# CA-675 || CLOUD TECHNOLOGY || ASSIGNMENT 1 – DATA ANALYSIS

| NAME | ANKIT SHARMA |
|---|---|
| STUDENT ID | 20211119 |
| EMAIL | ankit.sharma28@mail.dcu.ie |
| PROGRAMME OF STUDY | MSC. IN COMPUTING (DATA ANALYTICS) |
| GITHUB LINK | https://github.com/ankitapril/CA675-Assignment-1/tree/main |

**Task 1-Data Extraction**
We must acquire top 2,00,000 posts from stack exchange

**Task 2-Load Data**
Data is load with hive
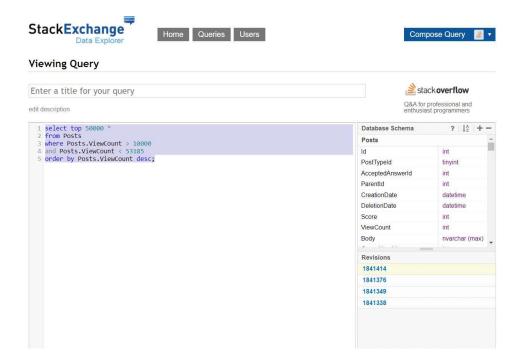
**Task 3-Query with Hive**

Top 10 posts by score? The top 10 users by post score? The number of distinct users, who used the word 'cloud' in one of their posts?

**Task 4- Calculate TF-IDF with Hive**

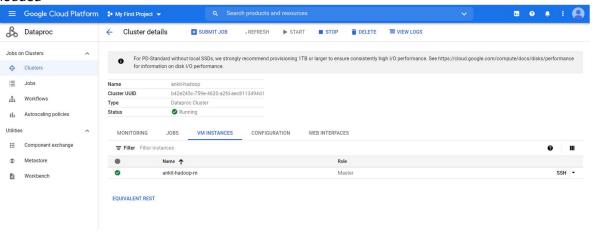Find Top 10 terms used for each of the top 10 users by post score

**Task 1-Extarcting Data**

- ➢ Extracting top 2,00,000 posts from Stack Exchange

- ➢ Data Acquisition from Stack Exchange Data Explorer (SEDE).

- ➢ We can download only 50,000 posts at one time. So, we will run 4 query to extract data

- ➢ Below are the queries
- • select top 50000 * from posts where Posts.ViewCount < 10030841 order by posts.ViewCount DESC;
- • select top 50000 * from Posts where Posts.ViewCount >62000  and Posts.ViewCount < 127042 order by Posts.ViewCount desc;
- • select top 50000 * from Posts where Posts.ViewCount >20000  and Posts.ViewCount < 74480 order by Posts.ViewCount desc;
- • select top 50000 * from Posts where Posts.ViewCount > 10000 and Posts.ViewCount < 53185 order by Posts.ViewCount desc;

**Task 2-Load Data**

For loading data. Cluster is created in GCP and cleaned csv file is uploaded and then data is loaded

**File Uploaded**



**Data loaded through Hive**



**Task 3- Query with Hive**

➢ Hive->create Database->ankit_hadoop->create table->hadoop_stack….(so on)

➢ 3.1-Top 10 Post by Score
- SELECT id, title, score FROM hadoop_stack ORDER BY score DESC LIMIT 10;



➢ 3.2-Top 10 Users by Post Score
- SELECT user_id, SUM(score) AS Tot_Sc FROM hadoop_stack GROUP BY user_id ORDER BY Tot_Sc DESC LIMIT 10

- SELECT COUNT (DISTINCT user_id) FROM hadoop_stack WHERE (LOWER(title) LIKE '%cloud%' OR LOWER(body) LIKE '%cloud%');



## Task 4-Calculate TF-IDF with Hive

Through Hive mall TF-IDF is calculated

Jar file is added



➢ Below are the queries to calculate TF-IDF

- create temporary macro max2(x INT, y INT) if(x>y,x,y);
- create temporary macro tfidf(tf FLOAT, df_t INT, n_docs INT) tf * (log(10, CAST(n_docs as FLOAT)/max2(1,df_t)) + 1.0);
- create table distOwnerIDs as SELECT user_id, SUM(score) AS TotalScore FROM hadoop_stack GROUP BY user_id ORDER BY TotalScore DESC LIMIT 10;
- create table mainUSRData as Select HT.user_id,title from hadoop_stack HT JOIN distOwnerIDs DO on HT.user_id = DO.user_id
- create or replace view mainUSRView as select user_id, eachword from mainUSRData LATERAL VIEW explode(tokenize(title, True)) t as eachword where not is_stopword(eachword);

- create or replace view tempView as select user_id, eachword, freq from (select user_id, tf(eachword) as word2freq from mainUSRView group by user_id) t LATERAL VIEW explode(word2freq) t2 as eachword, freq;
- create or replace view tfFinalView as select * from (select user_id, eachword, freq,rank() over (partition by user_id order by freq desc) as rn from tempView as t) as t where t.rn<=10 ;
- select * from tfFinalView;

```
hive> create or replace view tfFinalView as select * from (select user_id, eachword, freq,rank() over (partition by owneruser_id order by freq desc) as rn from tempView as t)
;
FAILED: SemanticException Failed to breakup Windowing invocations into Groups. At least 1 group must only depend on input columns. Also check for circular dependencies.
Underlying error: org.apache.hadoop.hive.ql.parse.SemanticException: Line 1:111 Invalid table alias or column reference 'owneruser_id': (possible column names are: user_id, eac
hive> create or replace view tfFinalView as select * from (select user_id, eachword, freq,rank()  over (partition by user_id order by freq desc) as rn from tempView as t) as t
OK
Time taken: 0.207 seconds
hive> select * from tfFinalView;
Query ID = ankitrich26_20211028101023_71c280ab-c1f3-41cc-860f-24262a5b56d4
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1635412257105_0005)

--------------------------------------------------------------------------------
        VERTICES      MODE      STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
--------------------------------------------------------------------------------
Map 1 .......... container    SUCCEEDED     1          1        0        0       0       0
Reducer 2 ...... container    SUCCEEDED     1          1        0        0       0       0
Reducer 3 ...... container    SUCCEEDED     1          1        0        0       0       0
--------------------------------------------------------------------------------
VERTICES: 03/03  [==========================>>] 100%  ELAPSED TIME: 6.68 s

OK
4883    python   0.038251366    1
4883    process  0.021857923    2
4883    ruby     0.021857923    2
4883    git      0.016393442    4
4883    table    0.016393442    4
4883    rails    0.016393442    4
4883    rename   0.010928961    7
4883    quotes   0.010928961    7
4883    branch   0.010928961    7
4883    list     0.010928961    7
4883    windows  0.010928961    7
4883    local    0.010928961    7
4883    write    0.010928961    7
4883    vs       0.010928961    7
4883    possible         0.010928961    7
4883    variable         0.010928961    7
4883    difference       0.010928961    7
4883    find     0.010928961    7
4883    style    0.010928961    7
6068    sql      0.026666667    1
```

```
6068    java     0.01777778     5
6068    make     0.017777778    5
9951    git      0.0390625      1
9951    python   0.0390625      1
9951    file     0.03125 3
9951    get      0.0234375      4
9951    javascript       0.0234375      4
9951    string   0.0234375      4
9951    interpreter      0.015625       7
9951    java     0.015625       7
9951    using    0.015625       7
9951    way      0.015625       7
9951    current  0.015625       7
9951    script   0.015625       7
9951    branch   0.015625       7
9951    android  0.015625       7
9951    dictionary       0.015625       7
49153   php      0.047058824    1
49153   using    0.04235294     2
49153   java     0.030588236    3
49153   javascript       0.030588236    3
49153   get      0.023529412    5
49153   array    0.016470589    6
49153   jquery   0.014117647    7
49153   file     0.014117647    7
49153   string   0.011764706    9
49153   class    0.009411764    10
51816   python   0.08423913     1
51816   wpf      0.024456521    2
51816   get      0.019021738    3
51816   string   0.016304348    4
51816   list     0.016304348    4
51816   class    0.013586956    6
51816   c        0.013586956    6
51816   function         0.013586956    6
51816   value    0.013586956    6
51816   index    0.010869565    10
51816   values   0.010869565    10
51816   vs       0.010869565    10
51816   use      0.010869565    10
63051   vs       0.022435898    1
63051   bash     0.022435898    1
63051   python   0.012820513    3
63051   java     0.009615385    4
63051   way      0.009615385    4
63051   file     0.009615385    4
63051   find     0.009615385    4
63051   get      0.009615385    4
```

**References**:

1.) https://www.tutorialspoint.com/hive/index.htm
2.) Stack Exchange