# Forest Fire Predictive Analytics with AWS Cloud

Team 9: Anna Chow, Hemang Behl, Jason Gonsalves, Kevin Chuang, Richita Das

## 1. What is the project idea /What the application does?

Wildfires can cause devastating destruction and cost millions of dollars in damage, especially with increasing human development near wilderness or rural areas. The ability to predict wildfires would be a preemptive measure to prevent and/or manage wildfires. Factors affecting wildfires would be explored and used to gauge the probability of a wildfire occurring. Based on the location and the corresponding land type, a machine learning model will be trained to identify the major contributing factors of the fire. The whole process will be carried out in cloud where the model can directly consume the source file.

## 2. What are the technologies used?

- Python + libraries (numpy, sklearn, scipy, matplotlib, Flask, boto3, sagemaker) + IDE (Jupyter Notebook/Spyder)
- SQL
- AWS S3
- AWS Glacier
- AWS Glue
- AWS Athena
- AWS Redshift
- AWS QuickSight
- AWS SageMaker
- AWS CloudWatch Logs

## 3. Feature list of the dataset:

- FOD_ID = Global unique identifier.
- FPA_ID = Unique identifier that contains information necessary to track back to the original record in the source dataset.
- SOURCE_SYSTEM_TYPE = Type of source database or system that the record was drawn from (federal, nonfederal, or interagency).
- SOURCE_SYSTEM = Name of or other identifier for source database or system that the record was drawn from.
- NWCG_REPORTING_AGENCY = Active National Wildlife Coordinating Group (NWCG) Unit Identifier for the agency preparing the fire report
- NWCG_REPORTING_UNIT_ID = Active NWCG Unit Identifier for the unit preparing the fire report.

- NWCG_REPORTING_UNIT_NAME = Active NWCG Unit Name for the unit preparing the fire report.
- SOURCE_REPORTING_UNIT = Code for the agency unit preparing the fire report, based on code/name in the source dataset.
- SOURCE_REPORTING_UNIT_NAME = Name of reporting agency unit preparing the fire report, based on code/name in the source dataset.
- LOCAL_FIRE_REPORT_ID = Number or code that uniquely identifies an incident report for a particular reporting unit and a particular calendar year.
- LOCAL_INCIDENT_ID = Number or code that uniquely identifies an incident for a particular local fire management organization within a particular calendar year.
- FIRE_CODE = Code used within the interagency wildland fire community to track and compile cost information for emergency fire suppression
- FIRE_NAME = Name of the incident, from the fire report (primary) or ICS-209 report (secondary).
- ICS_209_INCIDENT_NUMBER = Incident (event) identifier, from the ICS-209 report.
- ICS_209_NAME = Name of the incident, from the ICS-209 report.
- MTBS_ID = Incident identifier, from the MTBS perimeter dataset.
- MTBS_FIRE_NAME = Name of the incident, from the MTBS perimeter dataset.
- COMPLEX_NAME = Name of the complex under which the fire was ultimately managed, when discernible.
- FIRE_YEAR = Calendar year in which the fire was discovered or confirmed to exist.
- DISCOVERY_DATE = Date on which the fire was discovered or confirmed to exist.
- DISCOVERY_DOY = Day of year on which the fire was discovered or confirmed to exist.
- DISCOVERY_TIME = Time of day that the fire was discovered or confirmed to exist.
- STAT_CAUSE_CODE = Code for the (statistical) cause of the fire.
- STAT_CAUSE_DESCR = Description of the (statistical) cause of the fire.
- CONT_DATE = Date on which the fire was declared contained or otherwise controlled (mm/dd/yyyy)
- CONT_DOY = Day of year on which the fire was declared contained or otherwise controlled.
- CONT_TIME = Time of day that the fire was declared contained or otherwise controlled (hhmm where hh=hour, mm=minutes).
- FIRE_SIZE = Estimate of acres within the final perimeter of the fire.
- FIRE_SIZE_CLASS = Code for fire size based on the number of acres within the final fire perimeter expenditures (A=greater than 0 but less than or equal to 0.25 acres, B=0.26-9.9 acres, C=10.0-99.9 acres, D=100-299 acres, E=300 to 999 acres, F=1000 to 4999 acres, and G=5000+ acres).
- LATITUDE = Latitude (NAD83) for point location of the fire (decimal degrees).
- LONGITUDE = Longitude (NAD83) for point location of the fire (decimal degrees).

- ■ OWNER_CODE = Code for primary owner or entity responsible for managing the land at the point of origin of the fire at the time of the incident.
- ■ OWNER_DESCR = Name of primary owner or entity responsible for managing the land at the point of origin of the fire at the time of the incident.
- ■ STATE = Two-letter alphabetic code for the state in which the fire burned (or originated), based on the nominal designation in the fire report.
- ■ COUNTY = County, or equivalent, in which the fire burned (or originated), based on nominal designation in the fire report.
- ■ FIPS_CODE = Three-digit code from the Federal Information Process Standards (FIPS) publication 6-4 for representation of counties and equivalent entities.
- ■ FIPS_NAME = County name from the FIPS publication 6-4 for representation of counties and equivalent entities.

**Features post feature engineering:**
- ● FIRE_YEAR  (long)
- ● STAT_CAUSE_DESCR (string)
- ● FIRE_SIZE (double)
- ● FIRE_SIZE_CLASS (string)
- ● LATITUDE  (double)
- ● LONGITUDE (double)
- ● STATE (string)
- ● COUNTY (long)
- ● DISCOVERY_DATE (double)
- ● CNT_DATE (double)

# 4. Sample Demo Screenshots

## AWS Athena

SQL Queries

SELECT count(fod_id) as count, state  FROM sampledb.test_266 GROUP BY state;

Results

| | count | state |
|---|---|---|
| 1 | 79 | MT |
| 2 | 8 | WA |
| 3 | 41 | ID |
| 4 | 28 | CO |
| 5 | 34 | MN |
| 6 | 16 | TX |
| 7 | 196 | AZ |
| 8 | 18 | MO |
| 9 | 3 | SC |
| 10 | 25 | NV |
| 11 | 5 | AR |
| 12 | 8 | FL |
| 13 | 0 | STATE |
| 14 | 23 | UT |
| 15 | 121 | CA |
| 16 | 121 | OR |
| 17 | 10 | WY |
| 18 | 12 | LA |
| 19 | 2 | OK |
| 20 | 111 | NM |
| 21 | 8 | NC |
| 22 | 1 | NE |
| 23 | 23 | SD |

SELECT count(fod_id) as count, county  FROM sampledb.test_266 GROUP BY county;

Results

| | count | county |
|---|---|---|
| 1 | 33 | 37 |
| 2 | 17 | 33 |
| 3 | 3 | 87 |
| 4 | 19 | 21 |
| 5 | 10 | 29 |
| 6 | 22 | 25 |
| 7 | 8 | 75 |
| 8 | 33 | 3 |
| 9 | 5 | 63 |
| 10 | 3 | 59 |
| 11 | 2 | 67 |
| 12 | 69 | 17 |
| 13 | 22 | 47 |
| 14 | 3 | 55 |
| 15 | 14 | 51 |
| 16 | 2 | 91 |
| 17 | 1 | 95 |
| 18 | 1 | 213 |
| 19 | 1 | 165 |
| 20 | 3 | 510 |
| 21 | 1 | 81 |
| 22 | 2 | 93 |
| 23 | 1 | 89 |

select max(count) as c, state from
(SELECT state,count(fod_id) as count FROM sampledb.test_266 GROUP  BY state order by count desc) group by state order by c desc limit 1  ;

Results

| | c | state |
|---|---|---|
| 1 | 196 | AZ |

select * from
(SELECT state,count(fod_id) as count FROM sampledb.test_266 GROUP  BY state order by count
limit 2) where state <> 'STATE';

**Results**

| | state | count |
|---|---|---|
| 1 | NE | 1 |

select max(stat_cause_descr) as cause,count(stat_cause_descr) as number from
sampledb.test_266 where state='CO' ;

**Results**

| | cause | number |
|---|---|---|
| 1 | Miscellaneous | 28 |

select max(stat_cause_descr) as cause,count(stat_cause_descr) as number from
sampledb.test_266 where state='AZ' ;

**Results**

| | cause | number |
|---|---|---|
| 1 | Smoking | 196 |

select max(stat_cause_descr) as cause,count(stat_cause_descr) as number from
sampledb.test_266 where state='MT' ;

**Results**
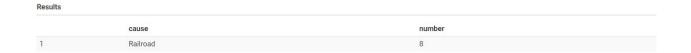
| | cause | number |
|---|---|---|
| 1 | Railroad | 79 |

select max(stat_cause_descr) as cause,count(stat_cause_descr) as number from
sampledb.test_266 where state='UT' ;

**Results**
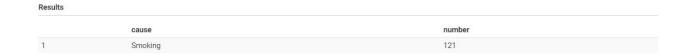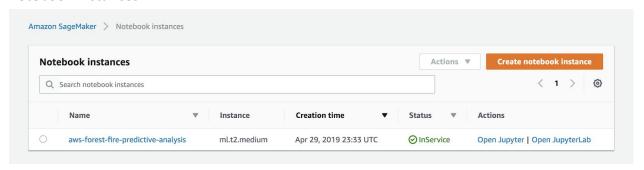
| | cause | number |
|---|---|---|
| 1 | Smoking | 23 |

select max(stat_cause_descr) as cause,count(stat_cause_descr) as number from
sampledb.test_266 where state='WA' ;

**Results**

| | cause | number |
|---|---|---|
| 1 | Railroad | 8 |

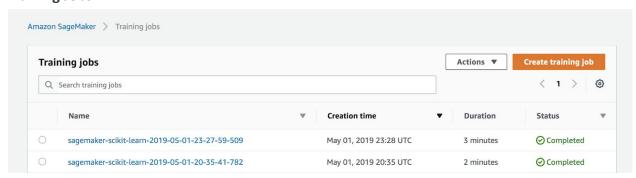select max(stat_cause_descr) as cause,count(stat_cause_descr) as number from sampledb.test_266 where state='OR' ;
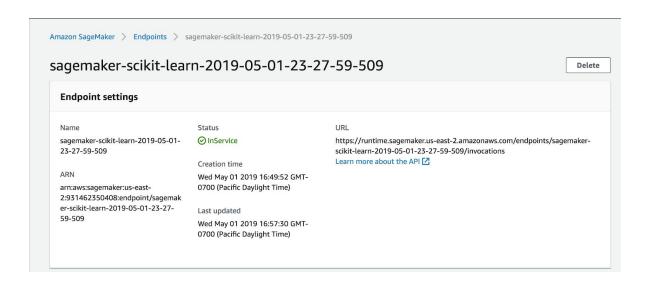
**Results**

| | cause | number |
|---|---|---|
| 1 | Smoking | 121 |

## AWS SageMaker

### Notebook Instances

Amazon SageMaker > Notebook instances

**Notebook instances**          Actions ▼   **Create notebook instance**

🔍 Search notebook instances                        ‹ 1 › ⚙

| | Name | Instance | Creation time | Status | Actions |
|---|---|---|---|---|---|
| ○ | aws-forest-fire-predictive-analysis | ml.t2.medium | Apr 29, 2019 23:33 UTC | ⊘ InService | Open Jupyter \| Open JupyterLab |

### Training Jobs

Amazon SageMaker > Training jobs

**Training jobs**          Actions ▼   **Create training job**

🔍 Search training jobs                        ‹ 1 › ⚙

| | Name | Creation time | Duration | Status |
|---|---|---|---|---|
| ○ | sagemaker-scikit-learn-2019-05-01-23-27-59-509 | May 01, 2019 23:28 UTC | 3 minutes | ⊘ Completed |
| ○ | sagemaker-scikit-learn-2019-05-01-20-35-41-782 | May 01, 2019 20:35 UTC | 2 minutes | ⊘ Completed |

### SageMaker Endpoint (Deployed Model)

Amazon SageMaker > Endpoints > sagemaker-scikit-learn-2019-05-01-23-27-59-509

# sagemaker-scikit-learn-2019-05-01-23-27-59-509

Delete

## Endpoint settings

**Name**
sagemaker-scikit-learn-2019-05-01-23-27-59-509

**ARN**
arn:aws:sagemaker:us-east-2:931462350408:endpoint/sagemaker-scikit-learn-2019-05-01-23-27-59-509

**Status**
⊘ InService

**Creation time**
Wed May 01 2019 16:49:52 GMT-0700 (Pacific Daylight Time)

**Last updated**
Wed May 01 2019 16:57:30 GMT-0700 (Pacific Daylight Time)

**URL**
https://runtime.sagemaker.us-east-2.amazonaws.com/endpoints/sagemaker-scikit-learn-2019-05-01-23-27-59-509/invocations
Learn more about the API

# CloudWatch Logs

CloudWatch > Log Groups > /aws/sagemaker/Endpoints/sagemaker-scikit-learn-2019-05-01-23-27-59-509 > All streams

✖
ℹ **Try CloudWatch Logs Insights**
CloudWatch Logs Insights allows you to search and analyze your logs using a new, purpose-built query language. Click here to experience it. If you want to learn more, read the AWS blog or visit our documentation.

Expand all ● Row ○ Text

Filter events        all  30s **5m** 1h  6h  1d  1w  custom ▾

| | Time (UTC +00:00) | Message | Show in stream |
|---|---|---|---|
| | 2019-05-02 | | |
| | | *No older events found for the selected date range. Adjust the date range.* | |
| ▸ | 01:18:49 | 10.32.0.2 - - [02/May/2019:01:18:48 +0000] "GET /ping HTTP/1.1" 200 0 "-" "AHC/2.0" | AllTraffic/i-073d4778be439fa6b |
| ▸ | 01:18:54 | 10.32.0.2 - - [02/May/2019:01:18:53 +0000] "GET /ping HTTP/1.1" 200 0 "-" "AHC/2.0" | AllTraffic/i-073d4778be439fa6b |
| ▸ | 01:18:59 | 10.32.0.2 - - [02/May/2019:01:18:58 +0000] "GET /ping HTTP/1.1" 200 0 "-" "AHC/2.0" | AllTraffic/i-073d4778be439fa6b |
| ▸ | 01:19:04 | 10.32.0.2 - - [02/May/2019:01:19:03 +0000] "GET /ping HTTP/1.1" 200 0 "-" "AHC/2.0" | AllTraffic/i-073d4778be439fa6b |
| ▸ | 01:19:09 | 10.32.0.2 - - [02/May/2019:01:19:08 +0000] "GET /ping HTTP/1.1" 200 0 "-" "AHC/2.0" | AllTraffic/i-073d4778be439fa6b |
| ▸ | 01:19:14 | 10.32.0.2 - - [02/May/2019:01:19:13 +0000] "GET /ping HTTP/1.1" 200 0 "-" "AHC/2.0" | AllTraffic/i-073d4778be439fa6b |

# QuickSight Analysis

## Dashboard Setup

## QuickSight Dashboard Overview

CMPE 266 - Big Data Engineering & Analytics
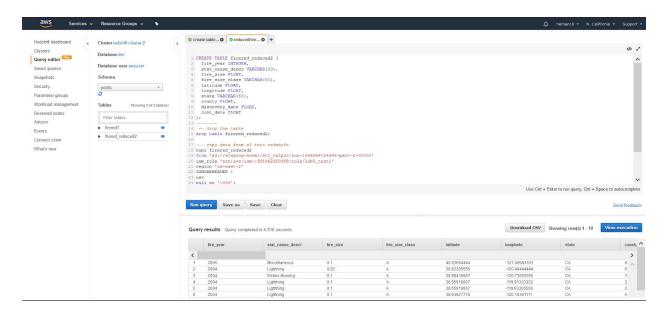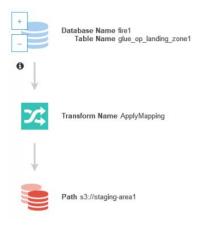
**RedShift Demo**

Querying data loaded from S3 (raw data)
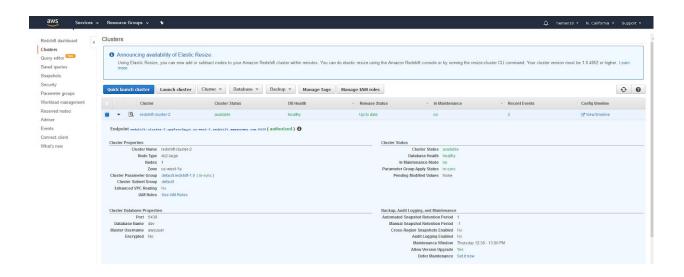


Querying data loaded from S3 (reduced data set)

**AWS Glue**
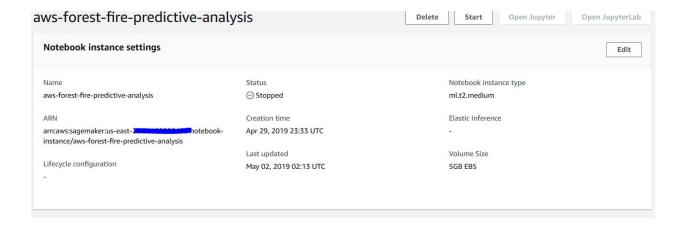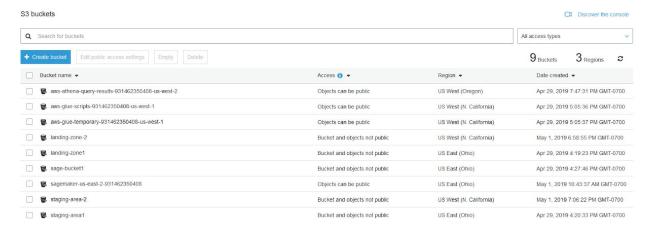


# 5. AWS Configurations Screenshots

## 5.1 Redshift Cluster

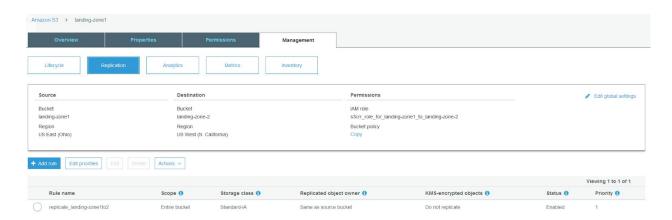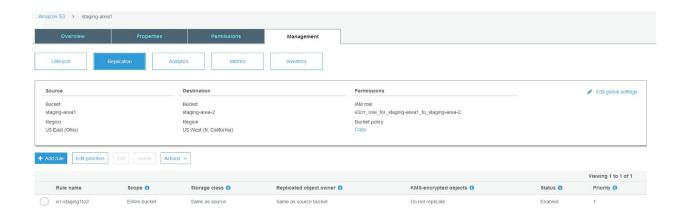## 5.2 Sagemaker



## 5.3 S3



## 5.3.1 Cross Region Replication

## 5.3.2 Transition to Glacier

## 5.4 Glue

## 5.4.1 Crawler

Crawlers > fire_crawler1

Run crawler    Edit

| | |
|---|---|
| Name | fire_crawler1 |
| Description | |
| Create a single schema for each S3 path | false |
| Security configuration | |
| Tags | - |
| State | Ready |
| Schedule | |
| Last updated | Mon Apr 29 16:50:38 GMT-700 2019 |
| Date created | Mon Apr 29 16:50:38 GMT-700 2019 |
| Database | fire1 |
| Table prefix | glue_op_ |
| Service role | service-role/AWSGlueServiceRole-glues3access |
| Selected classifiers | |
| Data store | S3 |
| Include path | s3://landing-zone1 |
| Exclude patterns | |
| Data store | S3 |
| Include path | s3://staging-area1 |
| Exclude patterns | |
| Data store | S3 |
| Include path | s3://sage-bucket1 |
| Exclude patterns | |

### Configuration options

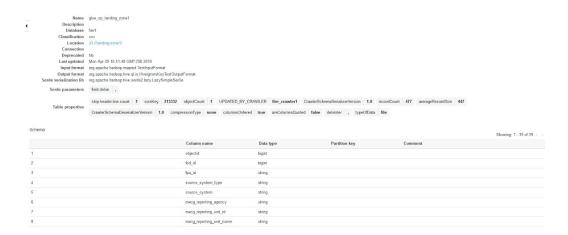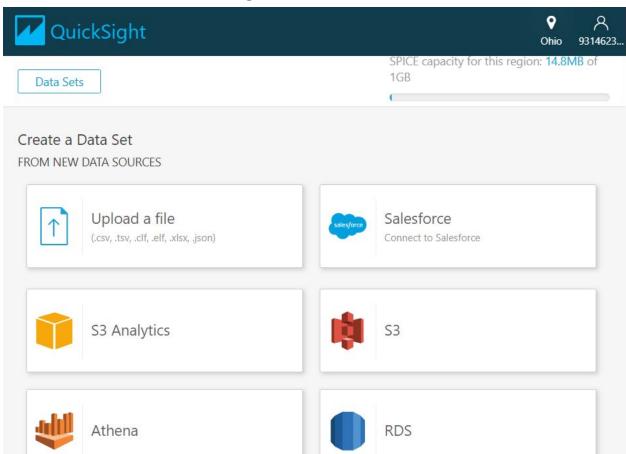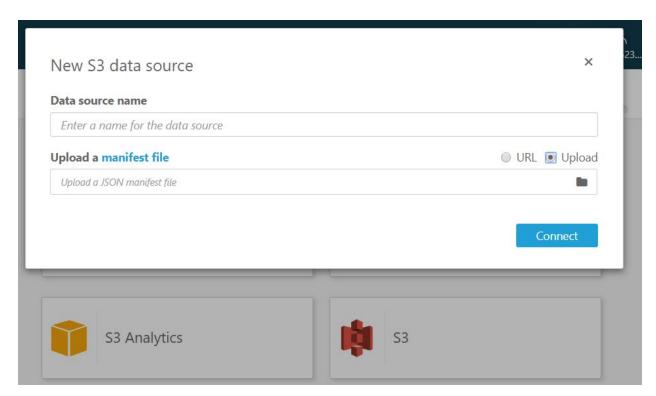| | |
|---|---|
| Schema updates in the data store | Update the table definition in the data catalog. |
| Object deletion in the data store | Mark the table as deprecated in the data catalog. |

## 5.4.2 ETL Mapping

## 5.4.3 Glue Table
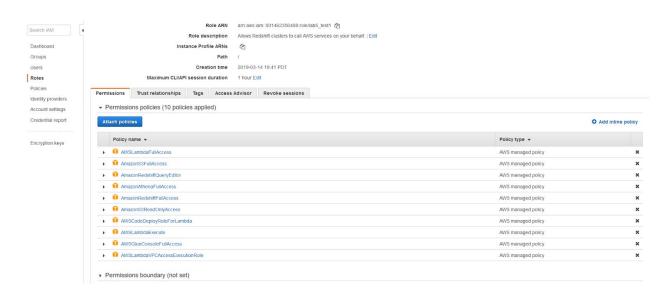
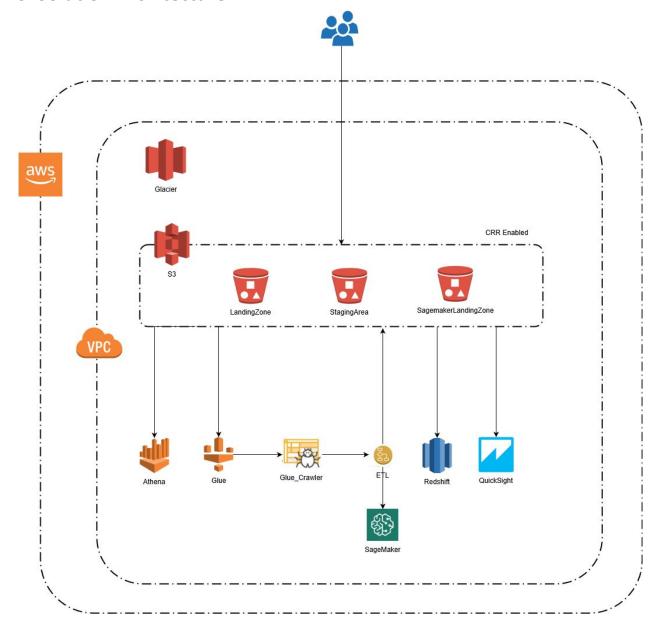

## 5.5 QuickSight

Connect an S3 bucket to QuickSight

## Manifest JSON file template and S3 file specifications

## 5.6 IAM

Policies:

# 6. Solution Architecture



- The main aim of the project is to develop a model which can identify and categorize the various fire incidents based on their cause.
- The data file is first uploaded onto the S3 bucket.
- We use Athena and Glue (using Glue crawlers) to read the data from the S3 buckets.
- Since the initial data has multiple columns, we perform feature engineering to extract only the required columns
- The data is cleaned using ETL functions and then the data is sent to the SageMaker where it is used to train the model

- A model has been developed using Python libraries to categorize the various fire causes.
- We are using QuickSight to analyze the data from the S3 buckets and create dashboards presenting the data in charts showing the various areas affected by fire. We created histograms depicting various information.
- We are using three S3 buckets, each in different regions to increase availability and reliability.

## 7. URL to GitHub.

https://github.com/k-chuang/aws-forest-fire-predictive-analytics

## 8. Future Scope

- Using Lambda to automate the flow so that as soon as the document lands in the S3 landing area, the other functions such as ETL jobs and training models are triggered automatically.
- Online Training (i.e. training on the fly for data as it comes in) for SageMaker
- Create an interactive frontend for the user to upload the document for which analysis is wanted.
- The user can also query the system via the frontend.