

# A Tag Means a Lot Than It is in a Folksonomic System

Ho Keung Tsoi  
Department of Computer Science  
Hong Kong Baptist University  
Kowloon Tong, Hong Kong  
hktsoi@comp.hkbu.edu.hk

## Abstract

*Folksonomic system allows users to use tags to describe items, these tags do not just exist in the form of textual description, they actually bear more meaning underneath, such as user preference. In this paper, we first show the distribution of preferences and semantic categories across a folksonomic system, and then develop a hybrid design to cope with the cold-start problem.*

*We speculate the semantic categories formed in user perspective and item perspective in a folksonomic system are different. They represent different preferences and meaning and are believed to be crucial in recommender algorithm design. Through a dimensionality reduction technique, the Latent Dirichlet Allocation, we demonstrate our speculation is correct.*

*In this regards, we design a hybrid strategy accordingly. Our strategy consists of two stages. First we enhance the user's tag profile by WordNet, so as to provide more sound information for later use. The second stage is to find a winning cluster, that maximizes the user's preference. The evaluation reveals our design outperform other existing approaches. This verifies our idea of leveraging various users' interests in the recommendation process, is capable of yielding a better result. Since this strategy can stimulate user's preference, it can enhance user experience as well as solving the cold-start problem.*

## 1 Introduction

How do you judge the person at first sight? By his appearance or by his dressing? You do not have too much information about this individual in the first meet, and hence you cannot do too much. This is also the case in providing recommendation to novel users, the lack of prior information hinder us from understanding the new comers.

Folksonomic system [23] has become popular and growing rapidly in recent years. A folksonomic system allows

Internet users to assign keywords – so called tags, to annotate resources. The role of these tags is to help users to manage, navigate and explore resources. Living examples of this system include Flickr<sup>1</sup>, Last.fm<sup>2</sup> or Delicious<sup>3</sup>.

Different analysis of tagging pattern and motivation in folksonomic systems have been done by peer researches. They show the types of tags used in the social tagging process can be classified in the categories of *Personal*, *Factual* and *Subjective*[3], and a semantic space of social tags will gradually be evolved from the folksonomic system, this semantic classification of tags formed by social tagging has some self-organizing characteristic[20]. As for the motivation of applying tags, study has shown that it is driven by the purpose of sharing and personal information management[15].

To this end, various tag recommendation algorithms have been deviated from these folksonomic systems. Tag recommendation can facilitate users to browse and search resources, as well as to manage and retrieve their own resources. Contemporary tag recommendation algorithms include cluster-based[19], memory-based[1], content-based[29] or collaborative filtering[27] approach. These approaches rely heavily on available prior information, such as rating, to find a matching relevant candidate to return to user. When a novel user is encountered that prior information is yet available, the recommender system struggles to generate a recommendation. This is know as the "cold-start" problem [18].

However, stimulating user's preference should be the primary goal of recommendation. If a user doesn't interested in a system, he will not use the system anymore, not to mention to accomplish the above tasks. But the aforementioned algorithms either suffering from the cold-start problem, or overlook the essence of user's preference, we hope to seek a solution that balances the two sides. In this paper, we design a cluster-based algorithm to give solution to the cold-start

---

<sup>1</sup><http://www.flickr.com>

<sup>2</sup><http://www.lastfm.com>

<sup>3</sup><http://delicious.com>

problem, while reserving user preference in the process.

The remainder of this paper is organized as follows. Section 2 is the literature review. In Section 3 we show the dataset we use throughout this paper. In Section 4, we present an analysis, and state the difference of semantic categories built in user- and item-space. In Section 5, a brief introduction of algorithms to be experimented and our strategy is presented. A comparison and evaluation of these algorithms are examined in Section 6. In Section 7 we summarize and discuss the algorithm and future works.

## 2 LITERATURE REVIEW

A folksonomy can be described as a four-tuple: a set of users,  $U$ ; a set of resources,  $R$ ; a set of tags,  $T$ ; and a set of assignments,  $A$ . The data in the folksonomy is denoted as  $D$  and is defined as:  $D = \langle U, R, T, A \rangle$ . The assignments,  $A$ , are represented as a set of triples containing a user, tag and resource defined as:  $A \subseteq \{ \langle u, r, t \rangle : u \in U, r \in R, t \in T \}$ . Therefore a folksonomy can be regarded as a tripartite hyper-graph with users, tags, and resources represented as nodes and the assignments represented as hyper-edges connecting one user, one tag and one resource[24].

As such, Graph Theory has always been adopted to provide recommendation in folksonomy. A graph-based ranking algorithm for interrelated multi-type object is proposed[14]. The task of Personalized Tag Recommendation is modeled as a "query and ranking" problem. When a user issues a tagging request, both the document and the user are treated as a part of the query. This algorithm ranks tags by considering both relevance to the document and preference of the user. Likewise, some authors are inspired by the algorithm *PageRank*, and use *authoritative* tags to enrich user query[9]. Each folksonomic user is maintained a profile in their approach, as well as a knowledge base consisting of two graphs called *Tag Resource Graph* and *Tag User Graph*. These graphs register the tags exploited in the folksonomy and the way they label involved resources, or the way they are registered in the user profiles. When user submit a query, *authoritative* tags are suggested to user and enrich user profile automatically. FolkRank [17], a enhancement of PageRank-like algorithm that takes into account the structure of folksonomies to search in the system.

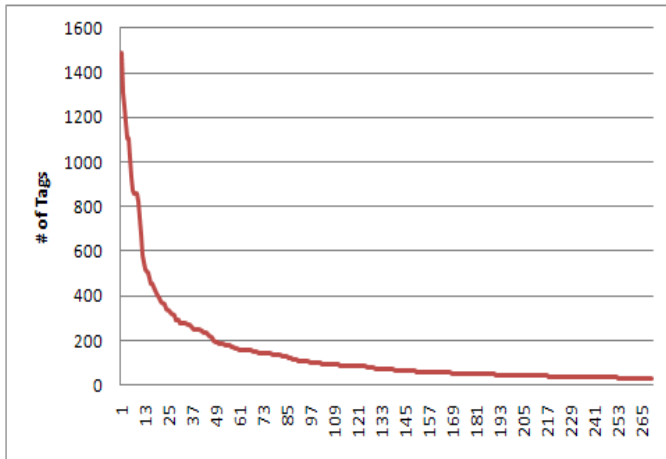
The idea of embedding content information in the recommendation process is not novel. The authors in [29] describe a movie recommendation system built purely on the keywords assigned to movies via collaborative tagging. Recommendations for the active user are produced by algorithms based on the similarity between the keywords of a movie and those of the tag-clouds of movies the user rates. According to [8], content-based recommender not only can recommend items, but also be used to infer user interests. They use a multivariate Poisson model for naive Bayes text

classification adapted to infer user profiles from both static content, as in classical content-based recommender, and tags provided by users to freely annotate items. The benefit of using content information includes solving the cold-start problem. Researches in [12] propose a probabilistic model for inferring the most probable tags from the text of the book. They combine a Relevance model developed for Information Retrieval, and the Collaborative Filtering approach, to generate tags from the content of books.

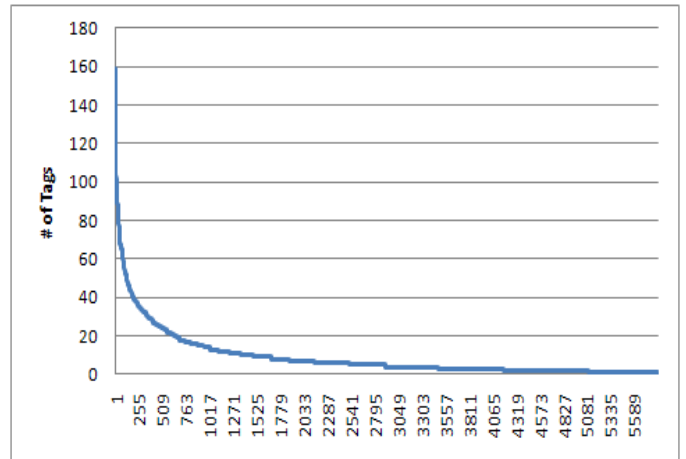
Besides content information, contextual information such as navigational pattern, and browsing behaviors also play a major role in recommendation process. [32] conduct a study to determine which context information sources can predict user's interests effectively. In particular, they evaluate *social*, *historic*, *task*, *collection*, and *user interaction*. Their result shows the context overlaps outperform any isolated source, and suggests designers can improve Website suggestion by these findings. [21] suggests items to users based on inferences made about user interests gleaned from their task environment, such as recently-viewed Web pages or the contents of active desktop applications. Underneath the obvious contextual information like the link structure of the Web, implicit connection in a folksonomic system can link related users together. Authors in [5] formulate user-induced links in collaborative tagging system as follows: if two documents that are maintained in the collection of the same user and/or assigned similar sets of tags can be considered as related from the perspective of the user. They then demonstrate that this kind of induced-link achieves much higher accuracy than existing hyperlinks. Contextual information also includes emotional context and the like. [13] present a SPA system, which elicits user's preference through a rich interaction through highly dynamic environments, networked game for example. This approach is particularly useful in social software systems, as it can easily acquire information via user activities and tasks. As shown, the task of generating recommendation is not limited to the scope of folksonomy.

Due to the textual nature of tags, each tag bears a semantic meaning. This leads to researches focusing on the semantic dimension to produce recommendation. WordNet dictionary [30] and ontologies from open linked data published on the Web [6] are the tools to support this task. To recommend items which are about similar contents, [10] find semantic relations between tags on different semantic sources and calculate their semantic similarity. In addition to applying semantic calculus directly on the tags, clustering tags at a higher level can also be done. [20] use self-organizing map to determine the tags' semantic dimension in social tagging system. This higher representation of tags gives more representative, and can amend the weaknesses of traditional methods such as experts or statistics.

Also because of this very nature of tags, the issue of re-



**Figure 1. Tag assignments distribution of each user**



**Figure 2. Tag assignments distribution of each movie**

dundancy and ambiguity has to be solved. The free annotation in a folksonomic system permits unsupervised tagging, and users can use tag in a way they desire. Consequently, a single tag has many different meanings or results in ambiguity and redundancy in which several tags have the same meaning. However, since traditional evaluation metrics such as Precision and Recall cannot account the effect of these phenomena, [11] use a cluster-based approach to define ambiguity and redundancy and provide evaluation on real world datasets. They show this evaluation strategy can more reveal the utility of tag recommender.

### 3 DATASET

The work of this paper is based on the MovieLens<sup>4</sup> dataset, and the movie's information obtained from IMDB<sup>5</sup>. MovieLens is the movie recommender system maintained by GroupLens Research. Each movie in the dataset has a link led to that movie's description on IMDB. Since we are interested in user tag profile, we trimmed the data as follows to counteract the effect of skewed distribution[28]. Figure 1 and Figure 2 show the distribution of our dataset. They reflect the long-tail phenomena in a folksonomic system – a majority of users/items use very few tags.

In the original dataset, we extracted a set of users who applied at least 30 tags (include duplicated tags). The set of tags belong to these users and the set of movies related to these users are considered. As a result, we have the following data.

For each of the related movie, we additionally crawled

**Table 1. Summary of the data set used**

	MovieLens
# users	271
# distinct tags	6,409
# movies	5,840

the movie's description from the associated IMDB link<sup>6</sup>. Keywords are then extracted from the textual description by comparing against the standard stop-word list.

### 4 THE DIVERGENCE OF SEMANTIC CATEGORIES

There are two approaches to discover semantic categories given a folksonomic system. One is to treat each user as a document, and the list of tags *assigned by* this user as words; alternatively, the list of tags *assigned to* an item can be regarded as a document, and thus the tags are words to describe this item. At this point, we have user-space and item-space. We further define that the semantic categories in user-space reveal the general preferences of users, while those in item-space reveal the objective description of items.

We come to these definitions due to the following observations. In user-space, a user assigns tag based on individual preference. A user might repeatedly use the same tag on different items, for instance. So the bag of tags used by this user can indicate his general preference. Whereas in the item-space, the item's tags are given by more than one user, and therefore those most frequently occurring tags can lighten the effect of individual bias, and hence objective. A common way to examine the semantic categories

<sup>4</sup><http://movielens.umn.edu/>

<sup>5</sup><http://www.imdb.com/>

<sup>6</sup>Example link: <http://akas.imdb.com/title/tt0114709/synopsis>

is through the technique of dimensionality reduction. Algorithms like Latent Semantic Indexing[31], Latent Dirichlet Allocation[7] and Self-Organizing Map[26] are capable of achieving this purpose. In particular, we rely on Latent Dirichlet Allocation (hereafter, LDA) for our subsequent analysis and algorithm development.

#### 4.1 The Analysis

We use LDA to examine the divergence of semantic categories among user-space and item-space. For illustrative purpose, we specify the number of latent topics to be 30. The result can be found in Table 2.

For each topic, we use the top three representative tags to represent the topic, the representativeness of a tag is in turn measured by the *term frequency*[10] w.r.t. the topic. It is seen that the semantic categories found in the two spaces not always agree with each other. For example, the cluster {action, Comedy, Drama} in user-space reflects the preference of watching comedic drama movie; the {Betamax, James Bond, 007} in item-space indicates the 007's series movies. This implies dissimilarity in the two spaces. The result is expectable in the sense that the semantic categories formed in user-space reflect the general preferences of users, and that in the item-space tells the factual dimension of resources. In the subsequent section, we will show how our algorithm addresses these findings.

### 5 TAG RECOMMENDATION

Based on the findings from previous analysis, we design a strategy to account for them. In this section, we present our approach, and show how we will deal with the issues. Then some of the example algorithms are selected and briefly described, so as to provide readers a brief understanding of the contemporary development of tag recommendations.

#### 5.1 Our Approach

Our goal is to provide recommendation to a novel user in a fashion such that maximizes the user's preferences. We reference to the design of existing tag recommendation algorithms, and try to combine their advantages into one. Our idea is to utilize user's tag profile to generate recommendation. By tag profile we mean the tags a user applied. A new user has a tag profile of length one as he first uses the system; this user has a tag profile of length two as he applies the second tag, so on and so forth for instance.

To begin with, we use LDA to discover the semantic categories in both user-space and item-space. Here the item-space contains not only the tags assigned by users, but also the keywords extracted from the movie's description link.

The keywords extraction is done by removing words from the standard stop-word list, and the remaining words become the keywords of that movie. Then the similarity of each cluster in user-space to each cluster in item-space are determined using *symmetric Jaccard coefficient* :

$$sim(C_{i^{th}}^{user}, C_{j^{th}}^{item}) = \frac{|T_i \cap T_j|}{|T_i \cup T_j|}$$

$$T_i \in C_{i^{th}}^{user}$$

$$T_j \in C_{j^{th}}^{item}$$

where  $T$  is the tagset,  $C_{i^{th}}^{user}$  denotes the  $i^{th}$  cluster in user-space, and  $C_{j^{th}}^{item}$  denotes the  $j^{th}$  cluster in item-space. This equivalent to bridge the gap between user-space and item-space. Whereas the traditional LDA concerns only either one, and neglect the other, which lead to the resulting clusters do not cover all the possible semantic categories.

Now, lets get back to the user's tag profile. For a new user, immediately after he has entered the first tag, we enrich his profile to five tags using WordNet [25]. We return  $N$  recommended tags which have the highest similarity value to the querying tag using Wu and Palmer metric[33], where  $N$  depends on the length of the current tag profile. Formally, we enrich the user's tag profile in the initial stage in the following manner:

$$N = \begin{cases} (5 - |Profile|), & \text{if } |Profile| < 5 \\ 0, & \text{if } |Profile| \geq 5 \end{cases}$$

After the initial enrichment, we have a tag profile of length at least five. The tag profile is then used for finding the most relevant cluster in the user-space, and the relevance is defined as the value of overlap between the tag profile and cluster. With the  $C_{i^{th}}^{user}$  located, we return both  $C_{i^{th}}^{user}$  and the most similar cluster in item-space  $C_{j^{th}}^{item}$  found in previous step.

The last step is to decide which cluster to be the final recommendation. Again, we use the maximum overlap for our measurement.

$$WinningCluster = \arg \max(C_{i^{th}}^{user} \cap Profile, C_{j^{th}}^{item} \cap Profile)$$

The rationale behind our strategy is to maintain the tag profile to have certain length, so that we can make use of this piece of information in the recommendation process. Using WordNet, we can enrich the profile with semantically correlated tags. The clusters in the user-space and item-space, on the other hand, represent users' general interests and objective factual information respectively. Choosing among these two groups in the last step captured the importance of user preference, as it has the largest degree of

**Table 2. The Latent Topics found in User-Space and Item-Space**

User-Space	Item-Space
action,Comedy,Drama	serial killer,martial arts,beautiful
boring,PG13,afternoon section	Owned,Crime,dvd
Can't remember,Friday night movie,Didn't finish	Nudity (Full Frontal - Notable),lesbian,Musical
dvd,DIVX,Want	Tumey's DVDs,imdb top 250,black and white
AFI 100,Disney,AFI 100 (Laughs)	library,gay,erlend's DVDs
girlie movie,Hitchcock	Bruce Willis,psychology,ghosts
seen more than once,overrated,James Bond	Oscar (Best Picture),documentary,Oscar (Best Cinematography)
Tumey's DVDs,USA film registry,Tumey's To See Again	based on a book,adapted from:book,Fantasy
less than 300 ratings,avi,violent	70mm,World War II,history
movie to see,National Film Registry,ClearPlay	comic book,holocaust,super-hero
classic,Criterion,history	classic,imdb top 250,National Film Registry
Nudity (Topless),Nudity (Topless - Brief),Nudity (Full Frontal - Notable)	anime,In Netflix queue,Japan
erlend's DVDs,Sven's to see list,based on book	Betamax,James Bond,007
Bibliothek,seen at the cinema,watched 2006	less than 300 ratings,To See,Sven's to see list
aliens,drugs,remake	boring,Johnny Depp,Adventure
70mm,Betamax,DVD-Video	action,aliens,Eric's Dvds
corvallis library,hw drama	dvd-r,library vhs,Scary Movies To See on Halloween
imdb top 250,netflix,oppl	dvd,sci-fi,Futuristmovies.com
on computer,funny,ohsoso	comedy,funny,chick flick
anime,need to own,breakthroughs	drama,biography,christmas
atmospheric,Golden Palm	erlend's DVDs,atmospheric,Criterion
Futuristmovies.com,documentary,space	directorial debut,time travel,seen more than once
owned,adapted from:book,based on a TV show	zombies,movie to see,Angelina Jolie
Johnny Depp,Brad Pitt,Arnold Schwarzenegger	movie to see,Below R,parody
In Netflix queue,Disney,Christmas	Disney,Animation,Pixar
ummarti2006,2.5,cars	remake,Brad Pitt,Based on a TV show
Oscar (Best Picture),psychology,toplist08	Nudity (Topless - Notable),VHS,Jackie Chan
based on a book,directorial debut,black and white	Can't remember,own,based on a play
World War II,Jack Nicholson,Clint Eastwood	ClearPlay,drugs,PG13
own,Eric's Dvds,Ei muista	Nudity (Topless),Nudity (Topless - Brief),netflix

agreement to the user's tag profile. The visual presentation of our idea is presented in Figure 3. The following subsections introduce the traditional approaches in tag recommendation.

## 5.2 Collaborative Filtering

In Collaborative Filtering[1, 27], an object is suggested to a user  $u$  if it was rated as relevant by a group of users having a profile similar to the one of  $u$ . The profile can be established by user's rating, and the relevance is measured by similarity metric such as cosine-similarity. Formally, similarity between users  $i$  and  $j$ , denoted by  $sim(i,j)$  is given by

$$sim(i, j) = \cos(\vec{i}, \vec{j}) = \frac{\vec{i} \cdot \vec{j}}{\|\vec{i}\|_2 * \|\vec{j}\|_2}$$

where  $\cdot$  denotes the dot-product of the two vectors. Variation of such approach includes Pearson Correlation[22], similarity-weighted average of the rating[4] are used for distance metric. This algorithm is intuitive and efficient, but sophisticated readers should spotted the problem of cold-start. Because of the nature of this algorithm relies heavily on user profile, it hardly provides recommendation to novel user or user who doesn't rate.

## 5.3 Association Rule Mining

Association rule mining finds interesting associations and correlation relationships among large set of data items. It has a form  $T_1 \longrightarrow T_2$ , where  $T_1$  and  $T_2$  are items (Tags in our case), this indicates  $T_1$  implies  $T_2$ . Association rules show attribute value conditions that occur frequently together in a given dataset. A typical and widely-used example of association rule mining is Market Basket Analysis[2]. The three key measures for association rules are support, confident and interest. The support is simply the number of transactions that include all items in the antecedent and consequent parts of the rule. i.e., an estimate of the joint probability  $P(item_1, item_2)$ . Confidence is an estimate of the conditional probability  $P(item_1|item_2)$ . Interest (a.k.a. Lift) is the ratio of Confidence to Expected Confidence ( $\frac{P(item_1, item_2)}{P(item_1)P(item_2)}$ ).

In the context of tag recommendation, if many resources with tags  $Tag_1$  are typically also annotated with tags  $Tag_2$ , then a new resource with tags  $Tag_1$  may also be meaningfully annotated with tags  $Tag_2$ [16]. But the nature of skewed distribution [28] of tags in a folksonmic system, prohibited the association rule to yield a better performance. If one prefers to maintain a large coverage of tags, the Confidence (so as the Support) parameters have to be lowered.

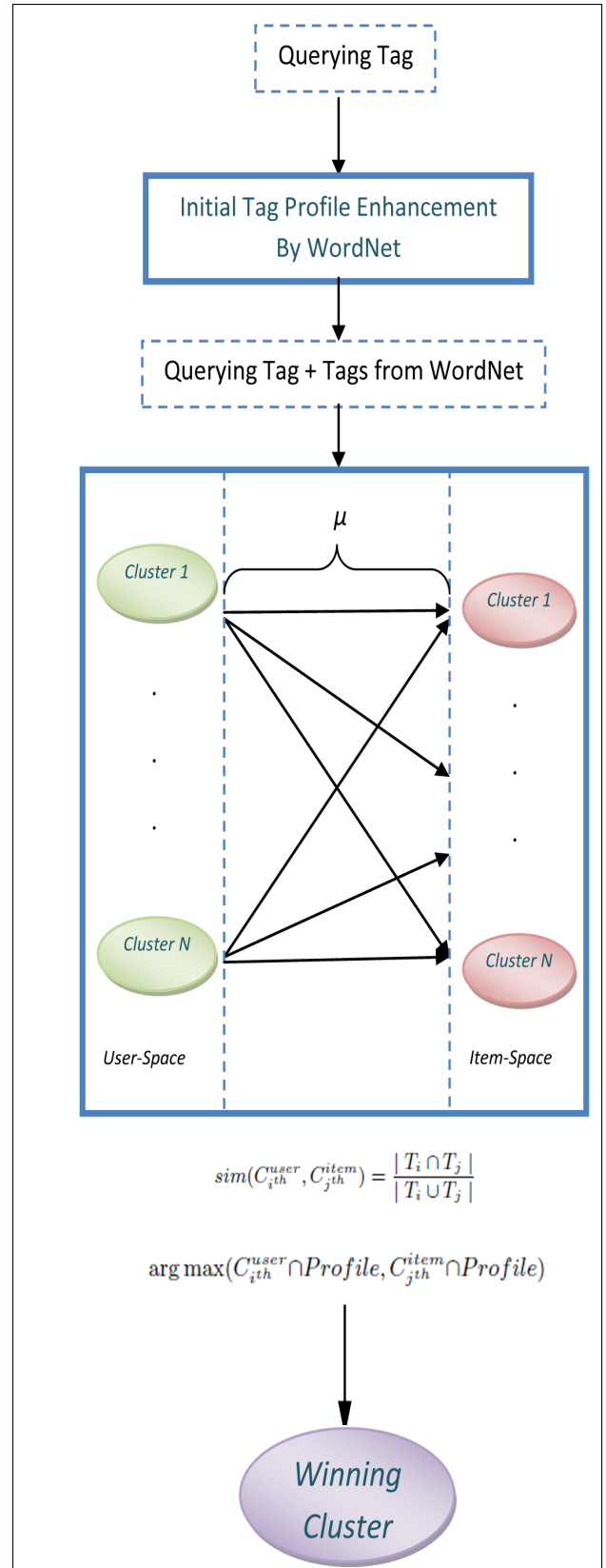


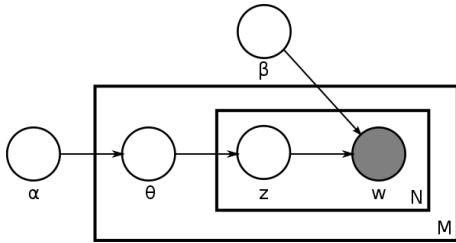
Figure 3. Graphical representation of our strategy

## 5.4 Latent Dirichlet Allocation

Latent Dirichlet allocation (LDA) is a generative probabilistic model of a corpus and it assumes there are  $k$  underlying latent topics. The basic idea is that documents are represented as random mixtures over latent topics, where each topic is characterized by a multinomial distribution over words. Using the terminology in our context, multiple users are annotating resources, and the resulting topics reflect a collaborative shared view of the resource and the tags of the topics reflect a common vocabulary to describe the resource.

The generative process of LDA can be formalized as follows:

1. Choose  $N \sim \text{Poisson}(\xi)$
2. Choose  $\theta \sim \text{Dir}(\alpha)$
3. For each of the  $N$  words  $w_n$ :
  - (a) Choose a topic  $z_n \sim \text{Multinomial}(\theta)$
  - (b) Choose a word  $w_n$  from  $p(w_n | z_n, \beta)$ , a multinomial probability conditioned on the topic  $z_n$ .



**Figure 4. Graphical model representation of LDA, adopted from [7]**

The parameter  $\theta$  indicates the mixing proportion of different topics in a particular resource.  $\alpha$  is the parameter of a Dirichlet distribution that controls how the mixing proportions  $\theta$  vary among different resources.  $\beta$  is the parameter of a set of multinomial distributions, each of them indicates the distribution of tags within a particular topic. Learning a LDA model from a collection of resources  $D = \{t_1, t_2, \dots, t_3\}$  involves finding  $\alpha$  and  $\beta$  that maximize the log likelihood of the data  $l(\alpha, \beta) = \sum_{d=1}^M \log P(w_d | \alpha, \beta)$ . This parameter estimation problem can be solved by the variational EM algorithm.[7]

## 5.5 WordNet

WordNet is an online lexical database designed for use under program control. English nouns, verbs, adjectives,

and adverbs are organized into sets of synonyms, each representing a lexicalized concept [25].

To utilize WordNet as a support tool for the tag recommendation[10], each tag is associated with a semantic knowledge and the recommendation is produced by returning tags with the highest similarity value to the querying tag, where the similarity between tags can be measured by the semantic similarity using the formula proposed by Wu and Palmer [33]. The advantage of this approach is that the recommended tags are lexically correlated to the querying tag, but its downside is that it takes no user preference into account.

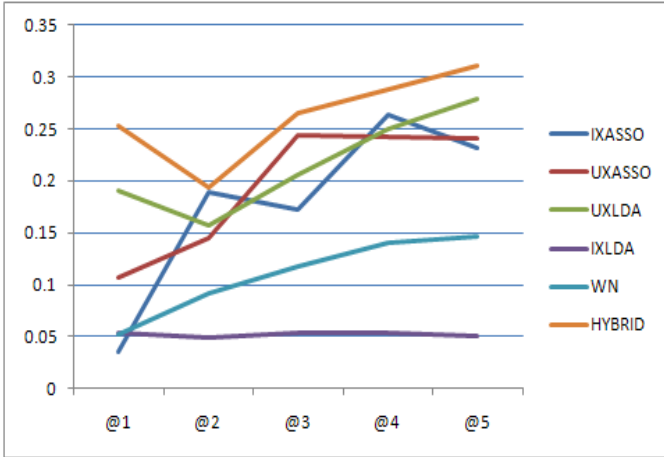
## 6 COMPARISON AND EVALUATION

In particular, we evaluate Association Rule Mining(*UXASSO*, *IXASSO*)[16], Latent Dirichlet Allocation (*UXLDA*, *IXLDA*)[7, 19], WordNet(*WN*)[33], and our hybrid one. The results are summarized in Figure 5 to Figure 6. The standard metrics in Information Retrieval, i.e. Precision, Recall are adopted. The horizontal axis of the graphs depict the number of tags in user profile. For each run, 200 iterations with the same length of tag profile are performed, and the averaged values are presented in the graphs.

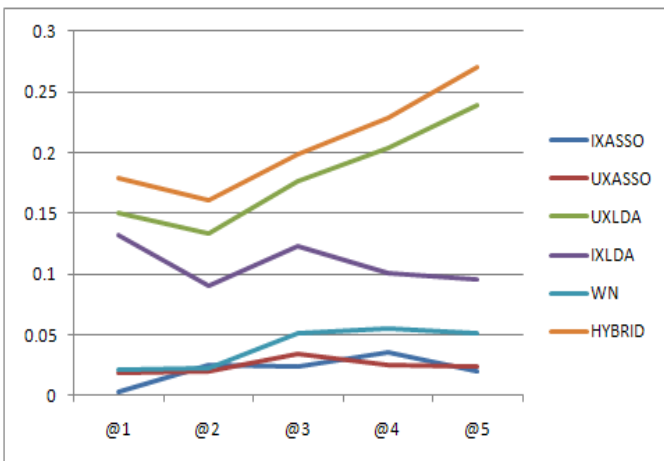
To demonstrate the effect on different semantic categories found in different spaces, we evaluate the same algorithm with two perspectives, namely user-space (hereafter, *UX*) and item-space (hereafter, *IX*). And because we are concerning about the cold-start problem, each user we withhold all but one of his tags, and withhold one tag less for another round until his tag profile length reaches five. This setup can help us to simulate the cold-start problem and to determine how the number of tags in the user profile influent the recommendation process.

Besides using the dataset as described in Section 3, we use the following parameter(s) in the experiment. In item-space association rule, the Support and Confidence are set to 10% and 80% respectively; whereas in user-space association rule, these values are 20% and 80%. Since we would like to obtain a considerable amount of rules, and hence larger coverage, we have to lower the Support values to achieve so in a skewed[28] dataset. The number of topics of LDA in both *ix* and *ux* are set to 500. As for the measurement of semantic distance in WordNet, we adopt the Wu and Palmer distance metric [33]. Table 3 summarized these settings.

From the graphs, it can be seen that the experimented algorithms can be classified into two distinct classes. One is sensitive to the current user tag profile, another one is not. We first go through the sensitive one, followed by the opposite group.



**Figure 5. Average precision of different algorithms with variable profile length**



**Figure 6. Average recall of different algorithms with variable profile length**

	Support	Confidence	No. of Topics	Similarity Metric
ixasso	10%	80%	/	/
uxasso	20%	80%	/	/
uxlda	/	/	500	/
ixlda	/	/	500	/
wn	/	/	/	Wu [33]
hybrid	/	/	500	Wu[33]

**Table 3. Parameters of different algorithms**

## 6.1 Sensitive to Tag Profile

All of the experimented algorithms are sensitive to the user’s tag profile, except the item-based Latent Dirichlet Allocation. These algorithms have a general trend in common, the averaged precision values increase as the the number of available tags in user’s tag profile increase. That is, if we know more about user’s interest, we can provide better recommendations.

In this paragraph, we would like to draw your attention to both *UXLDA* and our *HYBRID* approach. Interestingly, when you look at Figure 5 to Figure 6, you might find there are two lines running almost parallel to each other, and the one representing our strategy shifted upward. This indicates our design surpasses *UXLDA*, while the latter one in turn topped the rest of the experimented strategies. The advantages of *UXLDA* is that it emphasizes on the user’s preference. Users with variety of preferences able to find a cluster fitting their taste. As in the case in user-based Association Rule, recommending tags that stimulate user’s preference yield a better performance. Though *UXLDA* and *UXASSO* are doing similar tasks, *UXLDA* can get rid of the problem of low coverage and hence excel *UXASSO*.

Bear in mind that some individuals would prefer objective tags and others have their specific preferences, our *HYBRID* design further improves *UXLDA* by considering both factors in once. And the initial tag profile enhancement stage of our design plays an important role. It magnifies the user’s preference, and is crucial to our subsequent decision of what to deliver to user. The algorithm comes to a *dilemma point* when there are exactly two tags in the user’s tag profile, because at this stage, the two tags have equal weight, if these two tags have contradictory meaning, we cannot tell with confident that this user prefers either sides, and hence the performance dropped slightly. This is also the case for *UXLDA*. But this issue is rectified as the tag profile grows.

In the WordNet approach, the algorithm suggests *Top k* recommended tags to the user. The ranking is done by finding the most similar tags to the tags available in tag profile using Wu and Palmer[33] metric. We set the *k* to be five in our case. Assuming a user has a consistent preference, he



will use more or less the same set of tags to annotate objects. For example, if a user is optimistic, then it is likely for him to use tags such as 'great', 'funny', or 'happy' to describe an object. This explains the increasing precision values, because WordNet provides tags which are lexically correlated to the querying tags (user's preference).

Let's then take a look at the Association Rule. This algorithm generally outperforms WordNet, regardless user-space or item-space. In the beginning stage, when there is only one tag available, *UXASSO* performs better than *IXASSO*. This phenomenon is reasonable in the sense that the rules discovered in *UXASSO* reflect the user's preference. By suggesting tags with high Confidence, the probability of touching user's interest is relatively high. The low precision in *IXASSO* can be attributed to the fact that the rules generated from item-based transactions are generally objective, seldom biased towards individual's preference, which is in line with our assumption.

The distinction of these two approaches become blurred as more and more tags are available in the user's tag profile. This is especially the case when the number of tags reached five. We suppose this is caused by the low coverage in Association Rule. Because the distribution of tags is sparse, it is not easy to construct a rule given certain values of Confidence and Support. Only a small portion of tags, which are frequently co-used by the same user or co-exist in the same item, formed the basis of the rules. This leads to the result of low coverage no matter it is user-based or item-based. Consequently, the algorithm reaches a saturation point once the length of tag profile approach five, in this case for instance.

## 6.2 Insensitive to Tag Profile

The only experimented algorithm falls in this class, is the item-based Latent Dirichlet Allocation. As shown in Figure 5, the averaged precision values of this algorithm remain steady as the number of tags in user tag profile increases. As its name implies, the semantic categories or clusters formed in *IXLDA* reveal only the semantic categories of items, which is objective and descriptive in nature. It is good for discovering the semantic dimensions among items, but insufficient to stimulate user's interest.

As observed from the graphs, we interpret that there is only a few portion of overlapping between the user's preference (i.e. the tags in the user profile) and the item's semantic categories. The more availability of tags in the user's tag profile, the more obvious is the user's preference. However, there is no positive correlation between the length of tag profile and the averaged precision values. The precision values do not grow as the user's tag profile grows. We draw the conclusion that the objective information of items, or the factual categories of tags, merely occupying a small proportion of the whole set of user's preference, and this in-

formation is not enough to generate recommendations that fit all kinds of user preferences.

## 7 CONCLUSIONS

We demonstrated that taking the user's preference into account can improve the recommendation results. In a typical recommendation algorithm design, designers always focus on either user-based or item-based information, and overlook the difference between them. As we shown in our analysis, difference does exist.

We recognize the pattern found in user-based tag transactions as their preferences indicator. When a folksonomic system is mature, users in this system composite different niches, each representing certain preferences. An individual can find a niche to suit his interest. In contrast, the tag transactions in item-based are contributed by various users from different niches and thus preferences, resulting in avoiding a particular preference from dominating over others. We interpret these outcomes are objective, and unbiased. This follows the idea of classifying tags into *Personal*, *Subjective* and *Factual* categories. With *Personal*, *Subjective* classes belong to user-space, and *Factual* class belongs to item-space.

Upon verifying the dissimilarity among the user and item perspective, we examined the performance of our design, which takes the observations into account. We use a competitive strategy to find a winning cluster to user, that is, with cluster from user-space and item-space on hand, the winning cluster is the one which has the largest agreement with user tag profile. The recommendation generated in this way can leverage different interests, and therefore give a better result. The preliminary tag profile lengthening stage of our strategy enable us to maximize user's preference, which is essential to deal with the cold-start problem, as it doesn't have prior knowledge of the user.

Given the evident of different conceptual meaning found in user-space and item-space, as well as the benefit of considering them together, a simple tweak to the existing algorithms in this approach can outperform those only considered a single side, while the users of these systems can be more stimulating.

## References

- [1] G. Adomavicius and A. Tuzhilin. Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions. *IEEE Trans. on Knowl. and Data Eng.*, 17(6):734–749, 2005.
- [2] R. Agrawal, T. Imieliński, and A. Swami. Mining association rules between sets of items in large databases. In *SIGMOD '93: Proceedings of the 1993*

*ACM SIGMOD international conference on Management of data*, pages 207–216, New York, NY, USA, 1993. ACM.

- [3] H. S. Al-Khalifa and H. C. Davis. Towards better understanding of folksonomic patterns. In *HT '07: Proceedings of the eighteenth conference on Hypertext and hypermedia*, pages 163–166, New York, NY, USA, 2007. ACM.
- [4] X. Amatriain, N. Lathia, J. M. Pujol, H. Kwak, and N. Oliver. The wisdom of the few: a collaborative filtering approach based on expert opinions from the web. In *SIGIR '09: Proceedings of the 32nd international ACM SIGIR conference on Research and development in information retrieval*, pages 532–539, New York, NY, USA, 2009. ACM.
- [5] C.-m. Au Yeung, N. Gibbins, and N. Shadbolt. User-induced links in collaborative tagging systems. In *CIKM '09: Proceeding of the 18th ACM conference on Information and knowledge management*, pages 787–796, New York, NY, USA, 2009. ACM.
- [6] C. Bizer, T. Heath, K. Idehen, and T. Berners-Lee. Linked data on the web (ldow2008). In *WWW '08: Proceeding of the 17th international conference on World Wide Web*, pages 1265–1266, New York, NY, USA, 2008. ACM.
- [7] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, 2003.
- [8] M. de Gemmis, P. Lops, G. Semeraro, and P. Basile. Integrating tags in a semantic content-based recommender. In *RecSys '08: Proceedings of the 2008 ACM conference on Recommender systems*, pages 163–170, New York, NY, USA, 2008. ACM.
- [9] P. De Meo, G. Quattrone, and D. Ursino. A query expansion and user profile enrichment approach to improve the performance of recommender systems operating on a folksonomy. *User Modeling and User-Adapted Interaction*, 20(1):41–86, 2010.
- [10] F. Durao and P. Dolog. Extending a hybrid tag-based recommender system with personalization. In *SAC '10: Proceedings of the 2010 ACM Symposium on Applied Computing*, pages 1723–1727, New York, NY, USA, 2010. ACM.
- [11] J. Gemmell, M. Ramezani, T. Schimoler, L. Christiansen, and B. Mobasher. The impact of ambiguity and redundancy on tag recommendation in folksonomies. In *RecSys '09: Proceedings of the third ACM conference on Recommender systems*, pages 45–52, New York, NY, USA, 2009. ACM.
- [12] S. Givon and V. Lavrenko. Predicting social-tags for cold start book recommendations. In *RecSys '09: Proceedings of the third ACM conference on Recommender systems*, pages 333–336, New York, NY, USA, 2009. ACM.
- [13] G. Gonzalez, J. L. de la Rosa, M. Montaner, and S. Delfin. Embedding emotional context in recommender systems. In *ICDEW '07: Proceedings of the 2007 IEEE 23rd International Conference on Data Engineering Workshop*, pages 845–852, Washington, DC, USA, 2007. IEEE Computer Society.
- [14] Z. Guan, J. Bu, Q. Mei, C. Chen, and C. Wang. Personalized tag recommendation using graph-based ranking on multi-type interrelated objects. In *SIGIR '09: Proceedings of the 32nd international ACM SIGIR conference on Research and development in information retrieval*, pages 540–547, New York, NY, USA, 2009. ACM.
- [15] M. Heckner, M. Heilemann, and C. Wolff. Personal information management vs. resource sharing: Towards a model of information behaviour in social tagging systems. In *Int'l AAAI Conference on Weblogs and Social Media (ICWSM)*, San Jose, CA, USA, May 2009.
- [16] P. Heymann, D. Ramage, and H. Garcia-Molina. Social tag prediction. In *SIGIR '08: Proceedings of the 31st annual international ACM SIGIR conference on Research and development in information retrieval*, pages 531–538, New York, NY, USA, 2008. ACM.
- [17] A. Hotho, R. Jaschke, C. Schmitz, and G. Stumme. FolkRank: A ranking algorithm for folksonomies. In *Proc. FGIR 2006*, 2006.
- [18] A. B. J. A. Konstan, J. Riedl and J. Hellocker. Recommender systems: A grouplens perspective. In *Recommender Systems: Papers from the 1998 Workshop*. (AAAI Technical Report WS-98-08), pages 60C64. AAAI Press, 1998.
- [19] R. Krestel, P. Fankhauser, and W. Nejdl. Latent dirichlet allocation for tag recommendation. In *RecSys '09: Proceedings of the third ACM conference on Recommender systems*, pages 61–68, New York, NY, USA, 2009. ACM.
- [20] B. Li and Q. Zhu. The determination of semantic dimension in social tagging system based on som model. *Intelligent Information Technology Applications, 2007 Workshop on*, 1:909–913, 2008.
- [21] H. Lieberman. Letizia: An agent that assists web browsing. In *INTERNATIONAL JOINT CONFERENCE ON ARTIFICIAL INTELLIGENCE*, pages 924–929, 1995.

- [22] H. Ma, I. King, and M. R. Lyu. Effective missing data prediction for collaborative filtering. In *SIGIR '07: Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 39–46, New York, NY, USA, 2007. ACM.
- [23] A. Mathes. Folksonomies – cooperative classification and communication through shared metadata. *Computer Mediated Communication, LIS590CMC*, December, 2004.
- [24] P. Mika. Ontologies are us: A unified model of social networks and semantics. *Web Semant.*, 5(1):5–15, 2007.
- [25] G. A. Miller. Wordnet: a lexical database for english. *Commun. ACM*, 38(11):39–41, 1995.
- [26] B. S. Penn. Using self-organizing maps to visualize high-dimensional data. *Comput. Geosci.*, 31(5):531–544, 2005.
- [27] B. Sarwar, G. Karypis, J. Konstan, and J. Reidl. Item-based collaborative filtering recommendation algorithms. In *WWW '01: Proceedings of the 10th international conference on World Wide Web*, pages 285–295, New York, NY, USA, 2001. ACM.
- [28] S. Sen, J. Vig, and J. Riedl. Learning to recognize valuable tags. In *IUI '09: Proceedings of the 13th international conference on Intelligent user interfaces*, pages 87–96, New York, NY, USA, 2009. ACM.
- [29] M. Szomszor, C. Cattuto, H. Alani, K. O'Hara, A. Baldassarri, V. Loreto, and V. D. Servedio. Folksonomies, the semantic web, and movie recommendation. In *4th European Semantic Web Conference, Bridging the Gap between Semantic Web and Web 2.0*, 2007.
- [30] E. M. Voorhees. Using wordnet to disambiguate word senses for text retrieval. In *SIGIR '93: Proceedings of the 16th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 171–180, New York, NY, USA, 1993. ACM.
- [31] C.-P. Wei, C. C. Yang, and C.-M. Lin. A latent semantic indexing-based approach to multilingual document clustering. *Decis. Support Syst.*, 45(3):606–620, 2008.
- [32] R. W. White, P. Bailey, and L. Chen. Predicting user interests from contextual information. In *SIGIR '09: Proceedings of the 32nd international ACM SIGIR conference on Research and development in information retrieval*, pages 363–370, New York, NY, USA, 2009. ACM.
- [33] Z. Wu and M. Palmer. Verbs semantics and lexical selection. In *Proceedings of the 32nd annual meeting on Association for Computational Linguistics*, pages 133–138, Morristown, NJ, USA, 1994. Association for Computational Linguistics.