# Chapter I

# Leveraging User–specified Metadata to Personalize Image Search

**Kristina Lerman**
*USC Information Sciences Institute, USA*
**Anon Plangprasopchok**
*USC Information Sciences Institute, USA*

## ABSTRACT

The social media sites, such as Flickr and del.icio.us, allow users to upload content and annotate it with descriptive labels known as tags, join special-interest groups, etc. We believe user-generated metadata expresses user's tastes and interests and can be used to personalize information to an individual user. Specifically, we describe a machine learning method that analyzes a corpus of tagged content to find hidden topics. We then these learned topics to select content that matches user's interests. We empirically validated this approach on the social photo-sharing site Flickr, which allows users to annotate images with freely chosen tags and to search for images labeled with a certain tag. We use metadata associated with images tagged with an ambiguous query term to identify topics corresponding to different senses of the term, and then personalize results of image search by displaying to the user only those images that are of interest to her.

## INTRODUCTION

The rise of the Social Web underscores a fundamental transformation of the Web. Rather than simply searching for, and passively consuming, information, users of blogs, wikis and social media sites like del.icio.us, Flickr and digg, are creating, evaluating, and distributing information. In the process of using these sites, users are generating not only content that could be of interest to other users, but also a large quantity of metadata in the form of tags and ratings, which can be used to improve Web search and personalization.

Web personalization refers to the process of customizing Web experience to an individual user (Mobasher, 2000). Personalization is used by online stores to recommend relevant products to a particular user and to customize a user's shopping experience. It is used by advertising firms to target ads to a particular user. Search personalization has also been studied as a way to improve the quality of Web search (Ma, 2007) by disambiguating query terms based on user's browsing history or by eliminating irrelevant documents from search results.

Personalizing image search is an especially challenging problem, because, unlike documents, images generally contain little text that can be used for disambiguating terms. Consider, for example, a user searching for photos of "jaguars." Should the system return images of luxury cars or spotted felines to the user? In this context, personalization can help disambiguate query keywords used in image search or to weed out irrelevant images from search results. Therefore, if a user is interested in wildlife, the system will show her images of the predatory cat of South America and not of an automobile.

In this chapter we explore a novel source of evidence – user-generated metadata – that can be used to personalize image search results. We perform a case study of the technique on the social photo-

sharing site Flickr, which allows users to upload images and label them with freely-chosen keywords, known as tags. Tags are meant to help users organize content and make it searchable by themselves and others. In addition to describing and categorizing images, tags also capture user's photography interests. We use a machine learning method to find topics of a large corpus of tagged images returned by image search on Flickr. We then use the learned topics to match images to an individual user's interests. This appears to be a promising method for improving the quality of image search results.

## BACKGROUND

Traditionally, personalization techniques fall in one of two categories: collaborative-filtering or profile-based. The first, collaborative filtering (Breese, 1998; Schafer, 2007), aggregates opinions of many users to recommend new items to like-minded users. In these systems, users are asked to rate items on a universal scale. The system then analyses ratings from many users to identify those sharing similar opinions about items and recommends new items that these users liked. Netflix uses collaborative filtering to recommend movies to its subscribers. Amazon uses a similar technology to display other products that users who purchased a given product were also interested in. Since users are asked to rate items on a universal scale, the questions of how to design the rating system and how to elicit high quality ratings from users are very important. Despite the early concern that users lack incentives for making recommendations and, therefore, will be reluctant to make the extra effort, there is new evidence (Schafer, 2007) that this does not appear to be the case. It appears that, at the very least, users find value in a collaborative rating system as an extension of their memory.

The second class of personalization systems uses a profile of user's interests to target items for user's attention. The profile can be created explicitly by the user (Ma, 2007), or mined from data about user's behavior. Examples of the latter include data about user's Web browsing (Mobasher, 2000) and purchasing (Agrawal, 1994) behavior. One problem with this approach is that it is time-consuming for users to keep their explicit profiles current. Another problem is that while data mining methods have proven effective and commercially successful, in most cases they use proprietary data, which is not easily accessible to researchers.

Machine learning has played an increasingly important role in personalization. (Popescul, 2001) proposed a probabilistic generative model that describes co-occurrences of users and items of interest. In particular, the model assumes a user generates her topics of interest; then the topics generate documents and words in those documents if the user prefers those documents. The author-topic model (Rosen-Zvi , 2004) is also used to find latent topics in a collection of documents and group documents according to topic. If a user prefers one document (or topic), this method can be used to recommend other relevant documents. These models, however, do not carry any information about individual users, their tastes and interests. However, a recent work this area described a mixture model for collaborative filtering that takes into account users' intrinsic preferences about items (Jin, 2006). In this model, item rating is generated from both the item type and user's individual preference for that type. Intuitively, like-minded users provide similar ratings on similar types of items (e.g., movie genres). When predicting a rating of an item for a certain user, the user's previous ratings on other items will be used to infer a like-minded group of users, and then the "common" rating of that group is used in the prediction. This type of model can conceivably be adapted to social metadata and be used to personalize results of image search.

## LEVERAGING USER-GENERATED METADATA FOR PERSONALIZATION

The Web 2.0 has created an explosion not only in user-generated content, but also in user-generated metadata. This "data about data" is expressed in a number of ways on the Social Web sites: through tags (descriptive labels chosen by the user), ratings, comments and discussion about its, items that users mark as their favorite, and through the social networks users create and the special-interest groups they participate in. This metadata provides a wealth of information about individual user's tastes, preferences and interests. Social Web sites currently don't make much use of this data, except perhaps to target advertisement to individual users or groups. However, this data has the potential to transform how users

discover, process and use information. For example, Web browsing and search can be tuned to an individual user based on his or her expressed interests. Rather than requiring the user disambiguate query terms, e.g., through query expansion, in order to improve results of Web search, a personalization system would infer a user's meaning based on the rich trace of content and metadata the user has created. Such metadata could also filter the vast stream of new content created daily on the Web and recommend to the user only that content the user would find relevant or interesting. Personalization, recommendation and filtering are just some of the applications of user-generated metadata that have recently been explored by researchers.

## Issues, Controversies, Problems

In this chapter we focus on tags, although the analysis can be easily expanded to include other types of metadata, including social networks (Lerman et al., 2007). Tags are freely-chosen keywords users associate with content. Tagging was introduced as a means for users to organize their own content in order to facilitate searching and browsing for relevant information. The distinguishing feature of tagging systems is that they use an uncontrolled vocabulary, and that the user is free to highlight any one of the object's properties. From an algorithmic point of view, tagging systems offer many challenges that arise when users try to attach semantics to objects through keywords (Golder, 2006). These challenges are homonymy (the same tag may have different meanings), polysemy (tag has multiple related meanings), synonymy (multiple tags have the same meaning), and "basic level" variation (users describe an item by terms at different levels of specificity, e.g., "beagle" vs "dog"). Despite these challenges, tagging is a light weight, flexible categorization system. The growing amount of tagged content provides evidence that users are adopting tagging on Flickr (Marlow, 2006), Del.icio.us and other collaborative tagging systems. In a small case study we show how tags on the social photo-sharing site Flickr can be used to personalize results of image search.
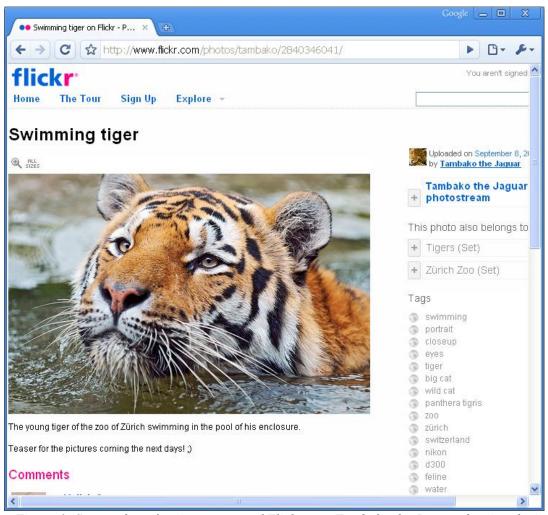
*Figure 1: Screen shot of an image page of Flickr user Tambako the Jaguar showing the image and the tags he attached to the image.*

Flickr consists of a collection of interlinked user, photo, tag and group pages. A typical Flickr photo page, shown in Figure 1, provides a variety of information about the image: who uploaded it and when, what groups it has been submitted to, its tags, who commented on the image and when, how many times the image was viewed or bookmarked as a "favorite." The user calling himself (user's may reveal their gender in their profile, as this user has chosen to do) "Tambako the Jaguar" posted a photograph of a swimming tiger at a Swiss zoo. To the right of the image is a list of keywords, tags, the user has associated with the image.[1] These tags include "tiger," "big cat," "wild cat," "panthera tigris," and "feline," all useful terms for describing this particular sense of the word "tiger." Clicking on a user's name brings up that user's photo stream, which shows the latest photos he uploaded, the images he marked as "favorite," and his profile, which gives information about the user, including a list of his social network (contacts) and groups he belong to. Clicking on the tag shows user's images that have been tagged with that keyword, or all public images that have been similarly tagged.
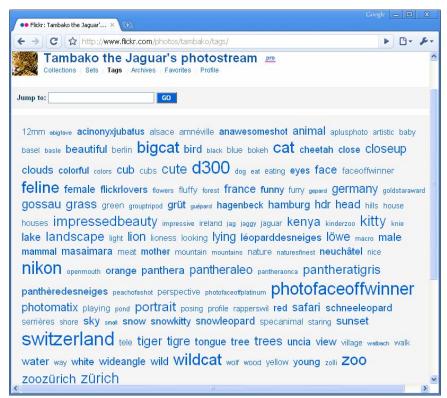
4

*Figure 2: Tag cloud view of the tags the owner of the image in Fig. 1 used to annotate his images. The bigger the font, the more frequently that tag was used by the user.*

Information about a user's photography tastes and interests is contained in the rich metadata he creates in his everyday activities on Flickr. He expresses these interests through the contacts he adds to his social networks, the groups he joins, the images of other photographers he marks as his favorite or comments on, as well as through tags he adds to his own images. Figure 2 shows a tag cloud view of the tags that "Tamboko the Jaguar" used to annotate his images on Flickr. The bigger the font, the more frequently that keyword was used. These tags clearly show that the user is interested in wildlife (bigcat, cat, lion, cheetah, tiger, tigre, wildcat) and nature (clouds, mountains) photography. They also show that he shoots with a Nikon (nikon, d300) and has traveled extensively in Europe (switzerland, germany, france) and parts of Africa (kenya). These interests are further reflected in the groups the user joined, which are listed on his profile page, that include such ad-hoc groups as "Horns and Antlers," "Exotic cats," "Cheetah Collection," and many others. In this work, we view group names just as we treat tags themselves. In fact, group names can be viewed as publicly agreed-upon tags.

Flickr allows users to search for photos that contain specified keywords in their descriptions (including titles) or tags. A user can search all public photos, or restrict the search to photos from her contacts, her own photos, or photos she marked as her favorite. Search results are by default displayed in reverse chronological order of being uploaded, with the most recent images on top. Another option is to display images by their "interestingness,"[2] with the most "interesting" images on top. Suppose a user is interested in wildlife photography and wants to see images of tigers on Flickr. As of September 9, 2008, the search of all public images tagged with the keyword "tiger" returned over 170,000 results. When arranged by "interestingness," the first few pages of results contain images of tigers, but also many irrelevant images of cats, kids, butterflies, flowers, and golf, as shown in Figure 3, and also sharks and screenshots of Mac OS X computer system.
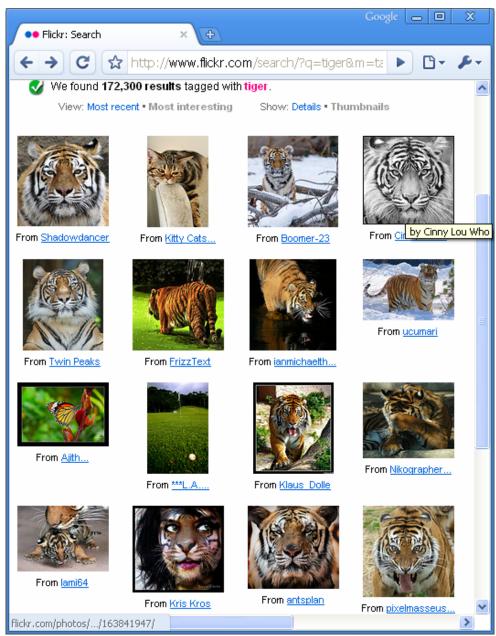
5

*Figure 3: Results of image search on Flickr for images tagged with "tiger"*

We assume that when a search term is ambiguous, the sense that the user has in mind is related to his or her interests. A wildlife photographer searching for "tiger" images is probably not interested in photographs of children with face paint. Similarly, a child photographer searching for pictures of "newborns" is most likely interested in images of human babies, not kittens or tiger cubs. In this chapter we show that we can improve the relevance of image search by personalizing image search results on Flickr. We use user-generated metadata, in the form of tags and the groups, for this purpose. Inferring personal interests from tags, however, is problematic, since this data is sparse (few tags per image) and noisy (idiosyncratic vocabulary use, synonyms, etc). Machine learning methods, which try to find statistical correlations in the data, directly address some of these challenges. In the section below, we describe a machine learning-based method that exploits information contained in user-generated metadata, specifically tags, to personalize image search results to an individual user.

## Probabilistic Model for Tag-based Personalization

We outline a probabilistic model that takes advantage of the images' tag and group information to discover latent topics contained in a set of images. If the dataset is a result of a search for images that have been tagged with the query term, the topics correspond to different senses of the query term. The users' interests can similarly be described by collections of tags they used to describe their own images. The latent topics found by the model can be used to personalize search results by finding images on topics that are of interest to the user.

We consider four types of entities in the model: a set of users $U = \{u_1, \ldots, u_n\}$, a set of images or photos $I = \{i_1, \ldots, i_m\}$, a set of tags $T = \{t_1, \ldots, t_o\}$, and a set of groups $G = \{g_1, \ldots, g_p\}$. A photo $i_x$ posted by user (image owner) $u_x$ is described by a set of tags $\{t_{x1}, t_{x2}, \ldots\}$ and submitted to several groups $\{g_{x1}, g_{x2}, \ldots\}$. This post could be viewed as a tuple $<i_x, u_x, \{t_{x1}, t_{x2}, \ldots\}, \{g_{x1}, g_{x2}, \ldots\}>$. We assume that there are n users, m posted photos and p groups in Flickr. Meanwhile, the vocabulary size of tags is q. In order to filter images retrieved by Flickr in response to tag search and personalize them for a user u, we compute the conditional probability $p(i|u)$, that describes the probability that the photo i is relevant to u based on her interests. Images with high enough $p(i|u)$ are then presented to the user as relevant images.

As mentioned earlier, users choose tags from an uncontrolled vocabulary according to their styles and interests. Images of the same subject could be tagged with different keywords although they have similar meaning. Meanwhile, the same keyword could be used to tag images of different subjects. In addition, a particular tag frequently used by one user may have a different meaning to another user. Probabilistic models offer a mechanism for addressing the issues of synonymy, homonymy and tag sparseness that arise in tagging systems.

We use a probabilistic topic model (Rosen-Zvi, 2004) to model user's image posting behaviour. As in a typical probabilistic topic model, topics are hidden variables, representing knowledge categories. In our case, topics are equivalent to image owner's interests. The process of photo posting by a particular user could be described as a stochastic process:

- User u decides to post a photo i.
- Based on user u's interests and the subject of the photo, a set of topics z are chosen.
- Tag t is then selected based on the set of topics chosen in the previous state.
- In case that u decides to expose her photo to some groups, a group g is then selected according to the chosen topics.
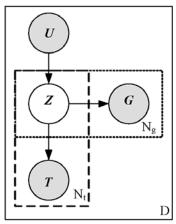


*Figure 4. Graphical representation for model-based information filtering. U, T, G and Z denote variables "User", "Tag", "Group", and "Topic" respectively. Nt represents a number of tag occurrences for a one photo (by the photo owner); D represents a number of all photos on Flickr. Meanwhile, $N_g$ denotes a number of groups for a particular photo.*

The process is depicted in a graphical form in Figure 4. We do not treat the image i as a variable in the model but view it as a co-occurrence of a user, a set of tags and a set of groups. From the process described above, we can represent the joint probability of user, tag and group for a particular photo as

$$p(i) = p(u_i, T_i, G_i) = p(u_i)\left(\prod_{n_t}\left[\sum_k p(z\,|\,u_i)p(t_i\,|\,z)\right]^{n_i(t)}\right)\cdot\left(\prod_{n_g}\left[\sum_k p(z\,|\,u_i)p(g_i\,|\,z)\right]^{n_i(g)}\right)$$

$n_t$ and $n_g$ are the numbers of all possible tags and groups respectively in the data set. Meanwhile, $n_i(t)$ and $n_i(g)$ act as indicator functions: $n_i(t)=1$ if an image i is tagged with tag t; otherwise, it is 0. Similarly, $n_i(g)=1$ if an image i is submitted to group g; otherwise, it is 0. k is the predefined number of topics. Note that it is straightforward to exclude photo's group information from the above equation simply by omitting the terms relevant to g.

In order to estimate parameters $p(z|u_i)$, $p(t_i|z)$, and $p(g_i|z)$, we define a log likelihood L=log(($\Pi_i$p(i))), which measures how the estimated parameters fit the observed data, in our case all the photos in the dataset. We use the EM algorithm (Dempster, 1977) to iterate between parameter estimates until the log likelihood for all parameter values converges. L is used as an objective function to estimate all parameters. In the expectation step (E-step), the joint probability of the hidden variable Z given all observations is computed from the following equations:

$$p(z\,|\,t,u) \propto p(z\,|\,u)\cdot p(t\,|\,z)$$

$$p(z\,|\,g,u) \propto p(z\,|\,u)\cdot p(g\,|\,z)$$

L cannot be maximized easily, since the summation over the hidden variable Z appears inside the logarithm. We instead maximize the expected complete data log-likelihood over the hidden variable, $E[L^c]$, which is defined as

$$E[L^c] \quad = \quad \sum_m \log(p(u)$$
$$+ \quad \sum_m\sum_t n_i(t)\cdot\sum_z p(z\,|\,u,t)\log(p(z\,|\,u)\cdot p(t\,|\,z))$$
$$+ \quad \sum_m\sum_g n_i(g)\cdot\sum_z p(z\,|\,u,g)\log(p(z\,|\,u)\cdot p(g\,|\,z))$$

Since the term $\Sigma\log(p(u)$ is not relevant to parameters and can be computed directly from the observed data, we discard this term from the expected complete data loglikelihood. With normalization constraints on all parameters, Lagrange multipliers $\tau$, $\rho$, $\psi$ are added to the expected log likelihood, yielding the following equation

$$H = E[L^c] \quad + \quad \sum_z \tau_z\left(1-\sum_t p(t\,|\,z)\right)$$
$$+ \quad \sum_z \rho_z\left(1-\sum_t p(g\,|\,z)\right)$$
$$+ \quad \sum_u \psi_u\left(1-\sum_z p(z\,|\,u)\right)$$

We maximize *H* with respect to *p(t/z)*, *p(g/z)*, and *p(z/u)*, and then eliminate the Lagrange multipliers to obtain the following equations for the maximization step:

$$p(t \mid z) \quad \propto \quad \sum_i n_i(t) \cdot p(z \mid t, u)$$

$$p(g \mid z) \quad \propto \quad \sum_i n_i(g) \cdot p(z \mid g, u)$$

$$p(z \mid u) \quad \propto \quad \sum_i \left( \sum_t n_i(t) \cdot p(z \mid t, u) + \sum_g n_i(g) \cdot p(z \mid g, u) \right)$$

The algorithm iterates between E and M step until the log likelihood for all parameter values converges. Additional details about model derivation and inference method can be found in (Lerman, 2007).

We can use the parameters inferred from the dataset to find the images *i* most relevant to the interests of a particular user u'. We do so by computing the conditional probability p(i|u'):

$$p(i \mid u') = \sum_z p(u_i, T_i, G_i \mid z) \cdot p(z \mid u'),$$

where $u_i$ is the owner of image i in the data set, and $T_i$ and $G_i$ are, respectively, the set of all the tags and groups for the image i.

We represent the interests of user u' as an aggregate of the tags she used in the past for tagging her own images. This information is used to to approximate p(z|u'):

$$p(z \mid u') \propto \sum_t n(t' = t) \cdot p(z \mid t)$$

where n(t'=t) is a frequency (or weight) of tag t' used by u'. Here we view n(t'=t) as proportional to p(t'|u'). Note that we can use either all the tags u' had applied to the images in her photostream, or a subset of these tags, e.g., only those that co-occur with some tag in user's images.

## Flickr Case Study

To show how user-generated metadata can be used to personalize image search results, we retrieved a variety of data from Flickr using their public API. We collected images by performing a single keyword tag search of all public images on Flickr. We specified that the returned images are ordered by their "interestingness" value, with most interesting images first. We retrieved the links to the top 4500 images for each of the search term. We indicate the possible senses of the query term below:

- **tiger**: (a) big cat ( e.g., Asian tiger), (b) shark (Tiger shark), (c) flower (Tiger Lily), (d) golfing (Tiger Woods), etc.
- **newborn**: (a) human baby, (b) kitten, (c) puppy, (d) duckling, (e) foal, etc.
- **beetle**: (a) a type of insect and (b) Volkswagen car

For each image in the set, we used Flickr's API to retrieve the name of the user who posted the image (image owner), and all the image's tags and groups.

We manually evaluated the top 500 images in each data set and marked each as relevant if it was related to the first sense (a) of the search term listed above, or not relevant, if the evaluator deemed it not relevant or could not understand the image well enough to judge its relevance.

| query | relevant | not relevant | precision |
|-------|----------|--------------|-----------|
| newborn | 412 | 83 | 0.82 |
| tiger | 337 | 156 | 0.67 |
| beetle | 232 | 268 | 0.46 |

*Table 1 – Number of the top 500 most "interesting" images in each search set that were deemed relevant to the first sense of the query term.*

The table above reports search precision within the 500 labeled images, as judged from the point of view of the searching users. Precision is defined as the proportion of relevant images within the top 500 images. Search precision on these sample queries is not very high due to the presence of false

positives – images not relevant to the sense of the search term the user had in mind. We do not compute search recall, or the proportion of all relevant images that are retrieved, since it is difficult for us to estimate how many images relevant to each search there are on Flickr.

Our objective is to personalize image search results; therefore, to evaluate our approach, we need to have users to whom the search results will be tailored. We identified four users who are interested in the first sense of each search term. For the **newborn** set, those users were one of the authors and three other contacts within that user's social network who are known to be interested in child photography. For the other data sets, the users were chosen from among the photographers whose images were returned by the tag search. We studied each user's profile, including group membership, user's statement, and user's photo stream, to confirm that the user was interested in the first sense of the search term. For each of the twelve users, we retrieved a list of all tags, with their frequencies, that these users have used to annotate their own images.

The model was trained separately on each set of 4500 images, with the number of topics fixed at ten. Computation of p(t|z) is central to the parameter estimation process, and it tells us something about how strongly a tag t contributes to a topic z. Table 1 shows the most probable 25 tags for some of the learned topics in the **tiger** dataset. Although the tag "tiger" dominates most topics, we can discern different themes from the other tags that appear in each topic. Thus, topic z3 is obviously about domestic cats, while topic z8 is about Apple computer products. Meanwhile, topic z2 is about flowers and colors ("flower," "lily," "yellow," "pink," "red"); topic z6 is about places ("losangeles," "sandiego," "lasvegas," "stuttgard,"), presumably because they have zoos. Topic z7 contains several variations of tiger's scientific name, "panthera tigris." This method also appears to identify related terms which can be used to expand the query. Topic z5, for example, gives synonyms "cat," "kitty," as well as the more general term "pet" and the more specific terms "kitten" and "tabby." It even contains the Spanish version of the word: "gatto." Recognizing ambiguity of tags, Flickr separates images tagged with some keyword into clusters, with images in each cluster related by meaning. For the tag "tiger",[3] for example, it finds four clusters. The first cluster is about wildlife in zoos, the second about Apple Computer products, and the third about orange flowers. The fourth cluster contains images invited to best-of groups and tagged with group names, such as "specanimal", "impressedbybeauty," etc. Although clustering appears to find different senses of ambiguous tags similar to our topic model approach, our framework has the added advantage that the learned topics (or more accurately, the learned probabilities) can be further used to personalize search results.

| z1 | z2 | z3 | z6 | z7 | z8 |
|---|---|---|---|---|---|
| tiger | tiger | tiger | tiger | nationalzoo | tiger |
| zoo | specanimal | cat | tigers | tiger | apple |
| animal | animal…lite | kitty | dczoo | sumatrantiger | mac |
| nature | abigfave | cute | tigercub | zoo | osx |
| animals | flower | kitten | california | nikon | macintosh |
| wild | butterfly | cats | lion | washingtondc | screenshot |
| tijger | macro | orange | cat | smithsonian | macosx |
| wildlife | yellow | eyes | cc100 | washington | desktop |
| ilovenature | swallowtail | pet | florida | animals | imac |
| cub | lily | tabby | girl | cat | stevejobs |
| siberiantiger | green | stripes | wilhelma | bigcat | dashboard |
| blijdorp | canon | whiskers | self | tigris | macbook |
| london | insect | white | lasvegas | panthera | powerbook |
| australia | nature | art | stuttgart | bigcats | os |
| portfolio | pink | feline | me | d70s | 104 |
| white | red | fur | baby | panthera...sumatrae | canon |
| dierentuin | flowers | animal | tattoo | dc | x |
| toronto | orange | gatto | endangered | sumatrae | ipod |
| stripes | eastern | pets | illustration | animal | computer |
| amurtiger | usa | black | ?? | 2005 | ibook |

| nikon...ggallery | impressed… | paws | losangeles | pantheratigris | intel |
|---|---|---|---|---|---|
| s5600 | tag2 | furry | portrait | nikond70 | keyboard |
| eyes | specnature | nose | sandiego | d70 | widget |
| sydney | black | teeth | lazoo | 2006 | wallpaper |
| cat | streetart | beautiful | giraffe | topv111 | laptop |

*Table 2 – Top 25 tags ordered by p(t|z) for some of the learned topics in the "tiger" dataset.*

We evaluated model-based personalization by using the learned parameters and the information about the interests of the selected users to compute p(i|u') for the top 500 (manually labeled) images in the set. Once images were ranked by how similar they are to user's interests, we calculated how many of the top-ranked x images were relevant to each user. From this number, we calculated the ==precision== of search, reported in Figure 5. The thick line in Figure 5 presents results of plain search, with images ranked by Flickr according to how "interesting" they are, while the thin dashed lines report precision of personalized search results for each of the users. As can be seen from the figure, most of the dashed lines are above the plain search line, indicating improve relevance for most users. The best results were for the **beetle** set. While fewer than half of the returned images were relevant to the "insect" sense of the word, personalization filtering pushed relevant images higher. In fact, for three of the four users, all of the top 100 images were deemed to be relevant. On the **newborn** set, personalization generally helped improve search results for all but user3. For two of the users, the top 200 of the filtered images were all relevant. Results were less impressive for the **tiger** set, where plain search outperformed filtered search for three of the four users. The four chosen users were all highly regarded photographers, not quite average Flickr users, and had wide ranging photography interests. The poor performance of personalization can probably be explained by these users' breadth of interests.
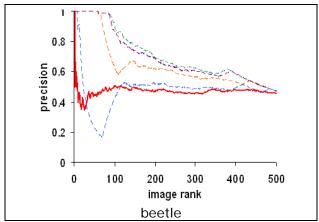


newborn

tiger

*Figure 5: Tag-based personalization results for tag search on Flickr for query words "newborn", "tiger", and "beetle". We picked four (different) users for each query that were interested in a single sense of the query term.*

## FUTURE RESEARCH DIRECTIONS

User-generated metadata is a rich source of information about user's tastes and preferences that can be leveraged to personalize information to an individual user. This personalization can be applied to browsing and search. In this chapter we explored the use of tags and groups (which were also viewed as publicly agreed-upon tags) for representing user's interests. In addition to tags, users express their interests in other ways, e.g., through the social networks they join and through the content they mark as their favorite. It is important to develop algorithmic approaches that combine multiple heterogeneous sources of metadata to succinctly represent user's information preferences.

The personalization method described in this chapter will fail if a user makes a query in a domain in which she has not previously expressed any interest. For example, suppose that a child portrait photographer wants to find beautiful mountain scenery. If she has never created tags relating to mountains landscape photography in general, the personalization method described above will fail. However, the Flickr community as a whole has generated a significant amount of data about nature and landscape photography and mountains in particular. Analysis of community-generated data can help the user discover mountain imagery the community has identified as being good. We need algorithms to mine community-generated metadata and knowledge to identify community-specific topics of interest, vocabulary, authorities within the communities and community-vetted content.

## CONCLUSION

In addition to creating content, users of Web 2.0 sites generate large quantities of metadata, or data about data, that describe their interests, tastes and preferences. These metadata, in the form of tags and social networks, are created mainly to help users organize and manage their own content. These types of metadata can also be used to target relevant content to the user through recommendation or personalization.

This chapter describes a machine learning-based method for personalizing results of image search on Flickr. Our method relies on metadata created by users through their everyday activities on Flickr, namely the tags they used for annotating their images and the groups to which they submitted these images. This information captures user's tastes and preferences in photography and can be used to personalize image search results to the individual user. We validated our approach by showing that it can be used to improve precision of image search on Flickr for three ambiguous terms: "newborn," "tiger,"

and "beetle." In addition to improving search precision, the tag-based approach can also be used to expand the search by suggesting other relevant keywords (e.g., "pantheratigris," "bigcat" and "cub" for the query "tiger").

## REFERENCES

Agrawal, R., & Srikant, R. (1994). Fast algorithms for mining association rules. In Bocca, J. B., Jarke, M.& Zaniolo, C. (Eds.), *Proceedings of the 20$^{th}$ Int. Conf. Very Large Data Bases, VLDB* (pp. 487—499). Morgan Kaufmann.

Breese, J., Heckerman, D.& Kadie, C. (1998). Empirical analysis of predictive algorithms for collaborative filtering. In *Proceedings of the 14th Annual Conference on Uncertainty in Artificial Intelligence* (pp. 43—52). San Francisco, CA: Morgan Kaufmann.

Dempster, A. P., Laird, N.M. & Rubin, D.B. (1977). Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society. Series B (Methodological) 39*(1), 1-38.

Golder, S.A. & Huberman, B.A.(2006). The structure of collaborative tagging systems. *Journal of Information Science 32*(2), 198-208.

Jin, R., Si, L., & Zhai, C. (2006) A study of mixture models for collaborative filtering. *Information Retrieval* 9(3):357–382.

Lerman, K., Plangprasopchok, A. & Wong, C. (2007). Personalizing Image Search Results on Flickr. In *Proceedings of AAAI workshop on Intelligent Techniques for Information Personalization.* Vancouver, Canada, AAAI Press.

Ma, Z., Pant, G.& Liu-Sheng, O.R. (2007). Interest-based personalized search. *ACM Trans. Inf. Syst. 25*(1).

Marlow, C., Naaman, M., boyd, d. & Davis, M. (2006). Ht06, tagging paper, taxonomy, flickr, academic article, toread. *Proceedings of Hypertext 2006*. New York: ACM.

Mobasher, B., Cooley, R. & Srivastava, J. (2000). Automatic personalization based on web usage mining. *Commun. ACM 43*(8), 142-151.

Popescul, A., Ungar, L., Pennock, D. & Lawrence, S. (2001). Probabilistic models for unified collaborative and content-based recommendation in sparse-data environments. In *17th Conference on Uncertainty in Artificial Intelligence* (pp. 437-444).

Rosen-Zvi, M., Griffiths, T., Steyvers, M. & Smyth, P. (2004). The author-topic model for authors and documents. In *Proceedings of the 20th conference on Uncertainty in artificial intelligence* (pp. 487—494). Arlington, Virginia, United States: AUAI Press.

Schafer, J., Frankowski, D., Herlocker, J. & Sen, S. (2007). Collaborative filtering recommender systems. *The Adaptive Web*, 291-324.

## ADDITIONAL READING SECTION

Crane, R. & Sornette, D. (2008) Viral, quality, and junk videos on youtube: Separating content from noise in an information-rich environment. *Proceedings of AAAI symposium on Social Information Processing (SIPS08)*, Menlo Park, CA, AAAI.

K. Lerman. (2007) Social information processing in social news aggregation. *IEEE Internet Computing: special issue on Social Search*, 11(6), pp.16-28.

Mika, P. (2005). Ontologies are us: A unified model of social networks and semantics. In *International Semantic Web Conference (ISWC-05)*.

Mislove, A., Gummadi, K.P., and Druschel, P. (2006) Exploiting social networks for internet search. *Proceedings of the 5th Workshop on Hot Topics in Networks (HotNets·S06)*.

Noll, M. G.  & Meinel, C. (2007). Web Search Personalization via Social Bookmarking and Tagging, *Proceedings of 6th International Semantic Web Conference (ISWC)*, Springer LNCS 4825, Busan, South Korea, pp. 367-380.

Perugini, S., Gonçalves, M. & Fox, E.A. (2004). Recommender systems research: A connection-centric survey. *Journal of Intelligent Information Systems, 23*(2), 107-143.

Plangprasopchok, A. & Lerman, K. (2007). Exploiting Social Annotation for Automatic Resource Discovery, *Proceedings of AAAI workshop on Information Integration (IIWeb-07)*.

## KEY TERMS & DEFINITIONS

**Social Media**:  a term that defines activities by which users create and publish content on the Web. Examples include Flickr, del.icio.us, Digg and many others.
**Social Web**:  an umbrella term that includes social media and social networking sites, like Facebook and MySpace.
**Machine learning**: a subfield of artificial intelligence that is concerned with algorithms and techniques for allowing computers to learn from data.
**Personalization**: algorithms and techniques that tailor content to individual users
**Image search**: a type of Web search that returns images matching a given (text) query
**Metadata**: 'data about data'
**Tag**: a freely-chosen keyword or term associated with content by the user

## ENDNOTES

[1] Although any user can tag the image unless specifically barred from doing so by the image owner, generally, only the image owner tags them.
[2] Flickr uses a proprietary algorithm to evaluate how "interesting" an image is based on the number of times it was viewed, commented on, marked as a favourite, among other factors.
[3] http://www.flickr.com/photos/tags/tiger/clusters/