

Online Recommender System for Radio Station Hosting: Experimental Results Revisited

Dmitry I. Ignatov

School of Applied Mathematics and
Information Science

National Research University Higher School of Economics
Moscow, Russia
dignatov@hse.ru

Taimuraz Abaev

National Research University Higher School of Economics
taymuraz.abaev@gmail.com

Sergey Nikolenko

Laboratory of Internet Studies

National Research University Higher School of Economics,
Steklov Institute of Mathematics at St. Petersburg of the RAS,
Saint Petersburg, Russia
sergey@logic.pdmi.ras.ru

Natalia Konstantinova

University of Wolverhampton, UK
n.konstantinova@wlv.ac.uk

Abstract—We present a new recommender system developed for the Russian interactive radio network *FMhost* based on a previously proposed model. The underlying model combines a collaborative user-based approach with information from tags of listened tracks in order to match user and radio station profiles. It follows an adaptive online learning strategy based on the user history. We compare the proposed algorithms and an industry standard technique based on singular value decomposition (SVD) in terms of precision, recall, and NDCG measures; experiments show that in our case the fusion-based approach shows the best results.

Keywords—*music recommender systems; interactive radio network; hybrid recommender system; information fusion; quality of service*

I. INTRODUCTION AND RELATED WORK

Music recommendation is an important direction in the field of recommender systems. Many recent works in this area have appeared at the International Society for Music Information Retrieval Conference (ISMIR) [1], Workshop on Music Recommendation and Discovery (WOMRAD) [2], [3], and the Recommender Systems conference (RecSys) [4]. Several broadcasting services, e.g., *last.fm*, *Yahoo!LaunchCast*, and *Pandora*, are known for their recommender systems and work on a commercial basis (however, the latter two do not operate in Russia). Despite many high quality works on different aspects of music recommendation, there are only a few studies devoted to online radio station recommender systems [5]. This work deals with the Russian online radio hosting service *FMhost*, in particular, its new hybrid recommender system.

Recently, the focus of computer science research dealing with music industry has shifted from music information retrieval and exploration [6], [7], [8] to music recommendation [9], [10]. It is not a new direction (see, e.g., [11]); however, it is now inspired by new capabilities of large online services that can provide not only millions of tracks for listening, but also thousands of radio stations to choose from at a single web site. In addition to this, social tagging is an important factor that lets us apply new recommender algorithms based on tag similarity [12], [13], [14].

A widely acclaimed public contest on music recommender algorithms, KDD Cup¹, was recently held by *Yahoo!*. In KDD Cup, track 1 was devoted to learning to predict users' ratings of musical items (tracks, albums, artists, and genres) where the items formed a taxonomy: tracks belong to albums, albums belong to artists, where albums and tracks are also tagged with genres. Track 2 aimed at developing learning algorithms for separating music tracks scored highly by specific users from tracks that have not been scored by them. It attracted a lot of attention to the problems that are both typical for recommender systems and specific for music recommendation: scalability issues, capturing the dynamics and taxonomical properties of the items [15]. Another large music recommender contest, the Million Songs Dataset Challenge², with open data was held in 2012 by the Computer Audition Lab at UC San Diego and LabROSA at Columbia University [16]. The core data consists of triples (*user, song, count*); it covers approximately 1.2 million users and more than 380,000 songs.

Current music recommendation trends reflect the advantages of hybrid approaches and call for user-centric quality measures [17]. For instance, the work [18] proposes a novel approach based on the so-called “forgetting curve” to evaluate “freshness” of predictions. The work [19] studies the problem of how much metadata one needs in music recommendation: a subjective evaluation of 19 users has revealed that pure content-based methods can be drastically improved with genre tags. Finally, the authors proposed a recommender approach that starts from an explicit set of music tracks provided by the user as evidence of his/her preferences and then infers high-level semantic descriptors for each track [20].

In [21], the authors proposed a music recommender system Starnet for social networking. It generates social recommendations based on positive ratings of friends, non-social recommendations based on positive ratings of other users in the network, and random recommendations. Hottabs is another interesting recently developed online music recommendation system is Hottabs [22], it is dedicated to the learning how

¹<http://kddcup.yahoo.com/>

²<http://labrosa.ee.columbia.edu/millionsong/>

to play guitar. Some authors aim at improving music recommender systems with semantic extraction techniques [23], [24]. In [25], the author describes a system of genre recommendation for music and TV programs, which can be considered as an alternative channel selector. The authors of [26] proposed a recommender system GroupFan which is able to aggregate preferences of a group of users to their mutual satisfaction.

Many online services (e.g., *last.fm* and *LaunchCast*) call their audio streams “radio stations”, however they are actually just playlists from a database of tracks which are based on a recommender system rather than real predefined channels. *FMhost*³, on the other hand, provides users with online radio stations in the classical meaning of this term: human DJs perform live, and a radio station actually represents a strategy or mood of a specific person (DJ) who plays his/her own tracks, performs contests etc. Thus, the problem we aim to solve differs from most of the work done in terms of music recommendation, and some of the challenges are unique. For instance, our system is scalable and it allows users to listen to music without being registered, and as a result our test dataset contains only 4266 registered users with a recorded listening history. The dataset contains 2206 radio stations with 24803 non-zero entries in the user–station matrix; it is small and sparse, so standard recommendation techniques (e.g., SVD that we compare our algorithms to) fail to achieve good prediction quality.

The paper is organized as follows. In Section II, we describe the online radio service *FMhost*. In section III, we propose a novel recommender model, two basic recommender algorithms, a third algorithm that combines them, and show the recommender system architecture. Quality of service (QoS) measurements, a comparison with an SVD-based approach, and some insights on *FMhost* user behaviour are discussed in Section IV. Section V concludes the paper.

II. ONLINE SERVICE FMHOST.ME

A. A concise online broadcasting dictionary

We begin by briefly introducing some basic domain terminology. A *chart* is a track rating of a particular radio station; for example, a rock chart shows a certain number (e.g., 10) of most popular rock tracks, ranked from the most popular (rank 1) to the least popular (rank 10) according to a survey. A *live performance* (or just *live* for short) is a performance with one or several *DJs* (*disk jockeys*) assigned to it. They perform using their own PCs, and the audio stream is being redirected from them to an Icecast server and then distributed to the users. The DJs may also have their own blog for each live, where people can interact with DJs who perform live. *LiquidSoap* is a sound generator that broadcasts audio files (*.mp3, *.aac etc.) into an audio stream, and *Icecast* is a retranslation server that redirects an audio stream from one source (e.g., *LiquidSoap*) to many receivers.

B. The *FMhost* project

FMhost is an interactive radio network. The portal allows users to listen and broadcast their own radio stations. There

are four user categories in the portal: (1) unauthorized user, (2) listener, (3) Disk Jockey (DJ), and (4) radio station owner.

User capabilities vary with their status. Unauthorized listeners can listen to any station but cannot vote or become DJs and cannot use the recommender system and the rating system. Listeners can vote for tracks, lives, and radio stations. They are allowed to use the recommender system and the rating system, and they can subscribe to lives, radio stations, or DJs. They also can be appointed to a live and become a DJ.

There are three types of broadcasting: (1) stream redirection from another server, (2) AutoDJ translation, and (3) live performance. Stream redirection applies when a radio station owner has a separate dedicated server, uses *FMhost* as a broadcasting platform, and broadcasts with his own sound generator, e.g., *SamBroadcaster* (<http://spacial.com/sam-broadcaster>), *LiquidSoap* (<http://savonet.sourceforge.net/>) etc. AutoDj is a special option that allows the users to play music directly from the *FMhost* server. Every radio owner gets some space where he/she can download tracks, and then *LiquidSoap* generates the audio stream and the Icecast (<http://www.icecast.org/>) server redirects it to the listeners. Usually the owner sets a schedule for his radio station. Live performances are done by DJs. Everyone who has performed live at least once can be called a DJ. He or she can also be added to a radio station crew. Moreover, a DJ can perform lives at any station, not only on his own station where he or she is in a crew.

FMhost was the first project of its kind in Russia, starting in 2009. Following *FMhost*'s success, there now exist several radio broadcasting portals: <http://frodio.com/>, <http://myradio24.com/>, <http://www.radio-hoster.ru/>, <http://www.taghosting.ru/>, <http://www.economhost.com/>, and even <http://fmhosting.ru/>. In late 2011, *FMhost* was taken down for a major redesign of both codebase and recommender system architecture. In this paper, we describe the results of this upgrade.

The previous version of the recommender system experienced several problems, including tag discrepancy and personal tracks without tags at all. A survey conducted by *FMhost* with about a hundred respondents showed that more than half of them appreciated the previous version of our recommender system, and more than 80% of the answers were either positive or neutral (see Table 1 in [27]); nevertheless, the new recommender model and algorithms provide even better recommendations and make even less prediction mistakes.

C. *FMhost* conceptual improvements

The new version features a more complex system of user interactions. Every radio station has an owner who is not just a name but also has the ability to assign DJs for lives, prepare radio schedule, and assign lives and programs. A new broadcasting panel allows the DJs to play tracks with new features such as an equalizer or fading between tracks. A new algorithm for the recommender system, a new rating system, and a new chart system will also be launched in the new version. The rating system has been developed to rank radio stations and DJs according to their popularity and quality of work. A new core is being implemented, and a new concept of *LiquidSoap* and *Icecast* is being designed. The new system fixes all problems that were identified in the previous version.

³<http://host.fm/en/>

III. MODELS, ALGORITHMS AND RECOMMENDER ARCHITECTURE

A. Input data and general structure

Our model is based on three data matrices. The first matrix $A = (a_{ut})$ tracks the number of times user u visits radio stations with a certain tag t . Each radio station r broadcasts audio tracks with a certain set of tags T_r . The sets of all users, radio stations, and tags are denoted by U , R , and T respectively. The second matrix $B = (b_{rt})$ contains how many tracks with a tag t a radio station r has played. Finally, the third matrix $C = (c_{ur})$ contains the number of times a user u has visited a radio station r . For each of these three matrices, we denote by v^A , v^B , and v^C the respective vectors containing sums of elements: $v^A = \sum_{t \in T} a_{ut}$, $v^B = \sum_{t \in T} b_{rt}$, and $v^C = \sum_{r \in R} c_{ur}$. We also introduce for each matrix A , B , C the corresponding frequency of visits matrices A_f , B_f , and C_f ; the frequency matrix is a result of normalization the matrix with the respective vector of visits, e.g., $A_f = (a_{ut} \cdot (v_u^A)^{-1})$. Our model is not static: the matrices A , B , and C change after a user u visits a radio station r with a tag t , i.e., each value a_{ut} , b_{rt} , and c_{ur} is incremented by 1 after this visit.

The model consists of three main blocks: Individual-Based Recommender System (IBRS) model, Collaborative-Based Recommender System (CBRS) model, and Fusion Recommender System (FRS) that aggregates the results of the first two. Each model has its own algorithmic implementation.

B. IBRS

The **IBRS** model uses matrices A_f and B_f and aims to provide a particular user $u_0 \in U$ with top N recommendations represented mathematically by a special structure $Top_N(u)$. Formally, $Top_N(u_0)$ is a triple $(R_{u_0}, \preceq_{u_0}, \text{score})$, where R_{u_0} is a set of at most N radio stations recommended to a particular user u_0 , \preceq_{u_0} is a well-defined quasiordering (reflexive, transitive, and complete) on the set R_{u_0} , and score is a function which maps each radio station r from R_{u_0} to $[0, 1]$.

The algorithm computes the l_1 -norm (Manhattan) distance between a user u_0 and a radio station r : $d(u_0, r) = \sum_t t \in T |a_{u_0 t} - b_{rt}|$. Then distances between the user u_0 and all radio stations $r \in R$ are computed, and the algorithm constructs the relation \preceq_{u_0} according to the following rule: $r_i \preceq r_j$ iff $d(u_0, r_i) \leq d(u_0, r_j)$. The function score operates on R_{u_0} according to the following rule:

$$\text{score}(r_i) = 1 - d(u_0, r_i) / \max_{r_j \in R} d(u_0, r_j).$$

Finally, after selecting N radio stations with N largest values in R_{u_0} , we have a ranked list $Top_N(u_0)$ of radio stations recommended to the user u_0 . In case there are several elements with the same score (rank) so that $Top_N(u)$ is not uniquely defined, we simply choose the first elements according to some arbitrary ordering (e.g., lexicographically by their names).

As shown in Fig. 1, our simplified model takes into account only “listened tracks” but the previously proposed one [27] also deals with “liked tracks”, “liked radio stations”, and “favorite radio stations”.

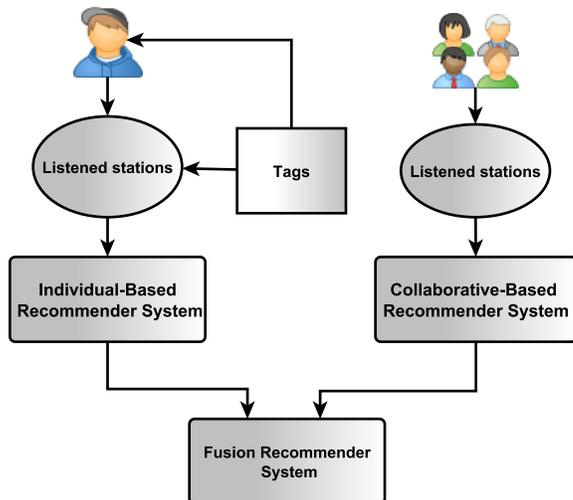


Fig. 1. The recommender system architecture

C. CBRS

The **CBRS** model is based on the C_f matrix (normalised number of times a user u has visited a radio station r). This matrix also yields a vector n^C which stores the total number of stations listened by each user $u \in U$. This vector also changes over the time, and this value is used as a threshold to transform matrix C_f to similarity matrix D via cosine similarity between users i and j :

$$\text{sim}(i, j) = \frac{c_{fi} \cdot c_{fj}}{\sqrt{\sum_{r \in R} c_{fir}^2} * \sqrt{\sum_{r \in R} c_{fjr}^2}}$$

After computing D , the algorithm constructs the list $Top_k(u_0) = (U_{u_0}, \preceq_{u_0}, \text{sim})$ of k users similar to the target user u_0 who awaits recommendations. We define the set of all radio stations user u_0 has listened to by $L(u_0) = \{r | c_{fur} = 0\}$. In a similar way, we define

$$\begin{aligned} Top_N(u_0) &= (R_{u_0}, \preceq_{u_0}, \text{score}), \\ \text{where } \text{score}(r) &= \text{sim}(u^*) \cdot c_{fu^*r} \\ \text{and } u^* &= \arg \max_{u \in U_{u_0}, r \in U/L(u_0)} \text{sim}(u) \cdot c_{fur}. \end{aligned}$$

Note that score takes values in the range of $[0, 1]$. The problem of choosing exactly N best stations is solved in the same way as in the IBRS submodel.

D. FRS

After IBRS and CBRS have finished their work and provides the results, we have two ranked lists of recommended stations $Top_N^I(u_0)$ and $Top_N^C(u_0)$ for a target user u_0 from IBRS and CBRS respectively. The **FRS** submodel proposes a simple solution for aggregating these lists into the final recommendation structure $Top_N^E(u_0) = (R_{u_0}^E, \preceq_{u_0}^E, \text{score}^E)$. For every $r \in R_{u_0}^C \cup R_{u_0}^I$, the function $\text{score}^E(r)$ maps r to the weighted sum

$$\beta \cdot \text{score}^C(r) + (1 - \beta) \cdot \text{score}^I(r),$$

where $\beta \in [0,1]$, $\text{score}^C(r) = 0$ for all $r \notin R^C$, and $\text{score}^I(r) = 0$ for all $r \notin R^I$. The algorithm adds the best N radio stations according to this criterion to the set $R_{u_0}^C$.

IV. QUALITY OF SERVICE ASSESSMENT

To evaluate the quality of the developed system, we propose a variant of cross-validation technique [28]. We represent the dataset with an object-attribute table (binary relation) $T \subseteq U \times I$, where uTi iff user $u \in U$ used (purchased, watched, listened etc.) item $i \in I$. To evaluate the quality of recommendations in terms of precision and recall, we split the initial user set U into training and test subsets U_{train} and U_{test} , with test set smaller than the training set, e.g., with a 20/80 proportion. Recommendation precision and recall are evaluated on the test set, and this part of the algorithm is similar one step of conventional cross-validation. Then each user vector u from U_{test} is divided into two parts which consist of evaluated items I_{visible} and items I_{hidden} which we have intentionally hidden. Note that in the existing literature, the proportion between size of I_{visible} and I_{hidden} is not discussed even in similar schemes [29]. Then, for example, a user-based algorithm makes recommendations according to similarity between users from the test and training sets. Each user from U_{test} gets recommendations as a set of fixed size $r_n(u) = \{i_1, i_2, \dots, i_n\}$. Precision and recall are defined as

$$\text{Recall} = \frac{|r_n(u) \cap u^I \cap I_{\text{hidden}}|}{|u^I \cap I_{\text{hidden}}|},$$

$$\text{Precision} = \frac{|r_n(u) \cap u^I \cap I_{\text{hidden}}|}{|r_n(u) \cap I_{\text{hidden}}|},$$

where u^I is a set of all items from I used by u . These measures are calculated for each user and then averaged. The experiment can be repeated several times, e.g. 100, for different test and training set splits, and then the values are averaged again. In addition, one can select the set I_{hidden} at random, but with a specific proportion, e.g. 20%. The idea behind the method comes from traditional cross-validation, but in case of recommender systems some modifications are needed. We used a modified m -fold cross-validation, which is performed by splitting the initial set into m disjoint subsets, where each subset is used as a test set and the other subsets are considered as training ones. To evaluate ranking quality based on the number of performed listenings we use Normalized Discounted Cumulative Gain (NDCG).

Before we proceed to the detailed description of the procedure, we discuss some important aspects of the *FMHost* data that we have mined.

A. Basic statistics

The dataset contains following entities: users, tags, radio stations, and tracks. We have selected users with at least one tag in the profile, and then all the tags related to the selected users. At the next step, we have chosen radio stations that the selected users have listened to or with selected tags in the radio station profiles. Finally, we have chosen tracks related to at least one of the selected radio stations and to at least one tag. The resulting dataset contains 4266 users, 3618 tags, 2209 radio stations, and 4165 tracks. The corresponding matrices have the following number of nonempty entries: 38504 in the

user–tag matrix, 18539 in the radio station–tag, 24803 in the user–radio station, 18781 in the track–tag, and 22525 in the radio station–track matrix.

It is a well-known fact that social networking data often follows the so called power law distribution [30]. In order to choose which number of active users or radio stations we have to take into account for making recommendations, we performed a simple statistical analysis of the user and radio station activity. Around 20% of the users (only registered ones) were analysed.

TABLE I. BASIC PARAMETERS OF THE USER AND RADIO VISITS DATASETS, WITH POWER-LAW FITS AND THE CORRESPONDING p -value .

Dataset	n	$\langle x \rangle$	σ	x_{max}	\hat{x}_{min}	$\hat{\alpha}$	n_{tail}	p -value
User dataset	4187	5.86	12.9	191	12 ± 2	2.46(0.096)	117	0.099
Radio dataset	2209	11.22	60.05	1817	46 ± 11	2.37(0.22)	849	0.629

Table I shows p -values of statistical tests performed with the *Matlab* package introduced in [30]. It shows that the power law does fit the radio station dataset, and the probability to make an error by ruling out the null hypothesis (no power law) is about 0.1 for the user dataset. Thus, radio station visits dataset is more likely to follow the power law than the user visits dataset, but we should take it into account for both datasets; Fig. 2 shows how the power law actually fits our data.

This analysis implies useful consequences according to the well-known “80:20” rule $W = P^{(\alpha-2)/(\alpha-1)}$, which means that the fraction W of the wealth is in the hands of the richest P of the population. In our case, 50% of users make 80% of all radio station visits, and 50% of radio stations have 83% of all visits (which is actually a rather flat distribution compared to most services). Thus, if the service tends to take into account only active stations and users, it can cover 80% of all visits by considering 50% of their active audience. However, new radio stations still deserve to be recommended, so this rule can only be applied to the user database.

B. Quality assessment

For quality assessment in the IBRS subsystem, we count average precision and recall on the set $R_N \subset R$, where N is the number of randomly “hidden” radio stations. We suppose that for every station $r \in R_N$ and every user $u \in U$ the algorithm does not know whether the radio stations were visited, and we change A_f and R accordingly. Then IBRS attempts to recommend Top- N radio stations for the modified matrix A_f . Top- N average precision and recall are computed as follows:

$$\text{Precision} = \frac{\sum_{u \in U} \frac{|R_u^I \cap L_u \cap R_N|}{|L_u \cap R_u^I|}}{|U|},$$

$$\text{Recall} = \frac{\sum_{u \in U} \frac{|R_u^I \cap L_u \cap R_N|}{|L_u \cap R_N|}}{|U|}.$$

To deal with CBRS, we use a modification of the leave-one-out technique. At each step of the procedure for a particular user u , we “hide” all radio stations $r \in R_N$ by setting $c_{fur} =$

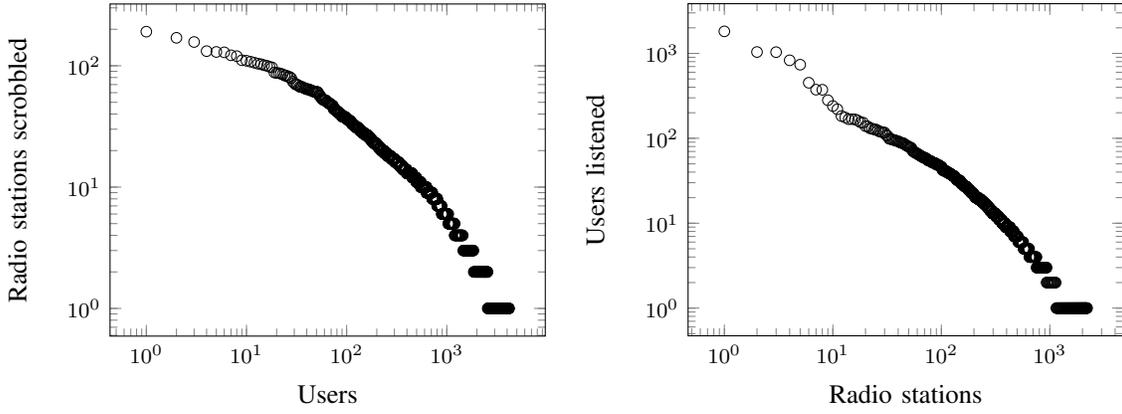


Fig. 2. Power law at FMHost. On the left, users sorted by the number of radio stations they have listened to. On the right, radio stations sorted by the number of their users. Both graphs (on log scale) show power law dependencies.

0. Then we perform CBRS algorithm assuming that $c_{fu'r}$ is unchanged for $u' \in U/u$ and then compute

$$\text{Precision} = \frac{\sum_{u \in U} \frac{|R_u^C \cap L_u \cap R_N|}{|L_u \cap R_u^C|}}{|U|},$$

$$\text{Recall} = \frac{\sum_{u \in U} \frac{|R_u^C \cap L_u \cap R_N|}{|L_u \cap R_N|}}{|U|}.$$

To tune the FRS system, we can use a combination of these two procedures trying to find the optimal β as

$$\beta^* = \arg \max_{\beta} \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{(\text{Precision} + \text{Recall})}.$$

After the launch of the system, we suppose to be able to collect enough reliable statistics in about one month of active operation in order to tune β and choose appropriate similarity and distance measures and thresholds. We suppose that the resulting system will provide reasonably accurate recommendations using only a single (last) month of user history and only 50% of the most active users. For quality assessment during the actual operation, we will compute Top-N precision and recall measures as well as NDCG. In addition, online surveys can be launched to assess user satisfaction with the new RS system.

C. Experimental results

The IBRS algorithm uses the user–tag matrix A and the radiostation–tag matrix B . As a result, we have a matrix of predicted scores that contains recommendations of radiostations for users. Figure 6 shows the precision, recall, and NDCG measure of IBRS versus the size of recommendation list. For the first recommended item, precision is about 30% and then it rapidly drops as Top- N grows to 5-10 elements. Then it becomes close to 1% while N goes to 100. Recall for the first recommended radiostation is 50%, and then recall value slowly increases as Top- N grows. NDCG slowly increases while Top- N grows. Our *Matlab* implementation of IBRS runs for about 80s for all users.

The CBRS algorithm uses the user–radiostation matrix C . To evaluate its quality, we used 3×3 -fold cross-validation

which we described in the beginning of Section IV. Therefore, as a result we have 9 different partitions of the initial data into training and test set, and all measures are averaged by all partitions. Figure 6 shows that CBRS precision is smaller than IBRS precision for small Top- N size: for the first recommended radiostation it is about 7%. However, for large Top- N size it goes down to 2% ($N = 100$), which is better than IBRS precision. In terms of recall CBRS also loses: recall of the first recommended radiostation is about 10%, and then for higher values of N it becomes closer to IBRS recall. NDCG of CBRS is strictly less than the NDCG of IBRS, which shows that IBRS is a better ranker than CBRS. The testing time for the entire cross-validation procedure is about 50 minutes.

The hybrid recommender system FRS uses output matrices with scores of recommended radiostations for IBRS and CBRS. To make the final ranking, it employs a weighted sum of the IBRS and CBRS output matrices. Parameter β is tuned for each Top- N size. To this end, for β from 0 to 1 with step 0.05 we have calculated the resulting weighted matrix; for this matrix, we calculate precision, recall, the F-measure, and NDCG. The procedure is repeated for N ranging from 1 to 100. The value of β is then tuned to maximize one of the measures. By maximizing the F-measure, we improve the general quality of the output but nearly ignore the ranking. Maximizing NDCG takes into account ranking. The time needed to tune β in the range from 0 to 1 with step 0.05 in our implementation is 200s. It takes up to 3s to calculate the final recommendation matrix with the chosen β .

When we maximize the F-measure with respect to β , for $N = 1$ IBRS has a slightly higher weight, about 0.57, but as Top- N grows to 30-35 radio stations, the parameter β smoothly decreases, thus increasing CBRS contribution; for $N > 35$ β drops to 0 (see Fig. 3). These results can be explained by the fact that IBRS provides higher precision and recall than CBRS, but as N grows the F-measure of CBRS is getting higher than for IBRS. As a result, the F-measure for FRS is not less than F-measure of either IBRS or CBRS individually, and FRS performs efficiently for every tested recommender list size. In particular, for N on the order of 15-20 radiostations the F-measure of FRS is higher than for the best basic method (IBRS) by 2-3%. For larger size of the top lists N , the FRS F-measure is close to CBRS.

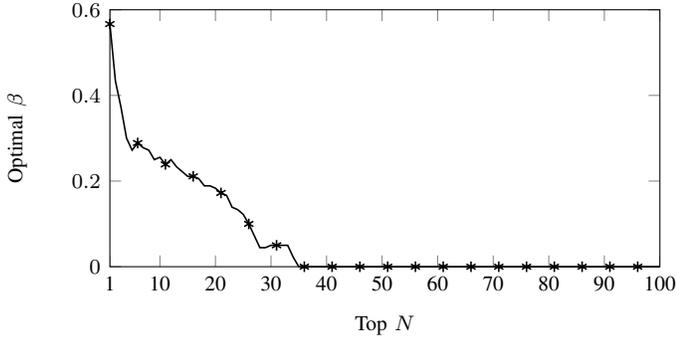


Fig. 3. Tuning the parameter β to maximize F -measure

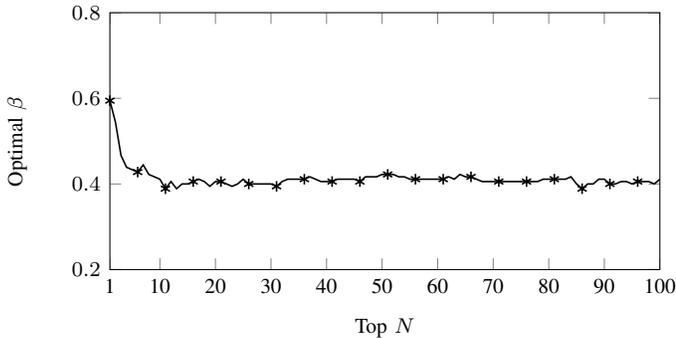


Fig. 4. Tuning the parameter β to maximize NDCG

NDCG maximisation by the parameter β prefers IBRS at first: the parameter begins with 0.6. But then in the range of N from 2 to 10 radiostations β decreases to 0.4 and stabilizes, oscillating for larger values of N around 0.40-0.41. This implies that CBRS contributes slightly more to the final recommendation (see Fig. 4). When we tune β to maximize NDCG, the weighted sum approach performed better than IBRS (which is better than CBRS with respect to NDCG) by 4-5% for all tested output sizes.

For comparison, we have also implemented a standard collaborative filtering model based on singular value decomposition (SVD) [31], [32], [33], [34], [35]. In this model, an item's rating is approximated as a scalar product of user and item feature vectors (plus some baseline predictors). In our setting, lacking explicit ratings, we have applied SVD to the matrix of user listenings to radio stations. Due to the power law dependencies we have discovered (see Section IV-A), we have used the logarithm of the number of listenings (plus one) in the model:

$$\log(\text{listen}_{u,r} + 1) \sim \mu + b_u + b_r + v_u^\top v_r,$$

where $\text{listen}_{u,r}$ is the number of times user u listened to radio station r , μ is the general mean, b_u and b_r are the baseline predictors for the user u and the station r respectively, and v_u and v_r are the vectors of the user and station features respectively.

It was clear that there is not enough data and the matrix is too sparse for SVD. This was also supported by our experiments: we did not see any improvement at all as the

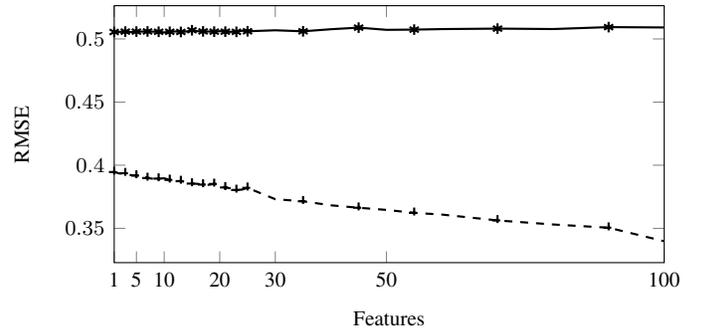


Fig. 5. Root mean squared error in SVD experiments as a function of the number of features. Solid line denotes error on the validation set; dashed line, error on the training set.

number of SVD features grew, the quality on the validation set was stable all the way from a single feature to about 100 and then started showing clear signs of overfitting; this is depicted in Fig. 5. Results of the main experiments also supported this observation: in Fig. 6, SVD clearly loses to the methods proposed in this work.

V. CONCLUSION AND FURTHER WORK

In this work, we have described the underlying models, algorithms, and system architecture of the new improved *FMHost* service and tested it on the available real dataset. We hope that the developed algorithms will help a user to find relevant radio stations for listening. In the future optimization and tuning, special attention should be paid to scalability issues and user-centric quality assessment.

By using bimodal cross-validation, we have built a hybrid algorithm FRS tuned to maximize either F -measure or NDCG for various values of N and β . The FRS algorithm performs better than the three other approaches, namely IBRS, CBRS, and SVD, both in terms of F -measure and in terms of NDCG.

According to the NDCG@n measure, IBRS is strictly better than CBRS, so the former one is a better ranker. The maximal value of F -measure is obtained for the Top- N of size 15.

Surprisingly, in our experiments the state-of-the-art SVD-based technique performed poorly in comparison to our proposed algorithms. This can be explained by the small size and sparseness of our dataset. We hypothesize that the methods described in this work will suit datasets with similar properties. Matrix factorization techniques remain a reasonable tool to increase scalability, but they have to be carefully adapted and assessed, taking into account the folksonomic nature of the track tags and rather disappointing results of the SVD-based technique. Another important issue is connected to the triadic relational nature of the data (users, radio stations or tracks, and tags), which constitutes the so called *folksonomy* [36], a fundamental data structure in resource-sharing systems with tags. As shown in [37], this data can be successfully mined by means of triclustering, so we also plan to build a tag-based recommender system by means of triclustering.

Acknowledgments: We would like to thank Rustam Tagiev and Mykola Pechenizkiy for their comments and Vasily Zakharchuk and Andrey Konstantiov for their very important work

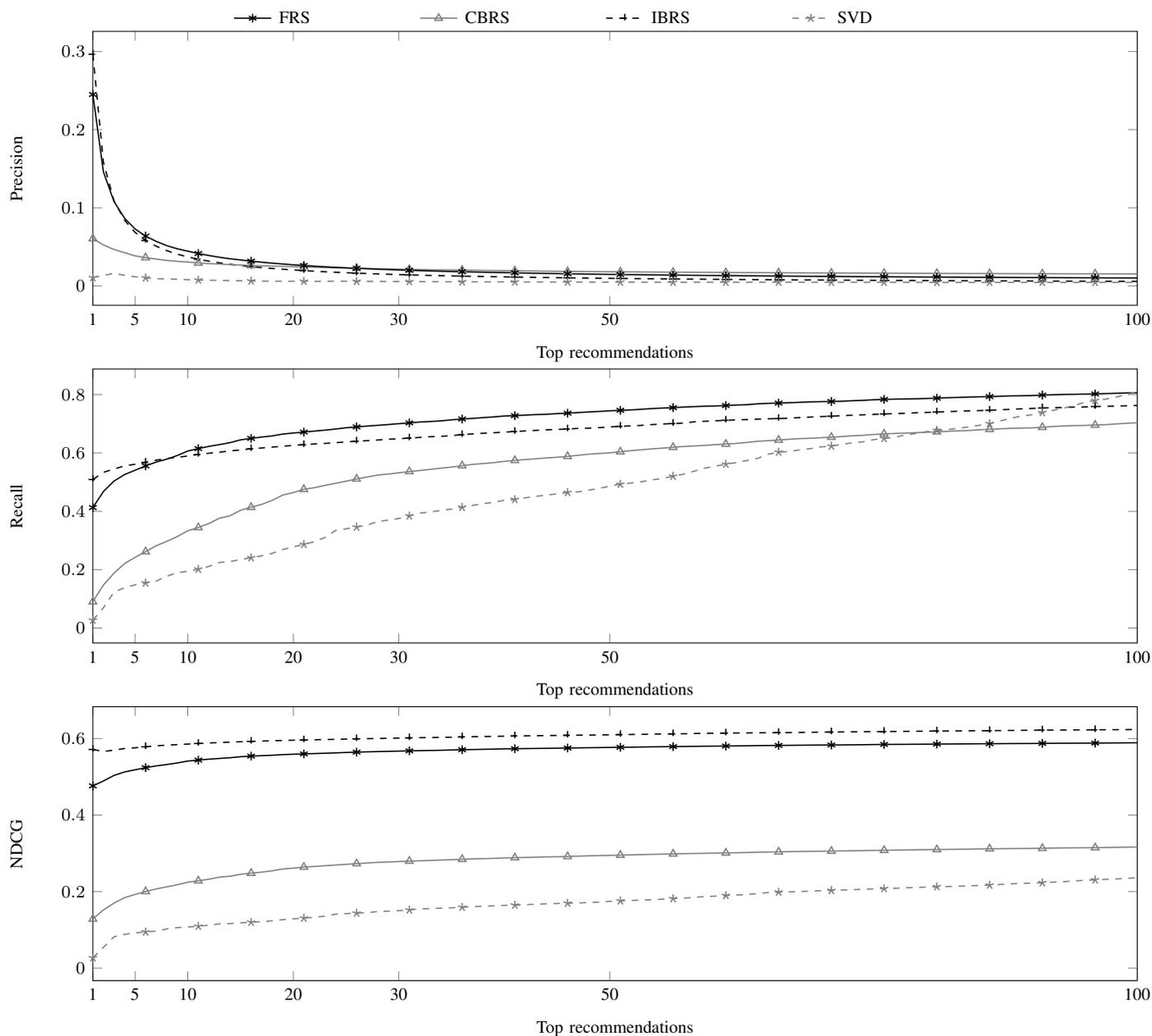


Fig. 6. Experimental results for four algorithms: IBRS, CBRS, FRS, and SVD. Graphs, top to bottom: precision, recall, and NDCG as a function of the number of top recommendations predicted for a user (averaged over all users).

at the previous stages of this project. Our research was done in the framework of the Basic Research Program at the National Research University Higher School of Economics in 2014, at the Laboratory of Intelligent Systems and Structural Analysis (Moscow) and Laboratory of Internet Studies (St. Petersburg). Dmitry Ignatov was supported by Russian Foundation for Basic Research (grant # 13-07-00504).

REFERENCES

- [1] A. de Souza Britto Jr., F. Gouyon, and S. Dixon, Eds., *Proceedings of the 14th International Society for Music Information Retrieval Conference, ISMIR 2013, Curitiba, Brazil, November 4-8, 2013*, 2013.
- [2] A. Anglade, C. Baccigalupo, N. Casagrande, Ò. Celma, and P. Lamere, "Workshop report: Womrad 2010," in *RecSys*, X. Amatriain, M. Torrens, P. Resnick, and M. Zanker, Eds. ACM, 2010, pp. 381–382.
- [3] A. Anglade, O. Celma, B. Fields, P. Lamere, and B. McFee, "Womrad: 2nd workshop on music recommendation and discovery," in *Proceedings of the fifth ACM conference on Recommender systems*, ser. RecSys '11. New York, NY, USA: ACM, 2011, pp. 381–382. [Online]. Available: <http://doi.acm.org/10.1145/2043932.2044013>
- [4] Q. Yang, I. King, Q. Li, P. Pu, and G. Karypis, Eds., *Seventh ACM Conference on Recommender Systems, RecSys '13, Hong Kong, China, October 12-16, 2013*. ACM, 2013.
- [5] N. Aizenberg, Y. Koren, and O. Somekh, "Build your own music recommender by modeling internet radio streams," in *WWW*, A. Mille, F. L. Gandon, J. Misselis, M. Rabinovich, and S. Staab, Eds. ACM, 2012, pp. 1–10.

- [6] O. Hilliges, P. Holzer, R. Klüber, and A. Butz, "Audioradar: A metaphorical visualization for the navigation of large music collections," in *Smart Graphics*, ser. Lecture Notes in Computer Science, A. Butz, B. D. Fisher, A. Krüger, and P. Olivier, Eds., vol. 4073. Springer, 2006, pp. 82–92.
- [7] D. F. Gleich, M. Rasmussen, K. Lang, and L. Zhukov, "The world of music: User ratings; spectral and spherical embeddings; map projections," Online report, 2006.
- [8] D. F. Gleich, L. Zhukov, M. Rasmussen, and K. Lang, "The World of Music: SDP Embedding of High Dimensional data," in *Information Visualization 2005*, 2005, interactive Poster.
- [9] K. Brandenburg, C. Dittmar, M. Gruhne, J. Abeer, H. Lukashevich, P. Dunker, D. Grtner, K. Wolter, and H. Grossmann, "Music search and recommendation," in *Handbook of Multimedia for Digital Entertainment and Arts*, B. Furht, Ed. Springer US, 2009, pp. 349–384.
- [10] Ö. Celma, *Music Recommendation and Discovery - The Long Tail, Long Fail, and Long Play in the Digital Music Space*. Springer, 2010.
- [11] P. Avesani, P. Massa, M. Nori, and A. Susi, "Collaborative radio community," in *Proceedings of the Second International Conference on Adaptive Hypermedia and Adaptive Web-Based Systems*, ser. AH '02. London, UK, UK: Springer-Verlag, 2002, pp. 462–465. [Online]. Available: <http://dl.acm.org/citation.cfm?id=647458.728391>
- [12] P. Symeonidis, M. M. Ruxanda, A. Nanopoulos, and Y. Manolopoulos, "Ternary semantic analysis of social tags for personalized music recommendation," in *ISMIR*, J. P. Bello, E. Chew, and D. Turnbull, Eds., 2008, pp. 219–224.
- [13] A. Nanopoulos, D. Rafailidis, P. Symeonidis, and Y. Manolopoulos, "Musicbox: Personalized music recommendation based on cubic analysis of social tags," *IEEE Transactions on Audio, Speech & Language Processing*, vol. 18, no. 2, pp. 407–412, 2010.
- [14] Y.-H. Yang, D. Bogdanov, P. Herrera, and M. Sordo, "Music retagging using label propagation and robust principal component analysis," in *WWW (Companion Volume)*, A. Mille, F. L. Gandon, J. Misselis, M. Rabinovich, and S. Staab, Eds. ACM, 2012, pp. 869–876.
- [15] N. Koenigstein, G. Dror, and Y. Koren, "Yahoo! music recommendations: modeling music ratings with temporal dynamics and item taxonomy," in *Proceedings of the fifth ACM conference on Recommender systems*, ser. RecSys '11. New York, NY, USA: ACM, 2011, pp. 165–172. [Online]. Available: <http://doi.acm.org/10.1145/2043932.2043964>
- [16] B. McFee, T. Bertin-Mahieux, D. P. W. Ellis, and G. R. G. Lanckriet, "The million song dataset challenge," in *WWW (Companion Volume)*, A. Mille, F. L. Gandon, J. Misselis, M. Rabinovich, and S. Staab, Eds. ACM, 2012, pp. 909–916.
- [17] O. Celma and P. Lamere, "Music recommendation and discovery revisited," in *Proceedings of the fifth ACM conference on Recommender systems*, ser. RecSys '11. New York, NY, USA: ACM, 2011, pp. 7–8. [Online]. Available: <http://doi.acm.org/10.1145/2043932.2043936>
- [18] Y. Hu and M. Oghara, "Nexttone player: A music recommendation system based on user behavior," in *ISMIR*, A. Klapuri and C. Leider, Eds. University of Miami, 2011, pp. 103–108.
- [19] D. Bogdanov and P. Herrera, "How much metadata do we need in music recommendation? a subjective evaluation using preference sets," in *ISMIR*, A. Klapuri and C. Leider, Eds. University of Miami, 2011, pp. 97–102.
- [20] D. Bogdanov, M. Haro, F. Fuhrmann, A. Xambó, E. Gómez, and P. Herrera, "Semantic audio content-based music recommendation and visualization based on user preference examples," *Inf. Process. Manage.*, vol. 49, no. 1, pp. 13–33, 2013.
- [21] C. S. Mesnage, A. Rafiq, S. Dixon, and R. P. Brixtel, "Music discovery with social networks," in *Workshop on Music Recommendation and Discovery 2011*. Workshop on Music Recommendation and Discovery, Oct. 2011, pp. 1–6. [Online]. Available: http://ceur-ws.org/Vol-793/womrad2011_complete.pdf
- [22] M. Barthet, A. Anglade, G. Fazekas, and R. Kolozali, Sefki Macrae, "Music recommendation for music learning: Hotttabs, a multimedia guitar tutor," in *Workshop on Music Recommendation and Discovery 2011*. Workshop on Music Recommendation and Discovery, Oct. 2011, pp. 7–13. [Online]. Available: http://ceur-ws.org/Vol-793/womrad2011_complete.pdf
- [23] I. Tatlı and A. Birtürk, "Using semantic relations in context-based music recommendations," in *Workshop on Music Recommendation and Discovery 2011*. Workshop on Music Recommendation and Discovery, Oct. 2011, pp. 14–17. [Online]. Available: http://ceur-ws.org/Vol-793/womrad2011_complete.pdf
- [24] P. Knees and M. Schedl, "Towards semantic music information extraction from the web using rule patterns and supervised learning," in *Workshop on Music Recommendation and Discovery 2011*. Workshop on Music Recommendation and Discovery, Oct. 2011, pp. 18–25. [Online]. Available: http://ceur-ws.org/Vol-793/womrad2011_complete.pdf
- [25] I. Knopke, "The importance of service and genre in recommendations for online radio and television programmes," in *Workshop on Music Recommendation and Discovery 2011*. Workshop on Music Recommendation and Discovery, Oct. 2011, pp. 26–29. [Online]. Available: http://ceur-ws.org/Vol-793/womrad2011_complete.pdf
- [26] G. Popescu and P. Pu, "Probabilistic game theoretic algorithms for group recommender systems," in *Workshop on Music Recommendation and Discovery 2011*. Workshop on Music Recommendation and Discovery, Oct. 2011, pp. 7–12. [Online]. Available: http://ceur-ws.org/Vol-793/womrad2011_complete.pdf
- [27] D. I. Ignatov, A. V. Konstantinov, S. I. Nikolenko, J. Poelmans, and V. Zaharchuk, "Online recommender system for radio station hosting," in *BIR*, ser. Lecture Notes in Business Information Processing, N. Aseeva, E. Babkin, and O. Kozyrev, Eds., vol. 128. Springer, 2012, pp. 1–12.
- [28] D. I. Ignatov, J. Poelmans, G. Dedene, and S. Viaene, "A New Cross-Validation Technique to Evaluate Quality of Recommender Systems," in *PerMin*, ser. LNCS, M. K. Kundu, S. Mitra, D. Mazumdar, and S. K. Pal, Eds., vol. 7143. Springer, 2012, pp. 195–202.
- [29] P. Symeonidis, A. Nanopoulos, A. N. Papadopoulos, and Y. Manolopoulos, "Nearest-biclusters collaborative filtering based on constant and coherent values," *Inf. Retr.*, vol. 11, no. 1, pp. 51–75, 2008.
- [30] A. Clauset, C. R. Shalizi, and M. E. J. Newman, "Power-law distributions in empirical data," *SIAM Rev.*, vol. 51, no. 4, pp. 661–703, Nov. 2009.
- [31] R. M. Bell and Y. Koren, "Scalable collaborative filtering with jointly derived neighborhood interpolation weights," in *Proceedings of the 7th IEEE International Conference on Data Mining*. Omaha, Nebraska, USA: IEEE Computer Society, 2007, pp. 43–52.
- [32] —, "Lessons from the Netflix Prize challenge," *SIGKDD Explorations*, vol. 9, no. 2, pp. 75–79, 2007.
- [33] Y. Koren, "Factorization meets the neighborhood: a multifaceted collaborative filtering model," in *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. Las Vegas, Nevada, USA: IEEE Computer Society, 2008, pp. 426–434.
- [34] Y. Koren, R. M. Bell, and C. Volinsky, "Matrix factorization techniques for recommender systems," *IEEE Computer*, vol. 42, no. 8, pp. 30–37, 2009.
- [35] Y. Koren and R. M. Bell, "Advances in collaborative filtering," in *Recommender Systems Handbook*, F. Ricci, L. Rokach, B. Shapira, and P. B. Kantor, Eds. Springer, 2011, pp. 145–186.
- [36] T. Vander Wal, "Folksonomy Coinage and Definition," 2007, <http://vanderwal.net/folksonomy.html> (accessed on 12.03.2012). [Online]. Available: <http://vanderwal.net/folksonomy.html>
- [37] D. I. Ignatov, S. O. Kuznetsov, J. Poelmans, and L. E. Zhukov, "Can triconcepts become triclusters?" *Int. J. General Systems*, vol. 42, no. 6, pp. 572–593, 2013.
- [38] A. Klapuri and C. Leider, Eds., *Proceedings of the 12th International Society for Music Information Retrieval Conference, ISMIR 2011, Miami, Florida, USA, October 24-28, 2011*. University of Miami, 2011.
- [39] A. Mille, F. L. Gandon, J. Misselis, M. Rabinovich, and S. Staab, Eds., *Proceedings of the 21st World Wide Web Conference, WWW 2012, Lyon, France, April 16-20, 2012 (Companion Volume)*. ACM, 2012.