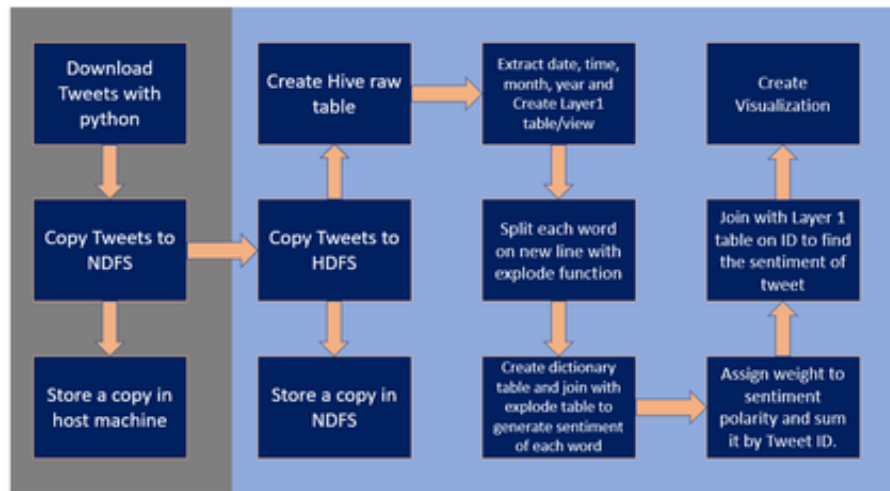# Project 01 - Twitter Sentiment Analysis in Hive

**This project will use Apache Hive to analyze the sentiment of incoming twitter data. The pipeline would be as follows.**



## Getting Data from Twitter API:

**The tweet generation code has been written in Python so it is necessary to install it on host machine.**
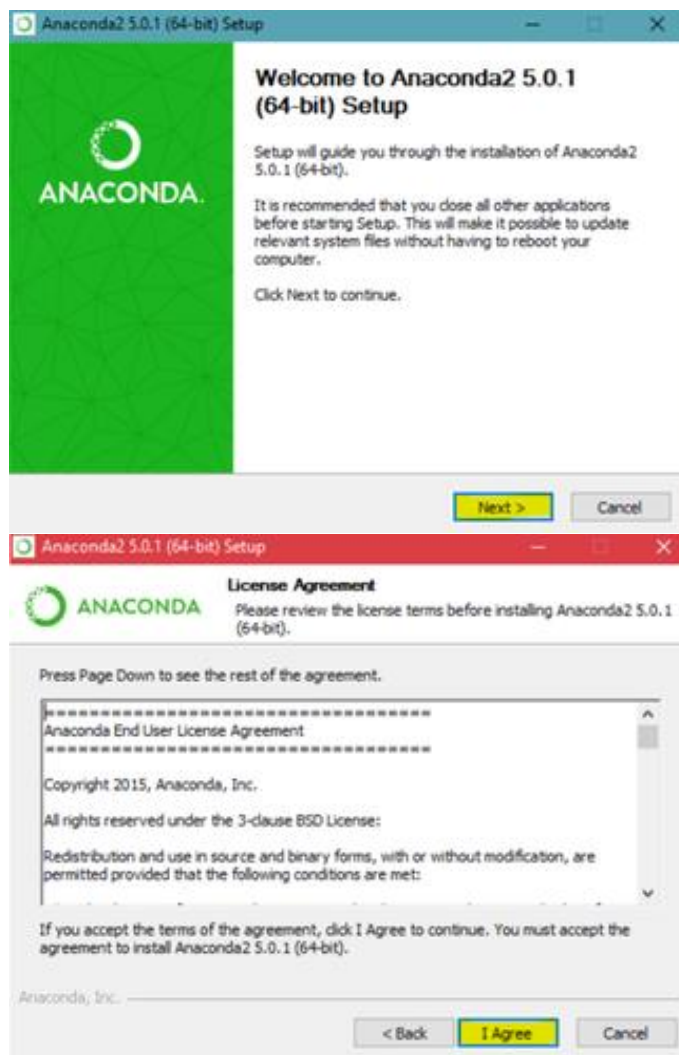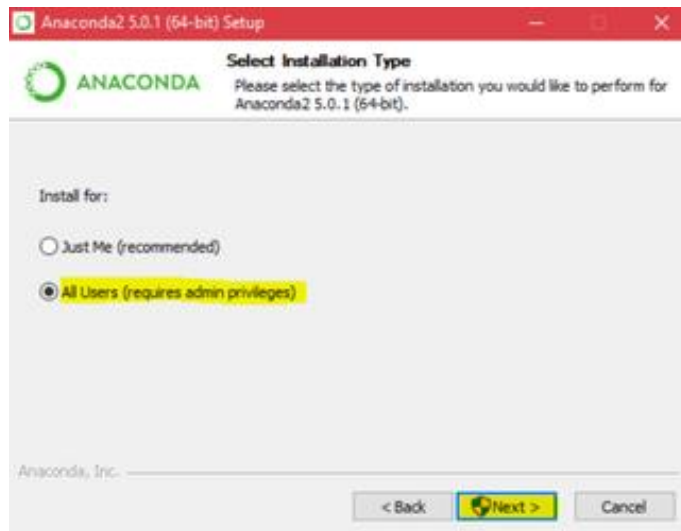
**Install Python (Anaconda):**

**The tweet generation code has been written in Python so it is necessary to install it on host machine. We'll install Anaconda which is a complete package for Python that provides a stable version of Python and a complete IDE and notebook for python coding. In addition, it also provides tools for debugging and testing making it a go-to application for python development. Here's the download link for Anaconda's website. Choose Python 3.8 installer according to your machine specs (32/64 bit).**
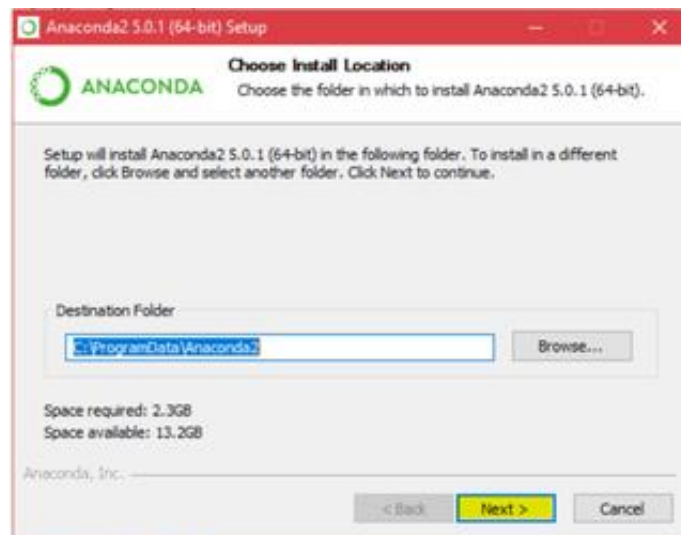**https://www.anaconda.com/products/individual**

**Anaconda Installers**

| Windows ⊞ | MacOS  | Linux △ |
|---|---|---|
| Python 3.8 | Python 3.8 | Python 3.8 |
| 64-Bit Graphical Installer (466 MB) | 64-Bit Graphical Installer (462 MB) | 64-Bit (x86) Installer (550 MB) |
| 32-Bit Graphical Installer (397 MB) | 64-Bit Command Line Installer (454 MB) | 64-Bit (Power8 and Power9) Installer (290 MB) |

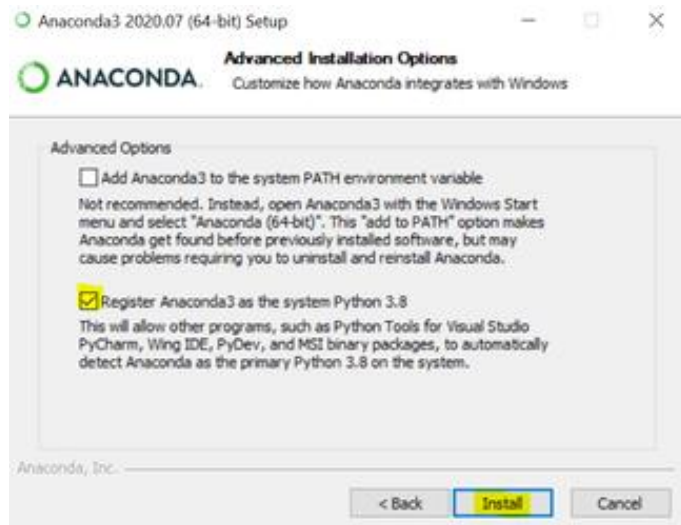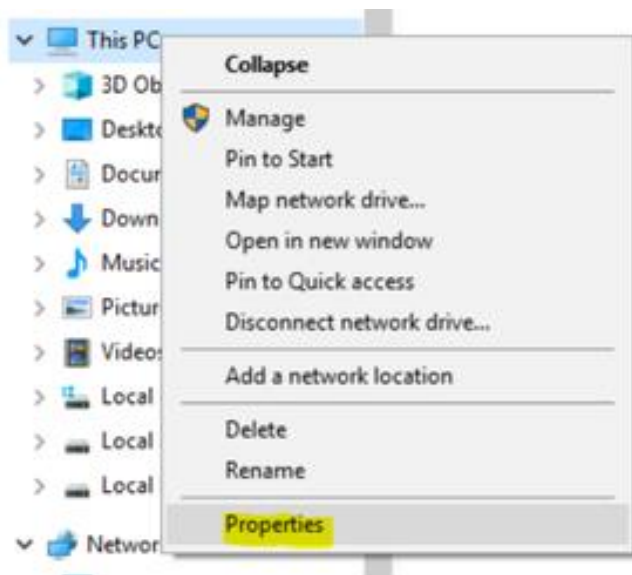**After downloading, open the installer click next and follow the process:**

**Select the folder you want to install Anaconda in and click next. Take note of the installation path entered in the Destination Folder as we will need it later on. Its best to copy the path in a notepad file for later use.**
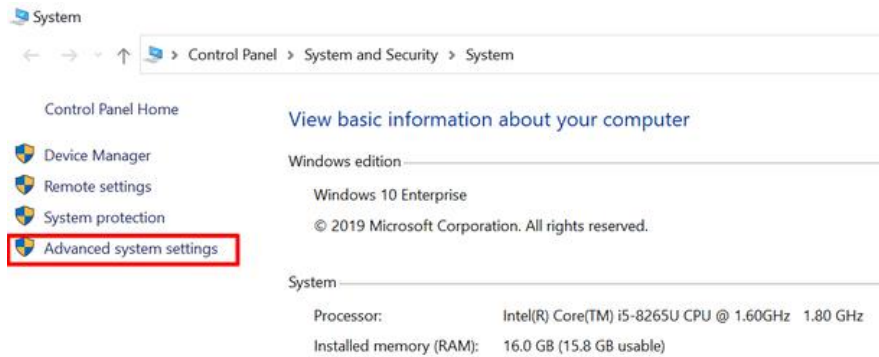


**In the Advanced options step for the make sure the Add Anaconda to my PATH environment variable is unchecked as below. This is an option but as the warning below it states, it can cause problems. Check the other option and click install.**
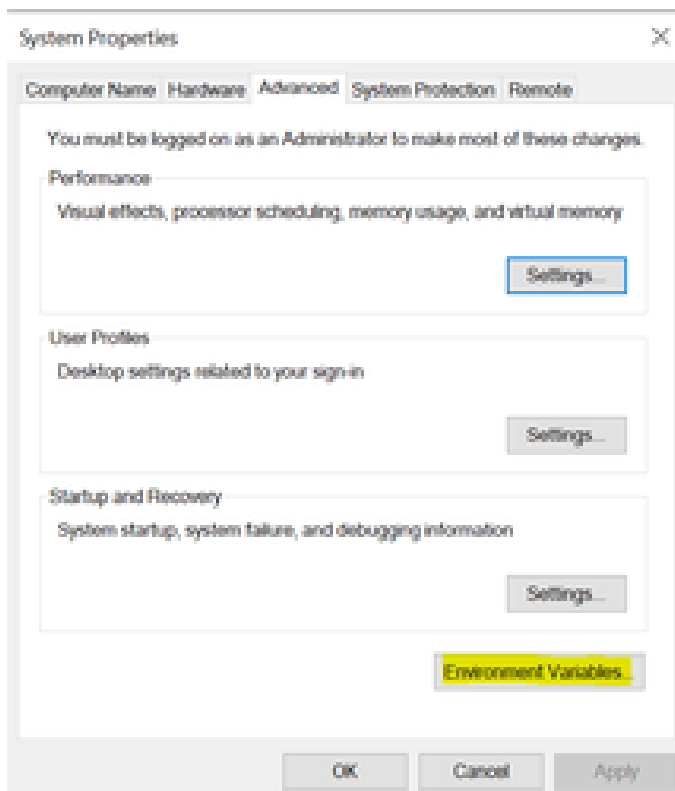
**The installation will be done after that. Now we need to add the path for Anaconda ourselves as we unchecked the option in the previous point. Open file explorer and on the left side and right click the computer name (usually named This PC but can be different). This will open the system's properties.**
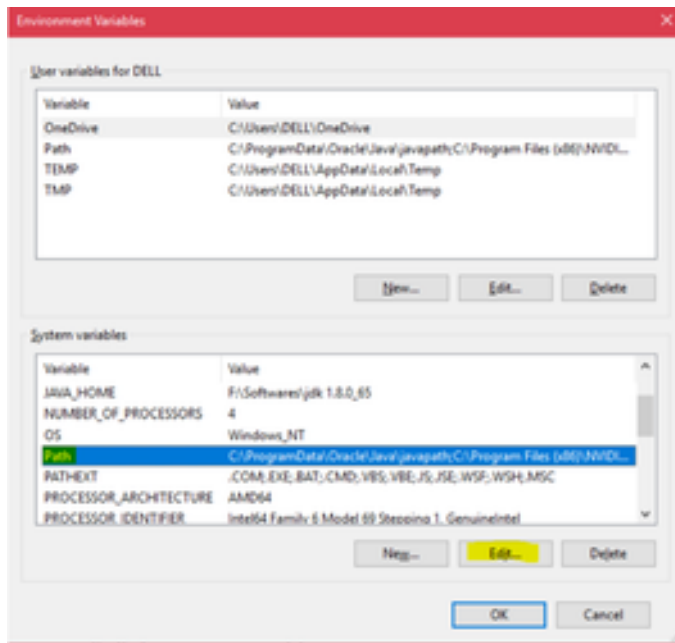


**Click on Advanced System Settings on the left hand side of the system properties window. A new window will open.**
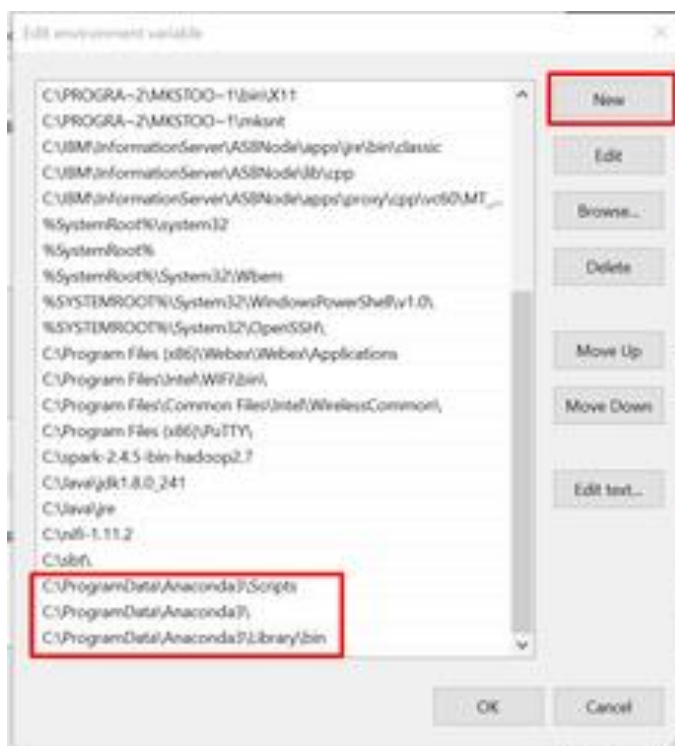
**Click on Environment Variables button. A new window will open.**



**Find the Path variable in the System variables tab, select it and click Edit. A new window will open.**

Click the New button and add these three paths one by one as shown below. Make sure that you DON'T delete previously set environment variables. Then click OK. Click OK again on the Environment Variables window and OK once again on the System Properties window. Now the path is set for Anaconda.
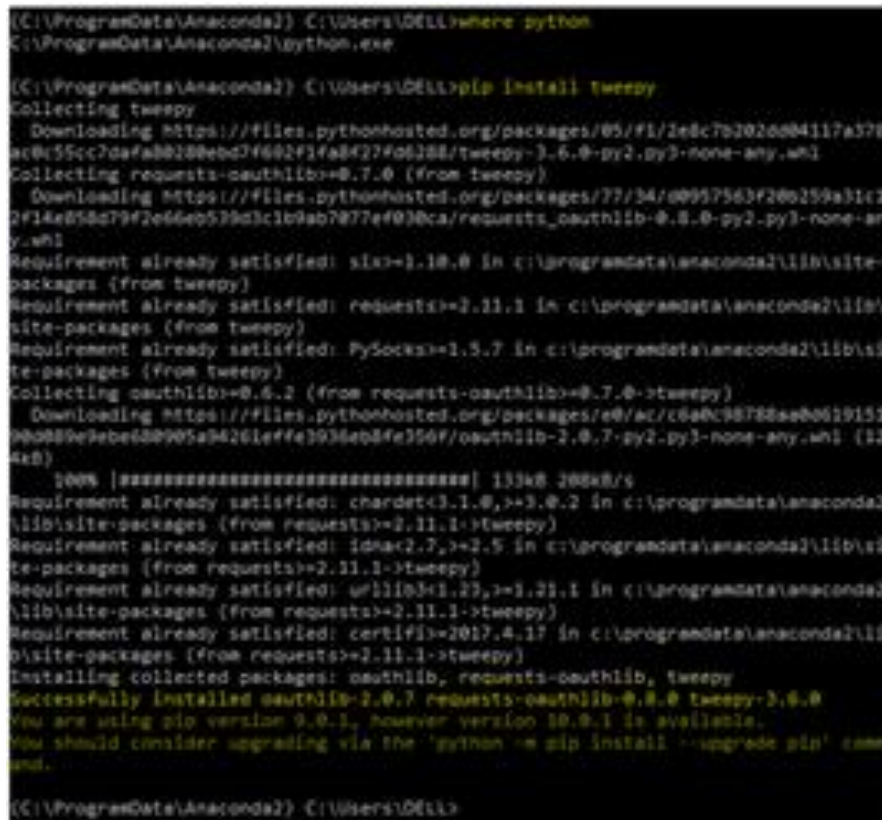
**Install Tweepy:**

**Next we need to install the tweepy library for Python. Tweepy provides a twitter API for Python and we can use it to fetch public tweets from Twitter.**
**Open the Anaconda Prompt. Enter the command *where python* to make sure python is installed correctly. The command should return the path of the python installation.**

```
where python
```

**Next enter the command *pip install tweepy* to install tweepy. Pip will automatically find and install tweepy.**

```
pip install tweepy
```



**Create .bat file to generate tweets:**

**The tweet fetching and transfer has to be done recursively. So, we'll create a .bat file for this as it will be helpful in automating the process later. This .bat file will be calling the .py file (will be shared with students).**

project_py_3.py

**Create a .bat file - just by creating a normal notepad file and saving it with .bat extension - in the folder where you placed your python file. <span style="color:orange">Please make sure not to put them on the desktop.</span> Open the batch file and paste the following code into it:**

```
cd /d "path-of-your-python-file"

"path-to-your-python-executable\python.exe" "path-to-pythonfile\name-of-python-file.py"

pscp -pw "your-vm-root-password" "path-to-where-your-json-filesare-created\*.json" root@<your-vm-ip-address>:/root/<folder-namewhere-tweets-are received> && move "path-to-where-your-json-filesare-created\*.json" "path-to-folder-where-json-files-are-to-bestored-after-processing"
```

**So, for example if my directory structure is as below:**

**Then my .bat file would be something like this (assuming I have placed the .py and .bat file in D:\BigData\Project\ )**

```
File Edit Format View Help
cd /d "D:\BigData\Project\"
C:\ProgramData\Anaconda3\python.exe "D:\BigData\Project\project_py_3.py"
pscp -pw "Myp@ss1234" "D:\BigData\Project\Received\*.json" root@192.168.10.124:/root/tweet_received && move "D:\BigData\Project\Received\*.json" "D:\BigData\Project\Processed\"
```

**You can double click the file to run it and it will start downloading tweets.**

# ETL Commands and Hive:

**Copy tweets to HDFS :**

**Running the .bat file will copy the tweets to you Linux (NDFS) as well - provided that your VM is running. Next you have to copy it on HDFS to analyze in Hive. For that we will use HDFS put command.**

```
hdfs dfs -put ./<folder-where-json-files-are-received-from-laptop>/*
/<destination-folder-in-hdfs> && mv ~<folder-where-json-files-arereceived-
from-laptop>/* ~/<folder-to-put-files-after-processing>
```

**for example if my NDFS and HDFS directories are as below:**

**Then my HDFS put command should be like:**



**Running the put command will copy all the tweets in tweet_received directory to HDFS' respective directory.**
**Next, We will create the staging table on top of this HDFS directory.**

**Creating Hive table:**

**Now we can create the hive table to load tweets in it. In the query below I have only selected some columns from the .json file to show that we can just load the columns we require later in the process. Use the query or modify it to create the table where tweets will be loaded.**

```
CREATE EXTERNAL TABLE IF NOT EXISTS raw_tweets
(
created_at string,
id string,
id_str string,
text string,
source string,
truncated string,
user_tw struct
```
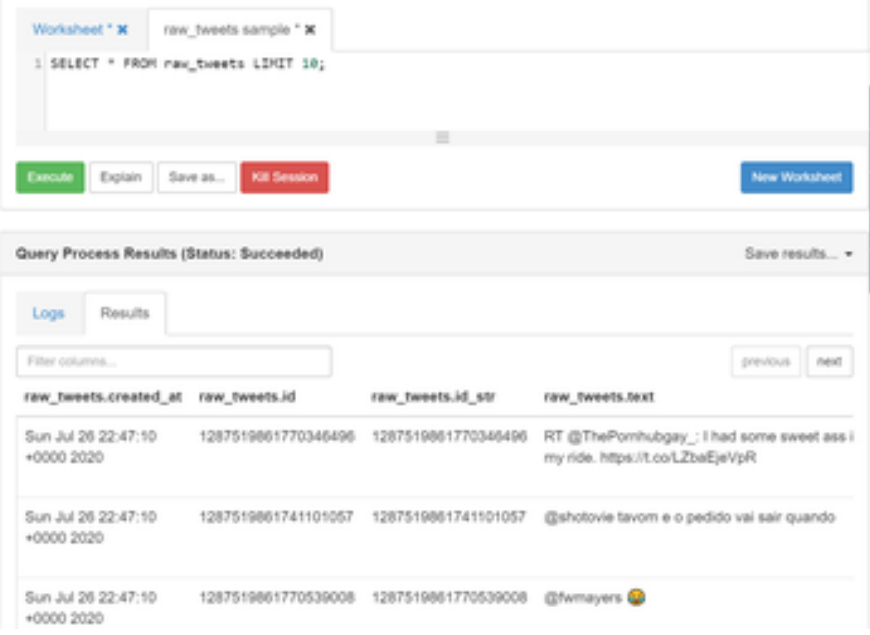
```
<
id:string,
id_str:string,
name:string,
screen_name:string,
location:string,
url:string,
description:string,
translator_type:string,
protected:string,
verified:string,
followers_count:string,
friends_count:string,
listed_count:string,
favourites_count:string,
created_at:string,
utc_offset:string,
time_zone:string
>
)
ROW FORMAT SERDE 'org.apache.hive.hcatalog.data.JsonSerDe'
STORED AS TEXTFILE
LOCATION '/apps/hive/warehouse/raw_tweets/' ;
```

**Executing this DDL will create the raw table on specified directory. Querying this table should show data as we have already copied it in previous step.**



**Create dictionary table :**

**Next we will create the dictionary table and for that we will use *dictionary.tsv* file. For that copy the dictionary.tsv file to HDFS. I have copied it to /tmp/dictionary/**

dictionary.tsv

## Create the table using the following DDL and load data into this table.

```
CREATE EXTERNAL TABLE my_dictionary (
    type string,
    length int,
    word string,
    pos string,
    stemmed string,
    polarity string
)
ROW FORMAT DELIMITED FIELDS TERMINATED BY '\t'
STORED AS TEXTFILE;

--Following will load data from HDFS directory into Hive table.
LOAD DATA INPATH '/tmp/dictionary/' OVERWRITE INTO TABLE
default.my_dictionary;
```

## Create L1, L2 and L3 views :

**To analyze the sentiment of a tweet, we need to break it down into words so that we can match the sentiment of words from dictionary table.**
**First, we will create Layer 1 view which will extract date elements, convert the text to lowercase and remove any new line characters.**

```
create or replace view default.Layer1 AS
select
created_at,
substr(created_at,27,4) as years,
substr(created_at,5,3) as months,
substr(created_at,9,2) as days,
substr(created_at,12,8) as times,
id,
lower(regexp_replace(text,'\n','')) as text
from default.raw_tweets;
```

| layer1.created_at | layer1.years | layer1.months | layer1.days | layer1.times | layer1.id | layer1.text |
|---|---|---|---|---|---|---|
| Sun Jul 26 20:54:05 +0000 2020 | 2020 | Jul | 26 | 20:54:05 | 1287491403384270848 | rt @umalambai @nokzen101 w this pain as a n phela this is ou https://t.co/qhh |
| Sun Jul 26 20:54:05 +0000 2020 | 2020 | Jul | 26 | 20:54:05 | 1287491403392655360 | rt @archaiospo πως θα χαρακτ την πεθερά σοι θηλυκό πετρετζ γιατί;- γιατί όλο μαγείρευε σκευ και όταν την αν |

**Layer 2 view will explode the tweet by ID and separate each word into new line.**

```
create or replace view default.Layer2 AS
select id, words
from default.Layer1
lateral view explode(split(text,'\\W+')) text as words
;
```

| layer2.id | layer2.words |
|---|---|
| 128749140338427O848 | rt |
| 128749140338427O848 | umalambanezn |
| 128749140338427O848 | nokzen101 |
| 128749140338427O848 | we |
| 128749140338427O848 | carry |
| 128749140338427O848 | this |
| 128749140338427O848 | pain |
| 128749140338427O848 | as |
| 128749140338427O848 | a |
| 128749140338427O848 | nation |

**Layer 3 view matches each word with the dictionary table picks whether polarity is negative or positive and assigns a value -1 or +1 respectively.**

```
Create or Replace view default.Layer3 AS
select
id, L2.words,
case d.polarity
       when 'negative' then -1
       when 'positive' then 1
       else 0 end
as polarity
from Layer2 L2 left outer join my_dictionary d
on L2.words=d.word
;
```

| layer3.id | layer3.words | layer3.polarity |
|---|---|---|
| 1287491403384270848 | rt | 0 |
| 1287491403384270848 | umalambanezn | 0 |
| 1287491403384270848 | nokzen101 | 0 |
| 1287491403384270848 | we | 0 |
| 1287491403384270848 | carry | 0 |
| 1287491403384270848 | this | 0 |
| 1287491403384270848 | pain | -1 |
| 1287491403384270848 | pain | -1 |
| 1287491403384270848 | pain | -1 |
| 1287491403384270848 | as | 0 |

**Finally, Layer 4 will sum up the polarity against a particular ID which generates the collective sentiment of the tweet with that ID. If you are comfortable with writing complex queries, you can also combine Layer 3 and Layer 4 views and just create one view.**

```
create or replace view default.sentiment as
select
  id,
  case
    when sum( polarity ) > 0 then 'positive'
    when sum( polarity ) < 0 then 'negative'
    else 'neutral' end as sentiment
 from layer3 l3 group by id;
```

| sentiment.id | sentiment.sentiment |
|---|---|
| 1285586508922023936 | positive |
| 1285586546079469569 | neutral |
| 1285586559463493634 | positive |
| 1285586579545821184 | neutral |
| 1285586585929551874 | neutral |

Now, We have the sentiment of each tweet in *default.sentiment* table/view grouped by each ID. If you wish to combine it with the Raw table or Layer 1 view/table to find the actual content of tweet. You can do so using the following query or create a final table using a CTAS statement.

```
SELECT
  L1.*, s.sentiment
FROM layer1 L1 LEFT OUTER JOIN sentiment s on L1.id = s.id
;
```
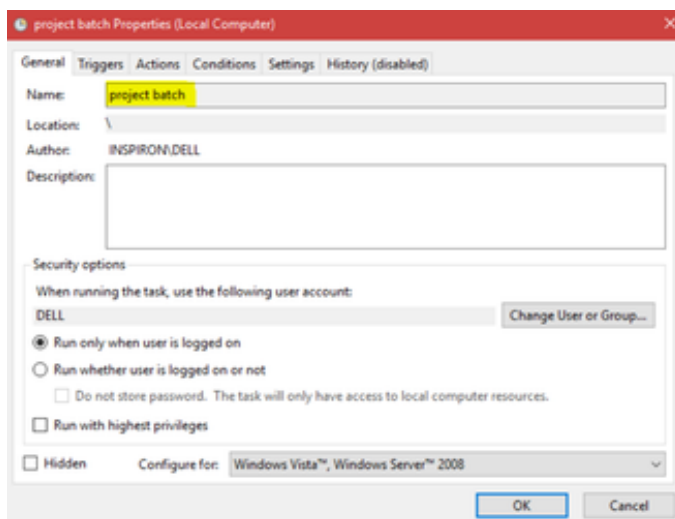
## Automation:

To automate this pipeline you just need to schedule the .bat file and HDFS put command. Rest will be done automatically as views will populate once underlying data does.

Schedule .bat file:

Open up Windows Task Scheduler. It is a default program of Windows and is already installed into your system. Click on Create Task on the right hand side of the window to create a new task:

**On the General tab enter the name you want to give to the task:**



**On the Triggers tab Click on New to create a new Trigger. Set scheduler to Daily, Repeat to one hour and indefinitely and check enabled as shown in the screenshot below. Click OK and the task will appear in the Triggers tab.**

**Go to the Actions tab and Click New. Choose action as Start a Program and give the path to the batch file in the field specified. Click Ok to save it.**



**On the Conditions tab uncheck everything. Click OK to create the Task**
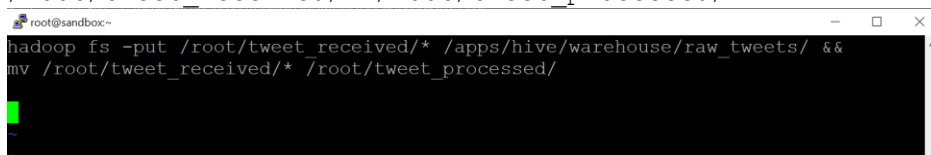
**According to the time you have set in the Triggers tab the task will activate and run the batch file. Make sure the IP address in the batch file is updated every time the VM is started as the IP address changes every time.**

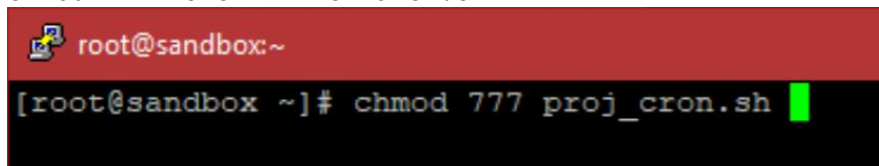**Schedule the HDFS file movement:**

**For the HDFS file movement we will use cron scheduler that is a built-in tool for Linux environment just like the windows scheduler.**
**First of all, create a .sh file using vi editor and paste the command that we used above to copy files from NDFS to HDFS and save this file. Give this files executable permissions as well.**

```
hadoop fs -put /root/tweet_received/* /apps/hive/warehouse/raw_tweets/ && mv
/root/tweet_received/* /root/tweet_processed/
```



```
chmod 777 <shell-file-name>.sh
```



**On the shell enter the command *crontab -e* to edit the cron list and add a new cronjob.**

**The text editor will open up. Enter I to get into insert mode and add the following line:**

```
10 * * * * /bin/bash -c "/<your-desired-file-name>.sh"
```



**The above line shows that your scheduled job will run at 10th minute of every hour.**
**Press Esc and enter :wq to save the cron job.**
**To view the installed cron job you can enter crontab -l in the shell.**

# Creating a visualization:

**Get creative and create a visualization by connecting Hive with Tableau using the ODBC connector. This task has been left to students so be creative.**