

to calculate Euclidean distance between an object j and cluster k ; $MaxEc(j,k)$ means the maximum Euclidean distance between an object j and cluster k .

2.2 Similarity computation

After grouping the items and users, we get a new item group rating matrix IG and a new user group rating matrix UG . We can use the item-based collaborative algorithm (IC) to calculate the similarity by using Pearson correlation-based similarity that is the most common measure for calculating the similarity. To make the correlation computation accurate we must first isolate the co-rated cases (i.e., cases where the users rated both i and j)

$$sim(i, j) = \frac{\sum_{u \in U} (R_{u,i} - \bar{R}_i)(R_{u,j} - \bar{R}_j)}{\sqrt{\sum_{u \in U} (R_{u,i} - \bar{R}_i)^2} \sqrt{\sum_{u \in U} (R_{u,j} - \bar{R}_j)^2}} \quad (2)$$

where $R_{u,i}$, $R_{u,j}$ is the rating given to item i and j by user u ; \bar{R}_i , \bar{R}_j is the mean rating given by user; and U is the total number of items.

After completion of calculating the item-based collaborative similarity, we get $sim_{IG}(i,j)$ the similarity of the new item group rating matrix. Thus we apply to calculate the user-based collaborative algorithm (UC). In UC, clustering is based on the attributes of user profiles and clustering result is treated as items. We apply Equation 2 to calculate $sim_{UG}(i,j)$ the similarity of new user group rating matrix, i and j mean the user and u mean the item, instead the original meaning.

In additional, we use Pearson correlation-based to calculate the similarity from item-rating matrix and user-rating matrix. We call $sim_I(i,j)$ and $sim_U(i,j)$ respectively. At last, the total user similarity is linear combination between $sim_{IG}(i,j)$ and $sim_I(i,j)$ Another is the total item similarity is linear combination between $sim_{UG}(i,j)$ and $sim_U(i,j)$.

$$simI(i, j) = sim_I(i, j) \times (1 - c) + sim_{IG}(i, j) \times c$$

where $simI(i,j)$ means the similarity between item i and j ; c means the combination coefficient; $sim_I(i,j)$ means that the similarity between item i and j , which is calculated from item-rating matrix; $sim_{IG}(i,j)$ means that the similarity between item i and j , which is calculated from item group rating matrix. Then make a linear combination of the similarity in user group rating matrix and the user-rating matrix as the above.

2.3 Collaborative prediction

The final step of offline mining process is to make a collaborative prediction. Prediction for an item is then computed by performing a weighted average of deviations from the neighbor's mean. The general formula for a prediction on item i of user u as follows:

$$P_{u,i} = \bar{R}_i + \frac{\sum_{k=1}^n (R_{u,k} - \bar{R}_k) \times simU(i, k)}{\sum_{k=1}^n |simU(i, k)|}$$

where $P_{u,i}$ represents the prediction for the user u on item i ; n means the total neighbors of item i ; $R_{u,k}$

means the user u rating on item i ; \bar{R}_i is the average ratings on item i ; $simU(i,k)$ means the similarity

between item i and its neighbor k ; \bar{R}_k means the average ratings on item k .

The score that the user u rated the item j can be predicted as $P_{u,j}$:

$$P_{u,j} = \bar{R}_j + \frac{\sum_{k=1}^n (R_{u,k} - \bar{R}_k) \times simI(j, k)}{\sum_{k=1}^n |simI(j, k)|}$$

3. Conclusions

We have proposed a new approach to combine collaborative and content-based filtering techniques. In offline mining, the collaborative sub-system provides two types of clusters, item-based and user-based, to overcome a shortage of ratings. Besides, in our approach, based on the information from user group rating matrix and item group rating matrix, we can make predictions for the new item and new user. Thus we can solve the cold start problem, sparsity, and scalability.

4. References

- [1] B.M. Sarwar, G. Karypis, J.A. Konstan, and J. Riedl, "Item-based Collaborative Filtering Recommendation Algorithms," *In Proc. 10th Int. WWW Conf.*, 2001, pp.285-295.
- [2] A. Jameson, J. Konstan, and J. Riedl, *AI techniques for personalized recommendation. Tutorial notes*, AAAI-02, Edmonton, Canada, 2002.
- [3] M. Claypool, A. Gokhale, T. Miranda, P. Murnikov, D. Netes, and M. Sartin, "Combining Content-Based and Collaborative Filters in and Online Newspaper," *Proc. ACM SIGIR '99 Workshop Recommender Systems: Algorithms and Evaluation*, Aug. 1999.
- [4] J. Han and M. Kamber, *Data mining: Concepts and Techniques*, Morgan-Kaufman, New York, 2001
- [5] G. Salton and C. Buckley, "Term-weight approaches in automatic retrieval," *Information Processing and Management*, vol. 24, no. 5, 1988, pp.513-523.