

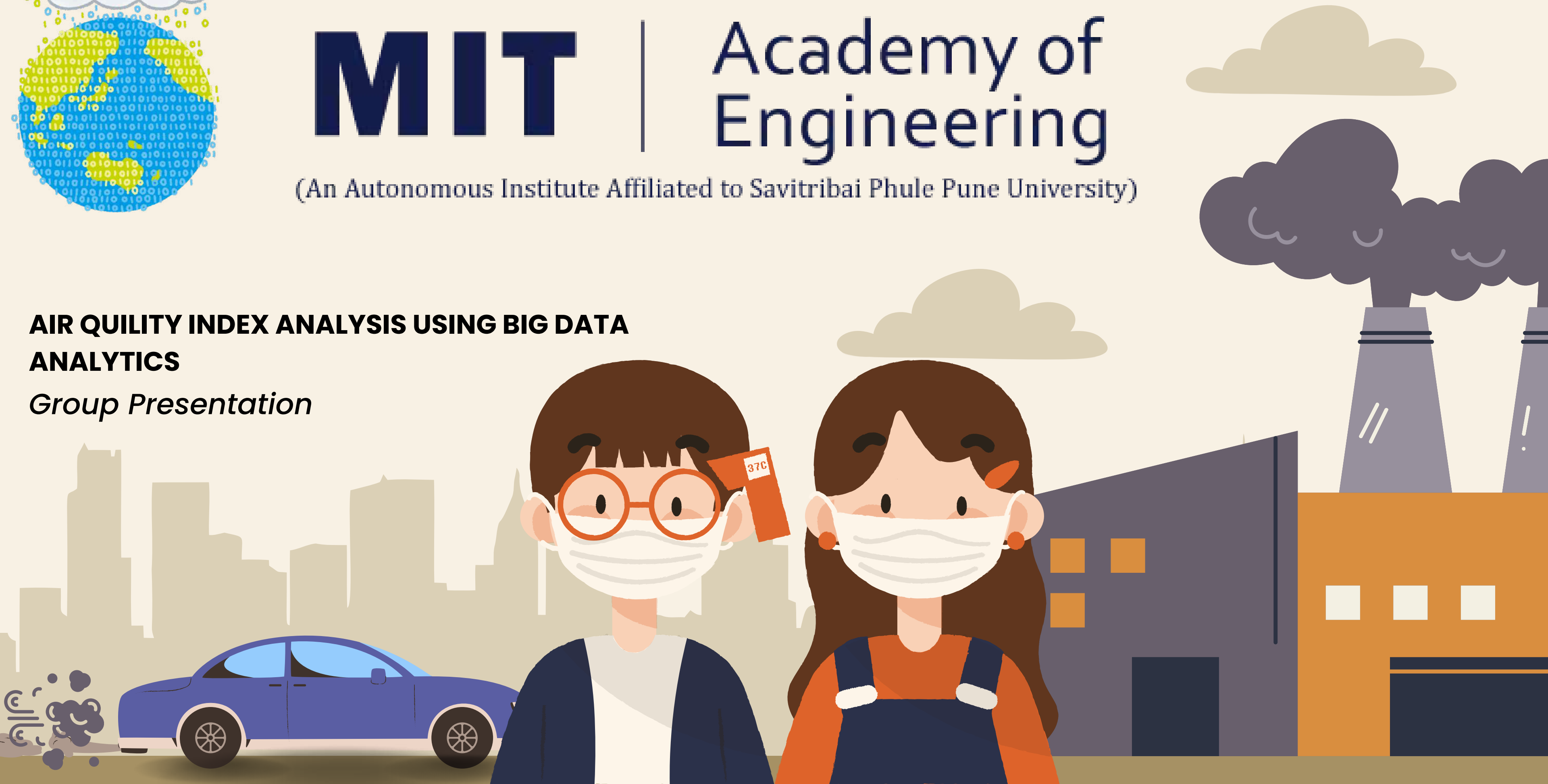


MIT | Academy of Engineering

(An Autonomous Institute Affiliated to Savitribai Phule Pune University)

AIR QUALITY INDEX ANALYSIS USING BIG DATA ANALYTICS

Group Presentation



GROUP Members name

Nirmal Chaturvedi	202201040210
Vedant Deokar	202201040164
Om Mukkawar	202201040193
Hridill Banik	202201040170

Project Guide- Mr. Pramod Dharmadhikari

Problem Statement:

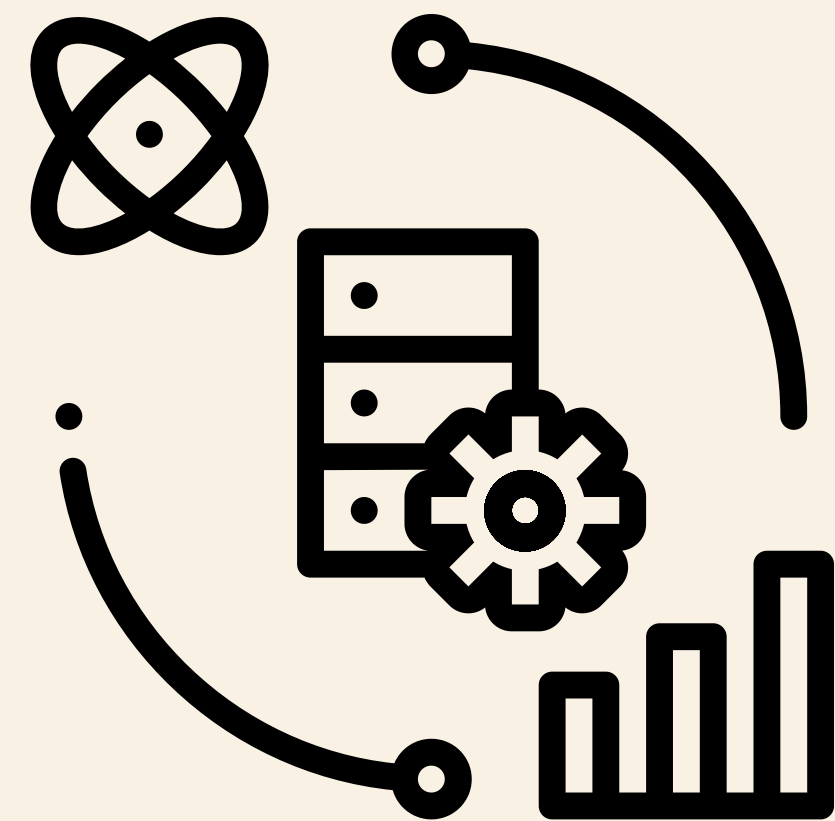
Due to increasing global pollution, air quality has become a major health concern worldwide. This project analyzes AQI data to identify the most polluted countries and categorize air quality levels to support better environmental decision-making.

dataset used - <https://www.kaggle.com/datasets/hasibalmuzdadid/global-air-pollution-dataset?resource=download>

1	Country	City	AQI Value	AQI Category	CO AQI Value	CO AQI Category	Ozone AQI Value	Ozone AQI Category	NO2 AQI Value	NO2 AQI Category	PM2.5 AQI Value	PM2.5 AQI Category		
2	Russian Federation	Praskoveyevskaya	51	Moderate	1	Good	36	Good	0	Good	51	Moderate		
3	Brazil	Presidente Dutra	41	Good	1	Good	5	Good	1	Good	41	Good		
4	Italy	Priolo Gargallo	66	Moderate	1	Good	39	Good	2	Good	66	Moderate		
5	Poland	Przasnysz	34	Good	1	Good	34	Good	0	Good	20	Good		
6	France	Punaauia	22	Good	0	Good	22	Good	0	Good	6	Good		
7	United States	Punta Gorda	54	Moderate	1	Good	14	Good	11	Good	54	Moderate		
8	Germany	Puttlingen	62	Moderate	1	Good	35	Good	3	Good	62	Moderate		
9	Belgium	Puurs	64	Moderate	1	Good	29	Good	7	Good	64	Moderate		
0	Russian Federation	Pyatigorsk	54	Moderate	1	Good	41	Good	1	Good	54	Moderate		
1	Egypt	Qalyub	142	Unhealthy	3	Good	89	Moderate	9	Good	142	Unhealthy for Sensitive Groups		
2	China	Qinzhou	68	Moderate	2	Good	68	Moderate	1	Good	58	Moderate		
3	Netherlands	Raalte	41	Good	1	Good	24	Good	6	Good	41	Good		
4	India	Radaur	158	Unhealthy	3	Good	139	Unhealthy for Sensitive Groups	1	Good	158	Unhealthy		
5	Pakistan	Radhan	158	Unhealthy	1	Good	50	Good	1	Good	158	Unhealthy		
6	Republic of Moldova	Radovis	83	Moderate	1	Good	46	Good	0	Good	83	Moderate		
7	France	Raismes	59	Moderate	1	Good	30	Good	4	Good	59	Moderate		
8	India	Rajgir	154	Unhealthy	3	Good	100	Unhealthy for Sensitive Groups	2	Good	154	Unhealthy		
9	Italy	Ramacca	55	Moderate	1	Good	47	Good	0	Good	55	Moderate		
0	United States	Phoenix	72	Moderate	1	Good	4	Good	23	Good	72	Moderate		
1	India	Phulabani	161	Unhealthy	2	Good	71	Moderate	0	Good	161	Unhealthy		
2	Poland	Piaseczno	28	Good	1	Good	28	Good	2	Good	28	Good		
3	India	Pimpri	118	Unhealthy	2	Good	30	Good	2	Good	118	Unhealthy for Sensitive Groups		
4	Brazil	Pindobacunga	33	Good	0	Good	10	Good	1	Good	33	Good		
5	China	Pingyin	150	Unhealthy	3	Good	95	Moderate	6	Good	150	Unhealthy		
6	Brazil	Pinheiral	154	Unhealthy	5	Good	0	Good	13	Good	154	Unhealthy		
7	India	Piravam	81	Moderate	1	Good	24	Good	1	Good	81	Moderate		
8	United States	Pittsburg	67	Moderate	1	Good	15	Good	3	Good	67	Moderate		



Hadoop / HDFS Operations



Hadoop Operations

hdfs dfs -ls /

```
[cloudera@quickstart ~]$ hdfs dfs -ls /
Found 6 items
drwxrwxrwx   - hdfs  supergroup          0 2017-10-23 09:15 /benchmarks
drwxr-xr-x   - hbase supergroup          0 2025-11-20 00:54 /hbase
drwxr-xr-x   - solr  solr                0 2017-10-23 09:18 /solr
drwxrwxrwt   - hdfs  supergroup          0 2025-11-19 13:27 /tmp
drwxr-xr-x   - hdfs  supergroup          0 2017-10-23 09:17 /user
drwxr-xr-x   - hdfs  supergroup          0 2017-10-23 09:17 /var
[cloudera@quickstart ~]$
```

hdfs dfs -ls /user/cloudera/airquality2

```
[cloudera@quickstart ~]$ hdfs dfs -ls /user/cloudera/airquality2
Found 1 items
-rw-r--r--   1 cloudera cloudera    1630831 2025-11-19 13:35 /user/cloudera/airquality2/airquality.csv
[cloudera@quickstart ~]$
```

hdfs dfs -cat /user/cloudera/airquality2/airquality.csv | head

```
[cloudera@quickstart ~]$ hdfs dfs -cat /user/cloudera/airquality2/airquality.csv | head
Country,City,AQI Value,AQI Category,CO AQI Value,CO AQI Category,Ozone AQI Value,Ozone AQI Category,N02 AQI Value,N02 AQI Category,PM2.5 AQI Value,PM2.5 AQI Category
Russian Federation,Prskoveya,51,Moderate,1,Good,36,Good,0,Good,51,Moderate
Brazil,Presidente Dutra,41,Good,1,Good,5,Good,1,Good,41,Good
Italy,Priolo Gargallo,66,Moderate,1,Good,39,Good,2,Good,66,Moderate
Poland,Przasnysz,34,Good,1,Good,34,Good,0,Good,20,Good
France,Punaauia,22,Good,0,Good,22,Good,0,Good,6,Good
United States of America,Punta Gorda,54,Moderate,1,Good,14,Good,11,Good,54,Moderate
Germany,Puttlingen,62,Moderate,1,Good,35,Good,3,Good,62,Moderate
Belgium,Puurs,64,Moderate,1,Good,29,Good,7,Good,64,Moderate
Russian Federation,Pyatigorsk,54,Moderate,1,Good,41,Good,1,Good,54,Moderate
cat: Unable to write to output stream.
[cloudera@quickstart ~]$
```

Copy + Move example

```
cat > /dev/null <& echo "enable to write to output stream."
[cloudera@quickstart ~]$ hdfs dfs -cp /user/cloudera/airquality2/airquality.csv /user/cloudera/airquality2/airquality_copy.csv
[cloudera@quickstart ~]$ hdfs dfs -ls /user/cloudera/airquality2
Found 2 items
-rw-r--r--  1 cloudera cloudera    1630831 2025-11-19 13:35 /user/cloudera/airquality2/airquality.csv
-rw-r--r--  1 cloudera cloudera    1630831 2025-11-20 01:12 /user/cloudera/airquality2/airquality_copy.csv
[cloudera@quickstart ~]$ hdfs dfs -mkdir /user/cloudera/airquality_backup
[cloudera@quickstart ~]$ hdfs dfs -mv /user/cloudera/airquality2/airquality_copy.csv /user/cloudera/airquality_backup/
[cloudera@quickstart ~]$ hdfs dfs -ls /user/cloudera/airquality_backup
Found 1 items
-rw-r--r--  1 cloudera cloudera    1630831 2025-11-20 01:12 /user/cloudera/airquality_backup/airquality_copy.csv
[cloudera@quickstart ~]$ █
```

Size + count

hdfs dfs -du -h /user/cloudera/airquality2

hdfs dfs -count /user/cloudera/airquality2

```
1.6 M  1.6 M  /user/cloudera/airquality2/airquality.csv
[cloudera@quickstart ~]$ hdfs dfs -count /user/cloudera/airquality2
      1      1      1630831 /user/cloudera/airquality2
[cloudera@quickstart ~]$ hdfs dfs -du -h /user/cloudera/airquality2
1.6 M  1.6 M  /user/cloudera/airquality2/airquality.csv
[cloudera@quickstart ~]$ hdfs dfs -count /user/cloudera/airquality2
      1      1      1630831 /user/cloudera/airquality2
[cloudera@quickstart ~]$ █
```



Apache Pig

PIG OPERATION

Pig Script (Data Processing)

[cloudera@quickstart ~]\$ pig

Load Air Quality Dataset and Filter Valid AQI Records

```
grunt> airdata = LOAD '/user/cloudera/airquality2/airquality.csv' USING PigStorage(',')
>> AS (
>>   country:chararray,
>>   city:chararray,
>>   aqi_value:int,
>>   aqi_category:chararray,
>>   co_value:int,
>>   co_category:chararray,
>>   ozone_value:int,
>>   ozone_category:chararray,
>>   no2_value:int,
>>   no2_category:chararray,
>>   pm25_value:int,
>>   pm25_category:chararray
>> );
grunt>
grunt> air_clean = FILTER airdata BY aqi_category IS NOT NULL;
grunt> █
```

Countries with best air quality (max Good cities):

```
good = FILTER air_clean BY aqi_category == 'Good';
good_grp = GROUP good BY country;
good_count = FOREACH good_grp GENERATE group AS country, COUNT(good) AS good_cities;
sorted_good = ORDER good_count BY good_cities DESC;
top10_good = LIMIT sorted_good 10;
DUMP top10_good;
```

```
2025-11-20 01:28:51,319 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(Brazil,1125)
(Russian Federation,1025)
(United States of America,1001)
(Germany,717)
(Japan,432)
(France,414)
(Spain,372)
(United Kingdom of Great Britain and Northern Ireland,319)
(Italy,283)
(Netherlands,274)
grunt> █
```

Store result in HDFS

STORE top10_good INTO '/user/cloudera/pig_output/top10_good_countries' USING PigStorage(',');

```
at org.apache.hadoop.util.RunJar.main(RunJar.java:136)
2025-11-20 01:32:11,991 [main] INFO org.apache.hadoop.mapred.ClientServiceDelegate - Application state is completed. FinalApplicationStatus=SUCCEEDED. Redirecting to job history server
2025-11-20 01:32:12,304 [main] WARN org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Encountered Warning FIELD_DISCARDED_TYPE_CONVERSION_FAILED 5 time(s).
2025-11-20 01:32:12,304 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
grunt> quit
[cloudera@quickstart ~]$ hdfs dfs -get /user/cloudera/pig_output/top10_good_countries /home/cloudera/output/
[cloudera@quickstart ~]$ mv /home/cloudera/output/top10_good_countries/part-r-00000 /home/cloudera/output/pig_top10_good_countries.csv
[cloudera@quickstart ~]$ sudo cp /home/cloudera/output/pig_top10_good_countries.csv /media/sf_global_air_pollution_dataset.csv/
[cloudera@quickstart ~]$ █
```

Top 10 Most Polluted Cities (Highest AQI).

```
airdata = LOAD '/user/cloudera/airquality2/airquality.csv' USING PigStorage(',')
AS (country:chararray, city:chararray, aqi_value:int, aqi_category:chararray,
    co_value:int, co_category:chararray, ozone_value:int, ozone_category:chararray,
    no2_value:int, no2_category:chararray, pm25_value:int, pm25_category:chararray);
air_filter = FILTER airdata BY aqi_value IS NOT NULL;
sorted_aqi = ORDER air_filter BY aqi_value DESC;
top10_polluted = LIMIT sorted_aqi 10;
DUMP top10_polluted;
```

```
2025-11-20 01:47:19,837 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat: Total input paths to process: 1
(India,Ratangarh,500,Hazardous,0,Good,37,Good,0,Good,471,Hazardous)
(India,Nuh,500,Hazardous,1,Good,44,Good,1,Good,420,Hazardous)
(India,Phalodi,500,Hazardous,0,Good,34,Good,0,Good,467,Hazardous)
(India,Bilari,500,Hazardous,4,Good,158,Unhealthy,4,Good,457,Hazardous)
(United States of America,Durango,500,Hazardous,133,Unhealthy for Sensitive Groups,0,Good,53,Moderate,500,Hazardous)
(India,Gohana,500,Hazardous,1,Good,47,Good,1,Good,500,Hazardous)
(India,Pilani,500,Hazardous,1,Good,39,Good,0,Good,433,Hazardous)
(India,Jahangirpur,500,Hazardous,1,Good,49,Good,1,Good,493,Hazardous)
(India,Gopamau,500,Hazardous,1,Good,52,Moderate,1,Good,378,Hazardous)
(India,Kasganj,500,Hazardous,1,Good,61,Moderate,3,Good,476,Hazardous)
grunt>
```

Count by AQI Category (Good / Moderate / Unhealthy).

```
grp_cat = GROUP airdata BY aqi_category;
cat_count = FOREACH grp_cat GENERATE group, COUNT(airdata);
sorted_cat = ORDER cat_count BY $1 DESC;
DUMP sorted_cat;
```

```
2025-11-20 01:47:19,837 [main] INFO org.apache.hadoop.mapreduce.lib.output.FileOutputFormat: Total output paths to process: 1
(Good,9688)
(Moderate,9088)
(Unhealthy,2215)
(Unhealthy for Sensitive Groups,1568)
(Very Unhealthy,286)
(Hazardous,191)
(AQI Category,1)
grunt>
```

```
grunt> quit
[cloudera@quickstart ~]$ hdfs dfs -ls /user/cloudera/output/aqi_category_count_pig
Found 2 items
-rw-r--r--  1 cloudera cloudera      0 2025-11-20 01:55 /user/cloudera/output/aqi_category_count_pig/_SUCCESS
-rw-r--r--  1 cloudera cloudera    123 2025-11-20 01:55 /user/cloudera/output/aqi_category_count_pig/part-r-000000
[cloudera@quickstart ~]$ hdfs dfs -get /user/cloudera/output/aqi_category_count_pig/part-r-000000 \
> /home/cloudera/aqi_category_count_pig.csv
[cloudera@quickstart ~]$ sudo cp /home/cloudera/aqi_category_count_pig.csv \
> /media/sf_global_air_pollution_dataset.csv/
[cloudera@quickstart ~]$
```



HIVE OPERATION

Hive Operation

Start Hive Services

```
sudo service hive-metastore start
sudo service hive-server2 start
```

```
[cloudera@quickstart ~]$ sudo service hive-metastore start
Starting Hive Metastore (hive-metastore): [ OK ]
Hive Metastore is running [ OK ]
[cloudera@quickstart ~]$ sudo service hive-server2 start
Hive Server2 is running [ OK ]
[cloudera@quickstart ~]$ █
```

Connect to Hive Using Beeline

```
beeline -u "jdbc:hive2://localhost:10000/airquality" -n cloudera
```

Creating the table- *air_table1*

```
CREATE TABLE air_table1 (
  country STRING,
  city STRING,
  aqi_value INT,
  aqi_category STRING,
  co_value INT,
  co_category STRING,
  ozone_value INT,
  ozone_category STRING,
  no2_value INT,
  no2_category STRING,
  pm25_value INT,
  pm25_category STRING
)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
STORED AS TEXTFILE;
```

Display all records

```
SELECT COUNT(*) FROM air_table1;
```

```
INFO : 2025-11-20 02:14:51,149 Stage-1 map = 100%, reduce = 0%, Cumu
INFO : 2025-11-20 02:15:31,574 Stage-1 map = 100%, reduce = 100%, Cu
INFO : MapReduce Total cumulative CPU time: 27 seconds 230 msec
INFO : Ended Job = job_1763628883521_0020
INFO : MapReduce Jobs Launched:
INFO : Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 27.23 sec
INFO : Total MapReduce CPU Time Spent: 27 seconds 230 msec
INFO : Completed executing command(queryId=hive_20251120021313_831717.
INFO : OK
+-----+--+
| _c0 |
+-----+--+
| 23464 |
+-----+--+
1 row selected (134.813 seconds)
0: jdbc:hive2://localhost:10000/airquality> █
```

Netherlands	3	Kesteren	37	37	Good	Good	1	Good	31	Good
Belgium	3	Good	37	39	Good	Good	1	Good	26	Good
Azerbaijan	0	Good	39	90	Good	Moderate	1	Good	43	Good
Brazil	0	Good	90	42	Moderate	Good	1	Good	6	Good
United States of America	0	Good	42	36	Good	Good	0	Good	35	Good
Italy	0	Good	36	36	Good	Good	1	Good	42	Good
China	4	Good	76	39	Moderate	Good	2	Good	39	Good
Sri Lanka	1	Good	37	43	Good	Good	1	Good	17	Good
Netherlands	2	Good	43	55	Good	Moderate	1	Good	39	Good
Serbia	1	Good	55	61	Moderate	Moderate	1	Good	32	Good
Italy	0	Good	61	64	Moderate	Moderate	1	Good	22	Good
Greece	3	Good	64	78	Moderate	Moderate	1	Good	40	Good
Russian Federation	0	Good	78	36	Moderate	Good	1	Good	32	Good
Uganda	6	Good	36	74	Good	Moderate	3	Good	21	Good
Germany	0	Good	74	68	Moderate	Moderate	2	Good	29	Good
Pakistan	8	Good	68	178	Moderate	Unhealthy	2	Good	120	Unhea
lthy for Sensitive Groups	0	Good	178	30	Unhealthy	Good	1	Good	30	Good
Netherlands	2	Good	30	21	Good	Good	1	Good	32	Good
Belgium	4	Good	21	61	Good	Moderate	1	Good		

```
LOAD DATA INPATH '/user/cloudera/airquality/airquality.csv'
INTO TABLE air_table1;
```


HIVE SCRIPT (DATA PROCESSING)

0: jdbc:hive2://localhost:10000/airquality>

Countries with Maximum Good AQI Cities

```
SELECT country, COUNT(*) AS good_count
FROM air_table1
WHERE aqi_category = 'Good'
GROUP BY country
ORDER BY good_count DESC
LIMIT 10;
```

INFO : Completed executing command(queryId=hive_20231120022323_50eb4c16-012d-4910-0110-d341)

INFO : OK

country	good_count
Brazil	1125
Russian Federation	1025
United States of America	1001
Germany	717
Japan	432
France	414
Spain	372
United Kingdom of Great Britain and Northern Ireland	319
Italy	283
Netherlands	274

10 rows selected (199.934 seconds)

0: jdbc:hive2://localhost:10000/airquality>

Average AQI by Country + Category Combination

```
INSERT OVERWRITE DIRECTORY
'/user/cloudera/hive_output/avg_aqi_country_category'
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
SELECT country, aqi_category, AVG(aqi_value) AS avg_aqi
FROM air_table1
GROUP BY country, aqi_category
ORDER BY avg_aqi DESC;
```

```
hdfs dfs -ls
user/cloudera/hive_output/avg_aqi_country_category
hdfs dfs -cat
cloudera/hive_output/avg_aqi_country_category/000000
_0 | head
```

country	aqi_category	avg_aqi
United States of America	Hazardous	500.0
Russian Federation	Hazardous	500.0
India	Hazardous	450.62025316455697
Republic of Korea	Hazardous	421.0
Pakistan	Hazardous	419.46153846153845
Democratic Republic of the Congo	Hazardous	415.0
South Africa	Hazardous	403.0
Chile	Hazardous	358.0
Mexico	Hazardous	339.6666666666667
China	Hazardous	332.3333333333333
Uzbekistan	Hazardous	300.0
Angola	Very Unhealthy	285.0
Australia	Very Unhealthy	264.0
Democratic Republic of the Congo	Very Unhealthy	249.5
Brazil	Very Unhealthy	239.33333333333334
Chile	Very Unhealthy	235.75
South Africa	Very Unhealthy	233.6
Mexico	Very Unhealthy	233.42105263157896
India	Very Unhealthy	232.91603053435114
Pakistan	Very Unhealthy	227.45833333333334
Mauritania	Very Unhealthy	224.0
Indonesia	Very Unhealthy	223.0
Nepal	Very Unhealthy	222.33333333333334
Iran (Islamic Republic of)	Very Unhealthy	220.25
Namibia	Very Unhealthy	218.0
Russian Federation	Very Unhealthy	214.0
Bangladesh	Very Unhealthy	214.0
Senegal	Very Unhealthy	211.0
Nigeria	Very Unhealthy	210.66666666666666
China	Very Unhealthy	210.04444444444445
Malaysia	Very Unhealthy	209.0
	Very Unhealthy	202.0
Guinea-Bissau	Unhealthy	195.0
Bahrain	Unhealthy	188.0
United Republic of Tanzania	Unhealthy	183.0
Peru	Unhealthy	179.0
Democratic Republic of the Congo	Unhealthy	175.125
Chile	Unhealthy	173.0
Libya	Unhealthy	172.71428571428572
Saudi Arabia	Unhealthy	172.33333333333334
Angola	Unhealthy	172.0
Russian Federation	Unhealthy	172.0
Mexico	Unhealthy	171.56976744186048

Storing the output in HDFS

hdfs dfs -ls /user/cloudera/hive_output/avg_aqi_country_category

hdfs dfs -cat /user/cloudera/hive_output/avg_aqi_country_category/000000_0 | head

```
1 cloudera cloudera      17326 2025-11-20 02:38 /user/cloudera/hive_output/avg_aqi_country_category/000000_0
[cloudera@quickstart ~]$ hdfs dfs -cat /user/cloudera/hive_output/avg_aqi_country_category/000000_0 | head
United States of America,Hazardous,500.0
Russian Federation,Hazardous,500.0
India,Hazardous,450.62025316455697
Republic of Korea,Hazardous,421.0
Pakistan,Hazardous,419.46153846153845
Democratic Republic of the Congo,Hazardous,415.0
South Africa,Hazardous,403.0
Chile,Hazardous,358.0
Mexico,Hazardous,339.6666666666667
China,Hazardous,332.3333333333333
cat: Unable to write to output stream.
[cloudera@quickstart ~]$ █
```

Unhealthy Cities % Contribution by Country

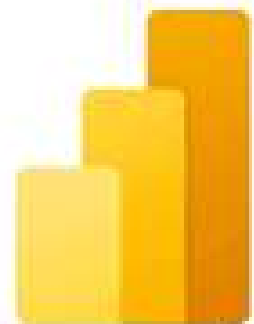
INSERT OVERWRITE DIRECTORY
'/user/cloudera/hive_output/unhealthy_percentage'
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
SELECT country,
ROUND((COUNT(*) * 100.0) / total.total_count, 2) AS
unhealthy_percentage
FROM air_table1
CROSS JOIN
(
SELECT COUNT(*) AS total_count
FROM air_table1
) total
WHERE aqi_category = 'Unhealthy'
GROUP BY country, total.total_count
ORDER BY unhealthy_percentage DESC;

country	unhealthy_percentage
India	4.69
China	1.25
Pakistan	0.98
Mexico	0.37
Indonesia	0.31
Brazil	0.19
Iran (Islamic Republic of)	0.17
South Africa	0.15
Bangladesh	0.11
Senegal	0.1
Nigeria	0.09
United States of America	0.08
Chile	0.06
Viet Nam	0.06
Uzbekistan	0.06
Venezuela (Bolivarian Republic of)	0.04
Philippines	0.04
Saudi Arabia	0.04
Dominican Republic	0.04
Libya	0.03
Uganda	0.03
Nepal	0.03
Malaysia	0.03
Kenya	0.03
Gambia	0.03
France	0.03
Egypt	0.03
Democratic Republic of the Congo	0.03
Angola	0.02
Afghanistan	0.02
Ethiopia	0.02
Oman	0.02
Yemen	0.02
United Republic of Tanzania	0.01
Mauritania	0.01
Mali	0.01

CO vs Ozone pollution Comparison per Country.

```
SELECT
  country,
  ROUND(AVG(co_value), 2) AS avg_co,
  ROUND(AVG(ozone_value), 2) AS avg_ozone
FROM air_table1
GROUP BY country
ORDER BY avg_co DESC;
```

country	avg_co	avg_ozone
Republic of Korea	27.0	0.0
South Africa	5.38	16.67
Democratic Republic of the Congo	5.29	27.14
Kingdom of Eswatini	4.67	15.67
Nigeria	3.81	24.85
Chile	3.8	10.05
Rwanda	3.62	22.08
China	3.42	88.32
Angola	3.15	22.7
Lesotho	2.8	14.0
Kuwait	2.67	135.67
Uganda	2.65	20.25
Burundi	2.62	21.38
El Salvador	2.58	13.64
Congo	2.58	23.83
Nepal	2.45	67.58
Bangladesh	2.44	46.2
Indonesia	2.42	46.57
Kenya	2.37	21.3
Viet Nam	2.37	46.23
Guatemala	2.33	11.16
Cameroon	2.33	17.05
Malaysia	2.23	36.5
Pakistan	2.07	89.14
Central African Republic	2.0	12.23
Bahrain	2.0	127.0
Algeria	1.92	47.17
Dominican Republic	1.82	17.63
Panama	1.78	18.04
Mexico	1.75	16.82
India	1.74	55.06
Haiti	1.72	22.94
Colombia	1.7	9.78
Chad	1.58	19.21

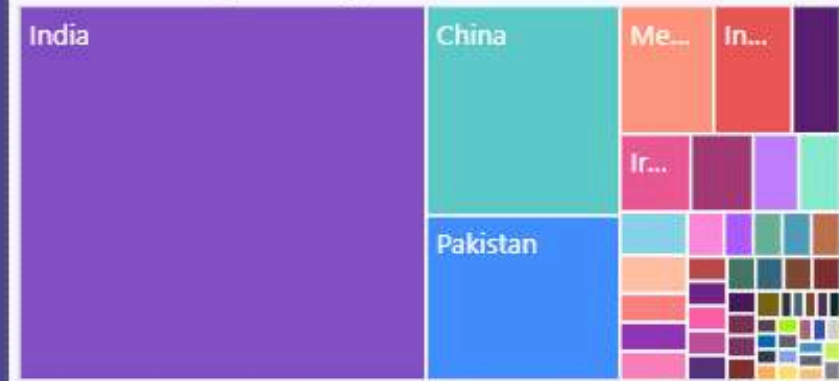


Power BI

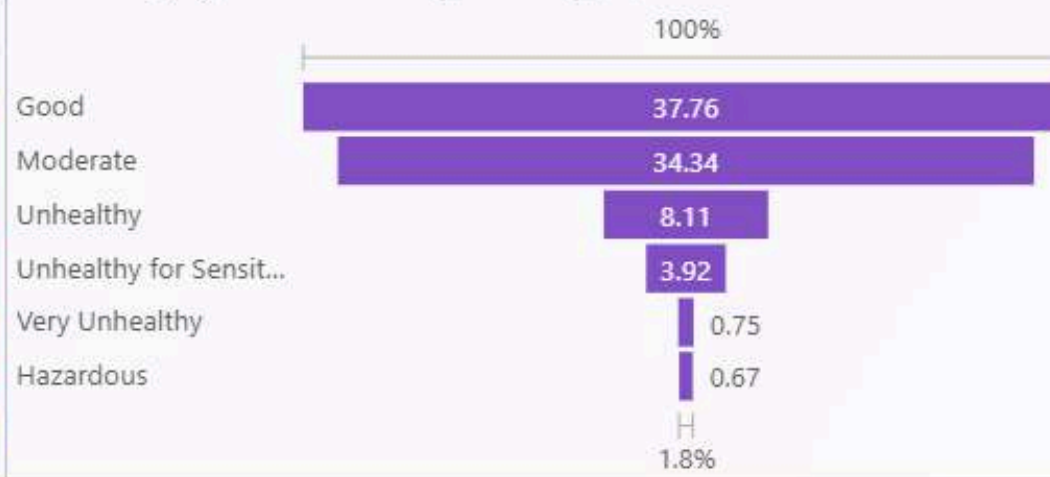
POWER BI

Air Pollution Trend Analysis

Sum of AQI by country



AQI Category Distribution by Country (%)



NUMBER OF GOOD AQI

18

NUMBER OF BAD AQI

9

country_name

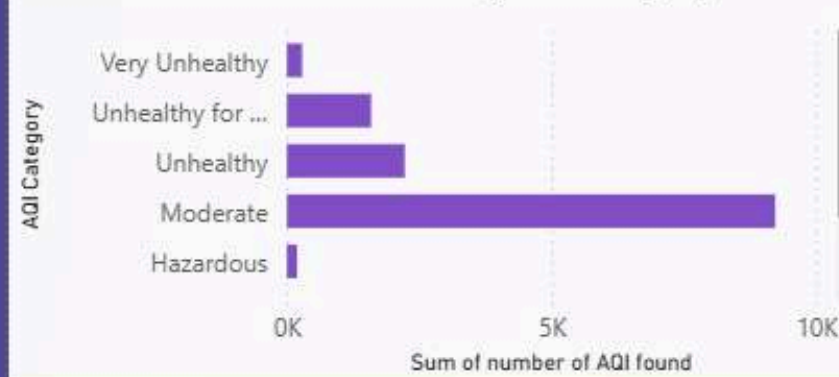
All

AQI CATEGORY

aqi.aqi_category

All

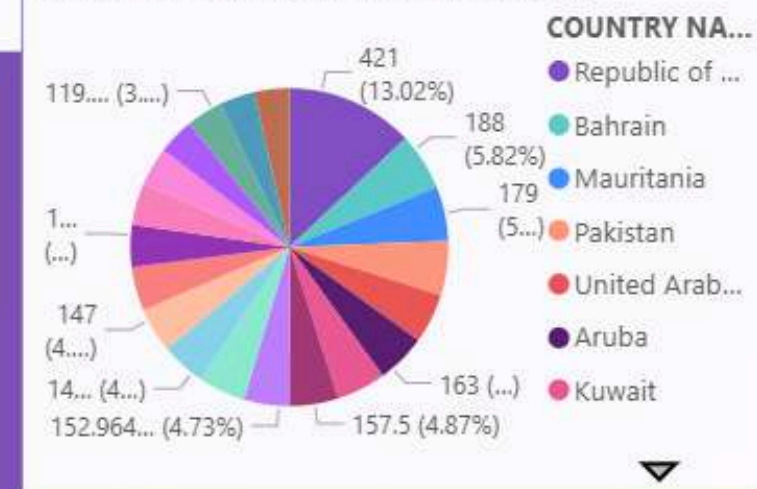
Sum of number of AQI found by AQI Category



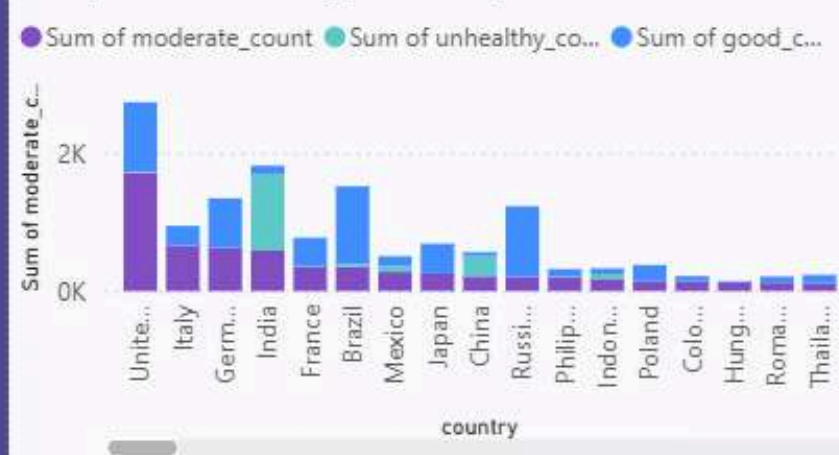
WORST AQI



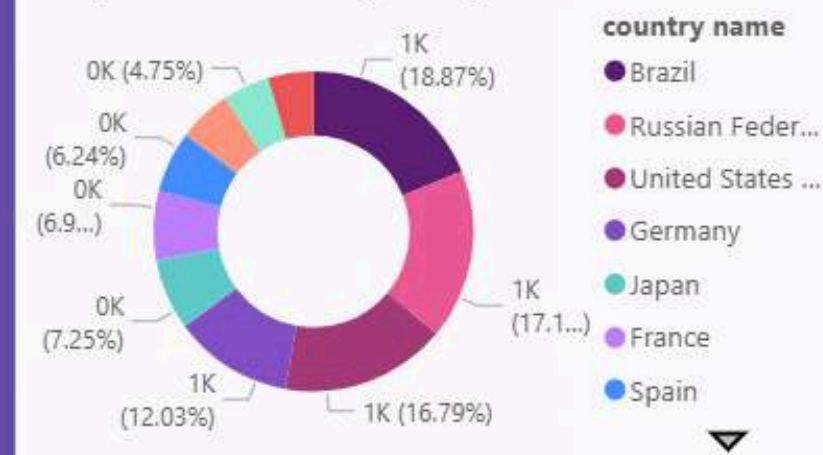
Global Air Quality Trends by Country



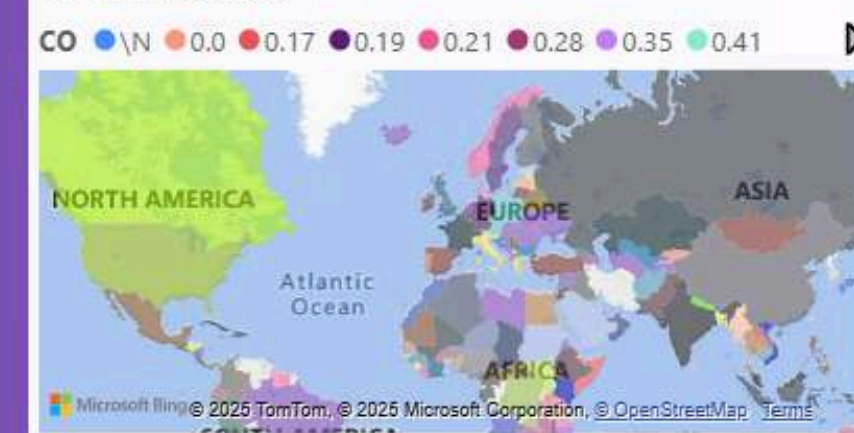
composition of the type of AQI present



composition of country with good AQI



CO Composition



THANK YOU

