**Music Recommendation System**

**Ceng History X - Iteration Report 2**

**April 12, 2014**

**Team Members:**

**1)** Asena Ok, 1746296, asenaok@gmail.com

**2)** Ayşe Aybüke Taşdirek, 1746353 aybuketasdirek@gmail.com

**3)** Birant Altınel, 1745744 birantaltinel@gmail.com

**4)** Hacer Nihal Tarkan, 1744291, tarkan.nihal@gmail.com

**Tasks for Iteration 2:**

| Planning Tasks | Responsible Person | Implementation | Is completed? |
|---|---|---|---|
| Developing Content Based Recommendation | All Team | No | This method is considered as inapplicable for our project |
| Improvement of recommendation algorithm with Neighborhood-Based Collaborative | All Team | Yes | Yes |
| Defining Evaluation Module Metrics | All Team | Yes | Yes |

**Table 1: Tasks for second iteration**

## 1. Why not Content-based?

In the first iteration, we've decided that tag based recommendation is not applicable for our data. Therefore, we analyzed the 'genre' of the songs to use while improving content-based approach. However, we just realized that 'genre' of the songs is not true. Please see the example below to observe the 'genres' are as irrelevant as we mentioned:

| ID | PROVIDERID | GENREID | CREATETIME | NAME |
|---|---|---|---|---|
| 3123440 | 7 | 22 | 2012-03-15 16:51:58.0 | Rolling In The Deep |
| 3383460 | 2 | 22 | 2013-09-02 21:06:15.0 | Ya Ali |
| 3383461 | 2 | 22 | 2013-09-02 21:06:15.0 | Erenler Cemi |
| 3383462 | 2 | 22 | 2013-09-02 21:06:16.0 | Kerbela'ya Yürüyorum |

**Table 2: Song in the same genre**

Also, nearly %10 of the songs has no genre information:



**Figure 1: Songs with null genre information**

**2. Neighborhood-Based (or Memory-Based) Filtering**

We also did some research on recommendation methods, and gave some reports on our researches. As we stated in our project proposal, neighbor-based filtering is a kind of collaborative filtering used in recommendation algorithms. In this method, to make a prediction, memory-based collaborative filtering algorithms use the whole or a sample of the user-item database. Users are grouped as people with similar interests. In these algorithms, identifying the neighbors of a new user is needed. Neighbors are the users who have the maximum number of similar interests with new user. After identification of neighbors, a prediction of recommendation for user can be generated. [1]

The neighborhood-based collaborative filtering algorithm is the most known memory-based collaborative filtering algorithm. This algorithm uses the following steps:

   I.    Calculating the similarity or weight between two users or two items is the first step and it is very important in memory based collaborative filtering algorithms. There are different methods to calculate the similarity. For example, in correlation-based similarity method the formula is given below.

$$w_{u,v} = \frac{\sum_{i \in I}(r_{u,i} - \bar{r}_u)(r_{v,i} - \bar{r}_v)}{\sqrt{\sum_{i \in I}(r_{u,i} - \bar{r}_u)^2} \sqrt{\sum_{i \in I}(r_{v,i} - \bar{r}_v)^2}}$$

Here $w_{u,v}$ is the similarity between two users u and v. On the other hand, it can be $w_{i,j}$ for the similarity between two items i and j. This method is also called Pearson correlation. $i \in I$ summations are over the items that both users u and v have rated. $\bar{r}_u$ is the average rating of commonly rated items of user u. Item based similarity calculation is similar to user based one. [1]

   II.   Producing a prediction for the active user. In the neighborhood-based collaborative filtering algorithm, weighed aggregate of the subset of nearest neighbors of the active user is used to produce a prediction.

$$P_{a,i} = \bar{r}_a + \frac{\sum_{u \in U}(r_{u,i} - \bar{r}_u) \cdot w_{a,u}}{\sum_{u \in U}|w_{a,u}|}$$

Here, $\bar{r}_a$ and $\bar{r}_u$ are the average ratings of user a and user u respectively. $w_{a,u}$ is the similarity between user a and user u. $u \in U$ summations are over the users who have rated on item i. [1]

III.   Generating top-N recommendations is step is to recommend N top-ranked items. There are user-based and item-based top-N recommendation algorithms. [1]

## 3. Evaluation Module

To calculate the accuracy of our recommendation algorithm, in light of our advisor's suggestion, we decided to use one of the most known metrics for recommendation engine evaluation: "precision metric". [2]

$$Precision = \frac{TP}{TP + FP}$$

Here in this formula, TP is the true positive value, which refers to the number of songs the user listened from the recommendation list of us. In below, the false positive value is, contrarily, the number of songs which the user has not listened from the recommendation list we provide to him/her.

In other words, we calculate the efficiency of our recommendation with looking at the ratio of the songs according to their listening states.

By referring the formula provided, we've acquired the below pseudo code:

```
float CalculatePrecision(){
    song[] userHistory;
    song[] recommendations;
    song[] intersection;
    recommendations <- getRecommendation();
    intersection <- intersect(userHistory, recommendations);
    float precision = intersection.length / recommendations.length;
    return precision;
}
```

**Figure 2: Pseudo code of precision calculation for evaluation module**
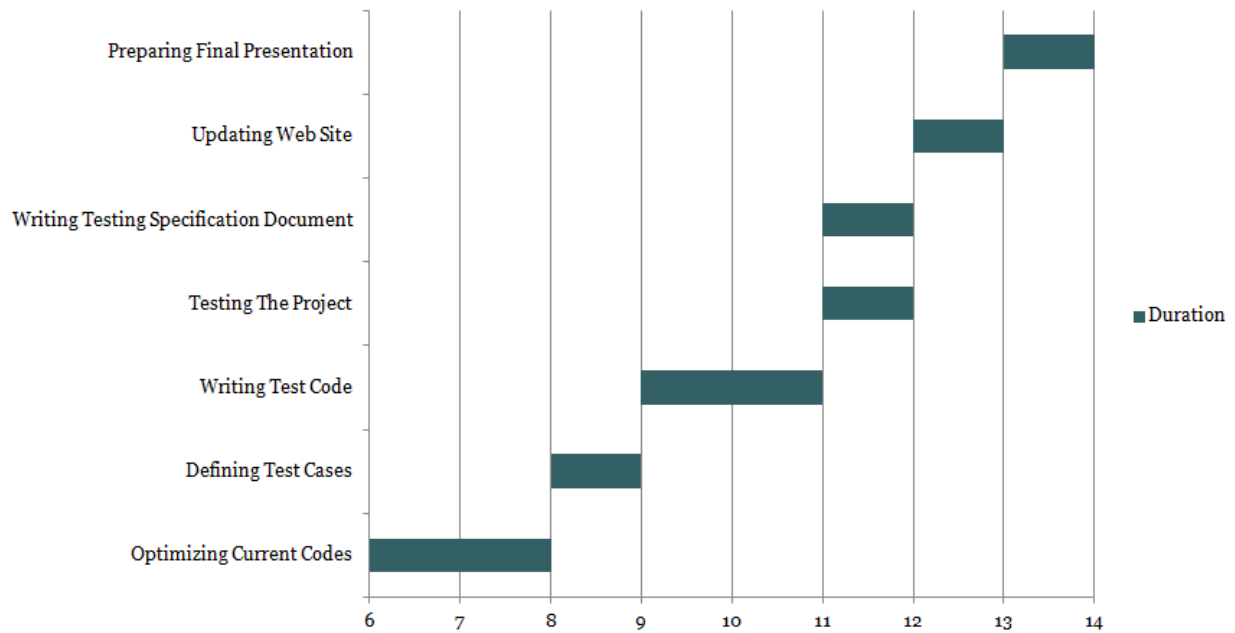
**Updated Gantt Chart for Our Future Works:**



**Figure 3: Gantt Chart of our future works**

**Finishing Evaluation Module:**

We will finish implementing our Evaluation Module.

**Solving the Inaccurate Rating Information In the DB:**

The rating values that are in the logs result in inaccurate calculations during the recommendation. For example it is very common in the logs that some users have the same rating value for all of the songs they have listened. Since this result in mathematically impossible calculations, we will either try to improve the rating information in the DB, or alter our algorithm to prevent erroneous calculations.

**Final Improvements on Collaborative Filtering:**

We will do our final touch on the recommendation algorithm and try to optimize our code for faster computation.

**Offline Similarity Calculations:**

Since our neighborhood-based collaborative filtering needs frequent user similarity values during computation, the system calculates hundreds, sometimes thousands of user similarity values. We will solve this by calculating the user similarity values between users offline and keep them in the database, so that the online part of the system will perform the tasks a lot faster.

**Optimizing DB structure:**

Since our recommendation and evaluation algorithms will be finalized in this iteration, we may update our database structure for a more efficient system.

**References:**

[1] Hong, Jun.(August 2009). A Survey of Collaborative Filtering Techniques. Advances in
        Artificial Inteligence. Retrieved from:
        http://www.hindawi.com/journals/aai/2009/421425/
[2] Ricci, Francesco. Rokach, Lior. Shapira, Brancha. Kantor, Paul. (2010). *Recommender
        Systems Handbook.* New York.