation reflects mostly the structural properties of the ontology.

The method we propose is extrinsic and task informed. We assess the adequacy of an ontology to support exploratory search (in the form of faceted search) by constructing a set of textual documents derived from the ontology and testing hypotheses on this dataset. This methodology is data-driven and automatically provides concrete feedback to the ontology maintainer.

3 Search Ontologies

The usage of an ontology in our project is motivated by the wish to improve the search experience, *i.e.*, we are interested in evaluating a search ontology as defined in the scope of ObIR (Ontology-based Information Retrieval).

The notion of semantic search refers to search techniques which go beyond the mere appearance of query words in possibly relevant documents, and aims to capture a deeper representation of the searched space and the knowledge embedded in it. Although search is widely used in the Internet, user satisfaction studies indicate general user dissatisfaction about irrelevant search results (low precision) or about obtaining too many results. The usage of an ontology can better support users expectations, at the cost of restricting the scope of a search engine to a specific domain. In our case, we investigate the entertainment domain. For such limited scope search, semantic technology helps the engine find more relevant documents by using links among concepts (*e.g.*, movies with the same actor, similar plot), cluster results along semantic attributes to improve navigation (faceted search)[13], and for conceptual indexing (search for “spy” and get “james bond”) [2].

In order to define the evaluation of a search ontology, we refer first to distinct types of search, which represent different types of information needs (following [14][15]):

- **Fact finding:** a precise set of results is requested. The number of retrieved documents is expected to be small (for instance, a specific movie in the entertainment domain). This may correspond to a “return visit” to a site or a short search session.
- **Exploration:** the user’s need is to obtain a general understanding of the search topic: high precision or recall is not required. For instance, the user explores a movies repository to find interesting movies according to his current mood or similarity with known movies.[5].
- **Comprehensive search:** the task is to find as many documents as possible on a given topic (high precision and recall), and to organize the resulting set in a synthetic manner. This task is also called “briefing”. Faceted search is also highly effective to support such a task.

According to these information need distinctions, Strasunskas and Tomassen [15] propose a set of evaluation measures for a search ontology:

- Generic quality evaluation: checks that the ontology is syntactically correct and that it is closely related to the domain.
- Search task fitness: a different measure is applied for each search task. Measures are taken with respect to a cluster of concepts. Fact-finding fitness for a cluster of concepts is a function of the number of instances, properties and data types of all concepts in the cluster. *Exploratory search fitness* is a function of the number of subclasses, and *Comprehensive search fitness* is a function of the number of object properties, sub- and super-classes and siblings. (In all cases, the numbers are divided by the number of concepts in the cluster).
- Search enhancement capability measures how useful the ontology is for query expansions, which improve recall and precision. Recall enhancement capability is a function of the number of labels, equivalent classes, intersections and unions of concepts in a given cluster. Precision enhancement capability is a function of the number of all OWL set operations [16], and of the number data and object properties of concepts in a given cluster.

Such metrics are useful to evaluate ontologies in the same sense that code complexity metrics are useful when developing software or readability index when assessing the quality of text. They correspond to what we call intrinsic measures above. These metrics capture the intuition that the search ontology properly supports the operation of a search engine. But these measures do not provide concrete feedback on the functional adequacy of the ontology to the domain. To illustrate the limitations of such intrinsic measures, it is possible to design an ontology to obtain high scores on all metrics with no knowledge of the domain, in a completely artificial manner, by optimizing the distribution of ontology instances across classes. To reuse the software development analogy, code complexity measures are useful to identify “bad code” (functions that are too long for example), but they do not help to assess the correctness or robustness of the code.

Beyond such metrics, we wish to define functional quality criteria for search ontologies. Gulla et al. [2] define the following desirable properties in a search ontology:

- Concept familiarity: the terminology introduced by the ontology is strongly connected to users terms in search queries.
- Document discrimination: the concept granularity in the ontology is compatible with the granularity used in users’ queries. This granularity compatibility allows good grouping of the search results according to the ontology concept hierarchy.
- Query formulation: the depth of the hierarchy in the ontology and the complexity and length of user queries should be compatible.
- Domain volatility: the ontology should be robust in the presence of frequent updates.

This classification of functional quality criteria is conceptually useful, but it does not provide a methodology or concrete tools to evaluate a given ontology. This is the task we address in this paper.

The evaluation methodology we introduce relies on the fact that given an ontology instance (in our domain, a movie), we can automatically retrieve large quantities of textual documents (movie reviews) associated to the instance. On the basis of this automatically acquired textual corpus, we can perform automatic linguistic analysis that determines whether the ontology reflects the information we mine in the texts.

Note that we focus on evaluating the ontology itself and its adequacy to the domain as a search ontology. However, we do not simulate the search process or measure specifically how the ontology affects steps in search operation (such as indexing, query expansion, result set clustering). Accordingly, the evaluation we suggest, although informed by the task (*i.e.*, we specifically evaluate a search ontology), is not a task-based evaluation (*i.e.*, we do not evaluate the ontology on a search benchmark).

4 Experimental Settings: An Ontology for Semantic Search in the Entertainment Domain

We illustrate our ontology evaluation method in the context of the entertainment domain. We first describe quantitatively the experiments we have run. Our project involves the semi-automatic acquisition of an ontology in the movies domain from semi-structured data sources (IMDb, Wikipedia and other similar sources). The objective of our project is to support exploratory search over a set of documents describing movies, actors and related information in the domain.

We first report on **intrinsic evaluation** metrics over the ontology we have been assessing: number of instances, relations, density. Such measures are domain-independent. Interpretation of these measures is eventually task-oriented: we compare the metrics with those established on “high-quality ontologies” in other domains. We use for this purpose the paradigm of OntoQA [17]. Following the definition of a *search ontology*, the ontology is not expected to have a deep hierarchical structure and complex (dense) relations. The basic metrics are illustrated in Table 1. Additional metrics (instance density, relation density) confirm the expectation that the search ontology we assess has a wide and shallow structure.

Table 1. Basic Measures of the Movies Ontology

Classes	33
Class instances	351,066
Relations	27
Relation instances	19
Movies	8,446
Persons	116,770

Extrinsic evaluation considers the two main search types we identified as our target scenario: fact finding and exploratory. In the first scenario, fact-finding

search, the user seeks precise results and knows what she should get, the main services expected from the ontology are:

- Produce high precision results and wide coverage for terms used in the queries
- Provide Named entity recognition functionality to allow fuzzy string matching and identify terminological variations.
- Identify anchors, *i.e.*, minimal facts that identify a movie (for example, its title, publication year, main actors, main keywords).

For the second scenario, exploratory search, precision and recall cannot be measured since the user does not know apriori what he expects to get. Different criteria have been proposed to assess the quality of an exploratory search system [18]. As mentioned above, we do not attempt a full task-based evaluation. We identify exact quality criteria for exploratory search that enable specific ways through which the ontology can improve the user experience. The services expected from the ontology are:

- Cluster instances by similarity to address the need to present large result-sets.
- Present result-sets using a faceted search GUI to provide efficient browsing and query refinement.
- Identify paths of exploration through which movies are identified (period, genre, actors, ...) to structure the sequence of queries.

Our task is to assess the adequacy of a specific ontology to provide the services listed above. To address this task, we adopt a corpus-based method: assume we have a corpus of textual documents associated to ontology instances. For example, for each movie instance in our ontology, we have a collection of texts. Our evaluation method translates tests on the ontology into tests on such an aligned textual corpus. We present next two specific tests illustrating this approach – to assess the ontology coverage and its classification adequacy.

5 Corpus-Anchored Ontology Evaluation

The first step of our method is to construct a corpus of documents aligned to the ontology instances. In our domain, we construct such a corpus automatically by mining movie reviews from the Web. We collected both professional, edited reviews taken from Robert Ebert’s Web site¹ and additional professional and users reviews published in the Metacritic Web site² and 13 similar Web sources. The key metadata we collect for each document is a unique identifier indicating to which movie the text is associated. The corpus we constructed for these experiments contains 11,706 reviews (of 3,146 movies). It contains 8.7M words, with an average of 749 words per review.

¹ <http://rogerebert.suntimes.com>

² <http://www.metacritic.com>

5.1 Assessing the Ontology Coverage

To assess the fitness of our ontology to support fact-finding search, we measured the named-entity coverage of the ontology, using the constructed text corpus as reference.

We first gathered a collection of potential named-entity labels in the corpus. In professional reviews, named entities are generally marked in the *html* source. Users' reviews are not edited nor formatted. For such reviews, we relied on the OpenCalais³ named entity recognizer to tag named entities in the corpus.

We then extracted all person names from the textual corpus and searched the labels for each entity in the ontology.

Results show that 74% of the named-entity that appear in professional reviews appear in our ontology. For user reviews (non-edited), the figure is 50%.

The main reasons for mismatches lay in orthography variations (such as accents or transliteration differences), mention of people not related to movie and aliasing or spelling variations (mostly in users reviews). We conclude that the coverage of people's names in ontology is satisfactory; however this test did not take into account variations in names and spelling that are expected.

To investigate terminological variation, we measured the ambiguity level of named-entity labels. By ambiguity, we refer to the possibility that a single name refers to more than one ontology instance. We also measured the level of terminological variation for each ontology instance – that is, given a single ontology instance (*e.g.*, an actor), how many variations of its name are found in the corpus. To identify variations in the text, we used the StringMetrics similarity matching library⁴. We experimented with the Levenstein, Jaro-Winkler and q-gram similarity measures. For example, using such similarity measures, we could match “Bill Jackson” with “William Jackson”.

For a version of the ontology that includes 117,556 instances referring to persons, taking into account surnames only, we found that 83% of the names are ambiguous. There are 18.57 variations on average for each ontology instance.

This simple exercise indicates how a textual corpus aligned with the ontology and mature language technology (named-entity recognition and flexible string similarity methods) allows us to measure a complex property of the ontology. This evaluation does not only provide a score for the ontology. It also indicates which specific named entities are used in the corpus, how often, which confusions can be expected when disambiguating query terms and how to specifically improve the terminology-related services provided by the ontology.

In the next section, we demonstrate how the more complex task of measuring the clustering adequacy of the ontology can also be assessed using text classification techniques.

³ <http://www.opencalais.com/>

⁴ <http://www.dcs.shef.ac.uk/sam/stringmetrics.html>

5.2 Assessing the Classification Fitness of an Ontology

As discussed above, the fitness of the ontology to support exploratory search is a function of the number of subclasses. We take this definition a step forward: the number of subclasses is valid if it produces a balanced view of the world domain (represented by the documents) and if the explicit characteristics of the hierarchy can be identified implicitly in the documents.

An ontology induces a hierarchical classification over its elements. Each class (*e.g.*, actor, genre) may be viewed as a dimension for classification of the texts that represent the domain. The ontology provides effective classification services if it meets two criteria:

- The Ontology classification is **useful** if the induced classification is well-balanced, enabling the user explore the dataset in an efficient manner (for exploratory purposes).
- The Ontology classification is **adequate** if the classification induced by the ontology is valid with respect to the domain, which is represented by texts.

Accordingly, we formulate the following hypothesis:

Hypothesis *If* the ontology indicates that some movies are “clustered” according to one of the dimensions, *then* documents associated to these movies should also be found to be associated by a text-classification engine that has been trained on the classification induced by the ontology.

The general procedure we performed to test this hypothesis is the following:

- Step 1:** Choose a dimension to test (we have tested genre, actors and keywords).
- Step 2:** Induce a set of categories (subsets of movies). The subclasses of this dimension and the films instantiated under each subclass defines a clustering of the movies. For example, if we evaluate the “genre” dimension, we cluster movies according to their genre property. In our ontology, this produces about 30 classes of movies (one for each genre value).
- Step 3:** Gather texts (from the reviews corpus, texts that were not used in the acquisition process of the ontology) related to these movies and form a collection (Text_{ij} , movie_i).
- Step 4:** Train a classifier on a subset of the texts (Text_{ik} , movie_i , category_i) where category_i is the category induced by the ontology.
- Step 5:** Test the trained classifier on withheld data (Text_{ij} , movie_i) and compute accuracy, precision and recall with respect to the category.

Hypothesis : Adequate classes yield high accuracy and F-measure on an instance-aligned corpus.

5.3 Parameters

There are several reasonable options to perform the text classification task in Step 4 above, with different methods of text representation and with different classifiers.

For text representation, we viewed texts as “bag of words”, *i.e.*, as unigrams, and represented each text as a Boolean vector in which each coordinate indicates the existence, or lack of existence, of a string in the text. We tested a few options of pre-processing on the texts and of selecting the features (the strings that we take into account when representing the text): with and without stemming⁵ and with and without filtering noise words; selecting features using Mutual Information (MI), or using TF/IDF; and with different numbers of features top 300 or 1000.

Mutual Information-based feature selection is inspired by [19] which shows that this method yields best results on text categorization by topic on a standard News corpus.

The feature selection methods we used are as follows: in TF/IDF, words with the highest values were chosen as features, for the entire corpus. In MI, the features with the highest mutual information associated with the class were chosen (a different set of features is used for every class).

For the classifying task, we used two methods: Support Vector Machines (SVM) (linear and quadratic) and Multinomial Naïve Bayes (MNB) as implemented in the Weka toolkit [20].

5.4 Results

We applied the classification procedure to the classification induced by the genre dimension. The classifiers were trained on the reviews corpus. We performed 5-fold cross-validation on the corpus.

The best text representation was established by testing the genre classifier on the task of classification of one class against all.

16 different experimental settings were tested:

- TF/IDF vs. MI.
- Vectors of size 300 vs. 1000 features.
- Stemmed words vs. Raw.
- Noise words filtered vs. no filtering.

For each possibility we tested both SVM and Naïve Bayes as classifiers.

Classification by Genre Genres, according to IMDb.com are defined to be “simply a categorization of certain types of art based upon their style, form, or content. Most movies can easily be described with certain umbrella terms, such as Westerns, dramas, or comedies”. The tested ontology includes 23 genre subclasses.

⁵ We used the classical Porter Stemmer for the experiment

We performed the classification process as described above, and found that the best combination is MI, 300 features, no stemming, noise filtering, and Naïve Bayes as classifier.

There are several ways to measure the performance of a binary classifier. A common measure in natural language processing tasks is the F-measure, that is defined by:

$$\frac{2TP^2}{2TP^2 + TP \times FP + TP \times FN}$$

Where TP is the number of elements in the 'positive' class that were correctly classified, FP is the number of the elements in the 'negative' class, falsely classified as positive, TN is the number of correctly classified negative elements and FN the number of elements in the positive class, classified as negative.

F-measure takes values between 0 (always mistaken) and 1 (always correct). For our task, we found Matthews correlation coefficient (MCC) more suitable. MCC is given by:

$$\frac{TP \times TN - FP \times FN}{\sqrt{(TP + FN)(TP + FP)(TN + FP)(TN + FN)}}$$

The values of MCC range between -1 (always wrong) and 1 (always correct).

In a typical classification scenario, the classifier *itself* is the object of evaluation (typically, several classifiers are compared on the same dataset), and hence the correlation between class size and F-measure is not an issue; in our case, it is the classes, not the classifier, that we want to evaluate. The advantage of MCC over F-measure comes from the fact that unlike F-measure, it is not affected by the "positive" class size. A random binary classifier, that does not consider the content but only the observed probability P of the positive class in the training set, is expected to yield an F-score of P . Its MCC is expected to be 0. Another point in favour of MCC is that it is symmetrical and thus more suitable to the case of classifying pairs of genres, in which there is no natural choice of the "positive" class.

The Average F-Measure is 0.49, and the average MCC - 0.43 (all results shown in Table 2).

The results indicate that some genres are very well defined (sport, family, documentary), while others cannot be recovered by analyzing the text of the reviews (history, war, biography, adult, short).⁶ While these figures provide a first assessment of the quality of each genre category, pair-wise classification provides finer-grained tests of the level to which pairs of genres can be distinguished. A subset of the results showing best and worst cases is shown in Table 3. We report MCC and error rate for these tests.

Pairwise classification was accurate: the overall error rate was less than 12.4%, and for no pair of classes was higher than 38.2%.

⁶ Specifically, the genres of music and musical are derived from the IMDb genres and are apparently confusing.

Table 2. MCC, F-Measure and error rates of classification engine One-vs-All, sorted by MCC

Genre	%Error	F-Measure	MCC
Sport	2.49	0.72	0.71
Family	5.59	0.63	0.60
Documentary	1.01	0.61	0.56
Adventure	11.88	0.63	0.56
Comedy	21.19	0.71	0.56
Thriller	20.06	0.68	0.54
Action	14.00	0.62	0.53
Sci-Fi	7.94	0.55	0.52
Animation	5.68	0.54	0.52
Horror	8.88	0.54	0.51
Fantasy	7.47	0.53	0.50
Music	3.85	0.50	0.48
Drama	22.18	0.84	0.47
Western	2.35	0.43	0.46
Crime	16.96	0.49	0.39
Romance	23.17	0.46	0.31
Musical	1.73	0.27	0.29
Mystery	12.50	0.35	0.28
History	5.30	0.29	0.28
War	5.97	0.26	0.23
Biography	6.25	0.23	0.22
Adult	6.43	0.19	0.18
Short	20.91	0.24	0.18

For comparison, we have tested a Baseline classifier: this was done by creating 25 random classes of 1,000 movies. We performed the same classification procedure. The results showed average F-measure lower than 0.16 and extremely low accuracy. This indicates that the corpus-anchored ontology evaluation method does not capture random patterns of text classification.

There can be overlap between two categories. For example, a movie can belong both to the genres of action and drama. In our experiment, we train and tested only on movies that belong to one genre of the pair being tested, but not the other.

6 The Case of Keywords

In this section, we suggest a methodology to deal with the situation where the ontology contains a vast number of classes which lack hierarchical structure. From the search perspective, this situation is undesirable since it is most likely unbalanced. From the evaluation perspective, applying the evaluation through classification directly on this set will not be effective, because text classification will not be effective on a large number of categories.

Table 3. Pairwise Classification of Genres

Pair	% Error	MCC
Sport - Fantasy	4.76	0.89
Crime - Family	7.60	0.82
Biography - Sci. Fi.	9.96	0.78
Drama - Animation	1.49	0.77
Horror - Thriller	7.15	0.27
Drama - War	12.23	0.23
History - Biography	38.28	0.23
Drama - History	0.51	-0.002

Keywords in IMDb and accordingly in our Entertainment ontology are based on what Szomszor *et al.* [21] call *free-for-all* tagging. Users can add new keywords, which are then moderated to prevent spamming. Overall, there are 10,529 unique keywords in the ontology. The relation between keywords and films is *many-to-many*: there are many keywords per movie, too many or too few (as few as *one*) movies per keyword and many of the keywords may be only weakly connected to the content of the movie.

Take for example the keywords associated with the movie “*Bonnie and Clyde*”: **bank robbery**, **celebrity criminal** may be good search terms; but **old woman**, **joke**, **face slap** or **marriage**, intuitively, do not characterize the desired plot. Bad search terms are either too general or too specific. Over-specific terms might be useful in searching for a particular movie (*e.g.*, “I’m looking for that movie with snakes and planes and forgot its name”) but, most likely, do not help exploratory search.

Hence, for search purposes, we wish to cluster keywords into more manageable and more useful exploratory dimensions. Once we cluster these keywords, we can use the same classification-based evaluation method to test the quality of these clusters.

We applied a clustering method based on Latent Dirichlet Allocation (LDA) [22] to cluster the keywords in our application. LDA is a generative probabilistic model for collections of discrete data. The main assumption of this model is that a document (represented as a bag-of-words) is a mixture of topics, and each word is generated by a topic with some probability. The mixture proportion of topics for each document as well as the topic to word generation probabilities words are learned from a given corpus in an unsupervised process. Three parameters (K , α , β) control the number of topics, the sparsity of the document to topics distribution and the sparsity of the topic to words distributions, respectively.

LDA has been used successfully in many textual data mining applications. In our application, we adopt LDA to learn possible topics from a set of keywords. The topics learned by LDA are clusters of keywords. To this end, we consider the list of keywords linked to a film to form a document, and perform LDA on all such documents.

Examples of the keyword clusters we acquired are shown in Fig.1 and Fig.2. We configured the LDA process to construct 100 clusters. On average, the obtained clusters contained 20 keywords.

6.1 Results

To test the validity of the clusters, we used the same method as described in Section 5.4. We applied the procedure to the classification induced by the genre dimension. The classifiers were trained on the reviews corpus.

As a baseline, we applied a very simple keyword clustering method - based on word similarity (two keywords are attached to the same cluster if they include a common word or words within a small edit distance. With such a simple clustering heuristic, the average F-measure obtained on text classification was extremely low (Average F-Measure = 0.07).

When using the LDA clustering method, results improved significantly:

Average F-measure is 0.48

Average MCC-measure is 0.28

Average error rate 0.26

Table 4 shows the results of the text classification method for some classes. While the classification results are not very high, they are much better than the baseline (range of 0.65 to 0.75 vs. 0.07). The MCC and F-measure values allow us to filter unreliable keyword clusters and to compare the potential of each keyword cluster to help in exploring the dataset.

class	F	MCC
10	0.74	0.27
72	0.71	0.25
83	0.70	0.21
24	0.66	0.19
90	0.66	0.38
71	0.66	0.39
45	0.66	0.21
36	0.65	0.43

Table 4. One-to-all Classification of Keywords

7 Conclusion and Future Work

We have presented a concrete ontology evaluation method based on the usage of a corpus of textual documents aligned with ontology instances. We have demonstrated how to operate such an evaluation in the case of an ontology in the entertainment domain used to improve a semantic search engine.

10 flashback-sequence mother-son-relationship hotel father-son-relationship restaurant face-slap premarital-sex bar car-accident hospital funeral los-angeles-california friendship drunk-scene beach blockbuster profanity title-spoken-by-character narration cemetery

72 based-on-novel character-name-in-title independent-film number-in-title acronym-in-title lost-film hobo kilt scottish-accent entire-title-is-capitalized-acronym clock-watcher party-lifestyle team-owner essex-wife wags team-captain wives-and-girlfriends aids once-upon-a-time-in-the-title m-a-s-h

83 based-on-novel character-name-in-title independent-film circus carnival clown amusement-park criminal-justice number-in-title acrobat gypsy midget roller-coaster fortune-teller sideshow hypnotism trapeze elephant carousel ferris-wheel

Fig. 1. LDA sets with highest F-measure

87 kids-and-family cartoon looney-tunes merrie-melodies australia bugs-bunny australian popeye porky-pig daffy-duck part-live-action chicken william-tell-overture woody-woodpecker sylvester pig screen-song breaking-the-fourth-wall duck anvil

86 boxing baseball sport soccer american-football football basketball coach boxer training olympics golf athlete college competition dandy stadium reverse-the-polarity-of-the-neutron-flow early-sound locker-room

36 native-american murder horse sheriff cowboy revenge gold outlaw spaghetti-western saloon ranch actor-shares-first-name-with-character bank-robbery desert cattle bandit texas shootout stagecoach arizona

Fig. 2. LDA sets with highest MCC-measure

We have first constructed an ontology-aligned textual corpus by developing a Web crawler of movie reviews. On the basis of this dataset, our first experiment measures the adequacy of the ontology to support fact-finding search. We have found specifically that our ontology has wide coverage but lacks support for ambiguity resolution and terminological variation handling. We use human-language technology to translate hypothesis on the ontology coverage into measures of properties of the textual

Our second experiment measures the adequacy of the ontology to support exploratory search. We have formulated hypotheses that capture the quality criteria of an exploratory search system, and tested these hypotheses on our ontology-aligned textual corpus. Specifically when testing the classification adequacy of our ontology along the “genre” dimension, we found that most of the genres in the ontology induce high-quality text classifiers - but some, such as sport and music) do not induce appropriate classifiers. This method provides specific feedback to the ontology maintainer.

Our tests indicate that the proposed text classification method provides useful feedback to information architects, even when dealing with extremely unstructured data such as clouds of unedited keywords.

Acknowledgments

This research is supported by Deutsche Telekom at the BGU T-Lab laboratories of Ben-Gurion University.

References

1. Stoica, E., Hearst, M.A., Richardson, M.: Automating creation of hierarchical faceted metadata structures. In Sidner, C.L., Schultz, T., Stone, M., Zhai, C., eds.: HLT-NAACL, The Association for Computational Linguistics (2007) 244–251
2. Gulla, J.A., Borch, H., Ingvaldsen, J.: Ontology learning for search applications. In: Proceedings of the 6th International Conference on Ontologies, Databases and Applications of Semantics. (2007)
3. Tunkelang, D.: Faceted Search. Synthesis Lectures on Information Concepts, Retrieval, and Services. Morgan and Claypool Publishers (2009)
4. Marti A. Hearst, Preston Smalley, C.C.: Faceted metadata for information architecture and search. (CHI 2006 Course)
5. White, R.W., Roth, R.A.: Exploratory Search: Beyond the Query-Response Paradigm. Synthesis Lectures on Information Concepts, Retrieval, and Services. Morgan and Claypool Publishers (2009)
6. Baeza-Yates, R., Ciaramita, M., Mika, P., Zaragoza, H.: Towards semantic search. In: Natural Language and Information Systems. (2008) 4–11
7. Hearst, M.A.: Uis for faceted navigation: Recent advances and remaining open problems. In: in the Workshop on Computer Interaction and Information Retrieval, HCIR 2008. (2008)
8. Moritz Stefaner, Sebastian Ferre, S.P.J.K., Zhang, Y.: User Interface Design. In: Dynamic Taxonomies and Faceted Search : Theory, Practice, and Experience. Volume 25 of The Information Retrieval Series. Springer (2009)
9. Burkhardt, F., Gulla, J.A., Liu, J., Weiss, C., Zhou, J.: Semi automatic ontology engineering in business applications. In: Proceedings of the 3rd International AST Workshop – Applications of Semantic Technologies. (2008)
10. Gangemi, A., Catenacci, C., Ciaramita, M., Lehmann, J.: Modelling ontology evaluation and validation. In: ESWC 2006. LNCS. Volume 4011. Springer, Heidelberg (2006) 140–154
11. Gomez-Perez, A.: Evaluation of ontologies. International Journal of Intelligent Systems **16** (2001) 391–409
12. Alani, H., Brewster, C.: Ontology ranking based on the analysis of concept structures. In: Proceedings of the 3rd International Conference on Knowledge Capture (K-Cap), Banff, Canada (2005) 51–58
13. Hearst, M.A.: Search User Interfaces. Cambridge University Press (2009)
14. Aula, A.: Query formulation in web information search. In: Proceedings of IADIS International Conference WWW/Internet. (2003) 403–410
15. Strasunskas, D., Tomassen, S.: Empirical in-sights on a value of ontology quality in ontology-driven web search. In: On the Move to Meaningful Internet Systems 2008: CoopIS, DOA, ODBASE, GADA, and IS. Springer-Verlag, Monterrey, Mexico (2008)
16. McGuinness, D.L., van Harmelen, F.: Owl web ontology language overview. Technical Report REC-owl-features-20040210, W3C (2004)
17. Tartir, S., Arpinar, I., Moore, M., Sheth, A., Aleman-Meza, B.: Ontoqa: Metric-based ontology quality analysis. In: Proceedings of Workshop on Knowledge Acquisition, Autonomous, Semantically Heterogeneous Data and Knowledge Sources. (2005)
18. White, R.W., Muresan, G., Marchionini, G., eds.: ACM SIGIR Workshop on “Evaluating Exploratory Search Systems”, Seattle (2006)

19. Dumais, S., Platt, J., Heckerman, D., Sahami, M.: Inductive learning algorithms and representations for text categorization. In: Proceedings of the Seventh international Conference on information and Knowledge Management, Bethesda, Maryland (1998) 2–7
20. Witten, I.H., Frank, E.: Data Mining: Practical machine learning tools and techniques. Second edn. Morgan Kaufmann, San Francisco (2005)
21. Szomszor, M., Cattuto, C., Alani, H., O’Hara, K., Baldassarri, A., Loreto, V., Servedio, V.D.: Folksonomies, the semantic web, and movie recommendation. In: Proceedings of the 4th European Semantic Web Conference, Bridging the Gap between Semantic Web and Web 2.0, Innsbruck, Austria (2007)
22. Blei, D.M., Ng, A.Y., Jordan, M.: Latent dirichlet allocation. *Journal of Machine Learning Research* **3** (2003) 993–1022