

In this project we prepare to do a user reliability score evaluate system. We have some problems need to discuss with you in data preparation stage.

**(Will it proper to use the formula below to do our data labelling work?)**

Step1 Data preparation

• Using yelp open source dataset(<https://www.yelp.com/dataset>) and do **data labelling** work.

**Explanation:** We have 16 attributes from our dataset as dictionary showing in appendix. Since we need to do score prediction but there's no score label in this dataset. We plan to calculate or labelling reliability score by using formula below:

$$Score = 1 - \frac{\sum_1^{RC} |RS_i - BAS|}{RC \cdot 4}$$

*RC* sum of one user's reviews

*RS* score the user marked to each merchant

*BAS* current score of each corresponding merchant

**This formula is aim to and more likely to calculate a average accuracy ratio.**

We can not just use the average score the user gave(average of RS) as the finally score.

Because reliability can not just be baed on user-side comments, we should take the bias into consideration. You can see this formula in a more abstract way:

**Score = 1-average percentage error= average accuracy= reliability score**

Example: here is list of a user who gave 5 comments on yelp:

<i>RS</i>	<i>BAS</i>
4.5	5
4	4.5
5	4
3	4.5
5	4.5

1.We will get a score like: score=1-(0.5+0.5+1+1.5+0.5)/5\*4=0.8

2.Take this score, multiply by 5( full score we can give or merchant can get on yelp), the average score this user gave to merchants is : 0.8 \*5=4.

3. the reason why we divide a constant number 4 in this formula is that the biggest deviation value we could get is (5-1) or(1-5) ,which absolute value of this deviation is 4.

## Appendix:

attribute dictionary

Attribute	Attribute Description
<i>user_id</i>	unique user id marked by 22 digits
<i>name</i>	first name of user
<i>review_count</i>	the total number of reviews they've written
<i>yelping_since</i>	years since user joined yelp
<i>friends</i>	list of user's friends marked by user_id
<i>useful</i>	number of useful votes sent by the user
<i>funny</i>	number of funny votes sent by the user
<i>cool</i>	number of cool votes sent by the user
<i>fans</i>	number of fans the user has
<i>elite</i>	the years the user was elite
<i>average_stars</i>	average rating of all reviews
<i>compliment_hot</i>	number of hot compliments received by the user
<i>compliment_more</i>	number of more compliments received by the user
<i>compliment_profile</i>	number of profile compliments received by the user
<i>compliment_cute</i>	number of cute compliments received by the user
<i>compliment_list</i>	number of list compliments received by the user