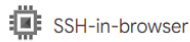


DATA CLEANING USING HIVE

This document covers all the steps taken to import the CSV to hive, clean, and export the CSV.

Step 1: Using DataProc create an Apache Hadoop cluster on Google Cloud Platform and Google Cloud Bucket to add the CSV dataset

Step 2: Launch the SSH session and import your CSV file using the GSUTIL URL to your local instance



```
Linux strikereport-m 5.10.0-26-cloud-amd64 #1 SMP Debian 5.10.197-1 (2023-09-29) x86_64

The programs included with the Debian GNU/Linux system are free software;
the exact distribution terms for each program are described in the
individual files in /usr/share/doc/*/copyright.

Debian GNU/Linux comes with ABSOLUTELY NO WARRANTY, to the extent
permitted by applicable law.
imgauravmehta@strikereport-m:~$ pwd
/home/imgauravmehta
imgauravmehta@strikereport-m:~$ mkdir strike_report
imgauravmehta@strikereport-m:~$ ls -l
total 4
drwxr-xr-x 2 imgauravmehta imgauravmehta 4096 Nov 30 19:38 strike_report
imgauravmehta@strikereport-m:~$ cd strike_report
imgauravmehta@strikereport-m:~/strike_report$ gsutil cp gs://strike_report/STRIKE_REPORTS.csv /^Cth/to/destination
imgauravmehta@strikereport-m:~/strike_report$ ^C
imgauravmehta@strikereport-m:~/strike_report$ gsutil cp gs://strike_report/STRIKE_REPORTS.csv /home/imgauravmehta/strike_report/.
Copying gs://strike_report/STRIKE_REPORTS.csv...
| [1 files][185.9 MiB/185.9 MiB]
Operation completed over 1 objects/185.9 MiB.
imgauravmehta@strikereport-m:~/strike_report$ ls -l
total 190400
-rw-r--r-- 1 imgauravmehta imgauravmehta 194965534 Nov 30 19:49 STRIKE_REPORTS.csv
imgauravmehta@strikereport-m:~/strike_report$
```

Step 3: Once the file is copied to the local instance copy the same file to the HADOOP file system

```
imgauravmehta@strikereport-m:~/strike_report$ hadoop fs -ls /
Found 3 items
drwxrwxrwt - hdfs hadoop 0 2023-11-30 19:33 /tmp
drwxrwxrwt - hdfs hadoop 0 2023-11-30 19:33 /user
drwxrwxrwt - hdfs hadoop 0 2023-11-30 19:33 /var
imgauravmehta@strikereport-m:~/strike_report$ hadoop fs -mkdir /user/strikereport
imgauravmehta@strikereport-m:~/strike_report$ hadoop fs -ls /user
Found 12 items
drwxrwxrwt - hdfs hadoop 0 2023-11-30 19:33 /user/dataproc
drwxrwxrwt - hdfs hadoop 0 2023-11-30 19:33 /user/hbase
drwxrwxrwt - hdfs hadoop 0 2023-11-30 19:33 /user/hdfs
drwxrwxrwt - hdfs hadoop 0 2023-11-30 19:33 /user/hive
drwxrwxrwt - hdfs hadoop 0 2023-11-30 19:33 /user/mapred
drwxrwxrwt - hdfs hadoop 0 2023-11-30 19:33 /user/pig
drwxrwxrwt - hdfs hadoop 0 2023-11-30 19:33 /user/solr
drwxrwxrwt - hdfs hadoop 0 2023-11-30 19:33 /user/spark
drwxr-xr-x - imgauravmehta hadoop 0 2023-11-30 19:51 /user/strikereport
drwxrwxrwt - hdfs hadoop 0 2023-11-30 19:33 /user/yarn
drwxrwxrwt - hdfs hadoop 0 2023-11-30 19:33 /user/zeppelin
drwxrwxrwt - hdfs hadoop 0 2023-11-30 19:33 /user/zookeeper
imgauravmehta@strikereport-m:~/strike_report$
```

```
imgauravmehta@strikereport-m:~/strike_report$ hadoop fs -copyFromLocal /home/imgauravmehta/strike_report/STRIKE_REPORTS.csv /user/strikereport/.
imgauravmehta@strikereport-m:~/strike_report$ hadoop fs -ls /user/strikereport
Found 1 items
-rw-r--r-- 1 imgauravmehta hadoop 194965534 2023-11-30 19:54 /user/strikereport/STRIKE_REPORTS.csv
imgauravmehta@strikereport-m:~/strike_report$
```

Step 4: Once the file is saved to our HADOOP file system, Initiate HIVE

```
imgauravmhta@strikerreport-mr:~/strike_report$ hive
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/lib/tez/lib/slf4j-reload4j-1.7.36.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/lib/hadoop/lib/slf4j-reload4j-1.7.36.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Reload4jLoggerFactory]
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/lib/tez/lib/slf4j-reload4j-1.7.36.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/lib/hadoop/lib/slf4j-reload4j-1.7.36.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Reload4jLoggerFactory]
Hive Session ID = 1ec1081c-f602-4ae9-8ba7-2726cdf63cdd

Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j2.properties Async: true
WARNING: An illegal reflective access operation has occurred
WARNING: Illegal reflective access by org.apache.hadoop.hive.common.StringInternUtils (file:/usr/lib/hive/lib/hive-common-3.1.3.jar) to field java.net.URI.string
WARNING: Please consider reporting this to the maintainers of org.apache.hadoop.hive.common.StringInternUtils
WARNING: Use --illegal-access=warn to enable warnings of further illegal reflective access operations
WARNING: All illegal access operations will be denied in a future release
Hive Session ID = e1754683-f018-4aa4-93c4-7077e57bc9f3
hive> 
```

Step 5: Create a new database in HIVE as air_accidents and use the same database to create the table

```
hive> create database air_accidents;
OK
Time taken: 0.534 seconds
```

```
hive> use air_accidents;
OK
Time taken: 0.06 seconds
hive> 
```

Step 6: Create an EXTERNAL TABLE in HIVE with the columns and their data types along with the delimiter passed as {,} to ensure the CSV file standards are maintained.

```
hive> CREATE EXTERNAL TABLE IF NOT EXISTS air_strike_data (
> INDEX NR INT,
> INCIDENT_DATE STRING,
> INCIDENT_MONTH INT,
> INCIDENT_YEAR INT, 'TIME' STRING,
> TIME_OF_DAY STRING,
> AIRPORT_ID STRING,
> AIRPORT STRING,
> LATITUDE DOUBLE,
> LONGITUDE DOUBLE,
> RUNWAY STRING,
> STATE STRING,
> FAAREGION STRING,
> LOCATION STRING,
> ENROUTE_STATE STRING,
> OPID STRING,
> OPERATOR STRING,
> REG STRING,
> FLT STRING,
> AIRCRAFT STRING,
> AMA STRING,
> AMO STRING,
> EMA STRING,
> EMO STRING,
> AC_CLASS STRING,
> AC_MASS INT,
> TYPE_ENG STRING,
> NUM_ENGS INT,
> ENG_1_POS INT,
> ENG_2_POS INT,
> ENG_3_POS INT,
> ENG_4_POS INT,
> PHASE_OF_FLIGHT STRING,
> HEIGHT INT,
> SPEED INT,
> DISTANCE INT,
> SKY STRING,
> PRECIPITATION STRING,
> AOS STRING,
> COST_REPAIRS INT,
> COST_OTHER INT,
> COST_REPAIRS_INFL_ADJ INT,
> COST_OTHER_INFL_ADJ INT,
> INGESTED_OTHER BOOLEAN,
> INDICATED_DAMAGE BOOLEAN,
> DAMAGE_LEVEL STRING,
> STR_RAD BOOLEAN,
> STR_ENG4 BOOLEAN,
> DAM_ENG4 BOOLEAN,
> ING_ENG4 BOOLEAN,
> STR_PROP BOOLEAN,
> DAM_PROP BOOLEAN,
> STR_WING_ROT BOOLEAN,
> DAM_WING_ROT BOOLEAN,
> STR_FUSE BOOLEAN,
> DAM_FUSE BOOLEAN,
> STR_LG BOOLEAN,
> DAM_LG BOOLEAN,
> STR_TAIL BOOLEAN,
> DAM_TAIL BOOLEAN,
> STR_LGHTS BOOLEAN,
> DAM_LGHTS BOOLEAN,
> STR_OTHER BOOLEAN,
> DAM_OTHER BOOLEAN,
> OTHER_SPECIFY STRING,
> EFFECT STRING,
> EFFECT_OTHER STRING,
> BIRD_BAND_NUMBER STRING,
> SPECIES_ID STRING,
> SPECIES STRING,
> REMARKS STRING,
> REMAINS_COLLECTED BOOLEAN,
> REMAINS_SENT BOOLEAN,
> WARNED BOOLEAN,
> NUM_SEEN INT,
> NUM_STRUCK INT,
> SIZE STRING,
> NR_INJURIES INT,
> NR_FATALITIES INT,
> COMMENTS STRING,
> REPORTED_NAME STRING,
> REPORTED_TITLE STRING,
> SOURCE STRING,
> PERSON STRING,
> LUPDATE STRING,
> TRANSFER STRING
> )
> ROW FORMAT DELIMITED
> FIELDS TERMINATED BY '{,'
> STORED AS TEXTFILE
> LOCATION '/user/strikerreport/STRIKE_REPORTS';
OK
Time taken: 0.564 seconds
hive> 
```

[illegible]

1. DROPPING COLUMNS
 - Dropping the columns that will not be required at all for analysis purposes especially with maximum NULL values.
Columns Dropped: Time, Runway, Location, Enroute_State, OPID (Operator Column already present), AOS, COST_REPAIRS, COST_OTHER, COST_REPAIRS_INFL_ADJ, COST_OTHER_INFL_ADJ, OTHER_SPECIFY, EFFECT_OTHER, BIRD_BAND_NUMBER, NUM_SEEN, NUM_STRUCK, NR_INJURIES, NR_FATALITIES, COMMENTS, REPORTED_NAME, REPORTED_TITLE, SOURCE, PERSON, LUPDATE, REASON
 - Dropping LATITUDE and LONGITUDES as we have AIRPORT name and AIRPORT ID.
 - Removing ENGINE_POSITIONS as well as this variable does not help in analysis purposes.
 - INCIDENT_DATE is yet another column that can be dropped as the values are inconsistent and INCIDENT_MONTH and INCIDENT_YEAR can help us for analysis purposes.

[illegible]

- Drop rows with airport as unknown or missing and create a new table as new_bird_strike using the query “CREATE TABLE new_bird_strike AS SELECT * FROM bird_strike WHERE AIRPORT != 'Unknown';”

```
hive> CREATE TABLE new_bird_strike AS SELECT * FROM bird_strike where airport != 'UNKNOWN';
Query ID = imgauravmehta_20231130235657_2ca8997e-7f00-429f-a867-e7995cf48632
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1701372766469_0008)

-----
VERTICES      MODE          STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
Map 1 ..... container      SUCCEEDED      4          4          0          0          0          0
-----
VERTICES: 01/01 [=====>]] 100% ELAPSED TIME: 17.17 s
-----
Moving data to directory hdfs://strikerreport-m/user/hive/warehouse/new_bird_strike
OK
Time taken: 17.629 seconds
hive> SELECT COUNT(*) FROM new_bird_strike;
OK
253393
Time taken: 0.119 seconds, Fetched: 1 row(s)
hive> □
```

- TIME_OF_DAY values with DUSK will be replaced with NIGHT and DAWN with DAW to remove data inconsistencies (Table name: birdstrike_1)

```
hive> CREATE TABLE birdstrike_1 AS SELECT *, CASE WHEN TIME_OF_DAY = 'Dusk' THEN 'Night' WHEN TIME_OF_DAY = 'Dawn' THEN 'Day' END AS TIME_OF_DAY FROM new_bird_strike;
FAILED: SemanticException [Error 10036]: Duplicate column name: time_of_day
hive> CREATE TABLE birdstrike_1 AS SELECT *, CASE WHEN TIME_OF_DAY = 'Dusk' THEN 'Night' WHEN TIME_OF_DAY = 'Dawn' THEN 'Day' END AS TIME_OF_DAY FROM new_bird_strike;
Query ID = imgauravmehta.20231201004113.f4e0f2c4-9113-446a-8d0a-b91ef276fc5
Total jobs = 1
Launching Job 1 out of 1
Status: Running (Executing on YARN cluster with App id application_1701372766469_0014)
```

VERTICES	MEKE	STATUS	TOTAL	COMPLETED	RUNNING	PENDING	FAILED	KILLED
Map 1	container	SUCCEEDED	6	6	0	0	0	0

```
VERTICES: 51/51 [=====] 100% ELAPSED TIME: 19.93 s
Moving data to directory hdfs://strike-report-m/user/hive/warehouse/birdstrike_1
OK
Time taken: 116.114 seconds
hive> select * from birdstrike_1 limit 5;
OK
1127070 5 2021 Day ELAX LOS ANGELES INTL CA AMP SKYWEST AIRLINES 322 22 04 A 4 D 2 Take-off Run 0 NULL 0 0v
erocrat Rain N false false false false false false false UNERS Unknown bird - small "Skywest/United Express FLT 5411 reported bird strike on RWY 25R during take off run near TWY B2. One (1)
lax false false false false false true false false false false false false false false false false false false false false false false false
small bird was recovered near the RWY 25R centerline at TWY B2 intersection. 04511 continued departure for PHX. We confirmed ID of remains. BIRD STRIKE NULL NULL NULL Yes *** Strike Report:
2021-05-16-153553-93 (Report ID: 1101761 [ Comment: ] ) Merged From 2021-05-16-50805-21 (Report ID: 1124799 [ Comment: ] ) on Aug 18 2021 1:53PM By mahalah.e.schan@usda.gov *** NULL NULL NULL Some Cloud
1127074 5 2021 Night EJFK JOHN F KENNEDY INTL NY AEA JETBLUE AIRWAYS 04A 07 23 01 A 4 D 2 Approach 1400 NULL NULL NULL
None N false false false false true false false false false false false false false false false false false false false false
lax false false false false false false Y Perching birds (y) "Please note the time stated in report reflects the time the strike was reported to JFK airport operations staff."
Upon my arrival to the gate NULL true true 1 REDACTED NULL
1127102 7 2021 Day RWY DETROIT METRO WAYNE COUNTY ARPT MI AGL FFA AIRLINES 188 17 22 04 A 4 D 2 Landing Roll 0 100 0 No
cloud N false false false false false false false false false false false false false false false false false false false false
lax false false false false false false false false false false false false false false false false false false false false
1127107 7 2021 RWY SAN FRANCISCO INTL ARPT CA AMP UNKNOWN NULL NULL NULL NULL NULL 0
lax false false false false false false false false false false false false false false false false false false false
lax false false false false false K3311 Cooper's hawk "(1) HAWK RECOVERED ON RUNWAY 10R NULL true false 1 REDACTED NULL NULL NULL 0
1127108 7 2021 RWY DETROIT METRO WAYNE COUNTY ARPT MI AGL UNKNOWN NULL NULL NULL NULL NULL 0
also false false false false false false false false false false false false false false false false false false false
lax false false false false false YH004 Horned lark Found during morning runway inspections. DAY true false NULL Small FAA Form 5200-7-E NULL
Time taken: 0.204 seconds, Fetched: 5 row(s)
hive>
```

Step 9: The table birdstrike_1 can now be exported to google bucket for further analysis. While exporting the file from HIVE to local instance, the file has multiple outputs.

The multiple output files can be combined into 1 single file using the `cat` command, the file can be now named `combined_output.csv`.

Using the `GSUTIL COPY` command, export `combined_output.csv` to Cloud Storage Bucket.

```
imgauravmehta@strike-report-m:~/strike_report$ ls
000000_0 000001_0 000002_0 000003_0 000004_0 000005_0
imgauravmehta@strike-report-m:~/strike_report$ ls -l
total 373456
-rw-r--r-- 1 imgauravmehta imgauravmehta 736364448 Dec 1 01:17 000000_0
-rw-r--r-- 1 imgauravmehta imgauravmehta 706041112 Dec 1 01:17 000001_0
-rw-r--r-- 1 imgauravmehta imgauravmehta 69221302 Dec 1 01:17 000002_0
-rw-r--r-- 1 imgauravmehta imgauravmehta 68377224 Dec 1 01:17 000003_0
-rw-r--r-- 1 imgauravmehta imgauravmehta 50723867 Dec 1 01:17 000004_0
-rw-r--r-- 1 imgauravmehta imgauravmehta 49846531 Dec 1 01:17 000005_0
imgauravmehta@strike-report-m:~/strike_report$ cat /home/imgauravmehta/strike_report/* > /home/imgauravmehta/strike_report/combined_output.csv
imgauravmehta@strike-report-m:~/strike_report$ ls
000000_0 000001_0 000002_0 000003_0 000004_0 000005_0 combined_output.csv
imgauravmehta@strike-report-m:~/strike_report$ less /home/imgauravmehta/strike_report/combined_output.csv
imgauravmehta@strike-report-m:~/strike_report$ gsutil cp /home/imgauravmehta/strike_report/combined_output.csv gs://strike_report
Copying file:///home/imgauravmehta/strike_report/combined_output.csv [Content-Type=text/csv]...
==> NOTE: You are uploading one or more large file(s), which would run
significantly faster if you enable parallel composite uploads. This
feature can be enabled by editing the
"parallel_composite_upload_threshold" value in your .boto
configuration file. However, note that if you do this large files will
be uploaded as "composite objects"
<https://cloud.google.com/storage/docs/composite-objects>, which
means that any user who downloads such objects will need to have a
compiled crcmod installed (see "gsutil help crcmod"). This is because
without a compiled crcmod, computing checksums on composite objects is
so slow that gsutil disables downloads of composite objects.

/ [1 files] [364.7 MiB/364.7 MiB]
Operation completed over 1 objects/364.7 MiB.
imgauravmehta@strike-report-m:~/strike_report$
```

Step 10: Verify the file in Google Bucket (strike_report)

strike_report

Location

us-east1 (South Carolina)

Storage class

Standard

Public access

Not public

Protection

None

OBJECTS

CONFIGURATION

PERMISSIONS

PROTECTION

LIFECYCLE

OBSERVABILITY

INVENTORY REPORTS

Buckets

>

strike_report

UPLOAD FILES

UPLOAD FOLDER

CREATE FOLDER

TRANSFER DATA

MANAGE HOLDS

DOWNLOAD

DELETE

Filter by name prefix only

Filter

Filter objects and folders

Show deleted data

<input type="checkbox"/>	Name	Size	Type	Created	Storage class	Last modified	Public access	Version history	Encryption	R
<input checked="" type="checkbox"/>	STRIKE_REPORTS.csv	185.9 MB	text/csv	Nov 30, 2023, 2:28:01 PM	Standard	Nov 30, 2023, 2:28:01 PM	Not public	—	Google-managed	⌵ ⌴ ⌵
<input checked="" type="checkbox"/>	combined_output.csv	364.7 MB	text/csv	Nov 30, 2023, 8:23:34 PM	Standard	Nov 30, 2023, 8:23:34 PM	Not public	—	Google-managed	⌵ ⌴ ⌵

The output file (`combined_output.csv`) can now be downloaded and loaded into Jupyter Notebook with Spark to perform further cleaning and analysis.