

Random graphs for performance evaluation of  
recommender systems\*

by

Szymon Chojnacki and Mieczysław A. Kłopotek

Institute of Computer Science, Polish Academy of Sciences  
Warsaw, Poland

**Abstract:** The purpose of this article is to introduce a new analytical framework dedicated to measuring performance of recommender systems. A standard approach is to assess the quality of a system by means of accuracy related statistics. However, the specificity of the environments in which recommender systems are deployed requires paying much attention to speed and memory requirements of the algorithms. Unfortunately, it is implausible to assess accurately the complexity of various algorithms with formal tools. This can be attributed to the fact that such analyses are usually based on an assumption of dense representation of underlying data structures. In real life, though, the algorithms operate on sparse data and are implemented with collections dedicated for them. Therefore, we propose to measure the complexity of recommender systems with artificial datasets that possess real-life properties. We utilize a recently developed bipartite graph generator to evaluate how the state-of-art recommender system behavior is determined and diversified by topological properties of the generated datasets.

**Keywords:** recommender systems, performance evaluation, random graphs, bipartite complex networks.

## 1. Introduction

Recommender systems are an important component of the Intelligent Web. These systems make information retrieval easier and push users from typing queries towards clicking at suggested links. We experience real-life recommender systems when browsing for books, movies, news or music. Such engines are an essential part of such websites as *Amazon*, *MovieLens* or *Last.fm*. Recommender systems are used to deal with the tasks that are typical for statistical classification methods. They fit especially the scenarios in which the number of attributes, classes or missing values is large. Classic data-mining techniques like

---

\*Submitted: October 2010; Accepted: June 2011.

*logistic regression* or *decision trees* are well suited to predict which category of news is the most interesting for a particular customer. Recommender systems are used to output more fine-grained results and point at concrete stories.

In recent years we have observed a surge of interest among the research community in recommender systems. One of the events that was responsible for this phenomenon was the Netflix Prize challenge<sup>1</sup>. The competition was organized by a large DVD retailer in US. The prize of 1 million dollars was awarded to the team that managed to improve RMSE (root mean standard error) of the retailer's Cinematch algorithm by more than 10%. The lesson we learned during the Netflix Prize is that the difference between the quality of simple methods and the sophisticated ones is not as significant as we could have expected. Moreover, in order to lower RMSE an ensemble of complex and computationally intensive methods has to be implemented. Even though the organizers made much effort to deliver realistic and huge data sets, the setting did not envision the problems that we need to face in diverse real-life recommender systems applications, such as:

- the *Cold Start* problem, i.e. an arrival of new users with short history (e.g. restricted to the last HTTP session)
- instant creation of new items (e.g. news, auction items or photos)
- real-time feedback from users about our performance.

These drawbacks were overcome during the Online Task of the Discovery Challenge organized as a part of the ECML 2009 (see Jäschke et al., 2009). The owners of the [www.BibSonomy.org](http://www.BibSonomy.org) bookmarking portal opened its interfaces to recommender systems taking part in the evaluation. Whenever a user of the BibSonomy was bookmarking a digital resource (a publication or a website) a query was sent to all the systems. The tag recommendation of a randomly selected system was displayed to the user. After the action a feedback with user's actions was sent to all systems. The systems could have been maintained during the challenge, because they were configured as web services. The results showed that all of the teams found it difficult to deliver majority of its recommendations within time constraint of 1 000 ms.

Our research was motivated by the above result and an observation that the development of recommender systems is limited by the fact that there are not enough possibilities to test the algorithms with various datasets. The data structure used by recommender systems is a sparse *user*  $\times$  *item* matrix with ratings. It is a hard exercise to generate randomly such matrices with predefined properties, resembling real-life situations, because of three reasons. Firstly, if we fix the number of users, items and rankings and try to place the rankings randomly in the matrix we obtain binomial distribution of the number of items rated by users (and vice versa). However, in real-life datasets the distributions are heavy-tailed. Secondly, simple random selection results in no correlations

---

<sup>1</sup><http://www.netflix.com>

among user's preferences. Such correlations exist in real datasets. Thirdly, if we generate one matrix with some desired properties and would like to add new users or items, we would probably lose the properties of the original matrix.

The problem has been taken up recently in Chojnacki and Kłopotek (2011). We proposed to look at the matrix with ratings as if it were a bipartite graph with nodes of both modalities representing users and items, respectively. A rating from the matrix is mapped onto an edge in the bigraph. We proposed an algorithm in which we can control not only simple statistics like numbers of users, items or rankings, but also obtain skewed distributions and correlations among users or items. Moreover, our random bigraph generator asymptotic properties were verified with formal and numerical tools and we can add users or items to the graph without losing the properties of the original datasets. The algorithm was obtained by adapting the advances in unipartite complex networks modeling onto a bipartite ground. We modified the preferential attachment model by Barabási and Albert (1999) with the extension proposed in Liu et al. (2002) and generalized the *surfing mechanism* introduced by Vázquez (2003).

In this paper we apply the generator to produce several random bigraphs with various properties and evaluate how these properties impact on the performance of the analyzed recommender systems. We analyze four features of the systems that in our opinion are responsible for the success of an algorithm in a real-life setting:

- time required to build a model from scratch
- memory consumption of the trained model
- latency of creating a recommendation
- time of updating the model with new ratings.

We considered six algorithms during the tests: UserBased, ItemBased, SlopeOne, UserThreshold, KnnItem and SVD (Owen et al., 2010). We relied on high-performance Mahout library<sup>2</sup>. Our attention was focused on five properties of artificial datasets: (1) size of the graph, (2) relative number of edges, (3) proportion of the number of users to the number of items, (4) clustering of edges, (5) distributions of node degrees of both modalities. We also checked the influence of two characteristics on the performance of the UserBased model, namely: similarity measure and the size of the neighborhood.

Throughout the article we show that our approach can be used to better understand the features of datasets that are responsible for the performance of recommender systems. We utilize this knowledge to show (1) which algorithms are best suited for various scenarios, (2) identify dataset features that are correlated with improving performance of some algorithms and worsening performance of others, (3) point at potential directions of improving the implementations of the algorithms.

---

<sup>2</sup><http://www.mahout.apache.org>

The rest of the article is organized as follows. In Section 2 we describe how recommender systems are evaluated. In Section 3 we outline the details of the applied random bigraph generator. The fourth section contains the results of extensive experiments. The last, fifth section is dedicated for the concluding remarks.

## 2. Performance of recommender systems

In our research we perceive performance in terms of real-life speed of using an algorithm. We omit analysis of statistical indicators such as accuracy, recall, f-measure, RMSE or lift charts. This is because, even though they are crucial in the process of selecting an algorithm to be deployed, it seems pointless to evaluate these measures within randomly generated datasets. One could argue that the usefulness of artificial datasets is questionable. And the best strategy is to evaluate all possible algorithms with one's real dataset and choose the most accurate one that gives recommendations within specified time constraints. We argue with this point of view. Based on our experience we are more likely to believe that the structure and topology of datasets changes rapidly and the performance may be affected by an appearance of outlying observations or unexpected growth in scale. It is hard to foresee all potential pitfalls and the need to evaluate the algorithms within a wide range of artificial data emerges. In order to justify our deduction we present in Fig. 1 the results of an online evaluation of the systems that participated in the social bookmarking Discovery Challenge organized as a part of ECML'09.

The algorithms were deployed as web services and latency was measured as the difference of time between sending a request and receiving a list of top five recommended tags. We participated in the evaluation and managed to lower the latency of our system to the level of 400 ms (see Chojnacki, Czerski and Kłopotek, 2010). It turned out to be an important improvement for the users of the BibSonomy and resulted in the highest rate of clicks among all the evaluated systems.

A good starting point when analyzing complexity of any algorithm is to assess its asymptotic properties. Such analyses are usually based on several simplifying assumptions. One of the assumptions that is virtually never met is the dense data structure representation assumption. For example vectors are either implemented as dense (e.g. `double[]`) or sparse arrays (e.g. `ArrayList<Double>`). Dedicated collections are used to find optimal trade-offs between memory consumption and speed. Moreover, several advanced issues such as Java `Object` overhead and its influence on garbage collectors performance need to be taken into account. A great discussion about these problems is contained in chapter 3 of Owen et al. (2010).

Nonetheless, theoretical analysis of complexity strengthens our general intuition about the upper bound of time needed for computations. A thorough analysis of wide range of recommender algorithms is provided in Jahrer, Töscher and

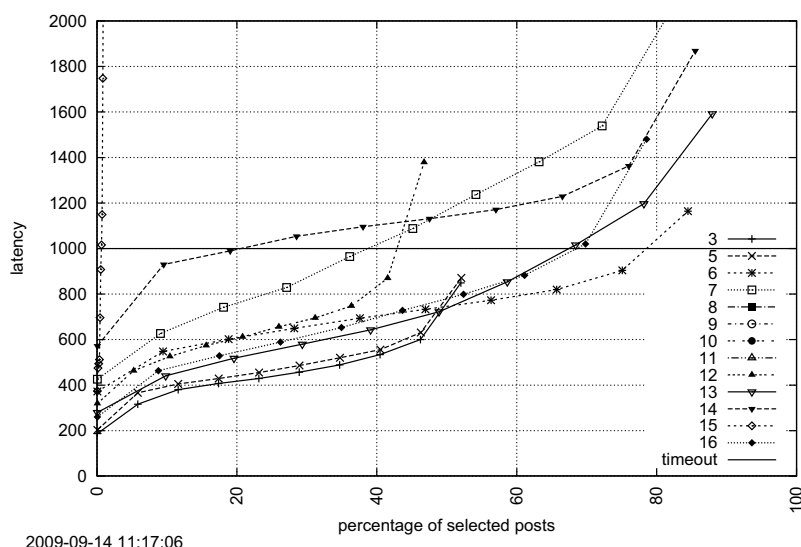


Figure 1. Latency of recommender systems evaluated during the Online Task of ECML-PKDD 2009 Discovery Challenge. None of the systems delivered recommendations to all queries. The points on curves show the percentiles of latency for each algorithm. For example the second most right point of algorithm 13 indicates that almost 90% of its recommendations were delivered within 1200ms. The same point tells us from the horizontal axis that these answers covered slightly below 80% of all queries. Source: <http://www.kde.cs.uni-kassel.de/ws/dc09/results/online/>

Legenstein (2010). In particular, it is proposed that training time, latency and memory consumption of SVD-based recommender are  $O(E)$ ,  $O(1)$ ,  $O(U + I)$ , where  $E$  is the number of ratings,  $U$  is the number of users and  $I$  is the number of items. The limitation of this kind of analysis is the fact that except for some rare situations it does not help us to decide which of two selected algorithms is expected to perform better on real data. An example of a set of soft rules meant to cope with such questions is drawn in Table 1.

In our experiments we did not find evidence to assert that UserBased algorithm performs significantly different than ItemBased when the proportion of the number of users to the number of items varies. We also found that there exist other factors than  $U$ ,  $I$  or  $E$  that are interpretable and impact on the performance in a coherent way. The results are discussed in detail in Sec. 4.

Table 1. Comparison of recommender systems, based on Owen et al. (2010)

Recommender	Key features
UserBased	Fast when number of users is relatively small
ItemBased	Fast when number of items is relatively small
SlopeOne	Recommendations and updates fast at runtime Requires large precomputations Suitable when number of items is relatively small
KnnItem	Good when number of items is relatively smaller
SVD	Good results Requires large precomputations

### 3. Bipartite random graph generator

In this section we describe an algorithm used to generate random bigraphs. The algorithm was introduced and described in detail in Chojnacki and Kłopotek (2011). In Sec. 3.1 we define the parameters of the algorithm, in Sec. 3.2 the properties of the generator are outlined.

#### 3.1. Generative procedure

The generative procedure consists of three steps: (1) new node creation, (2) edge attachment type selection and (3) running bouncing mechanism. The steps are run after an initialization of the bigraph. The procedure requires specifying eight parameters:

- $m$  – the number of initial loose edges with a user and an item at the ends
- $T$  – the number of iterations
- $p$  – the probability that a new node is a user  
( $1 - p$ ) is the probability that a new node is an item
- $u$  – the number of edges created by each new user
- $v$  – the number of edges created by each new item
- $\alpha$  – the probability that a new user's edge is being connected to an item with preferential attachment
- $\beta$  – the probability that a new item's edge is being connected to a user with preferential attachment
- $b$  – the fraction of preferentially attached edges that were created via a *bouncing mechanism*.

In the preferential attachment mechanism the probability that a node is drawn is linearly proportional to its degree. Opposite to the preferential attachment is random attachment, in which probability of selection is equal for all nodes. The model is based on an iterative repetition of three steps.

**Step 1** If a random number is greater than  $p$  create a new user with  $u$  loose edges, otherwise create a new item with  $v$  loose edges.

**Step 2** For each edge decide whether to join it to a node of the second modality randomly or with preferential attachment. The probability of selection preferential attachment is  $\alpha$  for new user and  $\beta$  for new item.

**Step 3** For each edge that is supposed to be created with preferential attachment decide if it should also be generated via a bouncing mechanism.

Bouncing is performed in three micro steps: (1) a random node is drawn from the nodes that are already joined with the new node, (2) a random neighbor of the drawn node is chosen, (3) a random neighbor of the neighbor is selected for joining with the new node. The bouncing mechanism was injected into the model in order to parametrize the level of transitivity in a graph. The transitivity is a feature of real datasets and in terms of recommender systems represent the correlations between items ranked by different users. In unipartite graphs transitivity is measured by the *local clustering coefficient*, which is calculated for each node as a number of edges among direct neighbors of the node divided by all possible pairs of the neighbors. In bipartite graphs the coefficient is always zero and is substituted by *bipartite local clustering coefficient (BLCC)*, Chojnacki and Klopotek (2011). Bipartite local clustering coefficient of node  $j$  takes values of one minus the proportion of second neighbors of the node to the potential number of the second neighbors of the node. The value of *BLCC* calculated for node  $j$  is given by:

$$BLCC_j = 1 - \frac{|N_2(j)|}{\sum_{i \in N_1(j)} (k_i - 1)}, \quad (1)$$

where  $|N_2(j)|$  stands for the number of the second neighbors of node  $j$ ,  $N_1(j)$  is a set of the first neighbors of node  $j$  and  $k_i$  is the degree of node  $i$ . The steps of the generator are depicted in Fig. 2.

### 3.2. Properties

One can see that after  $t$  iterations the bigraph consists of  $|U(t)| = m + pt$  users,  $|I(t)| = m + (1-p)t$  items, and  $|E(t)| = m + t(pu + (1-p)v)$  edges. Let us denote by  $\eta$  an average number of edges created during one iteration  $\eta = (pu + (1-p)v)$ . After relatively many iterations ( $t \gg m$ ) we can neglect  $m$ . In the presented model, an average user degree is:

$$\frac{|E(t)|}{|U(t)|} = \frac{m + t(pu + (1-p)v)}{m + pt} \approx \frac{\eta}{p}.$$

Analogously, an average item degree is:

$$\frac{|E(t)|}{|I(t)|} \approx \frac{\eta}{(1-p)}.$$

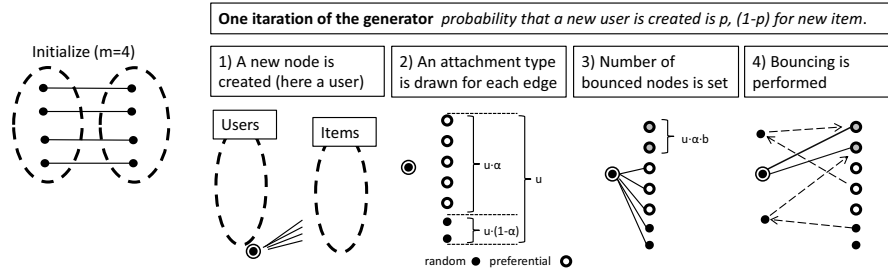


Figure 2. For each edge of a new node, that is to be connected with an existing node according to the preferential attachment mechanism, a decision is made whether to create it via a bouncing mechanism. In case of attaching new user node,  $u$  new edges are created. On average  $u \cdot \alpha$  edge endings are to be drawn preferentially and  $u \cdot \alpha \cdot b$  of them are to be obtained via bouncing from the nodes that are already selected

The values are time invariant, but depend on both  $u$  and  $v$ . In Figs. 3 to 5 three relations between model parameters and graph features are delineated.

## 4. Experiments

In order to evaluate the performance of the analyzed algorithms we generated 83 artificial bipartite graphs. The statistics describing the graphs are contained in Figs. 14 and 15 at the end of the paper. Each edge of the graph was augmented with a random integer from a set of possible rankings  $\{0, 1, 2, 3, 4, 5\}$ . After the last iteration (usually  $T = 10\,000$ ) one hundred more edges were created by running 100 steps for each graph with unchanged parameters. This enabled us to preserve asymptotic properties of the graphs within a set of rankings used to batch update of the models. The experiments were run in-memory within separate threads on a 64-bit Fedora operating system with the Quatro 2.66GHz Intel(R) Core(TM) i5 CPUs.

### 4.1. Evaluated systems

We evaluated six recommender algorithms implemented in the Mahout java library. Mahout contains highly efficient open-source implementations of machine-learning algorithms maintained by a vibrant community. It is powering several portals e.g. *SpeedDate*, *Yahoo! Mail*, *AOL* or *Mippin*. The algorithms are: GenericUserBasedRecommender (Herlocker et al., 1999), GenericItemBasedRecommender (Sarwar et al., 2001), SlopeOneRecommender (Lemire and MacLachlan, 2005), GenericUserBasedRecommender with the neighborhood defined by non-negative threshold similarity (Owen et al., 2010), KnnItemBasedRecommender (Bell and Koren, 2007) and SVDRecommender (Zhang et al., 2005).



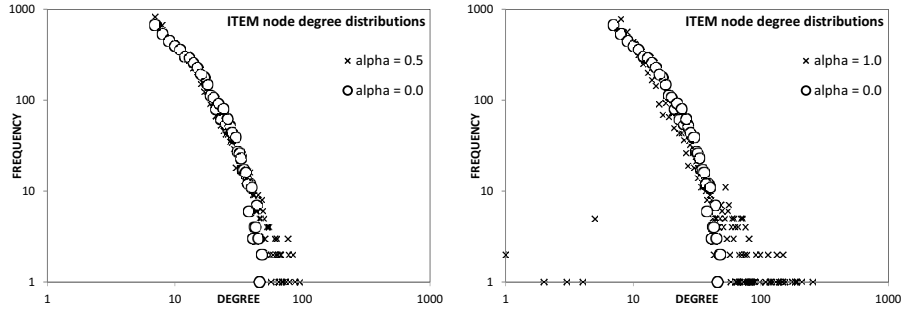


Figure 3. Parameters  $\alpha$  and  $\beta$  enable us to control the shape of node degree distributions. As the values approach unity, the shape becomes power-law, as the values tend to zero the shape tends to exponential

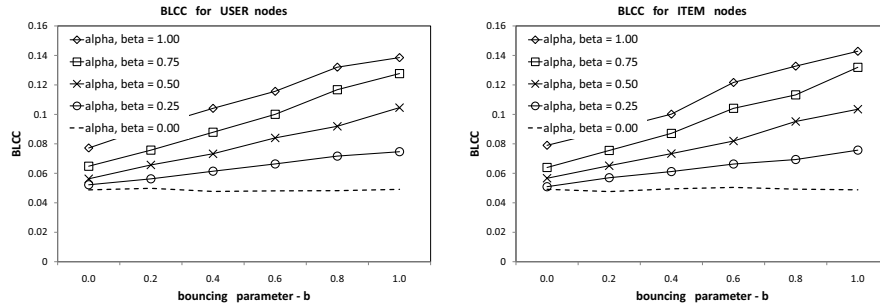


Figure 4. The growth of parameter  $b$  results in higher values of BLCC. It does not influence BLCC if  $\alpha = \beta = 0.0$

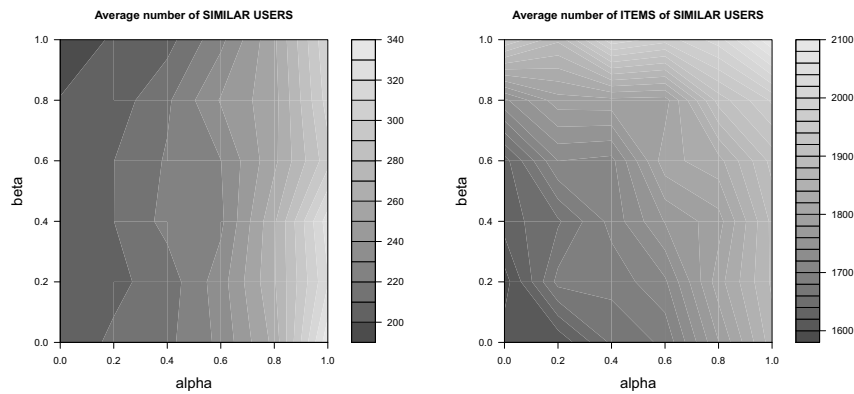


Figure 5. The average number of similar users (having ranked at least one item in common with a defined user) and their items depends on both  $\alpha$  and  $\beta$

The algorithms cover wide spectrum of approaches to the problems of *Collaborative Filtering*.

#### 4.2. Performance measures

In what follows we study the correlations between the performance of the six algorithms differentiated by five graph properties: size, density, proportion of users to items, clustering and shape of node degree distributions. We focus our attention on four performance statistics:

1. **BUILD** - time in milliseconds that is required to load whole bigraph from a text file and train a model, after this period of time the model is ready to create recommendations
2. **MEMORY** - memory consumption in megabytes of the built model, it was assessed as a difference between `Runtime.getRuntime().totalMemory()` and `Runtime.getRuntime().freeMemory()` after calling a garbage collector five times
3. **LATENCY** - an average time in milliseconds required to produce a recommendation for a sample of 500 users
4. **UPDATE** - time in millisecond of updating a model with 100 new ratings.

In case of UserBased recommender we set the size of the neighborhood to 200. SVD recommender was projected onto 10 factors and 200 iterations were run during training. The *Pearson Correlation* was used as a similarity measure. At the end we analyze specific parameters of the UserBased model.

#### 4.3. Scalability

In order to verify the ability of the analyzed algorithms to scale, we generated thirteen bigraphs. Each dataset was produced with the same parameters except of the number of iterations. The datasets are numbered from 10 to 22 in Fig. 14. The time of training the SVD recommender may be misleading<sup>3</sup>. For most systems, only training time grows linearly with the size of a dataset. Memory consumption and latency grows sublinearly. SlopeOne and SVD algorithms exhibit the poorest performance in terms of building time and memory consumption. These costs pay back in the phase of creating recommendations. KnnItem does not scale during training and is the slowest when creating recommendations. Except for SlopeOne all models refresh their structures immediately. Only the latency of SVD seems not to depend on the scale.

---

<sup>3</sup>In the chart (Fig. 6) the time of SVD building was divided by 1 000 to get comparable results. Such long time of building the model is a result of running 200 iterations of the gradient descent. However, even a single iteration takes on average  $1\,000/200 = 5$  times longer than building of the other models.

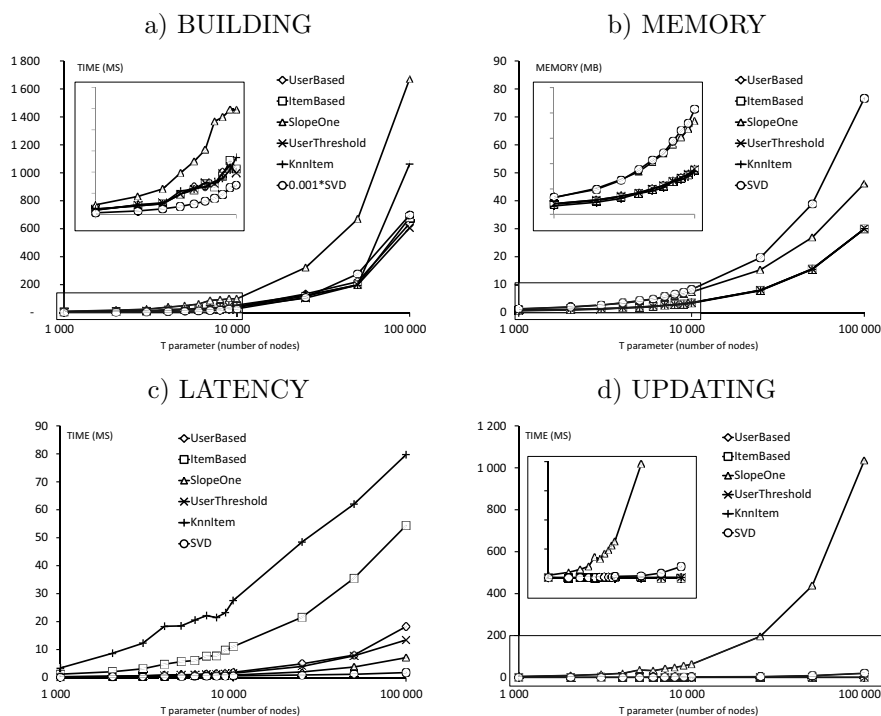


Figure 6. Performance of the recommenders conditioned by the size of bigraphs. On all charts horizontal axis has a logarithmic scale, allowing for even distribution of all observations

#### 4.4. Density

We generated 14 bigraphs to test how performance depends on density. The graphs are numbered from 23 to 36 in Fig. 14. All graphs have around the same number of nodes, but the number of edges varies between 30 100 and 240 100. The diversity was obtained by changing two parameters:  $u$  and  $v$ .

We omit the presentation of performance calculated for the first four bigraphs, in which  $u = v$ . All the four measures were growing steadily as we evaluated ( $u = v = 3$ ), ( $u = v = 6$ ), ( $u = v = 12$ ) and ( $u = v = 24$ ). This could have suggested that there exists a strong correlation between density and performance. However, when we compare five pairs of graphs with  $u \neq v$  (Fig. 7), the results become confusing. When setting  $u = 3$  and iterating over  $\{4, \dots, 15\}$  with  $v$ , the performance diminishes only slightly. If we freeze  $v = 3$  and iterate over  $u$ , we decrease the performance steadily. This suggests that the density on its own is not so much responsible for the performance<sup>4</sup>.

<sup>4</sup>Digging into this effect deeper we see that by increasing  $v$  we do not only increase the

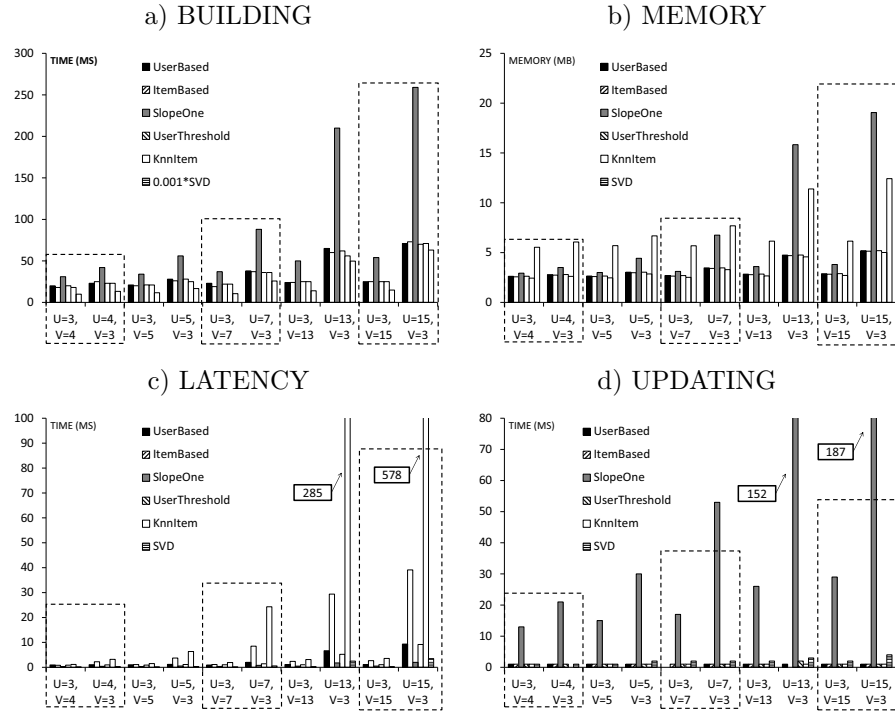


Figure 7. Performance of the recommenders versus the density of the graphs.

#### 4.5. Users to items proportion

Parameter  $p$  in the generator enables us to control the proportion of users to items. The parameter is interpreted as the probability that a new node is a user. We built nine graphs with constant numbers of nodes ( $= 10\ 200$ ) and edges ( $= 70\ 100$ ). The number of users varies between 1 166 and 9 082. The graphs are the first nine graphs in Fig. 14.

The relationship between  $p$  and performance is depicted in Fig. 8. Only the performance of SVD does not seem to depend on the proportion in none of four evaluations. Training time decreases as the relative number of users increases for all remaining models. The only model, in which memory consumption depends on the proportion is SlopeOne. It is very interesting that the lowest memory requirements are obtained for balanced graphs ( $p \approx 0.5$ ). The same non-monotonic relation can be observed in case of KnnItem and latency. The rest of the algorithms improve their latency as a relative number of items grows.

density, but also lower the variance of node degree distribution in user modality stronger than for the item modality. This shows that performance relies heavily on the distributions of node degrees, but with a different magnitude for both modalities.

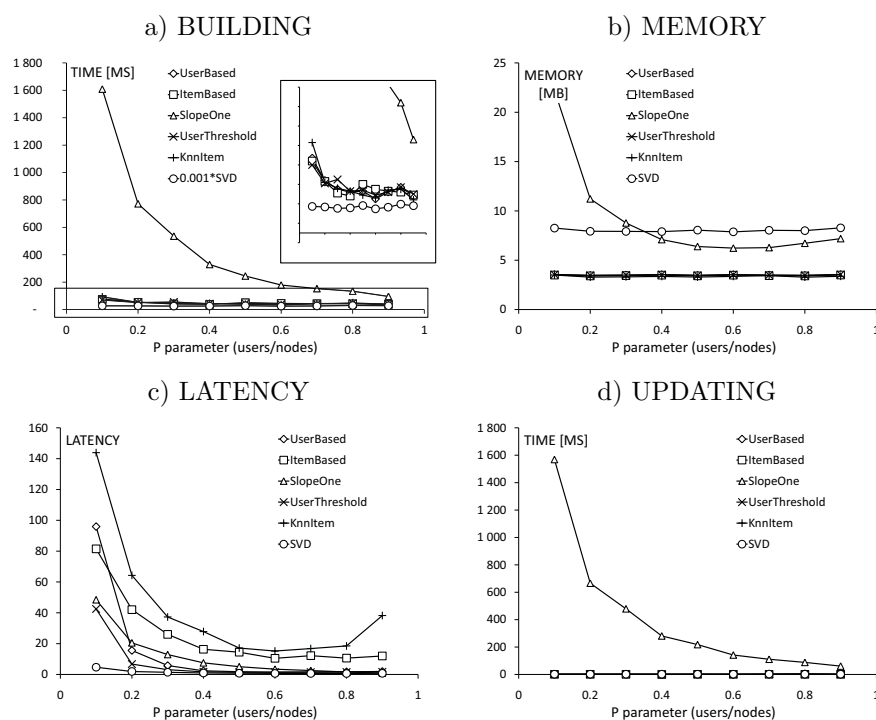


Figure 8. Performance modeled by the proportion of users to all nodes.

Even though most datasets have more items than users, the opposite situations can also happen. For example, the dataset evaluated during track II of the KDD Cup 2011 has more items than users. The fact that ItemBased recommender does not perform better than UserBased when there are relatively few items, Owen et al. (2010), raises our concern<sup>5</sup>.

#### 4.6. Clustering

Eleven graphs were generated to measure the influence of clustering on performance. The datasets are numbered from 37 to 47 in Figs. 14 and 15. The bigraphs were generated by changing the bouncing parameter  $b$  from 0.0 to 1.0 with 0.1 intervals. All graphs have 1 100 nodes and 120 100 edges.

The performance of SVD does not depend on bouncing. Memory consumption and time of building grow slightly in line with clustering only in case of SlopeOne model. The latency of creating a recommendation is correlated with

<sup>5</sup>This observation requires further investigation. In particular, we plan to verify if the claim that *UserBased algorithm is preferred over ItemBased when there are fewer users* is valid only when caching mechanisms are implemented

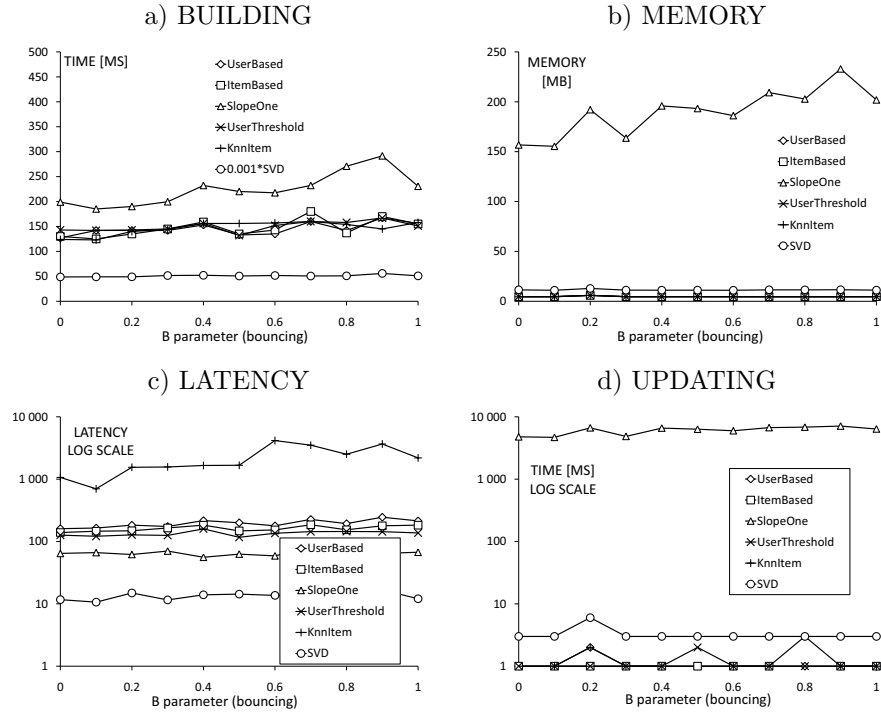


Figure 9. Performance dimensioned by various levels of clustering (transitivity)

clustering gently only in case of ItemBase, UserBase, KnnItem and UserThreshold algorithms. The relation between clustering and performance is in all variants weak and not stable.

#### 4.7. Shapes of degree distributions

We have shown in Chojnacki and Kłopotek (2011) that by changing  $\alpha$  and  $\beta$  we can output graphs with a mixture of the power-law distribution and the exponential distribution (compare Fig. 3). The two parameters enable us to control independently the shapes of degree distributions of both modalities. By increasing  $\alpha$  we can obtain graphs with constant number of users, items and edges, but growing average number of potentially similar users<sup>6</sup>. Moreover, an average number of distinct items of potentially similar users grows with both  $\alpha$  and  $\beta$  (compare Fig. 5). These results are consistent with the asymptotic Newman's formula, Newman et al. (2001). The formula is based on a local tree-like structure assumption and estimates the average number of second neighbors  $|N_2(j)|$

<sup>6</sup>We say that two users are potentially similar if they have rated at least one item in common.

of a random user  $j$  by means of the first moment of user degree distribution ( $\langle U \rangle$  and the first and the second moments of item degree distribution ( $\langle I \rangle$  and  $\langle I^2 \rangle$  respectively):

$$|N_2(j)| = \langle U \rangle \left( \frac{\langle I^2 \rangle}{\langle I \rangle} - 1 \right). \quad (2)$$

We have generated 36 bigraphs with all possible combinations of  $\alpha$  and  $\beta$  from the set of values  $\{0.0, 0.2, 0.4, 0.6, 0.8, 1.0\}$ . All the graphs have 10 000 nodes and 70 100 edges. They are numbered from 48 to 83 in Fig. 15. We observed that an average latency grows with either  $\alpha$  or  $\beta$  for all analyzed algorithms (Fig. 11). We also observed that building time and memory consumption is related to  $\alpha$  and  $\beta$  for all algorithms except for SVD and KnnItem. The performance of SlopeOne model is drawn in Fig. 10.

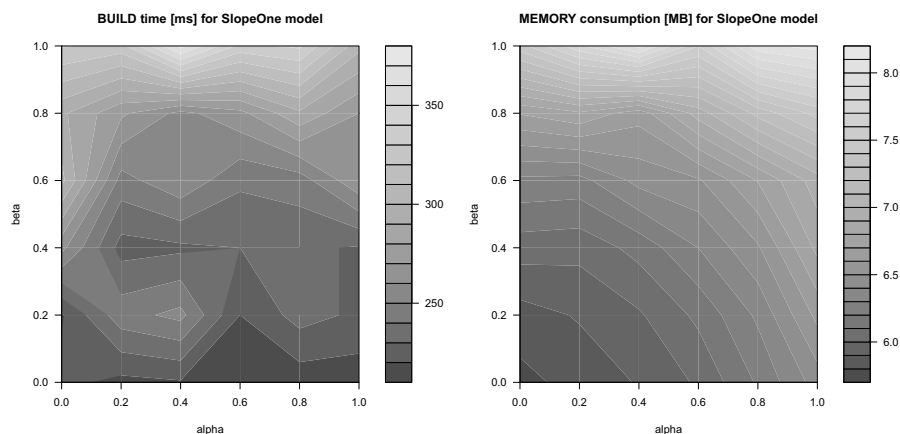


Figure 10. Training time and memory requirements of SlopeOne model

#### 4.8. Similarity measure

In the last two subsections we focus on UserBased algorithm and two parameters that are specific for it. The first parameter is the similarity measure and will be discussed in this subsection. The second parameter is the size of the neighborhood and will be described in next subsection.

The similarity between two vectors  $\mathbf{x} = (x_1, \dots, x_n)$  and  $\mathbf{y} = (y_1, \dots, y_n)$  can be defined in several ways. We analyzed all that are available in the Mahout library i.e. *Pearson similarity*, *Euclidean similarity*, *LogLikelihood similarity*, *Spearman similarity* and *Tanimoto similarity*. The third and the fifth are defined by set operations and can be easily used with *binary ratings*, when we only know if a user expressed a preference for an item or not. **Pearson similarity**

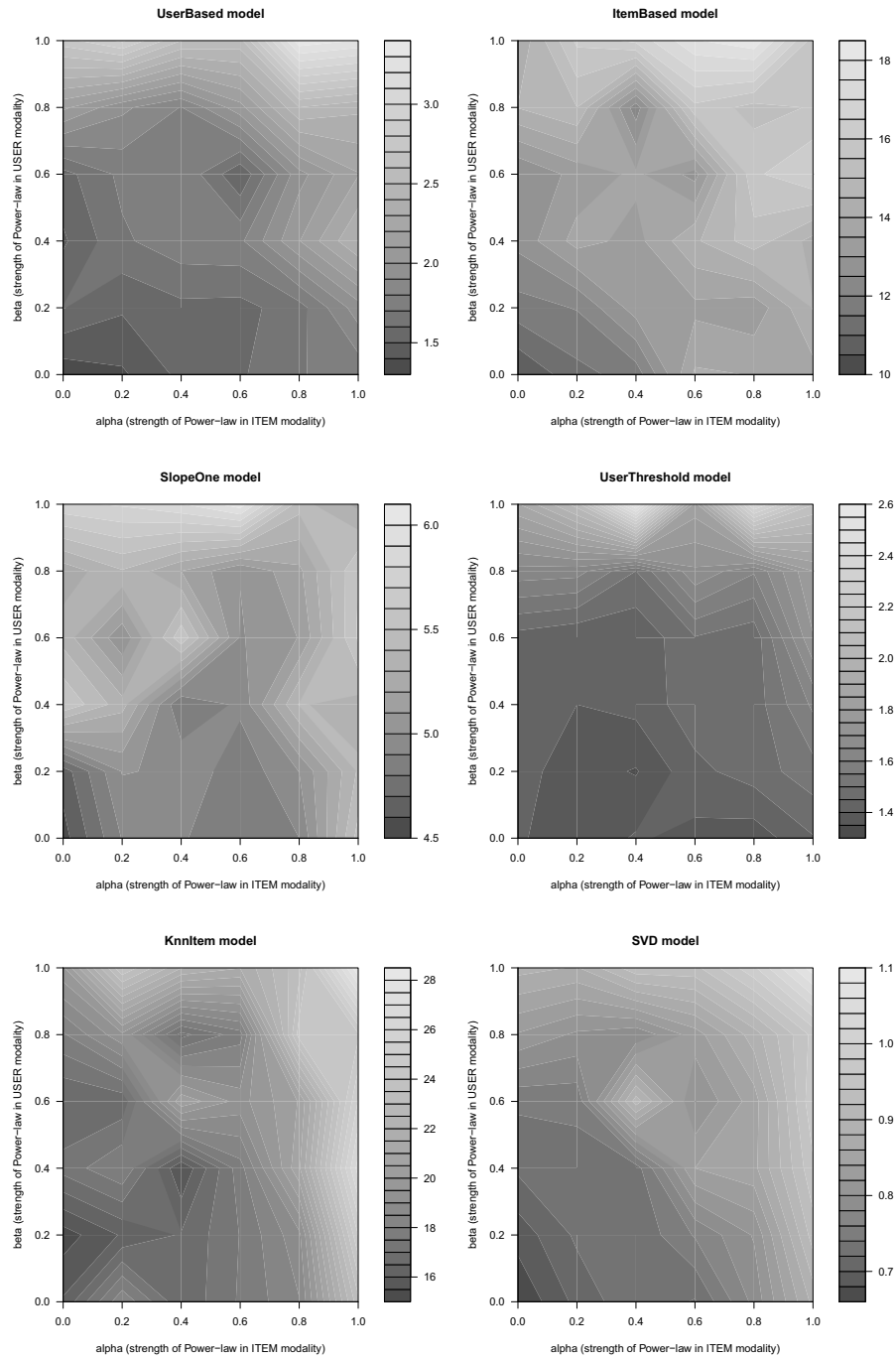


Figure 11. The influence of  $\alpha$  and  $\beta$  parameters on latency



is widely used and is related to the *Cosine similarity*. The only difference is that the former operates on centered data. **Euclidean similarity** is inversely proportional to the Euclidean distance between two vectors. **LogLikelihood similarity** is based on calculating four values: number of non-empty dimensions in both vectors, numbers of non-missing dimensions in the first and the second vector, and the number of all dimensions, Dunning (1993). **Spearman similarity** is calculated as Pearson similarity but  $x_i$  and  $y_j$  are substituted with their relative ranks i.e. the lowest value is 1, the second lowest is 2, and so on. **Tanimoto similarity** for binary data is calculated as a proportion of dimensions that are non-empty in both vectors to total number of non-empty dimensions in the vectors. Detailed definitions of all coefficients can be found in the javadoc API of Mahout.

We have observed in Fig. 12 that the latency of UserBased model depends on the selection of similarity measure. The fastest recommendations appear for the Pearson and Euclidean similarities. Tanimoto and Loglikelihood similarities are around three times slower. Spearman similarity is the slowest and does not scale well with the growth of density.

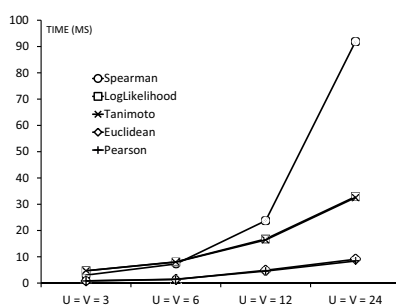


Figure 12. Dependence of latency on similarity metrics. Four analyzed graphs have the same number of nodes, they are numbered 23-26 in Fig. 14

#### 4.9. Size of the neighborhood

Size of the neighborhood is a parameter in UserBased recommender, which enables us to tune the accuracy of the algorithm. When a UserBased model is requested to deliver recommendations for a given user, two steps are performed. In the first step, the similarity of the user to all other users is calculated and only the most similar users are retained (the number is limited by the neighborhood). In the second step only items of the most similar users are weighted and the recommendation is selected among those items.

The results in Fig. 13 indicate that as long as graph density is below some threshold, the latency does not depend neither on density nor size of the neighborhood. This can be explained by the fact that the neighborhood parameter is triggered only when the number of similar users is greater than it.

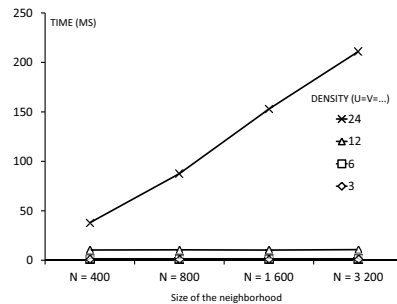


Figure 13. Latency of creating a recommendation in four graphs with varying density. UserBased algorithm run with four different levels of the neighborhood

## 5. Conclusion

In the article we have proposed to use random bipartite graphs to measure the performance of recommender systems. We have shown that the recently developed random bigraph generator, Chojnacki and Kłopotek (2011), can be used to generate a wide range of artificial datasets with predefined properties. The analytical framework can be used to compare various algorithms, but also to help us understand their complexity and point at potentially non-optimal implementations. We believe that the proposed methodology can be applied in various scenarios and settings. Further analyses need to be performed to understand the most intriguing results. In particular, this applies to the real relationship between the complexity of UserBased and ItemBased algorithms. Another issue is the emergence of U-shaped curve describing performance of various recommenders when the proportion of the number of users to the number of items is being changed. In the future we plan to evaluate how the performance of recommender systems is altered when they are implemented in a distributed environment.

## References

- BARABÁSI, A.L. and ALBERT, R. (1999) Emergence of Scaling in Random Networks. *Science*, **286** (5439), 509–512.
- BELL, R.M. and KOREN, Y. (2007) Scalable Collaborative Filtering with Jointly Derived Neighborhood Interpolation Weights. In: *Proc. of the 2007 Seventh IEEE International Conference on Data Mining*, IEEE Computer Society, 43–52.
- CHOJNACKI, S., CZERSKI, D. and KŁOPOTEK, M. (2010) Optimization of tag recommender systems in a real life setting. In: *3rd Conference on Human System Interaction*, Rzeszow, Poland. IEEE, 107–112.

- CHOJNACKI, S. and KLOPOTEK, M. (2011) Random Graphs for Bipartite Networks Modeling. *Journal of Control and Cybernetics*, submitted.
- DUNNING, T. (1993) Accurate Methods for the Statistics of Surprise and Coincidence. *Computational Linguistics*, **19** (1), 61–74.
- HERLOCKER, J.L., KONSTAN, J.A., BORCHERS, A. and RIEDL, J. (1999) An algorithmic framework for performing collaborative filtering. In: *Proc. of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval*, ACM Press, 230–237.
- JÄHRER, M., TÖSCHER, A. and LEGENSTEIN, R. (2010) Combining predictions for accurate recommender systems. In: *Proc. of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining (KDD '10)*, ACM, 693–702.
- JÄSCHKE, R., EISTERLEHNER, F., HOTHO, A. and STUMME, G. (2009) Testing and Evaluating Tag Recommenders in a Live System. In: D. Benz and F. Janssen, eds., *Workshop on Knowledge Discovery, Data Mining, and Machine Learning*. ACM, 44–51.
- LEMIRE, D. and MACLACHLAN, A. (2005) Slope One Predictors for Online Rating-Based Collaborative Filtering. In: *Proc. of SIAM Data Mining (SDM'05)*. ACM, 21–23.
- LIU, Z., LAI, Y-C. and DASGUPTA, P. (2002) Connectivity distribution and attack tolerance of general networks with both preferential and random attachments. *Physics Letters A*, **303** (5-6), 337–344.
- NEWMAN, M., STROGATZ, S. and WATTS, D. (2001) Random graphs with arbitrary degree distributions and their applications. *Phys. Rev. E*, **64** (2).
- OWEN, S., ANIL, R., DUNNING, T. and FRIEDMAN, E. (2010) *Mahout in action (MEAP)*. Manning Publication Co.
- SARWAR, B.M., KARYPIS, G., KONSTAN, J.A. and RIEDL, J. (2001) Item-based collaborative filtering recommendation algorithms. In: *Proc. of the 10th international conference on World Wide Web (WWW '01)*, ACM, 285–295.
- VÁZQUEZ, A. (2003) Growing network with local rules: Preferential attachment, clustering hierarchy, and degree correlations. *Phys. Rev. E*, **67** (5).
- ZHANG, S., WANG, W., FORD, J., MAKEDON, F. and PEARLMAN, J. (2005) Using singular value decomposition approximation for collaborative filtering. In: *Proc. of the Seventh IEEE International Conference on E-Commerce Technology*, IEEE Computer Society, 257–264.

No.	Generator's parameters								Graphs' properties				
	m	T	p	u	v	alpha	beta	b	users	items	edges	BLCC for users	BLCC for items
1	100	10 000	0.9	7	7	0.5	0.5	0.0	9 082	1 118	70 100	4%	23%
2	100	10 000	0.8	7	7	0.5	0.5	0.0	8 105	2 095	70 100	3%	10%
3	100	10 000	0.7	7	7	0.5	0.5	0.0	7 060	3 140	70 100	3%	6%
4	100	10 000	0.6	7	7	0.5	0.5	0.0	6 187	4 013	70 100	3%	4%
5	100	10 000	0.5	7	7	0.5	0.5	0.0	4 998	5 202	70 100	3%	3%
6	100	10 000	0.4	7	7	0.5	0.5	0.0	4 157	6 043	70 100	4%	3%
7	100	10 000	0.3	7	7	0.5	0.5	0.0	3 081	7 119	70 100	6%	3%
8	100	10 000	0.2	7	7	0.5	0.5	0.0	2 172	8 028	70 100	10%	3%
9	100	10 000	0.1	7	7	0.5	0.5	0.0	1 166	9 034	70 100	22%	4%
10	100	1 000	0.9	7	7	0.1	0.1	0.0	993	207	7 100	11%	40%
11	100	2 000	0.9	7	7	0.1	0.1	0.0	1 897	303	14 100	9%	39%
12	100	3 000	0.9	7	7	0.1	0.1	0.0	2 833	367	21 100	8%	39%
13	100	4 000	0.9	7	7	0.1	0.1	0.0	3 716	484	28 100	7%	34%
14	100	5 000	0.9	7	7	0.1	0.1	0.0	4 600	600	35 100	6%	30%
15	100	6 000	0.9	7	7	0.1	0.1	0.0	5 558	642	42 100	5%	31%
16	100	7 000	0.9	7	7	0.1	0.1	0.0	6 420	780	49 100	4%	27%
17	100	8 000	0.9	7	7	0.1	0.1	0.0	7 291	909	56 100	4%	24%
18	100	9 000	0.9	7	7	0.1	0.1	0.0	8 235	965	63 100	4%	24%
19	100	10 000	0.9	7	7	0.1	0.1	0.0	9 120	1 080	70 100	3%	22%
20	100	25 000	0.9	7	7	0.1	0.1	0.0	22 627	2 573	175 100	2%	13%
21	100	50 000	0.9	7	7	0.1	0.1	0.0	45 214	4 986	350 100	1%	8%
22	100	100 000	0.9	7	7	0.1	0.1	0.0	90 192	10 008	700 100	1%	5%
23	100	10 000	0.9	3	3	0.1	0.1	0.0	9 063	1 137	30 100	0%	4%
24	100	10 000	0.9	6	6	0.1	0.1	0.0	9 106	1 094	60 100	2%	17%
25	100	10 000	0.9	12	12	0.1	0.1	0.0	9 126	1 074	120 100	11%	43%
26	100	10 000	0.9	24	24	0.1	0.1	0.0	9 113	1 087	240 100	37%	71%
27	100	10 000	0.9	3	4	0.1	0.1	0.0	9 110	1 090	31 090	0%	5%
28	100	10 000	0.9	3	5	0.1	0.1	0.0	9 076	1 124	32 148	0%	5%
29	100	10 000	0.9	3	7	0.1	0.1	0.0	9 084	1 116	34 164	1%	6%
30	100	10 000	0.9	3	13	0.1	0.1	0.0	9 076	1 124	40 340	1%	8%
31	100	10 000	0.9	3	15	0.1	0.1	0.0	9 089	1 111	42 232	1%	9%
32	100	10 000	0.9	4	3	0.1	0.1	0.0	9 108	1 092	39 108	1%	8%
33	100	10 000	0.9	5	3	0.1	0.1	0.0	9 133	1 067	48 166	2%	12%
34	100	10 000	0.9	7	3	0.1	0.1	0.0	9 123	1 077	66 192	3%	21%
35	100	10 000	0.9	13	3	0.1	0.1	0.0	9 077	1 123	119 870	12%	40%
36	100	10 000	0.9	15	3	0.1	0.1	0.0	9 114	1 086	138 268	15%	47%
37	100	10 000	0.2	12	12	0.9	0.9	0.0	2 039	8 161	120 100	29%	16%
38	100	10 000	0.2	12	12	0.9	0.9	0.1	2 128	8 072	120 100	29%	16%
39	100	10 000	0.2	12	12	0.9	0.9	0.2	2 026	8 174	120 100	30%	17%
40	100	10 000	0.2	12	12	0.9	0.9	0.3	2 146	8 054	120 100	30%	16%
41	100	10 000	0.2	12	12	0.9	0.9	0.4	2 065	8 135	120 100	31%	18%

Figure 14. Random graphs generated to test the performance of recommender systems (part 1). Generator parameters and BLCC are defined in Sec. 3.1

No.	Generator's parameters								Graphs' properties				
	m	T	p	u	v	alpha	beta	b	users	items	edges	BLCC for users	BLCC for items
42	100	10 000	0.2	12	12	0.9	0.9	<b>0.5</b>	2 106	8 094	120 100	32%	18%
43	100	10 000	0.2	12	12	0.9	0.9	<b>0.6</b>	2 100	8 100	120 100	32%	17%
44	100	10 000	0.2	12	12	0.9	0.9	<b>0.7</b>	2 022	8 178	120 100	33%	18%
45	100	10 000	0.2	12	12	0.9	0.9	<b>0.8</b>	2 013	8 187	120 100	34%	18%
46	100	10 000	0.2	12	12	0.9	0.9	<b>0.9</b>	2 031	8 169	120 100	35%	19%
47	100	10 000	0.2	12	12	0.9	0.9	<b>1.0</b>	2 120	8 080	120 100	34%	18%
48	100	10 000	0.5	7	7	<b>1.0</b>	<b>1.0</b>	0.0	5 070	5 130	70 100	4%	4%
49	100	10 000	0.5	7	7	<b>1.0</b>	<b>0.8</b>	0.0	5 058	5 142	70 100	4%	4%
50	100	10 000	0.5	7	7	<b>1.0</b>	<b>0.6</b>	0.0	5 054	5 146	70 100	4%	4%
51	100	10 000	0.5	7	7	<b>1.0</b>	<b>0.4</b>	0.0	5 130	5 070	70 100	4%	4%
52	100	10 000	0.5	7	7	<b>1.0</b>	<b>0.2</b>	0.0	5 027	5 173	70 100	4%	3%
53	100	10 000	0.5	7	7	<b>1.0</b>	<b>0.0</b>	0.0	5 125	5 075	70 100	3%	3%
54	100	10 000	0.5	7	7	<b>0.8</b>	<b>1.0</b>	0.0	5 141	5 059	70 100	4%	4%
55	100	10 000	0.5	7	7	<b>0.8</b>	<b>0.8</b>	0.0	5 114	5 086	70 100	4%	4%
56	100	10 000	0.5	7	7	<b>0.8</b>	<b>0.6</b>	0.0	5 112	5 088	70 100	3%	3%
57	100	10 000	0.5	7	7	<b>0.8</b>	<b>0.4</b>	0.0	5 088	5 112	70 100	3%	3%
58	100	10 000	0.5	7	7	<b>0.8</b>	<b>0.2</b>	0.0	5 130	5 070	70 100	3%	3%
59	100	10 000	0.5	7	7	<b>0.8</b>	<b>0.0</b>	0.0	5 113	5 087	70 100	3%	3%
60	100	10 000	0.5	7	7	<b>0.6</b>	<b>1.0</b>	0.0	5 086	5 114	70 100	4%	4%
61	100	10 000	0.5	7	7	<b>0.6</b>	<b>0.8</b>	0.0	5 170	5 030	70 100	3%	3%
62	100	10 000	0.5	7	7	<b>0.6</b>	<b>0.6</b>	0.0	5 127	5 073	70 100	3%	3%
63	100	10 000	0.5	7	7	<b>0.6</b>	<b>0.4</b>	0.0	5 081	5 119	70 100	3%	3%
64	100	10 000	0.5	7	7	<b>0.6</b>	<b>0.2</b>	0.0	5 085	5 115	70 100	3%	3%
65	100	10 000	0.5	7	7	<b>0.6</b>	<b>0.0</b>	0.0	5 117	5 083	70 100	3%	3%
66	100	10 000	0.5	7	7	<b>0.4</b>	<b>1.0</b>	0.0	5 145	5 055	70 100	4%	4%
67	100	10 000	0.5	7	7	<b>0.4</b>	<b>0.8</b>	0.0	5 098	5 102	70 100	3%	3%
68	100	10 000	0.5	7	7	<b>0.4</b>	<b>0.6</b>	0.0	5 009	5 191	70 100	3%	3%
69	100	10 000	0.5	7	7	<b>0.4</b>	<b>0.4</b>	0.0	5 176	5 024	70 100	3%	3%
70	100	10 000	0.5	7	7	<b>0.4</b>	<b>0.2</b>	0.0	5 052	5 148	70 100	3%	3%
71	100	10 000	0.5	7	7	<b>0.4</b>	<b>0.0</b>	0.0	5 042	5 158	70 100	3%	3%
72	100	10 000	0.5	7	7	<b>0.2</b>	<b>1.0</b>	0.0	5 030	5 170	70 100	4%	4%
73	100	10 000	0.5	7	7	<b>0.2</b>	<b>0.8</b>	0.0	5 063	5 137	70 100	3%	3%
74	100	10 000	0.5	7	7	<b>0.2</b>	<b>0.6</b>	0.0	5 156	5 044	70 100	3%	3%
75	100	10 000	0.5	7	7	<b>0.2</b>	<b>0.4</b>	0.0	5 158	5 042	70 100	3%	3%
76	100	10 000	0.5	7	7	<b>0.2</b>	<b>0.2</b>	0.0	5 084	5 116	70 100	3%	3%
77	100	10 000	0.5	7	7	<b>0.2</b>	<b>0.0</b>	0.0	5 039	5 161	70 100	3%	3%
78	100	10 000	0.5	7	7	<b>0.0</b>	<b>1.0</b>	0.0	5 054	5 146	70 100	3%	3%
79	100	10 000	0.5	7	7	<b>0.0</b>	<b>0.8</b>	0.0	5 102	5 098	70 100	3%	3%
80	100	10 000	0.5	7	7	<b>0.0</b>	<b>0.6</b>	0.0	5 044	5 156	70 100	3%	3%
81	100	10 000	0.5	7	7	<b>0.0</b>	<b>0.4</b>	0.0	4 993	5 207	70 100	3%	3%
82	100	10 000	0.5	7	7	<b>0.0</b>	<b>0.2</b>	0.0	5 145	5 055	70 100	3%	3%
83	100	10 000	0.5	7	7	<b>0.0</b>	<b>0.0</b>	0.0	5 188	5 012	70 100	2%	3%

Figure 15. Random graphs generated to test the performance of recommender systems (part 2). Generator parameters and BLCC are defined in Sec. 3.1.

