

# Learning What People (Don't) Want

Thomas Hofmann

Department of Computer Science  
Brown University, Providence, RI 20912  
th@cs.brown.edu

**Abstract.** Recommender systems make use of a database of user ratings to generate personalized recommendations and help people to find relevant products, items, or documents. In this paper, we present a probabilistic, model-based framework for user ratings based on a novel collaborative filtering technique that performs an automatic decomposition of user preferences. Our approach has several benefits, including highly accurate predictions, task-optimized model learning, mining of interest groups and patterns, as well as a highly efficient and scalable computation of predictions and recommendation lists.

## 1 Introduction

Knowing what one really wants is far from obvious. In order to figure out which book to read, which movie to watch, or which Web site to visit, we often rely on advice given to us by other people. Yet, in many situations one would like to automate the process of recommending items by anticipating user preferences. For example, in an electronic commerce setting, it is in the best interest of the seller as well as the customer to automatically receive accurate recommendations. Similarly, in the case of news groups or Web communities one would like to have some form of automated support, since these user groups are often far too large to allow its members to directly share experiences or to interchange recommendation. And even if one is lucky enough to get a person's advice on a subject, the question of who's advice to trust on which issue always remains a crucial question.

*Recommender systems* have been proposed and utilized in this context to leverage existing user data (ratings, profiles, logs) and to help people share their evaluations. Prominent Internet sites like amazon.com or CDnow use such systems for personalized product recommendation and to complement the traditional content-oriented search functionality. The predominant technique applied in the academic as well as the commercial world is *collaborative filtering* [3] and the most popular methods are variations around the memory-based approach proposed in the GroupLens [9, 8] project. The underlying principle is extremely simple - for a given (so-called *active*) user, find people with similar interests and use their ratings and judgements to recommend new items. Technical details mainly deal with the question of how to compute the similarity between user profiles and of how to exactly perform the combination step.

In this paper, we propose a radically different approach to collaborative filtering which uses a sparse matrix decomposition technique called Probabilistic Latent Semantic Analysis [5] to automatically discover preference patterns in user profile data. In order to motivate our approach, we will first briefly discuss the dominant paradigm for collaborative filtering in Section 2 and point out what we think are its conceptual flaws. In Section 3, we will then present our approach which centers around a latent class model that can be varied and extended in several ways. Finally, Section 4 will present an experimental evaluation and illustration on the largest publicly available data set, the EachMovie data.

## 2 Memory-Based Approaches

The by far most popular class of algorithms for collaborative filtering relies on the following principle. Given an active user  $a$  and a database of user profiles, recommendations are determined by, (i) computing the similarity between the active user and all the users in the dataset, (ii) for each item, forming a weighted and properly normalized vote over all (or some) users in the database, where the weights reflect the similarity computed in (i).

More specifically, the Pearson correlation has been proposed in many landmark papers [8, 10], which computes the similarity between two user profiles as

$$w(a, u) = \sum_{y \in V_a \cap V_u} \frac{(v_{a,y} - \bar{v}_a)(v_{u,y} - \bar{v}_u)}{\sqrt{\sigma_a \sigma_u}}. \quad (1)$$

Here  $v_{u,y}$  denotes the (known) vote of user  $u$  on item  $y$ , while  $\bar{v}_u$  and  $\sigma_u$  refer to the mean and standard deviation of votes of user  $u$ , respectively. Based on these weights a predictive vote for the active user is computed as follows,

$$p_{a,y} = \frac{\sum_{u: y \in V_u} w(a, u) v_{u,y}}{\sum_{u: y \in V_u} |w(a, u)|}. \quad (2)$$

Variations of this basic methodology include the use of different correlation measures like the Spearman correlation [4], the thresholding of weights, and the inhomogeneous weighting of different items in computing the correlation between profiles. These methods are often called *memory-based* methods or *neighbor-based* methods, since predicted votes are computed directly from the database of user profiles.

What is the fundamental assumption underlying this class of algorithms? Most importantly in our opinion, it is the idea of a *constant similarity* measure between user profiles. Here “constant” refers to the fact that  $w(a, u)$  does not depend on the actual item  $y$  under consideration.<sup>1</sup> This is a very strong assumption, since it pre-supposes that in making a prediction for an active user, the

<sup>1</sup> Only the overall normalization in Eq. (2) is different for every item, although one can hardly consider this to be an “item-specific” weighting.

“degree of trust” in the rating of some other user does not depend on the specific item.

This can be best illustrated with an example: imagine we would like to recommend music to an active user  $A$  who has a strong interest in opera. In the database, we find a very “post-modern character”,  $B$ , who adores opera but at the same time likes Salsa rhythms. Should one recommend Salsa to the purist opera fan? Or should one effectively eliminate  $B$  from the neighbor set of  $A$ , although  $B$ ’s opera recommendations could be highly relevant? Similarly, in the opposite situation, imagine that “typically” opera and Salsa are negatively correlated and most people who like one genre dislike the other one. This would imply that one would either recommend items of one or the other type to  $B$ , but never both, dependent on the exact similarity weights and proportions.

Intuitively, what is missing in standard memory-based approaches is the possibility to model that fact that one person is a reliable recommender for another person with respect to a *subset of items*, but not necessarily for all possible items. This deficit is expected to become more relevant with increasing diversity of the item set and users’ interest patterns. In general, we strongly believe that the simple notion of similarity between one person and another person fails to capture the multi-dimensional nature of human preferences.

We see the solution to this problem in a method that couples recommendation with a decomposition of preferences, where each preference pattern models an interest or trait shared among a community of people. There will be a reciprocal characterization of user communities by sets of items and vice versa, along with the typical ratings that describe the nature of the relation, be it positive, neutral, or negative.

### 3 Decomposing Preferences

#### 3.1 Preview

We propose a probabilistic model-based approach to collaborative filtering that overcomes the major limitations of memory-based methods and offers a number of additional advantages:

- User profiles are explicitly decomposed into statistically significant *interest patterns*, each interest pattern typically dealing with a subset of the possible items. This allows to selectively share recommendations among users, i.e., effectively introduces a dependency on the specific item or type of item and thereby increases accuracy.
- The preference decomposition not only potentially improves the recommender accuracy, it also perform *data mining* by revealing hidden patterns that drive user ratings. This can be used for identifying user communities and for data visualization.
- The probabilistic preference model can be used to tailor recommender systems to the specific task at hand. For example, if each item comes with a specific cost or benefit (e.g., for a seller), one can optimize the recommendations in order to maximize the expected utility.

- Once the model is trained, the database is no longer required to make recommendations. A model-based approach thus avoids the need to have a potentially huge database available to make recommendations. This has advantages in terms of memory requirements as well as system speed.

### 3.2 Latent Class Models for User Ratings

There are two variants of probabilistic models that we investigate in this paper that correspond to different settings of the recommendation problem. (I) Predicting a user rating for an unknown item, i.e., estimating  $P(v|u, y)$  or  $p_{u,y} \equiv \sum_v v P(v|u, y)$  and (II) predicting a rating in conjunction with a selection, i.e., estimating  $P(v, y|u)$ . Which setting is more appropriate depends on the actual application and objective. For example, in an electronic commerce setting,  $P(y|u)$  may model the fact whether a user is likely to purchase a product, which may be an important part of the system, irrespective of the actual user satisfaction expressed by the rating  $v$ .

To be more formal, let us introduce a latent variable, denoted by  $z_{u,y}$ , for each (actual and potential) observation triplet  $(u, y, v_{u,y})$ , thus defining the complete data of an elementary rating event by  $(u, y, v_{u,y}, z_{u,y})$ . Intuitively, each quadruple is supposed to model the fact that a person  $u$  votes on item  $y$  with rating  $v_{u,y}$ , “because of”  $z_{u,y}$ . Each latent variable is intended to offer a hypothetical explanation for a rating that is itself not directly observable. For example, a high rating of a person on a piece of music, say “Fidelio”, might be explained by the fact that “Fidelio” is an opera and that the person is interested in this particular genre. However, an alternative explanation could be that a person likes works by a specific composer (Beethoven in our case) and the fact that “Fidelio” is an opera might be completely irrelevant. Yet a third person might like “Fidelio” because of its political message and its praise of freedom and love. Of course, it is not clear *a priori* which possible explanations should be considered. And even for a finite candidate set of causes the question remains to what cause is in effect in each individual case. Each item may be liked or disliked for different reasons and persons may have more than one interest, different aspects of their taste being relevant for different ratings. In summary, we are looking for a method that would automatically find potential causes, determine which subset of the causes are likely to be relevant for a specific item as well as for a specific person, and in each individual case assign a probability to the fact that a cause will be “active” for a given rating.

Assume for now that the number of potential causes is fixed beforehand, i.e., each variable  $z_{u,y}$  may take one out of a finite number of say  $K$  values. Then the joint probability of an observation triplet will be given by summing over all possible states of the latent variable,

$$P(u, y, v) = \sum_z P(u, y, v, z) = \sum_z P(y, v|z, u)P(z|u)P(u). \quad (3)$$

We make the crucial independence assumption that conditioned on the true cause  $z$  of a rating, the vote as well as the selection of the item is independent

of the users identity which implies that  $P(y, v|z, u) = P(y, v|z)$ . The resulting model can thus be written as

$$P(v, y|u) = \sum_z P(v, y|z)P(z|u) \quad (4)$$

where we have neglected the user probability  $P(u)$ . In the first case, where one conditions on both, the user as well as the item, one will get similarly

$$P(v|u, y) = \sum_z P(v|y, z)P(z|u) \quad (5)$$

These parameters have a very intuitive meaning:  $P(z|u)$  represents to what extent a user “participates” in a common interest pattern, or more precisely, which fraction of a user’s ratings are explained by hidden cause  $z$ .  $P(v, y|z)$  models which items and item/vote combinations are more likely or less likely to occur in the interest group described by  $z$ . Similarly,  $P(v|y, z)$  models the probability that a user as a member of interest group  $z$  will vote with  $v$  on item  $y$ .

One can also investigate models in which the role of users and items are reversed, in particular in the symmetric setting of estimating  $P(v|u, y)$ . Notice that in the typical case where the number of users exceeds the number of items, the dimensionality of the resulting model will be quite different.

The presented approach is closely related to sparse matrix decomposition techniques and the paper effectively extends a model called *Probabilistic Latent Semantic Analysis* which has been used mainly in the context of information retrieval [6]. More details about the conceptional foundations can be found in [5]. A simplified model and preliminary experiments for collaborative filtering have been presented in [7].

### 3.3 Expectation Maximization Algorithm

The Expectation Maximization (EM) algorithm is a standard method for maximum likelihood estimation in latent variable models. The EM algorithm performs two steps in alternation, an E-step, where the expected value for the unobserved variables are computed given the current parameters and a M-step, where these values are used to update the parameters. More details can be found in [5]. In our case the E-step computes the probability that a cause  $z$  is associated with an observation triplet  $(u, y, v_{u,y})$ . By Bayes’ rule one obtains

$$P(z|u, y, v_{u,y}) = \frac{P(z|u)\rho(v_{u,y}, z, y)}{\sum_{z'} P(z'|u)\rho(v_{u,y}, z', y)}, \quad (6)$$

where  $\rho(v_{u,y}, z, y) = P(v_{u,y}|z, y)$  in case (I) and  $\rho(v_{u,y}, z, y) = P(v_{u,y}, y|z)$  in case (II). Notice that the occurrence of a vote on  $y$  matters in the computation of the posterior probabilities for  $z$  according to (II), but not in case (I).

In the M-step one updates the parameters, i.e., re-estimates the conditional probabilities that define the model. The user participation probabilities are given

by

$$P(z|u) = \frac{\sum_{(y,v_{u,y})} P(z|u, y, v_{u,y})}{\sum_{z'} \sum_{(y,v_{u,y})} P(z'|u, y, v_{u,y})} = \frac{\sum_{(y,v_{u,y})} P(z|u, y, v_{u,y})}{|\{(u, y, v_{u,y})\}|}. \quad (7)$$

The update equations for  $P(v|y, z)$  and  $P(v, y|z)$ , respectively, depend on the parametric form that one chooses for those distributions. We assume that the rating scale is quantized, e.g., a rating  $v$  might be an integral number between 0 and 5 (as it is the case for the EachMovie data set). Other distributions, like a truncated normal distribution for a continuous rating variable are possible and pose no additional conceptual problem. In the quantized case one can work directly with multinomial/binomial probability distributions, leading to

$$P(v|y, z) = \frac{\sum_u P(z|u, y, v)}{\sum_{v'} \sum_u P(z|u, y, v')}. \quad (8)$$

Other choices are possible, for example a normal distribution with mean  $\mu_{y,z}$  and variance  $\sigma_{y,z}^2$  for each combination of items and hidden causes. In this case, the mean would simply be re-estimated in the M-step according to

$$\mu_{y,z} = \frac{\sum_v \sum_u v P(z|u, y, v)}{\sum_v \sum_u P(z|u, y, v)}. \quad (9)$$

The Gaussian model has shown slightly worse performance in the experiments, we have therefore focused on the multinomial/binomial model in the experiments.

## 4 Experimental Evaluation

### 4.1 Evaluation Criteria

A thorough empirical analysis of collaborative filtering algorithms has been presented in [1] and we have adapted most of the proposed evaluation criteria. The effectiveness of collaborative filtering techniques can be measured in various ways dependent on how the recommender system is used and how results are presented to the user.

In the case that individual items are presented to the user, we use two evaluation metrics. The first scores is the absolute deviation of the predicted and actual vote,

$$E_1(u, y) = |p_{u,y} - v_{u,y}|. \quad (10)$$

Moreover, for quantized votes, we also evaluate the score

$$E_0(u, y) = \begin{cases} 0, & \text{if } p_{u,y} = v_{u,y} \\ 1, & \text{else.} \end{cases} \quad (11)$$

which only measures whether or not the vote was correctly predicted.

Method	Absolute Deviation	Relative Improvement	Prediction Accuracy	Relative Improvement
Baseline	1.091	0.00	33.4	0.00
Memory-based	0.951	12.8%	35.3	6.3%
Model-based				
K=5	0.972	10.9%	39.3	17.6%
K=10	0.951	12.8%	39.8	19.1%
K=20	0.947	13.2%	39.9	19.5%
K=50	0.937	14.1%	40.1	20.0%
K=100	0.927	15.0%	40.6	21.6%
K=150	0.926	15.1%	40.7	21.9%
K=200	0.924	15.3%	40.8	22.2%
Model-based (Role of items and users reversed)				
K=20	0.983	9.9%	37.7	12.8%
K=50	0.975	10.6%	38.4	14.9%
K=100	0.973	10.8%	38.5	15.2%

**Table 1.** Results for predicting individual test ratings on the Each Movie data set.  $K$  denotes the number of latent states used to perform the decomposition.

For ranked lists of items we have to assign score to permutations of items. We denote a permutation by  $\tau$  and the rank of an item  $y$  with respect to  $\tau$  by  $\tau(y)$ , e.g., the top ranked item  $y$  will have  $\tau(y) = 1$ , the second item  $\tau(y) = 2$ , and so forth. Items with known votes are not included in the ranking. We then use the following rank score for  $\tau$ ,

$$R(u, \tau) = \sum_y 2^{-\frac{\tau(y)-1}{\alpha-1}} \max(v_{u,y} - \bar{v}, 0), \quad (12)$$

with  $\bar{v}$  denoting the overall mean vote. The rationale behind this score is that when presented with a ranked list of items, users will sift through the list starting at the top, until they find a relevant item or simply give up. The probability that a user will ever take notice of an item at rank  $r$  is modeled as an exponential distribution with a half-life constant  $\alpha$  (set to 4 in our experiments). The total score for a population of users is then measured by (cf. [1])

$$R = 100 \frac{\sum_u R(u, \tau_u)}{\sum_u \max_{\tau'} R(u, \tau')}. \quad (13)$$

## 4.2 Evaluation Protocols

In the evaluation experiments the observed user ratings are partitioned into a training set and a test set. For each user the training data is used to compute votes on unrated items as well as a ranked list of recommendations. For the split into training and test data, we have used a leave-one-out protocol, where

Method	Rank Score	Relative Improvement
Baseline	26.95	0.00
Memory-based	36.71	
Model-based		
K=20	44.64	65.6%
K=50	45.91	70.3%
K=100	45.98	70.6%

**Table 2.** Rank score results for ordered recommendation lists on the Each Movie data set.  $K$  denotes the number of latent states used to perform the decomposition.

we randomly leave out one of the ratings for each user who has at least two observed ratings. This has been called the *All but 1* protocol in [1].

Notice that the above ranking measure takes a very simple form in this case: If the hold-out rating is below the average, simply skip this user. Otherwise, compute the rank of the hold out “positive” item and evaluate the expression corresponding to this item. The maximum achievable score would be to have the test item always be in first place.

### 4.3 Data Set

Unfortunately, there are very few data sets available to the academic community to test recommender systems. This seems largely because of the potential commercial value of user data and because of privacy issues. The data set we have used in our experiments is the EachMovie data[2]. There are 1623 items in this data set and 40??? user profiles with a total of over 2 million ratings. The rating scale is discrete taking values from 0 to 5.<sup>2</sup>

### 4.4 Results

We have compared our statistical approach with the memory based approach using the Pearson correlation as a similarity measure. The results on predicting ratings for given items are summarize in Table 1. As one can see, the model-based approach achieves great performance gains in terms of absolute error as well as in terms of prediction accuracy. According to both scores, the model-based approach performs significantly better. The improvement grows with the number of states chosen for the latent class variables, but more or less levels out at around  $K = 100$ . Notice that we have used early stopping on validation data to avoid overfitting (cf. [5] for more details on complexity control issues).

Although we have not implemented other collaborative filtering techniques like Bayesian clustering or Bayesian networks, the comparison in [1] have demonstrated that memory-based methods achieve excellent results on this benchmark.

<sup>2</sup> The original ratings have been multiplied by a factor of 5.

This makes the results in Table 1 even more remarkable. It can also be seen that models with reversed role of users and items show a weaker performance.

The results for generating ranked recommendation lists, the more typical setting for recommender systems are shown in Table 2. Notice that in this case the model denoted by case (II) has been used, which not only predicts the outcome of the rating, but also whether or not a person is familiar with a particular item and/or has rated an item. The achieved performance gains are even more substantial than for predicting ratings of single items. The relative performance gain over the baseline method of ranking items by overall popularity is more than 70%.

#### 4.5 Performance Issues

The advantages of the model-based approach in terms of memory and computation time are also considerable. In terms of memory requirement, first notice that we only have to store the parameters  $P(v, y|z)$  or  $P(v, |y, z)$ , respectively, and can simply reconstruct the user-specific parameters  $P(z|u)$  as needed. In fact, it is simple to show that the optimization problem of finding  $P(z|u)$  reduces to a convex cross-entropy minimization problem. Since the number of users is typically much larger than the number of items, this can make a crucial difference. On the other hand, memory-based approaches suffer greatly from the fact that the data can not be compressed into a model, but has to be kept in main memory at recommendation time.

Similarly there are great savings in terms of computation time. In memory-based approaches, user correlations have typically to be performed on-line which is a very time consuming process. In comparison, the model-based approach only requires to perform of the order of  $K$  operations to compute the probability for a single rating or rating/item pair. Since the probabilities  $P(z|u)$  are typically very sparse, the actual savings will often be even larger for a slight sacrifice in accuracy. In our experiments, the model-based approach was typically more than 10-20 times faster compared to the memory-based approach. In practice, memory-based methods often have to subsample the user profile base to remain scalable to allow real-time recommendations for large databases. This is not necessary for our model-based approach.

The main disadvantage of the model-based approach is the computational burden of the model training stage. For the EachMovie data, a typical training run took between 5 and 30 minutes on a Pentium III w/800 MHz (dependent on the dimensionality of the model). Yet, the training can be performed off-line where computational resources are typically abundant or at least much less critical.

#### 4.6 Mining User Preferences

Finally, we would like to illustrate that the decomposition of user ratings may lead to the discovery of interesting patterns and regularities that describe user

interest as well as dislikes. To that extent we have visualized the items for each latent variable state by sorting them according to the popularity within an interest group as measured by  $P(y|z)$  (irrespective of the actual vote). The average vote  $\sum_v vP(v|y, z)$  is displayed in rectangular brackets as well. Figure 1 and 2 display the interest groups extracted by the model with  $K = 40$ , ordered according to the average “positiveness” of each group, computed as  $\sum_{y,v} vP(v|y, z)P(y|z)$ .

#### 4.7 Conclusion

We have presented a powerful method for collaborative filtering and mining of user data. The method achieves a very good recommendation and prediction accuracy, is highly scalable, and extremely flexible. Conceptionally, the decomposition of user preferences is a radically novel idea that clearly distinguishes this approach from traditional memory-based approaches.

#### References

1. J. S. Breese, D. Heckerman, and C. Kardie. Empirical analysis of predictive algorithms for collaborative filtering. In *Proceedings of the 14th Conference on Uncertainty in Artificial Intelligence*, pages 43–52, 1998.
2. [www.research.digital.com/src/eachmovie/](http://www.research.digital.com/src/eachmovie/).
3. D. Goldberg, D. Nichols, B. Oki, and D. Terry. Using collaborative filtering to weave an information tapestry. *Communications of the ACM*, 35(2):61–70, 1992.
4. J. Herlocker, Konstanm J., Al Borchers, and J. Riedl. An algorithmic framework for collaborative filtering. In *Proceedings of SIGIR'99*, 1999.
5. T. Hofmann. Probabilistic latent semantic analysis. In *Proceedings of the 15th Conference on Uncertainty in AI*, pages 289–296, 1999.
6. T. Hofmann. Probabilistic latent semantic indexingb. In *Proceedings of the 22nd ACM-SIGIR International Conference on Research and Development in Information Retrieval, Berkeley, California*, pages 50–57, 1999.
7. T. Hofmann and J. Puzicha. Latent class models for collaborative filtering. In *Proceedings of the International Joint Conference on Artificial Intelligence*, 1999.
8. J. Konstan, B. Miller, D. Maltz, J. Herlocker, L. Gordon, and J. Riedl. GroupLens: Applying collaborative filtering to Usenet news. *Communications of the ACM*, 40(3):77–87, 1997.
9. P. Resnik, I. Neophytos, M. Suchak, P. Bergstrom, and J. Riedl. GroupLens: An open architecture for collaborative filtering of netnews. In *Proceedings of ACM CSCW'94 Conference on Computer-Supported Cooperative Work*, pages 175–186, 1994.
10. U. Shardanand and P. Maes. Social information filtering: Algorithms for automatic 'word of mouth'. In *ACM CHI'95 Conference on Human Factors in Computing Systems*, pages 210–217, 1995.

Interest Group 1, *4.8*	Interest Group 2, *4.6*	Interest Group 3, *4.5*	Interest Group 4, *4.4*	Interest Group 5, *4.4*
<a href="#">Twister</a> [4.6*] [0.064]	<a href="#">Batman (1989)</a> [4.1*] [0.066]	<a href="#">Trainspotting</a> [5*] [0.038]	<a href="#">Dead Man Walking</a> [5*] [0.052]	<a href="#">The Santa Clause</a> [4.5*] [0.014]
<a href="#">Independence Day (...)</a> [4.9*] [0.061]	<a href="#">Apollo 13</a> [5*] [0.065]	<a href="#">Fargo</a> [5*] [0.033]	<a href="#">The Truth about Ca...</a> [4.3*] [0.039]	<a href="#">Casper</a> [4.5*] [0.014]
<a href="#">Toy Story</a> [4.9*] [0.057]	<a href="#">True Lies</a> [4.7*] [0.059]	<a href="#">Pulp Fiction</a> [5*] [0.028]	<a href="#">Get Shorty</a> [4.6*] [0.036]	<a href="#">Robin Hood: Men in...</a> [4.3*] [0.013]
<a href="#">Broken Arrow</a> [4.4*] [0.054]	<a href="#">Batman Forever</a> [4.1*] [0.054]	<a href="#">Clerks</a> [4.7*] [0.023]	<a href="#">Sense and Sensibil...</a> [5*] [0.035]	<a href="#">Tommy Boy</a> [4.5*] [0.013]
Interest Group 6, *4.3*	Interest Group 7, *4.3*	Interest Group 8, *4.2*	Interest Group 9, *4*	Interest Group 10, *3.9*
<a href="#">The Remains of the...</a> [4.5*] [0.047]	<a href="#">The Empire Strikes...</a> [4.7*] [0.032]	<a href="#">Pretty Woman</a> [4.3*] [0.059]	<a href="#">Sleepers</a> [4.2*] [0.015]	<a href="#">A Clockwork Orange...</a> [4.2*] [0.01]
<a href="#">The Piano</a> [4.7*] [0.043]	<a href="#">Raiders of the Los...</a> [4.7*] [0.03]	<a href="#">Mrs. Doubtfire</a> [4.3*] [0.059]	<a href="#">Jerry Maguire</a> [4.6*] [0.013]	<a href="#">Amadeus (1984)</a> [4.2*] [0.0098]
<a href="#">Like Water For Cho...</a> [4.7*] [0.043]	<a href="#">Star Wars</a> [4.9*] [0.026]	<a href="#">Ghost</a> [4.4*] [0.057]	<a href="#">The First Wives Cl...</a> [3.8*] [0.013]	<a href="#">Psycho (1960)</a> [4.3*] [0.0098]
<a href="#">Much Ado About Not...</a> [4.6*] [0.041]	<a href="#">Indiana Jones and ...</a> [4.5*] [0.025]	<a href="#">Sleepless in Seatt...</a> [4.4*] [0.055]	<a href="#">William Shakespear...</a> [4.5*] [0.011]	<a href="#">One Flew Over the ...</a> [4.5*] [0.0095]
Interest Group 11, *3.8*	Interest Group 12, *3.8*	Interest Group 13, *3.7*	Interest Group 14, *3.7*	Interest Group 15, *3.7*
<a href="#">Independence Day (...)</a> [4.1*] [0.05]	<a href="#">The Professional</a> [4.3*] [0.034]	<a href="#">Boys on the Side</a> [3.8*] [0.022]	<a href="#">The Usual Suspects</a> [4.4*] [0.063]	<a href="#">Searching For Bobb...</a> [4.2*] [0.031]
<a href="#">Twister</a> [3.8*] [0.043]	<a href="#">The Crow</a> [4.1*] [0.03]	<a href="#">Only You</a> [3.5*] [0.021]	<a href="#">Quiz Show</a> [3.9*] [0.052]	<a href="#">Snow White and the...</a> [3.9*] [0.03]
<a href="#">Mission: Impossibl...</a> [3.7*] [0.041]	<a href="#">Demolition Man</a> [3.4*] [0.028]	<a href="#">Something to Talk ...</a> [3.3*] [0.02]	<a href="#">Get Shorty</a> [3.8*] [0.052]	<a href="#">Pinocchio (1940)</a> [3.9*] [0.025]
<a href="#">Toy Story</a> [4.1*] [0.038]	<a href="#">Species</a> [3.1*] [0.028]	<a href="#">Comina, Comina</a> [3.4*] [0.018]	<a href="#">Pulp Fiction</a> [4.6*] [0.044]	<a href="#">Sabrina</a> [3.7*] [0.023]
Interest Group 16, *3.6*	Interest Group 17, *3.6*	Interest Group 18, *3.5*	Interest Group 19, *3.5*	Interest Group 20, *3.5*
<a href="#">Outbreak</a> [3.8*] [0.041]	<a href="#">Three Colors: Red</a> [4.2*] [0.033]	<a href="#">Mary Poppins (1964...</a> [4.1*] [0.028]	<a href="#">Leaving Las Vegas</a> [3.8*] [0.095]	<a href="#">Casper</a> [3.5*] [0.025]
<a href="#">Waterworld</a> [3.1*] [0.04]	<a href="#">Three Colors: Blue</a> [4*] [0.032]	<a href="#">Cinderella (1950)</a> [4*] [0.024]	<a href="#">Dead Man Walking</a> [4.1*] [0.078]	<a href="#">First Knight</a> [4*] [0.023]
<a href="#">Braveheart</a> [4.9*] [0.038]	<a href="#">Three Colors: Whit...</a> [3.9*] [0.028]	<a href="#">The Sound of Music...</a> [4.2*] [0.022]	<a href="#">The Birdcage</a> [3.3*] [0.076]	<a href="#">Junior</a> [3.5*] [0.021]
<a href="#">The Net</a> [3.4*] [0.038]	<a href="#">Barcelona</a> [3.4*] [0.02]	<a href="#">Dumbo (1941)</a> [3.9*] [0.022]	<a href="#">Sense and Sensibil...</a> [3.9*] [0.074]	<a href="#">The Santa Clause</a> [3.5*] [0.019]

Fig. 1. Interest groups number 1-20.

Interest Group 21, *3.5*	Interest Group 22, *3.4*	Interest Group 23, *3.1*	Interest Group 24, *3.1*	Interest Group 25, *3*
<a href="#">Smoke</a> [3.7*] [0.024]	<a href="#">Money Train</a> [3.2*] [0.017]	<a href="#">Jumanji</a> [3.5*] [0.026]	<a href="#">Jurassic Park</a> [3.2*] [0.08]	<a href="#">Jerry Maguire</a> [3.5*] [0.021]
<a href="#">Mighty Aphrodite</a> [3.3*] [0.018]	<a href="#">Terminal Velocity</a> [3.2*] [0.016]	<a href="#">Oliver and Company...</a> [3.1*] [0.023]	<a href="#">Forrest Gump</a> [3.5*] [0.079]	<a href="#">The People vs. Lar...</a> [3.5*] [0.016]
<a href="#">Fargo</a> [4.1*] [0.018]	<a href="#">The Shadow</a> [3.3*] [0.015]	<a href="#">Muppet Treasure Is...</a> [3.1*] [0.022]	<a href="#">The Fugitive</a> [3.6*] [0.066]	<a href="#">Mars Attacks!</a> [2.6*] [0.015]
<a href="#">Four Weddings and ...</a> [3.3*] [0.018]	<a href="#">Hard Target</a> [3.4*] [0.015]	<a href="#">Powder</a> [3.5*] [0.021]	<a href="#">Terminator 2: Judg...</a> [3.6*] [0.056]	<a href="#">The Lost World: Ju...</a> [2.8*] [0.015]
Interest Group 26, *2.8*	Interest Group 27, *2.7*	Interest Group 28, *2.7*	Interest Group 29, *2.6*	Interest Group 30, *2.2*
<a href="#">Dances With Wolves</a> [3.4*] [0.11]	<a href="#">Naked Gun 33 1/3: ...</a> [2.4*] [0.02]	<a href="#">I.Q.</a> [2.8*] [0.022]	<a href="#">Mrs. Doubtfire</a> [3*] [0.049]	<a href="#">Mulholland Falls</a> [2.1*] [0.018]
<a href="#">Batman (1989)</a> [2.4*] [0.11]	<a href="#">The Hudsoner Prox...</a> [3.3*] [0.019]	<a href="#">French Kiss</a> [2.9*] [0.019]	<a href="#">Pretty Woman</a> [2.9*] [0.047]	<a href="#">The Arrival (Shock...</a> [2.5*] [0.017]
<a href="#">Pulp Fiction</a> [3.4*] [0.1]	<a href="#">Hot Shots! Part De...</a> [2.3*] [0.019]	<a href="#">Nine Months</a> [2.5*] [0.019]	<a href="#">Ghost</a> [3*] [0.042]	<a href="#">Primal Fear</a> [2.9*] [0.017]
<a href="#">Apollo 13</a> [3.6*] [0.094]	<a href="#">So I Married an Ax...</a> [2.6*] [0.016]	<a href="#">Junior</a> [2.4*] [0.018]	<a href="#">The Mask</a> [2.5*] [0.042]	<a href="#">City Hall</a> [2.4*] [0.016]
Interest Group 31, *2.2*	Interest Group 32, *2*	Interest Group 33, *1.8*	Interest Group 34, *1.8*	Interest Group 35, *1.7*
<a href="#">E.T.: The Extrater... [2.6*] [0.01]</a>	<a href="#">Lord of Illusions</a> [1.6*] [0.011]	<a href="#">Sleepless in Seatt...</a> [1.8*] [0.017]	<a href="#">Toy Story</a> [2.4*] [0.05]	<a href="#">Striptease</a> [0.025*] [0.033]
<a href="#">The Sound of Music...</a> [2.3*] [0.0086]	<a href="#">Tales From the Hoo...</a> [1.6*] [0.0087]	<a href="#">The Firm</a> [1.8*] [0.015]	<a href="#">Mission: Impossibl...</a> [1.8*] [0.049]	<a href="#">Independence Day (...</a> [0.87*] [0.029]
<a href="#">Top Gun (1986)</a> [2.3*] [0.0086]	<a href="#">Mallrats</a> [2.4*] [0.0083]	<a href="#">Pretty Woman</a> [1.5*] [0.015]	<a href="#">Independence Day (...</a> [2.1*] [0.048]	<a href="#">The Cable Guy</a> [0.16*] [0.028]
<a href="#">Mary Poppins (1964...</a> [2.3*] [0.0083]	<a href="#">Wes Craven's New N...</a> [2*] [0.0082]	<a href="#">Dave</a> [2*] [0.015]	<a href="#">Twister</a> [1.8*] [0.043]	<a href="#">Barb Wire</a> [4.9e-005*] [0.025]
Interest Group 36, *1.1*	Interest Group 37, *0.68*	Interest Group 38, *0.39*	Interest Group 39, *0.16*	Interest Group 40, *0.16*
<a href="#">Super Mario Bros.</a> [0.11*] [0.017]	<a href="#">Mighty Morphin Pow...</a> [0.017*] [0.033]	<a href="#">Dumb and Dumber</a> [0.0025*] [0.038]	<a href="#">Kazaam</a> [0.028*] [0.014]	<a href="#">Tales From the Hoo...</a> [0.022*] [0.0075]
<a href="#">The Beverly Hillbi...</a> [0.34*] [0.016]	<a href="#">The Brady Bunch Mo...</a> [0.28*] [0.024]	<a href="#">Ace Ventura: Pet D...</a> [0.016*] [0.034]	<a href="#">Children of the Co...</a> [0.021*] [0.014]	<a href="#">Vampire in Brookly...</a> [1.3e-005*] [0.007]
<a href="#">Richie Rich</a> [0.22*] [0.015]	<a href="#">Mortal Kombat</a> [0.21*] [0.018]	<a href="#">Ace Ventura: When ...</a> [0.00067*] [0.033]	<a href="#">A Very Brady Seque...</a> [0.083*] [0.012]	<a href="#">The Baby-Sitters C...</a> [0.0063*] [0.007]
<a href="#">The Next Karate Ki...</a> [0.21*] [0.014]	<a href="#">The Bridges of Mad...</a> [0.015*] [0.018]	<a href="#">Waterworld</a> [0.034*] [0.028]	<a href="#">Halloween: The Cur...</a> [0.035*] [0.012]	<a href="#">Candyman: Farewell...</a> [0.0039*] [0.0065]

Fig. 2. Interest groups number 21-40.