

Matching Users With Common Attributes

Deniz Cakiroglu

1 Introduction

This paper describes an algorithm for matching users that have given set of attributes. It matches the users based on how many attributes they have in common. The attributes can be weighted. This algorithm was developed for matching users by their most listened artists and songs retrieved from Spotify. The implementation of this algorithm can be found [here](#).

2 Queuing Preference Changes

2.1 Algorithm of Queuing Changes

Prefs are received in a list for each user. For each element in the list, that pref's existing weight for given user is checked. If that pref exists and the weight is different, that pref is queued. If it doesn't exist, the pref is also queued.

Assume i , p , w stands for user id, pref id and weight respectively.

$A = \{ \text{most recent user prefs} \}$

with A_i representing the preference list of a given user

with A_{i_p} representing the specific preference weight in the list of a given user

This list is received externally, and used to populate Q , which after is used to update P .

$Q = \{ \text{prefs to be updated in the algorithm database} \}$

Each element in this set is a list containing $\{i, p, w\}$

$P = \{ \text{existing prefs for all users in database} \}$

with P_p representing the user list for a given preference

with P_{p_i} representing the weight of this preference for specific user

The operation to populate Q can be defined as:

$$(\forall_{i,p} \in (A \cup P))(A_{i_p} \cap \bar{P}_{p_i}) = Q$$

2.2 Implementation of Queuing Changes

Implementation of this operation can be found in the following function:

<https://github.com/DAB-Co/jam-server/blob/master/utils/algorithmEntryPoint.js#L46>

Note that this function needs to be called externally for each preference.

2.3 Complexity of Queuing Changes

The count of users is $|A|$. For each user with id i , there is $|A_i|$ preference updates. $\forall_i \in A(\max(A_i))$ gives the maximum amount of pref updates for a given user among all users.

So the complexity is $O((\forall_i \in A(\max(A_i))) * |A|)$. We will refer this as $|Q|$ for simplicity.

3 Applying Queued Changes to the graph

To apply the changes, we first subtract the existing weight for each given queued change. Otherwise, we would be adding new weight of a changed pref on top of the old one. After, subtraction is done, we add the changed pref.

$G = \{ \text{Graph containing total weight of all common prefs between each user with common prefs} \}$

G_{ij} : total weight of common prefs of user i and j

Before we start subtracting, we need to ensure every element in Q also exists in P . $(\forall_q \in (\bar{P}_p \cap Q_{q,p}), \forall_p \in P)(P_{p_{q.i}} : q.w)$

We add preference to the preference list, if it doesn't exist. This happens when a new preference that's not in the database is received.

The set of edges after subtraction

$(\forall_q \in Q, \forall_j \in (G_i \cap P_{q.i}), q.i \neq j)(G_{q.i_j} : G_{q.i_j} - P_{q.i}, G_{j_{q.i}} : G_{j_{q.i}} - P_{q.i})$

Re-adding the new weights

$(\forall_q \in Q, \forall_j \in (G_i \cap P_{q.i}), q.i \neq j)(G_{q.i_j} : G_{q.i_j} + q.e, G_{j_{q.i}} : G_{j_{q.i}} + q.w)$

3.1 Implementation of Applying Queued Changes to the graph

Implementation of this operation can be found in the following function:

<https://github.com/DAB-Co/jam-server/blob/master/utils/algorithmEntryPoint.js#L234>

3.2 Complexity of Applying Queued Changes

$|Q|$ is the length of the queue. Each element is applied to all its common members. The worst case here would be $|G|$. So the complexity is $O(|Q| * |G|)$.

One situation this happens is when every preference weight of every user is changed, but what is common remains the same.

4 Getting Matches

4.1 Algorithm of Getting Matches

4.1.1 Getting Biggest Neighbor

In theory, getting a match should be fairly simple.
 $(\forall_i \in G)(\max(G_i))$ gives j which is the match.

4.1.2 Previously matched neighbors

Assume we store already matched users in M and don't want duplicate matches.
 $M = \{ \text{previously matched users} \}$
 where M_i gives the list of users it has been matched
 where M_{i_j} exists, there has been a match between the two users

Then

$(\forall_i \in G \cup M)(\max(G_i \cap \overline{M}_i))$ to exclude the already matched users.

4.1.3 No unmatched neighbors

What if there are no matches left? In other words, what if a user has previously matched with all it's neighbors in the graph or has no neighbors? Or maybe all their neighbors have already matched with someone in this iteration of the algorithm. The neighbors only consist of other users that the given user has common attributes with.

$L = \{\text{users that couldn't match their neighbors, i.e. } G_i \cap \overline{M}_i = \emptyset\}$

For any given user in L , any other user that hasn't been matched with the given user must be non neighboring.

We match people in L with each other.

$\{\forall_{i,j} \in G(i \neq j, \nexists M_{i_j})\}$

$L_matched = \text{Set}()$

for $u1$ **in** L :

if $u1$ **in** $L_matched$:

continue

 found = False

for $u2$ **in** l :

if $u1 \neq u2$ **and** $u2$ **not in** $L_matched$ **and not** $M.u1.u2$:

$L_matched.add(u1)$

$L_matched.add(u2)$

```

        found = True
        M.u1.u2 = True
        M.u2.u1 = True
        break
    if not found:
        LL.add(u1)

```

4.1.4 No matches left this iteration

We can't guarantee everyone in L gets a match. It is possible that for a given user in L, all the other users have gotten a match in this iteration or given user has matched them already before.

$LL = \{ \text{user(s) that couldn't match with other users in L} \}$

We match the user(s) in LL with someone random among all users. Which means, it is possible for some users to get more than 1 match a day.

$\{\forall_i \in LL, \forall_j \in G, \exists M_{i_j}, \exists M_{j_i}\}$

4.1.5 Preventing too many matches after random matching

What if multiple users in LL end up matching the same user?

To prevent this, we hard code a match count.

We allow the user(s) in LL to match with their selected user only if the selected person had not been matched more than twice.

Why twice?

Because if only 1 user remains in LL with no matches due to all possible matches having matched once this iteration, it wouldn't be able to find a match.

This means, up to 2 matches a day are possible in certain situations for a few lucky users.

```

LL_matched = Set()
for u1 in LL:
    if u1 in LL_matched:
        continue

    for u2, _ in G:
        m1 = matched_today.has(u2) ? 1 : 0
        m2 = L_matched.has(u2) ? 1 : 0
        m3 = LL_matched.has(u2) ? 1 : 0
        if M.u1.u2 or (m1+m2+m3 > 2):
            continue

        M.u1.u2 = True
        M.u2.u1 = True
        LL_matched.add(u1)

```

`LL_matched.add(u2)`

4.2 Implementation of Getting Matches

The complete implementation of step 4 can be found in this function:

<https://github.com/DAB-Co/jam-server/blob/master/utils/algorithmEntryPoint.js#L55>

4.3 Complexity of Getting Matches

To get maximum match for a given node, we need to iterate all of it's neighbors value, after updating them in step 3.

Step 3: $(\forall_i \in G)(\max(G_i \cap \overline{M}_i))$

Maximum neighbors that a given user can have is $|G| - 1$.

This step is $O(|G|^2)$.

Worst case for L would be if all the users couldn't find a match, as explained above.

Worst case complexity of matching in L is $O(|L|^2) = O(|G|^2)$

Same goes for LL. If every user in L has already matched with each other, LL would contain all the existing users.

Worst case complexity of matching in LL is $O(|LL|^2) = O(|G|^2)$

The total worst case complexity comes to:

$3 * O(|G|^2) = O(|G|^2)$

5 Overall Complexity of the Algorithm

$O(|Q| + |Q| * |G_1| + |G|^2)$. G_1 represents the state of the graph before $|Q|$ is applied.

This comes to $|G|^2$.

6 Additional Notes

In the example implementation, constraints such as common languages and active users are also included. Adding any constraint set is pretty simple.

For example:

$LA = \{ \text{users with common languages} \}$

where LA_i is the set of other users that user i can speak with

where LA_{i_j} means user i and j understand each other.

So, $\forall_i(\max(G_i \cap \overline{M}_i))$ becomes $\forall_i(\max(G_i \cap \overline{M}_i \cap LA_i))$.

And this repeats in every step that involves matching. Any other constraint can be added this way.

Assuming infinite memory, this shouldn't affect the complexity since all sets can be stored with data structures that allow $O(1)$ lookup.

At the time this paper is written, the last added constraint was the active user set.