

NRT Analysis & Processing of Streaming Data.

Aim : Develop a robust, fault tolerant, scalable solution to process, analyze, clean and store streaming data from various streaming applications and devices.

Problems in traditional Big Data Systems : Traditional Big data solutions have solved various problems of harnessing voluminous data. But organizations have been facing various challenges such as,

- **Harnessing Streaming data** : Data has become the key driver for businesses and it is the need of the hour for these businesses to harness and effectively store and use this data. Data is generated every second from devices and applications and it is important for organizations to collect and store this streaming data.
- **Building Scalability, Robustness & Maintainability** : Scalability, Robustness and Maintainability of the system has been a challenge for organizations. Hence the solution developed to process, analyze and store streaming data must solve these three issues.
- **Performing NRT Analysis and Data warehousing on the cloud** : Big data technologies have been evolving and many organizations have adopted these solutions. Business insights can be gained from the data that has been flowing into the business systems. It is the need of the hour to have solutions providing Near Real time Analysis of data as that can be used to handle and take more effective business decisions.

Approach : We aim to solve the problems mentioned above by using cloud services and evolved Big data technologies such as Spark.

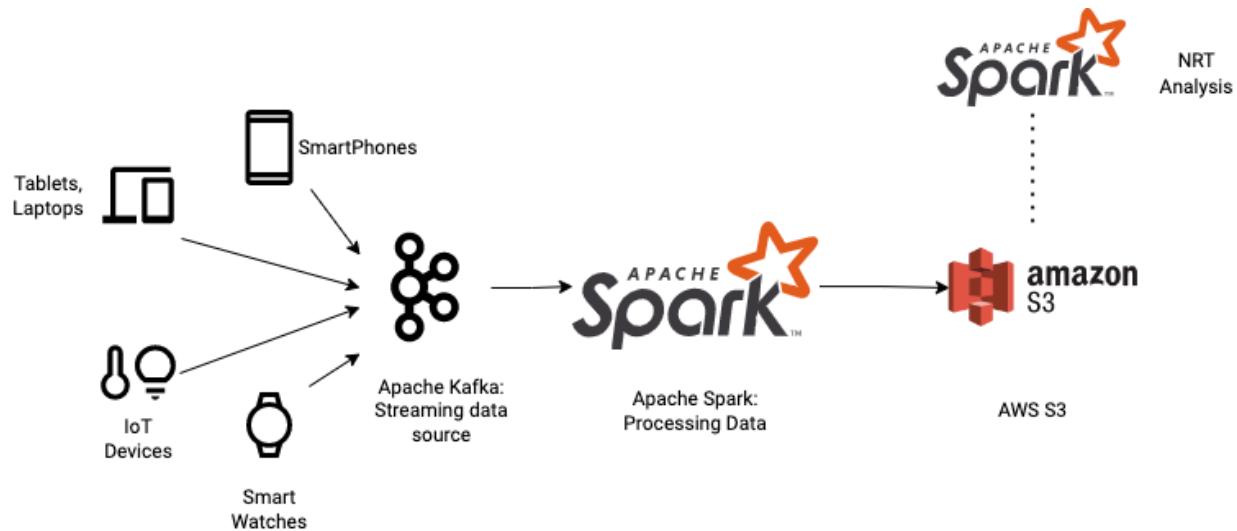
Functionalities :

- This solution would enable organizations & businesses to harness the streaming data in a cost effective manner by using managed services like Confluent (managed Apache Kafka), Databricks (managed Apache Spark) and AWS S3 (storage service).
- This solution would also enable organizations to perform NRT Analysis on streaming data from multiple sources.
- This streaming data would be cleaned and stored on cloud enabled data warehouses to be used later.

Technologies used :

- [Confluent](#) (managed [Apache Kafka](#))
- [Databricks](#) (managed [Apache Spark](#))
- AWS S3 (for data storage)

Architecture :



Producers : As seen in the Architecture diagram above there are various devices such as smartphones, smartwatches and IoT devices as well as various applications that are generating data every second. These devices and applications are referred to as producers hereafter. This data would be written to Kafka topics based on the device that generates the data. A Kafka topic would logically group similar data.

Apache Kafka : Apache Kafka is a distributed event store and stream-processing platform. Kafka uses a distributed paradigm to store data which enables scalability and fault tolerance for the solution. We use managed kafka i.e Confluent Kafka.

Data written to Kafka can be consumed 'n' consumer applications. Kafka supports data retention and can be configured from the cluster management page. (from 1 hour to forever)

Core capabilities of Kafka :

- High throughput
- Scalable
- Permanent Storage
- High Availability

Apache Spark : Apache Spark is a multi-language engine for executing data engineering, data science, and machine learning on single-node machines or clusters. Apache Spark supports processing streaming data as well as it is 100 times faster in memory and 10 times faster on disk than Hadoop MapReduce as it uses RAM for data processing while MapReduce writes the intermediate output to disk which is the reason for its slowness.

Here the consuming applications are spark jobs that read data from Kafka, clean it, and

perform NRT analysis on streaming data. Spark jobs are also programmed to persist data to AWS S3. This data persisted on S3 can be later used by the application team to perform analysis on this historical data.

Spark also supports external tables, wherein external SQL tables can be created in which the data lies on external cloud based object stores such as AWS S3 buckets or GCP Buckets and the table metadata lies on Databricks hive metastore.

For implementation of this project we have used Managed Spark provided by Databricks.

Core Capabilities of Spark :

- Batch and streaming data support.
- SQL based analysis.
- Machine Learning library support.
- Data Science.

Spark also supports various languages such as python, SQL, scala, java, R

AWS S3 : Amazon Simple Storage Service (Amazon S3) is an object storage service offering industry-leading scalability, data availability, security, and performance.

Data produced by devices and applications is consumed by Spark jobs. The data is cleaned, processed and written to S3 in JSON and CSV formats. External tables are then created on these S3 paths for application teams to be used later.

Deployment & Setup steps :

Git Repo : [ADBMS git repository](#)

Setting Up Confluent :

- Sign in to Confluent Cloud at <https://confluent.cloud>.
- Click Add cluster.
- On the Create cluster page, for the Basic cluster, select Begin configuration. This creates a Basic cluster, which supports single zone availability
- On the Region/zones page, choose a cloud provider, region, and select a single availability zone.
- Select Continue and set up payment on the next page.
- Specify a cluster name, review the configuration and cost information, and then select Launch cluster.

The screenshot shows the 'Create cluster' wizard on the 'Configuration & cost' tab. The cluster name is set to 'quickstart_cluster'. Below it, a table provides cost details per hour:

	Base cost	\$0 /hr
Write	\$0.1265	/GB
Read	\$0.1265	/GB
Storage	\$0.00015972	/GB-hour
Partitions	\$0.0046	/Partition-hour

Below the table, there are tabs for 'Configuration & cost' (selected), 'Usage limits', and 'Uptime SLA'. A note states: 'Settings marked with an asterisk (*) cannot be changed once you launch your cluster'. The cluster configuration section shows 'Basic' as the provider and 'us-west4' as the region. At the bottom are 'Go back', 'Review payment method', and 'Launch cluster' buttons.

- Now you can get started configuring apps and data on your new cluster.
- Once cluster has been created Navigate to API keys for your cluster generate Key and Secret and save it
- Use the "config.ini" from git and replace the <Server details>, <API Key> and <API secret> this will be used while running the mocked producer code.

Setting Up Producer (Mocked producers) :

- An application that is the source of the data stream is called a producer. In order to generate tokens or messages and further publish it to one or more topics in the Kafka cluster, we use Apache Kafka Producer.
- Steps to create a Kafka topic:
 - Sign in to Confluent Cloud at <https://confluent.cloud>.
 - From the navigation menu, click Topics, and then click Create topic.
 - In the Topic name field, type "weather" and then select Create with defaults.
 - The 'weathers' topic is created on the Kafka cluster and is available for use by producers and consumers.
 - Similarly, a 'energy consumption' topic was also created where energy consumption mock data was pushed.
- The producer code is present in 'tiktokStreamingData_producer.py', 'energyData_producer.py' and weatherStreaming_producer.py files.
- In these files, connection to the producer is first established using the config file passed via the command line.
- A producer instance is then created and mock data is then pushed to the producer using the producer.produce() function.
- These producers are just mocked sources of data to simulate the behavior of actual devices and applications.
- To run these producers we need Python 3 and some basic packages that can be installed using pip.
- Packages such as confluent_kafka and a few other basic packages are required.

Setting Up Databricks : Databricks is a paid service but a 14 day free trial is provided to first time users. But the infrastructure used for cluster and networking between nodes will be charged by the cloud provider.

- Start your free Databricks trial
 - [Databricks landing page](#)
 - [Databricks free trial](#)
- Create a cluster from the landing page by selecting “Compute” from the left side menu.

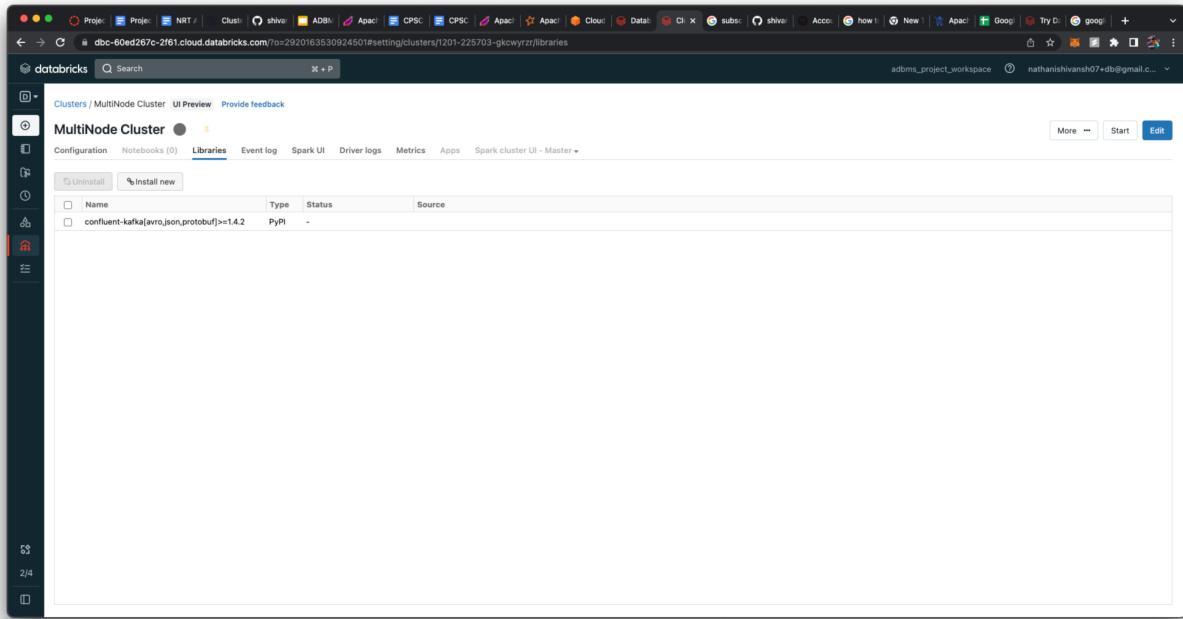
The screenshot shows the Databricks landing page. On the left, there is a sidebar with several options: Workspace, Repos, Recents, Data, Compute (which is highlighted with a red box), and Workflows. The main area displays various Databricks services: Notebook, Data import, Partner Connect, AutoML, and Transform data. Below these are sections for Last viewed, Release notes, and Blog posts. At the bottom, there are links for Partner Connect, Tasks Completed, and Menu options.

The screenshot shows the Databricks Compute page. It features a search bar at the top and a table below it. The table has columns for Name, Policy, Runtime, Active memory, Active cores, Active DBU / h, Source, and Creator. Two clusters are listed: '(default)basic-starter-cluster' and 'MultiNode Cluster'. At the top left of the table, there is a 'Create compute' button, which is highlighted with a red box. The table also includes filters for 'Created by me' and 'Accessible by me'.

Add the required cluster configurations and hit create cluster.

- Once the cluster is set up we need a notebook to write Spark Job/Code.
- From the left side menu click new -> click notebook and add the required details such as name, language that you wish to code in and cluster to which the job should be submitted.

- To connect to kafka we need the confluent kafka library.
- Open your cluster settings, navigate to the libraries tab, select PyPi and install the library "confluent-kafka[avro,json,protobuf]>=1.4.2"
- Once the library is installed this is how your cluster's library tab looks like.



- We use databricks secrets to store connection details for kafka from databricks, to set up kafka login details as secrets, use this [blog](#). Elseway you can directly use plain text for connecting to kafka. I.e. not from secrets but directly hard coded in the notebook. The connection details are once obtained while we set up a confluent cluster earlier.

Running the Application :

- Clone the github repo: [ADBMS Project](#)
- We have multiple examples in this repo.
- Navigate to the appropriate one that you wish to Run.
- Each example has a producer and a consumer file.
- Producer file is a plain python code to simulate streaming source.
- To run the Producer file the code expects a config file "config.ini" that has the kafka connection details.
- Run the producer code as shown below.

The screenshot shows the Confluent Cloud dashboard with the 'tiktok_category' topic selected. The 'Messages' tab is active, showing a list of produced messages. The terminal window on the right shows the command to run the producer and the resulting output of data being pushed to the Kafka topic.

```
kafka_managed_confluent --bash -T16x56
Shivansh-MBP:CPSC-531-ADBMS-Project shivansh$ cd tiktok_streaming_example/
Shivansh-MBP:tiktok_streaming_example shivansh$ cd ..
Shivansh-MBP:CPSC-531-ADBMS-Project shivansh$ cp ..\producer.py .
Shivansh-MBP:CPSC-531 shivansh$ ls
producer.py test1.py weatherStreaming.py
Shivansh-MBP:kafka_managed_confluent shivansh$ ls
test1.py
Shivansh-MBP:CPSC-531 shivansh$ ./producer.py test1.py weatherStreaming.py config.ini
1670722392.93971,Entertainment,3,0
Produced event to topic tiktok_category: key = 50000009 value = 1670722392.93971,Entertainment,3,0
1670722395.942857,Fitness/sports,21,0
Produced event to topic tiktok_category: key = 50000009 value = 1670722395.942857,Fitness/sports,21,0
5000001
1670722398.944313,Beauty/skin care,12,0
Produced event to topic tiktok_category: key = 50000001 value = 1670722398.944313,Beauty/skin care,12,0
50000001
1670722401.949785,Dance,11,0
Produced event to topic tiktok_category: key = 50000009 value = 1670722401.949785,Dance,11,0
1670722392.93971,Entertainment,3,0
1670722395.942857,Fitness/sports,21,0
1670722398.944313,Beauty/skin care,12,0
1670722401.949785,Dance,11,0
[1]: Stopped python3 tiktokStreamingData.py config.ini
Shivansh-MBP:kafka_managed_confluent shivansh$
```

- As soon as you run the file the code starts pushing data to kafka. You can check this by logging into the confluent dashboard and navigating to the topic that you are pushing data to and the messages produced start showing up here.
- Next up is running the Spark consuming application/job.
- Login to Databricks and open the notebook that we created during the setup phase.
- Paste the code from the "____consumer.py" file in each notebook cell.
- Run each cell by hitting "shift" and "enter".
- This is how the output looks like.

```

1 # writing the batched analysis to console
2 metrics.writeStream.format("console").outputMode("complete").start()

# 1 Spark Jobs
@15a8e107-ea35-4159-acef-a19481b7922 Last updated: 8 days ago

Dashboard Raw Data

Input vs. Processing Rate
records per second
0.3 rec/s
0.4 rec/s
Input rate: Processing rate
Batch Duration
Batch duration in seconds
16.7 s 16.3 s
16.7 s Average 16.3 s Latest
Aggregation State
102
Distinct keys
102

```

Out[4]: <pyspark.sql.streaming.query.StreamingQuery at 0x7f78b16aebe0>

Command completed

```

1 # writing streaming dataframe to s3 using writestream function which also ensures that new data is written using checkpoint concept.
2 tiktokDF.writeStream
3 .format("json")
4 .option("checkpointLocation", "s3://adlms.project.data/tiktok_streaming_data_checkpoint_1")
5 .option("path", "s3://adlms.project.data/tiktok_streaming_data_1")
6 .start()
7

```

- For the console output navigate to the cluster details page, navigate to Driver logs tab. Under it look for the Standard output.
- This is how the console log looks like with the expected output.

Log Type	Description	Size	Last Log
stdout	Standard output	86.53 KB	Latest log
stdout	Standard output	87.46 KB	Latest log
stdout	Standard output	112.26 KB	Latest log
stderr	Standard error	164 bytes	Latest log
stderr	Standard error	164 bytes	Latest log

Standard output

```

[5000004] Fitness/sports| 27.0| 1| 0|
[5000004] Home reno/DIY| 16.5| 2| 1|
only showing top 20 rows
-----
Batch: 15
-----
| key | category | avgWatchTime | totalViewed | likes |
+---+---+---+---+
| 5000000 | Dance | 12.0 | 1 | 0 |
| 5000000 | Fashion | 18.0 | 1 | 0 |
| 5000001 | Fitness/exc | 12.0 | 2 | 1 |
| 5000001 | Cooking/recipes | 6.5 | 2 | 0 |
| 5000001 | Dance | 14.0 | 1 | 0 |
| 5000001 | Fashion | 9.0 | 1 | 1 |

```

Standard error

- For the streaming example for energy and weather data we can run the producer and consumer codes similar to that of the TikTok example.
- Spark jobs also write data to S3. Below is a snapshot of S3

We're continuing to improve the S3 console to make it faster and easier to use. If you have feedback on the updated experience, choose [Provide feedback](#).

Objects (8)

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	checkpoint/	Folder	-	-	-
<input type="checkbox"/>	energy_consumption_checkpoint_mock/	Folder	-	-	-
<input type="checkbox"/>	energy_consumption_checkpoint/	Folder	-	-	-
<input type="checkbox"/>	energy_consumption_mock/	Folder	-	-	-
<input type="checkbox"/>	energy_consumption/	Folder	-	-	-
<input type="checkbox"/>	iot_data/	Folder	-	-	-
<input type="checkbox"/>	weather_data_checkpoint_mock/	Folder	-	-	-
<input type="checkbox"/>	weather_data_mock/	Folder	-	-	-

Amazon S3 > Buckets > adbms.project.data > energy_consumption/

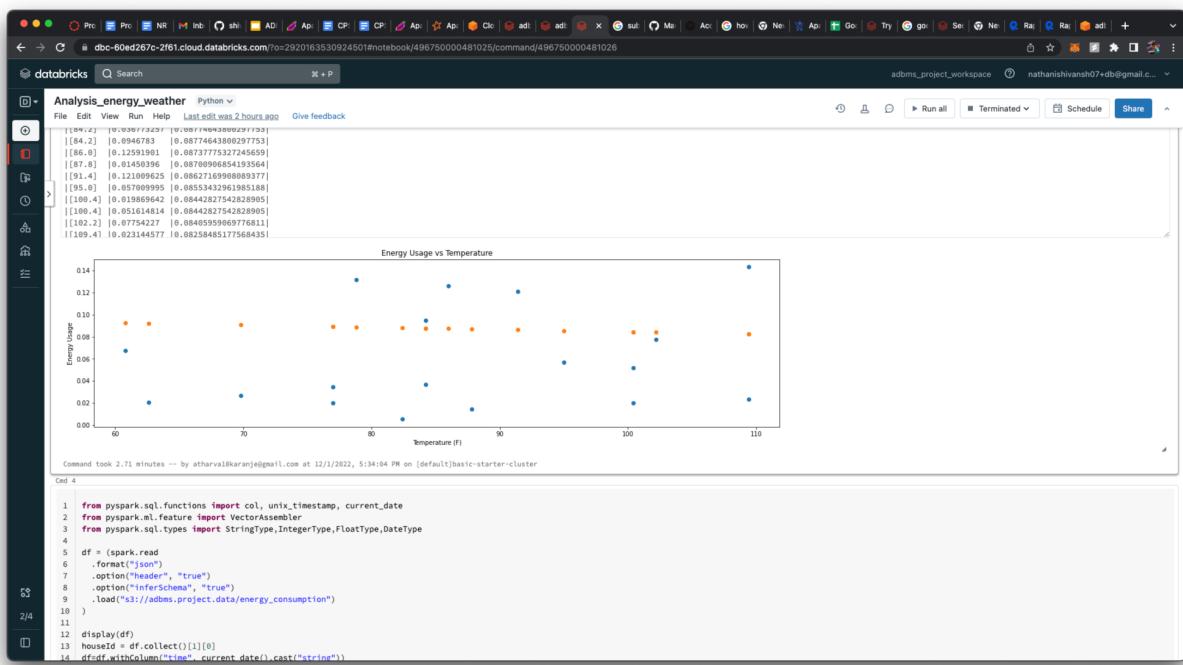
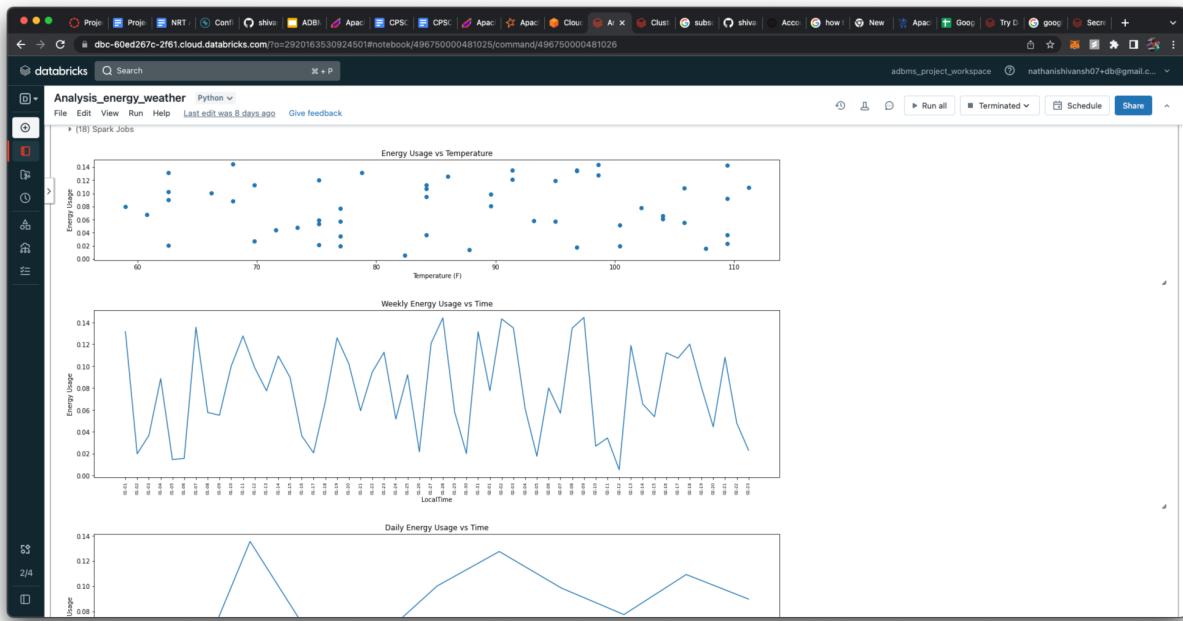
[Copy S3 URI](#)

Objects (24)

Objects are the fundamental entities stored in Amazon S3. You can use [Amazon S3 inventory](#) to get a list of all objects in your bucket. For others to access your objects, you'll need to explicitly grant them permissions. [Learn more](#)

<input type="checkbox"/>	Name	Type	Last modified	Size	Storage class
<input type="checkbox"/>	_spark_metadata/	Folder	-	-	-
<input type="checkbox"/>	part-00000-05f90ba1-5244-493a-b4d5-615cf2388569-c000.json	json	December 2, 2022, 14:24:13 (UTC-08:00)	2.3 KB	Standard
<input type="checkbox"/>	part-00000-11b87ff8-c962-4c18-a24a-6c02a2c040ef-c000.json	json	December 2, 2022, 14:24:31 (UTC-08:00)	648.0 B	Standard
<input type="checkbox"/>	part-00000-1b58e085-5069-4247-be5f-08af86f6962a-c000.json	json	December 2, 2022, 14:23:01 (UTC-08:00)	0 B	Standard

- Moving to the Visualization part of data we have a file by the name "Analysis_energy_weather.html" you can paste the contents to a new notebook and run the commands from the notebook. This is how the visualizations would look like.



The analysis code is present in the 'analysis_energy_weather' file and the code explanation is as follows:

- Analysis was performed on the weather and energy mock data that is generated.
- The data is initially stored in the Pyspark dataframes for analysis.
- The data is then processed and converted to a desirable format. (Ex. Data types are changes, dates are formatted, filtered based on a key column value, etc)
- Inner join is then performed on the dataframes based on the timestamp attribute to get the temperature and energy usage for a record at a particular timestamp.
- Scatter plots and line graphs are then generated to analyze and compare different attributes. The following graphs were created for analysis purpose:
 - Scatter plot of Energy usage vs temperature - This will help understand the energy consumption trend against the temperature.
 - Line graph for weekly energy usage vs time - This will help understand the energy usage trend during different days of a week.
 - Line graph for daily energy usage vs time - This will help understand the energy usage trend during different times of a day.
- Furthermore, prediction was also performed for the energy usage attribute using the linear regression model.
- For this, the data was initially split into 70% train and 30% test data.
- The linear regression model was trained based on the 'train data' and was tested against the 'test data' for evaluating the accuracy.
- The linear regression model will help predict energy usage for a household based on the temperature.

Conclusion: We have developed a robust, scalable and fault tolerant solution to harness streaming data by processing, analyzing it in NRT and storing it on cloud. This solution can be extended for any use case or organization where there is a need to collect, process and store streaming data.