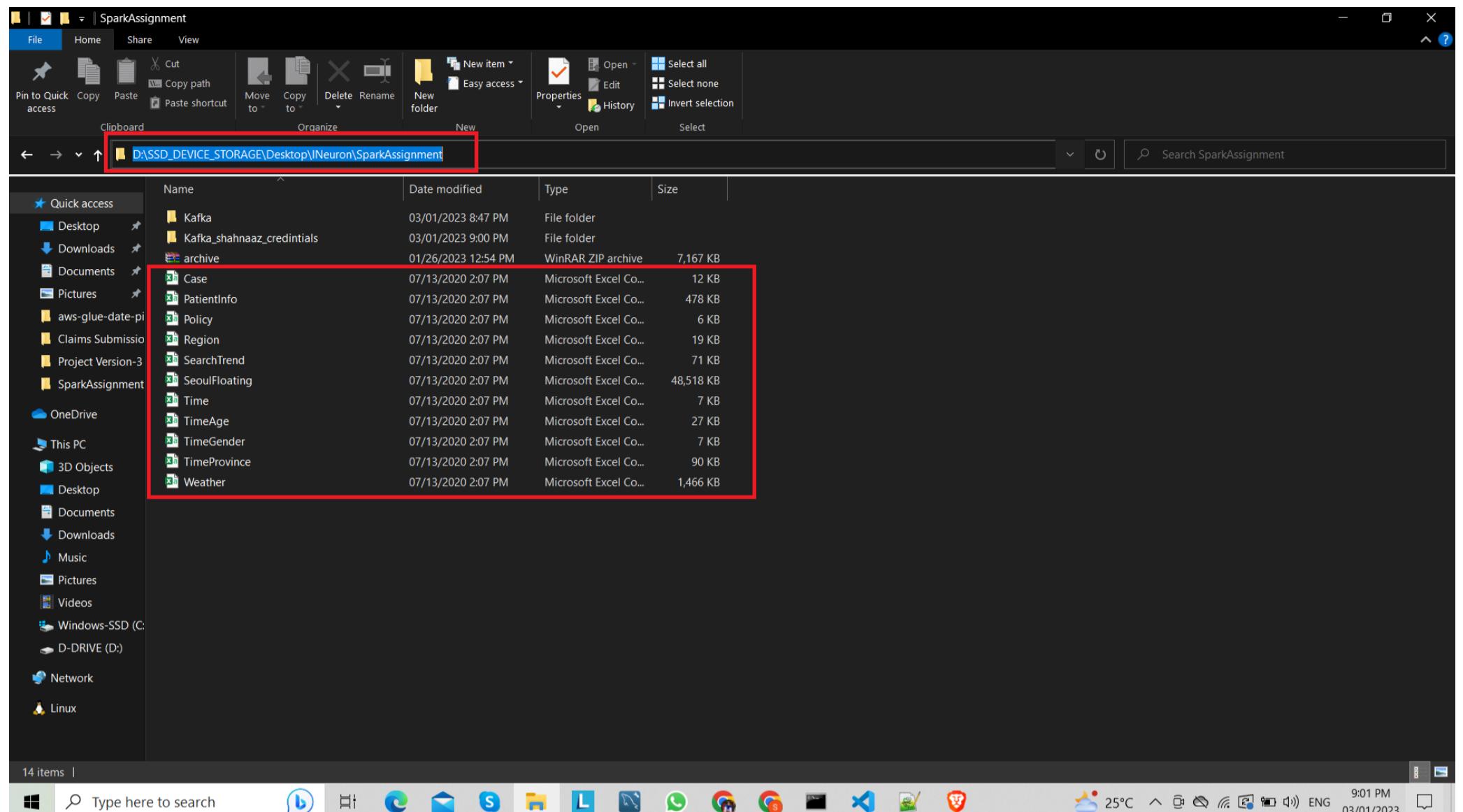


1. Download the data from the given URL :

<https://www.kaggle.com/datasets/kimjihoo/coronavirusdataset>



2. Create a producer with a python connector in confluent Kafka and stream your data.

The screenshot shows the Confluent Cloud Topics interface. The left sidebar includes options like Cluster Overview, Dashboard, Networking, API Keys, Cluster Settings, Stream Lineage, Stream Designer, Topics (which is selected), logDB, Connectors, Clients, and Schema Registry. The main area displays a table of topics. A red box highlights the 'patientinfo' topic in the list. The table columns include Topic name, Partitions, Production (last min), Consumption (last min), and Schema. Most topics have 'Edit schema' links, except for 'case' which has 'Set a schema'.

| Topic name | Partitions | Production (last min) | Consumption (last min) | Schema |
|---------------|------------|-----------------------|------------------------|--------------|
| case | 6 | -- | 0B/s | Edit schema |
| dim_customers | 6 | -- | -- | Set a schema |
| dim_dates | 6 | -- | -- | Set a schema |
| dim_employees | 6 | -- | -- | Set a schema |
| dim_locations | 6 | -- | -- | Set a schema |
| dim_products | 6 | -- | -- | Set a schema |
| fact_sales | 6 | -- | -- | Set a schema |
| patientinfo | 6 | 0B/s | 0B/s | Edit schema |
| policy | 6 | 0B/s | 0B/s | Edit schema |
| region | 6 | 0B/s | 0B/s | Edit schema |
| search_trend | 6 | 0B/s | 0B/s | Edit schema |
| SeoulFloating | 6 | 0B/s | 0B/s | Edit schema |
| Time | 6 | 0B/s | 0B/s | Edit schema |
| time_province | 6 | 0B/s | 0B/s | Edit schema |
| timeage | 6 | 0B/s | 0B/s | Edit schema |
| timegender | 6 | 0B/s | 0B/s | Edit schema |
| weather | 6 | 0B/s | 0B/s | Edit schema |

Topics - Confluent Cloud

confluent.cloud/environments/env-63ovw6/clusters/lkc-knoy36/topics/case/message-viewer

Stream Catalog LEARN 🔍 ⓘ ⚙️

case

Overview Messages Schema Configuration

Producers Consumers

Message fields

topic partition offset timestamp timestampType headers key value

case_id province city group

Description Add description

Tags Add tags to this topic

Add business metadata

Date created Mar 1 2023 12:21 PM

Date modified --

Retention time 1 week

Retention size Infinite

Number of partitions 6

Cleanup policy delete

JSON CSV Download

Topics - Confluent Cloud

confluent.cloud/environments/env-63ovw6/clusters/lkc-knoy36/topics/patientinfo/message-viewer

Stream Catalog LEARN 🔍 ⓘ ⚙️

patientinfo

Overview Messages Schema Configuration

Producers Consumers

Message fields

topic partition offset timestamp timestampType headers key value

patient_id sex age country

Description Add description

Tags Add tags to this topic

Add business metadata

Date created Mar 1 2023 7:02 PM

Date modified --

Retention time 1 week

Retention size Infinite

Number of partitions 6

Cleanup policy delete

JSON CSV Download

Topics - Confluent Cloud

confluent.cloud/environments/env-63ovw6/clusters/lkc-knoy36/topics/policy/message-viewer

Stream Catalog LEARN 🔔 ? ☰

policy

Overview Messages Schema Configuration

Producers Consumers

Message fields

Value

policy_id country type gov_policy detail

Bytes in/sec Bytes out/sec

Description Tags Date created Date modified Retention time Retention size Number of partitions Cleanup policy

JSON CSV Download

Explore Stream Lineage

Add description Add tags to this topic Add business metadata

Mar 1 2023 7:52 PM

1 week

Infinite

6

delete

Topics - Confluent Cloud

confluent.cloud/environments/env-63ovw6/clusters/lkc-knoy36/topics/region/message-viewer

Stream Catalog LEARN 🔔 ? ☰

region

Overview Messages Schema Configuration

Producers Consumers

Message fields

Value

code province city latitude longitude elementary_school_count

Bytes in/sec Bytes out/sec

Description Tags Date created Date modified Retention time Retention size Number of partitions Cleanup policy

25°C ENG 9:07 PM 03/01/2023

Add description Add tags to this topic Add business metadata

Mar 1 2023 7:58 PM

1 week

Infinite

6

delete

Topics - Confluent Cloud

confluent.cloud/environments/env-63ovw6/clusters/lkc-knoy36/topics/search_trend/message-viewer

Paused

Stream Catalog LEARN

search_trend

Overview Messages Schema Configuration

Producers Consumers

Message fields

Topic Partition Offset Timestamp

Value

Date Cold Flu Pneumonia

Description Tags Date created Date modified Retention time Retention size Number of partitions Cleanup policy

JSON CSV Download

Topics - Confluent Cloud

confluent.cloud/environments/env-63ovw6/clusters/lkc-knoy36/topics/SeoulFloating/message-viewer

Paused

Stream Catalog LEARN

SeoulFloating

Overview Messages Schema Configuration

Producers Consumers

Message fields

Topic Partition Offset Timestamp

Value

Date Hour Birth Year Sex

Description Tags Date created Date modified Retention time Retention size Number of partitions Cleanup policy

JSON CSV Download

Topics - Confluent Cloud

confluent.cloud/environments/env-63ovw6/clusters/lkc-knoy36/topics/Time/message-viewer

Stream Catalog LEARN 🔔 ? ☰

Time

Overview Messages Schema Configuration

Producers Consumers

Message fields

topic partition offset timestamp timestampType headers key value date time test negative

Description Add description

Tags Add tags to this topic

Add business metadata

Date created Mar 1 2023 8:22 PM

Date modified --

Retention time 1 week

Retention size Infinite

Number of partitions 6

Cleanup policy delete

JSON CSV Download

Topics - Confluent Cloud

confluent.cloud/environments/env-63ovw6/clusters/lkc-knoy36/topics/time_province/message-viewer

Stream Catalog LEARN 🔔 ? ☰

time_province

Overview Messages Schema Configuration

Producers Consumers

Message fields

topic partition offset timestamp timestampType headers key value date time province confirmed

Description Add description

Tags Add tags to this topic

Add business metadata

Date created Mar 1 2023 8:43 PM

Date modified --

Retention time 1 week

Retention size Infinite

Number of partitions 6

Cleanup policy delete

JSON CSV Download

Topics - Confluent Cloud

confluent.cloud/environments/env-63ovw6/clusters/lkc-knoy36/topics/timeage/message-viewer

Paused

Stream Catalog LEARN

timeage

Overview Messages Schema Configuration

Producers Consumers

Message fields

topic partition offset timestamp timestampType headers key value date time age confirmed

Description Add description

Tags Add tags to this topic

Add business metadata

Date created Mar 1 2023 8:27 PM

Date modified --

Retention time 1 week

Retention size Infinite

Number of partitions 6

Cleanup policy delete

JSON CSV Download

Topics - Confluent Cloud

confluent.cloud/environments/env-63ovw6/clusters/lkc-knoy36/topics/timegender/message-viewer

Paused

Stream Catalog LEARN

timegender

Overview Messages Schema Configuration

Producers Consumers

Message fields

topic partition offset timestamp timestampType headers key value date time sex confirmed

Description Add description

Tags Add tags to this topic

Add business metadata

Date created Mar 1 2023 8:37 PM

Date modified --

Retention time 1 week

Retention size Infinite

Number of partitions 6

Cleanup policy delete

JSON CSV Download

Topics - Confluent Cloud

confluent.cloud/environments/env-63ovw6/clusters/lkc-knoy36/topics/weather/message-viewer

mystartsearch Powered by Browse... New Tab YouTube Maps

CONFLUENT

HOME > ENVIRONMENTS > DEFAULT > CLUSTER_1 > TOPICS >

weather

Overview Messages Schema Configuration

Producers Consumers

Message fields

Topic Partition Offset Timestamp

Code Province Date Avg Temp Min Temp Max Temp Precipitation Max Wind Speed

30000 Gangwon-do 2016-01-18 -1.1 -9.7 4.5 0 4.0

13000 Gwangju 2016-01-18 -0.8 -4.4 3.2 4.1 4.0

14000 Incheon 2016-01-17 2.5 1 5.2 0 4.7

70000 Jeju-do 2016-01-16 7.5 4.2 11.1 0 4.3

Description Add description

Tags Add tags to this topic

Add business metadata

Date created Mar 1 2023 8:49 PM

Date modified --

Retention time 1 week

Retention size Infinite

Number of partitions 6

Cleanup policy delete

JSON CSV Download

Type here to search

25°C ENG 9:11 PM 03/01/2023

3. Consume your data through the python connector and dump it in mongo dB atlas.

Note: Here in the dataset you will be finding a multiple files you need to use all file for the Kafka and mongo dB

The screenshot shows the Confluent Cloud interface for managing connectors. On the left sidebar, under the 'Connectors' section, the 'MongoDBAtlasSinkConnector_0' is highlighted. The main panel displays the connector's details:

- MongoDBAtlasSinkConnector_0** (Running)
- Tasks**: 1, **Bytes/sec**: 0B/s
- Messages/sec**: 0, **Messages behind**: 0
- Overview** table:
 - Category: Sink
 - ID: lcc-j37k8q
 - Plugin name: MongoDB Atlas Sink

Below the connector details, there is a 'Connect with popular connectors' section featuring icons for Snowflake Sink, Google Cloud Storage Sink, Elasticsearch Service Sink, and MongoDB Atlas Source.

The screenshot shows the configuration page for the 'MongoDbAtlasSinkConnector_0'. The 'Topics' section has 'Topic names' set to 'case'. The 'Authentication' section has 'Connection host' as 'cluster0.wri9uwx.mongodb.net', 'Connection user' as 'mongodb', and 'Database name' set to 'kafka'. The 'Configuration' section has 'Input Kafka record value format' set to 'JSON_SR'. The 'Apply changes' button is visible at the bottom right.

Syedabdul's ... Access Manager Billing All Clusters Get Help Syedabdul

Project 0 Data Services App Services Charts

DEPLOYMENT Database SERVICES Triggers Data API Data Federation Search SECURITY Database Access Network Access Advanced Goto

SYEDABDUL'S ORG - 2022-12-22 > PROJECT 0 > DATABASES

ClusterO

Collections Overview Real Time Metrics Find Indexes Schema Anti-Patterns Aggregation Search Indexes

DATABASES: 1 COLLECTIONS: 1

+ Create Database Search Namespaces kafka case

STORAGE SIZE: 28KB LOGICAL DATA SIZE: 32.1KB TOTAL DOCUMENTS: 174 INDEXES TOTAL SIZE: 24KB

FIND FILTER { field: 'value' } OPTIONS Apply Reset

QUERY RESULTS: 161-174 OF 174

```
_id: ObjectId('6400d035e4b28f74a00fdcaa7')
infection_case: "etc"
province: "Sejong"
city: ""
latitude: ""
case_id: 1700006
confirmed: 1
group: false
longitude: ""
```

```
_id: ObjectId('6400d035e4b28f74a00fdcaa8')
infection_case: "Anyang Gunpo Pastors Group"
province: "Gyeonggi-do"
city: "Anyang-si"
latitude: "37.381784"
case_id: 2000011
```

Type here to search 26°C Mostly cloudy 10:15 PM 03/02/2023

Invensis Technologies Pvt. Ltd. - Mail - abdulqadir.syed - Network hcare - INV_TEAM - Slack Data | Cloud: MongoDB Cloud

Atlas Syedabdul's ... Access Manager Billing All Clusters Get Help Syedabdul

Project 0 Data Services App Services Charts

DEPLOYMENT Database SERVICES Triggers Data API Data Federation Search SECURITY Database Access Network Access Advanced Goto

SYEDABDUL'S ORG - 2022-12-22 > PROJECT 0 > DATABASES

ClusterO

Collections Overview Real Time Metrics Find Indexes Schema Anti-Patterns Aggregation Search Indexes

DATABASES: 1 COLLECTIONS: 2

+ Create Database Search Namespaces kafka case patientinfo

STORAGE SIZE: 712KB LOGICAL DATA SIZE: 1.69MB TOTAL DOCUMENTS: 5165 INDEXES TOTAL SIZE: 332KB

FIND FILTER { field: 'value' } OPTIONS Apply Reset

QUERY RESULTS: 1-20 OF MANY

```
_id: ObjectId('6400d327e4b28f74a00fdaba')
country: "Korea"
infection_case: "contact with patient"
city: "Jongno-gu"
sex: "male"
confirmed_date: "2020-01-30"
released_date: "2020-02-19"
contact_number: "17"
deceased_date: "NAN"
infected_by: "2002000001"
province: "Seoul"
symptom_onset_date: "NAN"
patient_id: 1000000003
state: "released"
age: "50s"
```

Invensis Technologies Pvt. Ltd. - Mail - abdulqadir.syed - Network - hc care - INV_TEAM - Slack - Data | Cloud: MongoDB Cloud

cloud.mongodb.com/v2/63a3f850eb7a144c3202b002#metrics/replicaSet/6400bb19de519e522e409747/explorer/kafka/policy/find

All Clusters Get Help Syedabdul

Atlas Syedabdul's ... Access Manager Billing

Project 0 Data Services App Services Charts

DEPLOYMENT Database SERVICES

Database PREVIEW

Trigger Data API Data Federation Search SECURITY

Database Access Network Access Advanced Goto

SYEDABDUL'S ORG - 2022-12-22 > PROJECT 0 > DATABASES

ClusterO

Collections Overview Real Time Metrics

DATABASES: 1 COLLECTIONS: 3

+ Create Database Search Namespaces kafka policy

TOTAL DOCUMENTS: 61 INDEXES TOTAL SIZE: 20KB

Find Indexes Schema Anti-Patterns Aggregation Search Indexes

INSERT DOCUMENT

FILTER { field: 'value' } OPTIONS Apply Reset

QUERY RESULTS: 1-20 OF MANY

```
_id: ObjectId('6400d543e4b28f74a00feee8')
end_date: "NAN"
country: "Korea"
gov_policy: "Special Immigration Procedure"
policy_id: 6
detail: "from Hong Kong"
type: "Immigration"
start_date: "2020-02-12"
```

```
_id: ObjectId('6400d543e4b28f74a00feee9')
end_date: "NAN"
country: "Korea"
gov_policy: "Emergency Use Authorization of Diagnostic Kit"
policy_id: 23
detail: "ath EUA"
type: "Health"
```

Invensis Technologies Pvt. Ltd. - Mail - abdulqadir.syed - Network - hc care - INV_TEAM - Slack - Data | Cloud: MongoDB Cloud

cloud.mongodb.com/v2/63a3f850eb7a144c3202b002#metrics/replicaSet/6400bb19de519e522e409747/explorer/kafka/region/find

All Clusters Get Help Syedabdul

Atlas Syedabdul's ... Access Manager Billing

Project 0 Data Services App Services Charts

DEPLOYMENT Database SERVICES

Database PREVIEW

Trigger Data API Data Federation Search SECURITY

Database Access Network Access Advanced Goto

SYEDABDUL'S ORG - 2022-12-22 > PROJECT 0 > DATABASES

ClusterO

Collections Overview Real Time Metrics

DATABASES: 1 COLLECTIONS: 4

+ Create Database Search Namespaces kafka region

TOTAL DOCUMENTS: 244 INDEXES TOTAL SIZE: 24KB

Find Indexes Schema Anti-Patterns Aggregation Search Indexes

INSERT DOCUMENT

FILTER { field: 'value' } OPTIONS Apply Reset

QUERY RESULTS: 1-20 OF MANY

```
_id: ObjectId('6400d66ee4b28f74a00fef26')
elderly_population_ratio: 15.38
code: 10000
province: "Seoul"
elementary_school_count: 607
city: "Seoul"
latitude: 37.566953
kindergarten_count: 830
university_count: 48
academy_ratio: 1.44
nursing_home_count: 22739
elderly_alone_ratio: 5.8
longitude: 126.977977
```

```
_id: ObjectId('6400d66ee4b28f74a00fef27')
elderly_population_ratio: 14.39
```

Syedabdul's org - 2022-12-22 > Project 0 > Databases

VERSION 5.0.15 REGION GCP Iowa (us-central)

Database PREVIEW

Services

Triggers

Data API

Data Federation

Search

Security

Database Access

Network Access

Advanced

Goto

DEPLOYMENT

Collections

Overview Real Time Metrics Collections Search Profiler Performance Advisor Online Archive Cmd Line Tools

kafka.search_trend

STORAGE SIZE: 108KB LOGICAL DATA SIZE: 176.39KB TOTAL DOCUMENTS: 1642 INDEXES TOTAL SIZE: 64KB

Find Indexes Schema Anti-Patterns Aggregation Search Indexes

search_trend

query results: 1-20 of many

```
_id: ObjectId('6400d83be4b28f74a00ff01b')
date: "2016-01-02"
pneumonia: 0.28826
coronavirus: 0.0089
cold: 0.13372
flu: 0.17135

_id: ObjectId('6400d83be4b28f74a00ff01c')
date: "2016-01-04"
pneumonia: 0.29008
coronavirus: 0.01145
cold: 0.17463
flu: 0.18626
```

Type here to search 24°C Mostly cloudy 10:41 PM 03/02/2023

Syedabdul's org - 2022-12-22 > Project 0 > Databases

VERSION 5.0.15 REGION GCP Iowa (us-central)

Database PREVIEW

Services

Triggers

Data API

Data Federation

Search

Security

Database Access

Network Access

Advanced

Goto

DEPLOYMENT

Collections

Overview Real Time Metrics Collections Search Profiler Performance Advisor Online Archive Cmd Line Tools

kafka.Time

STORAGE SIZE: 4KB LOGICAL DATA SIZE: 22.92KB TOTAL DOCUMENTS: 163 INDEXES TOTAL SIZE: 4KB

Find Indexes Schema Anti-Patterns Aggregation Search Indexes

Time

query results: 1-20 of many

```
_id: ObjectId('6400d905e4b28f74a00ff686')
date: "2020-01-29"
negative: 155
deceased: 0
test: 187
time: 16
confirmed: 4
released: 0

_id: ObjectId('6400d905e4b28f74a00ff687')
date: "2020-01-31"
negative: 245
deceased: 0
test: 312
time: 16
confirmed: 11
```

Syedabdul's ... Access Manager Billing All Clusters Get Help Syedabdul

Project 0 Data Services App Services Charts

DEPLOYMENT Database SERVICES Triggers Data API Data Federation Search SECURITY Database Access Network Access Advanced Goto

Atlas Syedabdul's ... Access Manager Billing All Clusters Get Help Syedabdul

SYEDABDUL'S ORG - 2022-12-22 > PROJECT 0 > DATABASES ClusterO

VERSION 5.0.15 REGION GCP Iowa (us-central)

Data Lake PREVIEW

Database Overview Real Time Metrics Collections Search Profiler Performance Advisor Online Archive Cmd Line Tools

+ Create Database Search Namespaces kafka timeage patientinfo polioy region search_trend timeage

TOTAL DOCUMENTS: 1089 INDEXES TOTAL SIZE: 4KB

Find Indexes Schema Anti-Patterns Aggregation Search Indexes

INSERT DOCUMENT

FILTER { field: 'value' } OPTIONS Apply Reset

QUERY RESULTS: 1-20 OF MANY

```
_id: ObjectId('6400d985e4b28f74a00ff72a')
date: "2020-03-02"
deceased: 0
time: 0
confirmed: 32
age: "0s"

_id: ObjectId('6400d985e4b28f74a00ff72b')
date: "2020-03-02"
deceased: 0
time: 0
confirmed: 169
age: "10s"
```

Type here to search 24°C Mostly cloudy 10:45 PM 03/02/2023

Invensis Technologies Pvt. Ltd. - Mail - abdulqadir.syed - Network Slack Data | Cloud: MongoDB Cloud

Atlas Syedabdul's ... Access Manager Billing All Clusters Get Help Syedabdul

SYEDABDUL'S ORG - 2022-12-22 > PROJECT 0 > DATABASES ClusterO

VERSION 5.0.15 REGION GCP Iowa (us-central)

Project 0 Data Services App Services Charts

DEPLOYMENT Database SERVICES Triggers Data API Data Federation Search SECURITY Database Access Network Access Advanced Goto

Atlas Syedabdul's ... Access Manager Billing All Clusters Get Help Syedabdul

SYEDABDUL'S ORG - 2022-12-22 > PROJECT 0 > DATABASES ClusterO

VERSION 5.0.15 REGION GCP Iowa (us-central)

Data Lake PREVIEW

Database Overview Real Time Metrics Collections Search Profiler Performance Advisor Online Archive Cmd Line Tools

+ Create Database Search Namespaces kafka timeage patientinfo polioy region search_trend timeage timegender

TOTAL DOCUMENTS: 242 INDEXES TOTAL SIZE: 4KB

Find Indexes Schema Anti-Patterns Aggregation Search Indexes

INSERT DOCUMENT

FILTER { field: 'value' } OPTIONS Apply Reset

QUERY RESULTS: 1-20 OF MANY

```
_id: ObjectId('6400d9fde4b28f74a00ffb6c')
date: "2020-03-04"
deceased: 20
sex: "male"
time: "0"
confirmed: 1996

_id: ObjectId('6400d9fde4b28f74a00ffb6d')
date: "2020-03-04"
deceased: 12
sex: "female"
time: "0"
confirmed: 3332
```

Syedabdul's org - 2022-12-22 > Project 0 > DATABASES

SYEDABDUL'S ORG - 2022-12-22 > PROJECT 0 > DATABASES

VERSION 5.0.15 REGION GCP Iowa (us-central)

Database PREVIEW

Services

Triggers

Data API

Data Federation

Search

Security

Database Access

Network Access

Advanced

Goto

DEPLOYMENT

Collections

Overview Real Time Metrics Collections Search Profiler Performance Advisor Online Archive Cmd Line Tools

kafka.time_province

STORAGE SIZE: 112KB LOGICAL DATA SIZE: 365.22KB TOTAL DOCUMENTS: 2771 INDEXES TOTAL SIZE: 100KB

Find Indexes Schema Anti-Patterns Aggregation Search Indexes

INSERT DOCUMENT

FILTER { field: 'value' }

OPTIONS Apply Reset

QUERY RESULTS: 1-20 OF MANY

```
_id: ObjectId('6400da4fe4b28f74a00fffc5f')
date: "2020-01-20"
deceased: 0
province: "Gwangju"
time: "16"
confirmed: 0
released: 0

_id: ObjectId('6400da4fe4b28f74a00fffc60')
date: "2020-01-20"
deceased: 0
province: "Chungcheongbuk-do"
time: "16"
confirmed: 0
released: 0
```

Type here to search 24°C Mostly cloudy 10:51 PM 03/02/2023

Syedabdul's org - 2022-12-22 > Project 0 > DATABASES

SYEDABDUL'S ORG - 2022-12-22 > PROJECT 0 > DATABASES

VERSION 5.0.15 REGION GCP Iowa (us-central)

Database PREVIEW

Services

Triggers

Data API

Data Federation

Search

Security

Database Access

Network Access

Advanced

Goto

DEPLOYMENT

Collections

Overview Real Time Metrics Collections Search Profiler Performance Advisor Online Archive Cmd Line Tools

kafka.weather

STORAGE SIZE: 2.71MB LOGICAL DATA SIZE: 6.09MB TOTAL DOCUMENTS: 26271 INDEXES TOTAL SIZE: 1.48MB

Find Indexes Schema Anti-Patterns Aggregation Search Indexes

INSERT DOCUMENT

FILTER { field: 'value' }

OPTIONS Apply Reset

QUERY RESULTS: 1-20 OF MANY

```
_id: ObjectId('6400dae3e4b28f74a0100733')
date: "2016-01-01"
precipitation: 0
avg_relative_humidity: 52.1
code: 11000
province: "Busan"
avg_temp: 5.3
min_temp: 1.1
most_wind_direction: 340
max_wind_speed: 7.4
max_temp: 10.9

_id: ObjectId('6400dae3e4b28f74a0100734')
date: "2016-01-01"
precipitation: 0
avg_relative_humidity: 70.5
```

The screenshot shows the MongoDB Atlas interface for the 'kafka.SeoulFloating' collection. The collection has 1085003 documents. Two sample documents are displayed:

```

_id: ObjectId('6400dbe6e4b28f74a0106dd7')
date: "2020-01-01"
hour: 0
province: "Seoul"
fp_num: 28880
city: "Dongjag-gu"
sex: "female"
birth_year: 20

_id: ObjectId('6400dbe6e4b28f74a0106dd8')
date: "2020-01-01"
hour: 0
province: "Seoul"
fp_num: 38300
city: "Gangseo-gu"
sex: "male"
...

```

4. Collect your data as a pyspark dataframe and perform different operations.

Note: Consider only three files for creating a data frame among all case, region and Time Province

- Read the data, show it and Count the number of records.
- Describe the data with a describe function.
- If there is any duplicate value drop it.
- Use limit function for showcasing a limited number of records.
- If you find the column name is not suitable, change the column name.[optional]
- Select the subset of the columns.
- If there is any null value, fill it with any random value or drop it.
- Filter the data based on different columns or variables and do the best analysis.

For example: We can filter a data frame using multiple conditions using AND(&), OR(|) and NOT(~) conditions. For example, we may want to find out all the different infection case in Daegu Province with more than 10 confirmed cases.

- Sort the number of confirmed cases. Confirmed column is there in the dataset. Check with descending sort also.
- In case of any wrong data type, cast that data type from integer to string or string to integer.
- Use group by on top of province and city column and agg it with sum of confirmed cases. For example
df.groupBy(["province", "city"]).agg(function.sum("confirmed"))
- For joins we will need one more file. you can use region file. User different different join methods, for example
cases.join(regions, ['province', 'city'], how='left') You can do your best analysis.

5. If you want, you can also use SQL with data frames. Let us try to run some SQL on the cases table.

For example:

```
cases.registerTempTable('cases table')
```

```
newDF = sqlContext.sql('select * from cases table where confirmed>100') newDF.show()
```

Here is a example how you can use df for sql now you can perform
various operations with GROUP BY, HAVING, AND ORDER BY

ABOVE MENTION QUERIES ARE IN ,Mongodb_Pyspark.ipynb file