

ON THE VALUE OF PURPOSE-ORIENTATION AND FOCUS ON LOCALS IN RECOMMENDING LEISURE ACTIVITIES

BEATRICE VALERI, FLORIAN DANIEL, FABIO CASATI

*DISI, University of Trento, Italy
via Sommarive 9, 38123 Povo, Trento, Italy
valeri@unitn.it*

FLORIAN DANIEL

*DISI, University of Trento, Italy
via Sommarive 9, 38123 Povo, Trento, Italy
daniel@unitn.it*

FABIO CASATI

*DISI, University of Trento, Italy
via Sommarive 9, 38123 Povo, Trento, Italy
casati@unitn.it*

Received (received date)

Revised (revised date)

Recommender systems are omnipresent today, especially on the Web, and the quality of their recommendations is crucial for user satisfaction. Contrary to most works on the topic, in this article the authors do not focus on the algorithmic side of the problem and instead study the importance of the data in input to the algorithms. They study the case of restaurant recommendations for locals and show that fine-tuned data and state-of-the-art algorithms easily outperform, for instance, TripAdvisor. The findings make a case for better-thought and purpose-tailored data collection techniques.

Keywords: Recommender systems, data collection, mobile recommendations, restaurants, TripAdvisor

Communicated by: to be filled by the Editorial

1. Introduction

Recommender systems are software systems that, given a set of items and a user, aim to predict the user's interest in the items and to suggest the user which items to inspect, use or buy in a given context. At the core of each recommender system there are two ingredients: first, the recommender *algorithms* that select candidate items; second, the *data* that provide the base for the recommendations [1]. In general, the better the algorithms and the more the data available, the better the recommendations.

Algorithms are the traditional focus of research. They can be split into two main classes: collaborative filtering [2] and content-based recommendation [3]. The former is based on logs of how users interact with items (e.g., if they read, rate, comment, like, buy items) and look for *users that behave similarly* to the target user; the latter is based on the descriptions of both

items and users (their profile and preferences) and look for items that have similar features to the ones the target user already liked in the past. Hybrid techniques bring both approaches together.

Data is often represented by ratings on items, and literature on recommender algorithms typically uses standard datasets for the assessment of algorithms, such as the 5-star ratings for movies by *MovieLens* [4] or *Netflix* [4, 5]. Other common rating scales are unary (like), binary (thumbs-up/thumbs-down) and 3-values scales (thumbs-up/neutral/thumbs-down). Ratings come as values with no extra information about the experience or context of the user that led to the judgments, leaving the interpretation of the ratings to the recommender algorithms. In our own work on the recommendation of leisure activities (drinking, dining, dancing) [6], we however noticed that the quality of recommendations on items (in our cases, restaurants and bars) strongly depends on the specific purpose of the activity: one place may be good for dining but not for drinking; another one may be good for a romantic dinner but not for one with friends. This kind of nuances is usually not captured by state-of-the-art recommender systems.

In addition, in the specific context of restaurants and dining, tourists have generally very limited knowledge of a new city and, therefore, most of the times rely on recommendations from family members or friends; if these are not available, tourists also like to rely on recommendations from *locals* [7]. They see *locals* as knowledgeable and trustworthy, since they know that locals know most of the available options, have been there themselves, and can provide personal recommendations for free. Again, this kind of knowledge from locals is typically lost in online tourist portals, which are mostly oriented toward and, hence, visited by tourists themselves. Good recommendation services specifically tailored to locals are still underrepresented.

These considerations raise the question whether data collected with i) special attention to locals and ii) information on specific usage purposes in mind can recommend restaurants with higher *precision* (that is, higher probability of recommending a place the user will actually like) than generic recommender systems, such as *TripAdvisor* for restaurants. In order to answer this question, this article studies the case of collaborative filtering algorithms and tests the following two hypotheses:

- H1: Data collected from locals and state-of-the-art, personalized recommendation algorithms produce recommendations of higher precision than generic recommender systems, such as *TripAdvisor*.
- H2: Recommendations computed from purpose-specific data outperform *TripAdvisor*.

The findings show that indeed taking these aspects into consideration can lead to improved recommendations with respect to the mainstream recommender in this domain. The findings also unveil insights that are particularly important to mobile recommender systems characterized by limited screen real estate.

2. Recommender Systems for Leisure

Collaborative filtering is very successful in the e-commerce sector, but for the leisure sector we have to take into consideration its specific characteristics: differently from item recommendation in e-commerce, *context* and *location* are really important. First of all, only places close

to the user can be experienced and people do not usually travel much to find them (within 14 miles from their house [8]). Other contextual information helps finding more interesting results: Baltrunas et al. [9] show that recommendation systems are able to increment user satisfaction by considering weather conditions, companions, time (season, weekday or time of the day) and familiarity with the area, along with other contextual information. Mobile devices provide support for context-aware recommendation systems. Thanks to their sensors, they can automatically collect contextual information, such as user position and therefore weather conditions, and also provide for proactive recommendations [10].

The need for support in searching interesting places is widely recognized and many services have been developed with this goal. TripAdvisor and Gogobot are two of the most used recommendation services in the travel sector, while Yelp and Google Local are more focused on locals, providing “yellow pages”-like services. These recommendation systems collect ratings using a 5-star scale or a variation of it, giving the users the possibility to express their experience with different nuances of satisfaction/dissatisfaction. They are also very popular: TripAdvisor counts 340 million unique monthly visitors [11], while Yelp counts 142 million unique monthly visitors [12]. The use of mobile devices has been growing steadily: 190 million people downloaded the TripAdvisor mobile app, and 50% of accesses to TripAdvisor and Yelp come from mobile devices (both smartphones and tablets) [11, 12].

Recently, Foursquare added local search functionality to its application [13], using both the check-in information of the location-based social network and user feedbacks, reviews and tastes to provide recommendations. Other services specifically focus on restaurants, such as The Fork (user-provided ratings) and Zagat (expert-provided ratings).

3. Background

In a recommender system, users express their opinions about items in form of ratings. *Items* can be anything users can experience and can have an opinion about, while a *user* is any person that has experienced some of the items the system is focused on. In this paper we focus on restaurants. These are physical establishments, so only people able to visit them can also experience them. From this perspective, a restaurant can have two kinds of customers: *locals* and *tourists*. *Locals* are people that live in the area, are familiar with the local cuisine and the restaurants, and can experience them several times. *Tourists* are visitors for business or leisure that generally have fewer chances to sample restaurants in a given area and are less familiar with the local cuisine.

When rating items, users commonly assign one rating per item, evaluating it according to their overall experience. *Multi-criteria* ratings, instead, ask users to add one rating for each of a set of predefined characteristics of the item. For example, in the case of restaurants, the criteria could be food quality, drink quality, service and popularity. This requires a user to consider the different aspects of her experience and give more ratings.

Orthogonally to this, a restaurant may be perceived differently depending on the *purpose* of the visit: the choice of a restaurant for a dinner with friends may differ from what we would choose for a romantic dinner or for a quick lunch. We already verified this hypothesis in our earlier research [6], where we also identified four main purposes: dinner with tourists, romantic dinner with the partner, dinner with friends and price/quality ratio (e.g., important for a lunch break).

4. Method

In this article we study whether recommendations based on *purpose-specific data* collected from *locals* outperform recommendations computed from the typical data collected from tourists by tourist portals (TripAdvisor).

4.1. Data Collection

We collected ratings using a 3-values thumb-up/thumbs-down scale for each of the purposes identified in [6] (dinner with tourists, dinner with partner, dinner with friends and lunch break): users can specify whether they like or don't like a restaurant, or are neutral about it. In general, the thumbs-up/thumbs-down leaves less space for controversy than using 5 stars: the user just has to think about whether the item is good or bad, without having to think about how good (or how bad). The neutral rating prevents forcing the user to like or dislike an item if it is considered borderline.

In May 2014 we collected ratings for 50 restaurants in Trento, Italy. We selected this list considering the most popular restaurants according to TripAdvisor that are located in the city center and easily reachable by everyone. We enrolled participants by distributing fliers to locals, sending emails to friends and colleagues, and also involving a small group of university students. The participants (114) were asked to rate restaurants for each of the purposes (the 4 identified above) at a time. The process produced a total of 4706 ratings, with 1529 ratings for “dinner with tourists”, 1113 ratings for “dinner with the partner”, 1112 ratings for “dinner with friends” and 952 ratings for “price/quality ratio”. The restaurants received a minimum of 4 ratings and a maximum of 112 ratings per purpose, while users added a minimum of 0 ratings and a maximum of 49 ratings per purpose, with an average of 11 ratings per purpose.

4.1.1. Recommendation Algorithms

Computing purpose-specific ratings poses challenges to the recommendation algorithms, as the algorithm has to work with multiple ratings per item per purpose. A first way to approach this multiplicity is to filter ratings to create one dataset for each purpose; in this way, only the information about the purpose the requester is interested in is used to compute recommendations. Another way is to merge all ratings from the different purposes and to compute aggregated ratings valid for all purposes, similarly to how multi-criteria ratings are handled by recommendation algorithms. A third solution is to learn user tastes using all collected data for all purposes and to compute ratings for each purpose individually; in this way, the whole information is used to extract taste features and to compute similarities between users or items, but only the ratings specific to a purpose for a user in a given instant of time are used for the prediction of ratings for unknown restaurants.

We followed this last solution. To handle the presence of 4 ratings per user-restaurant pair (one for each purpose), we split each user's ratings for a restaurant into 4 purpose-specific restaurant-purpose pair, resulting in 200 ($50 * 4$) items. In this way, all ratings can be considered in the computation of the model used by the algorithm (like building clusters for cluster-based collaborative filtering or computing matrix factorization for SVD), while only the restaurant-purpose pairs for the requested purpose are considered to build the rank when computing recommendations. To adapt the algorithms to this behavior, we only need

to extend them with a final filter of items by purpose.

For computing recommendations we selected four state-of-the-art, personalized, collaborative filtering algorithms implemented by Apache’s Mahout library (<http://mahout.apache.org>):

- *User-based collaborative filtering* identifies a requester’s neighbors (the users with similar tastes) and uses their ratings and the level of similarity with the requester to compute a prediction of the requester’s ratings for the items she does not know yet.
- *Cluster-based collaborative filtering* pre-groups users into clusters of users with similar tastes and averages the ratings of all users within each cluster to compute a prediction of the requester’s ratings for unknown items. We specifically use hierarchical clustering.
- *Slope One* is an item-based algorithm that leverages on the principle of “popularity differential,” that is, on how much one item is liked more than another. In order to predict the rating of an item, it considers information both from other items rated by the requester (and their ratings from other users) and from other users who rated the item (and their ratings to other items) [14].
- *SVD* is a matrix factorization algorithm that computes ratings out of features automatically extracted from a known, incomplete user-item matrix. The matrix is decomposed into a user-feature, a feature-item, and a feature-feature matrix. Rating predictions are computed as the product of the requester’s row, the feature-feature matrix, and the item’s column.

These algorithms have been selected as they are popular and simple, two properties that allow us to communicate better the effects of the data on recommendation quality. Other algorithms have been shown to perform similarly or even better under certain conditions, but our goal is more that of understanding and communicating the effect with widely known and adopted algorithms. Since all restaurants in our dataset are easily reachable by foot, user location and time (the usual contextual information) are not needed; we consider instead the purpose the requester is interested in.

4.2. Quality Metric

We compare algorithms based on their *precision* (since we don’t have full knowledge of the users’ interests they may have rated only a subset of the restaurants they actually know we cannot compute meaningful recall values). Given a user u , the *list* of computed recommendations, and the purpose p , we compare the performance of the algorithms using the following *precision* metric (following [15]):

$$Precision(u, p, list) = \frac{||Good(u, p, list)||}{||Good(u, p, list)|| + ||Bad(u, p, list)||} \quad (1)$$

where:

- $Good(u, p, list)$ = items in *list* that have been rated positively by user u for purpose p
- $Bad(u, p, list)$ = items in *list* that have been rated negatively by user u for purpose p

For the comparison, we split the users’ ratings for each purpose p into a training set (the ratings the algorithms can use to build the user profile) and a test set (the ratings used to compute the precision of recommendations) with a 70/30 proportion. We tested only users that had at least 6 ratings per selected purpose, leaving at least 2 ratings for testing (the

ceiling of the 30% split), and omitted items the users didn’t express any opinion for. Test ratings were randomly collected half from users’ positive ratings and half from their neutral or negative ratings (to test good and bad predictions). To make the test independent of the computed split of ratings, each query was repeated with 5 different random splits.

4.3. *Algorithms tuning and configuration*

Given this evaluation strategy, all algorithms underwent a dry run to configure them for best performance: we collected $N_p = 5$ recommendations for each purpose from each algorithm and averaged the precision of each list of recommendations (20 per user: 4 purposes by 5 training/test splits). *Slope One* has no parameters, so it was used as it is. *User-based collaborative filtering* depends on the used similarity metric, neighborhood strategy and neighborhood size. We tested Pearson correlation, log likelihood, Spearman correlation, Tanimoto coefficient, cosine similarity, Euclidean-distance-based similarity and Yule similarity. The best precision was obtained with neighborhood selected by similarity threshold, using Yule similarity and similarity threshold 0.3, with a precision of 76%. For *cluster-based collaborative filtering*, we used the same similarity metrics as for user-based collaborative filtering and identified the best configuration in log likelihood similarity and stopping condition expressed as fixed number of clusters, set to 3, with a precision of 70%. The best precision with *SVD* was obtained with 10 features and 30 iterations, with 65% of precision.

5. Results

5.1. *Aggregate precision*

For the comparison, we select the top N_p recommendations from TripAdvisor in the order proposed by TripAdvisor. We vary N_p from 2 to 15 to study the effect of the recommendation set size on precision and compute the precision of our recommender algorithms by averaging the results of all the purposes over 570 individual data points (114 users times 5 random splits per run) per purpose.

For a first assessment of the difference between the dataset underlying TripAdvisor and our own dataset, we compare the recommendations of TripAdvisor with a similar non-personalized, average-based recommendation algorithm using our dataset. This baseline algorithm computes the predictions of user ratings by computing the “lower bound of Wilson’s score confidence interval” (<http://www.evanmiller.org/how-not-to-sort-by-average-rating.html>). This formula computes a confidence interval for the average rating we would obtain if we had all ratings by the full population, starting from a sample of ratings. The lower bound tells “the item is liked at least that much.”

Our own dataset differs from TripAdvisor’s one in four key aspects: (i) 3-value vs. 5-value rating scales, (ii) purpose-based vs. generic ratings, (iii) locals vs. tourists, and (iv) small amount vs. large amount of ratings. Since we don’t have access to the actual dataset and algorithm used by TripAdvisor, we cannot distinguish the effects of each of these aspects, but we can still see in Figure 1 how TripAdvisor generally produces better recommendations than the baseline (except for $N_p = 5$). The key to this better performance most likely lies in the bigger amount of ratings TripAdvisor can rely on.

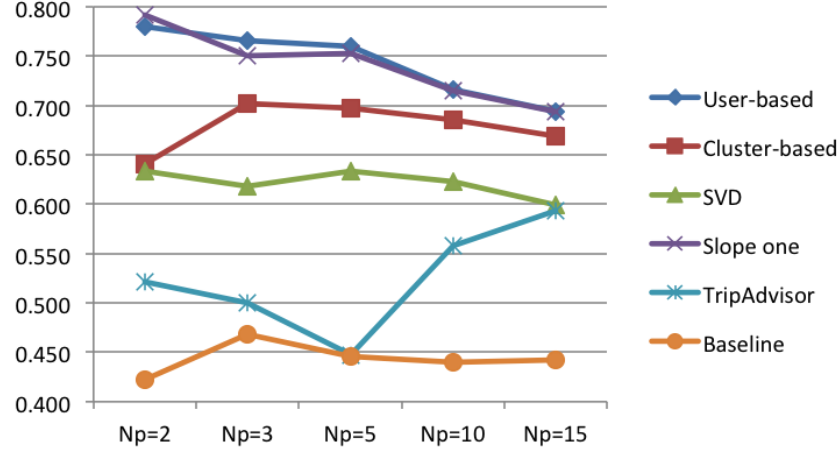


Fig. 1. Precision of the recommendation algorithms for varying Np .

Interestingly, if we now look at the precision of the personalized algorithms, we see that they all perform better than both TripAdvisor and the baseline. Slope One and User-based have the best precision and are very close to each other. Cluster-based is not far from the top recommenders, with only a distance of 2 percentage points in precision for higher Np , while SVD performs worse. TripAdvisor’s precision is highest for $Np = 15$, where it reaches the same precision of SVD, while it still is 10 percentage points lower than the best performance. This shows that as the size of the recommendation set grows, TripAdvisor has higher probability to contain good recommendations. In order to assess the expressive power of the charts in Figure 1, we took the precision values for $Np = 15$ and performed pair-wise t-tests. The tests confirm also statistically what is communicated by the chart visually: except for User-based/Slope one and TripAdvisor/SVD, all precision values are significantly different (p-value ≤ 0.0001 , α -level = 0.05, considering the precision of 1280 recommendation lists for each algorithm).

Overall, Figure 1 shows that the precision of the best algorithm between the chosen personalized algorithms (user-based collaborative filtering) is from 10 ($Np = 15$) to 31 ($Np = 5$) percentage points higher than that of TripAdvisor (from 17% to 68% in relative terms). This makes us accept our hypothesis H1: *data collected from locals and state-of-the-art, personalized recommendation algorithms produce recommendations of higher precision than TripAdvisor*.

This means that even though our dataset is significantly smaller than that of TripAdvisor, the focus on locals and personalization yield recommendations that are of significantly higher quality compared to recommendations computed with a generic algorithm from a much larger dataset. TripAdvisor’s restaurant rank is in fact built using a huge amount of reviews mostly by tourists and specifically focuses on recommending restaurants to tourists. Our experiment aims to understand how to recommend restaurants to locals and shows that locals are a special class of users that are simply more demanding than generic tourists.

We have to keep in mind that these results have been obtained by averaging the precision of purpose-based recommendations. TripAdvisor starts with a disadvantage since it is built for tourists, and its recommendations could be worse for other purposes (as we will see next).

5.2. Purpose-specific Precision

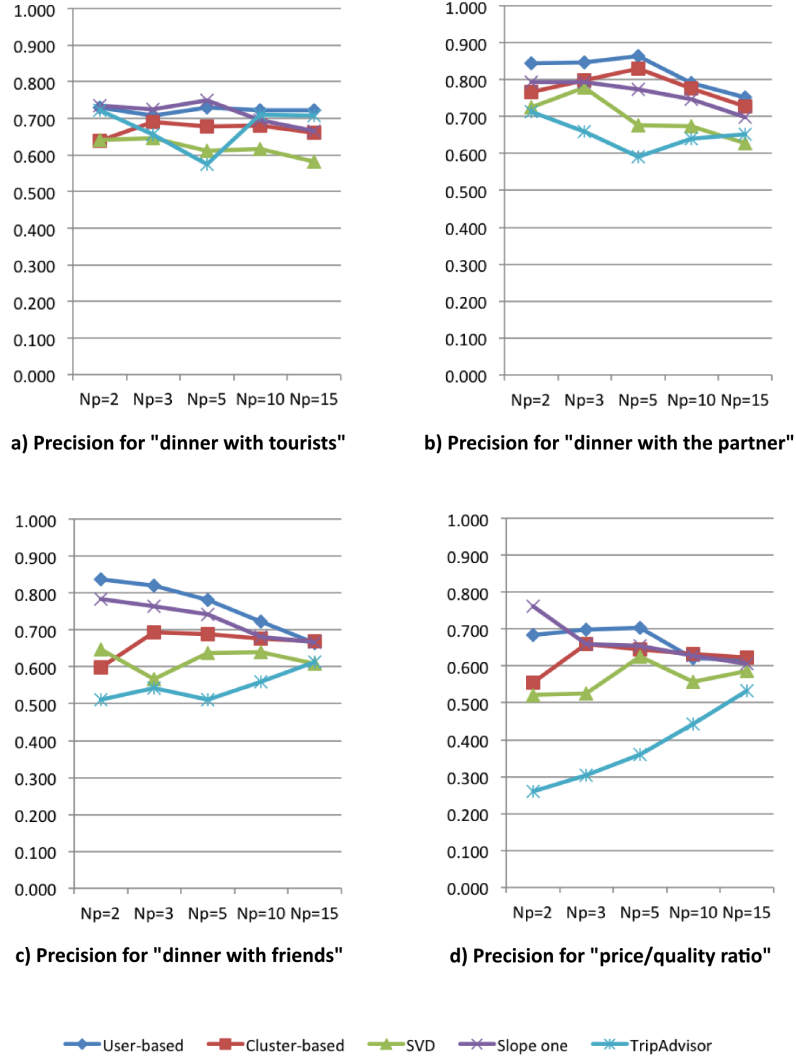


Fig. 2. Purpose-based precision for the five recommendation algorithms.

We now analyze the importance of purpose-specific ratings in recommending restaurants. In [6] we found that the restaurants perceived as good for bringing a tourist are similar to those for a romantic dinner with the partner, while the ones for going out with friends are very different and more related to the price/quality ratio. Next, we analyze concretely how the different recommenders behave depending on the purpose a user has in mind. The test setting of the experiments is the same as above, with the only difference that now we no longer aggregate results and instead keep purposes separated.

Figure 2 reports the precision graphs for each purpose. If we concentrate on TripAdvisor, we see that it provides good predictions for a dinner with tourists, while its precision decreases if the meal is to be consumed with the partner or friends, and it reaches its lowest value if a good price/quality ratio is the target (only 26% of precision for $Np = 2$). The personalized algorithms seem less affected by the purpose, with slightly higher precision for a dinner with the partner and slightly lower precision for the price/quality ratio. Slope One, User-based and Cluster-based collaborative filtering always outperform SVD.

These results clearly indicate that each purpose is different from the others, and algorithms that take care of these differences are able to build better recommendations than generic algorithms. TripAdvisor shows the best precision for $Np = 2$ and “dinner with tourists”, while the worst precision is obtained for $Np = 2$ and “price/quality ratio”, with a difference of 46 absolute percentage points. Purpose-based recommender algorithms have a more constant quality, with less difference between the best and the worst precision: for example, user-based collaborative filtering has the highest precision for $Np = 5$ and “dinner with the partner”, while the lowest one is for $Np = 15$ and “price/quality ratio”, with a difference of 25 absolute percentage points. This minor difference demonstrates a higher quality of purpose-based, personalized recommendations under all circumstances. This supports hypothesis H2 for the purposes dinner with partner, dinner with friends and price/quality ratio: *recommendations computed from purpose-specific data outperform TripAdvisor* for these purposes and may represent a strategic value for competitors of TripAdvisor that want to target locals instead of generic tourists.

TripAdvisor recommendations have instead a high precision for a dinner with tourists, and for this purpose their quality is in line with the ones computed with personalized recommendation algorithms. Given that these latter algorithms use data that stem from locals, this means that locals essentially agree with TripAdvisor on where to bring a tourist and where not. This, in turn, is a quality certificate for TripAdvisor for this specific purpose.

6. Discussion and Conclusion

Our experiments show that providing locals with restaurant recommendations is a tricky endeavor, because providing them with added value compared to generic tourist portals like TripAdvisor asks for advanced personalization, not only based on identity but also on purpose. The experiments further show that if data are collected with an eye on the purpose of the restaurant visit and from locals, even basic algorithms outperform generic recommendations. The improvement in recommendation quality thanks to tailored data is not only significant, but has a big effect size. These results are somewhat surprising, given that also more advanced and precise algorithms are available in the literature. The results however also show that TripAdvisor is still competitive in its own domain, i.e., recommendations for tourists.

At first glance, our comparison of the personalized algorithms with purpose-specific datasets and TripAdvisor may not seem fair: TripAdvisor does not specifically target locals; the available datasets are very different; and its underlying algorithms are not publically available and known. However, at the same time TripAdvisor is one of the key representatives of the state of the art in restaurant recommendations and one that nicely shows a one-size-fits-all approach that works for tourists. What we show in this article is that there is still huge space for improvement and businesses if the focus is shifted from a generic audience to locals.

Yet, doing so really requires a thorough planning of how to collect the necessary data and how to tailor it to the needs of locals. There is no shortcut solution to good data collection (e.g., crawling TripAdvisor or similar), and each new application will have to collect its own data according to its specific needs.

The results of our experiments also reveal another, slightly hidden message that is of particular importance to the world of mobile recommender systems: Mobile devices have typically small screens and are often used in situations in which the user cannot pay full attention to the device. This means that the user can see only few recommended items at a time and may not be willing or able to go through a long list of recommendations [16]. A mobile recommender system is thus particularly challenged even more than a desktop one to compute precise recommendations. The data in Figure 2 show that TripAdvisor performs particularly weakly for small result sets. The lesson is that simply porting a desktop version of a recommendation algorithm to a mobile recommender system may be dangerous, and personalization and data quality become even more important.

One limitation of the study is that our comparison of algorithms is based on the externally visible behavior of TripAdvisor. Its actual, internal algorithm and dataset are not made publicly available. Further, the algorithms we used were trained on the same dataset we also used for testing. TripAdvisor, on the one hand, could rely on a much bigger set of ratings for the training, but, on the other hand, the comparison was again based on our dataset of ratings. One difference between the two datasets is that TripAdvisor uses a 5-stars rating scale, while our dataset uses a 3-values thumb-up/thumbs-down scale. Understanding if the two rating scales affected our findings would require further analysis.

We see space for future investigations that study how to motivate users to provide more ratings (these are the real asset of recommender systems), how to identify fake, paid or differently biased ratings that are not trustworthy (e.g., provided by the restaurant owners themselves), and the effect of special offers and coupons.

References

1. G. Adomavicius, A. Tuzhilin (2005), *Toward the Next Generation of Recommender Systems: A Survey of the State-of-the-Art and Possible Extensions*, TKDD
2. P. Lops, M. De Gemmis and G. Semeraro (2011), *Content-based recommender systems: State of the art and trends*, Recommender Systems Handbook
3. Y. Koren and R. Bell (2011), *Advances in collaborative filtering*, Recommender Systems Handbook
4. F. Cacheda, V. Carneiro, D. Fernández, and V. Formoso (2011), *Comparison of collaborative filtering algorithms: Limitations of current techniques and proposals for scalable, high-performance recommender systems*, TWEB
5. J. Lee, M. Sun, and G. Lebanon (2012), *A comparative study of collaborative filtering algorithms*, arXiv preprint arXiv:1205.3193.
6. B. Valeri, M. Baez and F. Casati (2013), *Come Along: understanding and motivating participation to social leisure activities*, CG
7. R. Rompf, R. B. Dipietro, and P. Ricci (2005), *Locals' involvement in travelers' informational search and venue decision strategies while at destination*, J. Travel Tour. Mark.
8. T. Horozov, N. Narasimhan and V. Vasudevan (2006), *Using location for personalized POI recommendations in mobile environments*, SAINT
9. L. Baltrunas, B. Ludwig, S. Peer, and F. Ricci (2011), *Context-aware places of interest recommendations for mobile users*, Design, User Experience, and Usability. Theory, Methods, Tools and Practice

10. D. G. Vico, W. Woerndl, and R. Bader (2011), *A study on proactive delivery of restaurant recommendations for android smartphones*, RecSys Workshop on Personalization in Mobile Applications
11. *TripAdvisor Fact Sheet*, http://www.tripadvisor.com/PressCenter-c4-Fact_Sheet.html
12. *Yelp Fact Sheet*, <http://www.yelp.com/factsheet>
13. J. LinkLater (2014), *Foursquare: The New And Improved Yelp*, blog.sweetiq.com
14. D. Lemire and A. Maclachlan (2005), *Slope One Predictors for Online Rating-Based Collaborative Filtering*, SDM
15. G. Groh and C. Ehmig (2007), *Recommendations in taste related domains: collaborative filtering vs. social filtering*, GROUP
16. F. Ricci (2010), *Mobile recommender systems*, Inf. Technol. Tour.