# Mining Context-Aware Significant
# Travel Sequences from Geotagged Social Media

**Abdul Majid, Ling Chen, Hamid Turab Mirza, Ibrar Hussain & Gencai Chen**

College of Computer Science
Zhejiang University
Hangzhou, 310027, P.R. China
{majid, lingchen, hamid306, ibrar & chengc}@zju.edu.cn

## Introduction

Geotagged photos of users on social media site, i.e., Flickr provide plentiful location-based data, which has been exploited for location-based services, such as mapping geotags to places (Kennedy et al. 2007), and recommendation of personalized landmarks (Shi et al. 2011). As users' preferences to visit a location or multiple locations in a certain sequence could be affected by their current temporal, and weather context. Existing methods addressed queries either with free of context constraints or with a few dimensions of context. This paper considers the problem of mining context-aware significant semantic travel sequences from geotagged photos.

Our major contributions in this paper include: 1) A method that leverages the collective wisdom of people from community contributed geotagged photos collection to provide a set of travel sequences that are *significant* and match the user's current *context*. Our method categorizes context data to support complex context based queries and defines a reasonable function to generate a weighing of locations' significance that are used to score the travel sequences. 2) Evaluation of the performance achieved against an actual large-scale geotag dataset held by Flickr.

## The Proposed Method

We find the tourist locations using spatial proximity of photos and enrich the aggregated locations with semantic using textual tags annotated to photos in combination with information provided by online Web services. Profiles of locations are built to describe the contexts in which they have been visited. For temporal context, the temporal tags annotated to photos are exploited. Whereas, to derive weather context weather Web services are queried to retrieve weather conditions. We exploit the relationships among users, locations, and location categories to define significance of each location using a user-expertise model. Temporal information is analyzed to map and construct user travel sequences that define relationships between locations that users have visited. We give travel sequence recommendation to a new user based on context and significance score of locations.

**Finding Tourist Locations:** To find highly photographed locations as the tourist places in a city, P-DBSCAN (Kisile-

vich, Mansmann, and Keim 2010) is employed to cluster photos using their spatial proximity. Given a collection of photos $P$, The output of a P-DBSCAN is a set of locations $L = \{l_1, l_2, \cdots, l_n\}$. Each element $l = (P_l, g_l)$, where $P_l$ is a cluster of photos and $g_l$ is the geographical coordinates to represent the centroid of location $l$ and is computed from group of geotags annotated to photos in the cluster $P_l$.

**Semantic Annotations:** Textual tags of photos in cluster $P_l$ are utilized as in (Kennedy et al. 2007) to enrich location $l$ with *"name"*. To infer the *"category"* of location $l$, Web service such as *www.google.com/places* is queried to provide metadata (*name* and *type*) of Point of Interests (POIs) that are within the area with radius $r = 200$ meters from coordinates $g_l$. Google Places supports 126 types to describe POIs for search queries. We further generalize these types into 6 categories, i.e., education, shopping, religious, food, transportation, cultural, and entertainment. We select the *"category"* that is highest in frequency in the list of types associated with returned POIs as *"category"* of $l$. After annotation, a tourist location can be represented as $l = (P_l, g_l, name, category)$.

**Profiling Locations:** To build the profile of a location $l = (P_l, g_l, name, category)$, first we identify the visits made to this location by different users. Photos $P_l$ are used to infer the set of visits $V_l$ for location $l$. As a user can take more than one photo in a same visit, each user's photos are sorted using photos taken time. If the difference between the taken times of two consecutive photos is less than a visit duration threshold $visit_{th}$, both photos are considered belong to the same visit. We use the median of time-stamps of photos that belong to visit $v$ as the visit time $v.t$. This $v.t$ is used to retrieve weather conditions $w$, when visit $v$ was made by user $u$ at location $l$. Weather Services normally publish weather data at hourly, daily or monthly level that contain different variable like temperature, precipitation, etc., to describe the weather conditions. We define a context abstraction strategy to obtain abstract context concepts from raw contexts, i.e., time stamp and weather variable. E.g., the raw context $(21:30, 25C^o)$ can be abstracted to (night, warm).

Given the set of visits $V_l$ belongs to location $l$ with associated context concepts, the context concepts are considered *"popular"* that are higher in frequency. For example, $p(l.w) = (warm, sunny)$ depicts that location $l$ has been popularly visited in warm and sunny weather conditions. Af-

ter building the profiles of all locations, we maintain a locations database $LDB = \{l_1, l_2, \cdots, l_n\}$.

**Finding Users' Travel Sequences:** A travel sequence $s = \{l_1, l_2, \cdots, l_n\}$ can be taken as a trip made by a user to visit a sequence of locations in a temporal order, i.e., $l_{i+1}.t > l_i.t$, where $i = (1..n)$ represents the position of location in sequence $s$. To extract the travel sequences for each user: 1) The time-stamps annotated to user's contributed photos are exploited to sort the photos in order to yield his/her traveling history; 2) We split travel history into travel sequences if the difference in the time stamp of two consecutive photos is greater than a given threshold $trip_{dur}$; 3) Each photo is replaced with its corresponding tourist location from $LDB$. If two consecutive photos represent the same location then only one photo is taken into consideration. After the extraction of travel sequences for all users, we build a travel sequence database $SDB = \{S_1, S_2, \cdots, S_n\}$ where $S_i$ represents the set of trips made by user $i$.

**Building User-Location-Category Graph:** We organize three entities (users, locations, and location categories) and relationships (visits) among these entities into a meaningful data structure, i.e., user-location-category tripartite graph $G_{ULC} = (U; L; C; E_{UL}; W_{UL}; E_{UC}; W_{UC}; E_{LC})$, where $U$, $L$ and $C$ are nodes to represent users, locations and location categories respectively. $E_{UL}$ and $W_{UL}$ are sets of edges and edge weights between $U$ and $L$ to represent users' visits and the number of visits to particular locations. $E_{UC}$ and $W_{UC}$ are sets of edges and edge weights between $U$ and $C$ to represent users' visits and the number of visits to particular location categories. $E_{LC}$ are edges between $L$ and $C$ to describe the categories of locations.

**Mining Significance of Locations:** Significance of locations is mined using a user-expertise model, i.e., generate a weighing of location significance through the number of user visits to specific location categories. (1) Given $m$ users and $n$ categories, we build an $m \times n$ adjacency matrix $M_{UC}$. Each entry in $M_{UC}(p, q)$, depicts the experience of user $u_p$ in location category $c_q$. (2) To capture the relationship between users and locations, given $m$ users and $k$ locations, a $m \times k$ adjacency matrix $M_{UL}$ is built. So in matrix $M_{UL}$, $I_j = \{i\}$ where $M_{UL}(i, j) \neq 0$, is the set of indices of users who have visited location $j$. Significance score of each location $l_j(j = 1 \cdots k)$, which is of location category $c_q(q = 1 \cdots n)$ is computed by: $g(l_j) = \sum_{i \in I_j} M_{UC}(i, q)$.

**Recommendations:** Processing of context aware query $Q(t, w)$ made by user $u_p$ proceeds as a two-step approach: an initial filtering step retrieves travel sequences belong to target city from $SDB$ that meet the contextual constraints given in the query, thus producing a filtered set of sequences $S'$. In the second step: significance score for each travel sequence $s \in S'$ is computed by aggregating the significance score of locations it contains as: $r(s)) = \frac{1}{n} \sum_{i=1}^{n} g(l_i)$. Next the travel sequences in $S'$ are ranked using score $r$ and top $k$ travel sequences are returned as query result.

## Evaluation

**Data Collection:** Public API of Flickr (*www.flickr.com*) was used to collect 736383 geotagged photos that were taken

Table 1: Performance Comparison

|  | FrequentRank | ClassicRank | Our Method |
|---|---|---|---|
| nDCG@5 | 0.713 | 0.770 | 0.853 |
| nDCG@10 | 0.735 | 0.804 | 0.892 |

in six different cities of China between January 2001 and July 2011. Historical weather data of these cities was collected using public API of an online weather Web service *www.wunderground.com*.

**Baseline Methods:** As baseline methods, we used (a) Frequent-Rank (Yin et al. 2011), that employs PrefixSpan (Pei et al. 2001) to extract frequent sequential patterns from $SDB$ whose frequencies are not smaller than the given minimum support threshold $min_s = 2$ and (b) Classic-Rank, a method as given in (Zheng et al. 2009) to score the travel sequences by exploiting reenforcement relationships between users and locations from matrix $M_{UL}$.

**Results:** A large quantity of user queries were simulated with different contextual constraint settings. Four experts (tour guides) who are well familiar with these cities were requested to evaluate the output of our and baseline methods and provide the feedback. They were asked to evaluate top 5 and top 10 travel sequences using three scores, i.e., Very much interesting (3), interesting (2), neutral (1), not interesting (0). Table 1 depicts the effectiveness of our and baseline methods in terms of normalized discounted cumulative gain (nDCG@5 and nDCG@10). Results show that, as compared to baseline methods, our approach made recommendations that are more relevant (paired t-test with $p < 0.05$).

## References

Kennedy, L.; Naaman, M.; Ahern, S.; Nair, R.; and Rattenbury, T. 2007. How flickr helps us make sense of the world: context and content in community-contributed media collections. In *Proc. of Multimedia'07*, 631–640.

Kisilevich, S.; Mansmann, F.; and Keim, D. A. 2010. P-dbscan: a density based clustering algorithm for exploration and analysis of attractive areas using collections of geo-tagged photos. In *Proc. of COM.Geo*.

Shi, Y.; Serdyukov, P.; Hanjalic, A.; and Larson, M. 2011. Personalized landmark recommendation based on geotags from photo sharing sites. In *Proc. of ICWSM*.

Yin, Z.; Cao, L.; Han, J.; Luo, J.; and Huang, T. S. 2011. Diversified trajectory pattern ranking in geo-tagged social media. In *Proc. of SDM*, 980–991.

Pei, J.; Han, J.; Mortazavi-Asl, B.; Pinto, H.; Chen, Q.; Dayal, U.; and Hsu, M. 2001. Prefixspan: Mining sequential patterns by prefix-projected growth. In *Proc. of ICDE*, 215–224.

Zheng, Y.; Zhang, L.; Xie, X.; and Ma, W.-Y. 2009. Mining interesting locations and travel sequences from gps trajectories. In *Proc. of WWW'09*, 791–800.