## Detecting Sentiment Label for amazon customer reviews

In the first step of my analysis, I have used an exising data base with course name dsci5350 and I have created a table with name as 'amazon_rp0477' with the column names and data types as shown below in the query. As there is a chance of displaying null values for the column headers, I have set the table properties to skip the header line . In the second query, I have ran show tables to see whether the table has been created or not.



Then I have loaded the table with the file that I saved on my local desktop using 'load data local inpath 'Desktop/Amazon.txt' into table amazon_rp0477'. Using this command I have loaded the data into the table.

Now, I'm checking whether the file that is sitting on desktop has been loaded into amazon_rp0477 table or not. From the below screenshot, it is evident that amazon.txt file is loaded into amazon_rp0477.



Below is the data that can be viewed when clicked on the amazon.txt file.

I have created another table for dictionary.txt file as dictionary_rp0477





The above screenshot displays all the data that is contained in the file or loaded from the dictionary file into the table which we created.

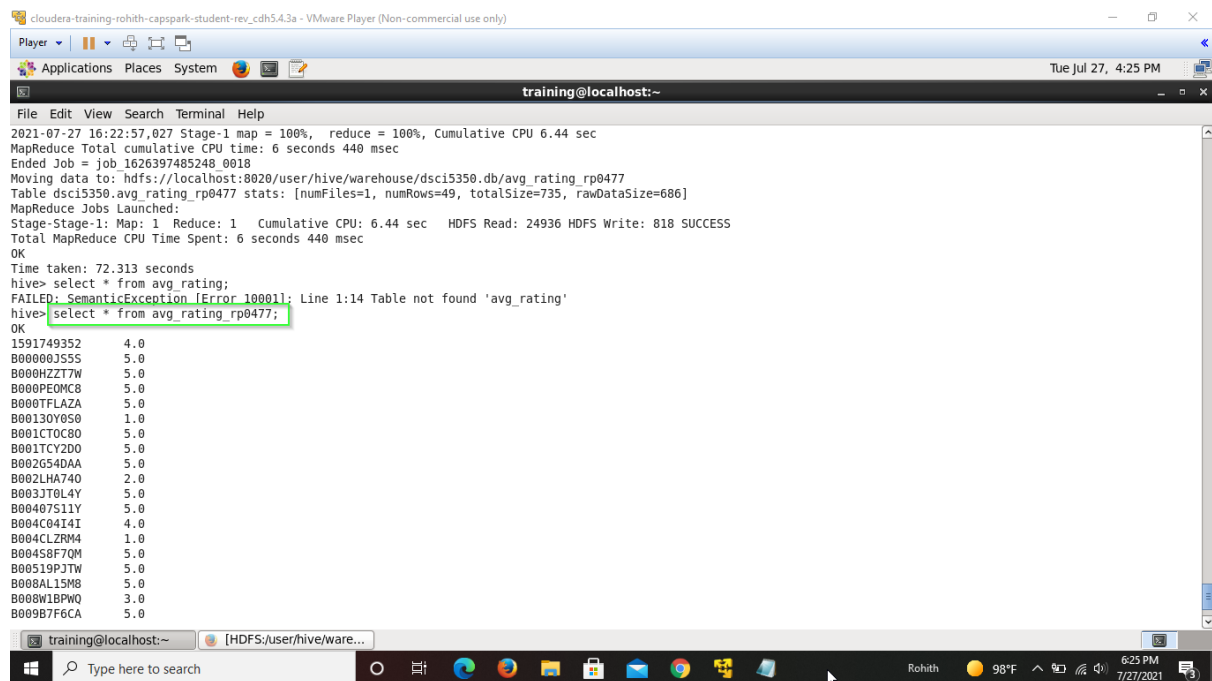2. A) Now, In order to find the products that has the average rating less than 2; we need to create a table that has average of the star rating column. So we created another table avg_rating_rp0477 with product id and their avg ratings.

The above query created a table for average reviews.



Here, we can see all the data i.e. with the product_id and avg_rating

In the above screenshot, I have ran a query to find the average rating that is less than 2 and could see that there are 5 products that have rating less than 2.

2. B) Here, we are asked to display the total number of products that are reviewed per day. I have ran a sql query to display and we have got 49 total number of products that are reviewed per day.



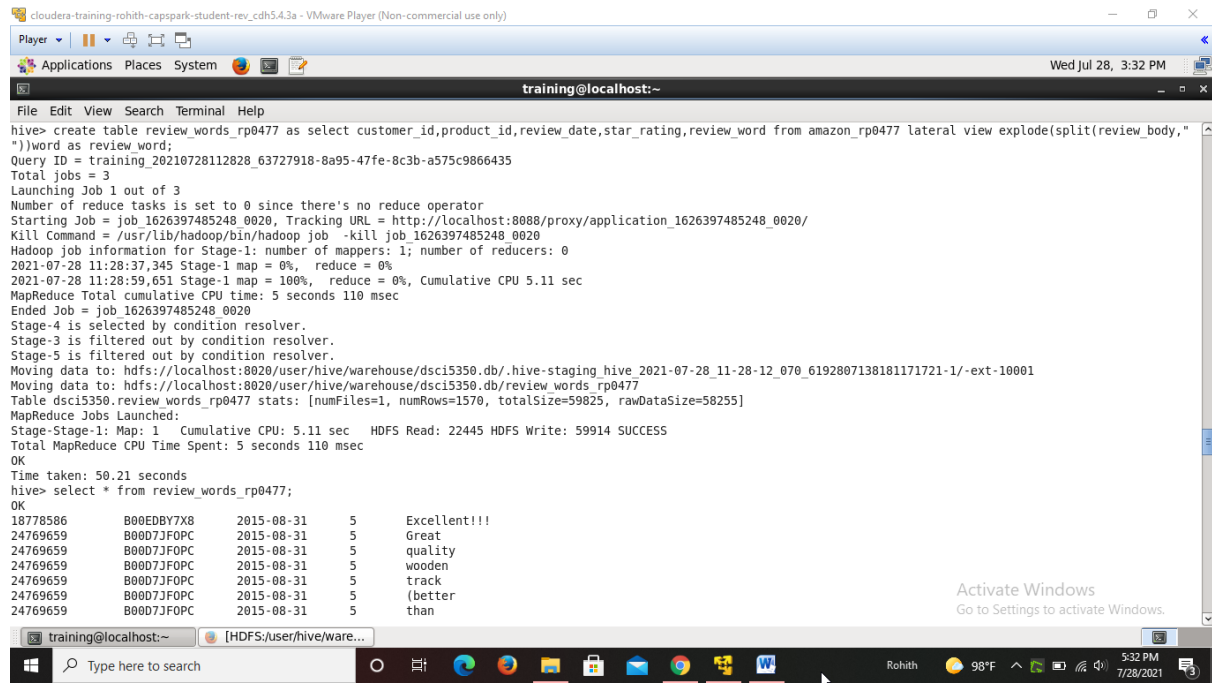2. C) Here we are asked to use sentiment analysis and display the top 5 good reviews amazon received for their products. In addition to this, we need to include customer_id, product_id, review_date, star_rating, total score and sentiment label in a new table.

Initially we need to have two new columns that as per the question. We need to find the total score and sentiment label using the rest of the columns. In order to find the total score, first we need to find the individual score of each word and then add them to sum up for the total score. Hence we can run a split query that basically splits all the words from a review body or sentence. Below is the query that is needed to run the query.
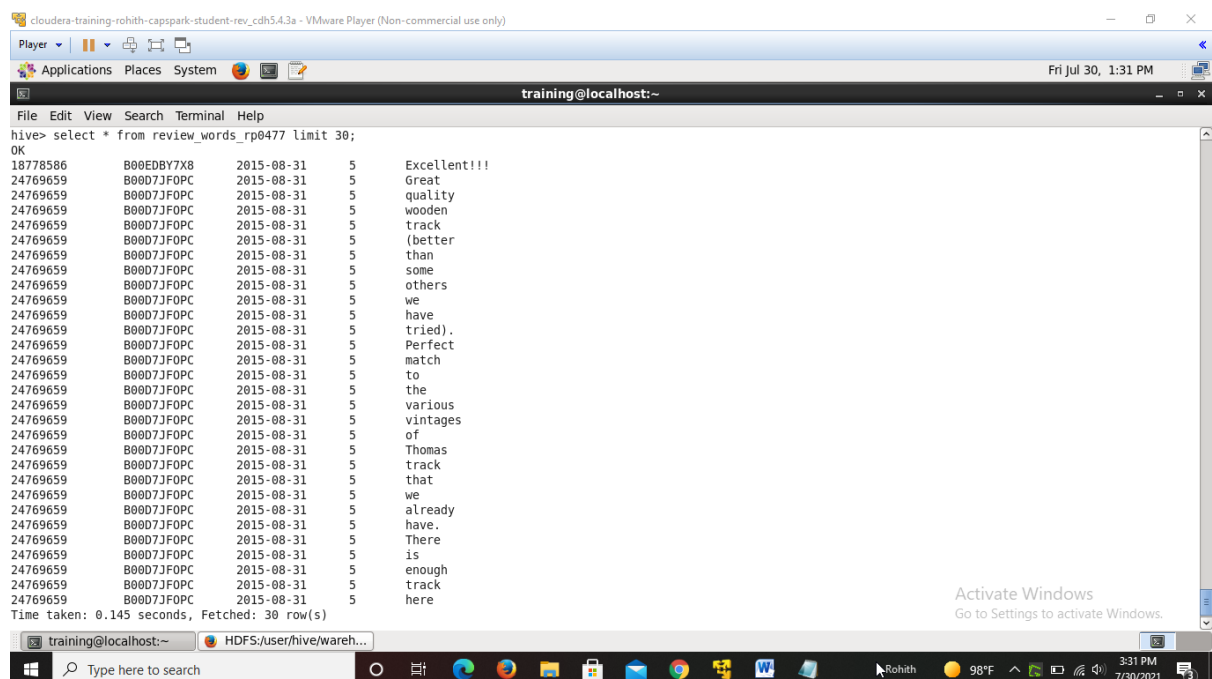


After creating a table using the query mentioned above, we display the data using select command and the above screenshot displays the data. In the below screenshot we can see the top 30 data as we have used limit 30

Now, we need to match the score of the review with the dictionary file and find the total score, hence I have created a table score_rp0477. In the below screenshot we can see the query that finds the total score for the review_words.





Now, we need to create a sentiment table that has a column total_score i.e. sum of all indiviual row score for the same product is calculated as total_score. Hence I have ran a query shown below and

output can also be seen in the below screenshot.



Now, we can see a column that is added beside the star rating i.e. total_score it received. Then in order to do sentiment analysis, I have ran a case statement using conditions if total_score > 0, display positive, total_score < 0 display negative else display neutral. In the below screenshot it is evident of running a query.



In the question it is asked to display the top 5 results for the good reviewed products.

In the above screenshot, we can see the results of top 5 good reviews that are received for the amazon products.

3. From the output of 2.c, we can infer that out of total reviews we have got only two negative reviews based on total score. Basically total score represents how negative or positive the total review is about. The drawback of this approach is whatever dictionary we tend to use has a chance of missing words that are present in the reviews. The words that are not present in the dictionary become useless and meaningless in the analysis. In the given dictionary file, we can see that same word is present with different verbs such as disregard, disregarded, disregarding, performance or accuracy might increase when for a single word for all variations instead of multiple words for a single variation. Though we add certain words based on our research and give a proper score to it, it works to certain extent but if the data is huge, we cannot read the data manually and assign a score to the words individually. So the best approach would be to have an individual dictionary type for different industries. For a retail industry like amazon, flipkart, alibaba, there should be certain different type of dictionary that helps to maximize the potential of the sentimental analysis. In addition to this we have got many speech processing techniques for analysis. The text can be converted to speech and processed for sentimental analysis to yield a better result.