Measuring Surprise in Recommender Systems

Marius Kaminskas Insight Centre for Data Analytics University College Cork, Ireland marius.kaminskas@insight-centre.org Derek Bridge Insight Centre for Data Analytics University College Cork, Ireland derek.bridge@insight-centre.org

ABSTRACT

A lot of current research on recommender systems focuses on objectives that go beyond the accuracy of the recommendations; for instance, ensuring that the list of recommended items is diverse. In this work we explore a particular beyondaccuracy objective — serendipity. Existing approaches to measuring serendipity rely on comparing the produced recommendations against a baseline recommender system. We take the first step toward a metric that would allow objectively measuring the surprise (or unexpectedness) of recommendations, without comparing them against an alternative system. We propose two ways to measure surprise, which constitutes the core component of serendipitous recommendations.

Through offline experiments we compare three state-ofthe-art recommendation algorithms in terms of their ability to generate surprising recommendations. For one of the suggested metrics, the results validate the intuition that a matrix factorization approach generates the most accurate but also the least surprising recommendations, while a userbased neighbourhood approach performs best in terms of surprise.

Categories and Subject Descriptors

H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval—Information Filtering

General Terms

Experimentation, Measurement

Keywords

Recommender systems; evaluation metrics; beyond accuracy; serendipity

1. INTRODUCTION

The goal of recommender systems (RS) research has traditionally been focused on accurately predicting users' ratings

ACM RecSys Workshop on Recommender Systems Evaluation: Dimensions and Design, October, 2014, Foster City, Silicon Valley, USA Copyright 20XX ACM X-XXXXX-XX-XXX/XX ...\$15.00. for unseen items. However, accuracy is not the only important objective of recommendation [12]. In recent years the focus of RS research has shifted to such objectives as ensuring the recommended items are novel and the set of items is diverse [17]. These objectives are of particular importance since in real life systems users are most likely to consider only a small set of recommendations. It is therefore crucial to make sure that this set of items is as interesting and engaging as possible.

In this paper we investigate serendipity — an objective denoting the ability of a system to generate recommendations that are both surprising and relevant [9]. Although serendipity is frequently mentioned in the research literature as an important complement to recommendation accuracy, until now there have been few attempts to formally define and measure it. A common intuition in the RS literature suggests that, of the state-of-the-art recommendation techniques, matrix factorization approaches and the item-based neighbourhood technique are less likely to make serendipitous recommendations than the user-based neighbourhood technique [6]. However, to the best of our knowledge, this intuition has not been confirmed by experimental evaluation.

Existing approaches to measuring serendipity rely on comparing the produced recommendations against a baseline recommender system [13, 8, 1]. We propose metrics that do not require two systems, but allow evaluating the surprise of recommendations produced by any existing techniques.

In the rest of this paper, we first review research that addresses the serendipity objective and position our work with respect to the previous efforts. We then propose two alternative ways of measuring surprise, which constitutes the core component of serendipity. Our main and final contribution lies in comparing the performance of three state-ofthe-art recommendation algorithms in terms of their ability to generate serendipitous recommendations. Our findings exemplify the trade-off between accuracy and serendipity, and provide a basis for further research on measuring and optimizing recommendation serendipity.

2. RELATED WORK

The term *serendipity*, meaning a "happy and unexpected discovery by accident", was coined in the 18^{th} century [2]. In the information retrieval (IR) literature, Toms [16] was among the first to recognize the importance of facilitating serendipity in information systems. Other IR studies have also analyzed the process of serendipitous information discovery and the paradox of designing for unexpected results [7, 11]. In the RS literature, Herlocker et al. [9] in-

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

formally defined serendipitous recommendations as "surprisingly interesting items the user might not have otherwise discovered".

2.1 Increasing serendipity

The first attempts to increase the serendipity of retrieved results were reported in the IR literature, e.g., Campos and Figueiredo designed a software agent to support serendipitous information discovery through web crawling [4]. More recently, André et al. [2] suggested viewing serendipity as a combination of chance discovery and the usefulness of discovered information. The authors provided guidelines to design information systems with better support for both components of serendipity: supporting chance encounters and enhancing the user's ability to recognize serendipitous content.

In the RS literature, Iaquinta et al. were among the first to optimize a recommender system for serendipity [10]. They considered items with low similarity to a user's profile as potentially serendipitous recommendations and modified a content-based system to include such items alongside the standard recommendations. Onuma et al. [15] proposed a graph-based algorithm for supporting surprising recommendations. The authors introduced the idea of computing a 'bridge score' of item nodes in the user-item bipartite graph. Nodes connecting separate interconnected areas in the graph receive high bridge scores as they bridge different subspaces in the item information space. The bridge score may be combined with an item relevance score when generating recommendations.

More recently, Oku and Hattori [14] presented a system that induced possibly serendipitous recommendations by selecting items whose content is a mixture of the content features of two items from the user's profile. Adamopoulos and Tuzhilin [1] presented an approach to generate unexpected recommendations by maximizing a utility function which combines an item's relevance and its distance from a set of 'expected' items. The set of expected items includes all items rated (i.e., seen) by the user and items similar to the seen ones (similarity is computed using item content features).

2.2 Measuring serendipity

Although serendipity has been frequently mentioned in the IR and RS literature, few works have provided formal serendipity definitions. This is not surprising, as the notion of an item being *surprising* or *unexpected* is difficult to measure without asking the user's opinion — and the notion that the item is *delightful* is probably impossible to capture. It is commonly agreed that serendipity consists of two components – surprise and relevance [9]. While relevance is easy to measure using ratings or other forms of user feedback, surprise (or unexpectedness) of recommendations is difficult to capture without relying on live user studies which are expensive to conduct.

As discussed above, previous approaches to increase serendipity relied on various heuristics to generate more surprising recommendations, e.g., an item being different from a user's profile [10, 1], an item being connected to distinct area in a user-item graph [15], or an item having a mixture of two input items' features [14]. However, when evaluating the quality of the results produced by their heuristics, a common practice among the authors is comparing the generated recommendations with recommendations produced by a primitive baseline system (i.e., one which is not optimized for serendipity). This approach to measuring serendipity was proposed by Murakami et al. [13], who argued that a primitive method produces easily predictable items, while the goal of a serendipitous recommender is to suggest items that are difficult to predict. This idea was later adopted by Ge et al. [8], who proposed a formulation of serendipity that combines this notion of unexpectedness with item relevance.

This comparative approach to serendipity measurement has a few drawbacks. First, it is sensitive to the choice of the baseline recommender system, and second, since it requires having two recommender systems – one optimized for serendipity and the other not –, it does not allow evaluating the serendipity of existing state-of-the-art recommendation approaches.

In this work, we take the first step toward a metric that would allow objectively measuring the surprise (or unexpectedness) of recommendations, without comparing them against an alternative system. We adopt the idea of measuring an item's distance from the user's profile.

As mentioned before, an item's distance from the set of items seen by the user has been exploited when optimizing for serendipity [1]. Also, Vargas and Castells, in their framework for measuring diversity and novelty, defined a metric that is based on computing an item's distance from a user's profile (although the authors do not use the term 'serendipity', but present it as a personalized novelty metric) [17].

Such a formulation of serendipity requires a distance measure between an item and a set of items. While the previous works used the *average* distance computed on item ratings or content features [17, 1], we hypothesize that averaging an item's distance from the user's profile items is not the best way to measure the item's surprise. Intuitively, a user perceives a recommendation as surprising if it is sufficiently different from *any* item the user has seen before. This information may be lost when averaging the computed distances, especially if the user has been exposed to diverse items. Therefore, we propose using the *lower bound* item distance from the user's profile items as an indicator of surprise. We believe this to be a reasonable estimate of how users perceive unexpected recommendations with respect to their past interactions with the system.

Note that we do not consider item relevance in our metric definitions. While relevance is also important to serendipity, we leave it to be measured separately by accuracy metrics. Furthermore, similarly to Adamopoulos and Tuzhilin [1], we also do not require a surprising item to be novel, although in principle, if an item is dissimilar to a user's profile, it is likely to be unknown to the user.

3. SURPRISE METRICS

We propose two alternative definitions of surprise — one based on users' rating behaviour and the other based on item content descriptors. Given the target item and the user's profile (a set of items rated by the user), both metrics produce a score that indicates the level of surprise the target item brings to the user.

3.1 Co-occurrence-based surprise

The first definition is based on the probability for the item to be seen (i.e., rated) together with the items in the user's profile. To measure the pairwise co-occurrence of items we employed normalized point-wise mutual information (PMI) [3], which measures the probability of observing specific outcomes of two independent random variables together. Given a pair of items i and j, we compute their PMI value as:

$$PMI(i,j) = \log_2 \frac{p(i,j)}{p(i)p(j)} / -\log_2 p(i,j)$$
(1)

where p(i) and p(j) represent the probabilities for the items to be rated by any user, and p(i, j) is the probability for the same user to rate both items. PMI values range from -1(in the limit) to 1, with -1 meaning the two items are never rated together, 0 signifying independence of the items, and 1 meaning complete co-occurrence of the items.

In order to measure the surprise of a recommended item i, we compute its PMI with each item in the user's profile, and take the maximum value as the overall surprise value:

$$S_{\text{co-occ}}(i) = \max_{i \in P} \text{PMI}(i, j)$$
(2)

where P is the user's profile (i.e., her set of rated items).

Since higher values of PMI(i, j) signify higher co-occurrence of items i and j (and therefore low surprise of seeing the two items together), taking the maximum value represents the lower bound of the surprise perceived by the user when item i is recommended.

Furthermore, to enable a comparison with the ideas in [17, 1], we conducted experiments with a variant where, instead of relying on the lower bound of surprise, we take the *average* PMI value across items in the user's profile:

$$S_{\text{co-occ}}^{avg}(i) = \frac{\sum_{j \in P} \text{PMI}(i, j)}{|P|}$$
(3)

For both $S_{\text{co-occ}}$ and $S_{\text{co-occ}}^{avg}$, lower values signify better surprise results.

Limitations.

We note that the co-occurrence-based definition of surprise metric may be sensitive to rare items, since the pointwise mutual information measure is known to be biased toward rare item pairs [3].

3.2 Content-based surprise

Our second surprise metric is based on distance applied to item content labels. We employed the complement of Jaccard similarity metric for comparing the items:

$$\operatorname{dist}(i,j) = 1 - \frac{L_i \cap L_j}{L_i \cup L_j} \tag{4}$$

where L_i and L_j are sets of labels describing items i and j.

To measure the surprise of a recommended item i, we take the minimum distance between the target item i and the user's profile items:

$$S_{\text{cont}}(i) = \min_{j \in P} \operatorname{dist}(i, j)$$
(5)

Taking the minimum distance value as the overall surprise represents the lower bound of how surprising the item is with respect to the seen items.

Similarly to the co-occurrence-based surprise, we also evaluated in our experiment a variant that takes the average distance:

$$S_{\text{cont}}^{avg}(i) = \frac{\sum_{j \in P} \operatorname{dist}(i, j)}{|P|}$$
(6)

For both S_{cont} and S_{cont}^{avg} , higher values signify better surprise results.

Limitations.

Note that the content-based metric is sensitive to the quality of item content labels. Low-quality labels may fail to capture the difference between items. Therefore, particular care should be taken when applying the metric to datasets with noisy data, e.g., user-generated labels.

4. EXPERIMENTS

4.1 Datasets

We tested the proposed surprise metrics on two benchmark datasets for offline recommender system evaluation the MovieLens 1M dataset and LastFM 1K dataset.

The MovieLens dataset contains ~ 1 million ratings, 6040 users and 3706 movies. The movies are annotated using a vocabulary of 18 genres (on average 1.65 genres per movie). To obtain richer content descriptors for the movies, we additionally scraped IMDb plot keywords for each movie, which resulted in an average of 81 labels per movie.

The LastFM dataset contains the listening events for 992 users and more than 100K artists. As the dataset is extremely sparse, we cleaned the set of artists by leaving only those for which we could obtain at least three LastFM tags and discarding artists that were listened to by fewer than 20 users. This resulted in 992 users and 7280 artists, with a total of \sim 500K ratings. The listening frequencies of the artists were transformed into ratings from 1 to 5 using the standard approach of converting implicit feedback into numerical ratings [17]. Since LastFM content labels for artists are user-generated (and not editorial as in the case of IMDb), to avoid noisy data we retrieved a maximum of the 10 most popular labels for every artist. This resulted in 9.2 labels per artist on average.

We used three state-of-the-art recommendation algorithms in our experiments: PureSVD — a matrix factorization approach (with 50 factors) that has been shown to perform well for top-N recommendation tasks [5] — and two standard neighbourhood methods — user-based and item-based collaborative filtering (with a neighbourhood size of 50) [6].

4.2 Evaluation methodology

In recent years, rating-based accuracy metrics for offline RS evaluations have been replaced by precision-oriented metrics that more closely reflect the users' interaction with the system — considering only a small set of top-ranked recommendations, ignoring the lower-ranked items.

In accordance with these state-of-the-art evaluation strategies, in this work we adopt the 'one plus random' methodology [5]. The methodology is based on splitting the dataset into a training set M and probe set P, constructing the test set T by selecting one highly-rated item per user from the probe set. Then, for each user u we make predictions for 1000 unrated items plus the one test item. The set of 1001 items is ranked according to the recommender's predicted score and the top-N recommendations are selected (we used N=10 in our experiments). If the test item is among the top-N items, we have a hit. The overall performance of the system — recall — is calculated as the ratio of the number of hits over the total number of test users. The underlying assumption that the 1000 unseen items are irrelevant is clearly undervaluing the performance, as certain items among the 1000 may be actually relevant for the user.

We believe this methodology to be appropriate when measuring the surprise of recommendations as it involves items the user has not discovered (i.e., unrated items), whereas other offline evaluation strategies employ items the user has had no trouble discovering (i.e., already rated items).

As described in Section 3, given an item and the user's profile, the co-occurrence-based and content-based surprise metrics produce a score in [-1; 1] and [0; 1] respectively. To obtain a single surprise value for a user's top-N recommendation list, we average the surprise values of all recommended items.

We computed the results using 5-fold cross validation with 80%-20% training/probe set split.

5. RESULTS AND DISCUSSION

As a first step in the evaluation, to make sure that the two surprise metrics are not redundant, we compared them using the following procedure: for each user in the dataset, we generated a list of top-10 recommendations using the PureSVD matrix factorization approach. We then ranked the 10 items using each metric and computed Spearman's rank correlation for the two rankings. Averaging the results over all users we obtained values close to zero for both MovieLens and LastFM datasets. This indicates that the two metric definitions capture different aspects of surprise.

In accordance with this finding, the two surprise metrics give different results in our main experiment.

5.1 Surprise value comparison

Table 1 lists the average recall and surprise values across test users for the evaluated algorithms — matrix factorization approach (MF), user-based approach (UB), and itembased approach (IB) — on both datasets. All reported metrics produce significantly different means when applied to different algorithms (p < 0.001 in a one-way ANOVA test). To measure which algorithms perform better/worse than the others, we applied the two-tailed t-test to each pair of algorithms.

Table 1: Recall and surprise results for the three tested algorithms. For each metric, values marked with $^+/^-$ are significantly better/worse than the other two values with p < 0.001.

| | Mc | vieLens | 1M | LastFM 1K | | | |
|---------------------------|-------------|-------------|-----------|-------------|-------------|-------------|--|
| Metric | MF | UB | IB | MF | UB | IB | |
| Recall | 0.334^{+} | 0.022^{-} | 0.065 | 0.427^{+} | 0.051 | 0.059 | |
| $S_{\text{co-occ}}$ | 0.349 | 0.349 | 0.354^- | 0.379^{+} | 0.404 | 0.407^{-} | |
| $S_{\rm cont}$ | 0.915^{-} | 0.928^{+} | 0.927 | 0.480^{-} | 0.643^{+} | 0.613 | |
| $S^{avg}_{\text{co-occ}}$ | 0.177^{-} | 0.146 | 0.097^+ | 0.141^{-} | 0.012^{+} | 0.049 | |
| $S_{\rm cont}^{avg}$ | 0.976^{-} | 0.981^+ | 0.980 | 0.887^{-} | 0.938^+ | 0.928 | |

As expected, MF produces the most accurate recommendations, outperforming both neighbourhood approaches.

For the content-based surprise metric (where higher values are better), the results are consistent for both minimum distance (Eq. 5) and average distance (Eq. 6) metric variants — MF performs significantly worse than both UB and IB, while UB performs better than IB.

On the other hand, for the co-occurrence-based surprise metric (where lower values are better), results differ for the minimum distance (Eq. 2) and average distance (Eq. 3) variants. In case of the minimum distance $S_{\text{co-occ}}$ variant, the MF approach performs as well as the UB and IB approaches, even outperforming both of them on the LastFM dataset. Contrastingly, the average distance variant $S_{\text{co-occ}}^{avg}$ shows MF to be outperformed by both UB and IB techniques.

An explanation for results of the co-occurrence-based surprise metric may lie in its sensitivity to rare items. In our experiments, the matrix factorization approach tends to rank popular items higher than both neighbourhood approaches, which can influence the values of point-wise mutual information. When one of the items in PMI(i, j) is rare, the normalized PMI metric can take extreme values — close to 1 if the two items frequently co-occur in the dataset, or close to -1 if they almost never occur together. Consequently, the minimum distance metric variant (which takes the maximum PMI value as the lower bound of an item's distance to a user's profile) tends to produce higher values when applied to rare items, while the average distance metric variant tends to produce low overall scores as it gets reduced by the negative PMI values. More experiments are needed to confirm this intuition. If confirmed, the $S_{\text{co-occ}}$ metric could be modified to avoid a bias toward rare items [3].

5.2 The impact of user's profile size

The size of a user's profile is an important factor to consider when measuring surprise of recommendations. Intuitively, the more items a user has been exposed to, the more difficult it is to surprise the user. To test this intuition, we plotted the surprise values obtained for individual users against their profile size. For both surprise metrics, we plotted two variants of results — one obtained using the minimum distance from the user's profile, and the other using the average distance from the profile. In both cases we measured the surprise of the recommendations generated using the most accurate recommendation method — the *PureSVD* approach.

Results for the minimum distance metric variant for both datasets are shown in the left half of Figure 1. As the user profile size increases, surprise values for the $S_{\text{co-occ}}$ metric tend to increase and the values for the S_{cont} metric decrease, i.e., in both cases the surprise of recommendations is decreasing. This finding confirms the intuition that given a large user's profile, there is a higher chance of finding an item similar to the recommended item, therefore, the chance of surprising the user is smaller.

The same trend is not observed when computing surprise using the average item's distance from a user's profile (right half of Figure 1), as the surprise values are independent of the profile size. We therefore believe that the minimum distance is more appropriate when measuring an item's surprise with respect to a user's profile, since averaging item distances results in losing information about the user's past interactions with the system.

5.3 Qualitative analysis

For a qualitative insight into the results, we looked at the recommendations generated for individual users. Table 2 shows for one user from the LastFM dataset 10 items from her profile and the top-10 recommendations generated using each algorithm. A subset of content labels is displayed for



Figure 1: Profile size influence on surprise values.

each item. (The full set of labels is omitted due to space limitations.) As can be seen from the content labels, all three algorithms generate items that are related to the user's musical preferences — independent rock and electronic music. However, user-based recommendations include more artists that may be surprising discoveries for the user, e.g., a jazz musician and a string quartet playing covers of popular songs in classical style. The item-based approach generates one such recommendation — a heavy metal band playing rockabilly songs.

6. CONCLUSIONS AND FUTURE WORK

In this paper, we investigated two metrics – co-occurrencebased (S_{co-occ}) and content-based (S_{cont}) – for measuring surprise in recommender systems. We evaluated each metric with three state-of-the-art recommenders on two benchmark datasets.

For both co-occurrence- and content-based metrics, two variants were evaluated — one using the minimum item's distance from the user's profile items as an indicator of surprise, and the other using the average distance from the profile items.

Using the $S_{\rm cont}$ surprise metric, the experimental results confirm that there is a trade-off between recommendation accuracy and serendipity [8] and the intuition that the matrix factorization approach, while being the most accurate, produces the least surprising recommendations with respect to the user's profile. Conversely, the user-to-user collaborative filtering approach produces the most surprising recommendations.

The above finding is confirmed for both $S_{\rm cont}$ metric variants — the minimum distance and average distance. How-

ever, the analysis of user's profile size influence on surprise values indicates that some information may be lost when using the average distance metric variant. Therefore, we believe that using the minimum distance to profile items is the better approach for measuring surprise.

The findings of $S_{\rm cont}$ are not fully confirmed using the alternative, co-occurrence-based surprise metric $S_{\rm co-occ}$. We believe the reason for this lies in the metric's sensitivity to rare items. However, the metric deserves further investigation, since in contrast to the content-based metric, it does not rely on item content labels and may capture a different aspect of surprise than the content-based metric.

The next steps in our research include measuring how the state-of-the-art recommendation algorithms perform with respect to other beyond-accuracy objectives, such as diversity, novelty, and coverage, in addition to serendipity. Additionally, we want to investigate the correlations and tradeoffs between these objectives. Finally, we would like to conduct a live user study, since no evaluation of beyondaccuracy objectives can be complete without obtaining feedback from real users.

7. ACKNOWLEDGMENTS

This publication has emanated from research conducted with the financial support of Science Foundation Ireland (SFI) under Grant Number SFI/12/RC/2289.

8. **REFERENCES**

[1] P. Adamopoulos and A. Tuzhilin. On unexpectedness in recommender systems: Or how to expect the unexpected. In *Procs. of the Workshop on Novelty and*

Table 2: A sample user profile and the generated recommendations from the LastFM dataset. The light-shaded item is the recovered hidden test item. The dark-shaded items are musicians belonging to music genres not present in the user's profile — jazz, classical, and heavy metal music.

| Profile | PureSVD (MF) | S _{co-occ} S _{cont} | User-based (UB) | S _{co-occ} S _{cont} | Item-based (IB) | S _{co-occ} S _{cont} |
|--|---|--|---------------------------------------|---------------------------------------|--|---------------------------------------|
| Death Cab for Cutie | Bloc Party | $0.318 \mid 0.333$ | Liam Finn | $0.357 \mid 0.750$ | Gram Parsons | $0.497 \mid 0.571$ |
| indie rock, alternative, emo, pop | indie rock, alternative, britpop | | indie, singer-songwriter, folk, pop | | country, singer-songwriter, folk, rock | |
| John Mayer | The Smiths | $0.381 \mid 0.333$ | Quietdrive | $0.395 \mid 0.571$ | Ralph Myerz | $0.390 \mid 0.667$ |
| singer-songwriter, acoustic, rock, blues | indie, new wave, alternative, britpop | | pop, punk, rock, alternative, emo | | electronic, downtempo, chillout, funk | |
| Damien Rice | Metric | $0.379 \mid 0.461$ | Jan Garbarek | $0.446 \mid 0.889$ | To My Boy | $0.365 \mid 0.750$ |
| singer-songwriter, acoustic, indie, folk | indie rock, female, alternative | | jazz, saxophone, instrumental | | new rave, electronic, indie, power pop | |
| Travis | M83 | $0.377 \mid 0.667$ | Steve Aoki | $0.379 \mid 0.571$ | The Enemy | $0.466 \mid 0.333$ |
| britpop, indie rock, alternative | electronic, shoegaze, post-rock, indie | | electronic, techno, indie, dance | | indie rock, alternative, britpop, punk | |
| The New Pornographers | Eels | $0.326 \mid 0.333$ | Girl Talk | $0.347 \mid 0.667$ | Volbeat | $0.321 \mid 0.750$ |
| indie rock, alternative, power pop | indie, alternative, rock, pop, american | | electronic, dance, hip-hop, mashup | | heavy metal, rockabilly, rock and roll | |
| The Postal Service | Pink Floyd | $0.264 \mid 0.571$ | Bloc Party | $0.318 \mid 0.333$ | Seasick Steve | $0.458 \mid 0.571$ |
| indie, electronic, pop, alternative, emo | rock, progressive, psychedelic, classic | | indie rock, alternative, britpop | | blues, rock, folk, country | |
| Toby Keith | Ratatat | $0.337 \mid 0.462$ | Ralph Myerz | $0.390 \mid 0.667$ | Elf Power | $0.422 \mid 0.333$ |
| modern country, male, pop, american | electronic, instrumental, indie | | electronic, downtempo, chillout, funk | | elephant 6, indie rock, lo-fi | |
| The Hold Steady | Green Day | $0.246 \mid 0.571$ | Olivia Ruiz | $0.394 \mid 0.750$ | Bloc Party | $0.318 \mid 0.333$ |
| indie rock, alternative, post-punk | punk rock, alternative, pop punk | | french, chanson francaise, female | | indie rock, alternative, britpop | |
| Clap Your Hands Say Yeah | Girl Talk | $0.347 \mid 0.667$ | Love And Rocket | $s 0.355 \mid 0.182$ | Biffy Clyro | $0.341 \mid 0.462$ |
| indie rock, alternative, indie pop, lo-fi | electronic, dance, hip-hop, mashup | | post-punk, new wave, gothic, 80s | | alternative, indie rock, emo | |
| The View | MSTRKRFT | $0.329 \mid 0.571$ | The String Quart | et $0.258 \mid 0.889$ | Boys Night Out | $0.345 \mid 0.571$ |
| indie rock, britpop, garage rock, punk electronic, dance, house, indie | | instrumental, classical, covers, strings | | emo, screamo, rock, hardcore | | |
| Average values: | | $0.330 \mid 0.497$ | | $0.364 \mid 0.627$ | | $0.392 \mid 0.534$ |

Diversity in Recommender Systems (DiveRS 2011), pages 11–18, 2011.

- [2] P. André, M. Schraefel, J. Teevan, and S. T. Dumais. Discovery is never by chance: designing for (un)serendipity. In *Procs. of the 7th ACM Conference* on *Creativity and Cognition*, pages 305–314, 2009.
- [3] G. Bouma. Normalized (pointwise) mutual information in collocation extraction. In Procs. of the Conference of the German Society for Computational Linguistics and Language Technology, pages 31–40, 2009.
- [4] J. Campos and A. D. de Figueiredo. Searching the unsearchable: Inducing serendipitous insights. In Procs. of the Workshop Program at the 4th International Conference on Case-Based Reasoning, pages 159–164, 2001.
- [5] P. Cremonesi, Y. Koren, and R. Turrin. Performance of recommender algorithms on top-n recommendation tasks. In *Procs. of the 4th ACM Conference on Recommender Systems*, pages 39–46, 2010.
- [6] C. Desrosiers and G. Karypis. A comprehensive survey of neighborhood-based recommendation methods. In F. Ricci et al., editors, *Recommender Systems Handbook*, pages 107–144. Springer, 2011.
- [7] A. Foster and N. Ford. Serendipity and information seeking: an empirical study. *Journal of Documentation*, 59(3):321–340, 2003.
- [8] M. Ge, C. Delgado-Battenfeld, and D. Jannach. Beyond accuracy: evaluating recommender systems by coverage and serendipity. In Procs. of the 4th ACM Conference on Recommender Systems, pages 257–260, 2010.
- [9] J. L. Herlocker, J. A. Konstan, L. G. Terveen, and J. T. Riedl. Evaluating collaborative filtering recommender systems. ACM Transactions on Information Systems, 22(1):5–53, 2004.
- [10] L. Iaquinta, M. de Gemmis, P. Lops, G. Semeraro, M. Filannino, and P. Molino. Introducing serendipity in a content-based recommender system. In *Procs. of*

the 8th International Conference on Hybrid Intelligent Systems, pages 168–173, 2008.

- [11] A. McBirnie. Seeking serendipity: the paradox of control. Aslib Proceedings, 60(6):600–618, 2008.
- [12] S. M. McNee, J. Riedl, and J. A. Konstan. Being accurate is not enough: how accuracy metrics have hurt recommender systems. In *CHI'06 Extended Abstracts on Human Factors in Computing Systems*, pages 1097–1101, 2006.
- [13] T. Murakami, K. Mori, and R. Orihara. Metrics for evaluating the serendipity of recommendation lists. In Procs. of the 2007 Conference on New Frontiers in Artificial Intelligence, pages 40–46, 2008.
- [14] K. Oku and F. Hattori. Fusion-based recommender system for improving serendipity. In Procs. of the Workshop on Novelty and Diversity in Recommender Systems (DiveRS 2011), pages 19–25, 2011.
- [15] K. Onuma, H. Tong, and C. Faloutsos. Tangent: a novel, 'surprise me', recommendation algorithm. In Procs. of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining, pages 657–666, 2009.
- [16] E. G. Toms. Serendipitous information retrieval. In Procs. of 1st DELOS Workshop on Information Seeking, Searching and Querying in Digital Libraries, pages 11–14, 2000.
- [17] S. Vargas and P. Castells. Rank and relevance in novelty and diversity metrics for recommender systems. In Procs. of the 5th ACM Conference on Recommender Systems, pages 109–116, 2011.