

Here is what I think.

I just want to try to normalize the global effect like what this article said.

(<http://blog.echen.me/2011/10/24/winning-the-netflix-prize-a-summary/> , the first part)

What we need do is having a preprocessing program for the whole data.

We need 3 statistics to get the new dataset. If I use  $x_{ij}$  to represent the rate for user  $i$  giving the movie  $j$ . (we only need the  $x_{ij}$  that have been given!) And we have totally  $n$  users and  $m$  movies. Then we need:

1.

$$\bar{x} = \frac{1}{\text{number of rate}} \sum_{i=1}^n \sum_{j=1}^m x_{ij}$$

2.

$$\bar{x}_i = \frac{1}{\text{number of moive user } i \text{ rated}} \sum_{j=1}^m x_{ij} \text{ for } \forall \text{ user } i$$

3.

$$\bar{x}_j = \frac{1}{\text{number of user who rated moive } j} \sum_{i=1}^n x_{ij} \text{ for } \forall \text{ movie } j$$

Then we want to make the dataset that just show the specific interaction between user  $i$  and movie  $j$ . we call this  $e_{ij}$

$$e_{ij} = x_{ij} + \bar{x} - \bar{x}_i - \bar{x}_j$$

So now we are focus on think about the similarity between  $e$ . The code will be the same, we just predict the new  $e_{ij}$ . And when we get  $e_{ij}$ , we can get  $x_{ij}$  back as,

$$x_{ij} = e_{ij} - \bar{x} + \bar{x}_i + \bar{x}_j$$