

# **CSEE5590/490: Big Data Programming**

## **Project Increment 2**

**Instructor - Zeenat Tariq**

**Team Number :- 6**

### **1. Project Title and Team Members**

International Student Data Analysis using Twitter API

- Tanvi Jain
- Saikumar Reddy Papagari
- Amarnadha Reddy Ankireddypalli
- Thotakura Naga Mounika

### **2. Introduction**

This project focuses on extraction of dataset (Twitter Data) from the Twitter API regarding all the information related to international students and using different tools like Hive, Hue pySpark, Solr to show different aspects of the data in a visualization form.

### **3. Background**

Twitter sentiment analysis is done in different ways before, using panda is the most common way. We will be using different hashtags which are related to international students and events related to international students to get a hold of better results based on those hashtags and more information about what events really took place and different aspects they

bring out. Here are some references of projects people have done implementing twitter analysis with big data, we are going to follow some aspects of these projects but our dataset is new and is extracted only for this project, so we could only find references related to those projects which have totally different dataset and completely different objectives.

1. <https://github.com/dsuarez993/bigdata-realtime-twitter-analysis>
2. <https://github.com/6vedant/TwitterAnalyticsHadoop>
3. <https://www.toptal.com/apache/apache-spark-streaming-twitter>

#### **4. Goals and Objectives:**

- **Motivation**

Studying abroad is a journey of education and discovery. There are currently over 1 million international students from more than 220 countries, coming to the United States annually. The individuals from this group belong to the international student community and we came to the United States to get a higher education. There are numerous situations where we are relied upon to follow those standards which the overall US citizens are not expected to follow in the long run since we have a place with the foreigner gathering. So, this gave us the establishment for this task and we chose to feature the fundamental information like percentage of students going to the United States for education, the probability of getting a work VISA, Immigration rules change

for F-1 Visa during COVID, Jobs for international students during covid(H1B sponsorship).

- **Significance**

Big data tools help to analyze the huge data which helps to provide efficient results. The sentimental analysis provides a brief understanding of various challenges that international students are currently facing and impact of covid-19 on the visa assurance.

- **Objectives**

The objective of this project is to get twitter data extraction using Twitter Data Analysis and then cleaning the data, performing sentimental analysis and importing files into Hadoop where we used different tools like Hive,Solr and pySpark.

- **Features**

The main feature of the project is to collect the Real timed tweets from the twitter API, also by performing the ETL which means we preprocess the data and extract the necessary data and then we load the extracted data in our HIVE.

Performed topic modelling using LDA and gensim model, and done visualization of the top 4 topics by t-SNE visualization.

Performed word cloud visualization using pyspark.

## 5. For all the features developed for this increment, write a documentation describing the design, implementation, testing, and deployment (including the precise descriptions and screenshots).

### •Dataset

We collected the data using twitter API using developer account and API keys. We have used different hashtags to get data related to 3 different scenarios as follows :-

#### 1. Generic International Student Data (**#F1visa**)

**#intlstudents #internationalstudents #studyinUSA)**

#### 2. Immigration rules change for F1 Visa during COVID in 2020 (**#AbolishICE**)

#### 3. Jobs for international students during covid (H1B sponsorship) (**#H1B, #h1bjobs**)

From the extraction we have the information about the tweet itself like this :-

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S
	Tweet Id	Text	Name	Screen Name	UTC	Created At	Favorites	Retweets	Language	Client	Tweet Ty	URLs	Hashtags	Mentions	Media Ty	Media UR			
1	1374879760	Science for the internation PRONINIB			2021-03-25T Thu Mar 25	0	0	en	<a href="#>Tweet	https://pronib	7	0	photo		https://pbs.twimg.com/media/ExSNpzrVoAMDixh.jpg				
2	1374877391	RT @marialuisapaul : Adriana Pérez adrianamper			2021-03-25T Thu Mar 25	0	0	en	<a href="#>Retweet	https://twitt	2	0							
3	1374876321	#internationalstudents	Maria Luisa I	marialuisaipa	2021-03-25T Thu Mar 25	1	1	en	<a href="#>Tweet	https://twitt	2	0							
4	1374876193	Today we had a special vis LEAP Intensi	LEAPWSU		2021-03-25T Thu Mar 25	0	2	0	en	<a href="#>Tweet		5	1 photo		https://pbs.twimg.com/media/ExSKDpWQAUbB-w.jpg				
5	1374872718	Education #schoolcounsel pigrantandas			2021-03-24T Wed Mar 24	0	0	und	<a href="#>Tweet		4	0 photo		https://pbs.twimg.com/media/ExS4yUWQAi9Nz.jpg					
6	1374863816	RT @ISClancer: Looking t	Crystal Kolrc	PROCrystal	2021-03-24T Wed Mar 24	0	0	en	<a href="#>Retweet	https://l.ead	5	0 photo		https://pbs.twimg.com/media/ExRFQdxXMagMEI.jpg					
7	1374856119	When it comes to choosing Lurnable.com lurnableedu			2021-03-24T Wed Mar 24	0	0	en	<a href="#>Tweet	https://bit.ly	7	0							
8	1374854640	Business classes are back! HouseToGroh housetogrow			2021-03-24T Wed Mar 24	0	0	en	<a href="#>Tweet	https://www	5	0 photo		https://pbs.twimg.com/media/ExR29XVEAQGzmA.jpg					
9	1374831723	RT @ISClancer: Looking t	PasadenCIT	PClancer	2021-03-24T Wed Mar 24	0	0	en	<a href="#>Retweet	https://l.ead	5	0 photo		https://pbs.twimg.com/media/ExRFQdxXMagMEI.jpg					
10	1374828718	Looking to find more infor ISC PCC			2021-03-24T Wed Mar 24	2	2	en	<a href="#>Tweet	https://l.ead	5	0 photo		https://pbs.twimg.com/media/ExRFQdxXMagMEI.jpg					
11	1374829002	@WVWise's free U.S. Adm	Dr. Kat Cohei draketchen		2021-03-24T Wed Mar 24	0	0	en	<a href="#>Tweet	https://go.st	5	0 photo		https://pbs.twimg.com/media/ExRFQdxXMagMEI.jpg					
12	1374824931	Now is the perfect time to StudentRoom StudentRoom			2021-03-24T Wed Mar 24	0	0	en	<a href="#>Tweet	http://ow.ly	2	1			https://pbs.twimg.com/media/ExRb9QjUAAxRfZ.jpg				
13	1374819046	SMARTS-UP, a mobility sdlopandrew	plopandrew		2021-03-24T Wed Mar 24	0	0	en	<a href="#>Tweet	http://plopai	14	0							
14	1374819787	RT @AlexFroger : Interna Laurence BO	LaurenceBoil		2021-03-24T Wed Mar 24	0	0	en	<a href="#>Retweet	https://lnkd.	5	0							
15	1374811992	Fulbright Scholarships in	plopandrew	plopandrew	2021-03-24T Wed Mar 24	0	0	en	<a href="#>Tweet	http://plopai	17	0							

The information about the user as well -

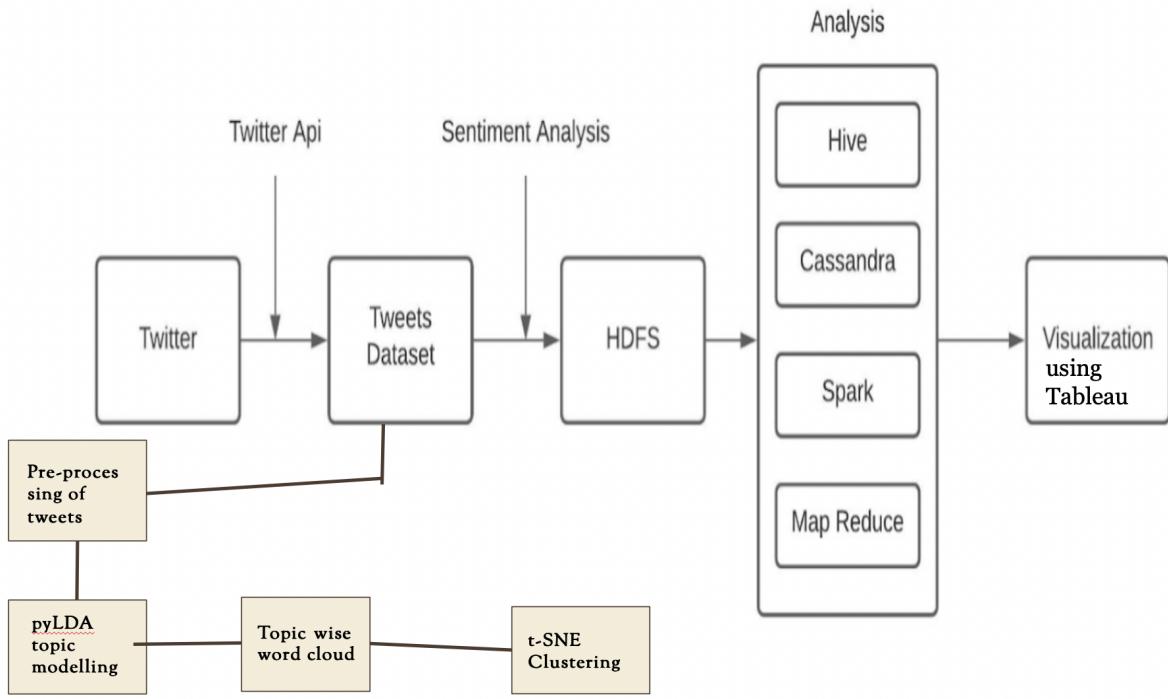
	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
1	User Id	Name	Screen Name/ UTC	Created At	Followers	Following	Favorites	Tweets	Lists	Bio	Location	URL	Verified	Default	Profile	
2	171719416	PRONIN Inte PRONINIB	2010-07-28T Wed Jul 28 0	173	315	210	349			1 We've been	Sydney, New http://www.	FALSE	FALSE			
3	1158836808	Adriana Pérez adrianamper	2019-08-06T Tue Aug 06 2	324	582	6690	1508			1 Guayaquil& Perla del Pac https://ndse		FALSE	TRUE			
4	10665229851	Maria Luisa I marialuisap	2018-11-25T Sun Nov 25 C	323	388	2156	1515			2 Talent & Incl Caracas, Ven https://mari		FALSE	FALSE			
5	1920828079	LEAP Intensi LEAPWSU	2013-09-30T Mon Sep 30	452	253	1784	3770			5 Wright State Dayton, OH	http://www.	FALSE	FALSE			
6	1227615420	plgrantandas plgrantandas	2020-02-12T Wed Feb 12	8	160	9	135			0 International	http://www.	FALSE	TRUE			
7	198577337	Crystal Kolrc IPROCystal	2010-10-04T Mon Oct 04 :	65	120	144	1521			2 Lies, damn li		FALSE	TRUE			
8	1077204432	Lurnable.com lurnableedu	2013-01-10T Thu Jan 10 1	39	14	14	793			3 Lurnable is tl	http://www.	FALSE	FALSE			
9	89557709501	HouseToGrow housetogrow	2017-08-10T Thu Aug 10 C	33	306	4	412			1 #Charity that Sydney, New http://house		FALSE	FALSE			
10	78736148	PasadenaCity PCClancer	2009-09-30T Wed Sep 30	6046	338	10010	10033			166 The official T Pasadena, C http://pasad		FALSE	FALSE			
11	386812395	ISC PCC	ISCLancer	2011-10-07T Fri Oct 07 23	59	9	34	524		0 The official T Pasadena, C http://www.	FALSE	TRUE				
12	1275139023	StudentRoom StudentRoom	2020-06-22T Mon Jun 22 :	6	16	4	65			0 Building the		FALSE	TRUE			
13	73586218181	Dr. Kat Cohei drkatcohen	2016-05-26T Thu May 26	11020	531	44957	21727			154 Founder and New York, N http://IvyWi		FALSE	FALSE			

It has all the different fields, which we will filter and use later using big data tools to perform queries and visualization.

These are some of the features and the description of those features :-

Feature	Description
id	Unique id of user
created at	Tweet created time stamp in UTC
text	UTF-8 text Tweet data
source	Type of device used to post the tweet
Name	Name of the user
Place	Geo location of user at time of tweet posted
Screen_name	Profile or Screen name of the user
Lang	Language opted by user
user_location	Location of the user
user_followers_count	Count of followers of user
user_friends_count	Count of friends of user
user_favourites_count	Count of people favorited the tweet
reply_count	Count of replies to the tweet
retweet_count	Count of retweets of the tweet

## •Detail design of Features with Project Workflow



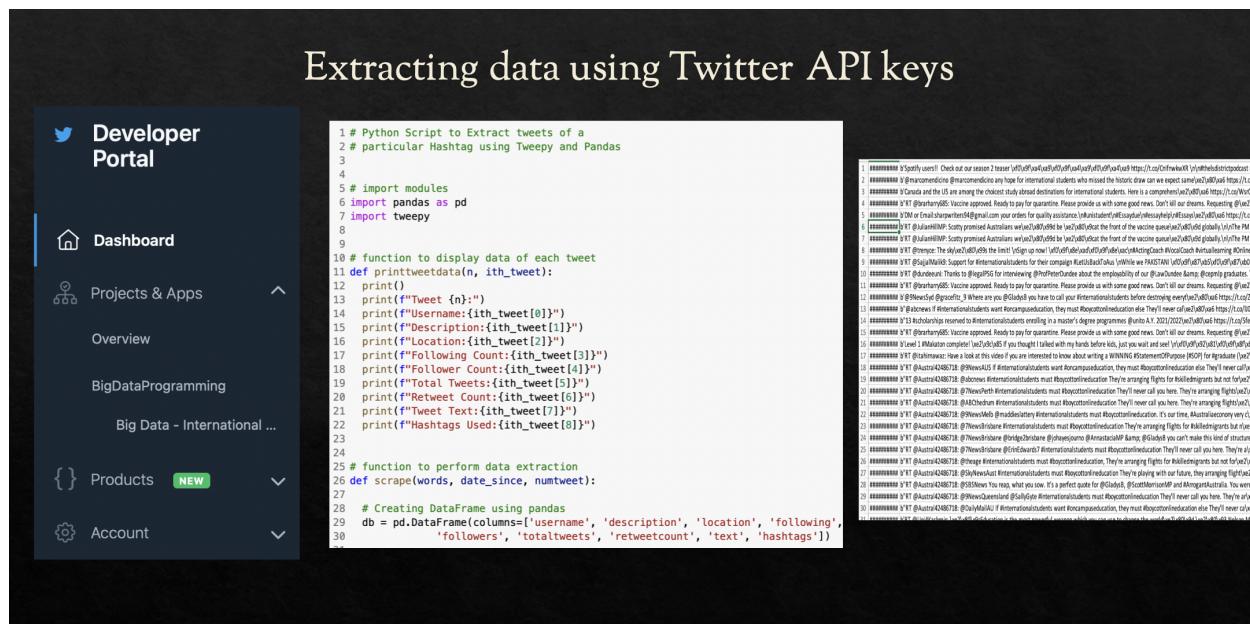
## •Analysis (Details about data)

We have extracted dataset using twitter API, these datasets are extracted using different hashtags focusing on different events related to international students, dividing the events and dataset is crucial since we want this project to be informational and focus on different aspects through different datasets. We will use sentiment analysis on the text (tweet) and use different queries to extract important features and visualize them in Tableau.

## •Implementation (Using Hadoop and Cassandra for Analysis)

**Data extraction :-**

## Extracting data using Twitter API keys



## Preprocessing of the data

## Preprocessing

```

1 def remove_pattern(input_txt, pattern):
2     r = re.compile(pattern)
3     for i in r:
4         input_txt = re.sub(i, '', input_txt)
5
6     return input_txt
7
8 # remove twitter handles (@user)
9 train['tidy_tweet'] = np.vectorize(remove_pattern)(train['tweets'], '@[\\w]*')
10
11 # remove special characters, numbers, punctuations
12 train['tidy_tweet'] = train['tidy_tweet'].str.replace("[^a-zA-Z]", " ")
13
14 #removing words less than length 2
15 train['tidy_tweet'] = train['tidy_tweet'].apply(lambda x: ' '.join([w for w in x.split() if len(w)>2]))
16
17 tokenized_tweet = train['tidy_tweet'].apply(lambda x: x.split())
18 tokenized_tweet.head()
19
20 [any, hope, for, international, students, who, ...
21 [Canadian and the, are, among, the, newest, ...
22 [and, the, opportunity, to, pay, their, ...
23 [Email, sharpwriters, com, your, orders, for, ...
24 [Scotty, promised, Australians, cat, the, froh...
Name: tidy_tweet, dtype: object
25
26
27 #stemming
28 from nltk.stem.porter import *
29 stemmer = PorterStemmer()
30
31 tokenized_tweet = tokenized_tweet.apply(lambda x: [stemmer.stem(i) for i in x]) # stemming
32 tokenized_tweet.head()

```

	tweets	tidy_tweet
0	b@marcomendicino @marcomendicino any hope for...	ani hope for intern student who miss the histo...
1	b@Canada and the US are among the choicest stu...	canada and the are among the choicest studi ab...
2	b@RT @brarharry685: Vaccine approved. Ready to...	vaccin approv readi pay for quarantin pleas pr...
3	b@DM or Email:sharpwriters94@gmail.com your or...	email sharpwrit com your order for qualiti ass...
4	b@RT @JulianHillMP: Scotty promised Australian...	scotti promis australian cat the front the vac...

Solr

## 1. Creation/generation of instance & Collection:

```
cloudera@quickstart:~
```

```
File Edit View Search Terminal Help
```

```
Documents Musical_instruments_reviews.csv workspace
Downloads newacad.java
```

```
[cloudera@quickstart ~]$ solrctl instancedir --generate /tmp/bdp_6
[cloudera@quickstart ~]$ ls /tmp/bdp_6/conf
admin-extra.html _schema_analysis_synonyms_english.json
admin-extra.menu-bottom.html schema.xml
admin-extra.menu-top.html scripts.conf
clustering solrconfig.xml
currency.xml solrconfig.xml.secure
elevate.xml spellings.txt
lang stopwords.txt
mapping-FoldToASCII.txt synonyms.txt
mapping-ISOLatin1Accent.txt update-script.js
protwords.txt velocity
_rest_managed.json xslt
_schema_analysis_stopwords_english.json
```

```
[cloudera@quickstart ~]$ ls /tmp/bdp_6/conf/schema.xml
/tmp/bdp_6/conf/schema.xml
```

```
[cloudera@quickstart ~]$ gedit /tmp/bdp_6/conf/schema.xml
[cloudera@quickstart ~]$ solrctl instancedir --create bdp_6 /tmp/bdp_6
Uploading configs from /tmp/bdp_6/conf to quickstart.cloudera:2181/solr. This may take up to a minute.
[cloudera@quickstart ~]$ solrctl collection --create bdp_6
[cloudera@quickstart ~]$
```

2. Edit the schema.xml created with the instance generation inside the configuration folder to change the attributes based on the dataset given:

```

@ schema.xml X
<field name="Tweet Id" type="string" indexed="true" stored="true" required="true" multiValued="false" />
<field name="Text" type="string" indexed="true" stored="true" />
<field name="Name" type="string" indexed="true" stored="true" />
<field name="Screen Name" type="string" indexed="true" stored="true" />
<field name="UTC" type="string" indexed="true" stored="true" />
<field name="Created At" type="string" indexed="true" stored="true" />
<field name="Favorites" type="string" indexed="true" stored="true" />
<field name="Retweets" type="string" indexed="true" stored="true" />
<field name="Language" type="string" indexed="true" stored="true" />
<field name="Client" type="string" indexed="true" stored="true" />
<field name="Tweet Type" type="string" indexed="true" stored="true" />
<field name="URLs" type="string" indexed="true" stored="true" />
<field name="Hashtags" type="string" indexed="true" stored="true" />
<field name="Mentions" type="string" indexed="true" stored="true" />
<field name="Media Type" type="string" indexed="true" stored="true" />
<field name="Media URLs" type="string" indexed="true" stored="true" />
<field name="" type="string" indexed="true" stored="true" />
<field name="" type="string" indexed="true" stored="true" />
<field name="" type="string" indexed="true" stored="true" />

<!-- points to the root document of a block of nested documents. Required for nested
document support, may be removed otherwise
-->
<field name="_root_" type="string" indexed="true" stored="false"/>

<field name="sku" type="text_en_splitting_tight" indexed="true" stored="true" omitNorms="true"/>
<field name="name" type="text_general" indexed="true" stored="true"/>
<field name="manu" type="text_general" indexed="true" stored="true" omitNorms="true"/>
<field name="cat" type="string" indexed="true" stored="true" multiValued="true"/>
<field name="features" type="text_general" indexed="true" stored="true" multiValued="true"/>
<field name="includes" type="text_general" indexed="true" stored="true" termVectors="true" termPositions="true" termOffsets="true" />

<field name="weight" type="float" indexed="true" stored="true"/>
<field name="price" type="float" indexed="true" stored="true"/>
<field name="popularity" type="int" indexed="true" stored="true" />
<field name="inStock" type="boolean" indexed="true" stored="true" />

<field name="store" type="location" indexed="true" stored="true"/>

```

3. Open the solr browser in the web browser and select the create collection on the left side dropdown.

4. Set "Tweet Id" as the primary key.

```

<!-- Field to use to determine and enforce document uniqueness.
Unless this field is marked with required="false", it will be a required field
-->
<uniqueKey> Tweet Id </uniqueKey>

```

5. Select the document type to csv. Then copy paste all the data inside the dataset into the documents field and submit the document

Request-Handler (qt)  
/update

Document Type  
CSV

Document(s)

	Tweet Id	Text	Name	Screen Name	UTC	Created At	Favorites	Retweets	Language	Client	Tweet Type	URLs	Hashtags	Mentions	Media Type	Media URLs
1	1374864939853225986	"RT @edu_visa_global : Counted as one of the top study destinations. #USA has a multitude of high-ranking programs to choose from.														
2	1374864939853225986	Get 100% Fee Waiver for the upcoming Intake. Contact us for any queries at https://t.co/3OcKGmEEfd														
3	1374864939853225986	#StudyInUSA #Education #StudyAbroad #Students #Scholarships #ApplyNow														
4	1374864939853225986	https://t.co/atjvfoR3j5", Education														
5	1374864939853225986	World,education_24x7,2021-03-24T23:25:20.000Z,Wed Mar 24 23:25:20 +0000														
6	1374864939853225986	2021 0 0 en Retweet https://edu-visa.com/contact/ 7 0 photo https://nhs.twimg.com														

Commit Within  
1000

Overwrite  
true

Submit Document

## HIVE:

### Created the database and database table in hive and loaded the data into hive table for AbolishICE dataset:

```

hive> create table AbolishICE (Tweet_Id string,Text string,Name string,Screen_Name string,UTC string,Created_At string,Favorites string,Retweets string,Language string,Client string,Tweet_Type string,URLs string,Hashtags string,Mentions string,Media_Type string,Media_URLs string) row format delimited fields terminated by ',' stored as textfile;
OK
Time taken: 0.353 seconds

```

```

hive> load data local inpath '/home/cloudera/Downloads/bdphive/AbolishICE.csv' into table AbolishICE;
Loading data to table default.abolishice
Table default.abolishice stats: [numFiles=1, totalSize=645903]
OK
Time taken: 0.598 seconds
hive> select * from AbolishICE limit 2;
OK
Tweet_Id      Text       Name        Screen_Name    UTC        Created_At   Favorites     Retweets    Language     Client      Tweet_Type   URLs        Hashtags    Mentions    Media_Type   Media_URLs
1374864939853225986 "RT @abolishICE_mny : After living in the US for nearly his whole life Hieu and 32 other Vietnamese refugees were deported by @POTUS and @VP last week. Please support his re-entry by donating to this fundraiser set up by @viettlead and @MekongNYC. If you don't have Venmo let us know. #AbolishICE https://t.co/T71Ny7d80" carb         carbnim 2021-03-25T00:48:30.000Z Thu Mar 25 00:48:30 +0000 0          0          en          <a href="https://mobile.twitter.com" rel="nofollow">Twitter Web App</a> Retweet      1          4          photo
Time taken: 0.382 seconds, Fetched: 2 row(s)
hive>

```

## **Created the database and database table in hive and loaded the data into hive table for F1visa dataset: Visualized the output:**

```
hive> create table Flvisa (Tweet_Id string,Text string,Name string,Screen_Name string,UTC string,Created_At string,Favorites string,Retweets string,Language string,Client string,Tweet_Type string,URLs string,Hashtags string,Mentions string,Media_Type string,Media_URLs string) row format delimited fields terminated by ',' stored as textfile;
OK
Time taken: 0.071 seconds
hive> load data local inpath '/home/cloudera/Downloads/bdphive/Flvisa.csv' into table Flvisa;
Loading data to table default.flvisa
Table default.flvisa stats: [numFiles=1, totalSize=40034]
OK
Time taken: 0.257 seconds
hive> select * from Flvisa limit 3;
OK
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| Tweet_Id | Text | Name | Screen_Name | UTC | Created_At | Favorites | Retweets | Language | Client | Tweet_Type | URLs | Hashtags | Mentions | Media_Type | Media_URLs |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| 1374834250667716689 | "RT @Maikel159222858 : All categories of immigrant visa in Cuba #F1Visa #F1Visa #F4Visa #F2AVisa #IR1Visa #IR2Visa #IR3Visa #CR1Visa #CR2Visa #K1Visa #K2Visa #K3Visa #K4Visa #DV2020 #DV2021 #amp; #CFRP should be processed in our country until when? | NULL | | | | | | |
| Meibis Aguilar | AguilarMeibis | AguilarMeibis | 2021-03-24T21:23:23.000Z | Wed Mar 24 21:23:23 +0000 2021 | 0 | en | <a href="https://mobile.twitter.com" rel="nofollow">Twitter Web App</a> | Retwe |
| et | 20 | 0 | NULL |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
Time taken: 0.064 seconds, Fetched: 3 row(s)
hive>
```

## **Created the database and database table in hive and loaded the data into hive table for h1b dataset: Visualized the output:**

```
hive> create table h1b (Tweet_Id string,Text string,Name string,Screen_Name string,UTC string,Created_At string,Favorites string,Retweets string,Language string,Client string,Tweet_Type string,URLs string,Hashtags string,Mentions string,Media_Type string,Media_URLs string) row format delimited fields terminated by ',' stored as textfile;
OK
Time taken: 0.108 seconds
hive> load data local inpath '/home/cloudera/Downloads/bdphive/h1b.csv' into table h1b;
Loading data to table default.h1b
Table default.h1b stats: [numFiles=1, totalSize=591164]
OK
Time taken: 0.339 seconds
hive> select * from h1b limit 2;
OK
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| Tweet_Id | Text | Name | Screen_Name | UTC | Created_At | Favorites | Retweets | Language | Client | Tweet_Type | URLs | Hashtags | Mentions | Media_Type | Media_URLs |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| 1374834189610576 | "RT @mohammed_raheem7utu: drive back home ?? home turned soar please dont for brothers surgery https://t.co/Vhfd1vEB #Sanjose #muslim #MuslimLivesMatter #veinbust #Indian #H1B #f1 visa" | Al Hindi | ??? | | | | | |
| andreaslin101 | 2021-03-25T00:55:00.000Z | Thu Mar 25 00:55:00 +0000 2021 | 0 | 0 | en | <a href="http://twitter.com/download/iphone" rel="nofollow">Twitter for iPhone</a> | Tweet | https://www.gofundme.com/f/pl |
| ease-help-for-the-surgery-of-mohammed-raheem7utu source=twitterutm_medium=social&utm_campaign=_pd&share_sheet | 8 | 0 |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
Time taken: 0.107 seconds, Fetched: 2 row(s)
hive>
```

## **Created the database and database table in hive and loaded the data into hive table for intlstudents dataset: Visualized the output:**

```
hive> create table intlstudents (Tweet_Id string,Text string,Name string,Screen_Name string,UTC string,Created_At string,Favorites string,Retweets string,Language string,Client string,Tweet_Type string,URLs string,Hashtags string,Mentions string,Media_Type string,Media_URLs string) row format delimited fields terminated by ',' stored as textfile;
OK
Time taken: 0.009 seconds
hive> load data local inpath '/home/cloudera/Downloads/bdphive/intlstudents.csv' into table intlstudents;
Loading data to table default.intlstudents
Table default.intlstudents stats: [numFiles=1, totalSize=128138]
OK
Time taken: 0.211 seconds
hive> select * from intlstudents limit 3;
OK
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| Tweet_Id | Text | Name | Screen_Name | UTC | Created_At | Favorites | Retweets | Language | Client | Tweet_Type | URLs | Hashtags | Mentions | Media_Type | Media_URLs |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| 137480872675276689 | "RT @rus_students: This China policy silence towards #intlstudents and the way universities treat us make me feel like the whole COVID19 pandemic is our fault ?? #TakeUsBackToChina" | Mr.Chips????? | | | | | | | |
| I_mzainil | 2021-03-25T01:01:37.000Z | Thu Mar 25 01:01:37 +0000 2021 | 0 | 0 | en | <a href="http://twitter.com/download/android" rel="nofollow">Twitter for Android</a> | Retweet | 2 | 0 |
| 1374862856592617476 | "RT @clayhenesly: During a riveting #PIELive21 discussion on disruptive #edtech solutions driving #digital #recruitment @cialfoplatform's @rohanpasari notes ""engagement (w/ #intlstudents & counselors is not an event but a process." Resonates w/ seasonal #intlst students. More things change... https://t.co/dC4yVgsXTq" | Rohan Pasari | rohanpasari | 2021-03-24T23:17:03 +0000 2021 | 0 |
| en | <a href="http://twitter.com/download/iphone" rel="nofollow">Twitter for iPhone</a> | Retweet | 6 | 2 |
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
Time taken: 0.044 seconds, Fetched: 3 row(s)
hive>
```

## **Created the database and database table in hive and loaded the data into hive table for study in USA dataset: Visualized the output:**

```
hive> create table studyinusa (Text string,Name string,Screen Name string,UTC string,Created At string,Favorites string,Retweets string,Language string,Client string,Tweet Type string,URLs string,Hashtags string,Mentions string,Media Type string,Media URLs string) row format delimited fields terminated by ',' stored as textfile;
Time taken: 0.001 seconds
hive> load data local inpath '/home/cloudera/Downloads/bdphive/studyinusa.csv' into table studyinusa;
Loading data to table default.studyinusa...
Table default.studyinusa stats: [numFiles=1, totalSize=123163]
OK
Time taken: 0.283 seconds
+-----+-----+
| From| Id | Text | Name | Screen Name | UTC | Created At | Favorites | Retweets | Language | Client | Tweet Type | URLs | Hashtags | Mentions | Media Type | Media URLs |
+-----+-----+
1 1374889172675276803 "RT @rus_002: [REDACTED] During a riveting #PELIVE21 discussion on disruptive #edtech solutions driving #digital_recruitment but a process... Resonates w/ seasoned #intld students. The more things change... https://t.co/dc4yqgsxTq" Rohan Pasari rohanpasari 2021-03-24T23:17:08.000Z Wed Mar 24 23:17:08 +0000 2021 0
1 137482381329979971 "RT @UAdvocacy #EDTrust Thank you so much for participating in our chat and for advocating for #immigrant & #intldstudents! #HigherEdImmChat Presidents' Alliance on Higher Ed & Immigration PresImmAlliance 2021-03-24T20:43:36.000Z Wed Mar 24 20:43:36 +0000 2021 0
1 137482474138202173 "RT @rus_002: [REDACTED] During a riveting #PELIVE21 discussion on disruptive #edtech solutions driving #digital_recruitment but a process... Resonates w/ seasoned #intld folk. The more things change... https://t.co/dc4yqgsxTq" Claito cialfoplatforms 2021-03-24T20:36.000Z Wed Mar 24 20:36:36 +0000 2021 0
1 137480587695686273 "RT @clahenleys : During a riveting #PELIVE21 discussion on disruptive #edtech solutions driving #digital_recruitment but a process... Resonates w/ seasoned #intld folk. The more things change... https://t.co/dc4yqgsxTq" The PIE News ThePIENews 2021-03-24T19:38:00.000Z Wed Mar 24 19:38:00 +0000 2021 0
1 1374799944992972 "RT @clahenleys : During a riveting #PELIVE21 discussion on disruptive #edtech solutions driving #digital_recruitment but a process... Resonates w/ seasoned #intld folk. The more things change... https://t.co/dc4yqgsxTq" ThePIELive ThePIELive 2021-03-24T19:07:03.000Z Wed Mar 24 19:07:03 +0000 2021 0
1 1374795481357576 "RT @rus_002: [REDACTED] During a riveting #PELIVE21 discussion on disruptive #edtech solutions driving #digital_recruitment but a process... Resonates w/ seasoned #intld students. The more things change... https://t.co/dc4yqgsxTq" Clay Hensley clahenleys 2021-03-24T18:54:23.000Z Wed Mar 24 18:54:23 +0000 2021 5 4 en <a href="https://about.twitter.com/products/tweetdeck" rel="nofollow">Twitter for iPhone</a>" Tweet https://t.co/71KPR201 University Language Campus Commons 2021-03-24T18:07:00.000Z Wed Mar 24 16:07:00 +0000 2021
1 13745746383920994 "10 tips for succeeding at an American university: https://t.co/6tVU5GeW #intldstudents #college https://t.co/71KPR201 University Language Campus Commons 2021-03-24T16:07:00.000Z Wed Mar 24 14:07:00 +0000 2021
AA87nb.jnp
```

# Imported file into hdfs

```
[cloudera@quickstart Team_6]$ hadoop fs -put /home/cloudera/Downloads/tweets_data_03_21.csv /Team_6
```

```
[cloudera@quickstart ~]$ mkdir Team_6  
[cloudera@quickstart ~]$ cd Team_6
```

**Created the database and database table in hive and loaded the data into hive table for General dataset:**

```
[cloudera@quickstart Team_6]$ hive
```

Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j.properties

```
hive> CREATE TABLE tweets dataFl string, created_at string, id string, lang string, text string, source string, truncated string, in_reply_to_status_id string, in_reply_to_status_id_str string, in_reply_to_user_id string, in_reply_to_user_id_str string, in_reply_to_screen_name string, user string, geo string, coordinates string, place string, contributors string,retweeted_status string, is_quote_status string, quote_count string, reply_count string, retweet_count string, favorite_count string, entities string, favorited_string, retweeted_string, filter_level string, lang string, timestamp_ms string, display_text string, range_string, extended_tweet string, quoted_status_id string, quoted_status_id_str string, quoted_status_permalink string, extended_entities string, possibly_sensitive string, withheld_in_countries string) row format delimited fields terminated by '\t' stored as textfile;
OK
Time taken: 6.212 seconds
hive> load data local inpath '/home/cloudera/Downloads/tweets_data_03_21.csv' into table tweets_data;
OK
hive> select count(*) from tweets_data;
OK
Table default_tweets_data stats: [numFiles=1, totalSize=33422781]
OK
Time taken: 3.601 seconds
hive> select * from tweets_data limit 10;
OK
```

## Visualized the Output:

## Loading the studentdata into the hadoop

1. First imported the .txt files into cloudera and then created a directory with name “StudentData” and moved the imported files into the directory

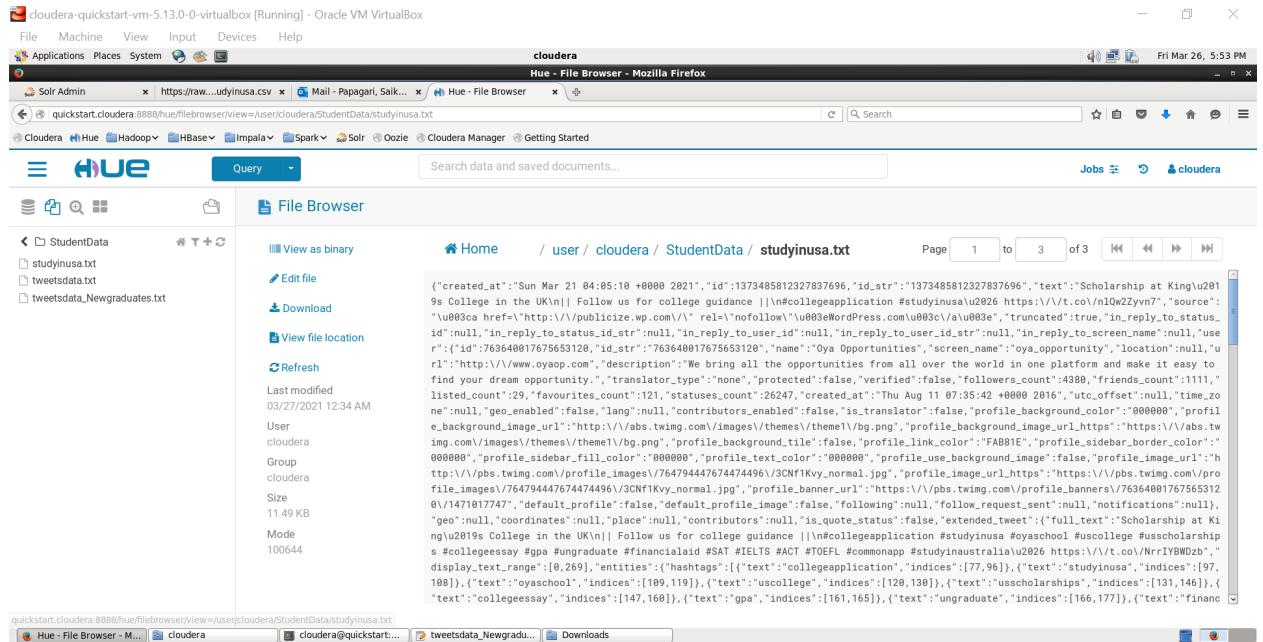


2. The output data which is loaded into Hue for tweetsdata\_Newgraduates.txt

3. Similarly, for the tweetsdata.txt

A screenshot of a Windows desktop showing the Cloudera Quickstart VM running in Oracle VM VirtualBox. The taskbar at the top has icons for File, Machine, View, Input, Devices, Help, Applications, Places, System, and a Cloudera Manager icon. The system tray shows the date as Fri Mar 26, 5:52 PM. The main window is the Hue File Browser in Mozilla Firefox, displaying a file named 'tweetsdata.txt' located at /user/cloudera/StudentData/tweetsdata.txt. The file content is a JSON object representing a tweet from Michelle Obama. The left sidebar shows a tree view of 'StudentData' containing 'studyinusa.txt', 'tweetsdata.txt', and 'tweetsdata\_Newgraduates.txt'. The top navigation bar includes links for Solr Admin, Mail - Papagari, Saik..., Hue - File Browser, and Cloudera Manager.

## 4. The output for file study.txt is as follows:



The screenshot shows the Hue File Browser interface. The left sidebar lists files: studyinusa.txt, tweetsdata.txt, and tweetsdata\_Newgraduates.txt. The main area displays the content of studyinusa.txt. The content is a JSON object representing a tweet from a user named 'udyinusa'. The tweet's text is a link to a scholarship opportunity at King's College in the UK. It includes various metadata fields like id, created\_at, and user information.

```
{"created_at": "Sun Mar 21 04:05:10 +0000 2021", "id": 1373485812327837696, "id_str": "1373485812327837696", "text": "Scholarship at King's College in the UK|| Follow us for college guidance ||\n#collegeapplication #studyinusa #u2026 https://t.co/n1QwZ2yvn7", "source": "\u0003<a href=\"http://publicize.wp.com/?\" rel=\"nofollow\">u083<#WordPress.com/u083</a\u0003", "truncated": true, "in_reply_to_status_id": null, "in_reply_to_status_id_str": null, "in_reply_to_user_id": null, "in_reply_to_user_id_str": null, "in_reply_to_screen_name": null, "user": {"id": 163640817675653126, "id_str": "163640817675653126", "name": "Oya Opportunities", "screen_name": "oya.opportunity", "location": null, "url": "http://www.oyaop.com", "description": "We bring all the opportunities from all over the world in one platform and make it easy to find your dream opportunity.", "translator_type": "none", "protected": false, "verified": false, "followers_count": 4388, "friends_count": 1111, "listed_count": 29, "favourites_count": 121, "statuses_count": 26247, "created_at": "Thu Aug 11 07:35:42 +0000 2016", "utc_offset": null, "time_zone": null, "geo_enabled": false, "lang": null, "contributors_enabled": false, "is_translator": false, "profile_background_color": "#000000", "profile_background_image_url": "https://abs.twimg.com/images/themes/theme1/bg.png", "profile_background_tile": false, "profile_link_color": "#FAB81E", "profile_sidebar_border_color": "#000000", "profile_sidebar_fill_color": "#000000", "profile_text_color": "#000000", "profile_use_background_image": false, "profile_image_url": "http://pbs.twimg.com/profile_images/76479447674474496/3CNFIKvy_normal.jpg", "profile_image_url_https": "https://pbs.twimg.com/profile_images/76479447674474496/3CNFIKvy_normal.jpg", "profile_banner_url": "https://pbs.twimg.com/profile_banners/763640017675653126/1471017747", "default_profile": false, "default_profile_image": false, "following": null, "follow_request_sent": null, "notifications": null}, "geo": null, "coordinates": null, "place": null, "contributors": null, "is_quote_status": false, "extended_tweet": {"full_text": "Scholarship at King's College in the UK|| Follow us for college guidance ||\n#collegeapplication #studyinusa #oyaop #uscollege #usscholarship #collegessay #gpa #ungraduate #financialaid #SAT #IELTS #ACT #TOEFL #commonapp #studyinaustralia #u2026 https://t.co/NrrIYWDzb", "display_text_range": [8, 269], "entities": [{"text": "#hashtag": [{"text": "#collegeapplication", "indices": [77, 96]}, {"text": "#studyinusa", "indices": [97, 108]}, {"text": "#oyaop", "indices": [109, 119]}, {"text": "#uscollege", "indices": [120, 130]}, {"text": "#usscholarships", "indices": [131, 146]}, {"text": "#collegessay", "indices": [147, 160]}], "text": "#ungraduate", "indices": [166, 177]}], "text": "#ungraduate", "indices": [166, 177]}], "text": "#financ
```

## Working with Spark using Scala:

The number of Tweets from various countries (Query)

```
scala> val q1 = sqlContext.sql("SELECT place.country,count(*) AS count FROM tweetDatatable GROUP BY place.country ORDER BY count DESC limit 10");  
q1: org.apache.spark.sql.DataFrame = [country: string, count: bigint]
```

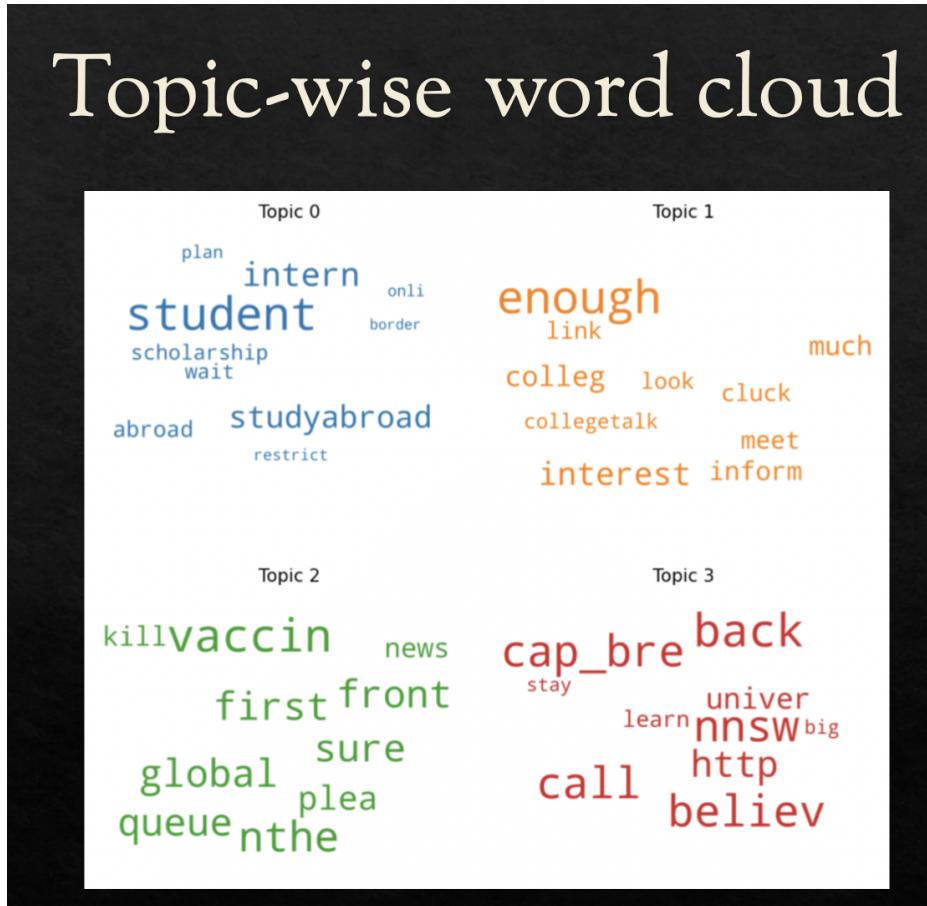
To retrieve the various languages in which tweets about web series are created.

The number of tweets about the web series in various languages was tallied.

```
scala> val q3=sqlContext.sql("select count(*) as count,lang as language from tweetDatatable where lang is not null group by lang");  
q3: org.apache.spark.sql.DataFrame = [count: bigint, language: string]
```

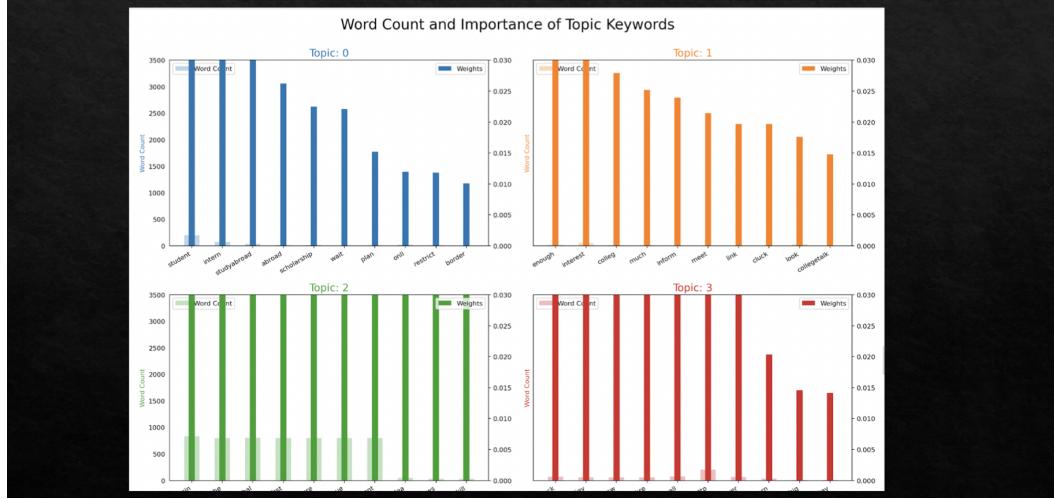
- Preliminary Results (Visualization of Results)

a.

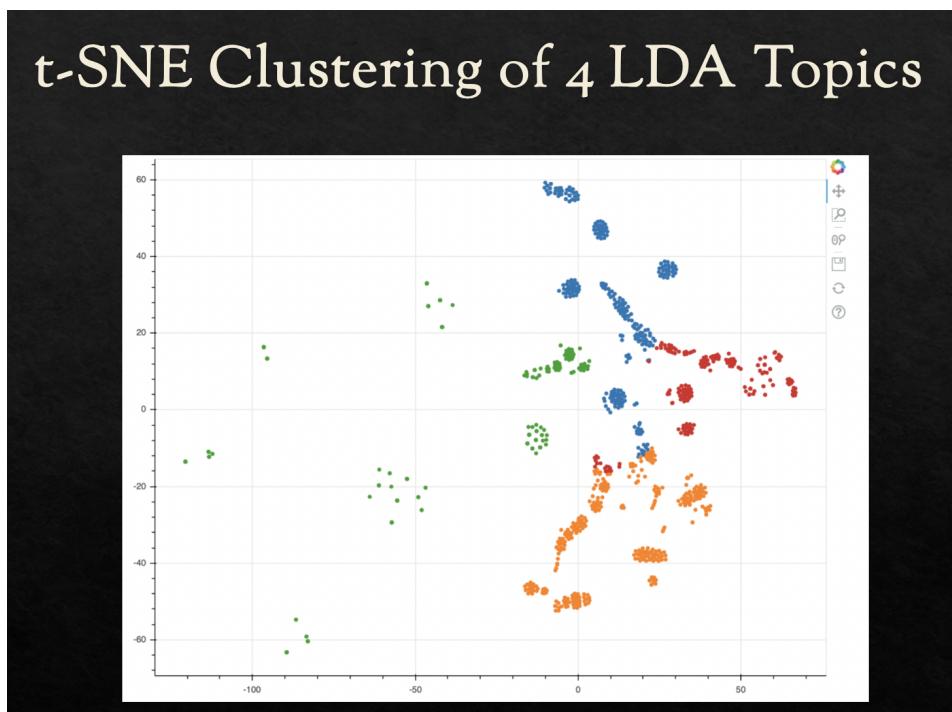


b.

## Importance of each word in the topic and word count

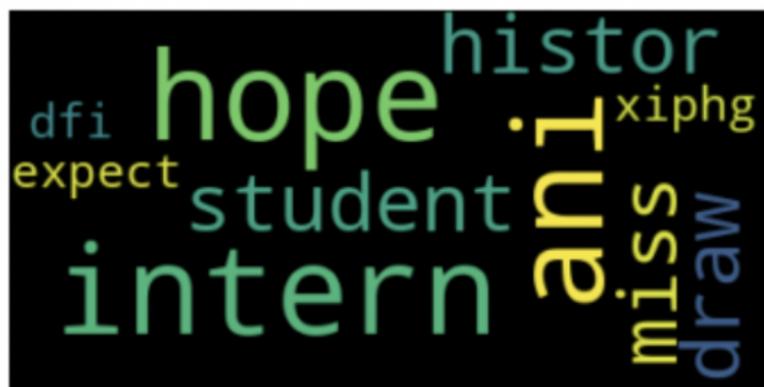


C.



d. Creating word cloud using pySpark

```
1 text = tidy_tweets[0]
2
3 # Create and generate a word cloud image:
4 wordcloud = WordCloud().generate(text)
5
6 # Display the generated image:
7 plt.imshow(wordcloud, interpolation='bilinear')
8 plt.axis("off")
9 plt.show()
```



- **Project Management**

- **Implementation status report**

- **Work completed:**

- **Description :-**

We have completed the extraction of dataset from twitter using twitter API, completed the tweets preprocessing part, sentiment analysis of those tweets. We have also started working on hadoop and pySpark and completed loading of the dataset in hive and creating tables, and creating pySpark and working on the tweets (text processing) aspect in spark to perform different analysing techniques. Loaded dataset into solr as well to work on the analysing aspect.

- **Responsibility (Task, Person)**

**Amar** - Data Extraction using tweepy, Creating collection and Instance in Solr, Editing schema.xml in accordance with dataset, Loading the dataset (.csv file) into Solr.

**Sai** - Importing the text files into Hadoop-Hue and exploring on working with Scala queries.

**Mounika** - Data extraction , loading table in Hive and documentation.

**Tanvi** - Twitter data collection from twitter API, preprocessing, (removing numbers, punctuations, words less than length 2, stemming, removing stopwords), using Gensim, Build the bigram and trigram models, LDA model for topic modelling, Word Cloud of Top N words in each topic, plot Word Count and Weights of Topic Keywords, t-SNE clustering of 4LDA topics, using pySpark for word cloud and text analysis/visualization.

- **Contributions (members/percentage)**

**Amar** - 20%

**Sai** - 20%

**Mounika** - 30%

**Tanvi** - 30%

- **Work to be completed**

- **Description**

We have to perform ETL process using spark's batch processing, and analysis using different pySpark codes to analyse the text data and different aspects of it, and also visualize them using Tableau. We have to

- **Responsibility (Task,**

## **Person)**

**Tanvi** - performing more analysis and text visualization and sentiment analysis in pySpark on the other 2 events well (AbolishICE and h1b), visualizing the important data using Tableau

**Mounika** - Implementing SQL and HIVE queries and using sqoop to transfer data from SQL to HDFS

**Amar** - Performing various queries in Solr, sentiment analysis using pySpark, visualisation using seaborn.

**Sai** - Data analysis using Scala for different datasets.

- **Issues/Concerns**

The big data tools and the tasks performed by increment 2 might alter a bit later on in the final phase of the project based on the time and knowledge constraints and what aspects will be important to showcase, if the project demands only hadoop, spark and solr, or addition of another big data tool for better data analysis, this can be only determined after moving forward with current plan and seeing how well it goes or if there is a need to change the plan.

**6. Once you are done with the above report, for story telling part of the project, address all the questions given in the following link for Increment 2 storytelling under Assignment 02section:**

### **1. Who**

The international students who are studying in the United States Of America.(or overseas).

The dataset is extracted using different hashtags to get data related to 3 different scenarios as follows :-

1. Generic International Student Data (**#F1visa #intlstudents #internationalstudents #studyinUSA**)
2. Immigration rules change for F1 Visa during COVID in 2020 (**#AbolishICE**)
3. Jobs for international students during covid (H1B sponsorship) (**#H1B, #h1bjobs**)

## **2. What**

Yes, the data set records the targeted events, activities, behaviors, etc. in Assignment 1. This is fundamentally about the variables. It records the username, the location, the tweets which tell us about what the users really think about the specific event that happened.

## **3. When**

The events take place on how people reacted to the challenges faced by International students during COVID like immigration rules for F1visa, jobs for international students during covid(H1B sponsorship).

## **4. Where means two things:**

These challenges can be faced by any international students in a foreign land, where they are a part of the immigrant community and do not have access to the perks that The citizens of that country get, which makes life a bit harder for even the very basic requirements of living.

## **5. Why means the possible causes and/origin of the problem.**

Due to the increase in the intake of international students over a period of time.

## **6. How:**

It happened during COVID, as a result there was a massive Unavailability of jobs(new graduates), many have lost their jobs due to pandemic, immigration rules change for F1 Visa during COVID.

## **References:**

1. [https://docs.cloudera.com/documentation/enterprise/6/6.3/topics/cm\\_mc\\_solr\\_service.html](https://docs.cloudera.com/documentation/enterprise/6/6.3/topics/cm_mc_solr_service.html)
2. <https://docs.cloudera.com/runtime/7.2.7/search-managing/topics/search-updating-the-schema-in-a-solr-collection.html>
3. <https://github.com/dsuarez993/bigdata-realtime-twitter-analysis>
4. <https://github.com/6vedant/TwitterAnalyticsHadoop>
5. <https://www.toptal.com/apache/apache-spark-streaming-twitter>
6. <http://spark.apache.org/>
7. <https://www.toptal.com/spark/introduction-to-apache-spark>
8. <https://hadoop.apache.org/>
9. <https://www.infoworld.com/article/3236869/what-is-apache-spark-the-big-data-platform-that-crushed-hadoop.html>

10.[https://towardsdatascience.com/extracting-data-from-twitte  
r-using-python-5ab67bff553a](https://towardsdatascience.com/extracting-data-from-twitter-using-python-5ab67bff553a)