

In a world where there is an information overload, it becomes imperative to have some kind of filtering in place to help users find relevant information. One of the ways to do is by profiling the user information, comparing the user profile with other profiles that have similar likes and dislikes and using it to recommend new information. Thus the user is able to find their preferred information from the internet more quickly and efficiently. These systems are called the "Recommender System". The main aim of these systems are to "Predict" what the user may like. For this to work, the system can either use the content or features of the item and match it with the user profile (Content based) or use the historical activities of the user to create a similarity matrix with other users or items and use that for prediction (Collaborative Filtering) or combine both approaches (hybrid). In this paper we will discuss more on the collaborative filtering approach.

Collaborative filtering has two main techniques: Memory based or Neighbourhood-based approach and Model-Based approach. A hybrid system can be built by combining both Neighbourhood-based and Model-based approach.

Memory or Neighbourhood based approach: The approach uses the fact that users with similarity in profiles will provide similar ratings to items or that similar items liked by a user will get similar ratings. The two types of Memory based approaches are

1. Item based collaborative approach
2. User based collaborative approach

Item-Based approach: In this the system will look at the ratings provided by a user and recommend items that have high ratings from users with similar interests. A good example to this would be, if you browse for a Midi-dress in an online apparel store, you will get a list of suggested items or "You may also like" list at the bottom or side that are midi-dresses viewed. The same behaviour can be noticed in Amazon or Goodreads as well. When you rate book "10", you will be recommended books with similar ratings.

Some advantages of using item-Based approach are:

1. The recommendation is a lot more consistent as it is based on what items were selected by the user and what is similar to the item selected and its rating.
2. Item based approach are much easier to explain as they are predicted based on the user's item selection and rating.
3. The predicted values are a lot more reliable.

User-Based approach: In this approach the system does not care about the ratings but instead clusters the users with similar tastes together. For example in YouTube Music, if the user listened to "Jazz" then the systems pairs the user with other users listening to jazz and the recommendation is based on what other users in the cluster are listening to. The relevance of the results will depend on how many choices the user makes.

Some advantages of using User-Based approach are:

1. The prediction and performance of the approach increases as the neighbourhood size increases.
2. There is a greater diversity to the recommendations.
3. Easier to implement than Item-Based approach

Model Based CF: In this approach models are developed using model based CF algorithms such as Bayesian CF, clustering CF, latent semantic CF, and CF using dimensionality reduction techniques. The dimensionality techniques are more popular in Model based approaches as it in general can be used to improve the robustness and accuracy of the memory based models. One big advantage of using this approach is that instead of having high dimensional matrix with sparse data, we will have smaller matrix there by using reduced space. This makes it highly scalable and easy to compute when combined with user based or item based approach.

The challenges of Collaborative filtering approach: Recommender systems are necessary to help in filtering the vast data available now and has shown that it can significantly increase a company's sales output. But, there are few main challenges that need to be addressed to ensure that they perform well.

1. **Data Sparsity and Cold Start:** The data produced by the web is huge but the data that are rated by users (to show preference) is very low. This makes building the user-item matrix difficult as the matrix will be sparse and thus leading to very weak or incorrect recommendation. This also occurs when there is a new user or Item added to the system (Cold start).

One way to overcome this can be if recommender system can start by displaying the latest trend. E.g. Amazon in its homepage displays the top 100 popular books or items. We can also request the user to provide their preferences as part of a short questionnaire. Last, option would be wait until the user searches for at least one item and then build recommendations based on all the groups/clusters the item belongs to. E.g. If the user searches for "Titanic", we can place the user in clusters such as "Romantic Movies", "Kate Winslet", "Leonardo Di Caprio", "James Cameron" etc.. And recommend the top trending ones in these clusters.

There were few techniques to deal with the cold start and the earliest one was content-boosted CF algorithm as a hybrid CF. A model based CF algorithm was later developed called TAN-ELR. PIP (Proximity-Impact-Popularity) was a method that used a heuristic measure for collaborative filtering. PIP did address issues with the cosine method and PCC. Another method called MJD (Mean-Jaccard-Difference) was developed. It combined three measures which were mean squared difference with the ratings, Jaccard similarity measure and the difference measure of ratings between two users. All the measures obtain their weight based on neural network learning. PSS (Proximity-Significance-singularity) was developed to overcome the shortcomings of PIP. It uses the proximity, significance and singularity aspects of the user ratings. By combining the global preference with the local context information the cold start problem can be overcome significantly.

2. **Scalability:** As the number of users and items grow the systems performance can get affected by using the traditional CF approach. Most of these systems require quick recommendations for all the online users and this means the computation engine has to perform at its optimal all the time.

There were few techniques used to overcome the scalability issue. C. Sneha, G. Varma in 2015 used Hadoop to implement user based CF algorithm which helps to address the issue of scalability. Dimensionality reduction technique can help with the scalability problem and combining that with incremental system as suggested by B. Sarwar, G. Karypis, J. Konstan, J.

Riedl, Fifth International Conference on Computer and Information Science. (2002) can help reduce the cost and complexity of the matrix factorization step. . C. Zeng, C.-X. Xing, L.-Z. Zhou, Proceedings of the 12th international conference on World Wide Web. (2003) proposed using matrix conversion method to convert the matrix from the user-item matrix to user-class matrix for similarity measure and an instance selection method to remove irrelevant instance and reduce the size of the training set. This way both scalability and sparsity issues were addressed. M. Papagelis, I. Rousidis, D. Plexousakis, E. Theoharopoulos, International Symposium on Methodologies for Intelligent Systems. (2005) suggested a technique that uses incremental update for user-user similarities to overcome scalability issue.

Moreover Model based CF proposed to save CPU and memory which makes it better performance and scalability wise by using Bayesian CF, clustering CF, latent semantic CF, and CF using dimensionality reduction techniques. But, the models as such can be very complex and difficult to maintain.

Conclusion:

There is a lot of scope for development in the improving the scalability, sparsity and cold start issues that the CF approach faces. Even though model based approach provides avenues to overcome these challenges, it has its own sets of issues being very complex, costly and difficult to maintain. On the other hand Memory based approach is easy to implement but allows for all the challenges mentioned above to pose a roadblock. I have not researched a hybrid system yet but can see advantages of using the best features from both approaches to overcome the challenges and provide more faster and accurate results.

References:

<https://medium.com/recommendation-systems/user-based-vs-item-based-collaborative-filtering-d40bb49c7060>

https://en.wikipedia.org/wiki/Collaborative_filtering

<https://www.iteratorshq.com/blog/an-introduction-recommender-systems-9-easy-examples/>

<https://goldberg.berkeley.edu/pubs/eigentaste.pdf>

https://en.wikipedia.org/wiki/Multi-armed_bandit

<https://www.be-terna.com/insights/recommendation-systems-in-e-commerce-whats-the-thing-youve-never-known-but-always-wanted-to>

https://www.researchgate.net/publication/320761149_Collaborative_Filtering_Recommender_System_Overview_and_Challenges?enrichId=rgreq-1f99186304522ff14f4adfafa968464-XXX&enrichSource=Y292ZXJQYWdlOzMyMDc2MTE0OTtBUzo4NTEwODM0MjYwODA3NjhAMTU3OTk0OTM1MTI2Mg%3D%3D&el=1_x_3&esc=publicationCoverPdf

