

Trend-Aware User Modeling with Location-Aware Trends

Marcel Kanta, Marián Šimko, Mária Bieliková

Institute of Informatics and Software Engineering,

Faculty of Informatics and Information Technologies , Slovak University of Technology,

Ilkovičova 3, 842 47 Bratislava, Slovakia

xkanta@is.stuba.sk, {simko, bielik}@fiit.stuba.sk

Abstract—Microblogs are a phenomenon of modern social media. As there is much real-time social information in there, they are candidates to be used as a source for mining important information enhancing user experience in variety of web applications, especially those related with content adaptation and recommendation. In this paper we deal with microblog-based user models. We propose trend-aware user model with location-aware trends, which focuses on location aspects and trends. It is a general model, which can be used in various domains. We evaluated the model in a domain of news recommendations and we showed that recommendation based on this model outperforms state-of-the-art approaches.

Keywords-microblog; Twitter; user modeling; trends;location

I. INTRODUCTION AND RELATED WORK

The Web brings people so much information that people face these days an information overload. There are several ways to cope with this problem; the most significant are recommender systems [9] and faceted search [1], which are present in various web applications and facilitate access to information and improve user experience on the Web. Web applications yet typically have their own data model, where the information is stored in a structured way with connections between entities. Most of those applications are personalized, so they try to capture the user characteristics, goals, interests or intents and they utilize it for recommendation. Recommender systems seek new ways how to improve their services and make user experience better. To achieve this goal, new sources of information have to be examined and utilized to make user model more precise.

Twitter is a highly significant source of social information and interaction. There is a constant concern in scientific research in microblogs, particularly in Twitter, covering wide scope of interests ranging from resource ranking to sentiment analysis [8][11][13]. In our work we use Twitter as a source for user modeling.

User modeling in open information spaces tends to rely on lightweight descriptions of subject domain [4][6]. There are several works dealing with user modeling based on or related to microblogging service Twitter and news recommendation [2][7][15].

Abel et al. in the work [2] created a framework for user modeling based on entities, topics or hash-tags. Tweet enrichment is also a part of this framework. In their approach they enrich tweets with entities/topics found in links users share in tweets. The result is the user model, which serves as a basis for making (news) recommendations for Twitter users. Their work was further enriched by considering trending topics in Twitter [7]. Provided recommendation of news was not only personal, but also trend-aware. Conclusions of the research of Gao et al. were that personalized recommendation is more important than trend-aware recommendation, but integrating trend-aware and personalized recommendation can improve recommendation results.

In our work we incorporate trend-awareness and personalization similarly to Gao et al. [7]. On top of that we use location-awareness to improve the results, thus the user model is more precise. The idea is based on the assumption that employing location of trends improves the quality of user model. In other words, we believe that applications that incorporate our proposed user model will have more precise results compared with applications employing traditional location-not-aware user models.

The rest of the paper is structured as follows. In section II we present our enhanced trend-aware user model. In section III the evaluation confirming our hypothesis is described. In section IV we sum up our work and provide conclusions

II. ENHANCED TREND-AWARE USER MODEL

Location-awareness is a natural phenomenon that we need to reflect in user modeling. Users are influenced by their context (including geolocation), which affects their decisions. This also means that users are interested in and access web documents containing topics, things, events that relate to or frequently occur in their nearest environment.

Our hypothesis is that location-awareness improves the quality of a user model. Location-aware approach is new aspect in personalized trend-aware news recommendation in microblogs. To our best knowledge, it was not exploited in any previous work before.

We formally define our user model by following the work of Gao et al. [7] and extending it with location aspects. We define the user model as follows

$$P(u) = (c, l, w(u, c, l)) | u \in U, c \in C, l \in L \quad (1)$$

where c stands for concept, w for weighting function, u for user and l for location. C , U , L represent a set of all considered concepts, users and locations, respectively.

In this definition, we capture concepts weights in relation to different locations they can be associated with. We introduce location l , which means that every concept and user belongs to quadtree region and its parent regions.

In addition to user model, we also extend definition of trend model introduced in [7]. We define location-aware trend model as follows:

$$T(I_j) = (c, l, w(I_j, c, l)) | c \in C, l \in L \quad (2)$$

where T is trend model for a time interval I_j , w is a certain weighting function, C and L are set of all concepts and locations related to time interval I_j , respectively.

Location-aware trend model is computed for every time interval I_j for each location l . The idea behind trends is that we model the characteristics of the trend, where the user is located in particular time and then we suppose users are influenced by those time and location-aware trends. It was shown that combination of user and trend models better describes user interests and reflects into improved recommendation [7]. In addition, incorporating trend model into users combined model solves cold-start problem, which emerges when we do not have any or enough information about a new user, who comes to the system.

The combined user model is defined as follows:

$$\vec{m}(I_j, u) = d * \vec{p}(u) + (1 - d) * \vec{t}(I_j) \quad (3)$$

where $p(u)$ is the user model and $t(I_j)$ is the regional trend model. In this model that is computed in every interval I_j for every user u , we compute combined model from equations (1) and (2). The parameter d is trend influence, a configurable parameter, where $d=1$ means that combined model consists only of user model and $d=0$ means combined model consists of trend model only.

We use TF-IDF and t-TF-IDF measures [7] for region and time as a weighting function w in equations (1) and (2). t-TF-IDF is a time-sensitive modification of standard TF-IDF. It uses temporal stability of concepts in form of computing a standard deviation of appearance of concepts in time quanta. Concepts that are more stable appear equally over every period and their weight is decreased in comparison with concepts that appear mainly in short period of time and then they disappear. Such concepts weight is increased in t-TF-IDF. In this way we use t-TF-IDF to capture trends.

A. Location Modeling

The principle of location awareness is that user model is modeled in regions. Weighting of concepts is done per region. Regions enabling location-awareness resemble divide and conquer strategy of algorithms that was proven effective over time. We use regions, where the computation

is done and the results are then aggregated. In location-aware user model it means that we compute weights of concepts in every quadtree region as seen in Figure 1. Note that aggregation is done only with regions and its parent node regions. That is because we compute some aspects of the user model in regions with different size and based on aggregation function we can also weight results in regions based on the importance of location aspect, so we can weight one size of region more to improve the user model even more. This region size weighting (aggregation function) can be calibrated based on feedback from real results to further improve the user model.

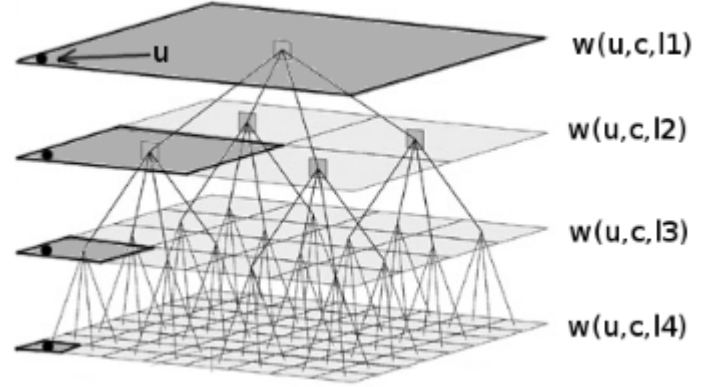


Fig. 1. Location-aware concept weighting based on quadtree regions

Our user model is defined as vectors of weighted concepts partitioned in PR-quadtree regions [5]. PR-quadtree is a tree that has one root and every node of the structure has 0 or 4 children. It is used in geographic information systems because of its advantages that were the reason why we chose this data structure over the others

- every region have similar number of entities,
- live partitioning,
- conserve space,
- parent region connection,
- fast location to region lookup.

Connection between child and its parent nodes is important to fast location lookup, when we want to add new entities (users) to a region map, or to be exact, in microblog domain, when the user tweeted from a location. We do not need to scan all the regions and then find those that matches, it is sufficient to traverse a tree from main node to its child using quadtree structure, this operation consumes $O(\log n)$ time. We can see an example of a visualisation of quadtree in Figure 2 that was created for users of Twitter with locality obtained by geotagging locality field in user profile. Those users were selected according to their high tweet count. It resembles the world map and its density of population. The parts of the map where many users are located (or associated with), are smaller regions; the map is more partitioned. It means for the user model that users living in more crowded quadtree regions have user model that is based on smaller surroundings. It is important that the quadtree is

actually a tree, thus every user is modeled in its closest surrounding region, but also in parent regions ending in the global node. That supports the idea that the user is affected by its surrounding with various dimensions. For example, there are news relevant only for one city, other news for country and some are independent on region, so there are no traces that news are significant in a particular region more than in other. Our location-aware model uses at most $M(\log n)$ times more data than the traditional model (n is maximal number of regions). It means that all the processing is a bit slower. When considering the improvements of user model we believe it is a very good trade-off.

III. EVALUATION: NEWS RECOMMENDATION

The location-aware enhancements we propose retain generality of the user model. However, since it was created while focusing on recommendation of web content (user links or news in a particular locality or global news), we evaluated the model with respect to this purpose. We use the defined user model for recommendation, i.e., we deal with a ranking problem, how to provide a user with ordered list of weighted links (to web content) based on their relevance to the user (with respect to time and location). We quantitatively evaluated our user model on the UMAP2011 Tweets dataset 1 acquired by Abel et al. [2] by performing a synthetic evaluation. The dataset contains 2 316 204 tweets posted by 1619 users. We performed the evaluation as a sequence of the following steps:

- User model acquisition
- News recommendation
- Results evaluation

In the first step we acquired user model by preprocessing tweets, enriching tweets based on link analysis and subsequent location-aware user model creation. Then we simulated recommendation of news for users in the dataset and evaluated the results by applying traditional information retrieval measures. Due to size of social networks the number of users and amount of content they produce is enormous. Hence the scalability of our algorithms had to be considered. We chose MapReduce programming model as a platform for evaluation. We used Google Hadoop 2 implementation and Hive for SQL-like syntax. In evaluation it showed to be a good decision, because it was about 30 times faster on provided cluster than single-threaded solution.

A. User model acquisition

Tweets preprocessing. In this step we obtained entities and topics, links, users and locations from tweets using custom JSON parser and semantic service OpenCalais 3. UMAP2011 Tweets dataset contains tweets description we used. We used web service OpenCalais for every tweet and the result obtained was JSON-formatted text containing semantic information about tweets. We parsed entities and topics from obtained texts. In order to obtain localities,

we got user identifiers contained in dataset and then we questioned Twitter API service for the Twitter user profile. We parsed text location from JSON output and we used batchgeo 4 to obtain user locality. Tweets enrichment. User tweets often point to web content that contains information potentially relevant for user model [3]. There are 1 066 929 links in the dataset we used. We obtained the text of those links and its topics and concepts using SemanticProxy 5 service. This service reads the content of links from tweet dataset, it filters header, footer, navigation and other irrelevant content from HTML and then it extracts the entities with probability and the topic of link. Those entities and topics were added to users tweets as an enrichment with a goal to improve the user model, because people are usually interested in content of links they tweeted, so it characterizes them better. Then we downloaded entities and topics from the fetched content using OpenCalais (Note that there is also UMAP 2011 News dataset, a dataset of news crawled from RSS of nyti.com, bbc.com and cnn.com: there are 77 860 news that consist of 1 896 328 entities. However, only 12k entities were linked with actual tweets using exact match, hence making this dataset insufficient.) User model creation. We assigned users a location they tweeted from. Users tweeted from all over the world; however, only 66 % of them were successfully geo-tagged by batchgeo. We further used only tweets and links from those users to show location aspect of user models. We assigned users with location to regions represented by PR-quadtrees using PR-quadtrees creation algorithm (see Location Modeling in section II). We decided to use at least 100 users per region and then we filtered out those containing less than 20, because too small regions would be useless (we would create regions based on too few users and regions would model particular users, not common local characteristics). We chose those parameters based on characteristics of the UMAP2011 Tweets dataset. Finally we had 53 regions in PR-quadtrees. When the model is being created, tweets with enriched metadata are arranged to regions based on location and time period. We used one week as time period. After weighting user models were created.

B. News Recommendation

We used a general synthetic evaluation approach used in machine learning, where we created user models from training data and then tested (evaluated) those models on other, test examples. In our work, we used tweets from first nine week periods of time for user model creation and the last one week period for recommendation. The same approach was employed by Geo et al. in [7], whom we want to compare with. We recommended top n links from testing dataset to our combined user models. Then we checked, if users actually posted a link that matched one of the recommended links. In related work authors typically use information about website accesses (e.g., logs of visited sites obtained from yahoo toolbar [12]). However, in the UMAP2011 datasets such information is not available. In fact, Gao et al. [7] used UMAP2011 News dataset that

is much smaller than all links, which we consider. They also linked news to tweets by utilizing similarity measure to find most relevant links for tweets instead of exact match between a tweet and a link, i.e., placing that link into the tweet. In our approach, we know exactly what links are contained in tweets so we do not rely on less accurate information. To generate recommendations we used cosine similarity that is commonly used in recommendation systems. Despite its simplicity this method is sufficient since our aim is to evaluate and compare models and not to devise most accurate recommender. In order to compute similarity (suitability) of a web page N_j with respect to the user model M_i , we use the following equation:

$$\text{sim}(M, N) = (c, \text{relevancy}(c, \text{url}))|c \quad (4)$$

where model of web document N is represented as a vector of concepts with relevancy retrieved from content of links by SemanticProxy. The recommendations were generated based on location-aware model, for every region the user belongs to. The smaller the region was, the more the user model was aware of local trends. We used an averaged recommendation from all recommendations from every region of a user. As our approach is based on location modeling based on composite regions (each region has its quadtree parent node regions), we capture local trends, but we are also aware of global trends. Since we obtained many recommendations for each user, we selected top- n to select only the most relevant web news documents for the user. The result set consisted of triplets: user, link and relevancy for each of 962 users (involved in both training and testing). The number of recommendations for each user varied based on actual value of n parameter. To evaluate the results, we used standard measures used for recommendation evaluation, such as precision at n ($P@n$), recall (R), F-measure ($F1$) and Mean Reciprocal Rank (MRR), which show various aspects of quality of the model:

$$P = 1R = 2F = 2 \quad (5)$$

where U is set of all users, $\text{RET}(u)$ is set of retrieved documents for user u and $\text{REL}(u)$ is set of relevant documents for user u . For mean reciprocal rank we borrowed definition from [14]:

$$\text{MRR} = 1 \quad (6)$$

where Q is a query. In our case, a query constitutes a user model involved in recommendation. In MRR we weight quality by rank (position) of first relevant document in ordered recommendation list. It is important to note that our evaluation and combined user model had more parameters (see Table I). When evaluating news recommendation we performed several simulations in order to determine influence of parameters.

Table

To compare the quality of model we create our proposed location-aware model and existing global model [7]. We used entities or topics as concept types from OpenCalais for user and trend models [2]. As a weighting function we used TF-IDF and t-TF-IDF [7]. We also experimented with different influence of trends on the combined model (parameter d), which defines ratio of user and trend model. Due to the finding of Gao et al. [7], user model in combined model is more important than trend model, so we focused the d parameter close to 1. In order to observe characteristics in relation to a number of recommendations there is a top- n parameter.

C. Results and discussion

We evaluated local and global combined model performing simulations with together 160 combinations of the parameters by applying the described measures. First, we compared influence of the parameters on the F-measure. The best and worst values of F-measure for particular setup are depicted in Figure 3. The most important parameter influencing the recommendation was trend influence d in combined user model. It was revealed that model based on users interests was 8 times better in average than model when only trends are considered according to F-measure. Entity-based models were twice as good as topic-based models. Number of recommendations $n=5$ was 1.5 times better than $n=100$ in average. t-TF-IDF improves the model as much as 4 % when compared with TF-IDF. Location awareness of models improved models by average by 2 %.

The presented F-measure values were averaged for various setups to give a basic picture of results. Particular models were much better, although some setups were worse than others. We focused on location awareness of model and analyzed those results in more detail. Tab. II shows a comparison based on Precision, Recall, F-measure and Mean Reciprocal Rank. Location-aware model improved in average Precision, Recall and F-measure increased by 2 %, while MRR decreased by 1.7 %. MRR measures the position of first relevant result. We suppose MRR was decreased because there were new relevant results in local model given for users that had 0 relevant results in top- n recommendation list based on global model. In case of location-aware models there were new results with lower position that decreased MRR measure. As F-measure was improved, decreased MRR does not necessarily mean worse model. MRR using all recommendations evaluated on models with parameters with best F-measure showed small improvements.

We found that the best model according to Precision, Recall and F-measure is entity-based location-aware t-TF-IDF user model ($d=1$). In Figure 4a we can see its

behavior when n parameter changes. We can observe a standard Precision and Recall pattern where Precision is decreased and Recall increased when n is increased. The best recommendation according to F-measure was obtained for $n=10$. The best location-aware improvement of the model was achieved for combined model ($d=0.8$) with t -TF-IDF as weighting function, topic concept type with $n=10$. It exceeds 27 % of improvement when considering Precision (see also Figure 4b). Trend models were improved more than user based models. Location-aware models were equal or better in 67.5 % than global models according to Precision, 77.5 % by Recall, 67.5 % by F-measure and 52.5 % by MRR. We found that Mean Reciprocal Rank was the best when $d=0.4$ in combined model. In [7] similar results were reported so we confirm these findings and we conclude that combined user model consisting of user model and trend model improves MRR measure. The best model yielded 4 % precision. It is important to note that there was a huge information overload in this dataset and we were recommending content to every user from 962 users. It is important to discuss limitations of evaluation we used. It could possibly result in even better results. Recommended content was marked as relevant only when there was exact match in URL. We recommended also other relevant content that was evaluated as irrelevant but in fact, it could be relevant as well. Twitter users often used URL shortener services (such as <http://bit.ly/>). We recommended URLs linked to relevant content, but shortened by shorteners. In our evaluation we did not link shortened URLs from more shorteners pointing to the same content. In addition, there was more content that was not exactly the same, but it was similar, e.g., story about some company mentioned in bbc.com, nyti.com and cnn.com. In our experiments we considered that this content was not the same (only one recommended link was evaluated as relevant). That means that the Precision and other measures were actually better than evaluated. However, our evaluation table

plan was consistent across all the models evaluated, so it was appropriate for comparison of those models.

IV. CONCLUSIONS

In this paper we proposed a location-aware user model for web content recommendation. We followed the work of Gao et al. [7] and researched how location aspects of both users and trends (represented by user and trend models, respectively) relate to the quality of combined user model and how it affects recommendation of web (news) content. We performed an evaluation of the combined user model with various parameters. We confirmed our hypothesis that location-awareness can improve the quality of model, as much as 27 % for best setup and 2 % in average. The best user models created when considering Precision, Recall, F-measure and Mean Reciprocal Rank were location-aware models using t -TF-IDF as weighting function, i.e., those

considering temporal characteristics reflecting trends. We also found that personalization based on user is 8 times better than one based on trend (in terms of F-measure). However, mix of user and trend in combined model can improve Mean Reciprocal Rank. We consider the results we obtained very reasonable. They show that location aspect in user modeling is very important especially in large scale systems such as nowadays very popular microblogs. We believe the importance of location-aware user modeling will be even more increased with the huge boom of smartphones and tablets with even higher support for location data production and utilization.