

# INTRODUCTION TO BIG DATA

## Phase III

Group 5

### Data Cleaning:

#### Run test.py

1. **Validity:** In the first step of data cleaning, we are checking in each table if the primary keys are being repeated or not. Primary keys across all the tables in our dataset are user\_id, business\_id, review\_id, category\_id. Our dataset data is represented in such a way that there can't be primary key duplicates in a table.
2. **Completeness:** We are checking if our dataset contains Null values in any table. This checks the completeness of the dataset.
3. **Accuracy:** Accuracy of the dataset is cross validated by checking the values of columns which have numerical values and checking if the value for each tuple is in the proper range as it should be. For example, the stars column values can range from 0 to max 5, it cannot go below 0 or above 5.
4. **Timeliness:** Timeliness is done by checking the date column of the dataset which shows that even the recent months data is contained in the dataset.

### Data Integration:

The following views were created:

1. CREATE MATERIALIZED VIEW business\_tips as select distinct b.name, t.tip from tip as t join review as r on t.business\_id = r.business\_id join business as b on b.business\_id = t.business\_id join users as u on t.user\_id = u.user\_id and t.user\_id = r.user\_id where t.complement\_count > 2 and r.stars > 2 and u.fans > 50

2. CREATE MATERIALIZED VIEW number\_of\_restaurants as select  
ad.postal\_code,ad.city, count(bc.business\_id) from category as c join  
business\_category as bc on c.id = bc.category\_id join address as ad on  
bc.business\_id = ad.business\_id where c.category\_type LIKE 'Restaurants'  
group by ad.postal\_code,ad.city, c.id order by count(bc.business\_id) DESC
3. CREATE MATERIALIZED VIEW Business\_name\_category AS SELECT  
bc.business\_id,name as Business\_name,Category\_type as Category FROM  
business as b JOIN business\_category as bc ON bc.business\_id = b.business\_id  
JOIN category as c ON c.id = bc.category\_id;
4. CREATE MATERIALIZED VIEW TopUsersTips AS SELECT b.name as business  
name,u.name as username,tip FROM users as u JOIN tip as t ON u.user\_id =  
t.user\_id JOIN business as b ON b.business\_id = t.business\_id WHERE  
review\_count>50 and fans=100;

View number 1 contains a business name and a tip for that business name. These printed tips are by users with compliment count greater than 2 and number of fans greater than 50. Also the ratings for a business is greater than 2.

View number 2 contains a postal code, the city name for that postal code and the number of business in that city with the category = 'Restaurants'

View number 3 contains all the business id's, their names and their category names.

View number 4 contains a business name, a user name and a tip by the user for that business.

## **ITEMSET MINING**

### **itemsetMining.py**

Using itemset mining on the yelp dataset we try to find out the number of users who have commonly given a tip to a business. For this the steps are as:

1. Build the best\_tips table from the tip table by pulling out user\_id, business\_id where the compliment count was greater than 0.
2. From the best\_tips we built the L1 lattice table by taking frequent itemsets of size 1 from best\_tips with a minimum support that the users in the itemsets have at least greater than 1 business in common.
3. Using this minimum support we find the lattice tables upto the nth level until the lattice level at the nth level is empty.

4. In the last step we take the  $n-1$  th lattice which will be non empty and join it with users table in order to display user names.

```
totale execution time 0.6700742244720459
users : Bella, Grace, MarVy, Sharr, Christalle, visited frequency :2
users : Mike, Grace, MarVy, Sharr, Christalle, visited frequency :2
users : Vi, Grace, MarVy, Sharr, Christalle, visited frequency :2
users : Mimi, Grace, MarVy, Sharr, Christalle, visited frequency :2
users : Terri, Christie, Joanna, Kevin, Amanda, visited frequency :2
users : Jordan, Grace, MarVy, Sharr, Christalle, visited frequency :2
users : Joe, Grace, MarVy, Sharr, Christalle, visited frequency :2
users : Jennifer, Grace, MarVy, Sharr, Christalle, visited frequency :2
```

**RDBMS OR Document Oriented for our Model:** Relational Databases works better for our dataset than the document-oriented databases as writing complex lattice queries in a Relational Database is much easier when compared to MongoDB. It makes sense as the lattice size keeps increasing, the queries get only bigger which is why RDBMS is better for our dataset. Also, writing code which generates queries dynamically can be done in RDBMS without any hassle. Document oriented would work better if the dataset is too large which is not the case with our dataset. Moreover, it is difficult to perform the cleaning process in MongoDB searching through each document for null values and since there exists no foreign key constraint in document oriented, it gets difficult to maintain consistent data which is required for itemset mining and data cleaning.