



INSTITUTE FOR ADVANCED
COMPUTING AND
SOFTWARE
DEVELOPME
NT AKURDI,
PUNE

Documentation On
“ US Federal Campaign Finance Data Analysis Using BigData Technologies ”
PG-DBDA MARCH 2022

Submitted By:
Group No: 06
UMESH BHARUDE 223313

Mr. Prashant Karhale
Center Coordinator

Mr. Akshay Tilekar
Project Guide

Contents

1.	Introduction	2
1.1	PROBLEM STATEMENT	
1.2	Abstract	
2.	TOOLS & TECHNOLOGIES	3
3.	Upload dataset from kaggle to EMR	4
4.	create the emr cluster and we choose services hadoop, spark and hive	5
5.	we copy the data from s3 bucket to emr	11
6.	Go to the pyspark on HDFS	12
7.	Read multiple csv from hdfs to spark	13
8.	we clean the data ,we remove the null column	15
9.	Save the cleaned data in S3 Bucket into parquet format	19
10.	Load the data from S3 to pyspark	19
11.	Create table and perform query in pyspark and store table in hive	20
12.	Tables which are stored in hive	27
13.	connectivity to power bi	28
14.	PowerBI Dashboard	32
15.	Requirements Specification	33
16.	References	34

1. Introduction

1.1 PROBLEM STATEMENT

To analyze the campaign finance & contributions (individuals & PACs) for all federal elections over the last 2 decades & a half.

1.2 Abstract

US FEDERAL CAMPAIGN FINANCE DATA ANALYSIS using AWS services like S3 Bucket, EMR, CFT, Spark, Hive and PowerBI. To analyze the historical data of US Federal Campaign's finance and contributions made individually and through PACs. I have determined the relationship between various features and their impact on contribution patterns for candidates. Performed ETL in Apache Spark and displayed respective visualizations in Power BI.

The dataset we have taken from kaggle.

Below is the following steps how we imported the data and store the shell script file in s3 bucket

2. TOOLS & TECHNOLOGIES

1. S3 → Storing raw data
2. Extract data from Kaggle and storing in Amazon S3 Bucket
3. CFT : Creating EMR

4.EMR :

Apache Hadoop Ecosystem:

- a. HDFS → Using Hadoop distributed file system for data storage
- b. Apache PySpark → Combine various tables for proper analysis and performing ETL
- c. Apache Hive → To process structured Data

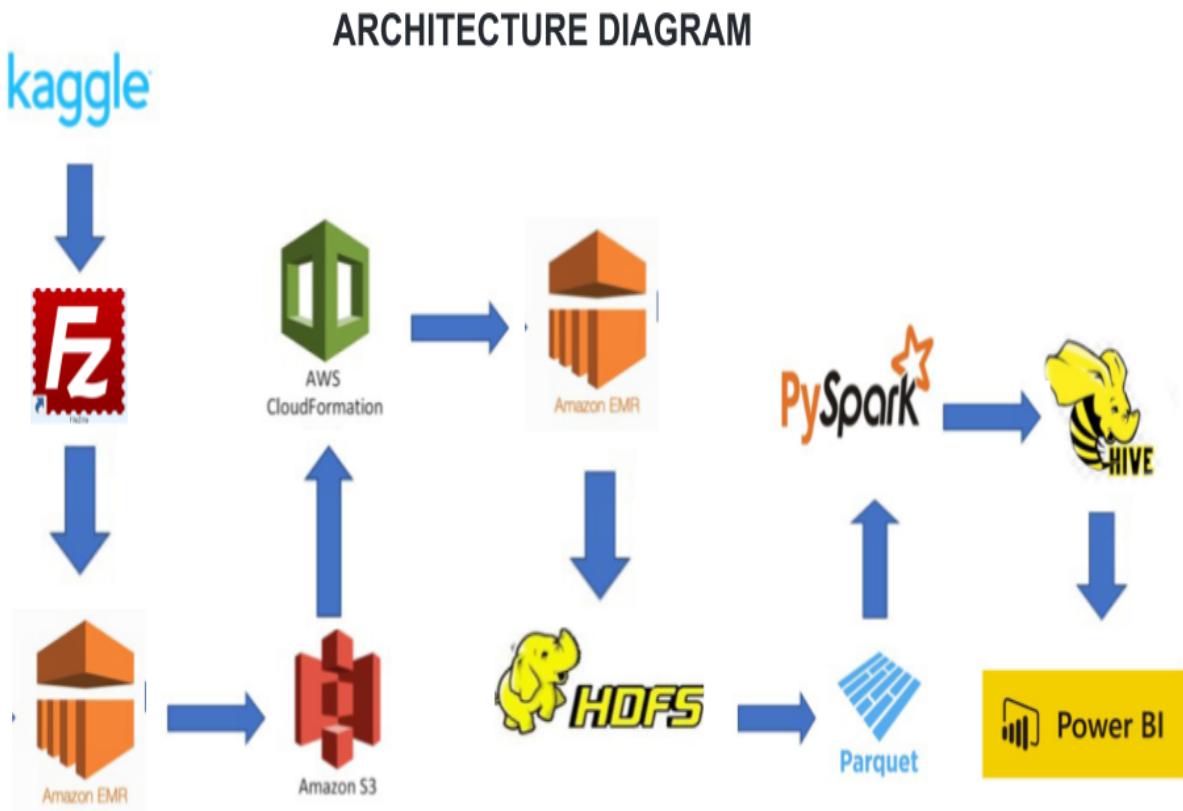
5. Languages to use:

- a. Python → ETL and Data Processing
- b. HQL → Creating Hive tables and loading data
- c. Shell Script → Automation

6. File Formats to be used:

- a. CSV
- b. Parquet

7. Power BI: Creating Dashboard

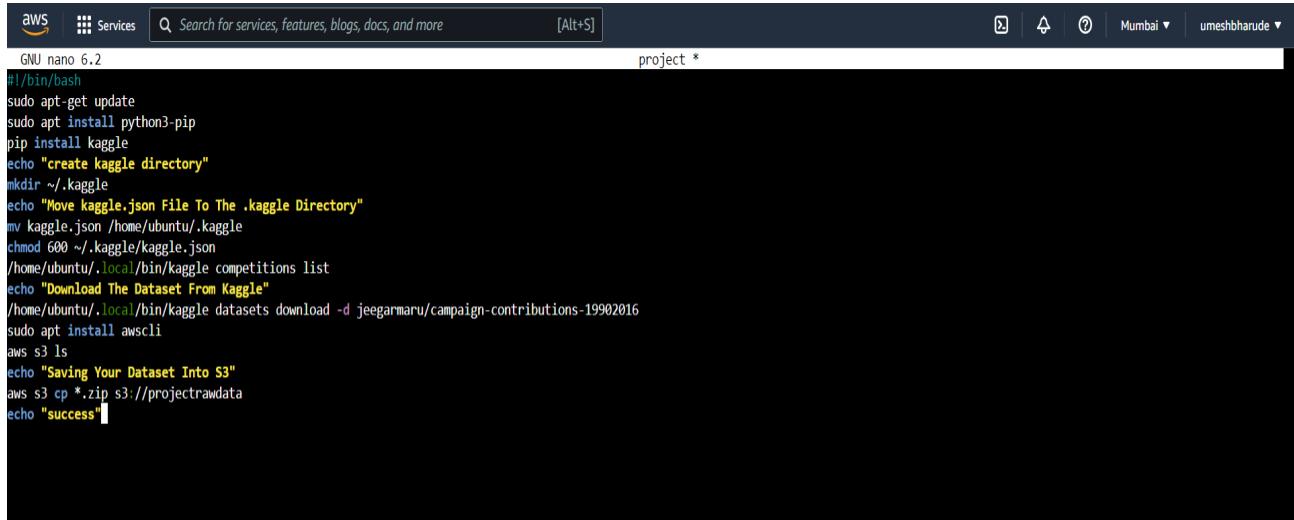


1) Upload dataset from kaggle to EMR:

- Go to the kaggle website.
- Click on Your profile button on the top right and then select My Account.
- Scroll down to the API section and click on the Create New API Token button.
- It will download the file kaggle.json & Save the file at a known location on your machine.
- Move the downloaded file to a location: `~/.kaggle/kaggle.json`.
- Create .kaggle folder in your home directory, you can create one using the command: `mkdir ~/.kaggle`
- Now move the downloaded file to this location using:
- `mv <location>/kaggle.json ~/.kaggle/kaggle.json`

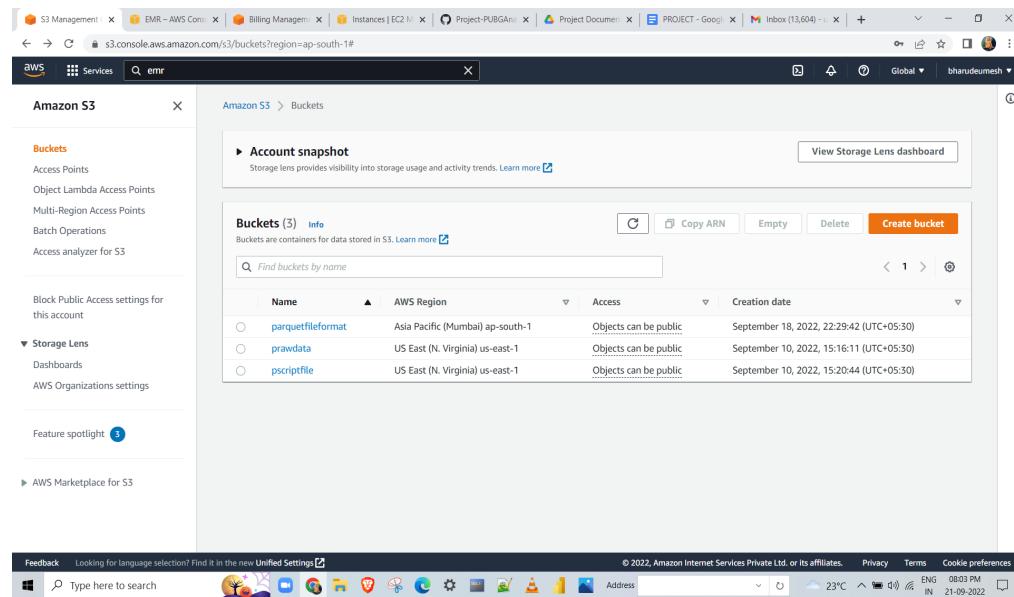
US Federal Campaign Finance Data Analysis Using BigData Technologies

- You need to give proper permissions to the file:
- chmod 600 ~/.kaggle/kaggle.json
- Then run the command as ...> *kaggle competitions list*



```
GNU nano 6.2
#!/bin/bash
sudo apt-get update
sudo apt install python3-pip
pip install kaggle
echo "Create kaggle directory"
mkdir ~/.kaggle
echo "Move kaggle.json File To The .kaggle Directory"
mv kaggle.json /home/ubuntu/.kaggle
chmod 600 ~/.kaggle/kaggle.json
/home/ubuntu/.local/bin/kaggle competitions list
echo "Download The Dataset From Kaggle"
/home/ubuntu/.local/bin/kaggle datasets download -d jeegarmar/campaign-contributions-19902016
sudo apt install awscli
aws s3 ls
echo "Saving Your Dataset Into S3"
aws s3 cp *.zip s3://projectrawdata
echo "success"
```

- We Run the above shell script
- The raw data is download and
- Then the raw data is stored in s3 bucket
- We also store the shell script file in s3 bucket



Amazon S3 > Buckets

Name	AWS Region	Access	Creation date
parquetfileformat	Asia Pacific (Mumbai) ap-south-1	Objects can be public	September 18, 2022, 22:29:42 (UTC+05:30)
prawdata	US East (N. Virginia) us-east-1	Objects can be public	September 10, 2022, 15:16:11 (UTC+05:30)
pscriptfile	US East (N. Virginia) us-east-1	Objects can be public	September 10, 2022, 15:20:44 (UTC+05:30)

US Federal Campaign Finance Data Analysis Using BigData Technologies

2)create the emr cluster and we choose services hadoop, spark and hive :

Software Configuration

- Hadoop 2.10.1
- JupyterHub 1.4.1
- Ganglia 3.7.2
- Hive 2.3.9
- JupyterEnterpriseGateway 2.1.0
- Mahout 0.13.0
- Oozie 5.2.1
- TensorFlow 2.4.1
- Zeppelin 0.10.0
- Tez 0.9.2
- HBase 1.4.13
- Presto 0.267
- Hue 4.10.0
- Spark 2.4.8
- Livy 0.7.1
- Flink 1.14.2
- Pig 0.17.0
- ZooKeeper 3.4.14
- Sqoop 1.4.7
- Phoenix 4.14.3
- HCatalog 2.3.9

Multiple master nodes (optional)

Use multiple master nodes to improve cluster availability. [Learn more](#)

AWS Glue Data Catalog settings (optional)

Use for Hive table metadata

Use for Spark table metadata

Edit software settings

Enter configuration Load JSON from S3

```
classification=config-file-name.properties=[myKey1=myValue1,myKey2=myValue2]
```

Steps (optional)

Node type	Instance type	Instance count	Purchasing option
Master	m5.xlarge	1 Instances	<input checked="" type="radio"/> On-demand <input type="radio"/> Spot <small>Use on-demand as max price</small>
Core	m5.xlarge	0 Instances	<input checked="" type="radio"/> On-demand <input type="radio"/> Spot <small>Use on-demand as max price</small>
Task	m5.xlarge	0 Instances	<input checked="" type="radio"/> On-demand <input type="radio"/> Spot <small>Use on-demand as max price</small>
+ Add task instance group			
Total core and task units		0 Total units	

US Federal Campaign Finance Data Analysis Using BigData Technologies

The screenshot shows the AWS EMR Create Cluster wizard, Step 3: General Cluster Settings. The interface includes tabs for Step 1: Software and Steps, Step 2: Hardware, Step 3: General Cluster Settings (which is selected), and Step 4: Security. The General Cluster Settings section contains fields for Cluster name (set to projectemr), Logging (unchecked), Termination protection (unchecked), and Tags. A Key field is present with the placeholder "Add a key to create a tag". The Operating System Options section shows Amazon Linux Release 2.0.20220719.0 selected, and there is an option to Enter an AMI ID. A checkbox for "Update all installed packages on reboot (recommended)" is checked. Navigation buttons at the bottom include "Cancel", "Previous", and "Next". The browser header shows the URL ap-south-1.console.aws.amazon.com/elasticmapreduce/home?region=ap-south-1#create-cluster.

US Federal Campaign Finance Data Analysis Using BigData Technologies

Create Cluster - Advanced Options

Step 1: Software and Steps Step 2: Hardware Step 3: General Cluster Settings Step 4: Security

Security Options

EC2 key pair: umeshbharude

Cluster visible to all IAM users in account

Permissions

Default Custom

Use default IAM roles. If roles are not present, they will be automatically created for you with managed policies for automatic policy updates.

EMR role: [EMR_DefaultRole](#) Use EMR_DefaultRole_V2

EC2 instance profile: [EMR_EC2_DefaultRole](#)

Auto Scaling role: [EMR_AutoScaling_DefaultRole](#)

► Security Configuration

► EC2 security groups

Create cluster

Amazon EMR

Cluster: projectemr Waiting Cluster ready to run steps.

Summary

ID: J-1LKEASN2PIBD
Creation date: 2022-09-21 16:07 (UTC+5:30)
Elapsed time: 3 minutes
After last step completes: Cluster waits
Termination protection: Off Change
Tags: -- View All / Edit
Master public DNS: ec2-43-205-124-245.ap-south-1.compute.amazonaws.com Connect to the Master Node Using SSH

Configuration details

Release label: emr-5.36.0
Hadoop distribution: Amazon 2.10.1
Applications: Hive 2.3.9, Spark 2.4.8
Log URI: --
EMRFS consistent view: Disabled
Custom AMI ID: --
Amazon Linux Release: 2.0.20220719.0 Learn more

Application user interfaces

Persistent user interfaces: Spark history server, YARN timeline server, Tez UI
On-cluster user Not Enabled Enable an SSH Connection
interfaces:

Network and hardware

Availability zone: ap-south-1a
Subnet ID: subnet-00e245c2312b55616
Master: Running 1 m5.xlarge
Core: --
Task: --
Cluster scaling: Not enabled
Auto-termination: Not enabled

Security and access

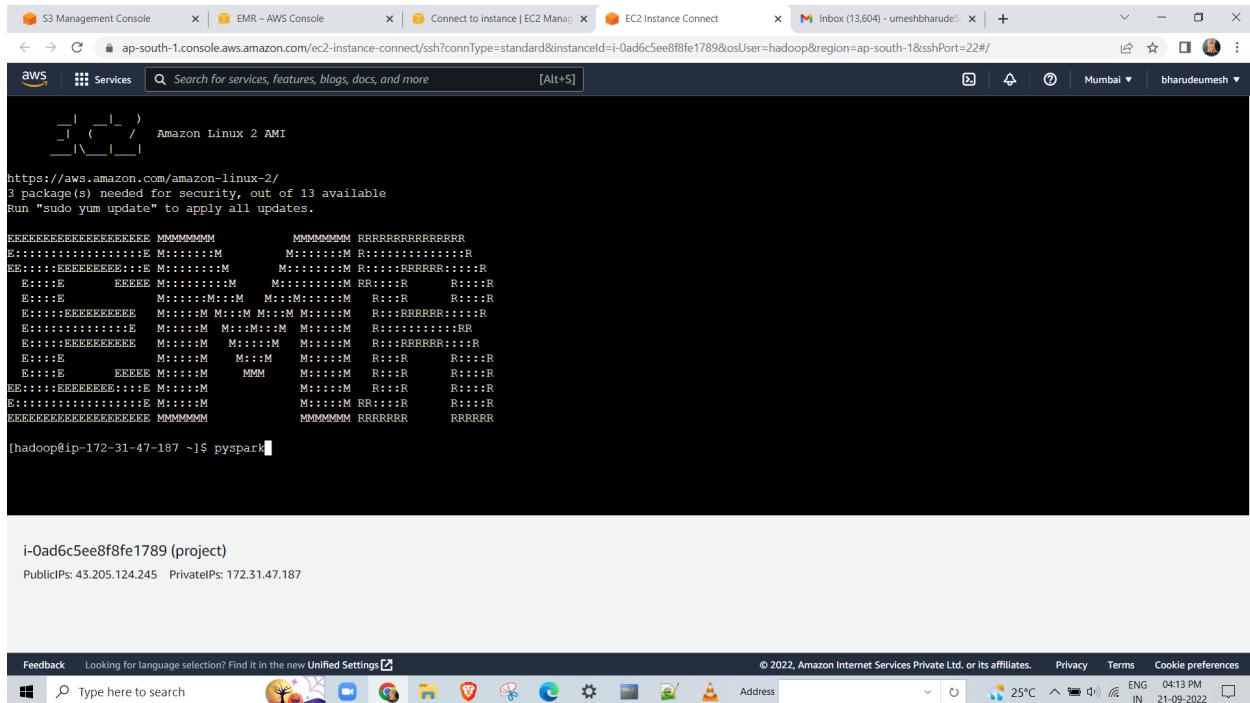
Key name: umeshbharude
EC2 instance profile: EMR_EC2_DefaultRole
EMR role: EMR_DefaultRole

US Federal Campaign Finance Data Analysis Using BigData Technologies

The screenshot shows the AWS Management Console interface. The top navigation bar includes tabs for S3 Management Console, EMR - AWS Console, Instances | EC2 Management Console, and Inbox (13,604) - umeshbharude. The main content area displays the EC2 Instances page, which lists one instance named "project" with the ID i-0ad6c5ee8f8fe1789. The instance is running, m5.xlarge, with 2/2 checks passed and no alarms. It is located in the ap-south-1a availability zone with a public IPv4 address of 43.205.124.245. Below this, a detailed view for instance i-0ad6c5ee8f8fe1789 is shown, including its summary, security groups, networking, storage, status checks, monitoring, and tags. The bottom portion of the screenshot shows the "Connect to instance" dialog for the same instance, where the user has selected the "EC2 Instance Connect" tab and entered the user name "hadoop". A note at the bottom of the dialog states: "Note: In most cases, the guessed user name is correct. However, read your AMI usage instructions to check if the AMI owner has changed the default AMI user name." The browser's address bar shows the URL ap-south-1.console.aws.amazon.com/ec2/home?region=ap-south-1#Instances\$instanceState=running.

IACSD-PG-DBDA-MARCH-22

US Federal Campaign Finance Data Analysis Using BigData Technologies



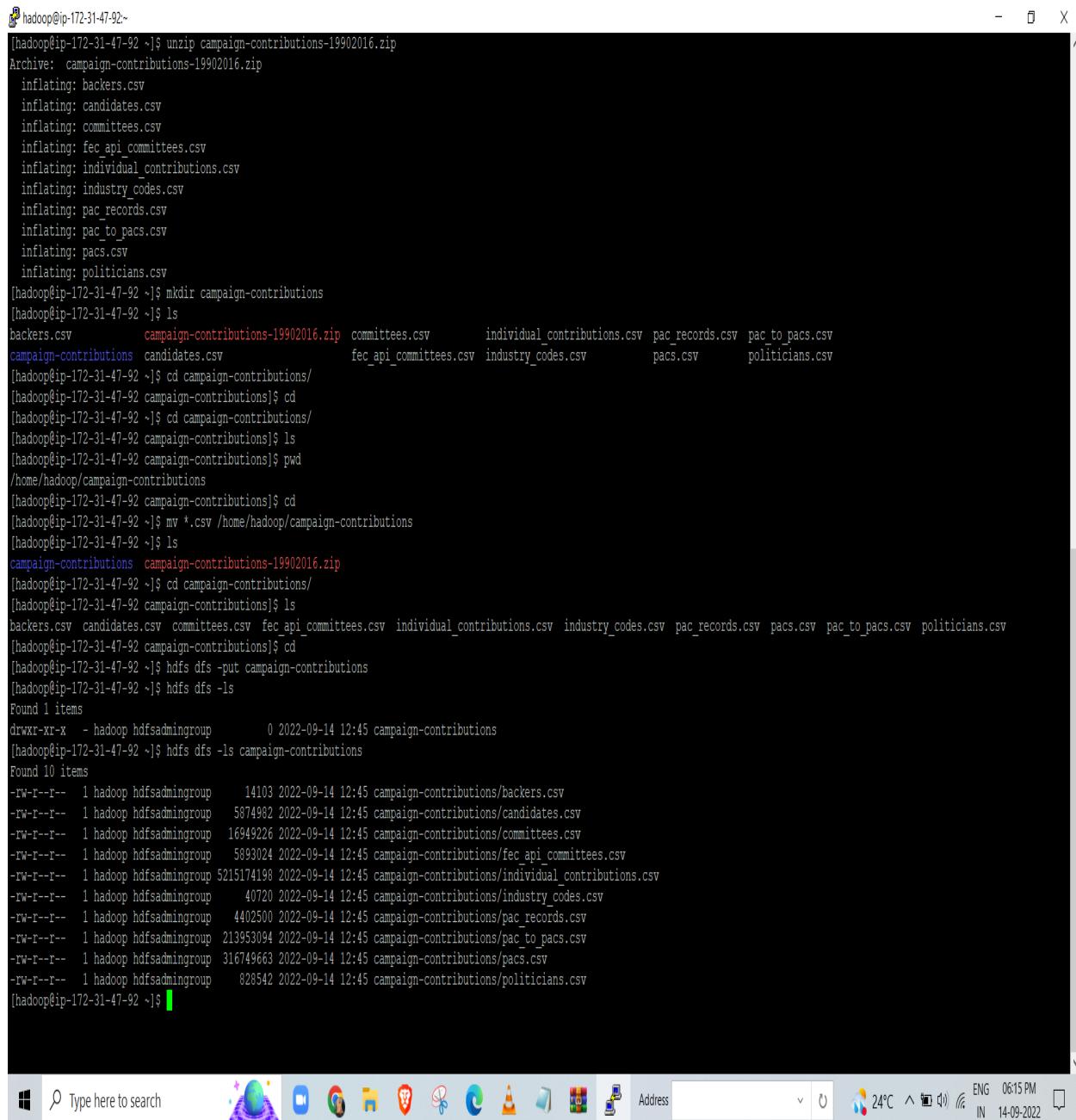
JACSD-PG-DBDA-MARCH-22

3)we copy the data from s3 bucket to emr:

By using the command

- aws s3 <path of bucket> ./
- Then we unzip the data and store data in hdfs by using command :
\$ hdfs dfs -put <directory name >

US Federal Campaign Finance Data Analysis Using BigData Technologies



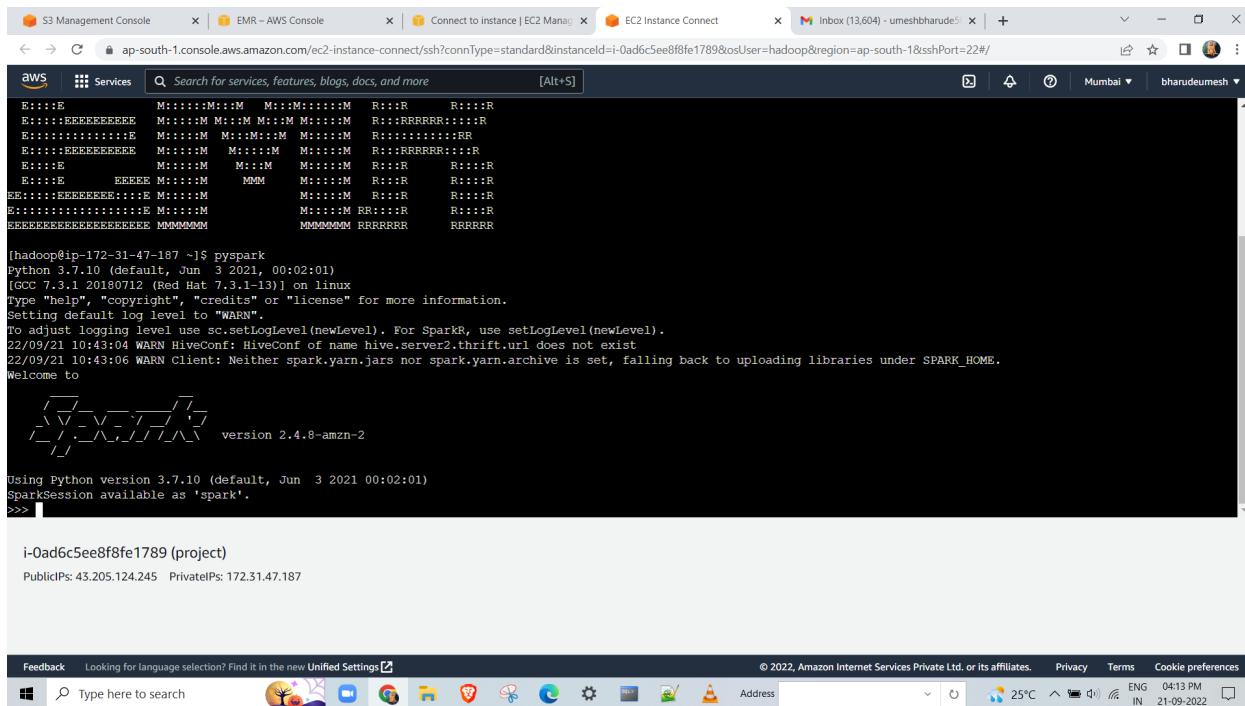
```
[hadoop@ip-172-31-47-92 ~]$ unzip campaign-contributions-19902016.zip
Archive: campaign-contributions-19902016.zip
  inflating: backers.csv
  inflating: candidates.csv
  inflating: committees.csv
  inflating: fec_api_committees.csv
  inflating: individual_contributions.csv
  inflating: industry_codes.csv
  inflating: pac_records.csv
  inflating: pac_to_pacs.csv
  inflating: pacs.csv
  inflating: politicians.csv
[hadoop@ip-172-31-47-92 ~]$ mkdir campaign-contributions
[hadoop@ip-172-31-47-92 ~]$ ls
backers.csv      campaign-contributions-19902016.zip  committees.csv      individual_contributions.csv  pac_records.csv  pac_to_pacs.csv
campaign-contributions  candidates.csv      fec_api_committees.csv  industry_codes.csv    pacs.csv        politicians.csv
[hadoop@ip-172-31-47-92 ~]$ cd campaign-contributions/
[hadoop@ip-172-31-47-92 campaign-contributions]$ cd
[hadoop@ip-172-31-47-92 ~]$ cd campaign-contributions/
[hadoop@ip-172-31-47-92 campaign-contributions]$ ls
[hadoop@ip-172-31-47-92 campaign-contributions]$ pwd
/home/hadoop/campaign-contributions
[hadoop@ip-172-31-47-92 campaign-contributions]$ cd
[hadoop@ip-172-31-47-92 ~]$ mv *.csv /home/hadoop/campaign-contributions
[hadoop@ip-172-31-47-92 ~]$ ls
campaign-contributions  campaign-contributions-19902016.zip
[hadoop@ip-172-31-47-92 ~]$ cd campaign-contributions/
[hadoop@ip-172-31-47-92 campaign-contributions]$ ls
backers.csv  candidates.csv  committees.csv  fec_api_committees.csv  individual_contributions.csv  industry_codes.csv  pac_records.csv  pacs.csv  pac_to_pacs.csv  politicians.csv
[hadoop@ip-172-31-47-92 campaign-contributions]$ cd
[hadoop@ip-172-31-47-92 ~]$ hdfs dfs -put campaign-contributions
[hadoop@ip-172-31-47-92 ~]$ hdfs dfs -ls
Found 1 items
drwxr-xr-x - hadoop hdfsadmin group 0 2022-09-14 12:45 campaign-contributions
[hadoop@ip-172-31-47-92 ~]$ hdfs dfs -ls campaign-contributions
Found 10 items
-rw-r--r-- 1 hadoop hdfsadmin group 14103 2022-09-14 12:45 campaign-contributions/backers.csv
-rw-r--r-- 1 hadoop hdfsadmin group 5874982 2022-09-14 12:45 campaign-contributions/candidates.csv
-rw-r--r-- 1 hadoop hdfsadmin group 16949226 2022-09-14 12:45 campaign-contributions/committees.csv
-rw-r--r-- 1 hadoop hdfsadmin group 5893024 2022-09-14 12:45 campaign-contributions/fec_api_committees.csv
-rw-r--r-- 1 hadoop hdfsadmin group 5215174198 2022-09-14 12:45 campaign-contributions/individual_contributions.csv
-rw-r--r-- 1 hadoop hdfsadmin group 40720 2022-09-14 12:45 campaign-contributions/industry_codes.csv
-rw-r--r-- 1 hadoop hdfsadmin group 4402500 2022-09-14 12:45 campaign-contributions/pac_records.csv
-rw-r--r-- 1 hadoop hdfsadmin group 213953094 2022-09-14 12:45 campaign-contributions/pac_to_pacs.csv
-rw-r--r-- 1 hadoop hdfsadmin group 316749663 2022-09-14 12:45 campaign-contributions/pacs.csv
-rw-r--r-- 1 hadoop hdfsadmin group 828542 2022-09-14 12:45 campaign-contributions/politicians.csv
[hadoop@ip-172-31-47-92 ~]$
```

4)Go to the pyspark on HDFS

Command ==> pyspark

IACSD-PG-DBDA-MARCH-22

US Federal Campaign Finance Data Analysis Using BigData Technologies



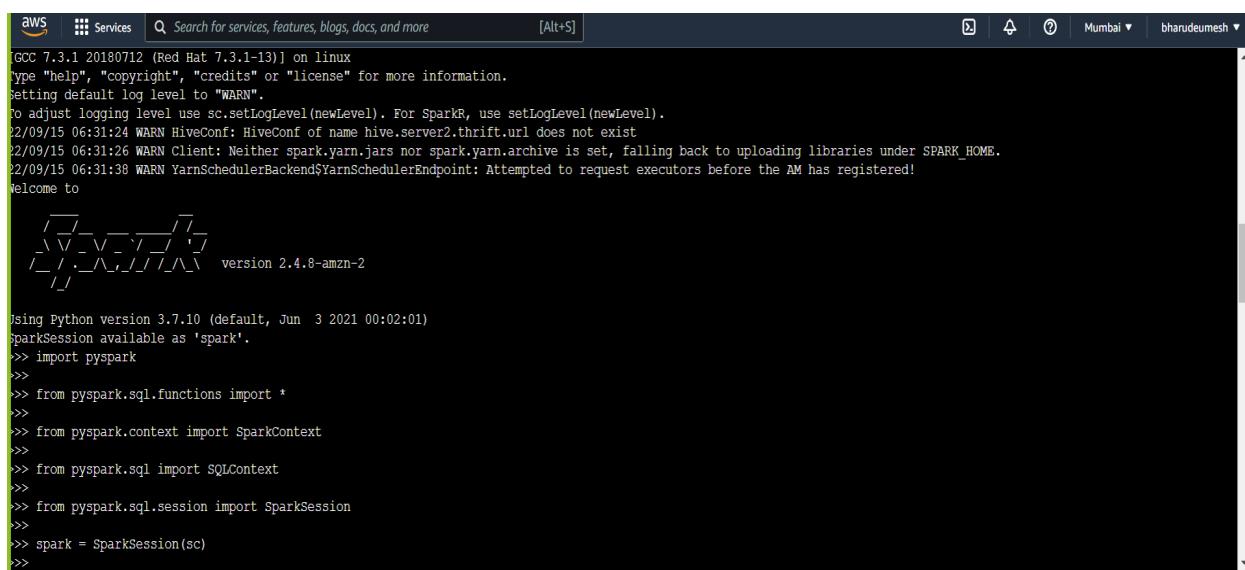
A screenshot of a terminal window titled "ap-south-1.console.aws.amazon.com/ec2-instance-connect/ssh?connType=standard&instanceId=i-0ad6c5ee8f8fe1789&osUser=hadoop®ion=ap-south-1&sshPort=22#". The terminal shows a series of ASCII art patterns (M's and R's) followed by a command-line session:

```
[hadoop@ip-172-31-47-187 ~]$ pyspark
Python 3.7.10 (default, Jun 3 2021 00:02:01)
[GCC 7.3.1 20180712 (Red Hat 7.3.1-13)] on linux
Type "help", "copyright", "credits" or "license" for more information.
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
22/09/21 10:43:04 WARN HiveConf: HiveConf of name hive.server2.thrift.url does not exist
22/09/21 10:43:06 WARN Client: Neither spark.yarn.jars nor spark.yarn.archive is set, falling back to uploading libraries under SPARK_HOME.
Welcome to
   ____/ \
  / \ \_ \_ \_ \_ \_ \_ \
 /_ \_ / . \_ \_ \_ \_ \_ \_ \
   \_ / \
version 2.4.8-amzn-2

Using Python version 3.7.10 (default, Jun 3 2021 00:02:01)
SparkSession available as 'spark'.
>>>
```

Below the terminal, the AWS navigation bar is visible, along with a search bar and a sidebar menu.

5) we install the libraries in pyspark:



A screenshot of a terminal window titled "ap-south-1.console.aws.amazon.com/ec2-instance-connect/ssh?connType=standard&instanceId=i-0ad6c5ee8f8fe1789&osUser=hadoop®ion=ap-south-1&sshPort=22#". The terminal shows a command-line session with library imports:

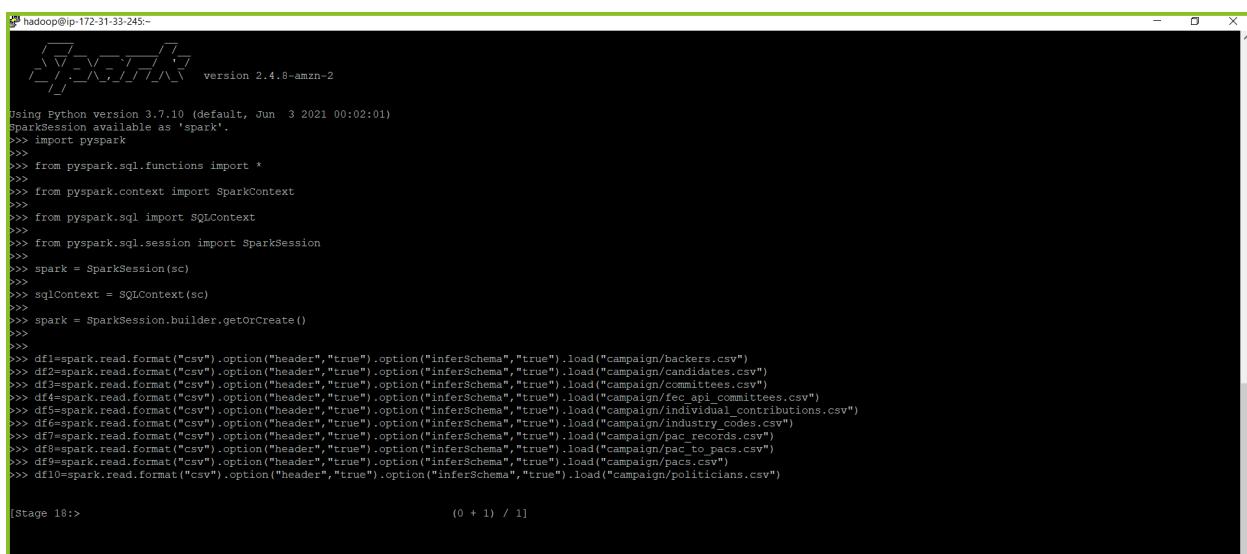
```
[GCC 7.3.1 20180712 (Red Hat 7.3.1-13)] on linux
Type "help", "copyright", "credits" or "license" for more information.
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
22/09/15 06:31:24 WARN HiveConf: HiveConf of name hive.server2.thrift.url does not exist
22/09/15 06:31:26 WARN Client: Neither spark.yarn.jars nor spark.yarn.archive is set, falling back to uploading libraries under SPARK_HOME.
22/09/15 06:31:38 WARN YarnSchedulerBackend$YarnSchedulerEndpoint: Attempted to request executors before the AM has registered!
Welcome to
   ____/ \
  / \ \_ \_ \_ \_ \_ \_ \
 /_ \_ / . \_ \_ \_ \_ \_ \_ \
   \_ / \
version 2.4.8-amzn-2

Using Python version 3.7.10 (default, Jun 3 2021 00:02:01)
SparkSession available as 'spark'.
>>> import pyspark
>>>
>>> from pyspark.sql.functions import *
>>>
>>> from pyspark.context import SparkContext
>>>
>>> from pyspark.sql import SQLContext
>>>
>>> from pyspark.sql.session import SparkSession
>>>
>>> spark = SparkSession(sc)
>>>
```

6)Read multiple csv from hdfs to spark

##creating dataframe :

```
df1=spark.read.format("csv").option("header","true").option("inferSchema","true")
).load("campaign/backers.csv")
df2=spark.read.format("csv").option("header","true").option("inferSchema","true")
).load("campaign/candidates.csv")
df3=spark.read.format("csv").option("header","true").option("inferSchema","true")
).load("campaign/committees.csv")
df4=spark.read.format("csv").option("header","true").option("inferSchema","true")
).load("campaign/fec_api_committees.csv")
df5=spark.read.format("csv").option("header","true").option("inferSchema","true")
).load("campaign/individual_contributions.csv")
df6=spark.read.format("csv").option("header","true").option("inferSchema","true")
).load("campaign/industry_codes.csv")
df7=spark.read.format("csv").option("header","true").option("inferSchema","true")
).load("campaign/pac_records.csv")
df8=spark.read.format("csv").option("header","true").option("inferSchema","true")
).load("campaign/pac_to_pacs.csv")
df9=spark.read.format("csv").option("header","true").option("inferSchema","true")
).load("campaign/pacs.csv")
df10=spark.read.format("csv").option("header","true").option("inferSchema","true")
").load("campaign/politicians.csv")
```



The screenshot shows a terminal window with the following content:

```

hadoop@ip-172-31-33-245:~$ 
version 2.4.8-amzn-2

Using Python version 3.7.10 (default, Jun  3 2021 00:02:01)
SparkSession available as 'spark'.
>>> import pyspark
>>>
>>> from pyspark.sql.functions import *
>>>
>>> from pyspark.context import SparkContext
>>> from pyspark.sql import SQLContext
>>> from pyspark.sql.session import SparkSession
>>> spark = SparkSession(sc)
>>> sqlContext = SQLContext(sc)
>>> spark = SparkSession.builder.getOrCreate()
>>>

>>> df1=spark.read.format("csv").option("header","true").option("inferSchema","true").load("campaign/backers.csv")
>>> df2=spark.read.format("csv").option("header","true").option("inferSchema","true").load("campaign/candidates.csv")
>>> df3=spark.read.format("csv").option("header","true").option("inferSchema","true").load("campaign/committees.csv")
>>> df4=spark.read.format("csv").option("header","true").option("inferSchema","true").load("campaign/individual_contributions.csv")
>>> df5=spark.read.format("csv").option("header","true").option("inferSchema","true").load("campaign/industry_codes.csv")
>>> df6=spark.read.format("csv").option("header","true").option("inferSchema","true").load("campaign/pac_records.csv")
>>> df8=spark.read.format("csv").option("header","true").option("inferSchema","true").load("campaign/pac_to_pacs.csv")
>>> df9=spark.read.format("csv").option("header","true").option("inferSchema","true").load("campaign/pacs.csv")
>>> df10=spark.read.format("csv").option("header","true").option("inferSchema","true").load("campaign/politicians.csv")

Stage 10:                                     (0 + 1) / 11

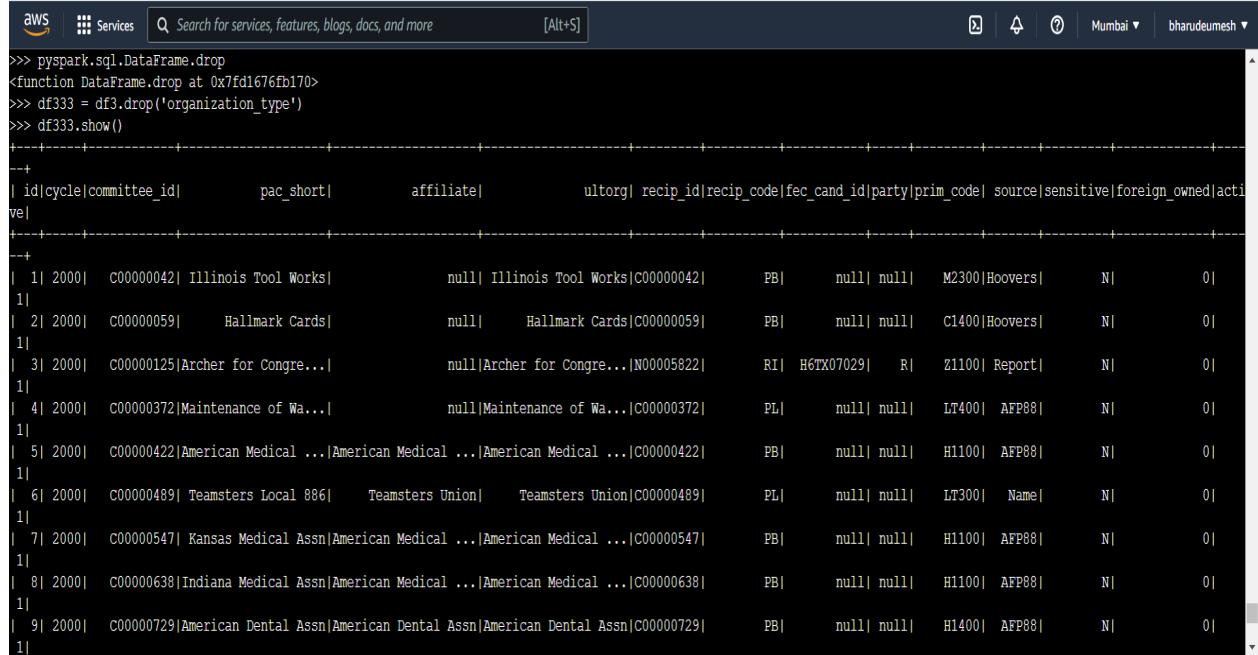
```

##creating temp view of tables:

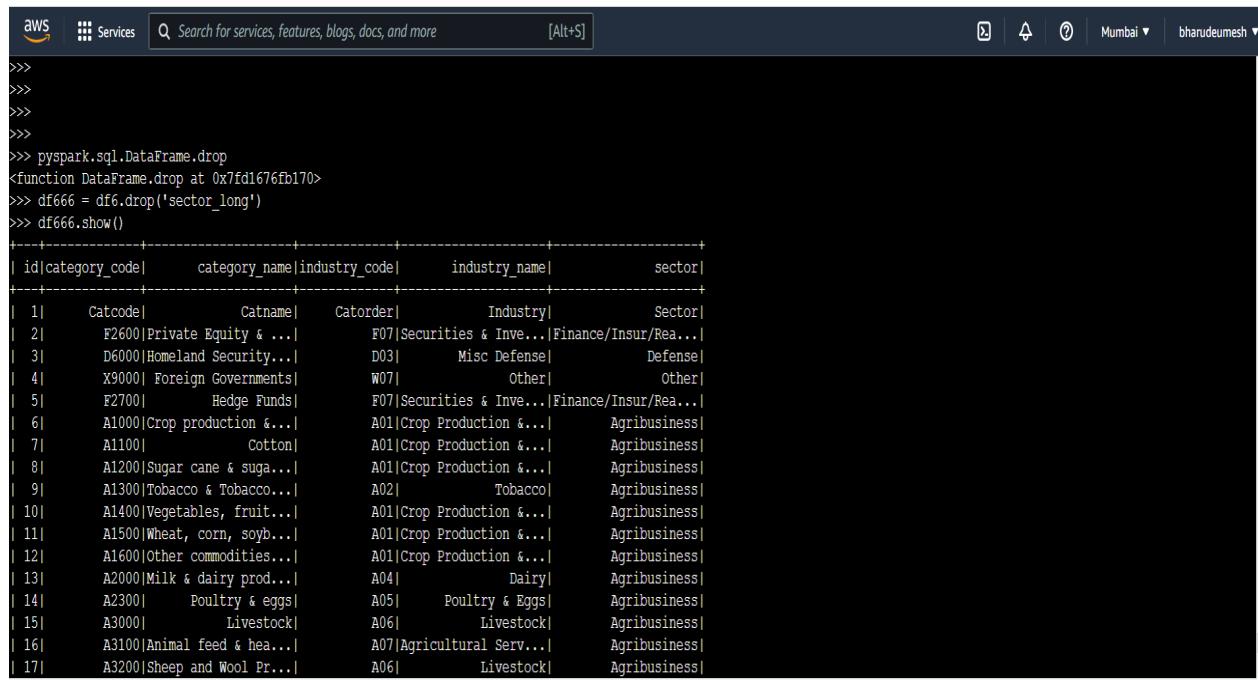
```
>>>
>>>
>>> df1.createOrReplaceTempView("backers")
>>> df2.createOrReplaceTempView("candidates")
>>> df3.createOrReplaceTempView("committees")
>>> df4.createOrReplaceTempView("fec_api_committees")
>>> df5.createOrReplaceTempView("individual_contributions")
>>> df6.createOrReplaceTempView("industry_codes")
>>> df7.createOrReplaceTempView("pac_records")
>>> df8.createOrReplaceTempView("pac_to_pacs")
>>> df9.createOrReplaceTempView("pacs")
>>> df10.createOrReplaceTempView("politicians")
>>> [REDACTED]
```

```
>>>
>>> df1=spark.read.format("csv").option("header","true").option("inferSchema","true").load("campaign/backers.csv")
>>>
>>>
>>> df1.createOrReplaceTempView("backers")
>>>
>>> spark.sql("select * from backers").show(10)
+---+-----+-----+-----+-----+
| id|    cid|          name|backer level|kickstarter|
+---+-----+-----+-----+-----+
| 1|N00003389|   Karan Dodia| leadership|      t|
| 2|N00009922|   Neal Foard| leadership|      t|
| 3|N00000528|   Kevin Pine|    senate|      t|
| 4|N00001093| Adam Edwards|    senate|      t|
| 5|N00007364|"Michael ""Pinski...|    senate|      t|
| 6|N00007364| Trevor von Stein|    senate|      t|
| 7|N00027605|   Susan Estep|    senate|      t|
| 8|N00002097| Wesley W. Driver|    senate|      t|
| 9|N00029303|   Betsy Wentzel|    senate|      t|
| 10|N00034580|   Chuck Dye|    senate|      t|
+---+-----+-----+-----+-----+
only showing top 10 rows
```

7) we clean the data ,we remove the null column:



```
>>> pyspark.sql.DataFrame.drop
<function DataFrame.drop at 0x7fd1676fb170>
>>> df333 = df3.drop('organization_type')
>>> df333.show()
+-----+-----+-----+-----+-----+-----+-----+-----+
| id|cycle|committee_id| pac_short| affiliate| ultorg| recip_id|recip_code|fec_cand_id|party|prim_code| source|sensitive|foreign_owned|acti
+-----+-----+-----+-----+-----+-----+-----+-----+
| 1| 2000| C00000042| Illinois Tool Works| null| Illinois Tool Works|C00000042| PB| null| null| M2300|Hoovers| N| 0|
| 2| 2000| C00000059| Hallmark Cards| null| Hallmark Cards|C00000059| PB| null| null| C1400|Hoovers| N| 0|
| 3| 2000| C00000125|Archer for Congre...| null|Archer for Congre...|N00005822| RI| H6TX07029| R| Z1100| Report| N| 0|
| 4| 2000| C00000372|Maintenance of Wa...| null|Maintenance of Wa...|C00000372| PL| null| null| LT400| AFP88| N| 0|
| 5| 2000| C00000422|American Medical ...|American Medical ...|American Medical ...|C00000422| PB| null| null| H1100| AFP88| N| 0|
| 6| 2000| C00000489| Teamsters Local 886| Teamsters Union| Teamsters Union|C00000489| PL| null| null| LT300| Name| N| 0|
| 7| 2000| C00000547| Kansas Medical Assn|American Medical ...|American Medical ...|C00000547| PB| null| null| H1100| AFP88| N| 0|
| 8| 2000| C00000638|Indiana Medical Assn|American Medical ...|American Medical ...|C00000638| PB| null| null| H1100| AFP88| N| 0|
| 9| 2000| C00000729|American Dental Assn|American Dental Assn|American Dental Assn|C00000729| PB| null| null| H1400| AFP88| N| 0|
+-----+
```



```
>>>
>>>
>>>
>>> pyspark.sql.DataFrame.drop
<function DataFrame.drop at 0x7fd1676fb170>
>>> df666 = df6.drop('sector_long')
>>> df666.show()
+-----+-----+-----+-----+
| id|category_code| category_name|industry_code| industry_name| sector|
+-----+-----+-----+-----+
| 1| Catcode| Catname| Catorder| Industry| Sector|
| 2| F2600|Private Equity & ...| F07|Securities & Inve...|Finance/Insur/Rea...
| 3| D6000|Homeland Security...| D03| Misc Defense| Defense|
| 4| X9000| Foreign Governments| W07| Other| Other|
| 5| F2700| Hedge Funds| F07|Securities & Inve...|Finance/Insur/Rea...
| 6| A1000|Crop production &...| A01|Crop Production &...| Agribusiness|
| 7| A1100| Cotton| A01|Crop Production &...| Agribusiness|
| 8| A1200|Sugar cane & suga...| A01|Crop Production &...| Agribusiness|
| 9| A1300|Tobacco & Tobacco...| A02| Tobacco| Agribusiness|
| 10| A1400|Vegetables, fruit...| A01|Crop Production &...| Agribusiness|
| 11| A1500|Wheat, corn, soy...| A01|Crop Production &...| Agribusiness|
| 12| A1600|Other commodities...| A01|Crop Production &...| Agribusiness|
| 13| A2000|Milk & dairy prod...| A04| Dairy| Agribusiness|
| 14| A2300| Poultry & eggs| A05| Poultry & Eggs| Agribusiness|
| 15| A3000| Livestock| A06| Livestock| Agribusiness|
| 16| A3100|Animal feed & hea...| A07|Agricultural Serv...| Agribusiness|
| 17| A3200|Sheep and Wool Pr...| A06| Livestock| Agribusiness|
+-----+
```

US Federal Campaign Finance Data Analysis Using BigData Technologies

```

aws Services Search for services, features, blogs, docs, and more [Alt+S]
>>> pyspark.sql.DataFrame.drop
<function DataFrame.drop at 0x7fd1676fb170>
>>> df888 = df8.drop('fec_occ_emp')
>>> df888.show()
+---+---+---+---+---+---+---+---+---+---+
| id|cycle|fec_rec_no| filer_id| donor_committee| contrib_lend_trans| city|state| zip|prim_code| date| amount|recipient_id|party| other_id|recip_c
ode|recip_prim_code|amend|report| pg| microfilm|type|real_code|source|
+---+---+---+---+---+---+---+---+---+---+
| 1| 2000| 1430708|C00000059| Hallmark Cards| null| null| | | C1400|2000-06-08 00:00:00| 1000.0| C00341867| null|C00341867| | | | | | |
| PI| J1200| A| N7 | P|2003598313| 24K| C1400| PAC | null| | | null| | | LT400|1999-03-16 00:00:00| 5000.0| C00190876| D|C00190876|
| 2| 2000| 62108|C00000372|Maintenance of Wa...| null| null| | | null| | | null| | | LT400|1999-07-31 00:00:00| 1000.0| C00147512| null|C00147512|
| PI| J2100| N1| M4 | P|99034461953| 24K| LT400| PAC | null| | | null| | | LT400|1999-07-31 00:00:00| 1000.0| C00147512| null|C00147512|
| 3| 2000| 397553|C00000372|Maintenance of Wa...| null| null| | | null| | | null| | | LT400|2000-06-08 00:00:00|35000.0| C00010603| D|C00010603|
| PI| J7500| N1| M8 | P|99034784637| 24K| LT400| PAC | null| | | null| | | LT400|2000-06-08 00:00:00|35000.0| C00010603| D|C00010603|
| 4| 2000| 1325828|C00000372|Maintenance of Wa...| null| null| | | null| | | null| | | LT400|2000-05-09 00:00:00|15000.0| C00010603| D|C00010603|
| DP| Z5200| N1| M7 | P|20035881685| 24K| LT400| PAC | null| | | null| | | LT400|2000-05-09 00:00:00|15000.0| C00010603| D|C00010603|
| 5| 2000| 1146337|C00000372|Maintenance of Wa...| null| null| | | null| | | null| | | LT400|2000-05-09 00:00:00|15000.0| C00010603| D|C00010603|
| DP| Z5200| A| M6 | P|20035732124| 24K| LT400| PAC | null| | | null| | | LT400|1999-11-21 00:00:00| 1092.0| C00010603| D|C00010603|
| 6| 2000| 825273|C00000372|Maintenance of Wa...| null| null| | | null| | | null| | | LT400|1999-11-21 00:00:00| 1092.0| C00010603| D|C00010603|
| DP| Z5200| A| M12| P|20035043806| 24Z| LT400| PAC | null| | | null| | | LT400|1999-11-21 00:00:00| 1092.0| C00010603| D|C00010603|
| 7| 2000| 1491090|C00000372|Maintenance of Wa...|SOUTH DAKOTA DEMO...| SIOUX FALLS| SD|57101| LT400|2000-10-05 00:00:00| 1000.0| C00160937| D|C00160937|
| DP| Z5200| N1| 12G| G|20990227762| 24K| LT400| PAC | null| | | null| | | LT400|2000-10-05 00:00:00| 1000.0| C00160937| D|C00160937|
| 8| 2000| 79737|C00000422|American Medical ...| null| null| | | null| | | null| | | H1100|1999-04-28 00:00:00|15000.0| C00042366| D|C00042366|
| DP| Z5200| N1| M5 | P|99034493582| 24K| H1100| PAC | null| | | null| | | H1100|1999-04-28 00:00:00|15000.0| C00042366| D|C00042366|
| 9| 2000| 1068453|C00000372|Maintenance of Wa...| null| null| | | null| | | null| | | LT400|2000-04-13 00:00:00| 550.0| C00160937| D|C00160937|
| DP| Z5200| N1| M5 | P|20035582436| 24K| LT400| PAC | null| | | null| | | LT400|2000-04-13 00:00:00| 550.0| C00160937| D|C00160937|

```



```

aws Services Search for services, features, blogs, docs, and more [Alt+S]
>>> pyspark.sql.DataFrame.drop
<function DataFrame.drop at 0x7fd1676fb170>
>>> df555 = df5.drop('ult_org')
>>> df555.show()
+---+---+---+---+---+---+---+---+---+---+
| id|cycle| fec_trans_id|contributor_id| contributor_name|recipient_id| org_name|real_code| date|amount|street| city|state| zip|re
cip_code|type|committee_id| other_id|gender|old_format_employer_occupation| microfilm| occupation| employer|source|
+---+---+---+---+---+---+---+---+---+---+
| 1| 2010|1010420120009248286| k00019100751| William G Crook| N00009638| Self-Employed| G0000|2011-03-01 00:00:00| 200| null| New York| NY|10022| | | | | |
| DP| 15 | C00431445| | | null|11931873404| SELF| SELF| GEN | null| | | null| | | null| | |
| 2| 2012|1010120130011206490| k0001544935| JONES, BERNHARDT N| C00042366| Retired| X1200|2012-10-05 00:00:00| 250| null| EUGENE| OR|97404|
| DP| 15 | C00042366| null| N| null|12021033713| RETIRED| N/A| Gen | null| | | null| | |
| 3| 2012|1010120130011206620| h30019228431| YAGGY, DUNCAN| C00042366| Yaggy Corp| Y4000|2012-10-16 00:00:00| 500| null| DURHAM| NC|27705|
| DP| 15E| C00042366|C00401224| M| null|12021053305| MANAGER| YAGGY CORPORATION| | null| | | null| | |
| 4| 2012|1010120130011206706| i3003429059| KELETI, STEVEN| C00042366| Judge Technical S...| Y4000|2012-10-17 00:00:00| 1000| null| MALDEN| MA|02148|
| DP| 15E| C00042366|C00401224| M| null|12021044111| ENGINEER|JUDGE TECHNICAL S...| | null| | | null| | |
| 5| 2012|1010120130011206819| m00021679750| HAWTHORNE, JOAN| C00042366|University of Nor...| H5100|2012-10-11 00:00:00| 500| null| GRAND FORKS| ND|58201|
| DP| 15E| C00042366|C00401224| F| null|12021043352| EDUCATOR|UNIVERSITY OF NOR...| Name | null| | | null| | |
| 6| 2012|1010120130011207076| h3001884327| HEDELMAN, SIDNEY| C00042366| Retired| X1200|2012-10-04 00:00:00| 200| null| MILL VALLEY| CA|94941|
| DP| 15E| C00042366|C00401224| N| null|12021043390| NOT EMPLOYED| N/A| Gen | null| | | null| | |
| 7| 2012|1010120130011207228| h30017191628| PRICE, LINDA N| C00042366| Retired| X1200|2012-10-17 00:00:00| 500| null| NAPA| CA|94558|
| DP| 15E| C00042366|C00401224| F| null|12021051206| RETIRED| N/A| Gen | null| | | null| | |
| 8| 2012|1010120130011207256| h3001636082| CORBITT, CRAIG| C00042366|Zelle, Hofmann et al| K1000|2012-10-17 00:00:00| 250| null| MILL VALLEY| CA|94941|
| DP| 15E| C00042366|C00401224| M| null|12021041752| ATTORNEY| ZELLE HOFMANN LLP Rept | null| | | null| | |
| 9| 2012|1010120130011207321| m0002107708| MCGAVICK, HUGH| C00042366| Attorney| K1000|2012-10-13 00:00:00| 200| null| PORTLAND| OR|97219|
| DP| 15E| C00042366|C00401224| M| null|12021050146| ATTORNEY| SELF-EMPLOYED| Gen | null| | | null| | |

```

US Federal Campaign Finance Data Analysis Using BigData Technologies

aws Services Search for services, features, blogs, docs, and more [Alt+S] Mumbai bharudeumesh

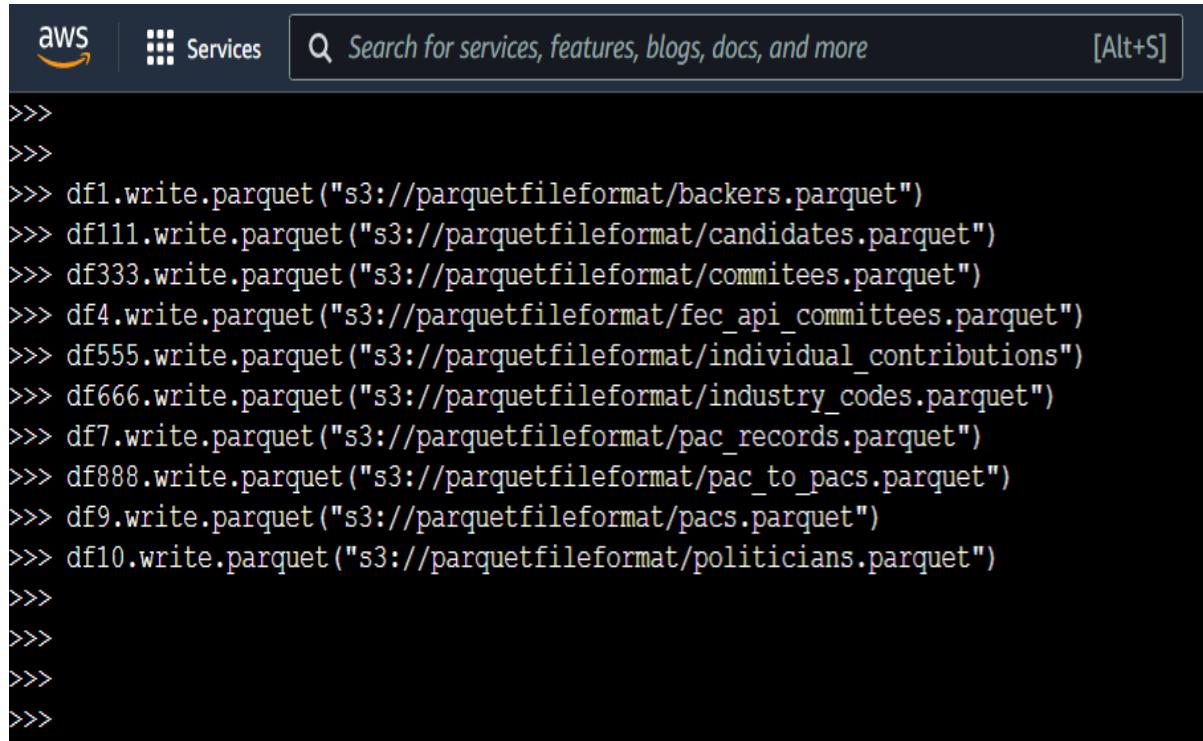
```
>>> pypark.sql.DataFrame.drop
>>> function DataFrame.drop at 0x7fd1676fb170>
>>> df555 = df555.drop('street')
>>> df555.show()
```

id	cycle	fec_trans_id	contributor_id	contributor_name	recipient_id	org_name	real_code	date	amount	city/state	zip	recip_code	
e_type	committee_id	other_id	gender	old_format_employer_occupation	microfilm	occupation		employer/source					
I	1	2010 1010420120009248286	k00019100751	William G Crook	N00009638	Self-Employed	G0000	2011-03-01 00:00:00	2001	New York	NY 110022	D	
E	15		C004314545		null 11931873404	SELF	SELF GEN						
I	2	2011 1010120130011206490	k0001544935	I JONES, BERNHARDT	NI	C000423661	Retired	X1200 2012-10-05 00:00:00	2501	EUGENE	OR 97404	D	
P	I	31	2011 1010120130011206620	h30019228431	YAGGY, DUNCAN	C000423661	RETIRED	N/A Gen					
P	I	15E	2011 1010120130011206706	C000423661 C004012241	MJ	null 12021053305	Yaggy Corp	Y4000 2012-10-16 00:00:00	5001	DURHAM	NC 27705	D	
I	41	2011 1010120130011206706	i3003429059	KELETI, STEVENI	C000423661 Judge Technical S...	JUDGE TECHNICAL S...	MANAGER	YAGGY CORPORATION	Y4000 2012-10-17 00:00:00	10001	MALDEN	MA 02148	D
P	I	15E	2011 1010120130011206819	C000423661 C004012241	MJ	null 12021044111	ENGINEER JUDGE TECHNICAL S...						
I	51	2011 1010120130011206819	m00021679758	HAWTHORNE, JOAN	C000423661 University of Nor...	University of Nor...	F1	H5100 2012-10-11 00:00:00	5001	GRAND FORKS	ND 58201	D	
P	I	15E	2011 1010120130011207076	C000423661 C004012241	F1	null 12021043352	EDUCATOR UNIVERSITY OF NOR... Name						
I	61	2012 1010120130011207076	h3001884327	HEDELMAN, SIDNEY	NI	C000423661	Retired	X1200 2012-10-04 00:00:00	2001	MILL VALLEY	CA 94941	D	
P	I	15E	2012 1010120130011207228	C000423661 C004012241	NI	null 12021043390	NOT EMPLOYED	N/A Gen					
I	71	2012 1010120130011207228	h30017191628	PRICE, LINDA	NI	C000423661	Retired	X1200 2012-10-17 00:00:00	5001	NAPA	CA 94558	D	
P	I	15E	2012 1010120130011207256	C000423661 C004012241	F1	null 120210120512061	RETIRED	N/A Gen					
I	81	2012 1010120130011207256	h3001636082	CORBITT, CRAIG	Zelle, Hofmann et al	C000423661 Zelle, Hofmann et al	K1000 2012-10-12 00:00:00	2501	MILL VALLEY	CA 94941	D		
P	I	15E	2012 1010120130011207321	C000423661 C004012241	MJ	null 12021014752	ATTORNEY	ZELLE HOFMANN LLP Rept	K1000 2012-10-13 00:00:00	2001	PORTLAND	OR 97219	D
I	91	2012 1010120130011207708	m0002107708	MCGAVICK, HUGH	C000423661	Attorney							

We renamed the column:

```
>>>
>>>
>>>
>>> df888.withColumnRenamed("pg","primary or general").show()
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| id|cycle|fec_rec_no|filer_id|donor_committee|contrib_lend_trans|city_state|zip|prim_code|date|amount|recipient_id|party|other_id|recip_c
ode|recip_prim_code|amend|report|primary_or_general|microfilm|type|real_code|source|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| 1| 2000|1430708|C00000059| Hallmark Cards| null| null| | C1400|2000-06-08 00:00:00| 1000.0| C00341867| null|C00341867| | |
| PII| J12001| AI_MT | P|2003598313| 24K| C1400| PAC | | | LT400|1999-03-16 00:00:00| 5000.0| C00190876| D|C00190876|
| PII| J21001| NI_MA | P|9903446195| 24K| LT400| PAC | | | LT400|1999-07-31 00:00:00| 1000.0| C00147512| null|C00147512|
| PII| 397553|C00000372|Maintenance of Wa...| null| null| null| | | LT400|2000-06-08 00:00:00|35000.0| C00010603| D|C00010603|
| PII| J75001| NI_M6 | P|9903478463| 24K| LT400| PAC | | | LT400|2000-05-09 00:00:00|15000.0| C00010603| D|C00010603|
| DPI| 1325828|C00000372|Maintenance of Wa...| null| null| null| | | LT400|1999-11-21 00:00:00| 1092.0| C00010603| D|C00010603|
| DPI| 252001| NI_M7 | P|20035881685| 24K| LT400| PAC | | | LT400|1999-07-31 00:00:00| 1000.0| C00147512| null|C00147512|
| PII| 51| 2000|114633|C00000372|Maintenance of Wa...| null| null| | | LT400|2000-05-09 00:00:00|15000.0| C00010603| D|C00010603|
| DPI| 252001| AI_M6 | P|20035732124| 24K| LT400| PAC | | | LT400|1999-11-21 00:00:00| 1092.0| C00010603| D|C00010603|
| DPI| 61| 2000|825273|C00000372|Maintenance of Wa...| null| null| | | LT400|1999-11-21 00:00:00| 1092.0| C00010603| D|C00010603|
| DPI| 252001| AI_M12| P|20035043806| 24Z| LT400| PAC | | | LT400|1999-11-21 00:00:00| 1092.0| C00010603| D|C00010603|
| PII| 71| 2000|1491090|C00000372|Maintenance of Wa...|SOUTH DAKOTA DEMO...|SIOUX FALLS| SD|57101| LT400|2000-10-05 00:00:00| 1000.0| C00160937| D|C00160937|
| DPI| 252001| NI_12G| G|209909227762| 24K| LT400| PAC | | | LT400|2000-10-05 00:00:00| 1000.0| C00160937| D|C00160937|
| PII| 81| 2000|79737|C00000422|American Medical ...| null| null| | | H1100|1999-04-28 00:00:00|15000.0| C00042366| D|C00042366|
| DPI| 252001| NI_M5 | P|9903449358| 24K| H1100| PAC | | | H1100|1999-04-28 00:00:00|15000.0| C00042366| D|C00042366|
```

8)Save the cleaned data in S3 Bucket into parquet format:



```
>>>
>>>
>>> df1.write.parquet("s3://parquetfileformat/backers.parquet")
>>> df111.write.parquet("s3://parquetfileformat/candidates.parquet")
>>> df333.write.parquet("s3://parquetfileformat/commitees.parquet")
>>> df4.write.parquet("s3://parquetfileformat/fec_api_committees.parquet")
>>> df555.write.parquet("s3://parquetfileformat/individual_contributions")
>>> df666.write.parquet("s3://parquetfileformat/industry_codes.parquet")
>>> df7.write.parquet("s3://parquetfileformat/pac_records.parquet")
>>> df888.write.parquet("s3://parquetfileformat/pac_to_pacs.parquet")
>>> df9.write.parquet("s3://parquetfileformat/pacs.parquet")
>>> df10.write.parquet("s3://parquetfileformat/politicians.parquet")
>>>
>>>
>>>
>>>
```

9)Load the data from S3 to pyspark:



```
>>>
>>> df1=spark.read.parquet("s3://parquetfileformat/backers.parquet")
>>> df2=spark.read.parquet("s3://parquetfileformat/candidates.parquet")
>>> df3=spark.read.parquet("s3://parquetfileformat/commitees.parquet")
>>> df4=spark.read.parquet("s3://parquetfileformat/fec_api_committees.parquet")
>>> df5=spark.read.parquet("s3://parquetfileformat/individual_contributions")
>>> df6=spark.read.parquet("s3://parquetfileformat/industry_codes.parquet")
>>> df7=spark.read.parquet("s3://parquetfileformat/pac_records.parquet")
>>> df8=spark.read.parquet("s3://parquetfileformat/pac_to_pacs.parquet")
>>> df9=spark.read.parquet("s3://parquetfileformat/pacs.parquet")
>>> df10=spark.read.parquet("s3://parquetfileformat/politicians.parquet")
>>>
>>>
>>>
```

```
>>> df1.count()
370
>>> df2.count()
67999
>>> df3.count()
157542
>>> df4.count()
45507
>>> df5.count()
24446421
>>> df6.count()
444
>>> df7.count()
44628
>>> df8.count()
1083525
>>> df9.count()
3539657
>>> df10.count()
26047
```

10)Create table and perform query in pyspark and store table in hive :

US Federal Campaign Finance Data Analysis Using BigData Technologies

```

aws Services Q Search for services, features, blogs, docs, and more [Alt+S]
>>> sqlContext = SQLContext(sc)
>>>
>>> spark = SparkSession.builder.getOrCreate()
>>>
>>> spark.sql("create table senate select name,backer_level from backers where backer_level = 'senate'")
22/09/18 17:50:01 WARN HiveConf: HiveConf of name hive.server2.thrift.url does not exist
DataFrame[]
>>> spark.sql("select * from senate").show()
+-----+-----+
|      name|backer_level|
+-----+-----+
| Kevin Pine|    senate|
| Adam Edwards|    senate|
| "Michael ""Pinsky...|    senate|
| Trevor von Stein|    senate|
| Susan Estep|    senate|
| Wesley W. Driver|    senate|
| Betsy Wentzel|    senate|
| Chuck Dye|    senate|
| David Cook|    senate|
| David Powell|    senate|
| Gena Mason|    senate|
| Ken Ledenbach|    senate|
| Micah Sifry|    senate|
| Milind Shah|    senate|
|Please do not use...|    senate|
| R.A.Plunkett|    senate|
| Roy Clymer|    senate|
+-----+-----+

```

```

>>> spark.sql("create table individual_contributions_details1 select contributor_name,gender,org_name,city,state,occupation from individual_contributions")
DataFrame[]
>>> spark.sql("select * from individual_contributions_details1").show()
+-----+-----+-----+-----+-----+
|contributor_name|gender|org_name|city|state|occupation|
+-----+-----+-----+-----+-----+
| BURSON, JOHN MR| M| Shapiro Burson LLP| ARLINGTON| VA| Attorney|
| CANDEAUX, PAT| N| Dr R Rust| EUGENE| OR| Retired|
| CANON, JOSEPH EDWIN| M| Tejon Exploration| ABILENE| TX| null|
| SIMS, CHARLES| M| Proskauer Rose| NEW YORK| NY| null|
| YOUNG 236, REXFOR...| M| null| NEWPORT NEWS| VA| null|
| JACOBS, CHARLES E| M| CE Jacobs Co| ALBANY| TX| null|
| PLASTER, DONALD J...| M| Compass Technology| VIRGINIA BEACH| VA| Computer Analyst|
| FREELAND, DOUGLAS| M| null| WRIGHT| WY| Maintenance Manager|
| BARTLETT, GAIL A| N| null| OKLAHOMA CITY| OK| null|
| O'BRIAN, FERN| F| Thompson Hine LLP| BETHESDA| MD| lawyer|
| FRITSCHE, WILLIAM...| M| Retired| VIRGINIA BEACH| VA| Retired|
| DELANEY, ROBERT| M| null| PELICANVILLE| TX| null|
| BLOOM, DAVID| M| Bloom Electric| YUKON| OK| null|
| SACCO, JOHN| M| Sierra Internatio...| BAKERSFIELD| CA| Managing Member|
| LAWRENCE, RICHARD...| M| Lawrence Firm| CINCINNATI| OH| null|
| LAMB, HERSH| N| null| VANCOUVER| WA| null|
| DEW, BRIAN| M| null| NEW ALBANY| OH| Electrical Contra...|
+-----+-----+-----+-----+-----+

```

```
>>> spark.sql("create table countof_backerslevel select backer_level,count(backer_level) from backers group by backer_level")
DataFrame[]
>>> spark.sql("select * from countof_backerslevel").show()
+-----+-----+
| backer_level|count(backer_level)|
+-----+-----+
|      house|        48|
|     upstart|        50|
|kickstarter_backer|      248|
|      senate|        22|
|   leadership|         2|
+-----+-----+
>>>
```

spark.sql("create table cycleof_candidates select b.cid,c.cycle,b.name,party from backers b join candidates c where c.cycle=1996 and b.cid=c.cid")

spark.sql("select * from cycleof_candidates").show()

```
+-----+-----+-----+
| cid|cycle|       name|party|
+-----+-----+-----+
|N00005165| 1996|    Warren|  R|
|N00002384| 1996| Alex Cox|  D|
|N00008999| 1996| Jorge Pinon|  R|
|N00000561| 1996|   Ainslie|  R|
|N00002247| 1996|   Nala R.|  R|
|N00007335| 1996| Mark Menke|  D|
|N00004541| 1996| lynn adler|  D|
|N00005195| 1996|   Gene Mason|  R|
|N00003333| 1996| Andrew DiSabatino|  D|
|N00002097| 1996| Wesley W. Driver|  D|
|N00000528| 1996|   Kevin Pine|  I|
|N00001743| 1996|   Steve King|  D|
|N00007063| 1996| Nicole Lauber|  R|
|N00001750| 1996| Steve Tuckner|  D|
|N00009082| 1996|Please do not use...|  R|
|N00009082| 1996|Please do not use...|  R|
|N00006863| 1996| Happy Elder Geek|  R|
```

```
|N00001305| 1996| Michael Maitlen| D|
|N00007458| 1996| Donald Simon| D|
|N00001802| 1996| Harper Reed| R|
+-----+-----+-----+
```

only showing top 20 rows

```
spark.sql("select cycle,first_last_party,party from candidates where party = 'R'")
```

```
+-----+-----+
|cycle| first_last_party|party|
+-----+-----+
| 2000|Richmond A Soluad...| R|
| 1996|Mark Alan Behnke (R)| R|
| 1996|Stephen Bonsal Yo...| R|
| 1996|Raymond J Clatwor...| R|
| 1996| Jim Ford (R)| R|
| 1996| Ernest J Istook (R)| R|
| 1996|Stephen Wayne Hof...| R|
| 1996|Olivia Coggin Eud...| R|
| 1996|Michael Minoru Fa...| R|
| 1996| Joseph Dr Smith (R)| R|
| 1996| James A McClure (R)| R|
| 2000|John A Holmes Jr (R)| R|
| 1996|James T Robinson (R)| R|
| 1996|William Jack Sext...| R|
| 1996| Gen Olson (R)| R|
| 1996|Berwick (Barry) P...| R|
| 1996| Mark E Souder (R)| R|
| 1996|Robert L Hammock (R)| R|
| 1996| Luis Acle Jr (R)| R|
| 1996| Jeffrey Bell (R)| R|
+-----+-----+
```

```
spark.sql("select party,count(party) from candidates where party<100 group by party").show()
```

```
+-----+
|party|count(party)|
+-----+
|  3|    5297|
+-----+
```

```
spark.sql("select party,count(party) from candidates group by party")
```

```
+-----+-----+
| party|count(party)|
+-----+-----+
| Jr (3)"|      4|
| U|    1273|
| D|   25778|
| R|   27941|
| L|   3053|
| 3|   5297|
| I|   4653|
+-----+-----+
```

```
spark.sql("select party,count(party) from candidates group by party having party = 'R'").show()
```

```
+-----+
|party|count(party)|
+-----+
| R|    27941|
+-----+
```

Q.display the contributor_name whose state is CA, occupation is retired.

```
spark.sql("select contributor_name, state , occupation from individual_contributions  
where state='CA' and occupation='retired'")
```

Q.display the cycle,committee type, contributor_name, city from fec_api_committees and individual_contributions.

```
spark.sql("select i.contributor_name, f.cycle, f.committee_type, i.city from  
fec_api_committees f join individual_contributions i  
where f.cycle=i.cycle")
```

Q:Display category_code,category_name from industry_codes by grouping the industry_code,industry_name,sector

```
spark.sql("select category_code,category_name from industry_codes group by  
industry_code,industry_name,sector")
```

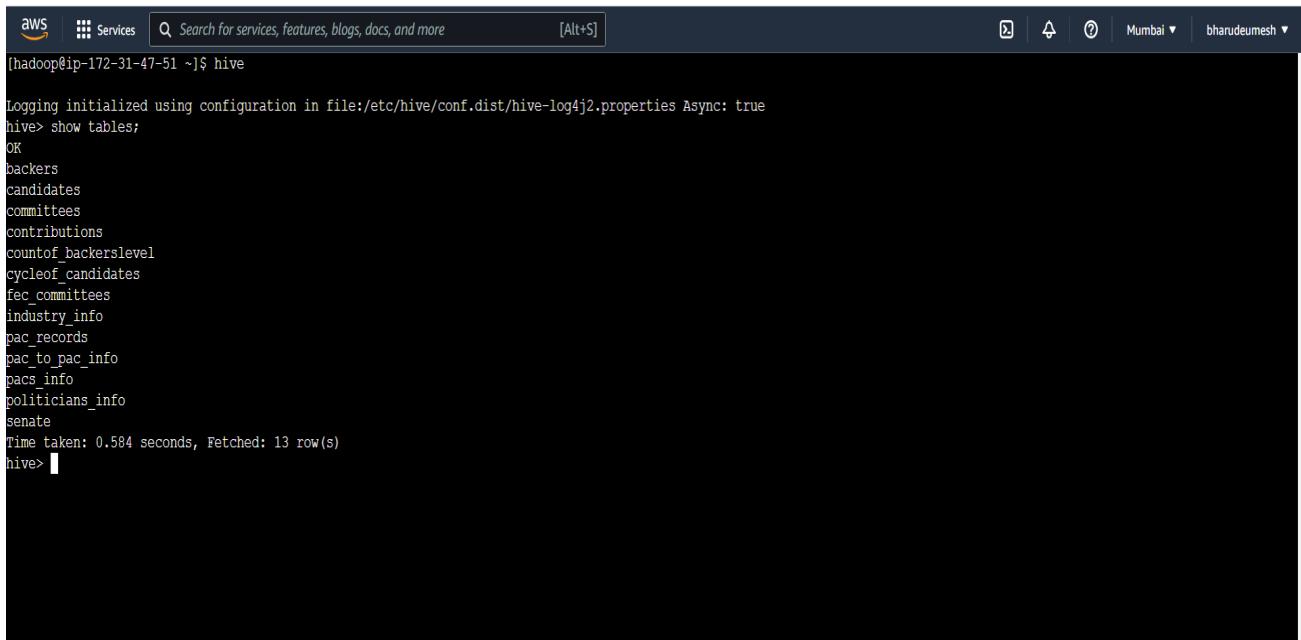
Q:which are the politicians contributed in different cycle

```
spark.sql("select p.cid,p.cycle,o.name as politician_name from pacs p join politicians o where p.cid=o.cid group by p.cycle")
```

```
>>> spark.sql("create table state_occupations select contributor_name, state , occupation from individual_contributions")
DataFrame[]
>>> spark.sql("select * from state_occupations").show()
+-----+-----+
| contributor_name|state|      occupation|
+-----+-----+
| BURSON, JOHN MR| VA|      Attorney|
| CANDEAUX, PAT| ORI|      Retired|
| CANON, JOSEPH EDWIN| TX|      null|
| SIMS, CHARLES| NY|      null|
| YOUNG 236, REXFOR...| VA|      null|
| JACOBS, CHARLES E| TX|      null|
| PLASTER, DONALD J...| VA| Computer Analyst|
| FREELAND, DOUGLAS| WY| Maintenance Manager|
| BARTLETT, GAIL A| OK|      null|
| O'BRIAN, FERN| MD|      lawyer|
| FRITSCHE, WILLIAM...| VA|      Retired|
| DELANEY, ROBERT| TX|      null|
| BLOOM, DAVID| OK|      null|
| SACCO, JOHN| CA| Managing Member|
| LAWRENCE, RICHARD...| OH|      null|
| LAMB, HERSHI| WA|      null|
| DEW, BRIAN| OH|Electrical Contra...|
| CHRISTENSEN, JAYNE| OK|      null|
| SANSARICQ, BERNARD| FL|      null|
| AYERS, HELEN MRS| IL|      null|
+-----+-----+
only showing top 20 rows
```

```
>>> spark.sql("select contributor_name, state , occupation from  
individual_contributions").show()  
+-----+-----+  
| contributor_name|state| occupation|  
+-----+-----+  
| BURSON, JOHN MR| VA| Attorney|  
| CANDEAUX, PAT| OR| Retired|  
| CANON, JOSEPH EDWIN| TX| null|  
| SIMS, CHARLES| NY| null|  
| YOUNG 236, REXFOR...| VA| null|  
| JACOBS, CHARLES E| TX| null|  
| PLASTER, DONALD J...| VA| Computer Analyst|  
| FREELAND, DOUGLAS| WY| Maintenance Manager|  
| BARTLETT, GAIL A| OK| null|  
| O'BRIAN, FERN| MD| lawyer|  
| FRITSCHE, WILLIAM...| VA| Retired|  
| DELANEY, ROBERT| TX| null|  
| BLOOM, DAVID| OK| null|  
| SACCO, JOHN| CA| Managing Member|  
| LAWRENCE, RICHARD...| OH| null|  
| LAMB, HERSH| WA| null|  
| DEW, BRIAN| OH|Electrical Contra...|  
| CHRISTENSEN, JAYNE| OK| null|  
| SANSARICQ, BERNARD| FL| null|  
| AYERS, HELEN MRS| IL| null|  
+-----+-----+  
only showing top 20 rows
```

9.Tables which are stored in hive:



A screenshot of a terminal window titled "aws Services". The search bar contains "Search for services, features, blogs, docs, and more [Alt+S]". The top right shows icons for a user profile, Mumbai, and bharudeumesh. The terminal output shows the results of a "show tables;" command in Hive:

```
[hadoop@ip-172-31-47-51 ~]$ hive
Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j2.properties Async: true
hive> show tables;
OK
backers
candidates
committees
contributions
countof_backerslevel
cycleof_candidates
fec_committees
industry_info
pac_records
pac_to_pac_info
pacs_info
politicians_info
senate
Time taken: 0.584 seconds, Fetched: 13 row(s)
hive>
```

10.connectivity to power bi:

From the created tables we are now connecting to power bi for better Visualization and understanding below are the queries with their respective visualization.

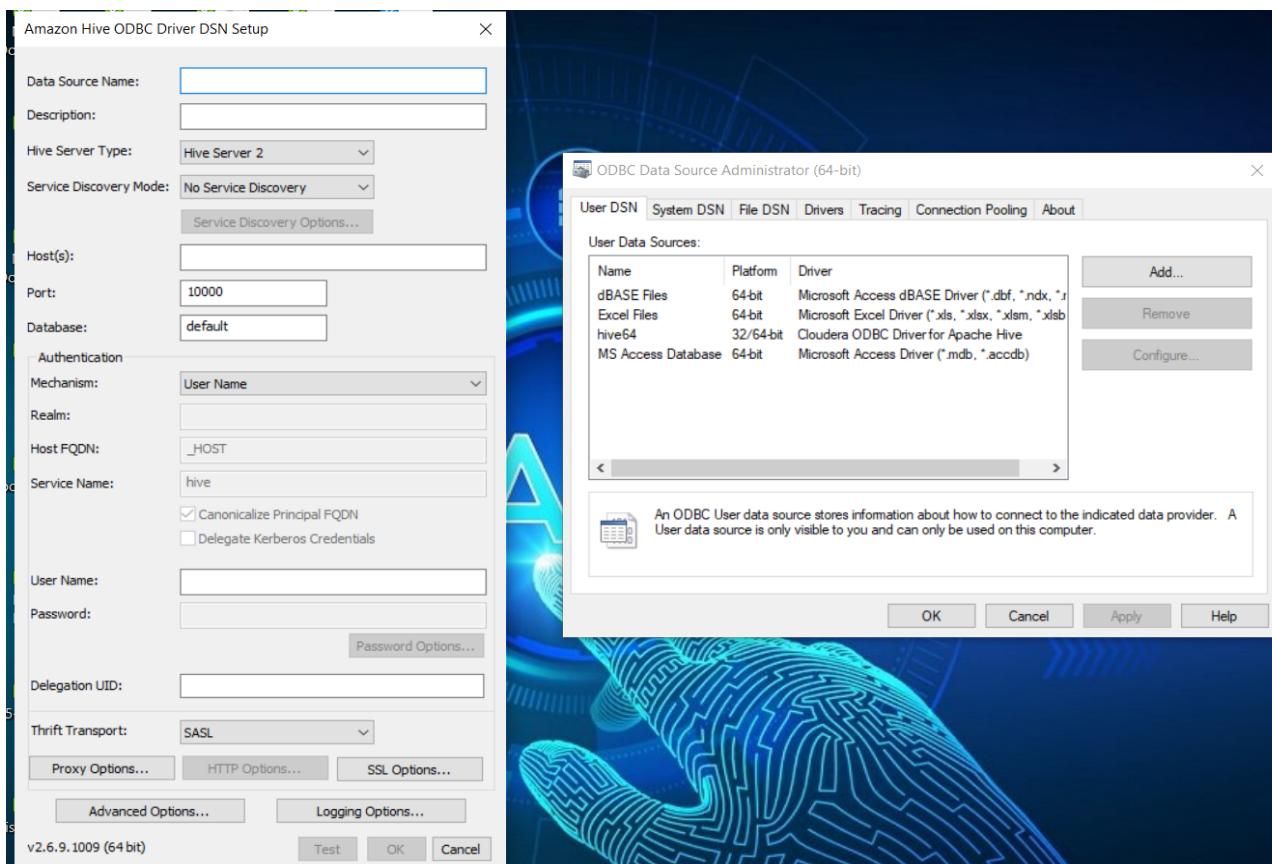
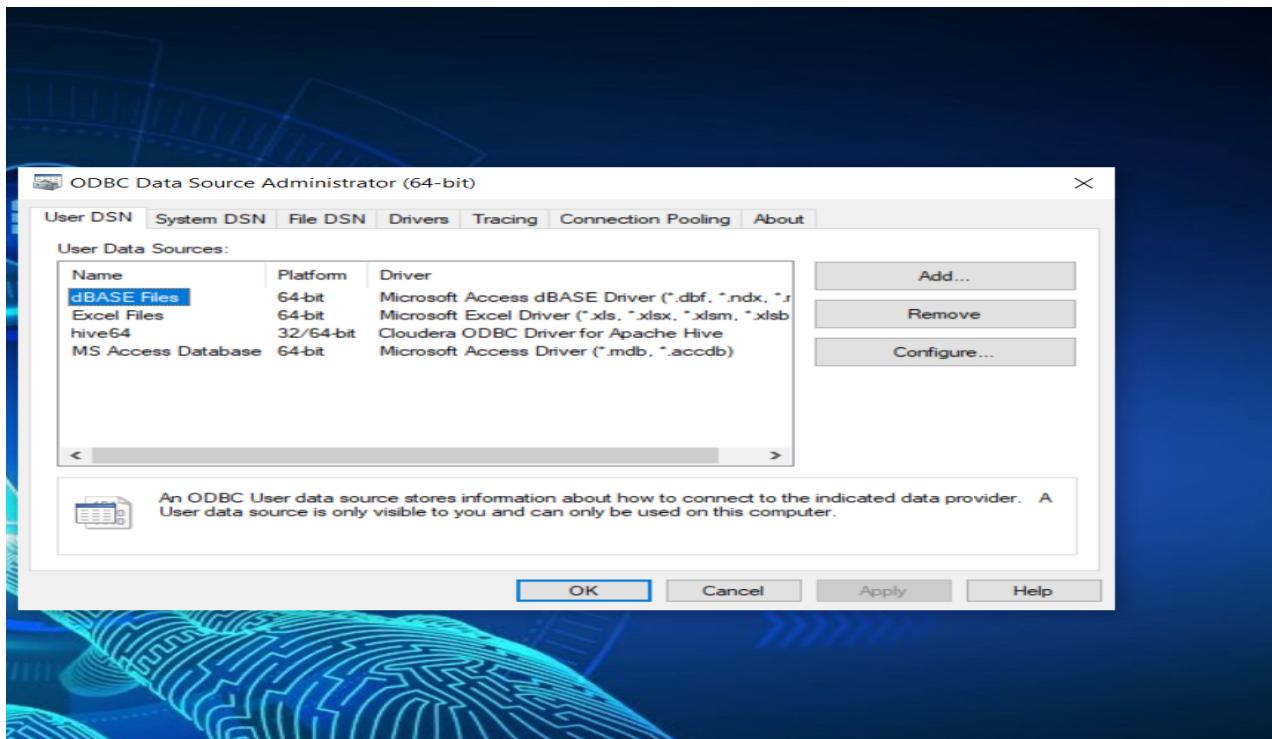
We connected power bi using ODBC DATA SOURCES by following steps

The screenshot shows the AWS EMR console with the search bar set to 'ec2'. On the left, there's a sidebar with navigation links for Amazon EMR, EMR Studio, EMR Serverless (with a 'New' button), EMR on EC2 (selected), Clusters, Notebooks, Git repositories, Security configurations, Block public access, VPC subnets, Events, and EMR on EKS. The main content area displays a cluster named 'projectemr' with the following details:

Name	ID	Status	Creation time (UTC+5:30)	Elapsed time	Normalized instance hours
projectemr	j-1J4FTPXMYH8CG	Waiting Cluster ready	2022-09-23 20:53 (UTC+5:30)	1 hour	0

Below the table, there's a 'Summary' section with details about the master public DNS (ec2-13-233-254-26.ap-south-1.compute.amazonaws.com) and termination protection (Off). It also shows tags and bootstrap actions. The 'Hardware' section indicates a single master node (m5.xlarge). At the bottom, there are 'View cluster details' and 'View monitoring details' buttons.

US Federal Campaign Finance Data Analysis Using BigData Technologies



US Federal Campaign Finance Data Analysis Using BigData Technologies

The image consists of three screenshots demonstrating the setup and use of ODBC for data analysis:

- Amazon Hive ODBC Driver DSN Setup:** A Windows dialog box showing the configuration of a new ODBC connection named "hive_connection". The "Hive Server Type" is set to "Hive Server 2" and "Service Discovery Mode" is set to "No Service Discovery". The "Test Results" panel shows a successful connection test.
- ODBC Data Source Administrator (64-bit):** A Windows system dialog showing the "User Data Sources" list. It includes entries for dBASE Files, Excel Files, hive64, and MS Access Database. The "hive64" entry is highlighted.
- Power BI Desktop - Get Data:** A screenshot of the Power BI interface showing the "Get Data" dialog. The "ODBC" connector is selected under the "Other" category. The "Visualizations" pane on the right shows various chart and report options.

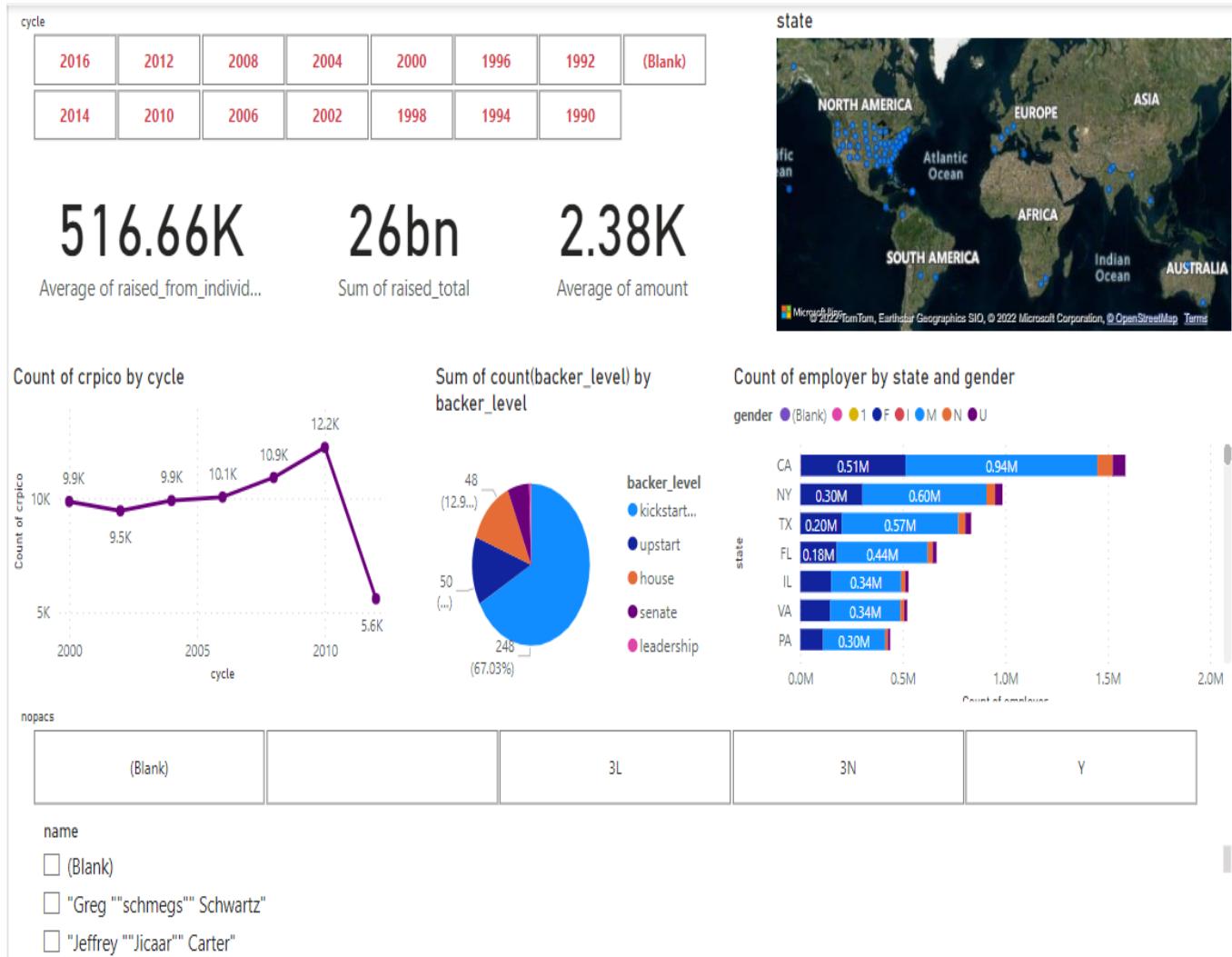
US Federal Campaign Finance Data Analysis Using BigData Technologies

The image consists of two screenshots of the Power BI Desktop application.

The top screenshot shows the 'From ODBC' dialog box. It has a dropdown menu 'Data source name (DSN)' set to 'hive_connection'. Below it is a link 'Get data from another source →'. The background shows the Power BI interface with various icons and a search bar.

The bottom screenshot shows the 'Navigator' feature. On the left, there's a tree view of data sources: 'ODBC (dsn=hive_connection) [1]', 'HIVE [1]', and 'default [13]'. Under 'default', several tables are listed, including 'backers' (which is selected), 'candidates', 'committees', 'contributions', 'countof_backerslevel', 'cycleof_candidates', 'fec_committees', 'industry_info', 'pac_records', 'pac_to_pac_info', 'pacs_info', 'politicians_info', and 'senate'. At the bottom of the Navigator window are buttons for 'Select Related Tables', 'Load', 'Transform Data', and 'Cancel'. A status message 'Preview is evaluating...' is visible in the center of the main workspace. The right side of the screen shows the 'Visualizations' pane with various chart and report options, and the 'Fields' pane.

11. PowerBI Dashboard:



12. Requirements Specification

Hardware Requirement

500 GB hard drive (Minimum requirement)

8 GB RAM (Minimum requirement)

PC x64-bit CPU

Software Requirement

Windows/Mac/Linux

Any Modern Web Browser like Google Chrome

To access aws account for creating emr etc

Power BI Desktop for creation dashboard

ODBC data sources for connecting hive to Power Bi

References

- US Federal Campaign Finance Data Analysis Using BigData Technologies Dataset

<https://www.kaggle.com/datasets/jeegarmaru/campaign-contributions-19902016>

<https://github.com/Priyanka743/Project-PUBGAnalysis>