

# Clustering Tags in Enterprise and Web Folksonomies

Edwin Simpson

edwin.simpson@hp.com

HP Labs, Filton Road, Stoke Gifford, Bristol, UK

## Abstract

Tags lack organizational structure limiting their utility for navigation. We present two clustering algorithms that improve this by organizing tags automatically. We apply the algorithms to two very different datasets, visualize the results and propose future improvements.

## Keywords

tagging, folksonomy, clustering, similarity, navigation

## Introduction and Motivation

Tags are simple, ad-hoc labels assigned by users to describe or annotate any kind of resource for future retrieval. Their flexibility means they are easy to add and they capture a user's perspective of resources. The tags added by a group of users form a FOLKSONOMY. Unfortunately, folksonomies are difficult to navigate if tags are presented as long lists. Also, different users refer to the same concepts using different tags (Golder & Huberman 2006) and often add tags of little use to others, e.g. "toRead".

It has been proposed that folksonomies contain nested groups of tags related to common topics (Heymann & Garcia-Molina 2006; Damme, Hepp, & Siorpaes 2007). This paper describes two algorithms that extract such topic groupings from TAG CO-OCCURRENCE data. Co-occurrences between tags occur when both tags are used with the same resource. The algorithms should produce evenly-sized and intuitive clusters for browsing.

We collected the first dataset from the social bookmarking service Delicious (<http://del.icio.us>), and the second from an internal bookmarking service, Labbies, used by a group of researchers at HPLabs. We obtained a subset of Delicious data by selecting all tags from users who have used the tag "dspace" during a 16 week period. The use of a common tag ensures there are some relationships between the tags in the dataset. Statistics of the datasets are given in Table 1.

## Tag Similarity Graphs

We can extract clusters from a similarity graph, where nodes represent tags and edge weights represent strength of similarity based on the number of tag co-occurrences.

Copyright © 2008, Association for the Advancement of Artificial Intelligence ([www.aaai.org](http://www.aaai.org)). All rights reserved.

| Dataset                | Labbies | Delicious |
|------------------------|---------|-----------|
| Number of Users        | 20      | 136       |
| Number of Bookmarks    | 1935    | 95155     |
| Number of Tags         | 2092    | 8012      |
| No. Tag Co-occurrences | 9526    | 61453     |

Table 1: Dataset Statistics

Extremely popular tags have high co-occurrence values with many other weakly related tags. To avoid this bias when calculating similarity, we calculate the NORMALIZED CO-OCCURRENCE, or NCO by normalizing the raw co-occurrence value relative to the popularity of two tags using the Jaccard index (Begelman, Keller, & Smadja 2006).

$$NCO = \frac{|A \cap B|}{|A \cup B|} \quad (1)$$

Here  $A$  is the set of documents tagged with tag  $a$ , and  $B$  the set of documents tagged with  $b$ . Other normalizations are also possible, such as cosine similarity.

## Clustering Algorithms

The algorithms we tested are hierarchical divisive clustering algorithms. Algorithm (1) is as follows, starting with a graph containing all tag relationships,  $G$ :

1. Count the number of clusters present in  $G$  by counting the disconnected sub-graphs.
2. Evaluate the current clustering using MODULARITY (Newman & Girvan 2004), a quality measure defined as:

$$modularity = Tre - ||e^2||, \quad (2)$$

where  $Tre$  is the fraction of edges that connect nodes in the same cluster, and  $||e^2||$  is the fraction of edges that would connect nodes in the same cluster if the clusters were marked randomly in the graph. If  $modularity >$  all previous modularities, set  $G_{maxMod} = G$ .

3. Remove the edge with the lowest NCO value from  $G$ .
4. Repeat process from step 1 until no edges are left. Heuristics may be used here to reduce the number of iterations, e.g. stop when a modularity exceeds a certain threshold.

