

**CS506-02**

**Distributed Cloud Computing**



**Final Project**

Name	Student Id
Naveen Sita Reddy Kovvuri	S02059263
Nithish Kumar Medabalimi	S02059513
Tanay Tammineni	S02065959
Sainath Reddy Mosali	S02061036
Lakshmi Kethineni	S02066249
Vidya Rani Vallala	S02066899

**Spring 2024 CS506-01 22619**

**PROFESSOR: DR. ZIPPING LIU**

**Southeast Missouri State University**

## **Project Name: Predicting the Effects of Earthquakes Using Cloud-Based Machine Learning**

### **Project Description:**

The goal of this project, "Predicting the Effects of Earthquakes Using Cloud-Based Machine Learning," is to create prediction models for earthquake features by utilizing PySpark and AWS SageMaker. Improving seismic event predicting and showcasing the effectiveness of cloud technology in handling and evaluating massive datasets are the main goals. The desired results include enhanced readiness for earthquakes and the capacity to react quickly to any seismic hazards by using reliable and accurate machine learning forecasts.

### **Scope:**

The scope of this project includes building and implementing machine learning models with a focus on seismic event prediction using PySpark and Amazon SageMaker. Large seismic data sets are handled efficiently by the project by utilizing AWS cloud services, such as S3 for data storage and EC2 for processing power. Machine learning processes are made scalable and reliable by using SageMaker for model deployment and use of PySpark for data preparation and model training. Additionally, this project examines how to integrate AWS Lambda for possible real-time data processing, which might improve the system's ability to react to seismic activity in real time.

### **Requirements:**

For this project, AWS's cloud computing services like EC2, S3, and Lambda, as well as PySpark for data processing and Amazon SageMaker for model deployment are required. For development and testing, Jupyter Lab is a crucial

tool, as are other Python libraries for data processing and show, such NumPy, pandas, and Matplotlib. High-performance computer resources are also required for the project to manage the massive amounts of data processing and analytics needed for earthquake prediction.

## **Architecture:**

**Data Management and Storage:** Raw and processed seismic data are stored on Amazon S3, which serves as the main data repository.

## **Data Pre-Processing and Model Development:**

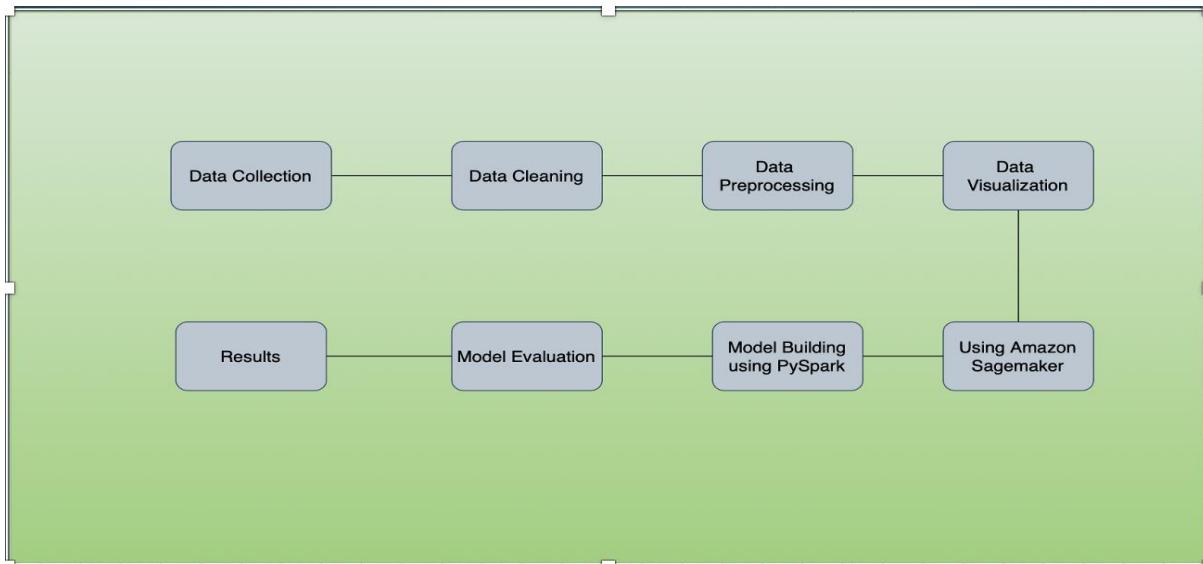
PySpark on AWS EC2: PySpark is run on EC2 instances to take use of its distributed computing resources for effective data processing and model training.

**Model Management and Deployment:** The machine learning models that have been trained go into production using AWS SageMaker.

## **Introduction:**

The goal of this project is to accurately predict earthquake characteristics by utilizing PySpark and Amazon SageMaker. The project makes optimal use of cloud computing to collect and analyze massive amounts of seismic data by deploying these technologies on AWS. Enhancing earthquake prediction, improving disaster response systems, and showcasing the effectiveness and scalability of cloud-based data processing are the main goals. This project is significant because it shows how modern methods of machine learning may be used in real-world situations to reduce the risks connected with natural catastrophes, which has major implications for preparedness and safety.

## Workflow Diagram:



## Dataset Description:

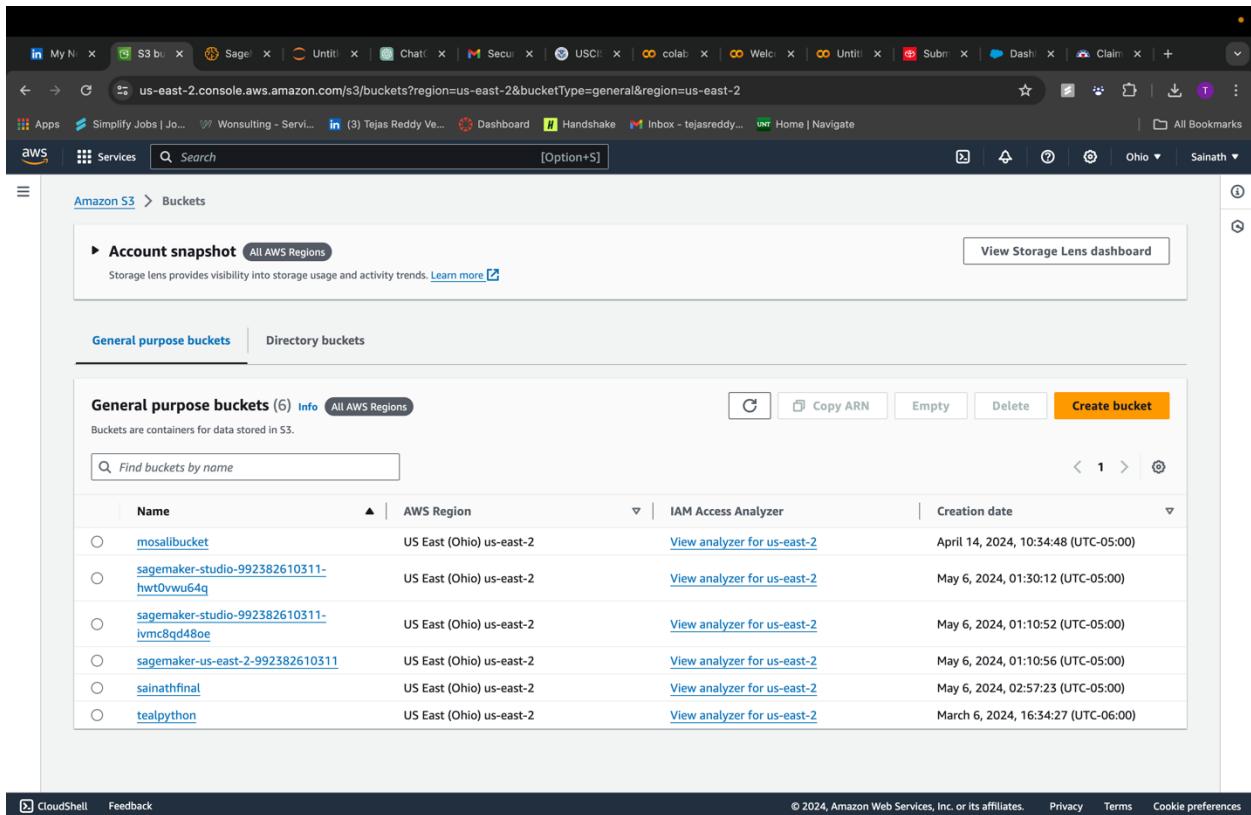
Variable Name	Description
<u>magNst</u>	Number of stations reporting the magnitude.
updated	last update of the record.
Depth	Depth of the earthquake in kilometers.
status	Event has been automatically reported or manually reviewed.
<u>magType</u>	The method or algorithm used to calculate the magnitude.
<u>horizontalError</u>	Errors in the measurement of the earthquake's horizontal position.
<u>dmin</u>	Horizontal distance from the epicenter to the nearest station (in degrees).
rms	Root Mean Square of the seismic wave, which is a measure of the accuracy of the network's earthquake location.
place	Descriptive location of the earthquake.
type	Type of seismic event.
gap	Gap between azimuthally adjacent stations.
status	Indicates if the event has been automatically reported or manually reviewed.
<u>locationSource</u>	The source that provided the location.
mag	Overall size or strength of the earthquake.
Time	Timestamp of the earthquake
Latitude, longitude	Coordinates of the earthquake epicenter.
<u>nst</u>	Number of seismic stations

## Cloud Technologies and Services Used:

This project provides a full featured cloud computing environment by using a variety of Amazon Web Services (AWS). **PySpark** is used to run complex algorithms for machine learning in an effective manner because to the scalable computing facilities provided by Amazon EC2. Amazon S3, which provides reliable and safe data storage solutions, handles data management and storage. With its capabilities for simple scaling and real-time analytics, **Amazon SageMaker** is a vital resource for installing, tracking, and controlling machine learning models. The project's ability to handle seismic data as it comes in is further improved by the optional use of AWS Lambda, which enables the execution of code in response to real-time data triggers.

## Project Implementation:

### Creating an S3 bucket



The screenshot shows the AWS S3 console interface. At the top, there is a navigation bar with various tabs and links. Below the navigation bar, the main content area has a header titled "Account snapshot" with a link to "View Storage Lens dashboard". There are two tabs at the top of the main content area: "General purpose buckets" (which is selected) and "Directory buckets". Below these tabs, there is a search bar labeled "Find buckets by name". The main content area displays a table of "General purpose buckets". The table has columns for "Name", "AWS Region", "IAM Access Analyzer", and "Creation date". There are 6 buckets listed:

Name	AWS Region	IAM Access Analyzer	Creation date
<a href="#">mosalibucket</a>	US East (Ohio) us-east-2	<a href="#">View analyzer for us-east-2</a>	April 14, 2024, 10:34:48 (UTC-05:00)
<a href="#">sagemaker-studio-992382610311-hwtovwu64q</a>	US East (Ohio) us-east-2	<a href="#">View analyzer for us-east-2</a>	May 6, 2024, 01:30:12 (UTC-05:00)
<a href="#">sagemaker-studio-992382610311-ivmc8qd48oe</a>	US East (Ohio) us-east-2	<a href="#">View analyzer for us-east-2</a>	May 6, 2024, 01:10:52 (UTC-05:00)
<a href="#">sagemaker-us-east-2-992382610311</a>	US East (Ohio) us-east-2	<a href="#">View analyzer for us-east-2</a>	May 6, 2024, 01:10:56 (UTC-05:00)
<a href="#">sainathfinal</a>	US East (Ohio) us-east-2	<a href="#">View analyzer for us-east-2</a>	May 6, 2024, 02:57:23 (UTC-05:00)
<a href="#">tealpython</a>	US East (Ohio) us-east-2	<a href="#">View analyzer for us-east-2</a>	March 6, 2024, 16:34:27 (UTC-06:00)

At the bottom of the page, there are links for "CloudShell", "Feedback", "© 2024, Amazon Web Services, Inc. or its affiliates.", "Privacy", "Terms", and "Cookie preferences".

using amazon sagemaker to create an instance to deploy ML models

setting up the domain to work in jupyter hub

	Name	Id	Status	Created on	Modified on
<input type="radio"/>	QuickSetupDomain-20240506T011069	d-cncvy4yqd4cy	<input checked="" type="checkbox"/> InService	May 06, 2024 06:10 UTC	May 06, 2024 06:16 UTC
<input type="radio"/>	QuickSetupDomain-20240506T013031	d-2z2pyqbp1a9y	<input checked="" type="checkbox"/> InService	May 06, 2024 06:30 UTC	May 06, 2024 06:33 UTC

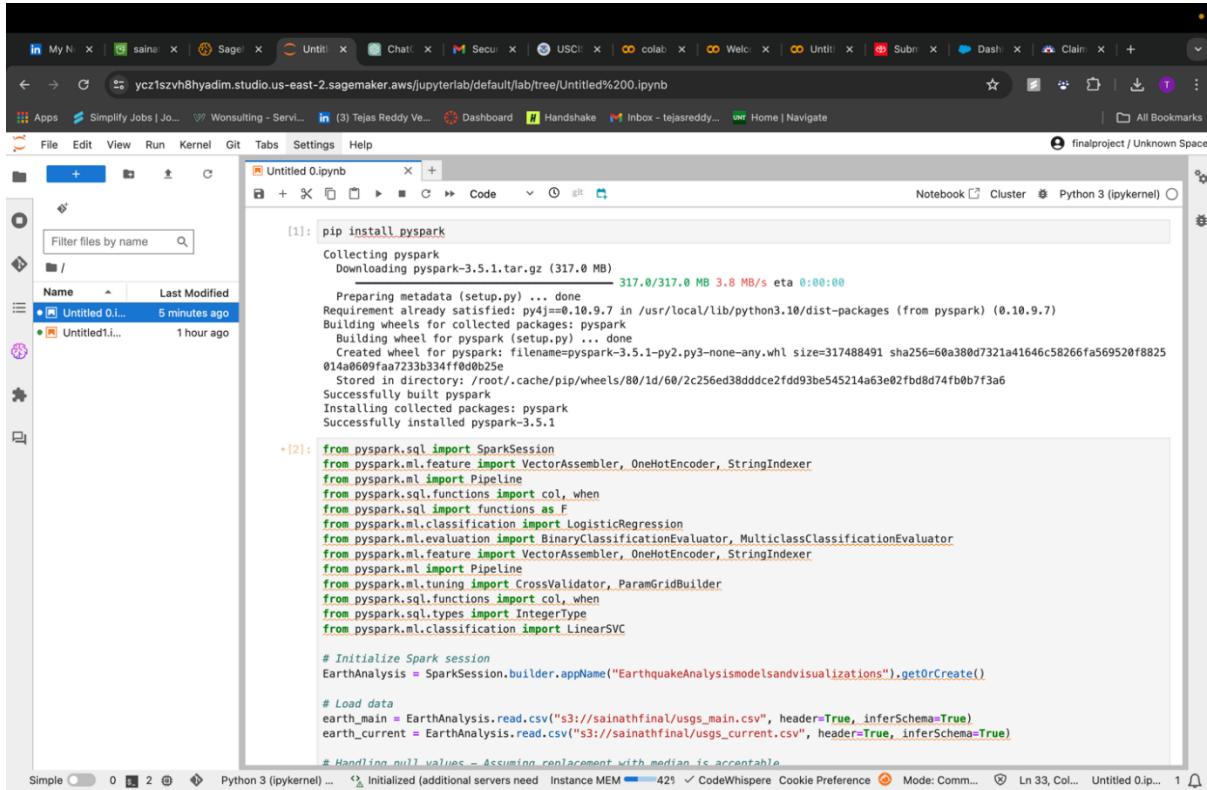
## launching the studio

The screenshot shows the 'Domain details' page for 'QuickSetupDomain-20240506T011069'. The left sidebar includes links for Getting started, Studio, Studio Lab, Canvas, RStudio, TensorBoard, Profiler, Admin configurations (Domains, Role manager, Images, Lifecycle configurations), SageMaker dashboard, Search, and JumpStart (Foundation models, Computer vision models, Natural language processing models). The main content area displays 'User profiles' with a table listing two users: 'finalproject' and 'default-20240506t011069', both created on May 06, 2024, at 06:20 UTC.

## open jupyter lab

The screenshot shows the 'JupyterLab' project settings page in SageMaker Studio. The left sidebar lists Applications (JupyterLab, RStudio, Canvas, Code Editor, Studio CLI), Home, Running instances, Data, Auto ML, Experiments, Jobs, Pipelines, Models, and JumpStart. The main area shows a 'Space Settings' section for a 'jupyter' project, which is private. It includes fields for Storage (5 GB), Lifecycle Configuration (No Script), and Attach custom EFS filesystem - optional (None). A message at the bottom indicates 'Opening project in a new tab'.

## Importing and installing required packages



```
[1]: pip install pyspark
Collecting pyspark
  Downloading pyspark-3.5.1.tar.gz (317.0 MB)
    Preparing metadata (setup.py) ... done
      Requirement already satisfied: py4j==0.10.9.7 in /usr/local/lib/python3.10/dist-packages (from pyspark) (0.10.9.7)
    Building wheels for collected packages: pyspark
      Building wheel for pyspark (setup.py) ... done
        Created wheel for pyspark: filename=pyspark-3.5.1-py2.py3-none-any.whl size=317488491 sha256=60a380d7321a41646c58266fa569520f8825
        Stored in directory: /root/.cache/pip/wheels/80/1d/60/c256ed38ddce2fd93be545214a63e02fb8d74fb0b7f3a6
      Successfully built pyspark
    Installing collected packages: pyspark
      Successfully installed pyspark-3.5.1

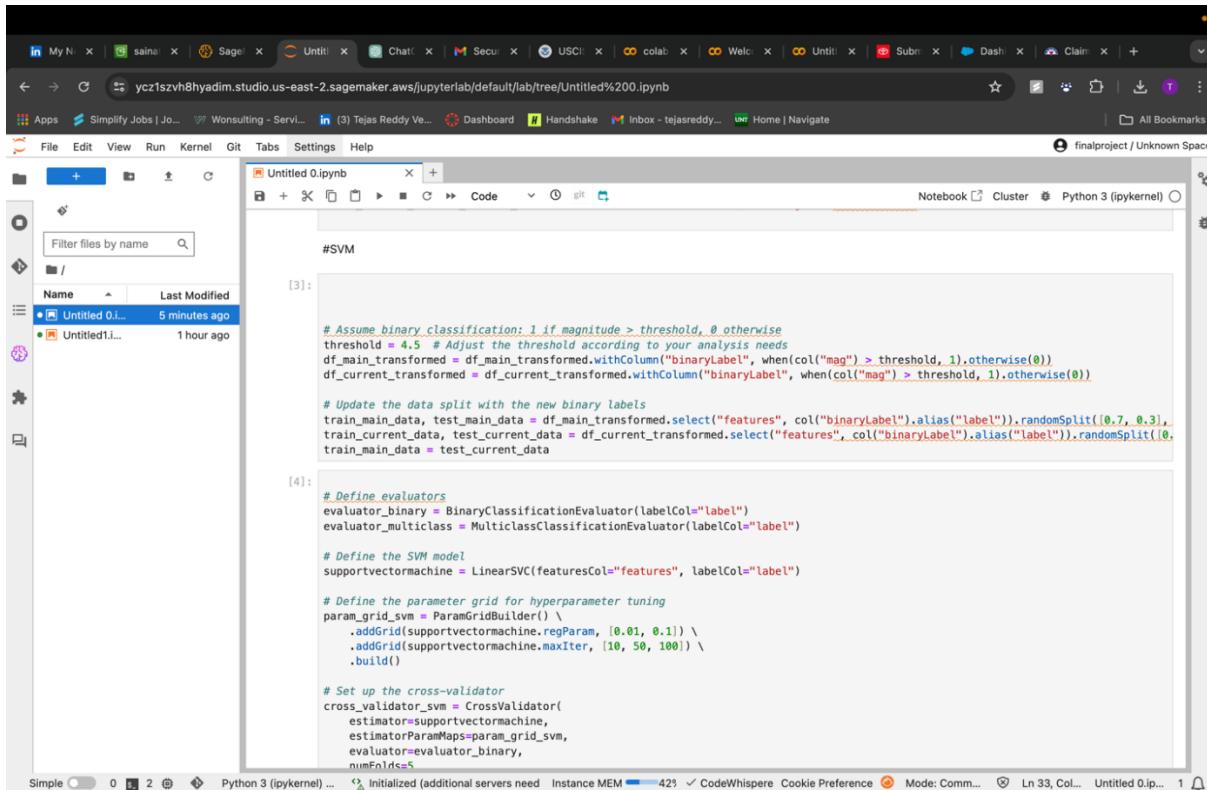
[2]: from pyspark.sql import SparkSession
from pyspark.ml.feature import VectorAssembler, OneHotEncoder, StringIndexer
from pyspark.ml import Pipeline
from pyspark.sql.functions import col, when
from pyspark.sql import functions as F
from pyspark.ml.classification import LogisticRegression
from pyspark.ml.evaluation import BinaryClassificationEvaluator, MulticlassClassificationEvaluator
from pyspark.ml.feature import VectorAssembler, OneHotEncoder, StringIndexer
from pyspark.ml import Pipeline
from pyspark.ml.tuning import CrossValidator, ParamGridBuilder
from pyspark.sql.functions import col, when
from pyspark.sql.types import IntegerType
from pyspark.ml.classification import LinearSVC

# Initialize Spark session
EarthAnalysis = SparkSession.builder.appName("EarthquakeAnalysismodelsandvisualizations").getOrCreate()

# Load data
earth_main = EarthAnalysis.read.csv("s3://sainathfinal/usgs_main.csv", header=True, inferSchema=True)
earth_current = EarthAnalysis.read.csv("s3://sainathfinal/usgs_current.csv", header=True, inferSchema=True)

# Handling null values - assuming replacement with median is acceptable
```

## retrieving data from



```
[3]:
#SVM

# Assume binary classification: 1 if magnitude > threshold, 0 otherwise
threshold = 4.5 # Adjust the threshold according to your analysis needs
df_main_transformed = df_main_transformed.withColumn("binaryLabel", when(col("mag") > threshold, 1).otherwise(0))
df_current_transformed = df_current_transformed.withColumn("binaryLabel", when(col("mag") > threshold, 1).otherwise(0))

# Update the data split with the new binary labels
train_main_data, test_main_data = df_main_transformed.select("features", col("binaryLabel").alias("label")).randomSplit([0.7, 0.3])
train_current_data, test_current_data = df_current_transformed.select("features", col("binaryLabel").alias("label")).randomSplit([0, 1])
train_main_data = test_current_data
test_main_data = test_current_data

[4]:
# Define evaluators
evaluator_binary = BinaryClassificationEvaluator(labelCol="label")
evaluator_multiclass = MulticlassClassificationEvaluator(labelCol="label")

# Define the SVM model
supportvectormachine = LinearSVC(featuresCol="features", labelCol="label")

# Define the parameter grid for hyperparameter tuning
param_grid_svm = ParamGridBuilder() \
    .addGrid(supportvectormachine.regParam, [0.01, 0.1]) \
    .addGrid(supportvectormachine.maxIter, [10, 50, 100]) \
    .build()

# Set up the cross-validator
cross_validator_svm = CrossValidator(
    estimator=supportvectormachine,
    estimatorParamMaps=param_grid_svm,
    evaluator=evaluator_binary,
    numFolds=5)
```

## s3 bucket and creating a SVM (support vector machine) model

```

f1_score_svm = evaluator_multiclass.evaluate(predictions_svm, {evaluator_multiclass.metricName: "f1"})
precision_svm = evaluator_multiclass.evaluate(predictions_svm, {evaluator_multiclass.metricName: "weightedPrecision"})
recall_svm = evaluator_multiclass.evaluate(predictions_svm, {evaluator_multiclass.metricName: "weightedRecall"})

# Print the results
print("Support Vector Machine (SVM) Results:")
print("Accuracy: (accuracy_svm)")
print("F1 Score: (f1_score_svm)")
print("Precision: (precision_svm)")
print("Recall: (recall_svm)")

# Generating and printing the confusion matrix
confusion_matrix_svm = predictions_svm.groupBy("label").pivot("prediction", [0, 1]).count()
confusion_matrix_svm.show()

# Best model parameters
best_model_svm = cv_model_svm.bestModel
print("Best Model Parameters for SVM:")
print(" - RegParam: (best_model_svm._java_obj.getRegParam())")
print(" - MaxIter: (best_model_svm._java_obj.getMaxIter())")

```

Support Vector Machine (SVM) Results:

Accuracy: 0.9230769230769231  
F1 Score: 0.9286510590858416  
Precision: 0.9387057387057387  
Recall: 0.9230769230769231

label	0	1
1	1	4
0	44	3

Best Model Parameters for SVM:  
- RegParam: 0.01  
- MaxIter: 10

## Results: SVM results

```

print("Recall: (recall)")

# Generating and printing the confusion matrix
confusion_matrix = predictions.groupBy("label").pivot("prediction", [0, 1]).count()
confusion_matrix.show()

# Best model parameters
best_model = cv_model.bestModel
print("Best Model Parameters:")
print(" - RegParam: (best_model._java_obj.getRegParam())")
print(" - MaxIter: (best_model._java_obj.getMaxIter())")
print(" - ElasticNetParam: (best_model._java_obj.getElasticNetParam())")

Accuracy: 0.9230769230769231
F1 Score: 0.9286510590858416
Precision: 0.9387057387057387
Recall: 0.9230769230769231

table

```

label	0	1
1	1	4
0	44	3

Best Model Parameters:  
- RegParam: 0.01  
- MaxIter: 10  
- ElasticNetParam: 0.0

#decision tree

```

[7]: from pyspark.ml.classification import DecisionTreeClassifier
      # Define the Decision Tree model
      decision = DecisionTreeClassifier(featuresCol="features", labelCol="label")
      # Define the parameter grid for hyperparameter tuning
      param_grid_dt = ParamGridBuilder()

```

## Design tree results

The screenshot shows a Jupyter Notebook interface with a single open cell titled "Untitled 0.ipynb". The code in the cell is as follows:

```

# Generating and printing the confusion matrix
confusion_matrix_dt = predictions_dt.groupBy("label").pivot("prediction", [0, 1]).count()
confusion_matrix_dt.show()

# Best model parameters
best_model_dt = cv_model_dt.bestModel
print("Best Model Parameters for Decision Trees:")
print(f"- MaxDepth: {best_model_dt._java_obj.getMaxDepth()}")
print(f"- MaxBins: {best_model_dt._java_obj.getMaxBins()}")


Decision Tree Results:
Accuracy: 0.8846153846153846
F1 Score: 0.8929765886287626
Precision: 0.9048840048840049
Recall: 0.8846153846153846

+-----+
|label| 0| 1|
+-----+
| 1| 2| 3|
| 0| 43| 4|
+-----+


Best Model Parameters for Decision Trees:
- MaxDepth: 5
- MaxBins: 60

[8]: # Print the decision tree's structure
print("Decision Tree Model Structure:")
print(best_model_dt.toDebugString)

```

The output of the code is displayed below the code cell, showing the decision tree's structure and its performance metrics.

## decision tree structure

```

Decision Tree Model Structure:
DecisionTreeClassificationModel: uid=DecisionTreeClassifier_7702185b03b5, depth=3, numNodes=7, numClasses=2, numFeatures=16
If (feature 10 <= 1.2185)
  Predict: 0.0
Else (feature 10 > 1.2185)
  If (feature 6 <= -176.73579999999998)
    Predict: 0.0
  Else (feature 6 > -176.73579999999998)
    If (feature 8 <= 17.5)
      Predict: 0.0
    Else (feature 8 > 17.5)
      Predict: 1.0

```

random forest results

```
Random Forest Results:
Accuracy: 0.8846153846153846
F1 Score: 0.8929765886287626
Precision: 0.9048840048840049
Recall: 0.8846153846153846
+-----+
|label| 0| 1|
+-----+
| 1| 2| 3|
| 0| 43| 4|
+-----+
```

Best Model Parameters for Random Forest:

- MaxDepth: 5
- NumTrees: 20

printing the importance of the features

```
# Print feature importances
importances = best_model_rf.featureImportances
print("Feature Importances:", importances)
# List of input columns to the VectorAssembler
input_columns = assembler_ins # This list should have been defined when you set up your VectorAssembler

# Creating a dictionary to map feature vector indices back to feature names
feature_index = {i: col for i, col in enumerate(input_columns)}

# Optionally, convert feature importances to a more readable format
# Assuming you have a list of feature names corresponding to indices in the features vector
feature_names = [feature_index[i] for i in range(len(feature_index))] # Reuse feature_index from previous examples
feature_importance_dict = dict(zip(feature_names, importances))
print("Readable Feature Importances:", feature_importance_dict)

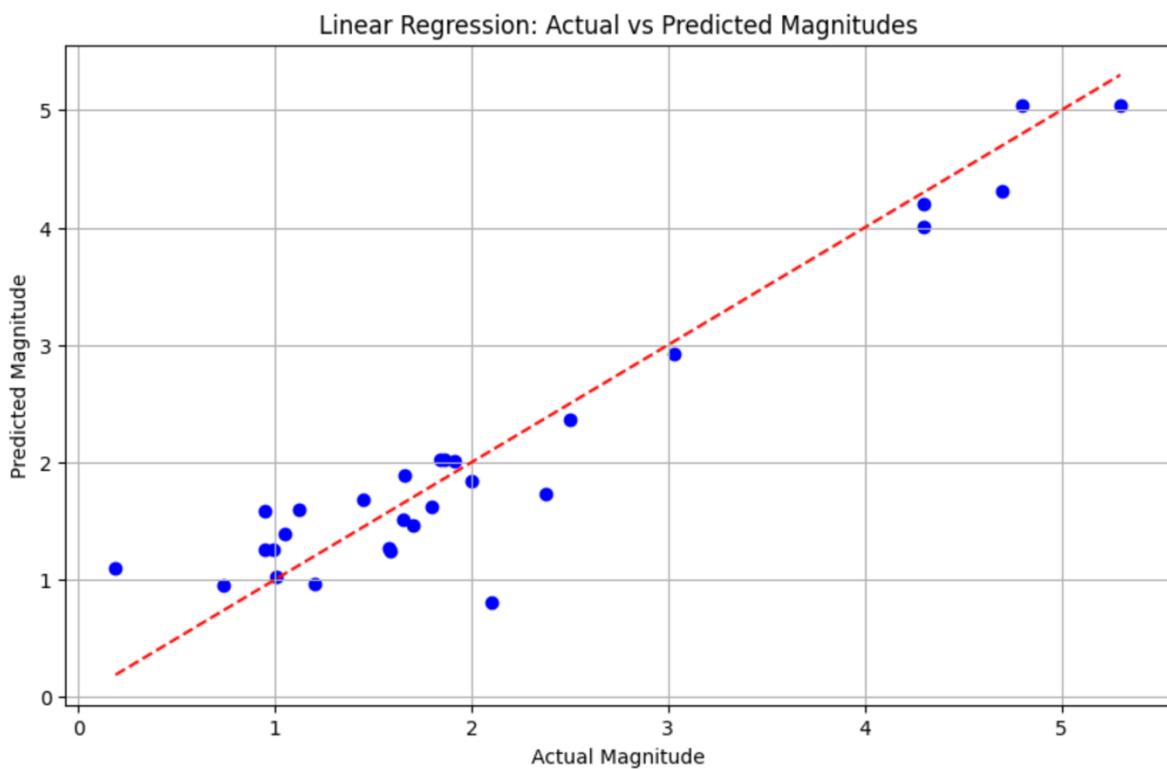
Feature Importances: (16,[0,2,5,6,7,8,9,10,11,12,13,14,15],[0.0060206460206460215,0.17621730073258485,0.12445300689048641,0.0638252926035082,0.0618767779534876,0.03896364929171424,0.06349353430548997,0.1596512358758316,0.0545426105365245,0.11474855478793751,0.029275085232850585,0.008660130718954244,0.09827217519128048])
Readable Feature Importances: {'magTypeOHE': 0.0060206460206460215, 'typeOHE': 0.0, 'latitude': 0.17621730073258485, 'longitude': 0.0, 'depth': 0.0, 'nst': 0.12445300689048641, 'gap': 0.0638252926035082, 'dmin': 0.0618767779534876, 'rms': 0.03896364929171424, 'horizontalError': 0.06349353430548997, 'depthError': 0.1596512358758316, 'magError': 0.0545426105365245, 'magNst': 0.11474855478793751}
```

## Linear Regression plots

prediction	label	features
5.04192489342752	4.8	[0.0,0.0,1.0,0.0,...]
4.009377612272793	4.3	[0.0,0.0,1.0,0.0,...]
4.316257388620885	4.7	[0.0,0.0,1.0,0.0,...]
5.038339507460163	5.3	[0.0,0.0,1.0,0.0,...]
4.2042868440534935	4.3	[0.0,0.0,1.0,0.0,...]

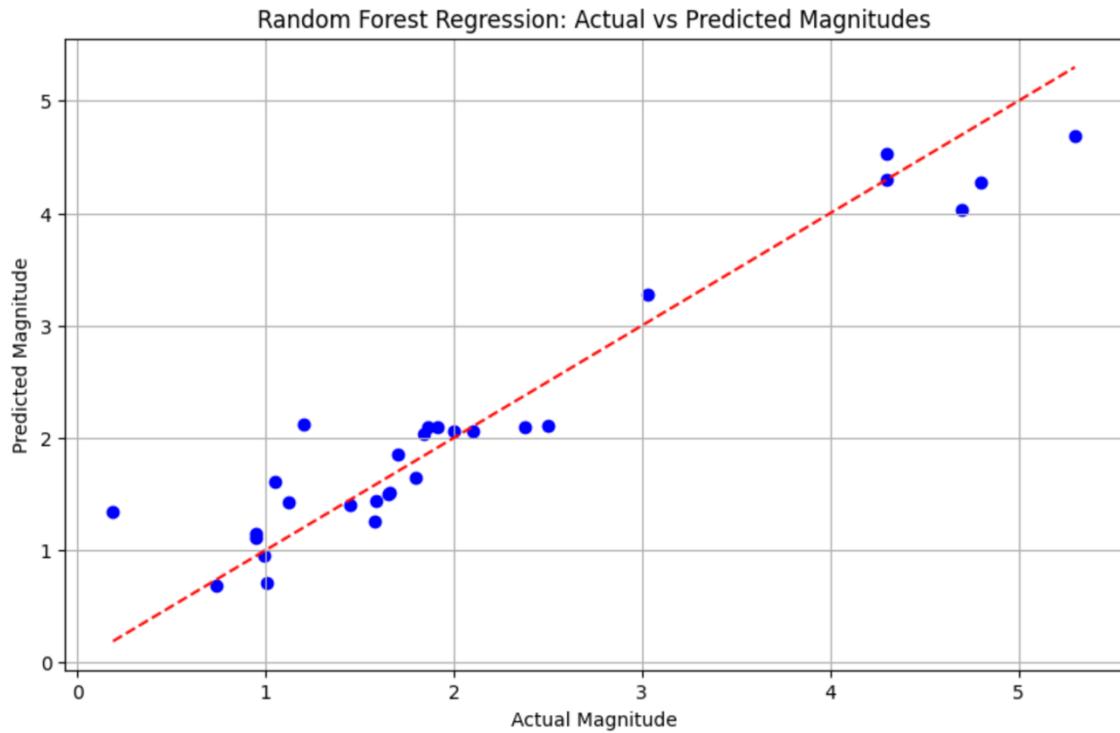
only showing top 5 rows

Root Mean Squared Error (RMSE) on test data = 0.410369

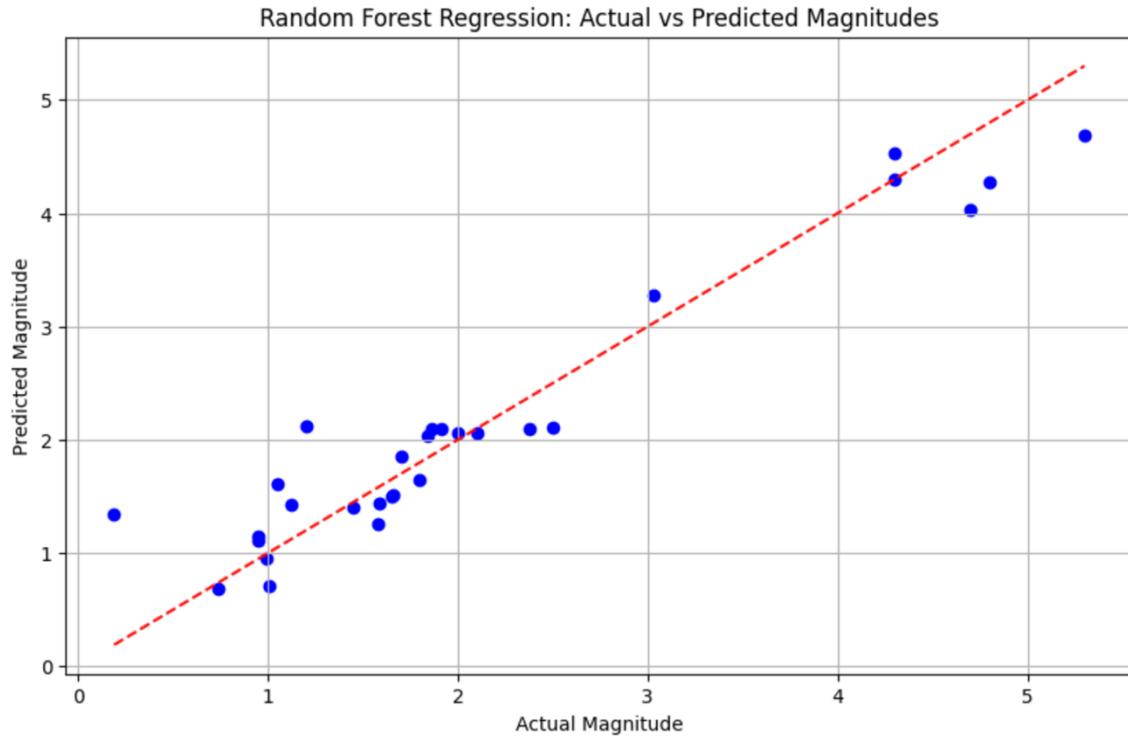


## Random forest plots

Root Mean Squared Error (RMSE) on test data with Random Forest = 0.396129



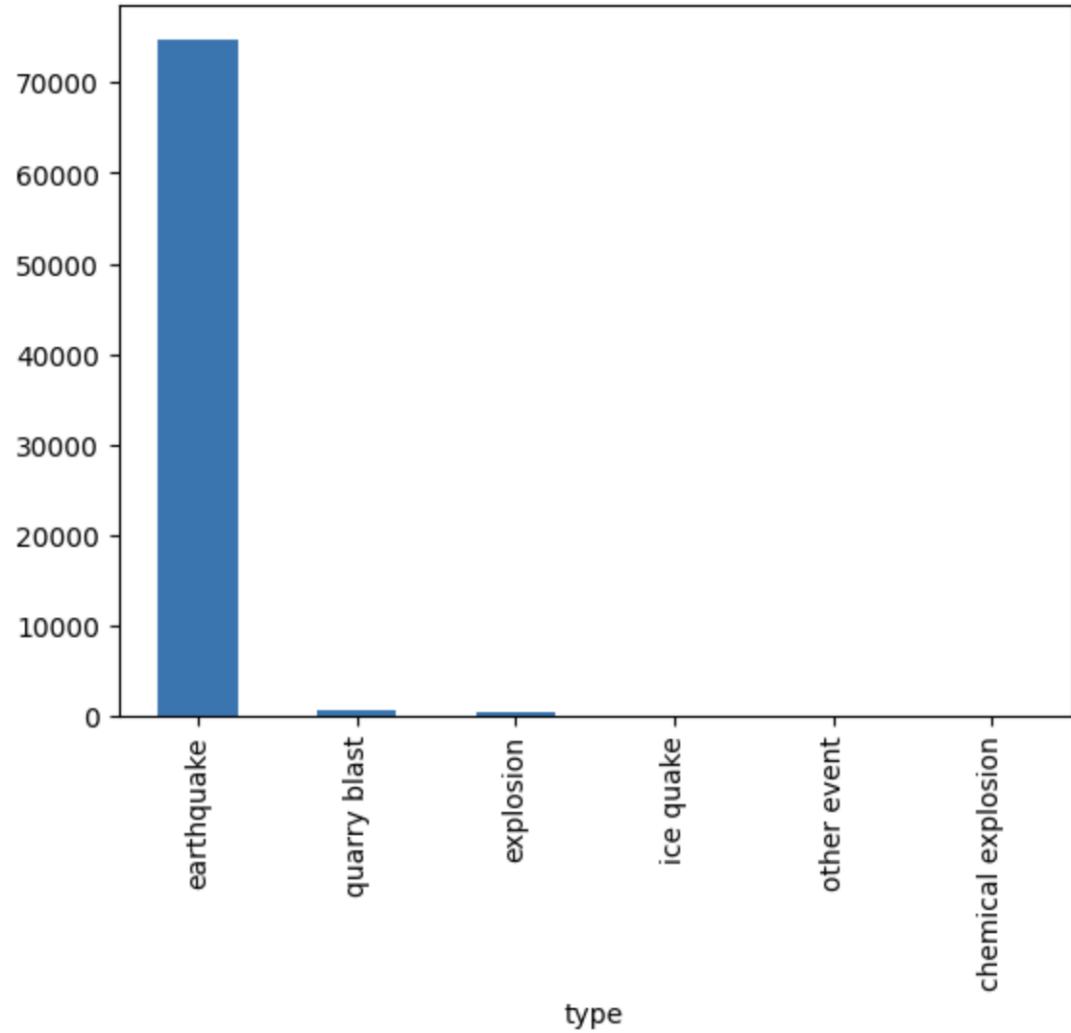
Root Mean Squared Error (RMSE) on test data with Random Forest = 0.396129



## Data Visualization:

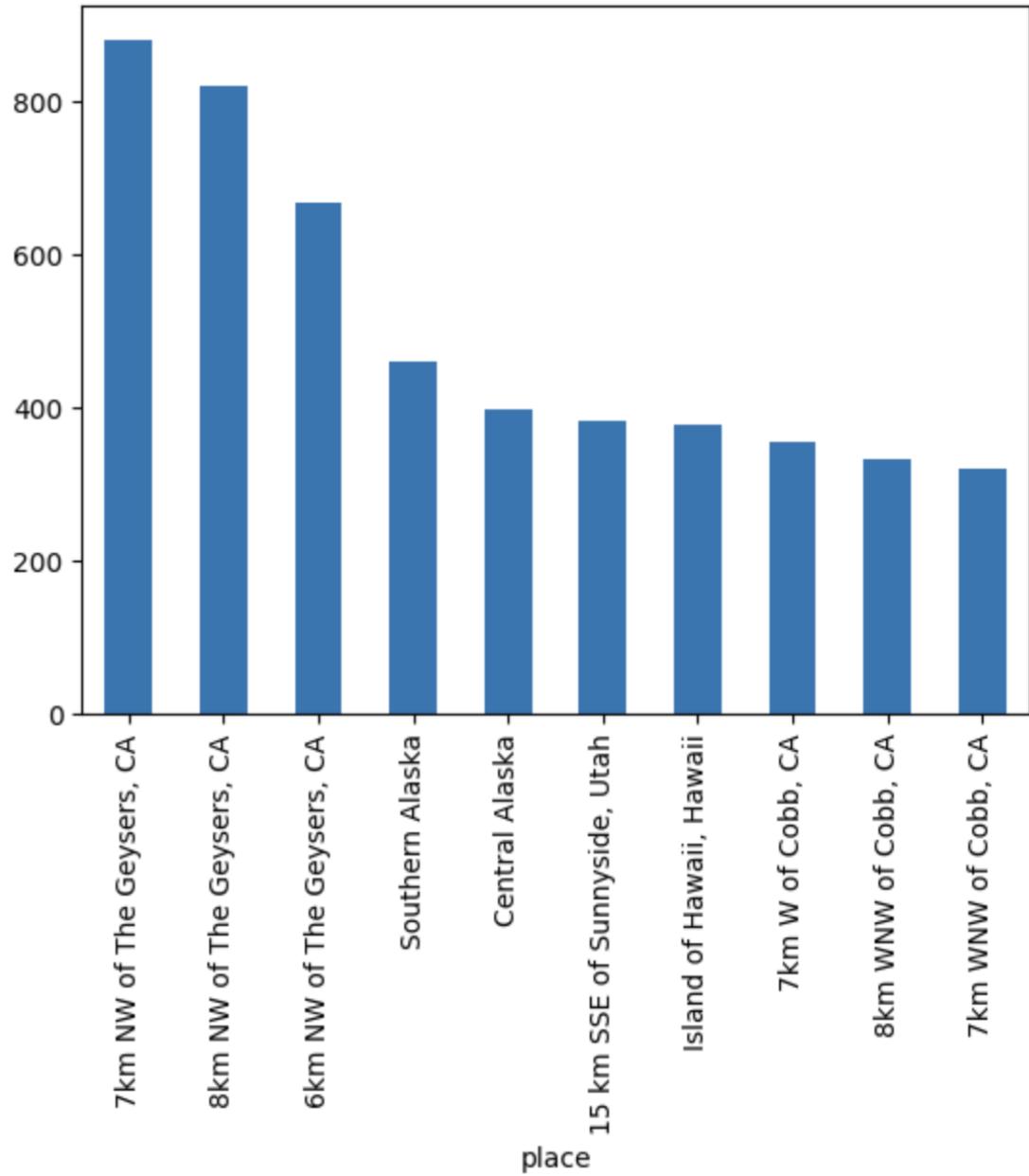
type of event

<Axes: xlabel='type'>

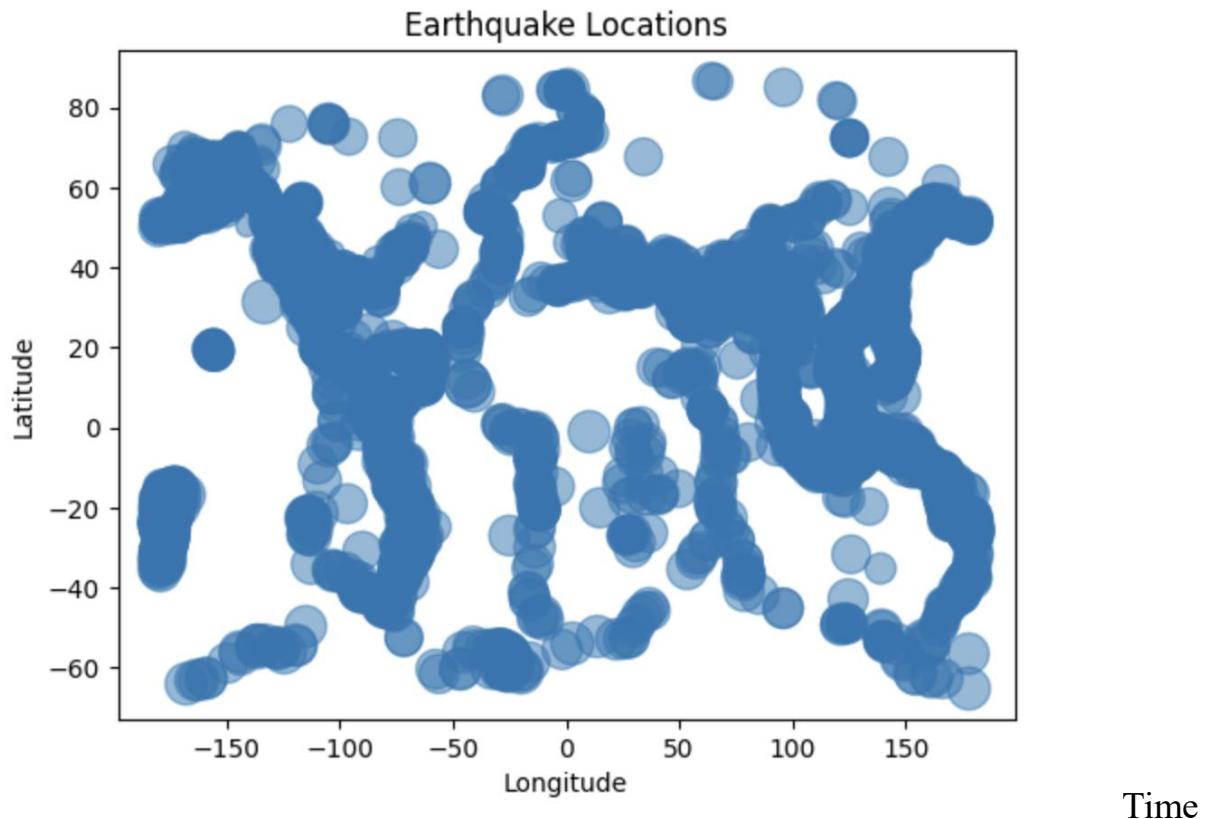


## place vs event

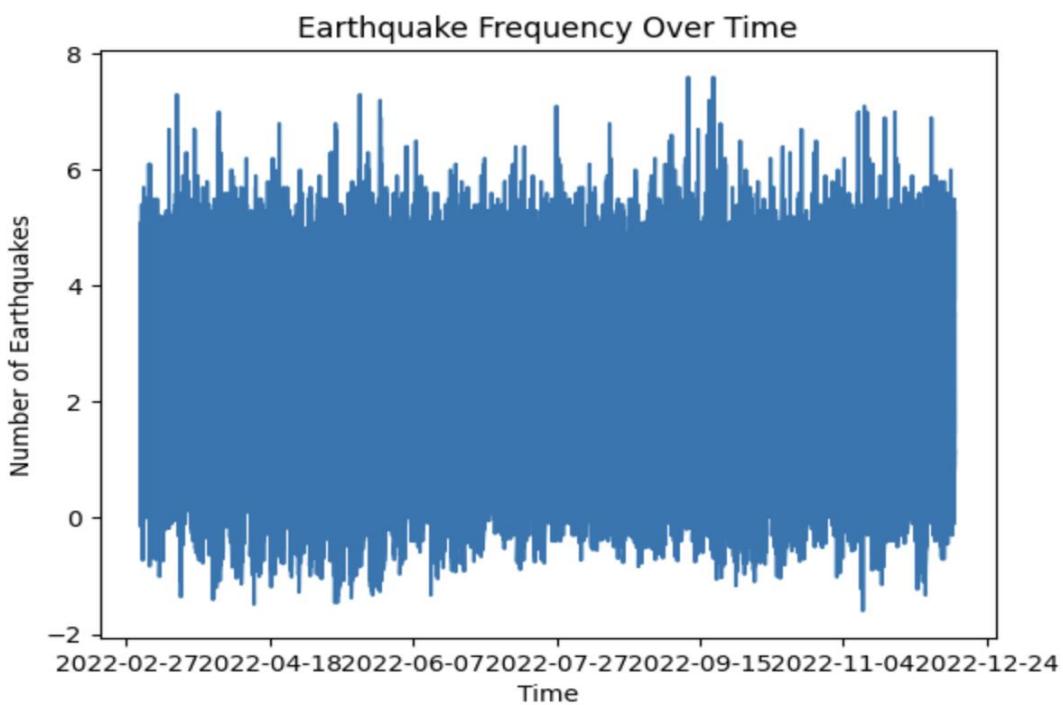
&lt;Axes: xlabel='place'&gt;



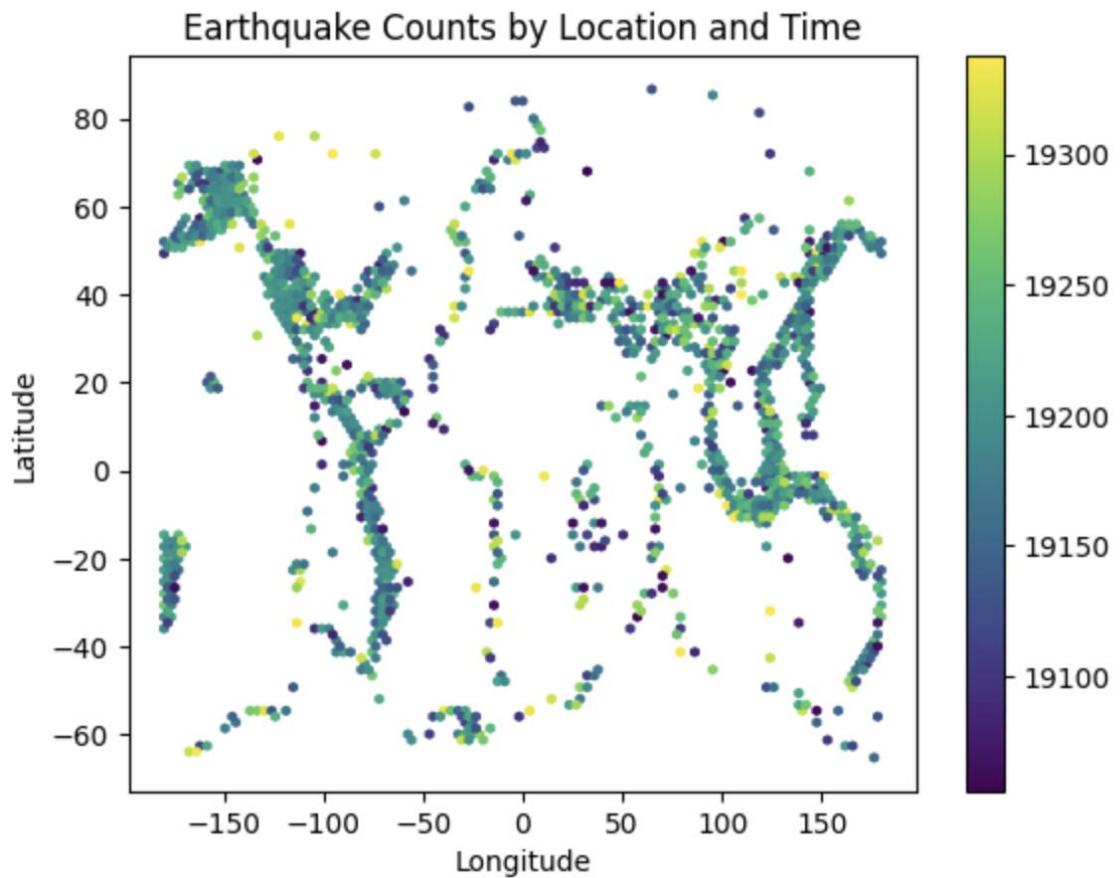
longitude vs latitude to show location worldwide



vs number of earthquakes



### Earthquake counts by location and time



## Dealing with null values

```
US_df.isnull().sum()
```

```
time          0
latitude      0
longitude     0
depth         0
mag           14
magType       14
place          0
type          0
state          0
severity       0
depth category 0
dtype: int64
```

```
null_mag = US_df[US_df['mag'].isnull()].index
US_df.drop(null_mag, inplace = True)
US_df.isnull().sum()
```

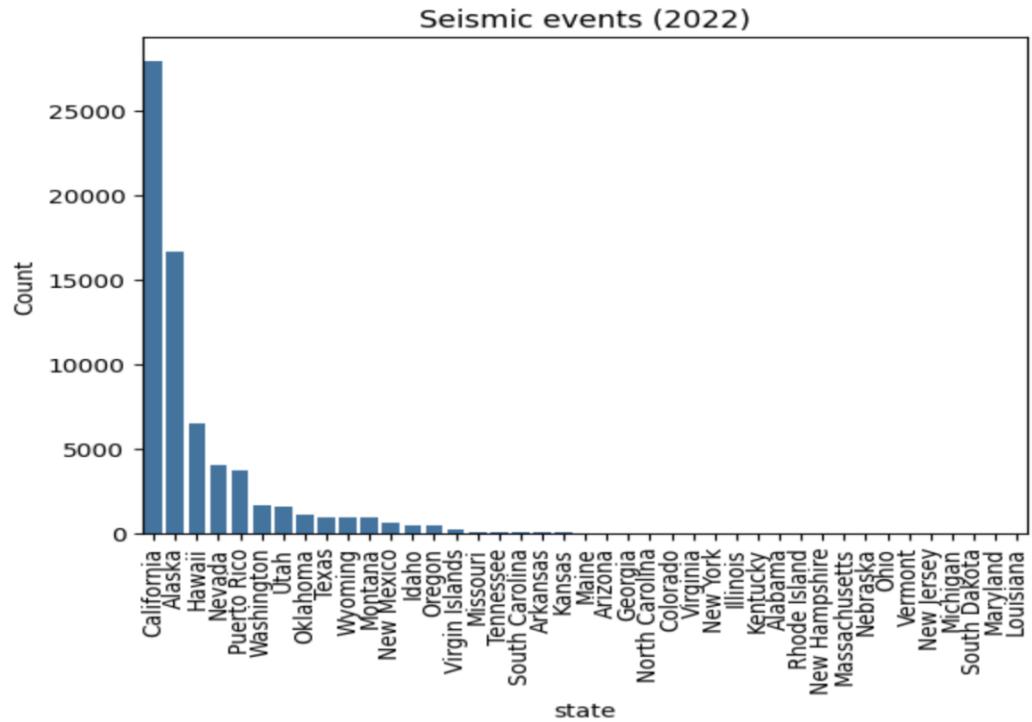
```
time          0
latitude      0
longitude     0
depth         0
mag           0
magType       0
place          0
type          0
state          0
severity       0
depth category 0
dtype: int64
```

## Original and changed values

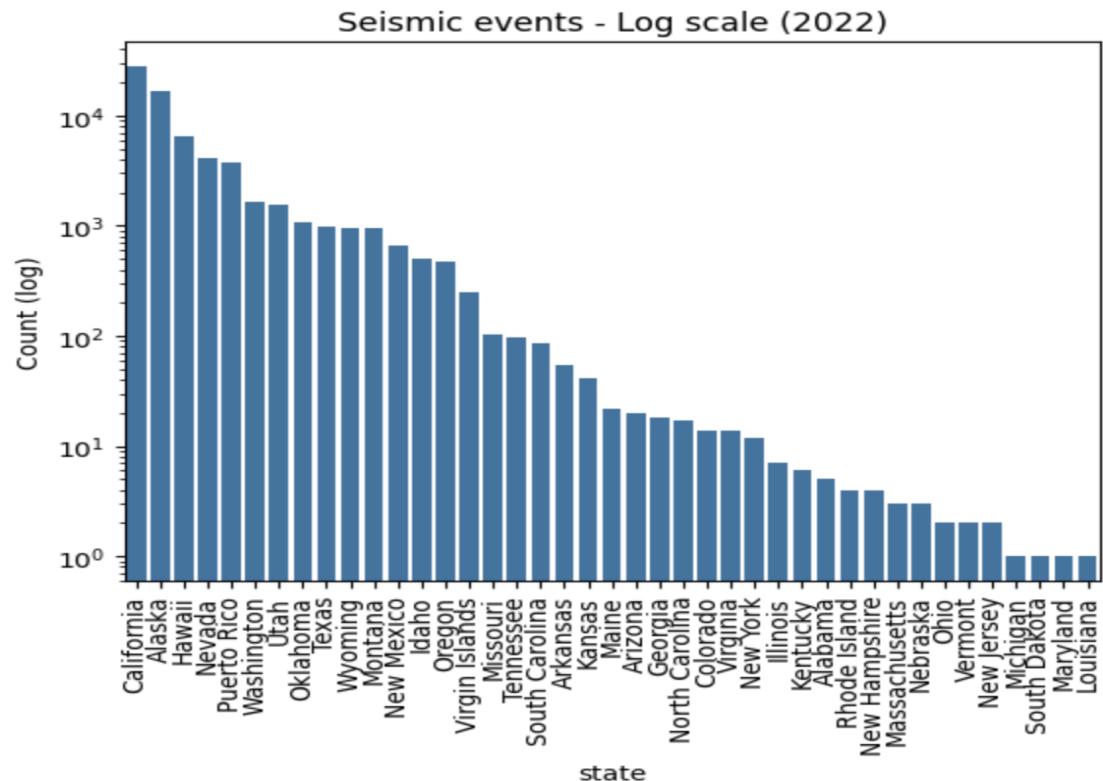
```
print('Original value counts for magnitude types:')
print(US_df['magType'].value_counts())
US_df['magType'].replace(['Ml', 'mww', 'Mb'], ['ml', 'mw', 'mb'], inplace = True)
print('Corrected value counts for magnitude types:')
print(US_df['magType'].value_counts())

Original value counts for magnitude types:
magType
ml      43344
md      24598
mb      244
mb_lg    125
mh       52
mw       32
mlv      25
mww      17
mwr      16
Ml       12
mlr      11
Mi       3
Mb       2
Name: count, dtype: int64
Corrected value counts for magnitude types:
magType
ml      43356
md      24598
mb      246
mb_lg    125
mh       52
mw       49
mlv      25
mwr      16
mlr      11
Mi       3
Name: count, dtype: int64
```

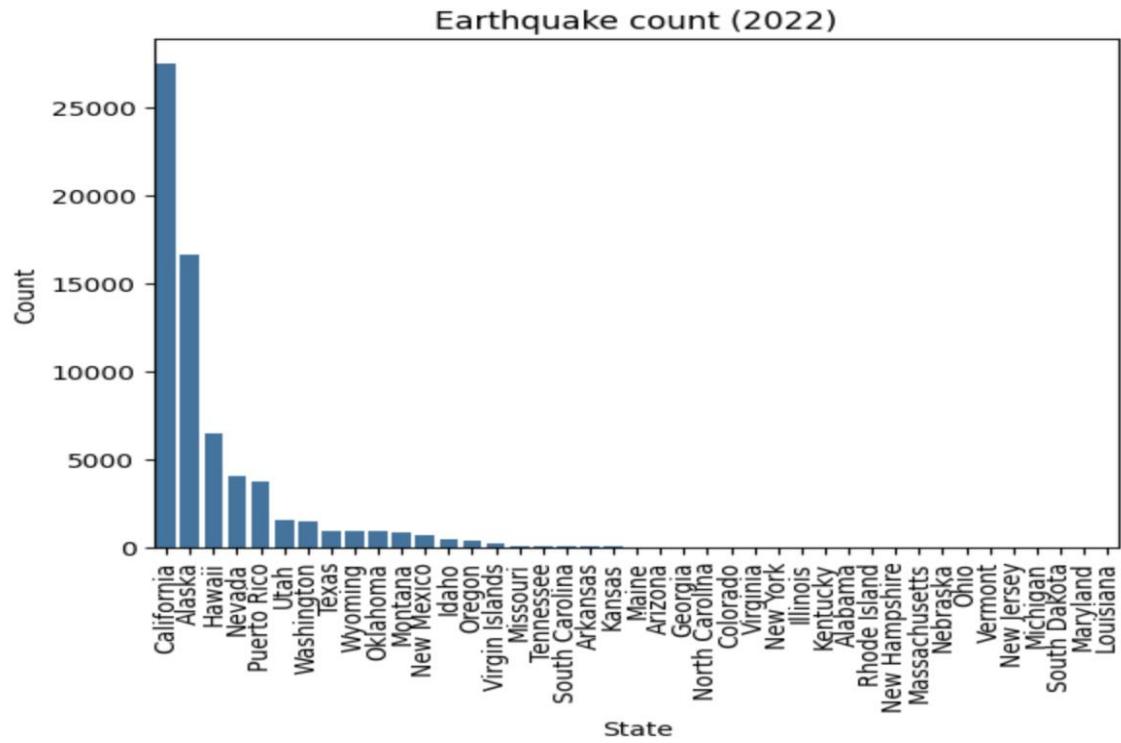
### State vs count



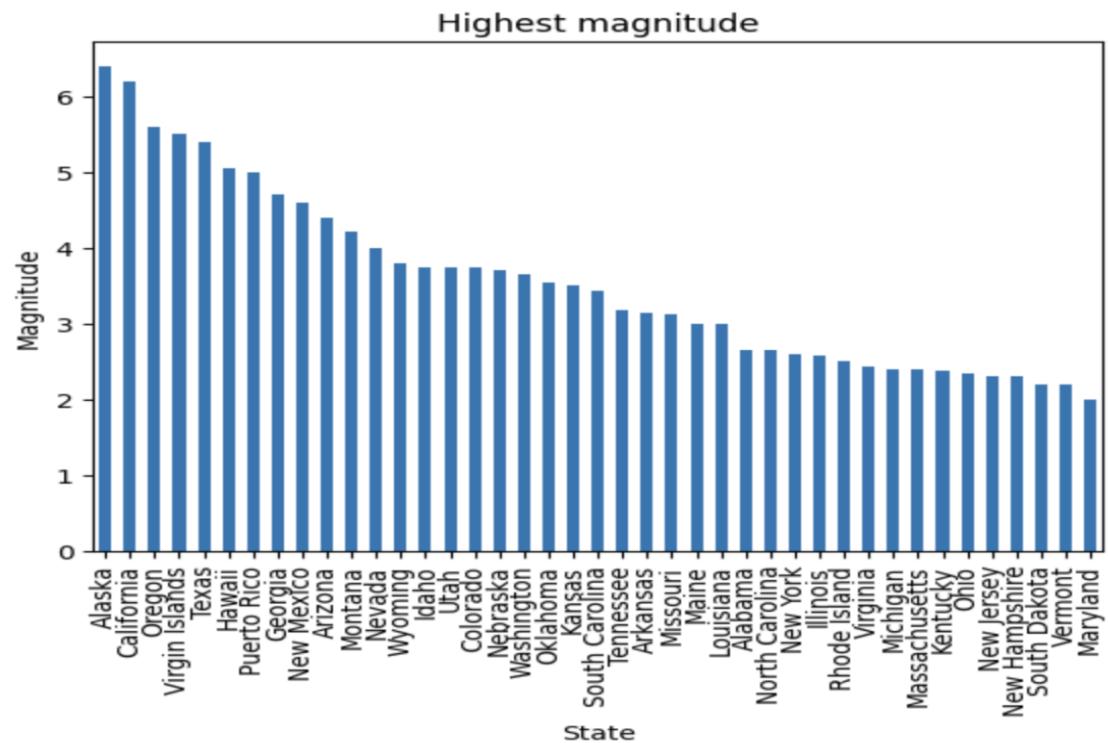
### States count on a log scale



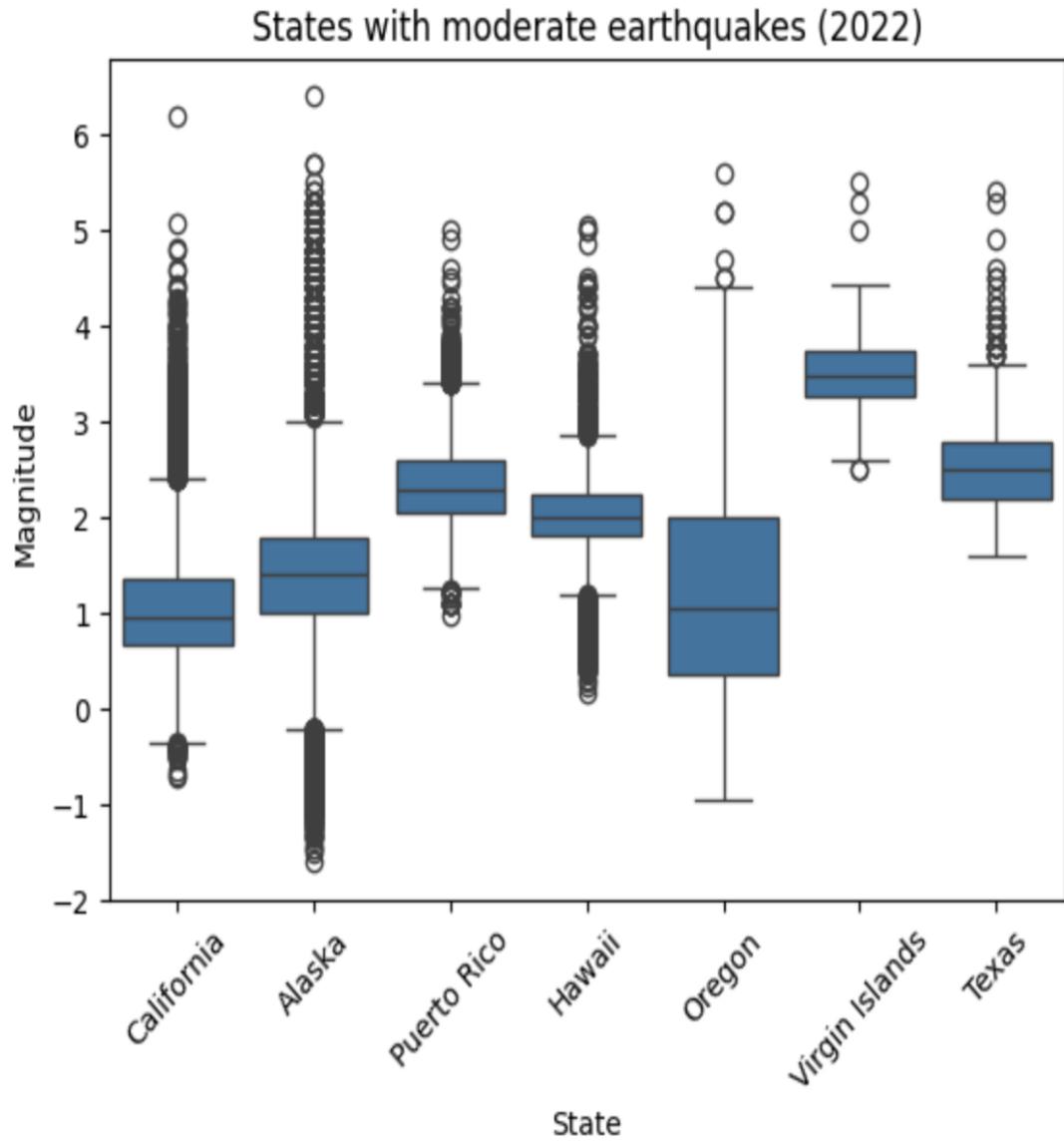
## Earthquake count in year 2022



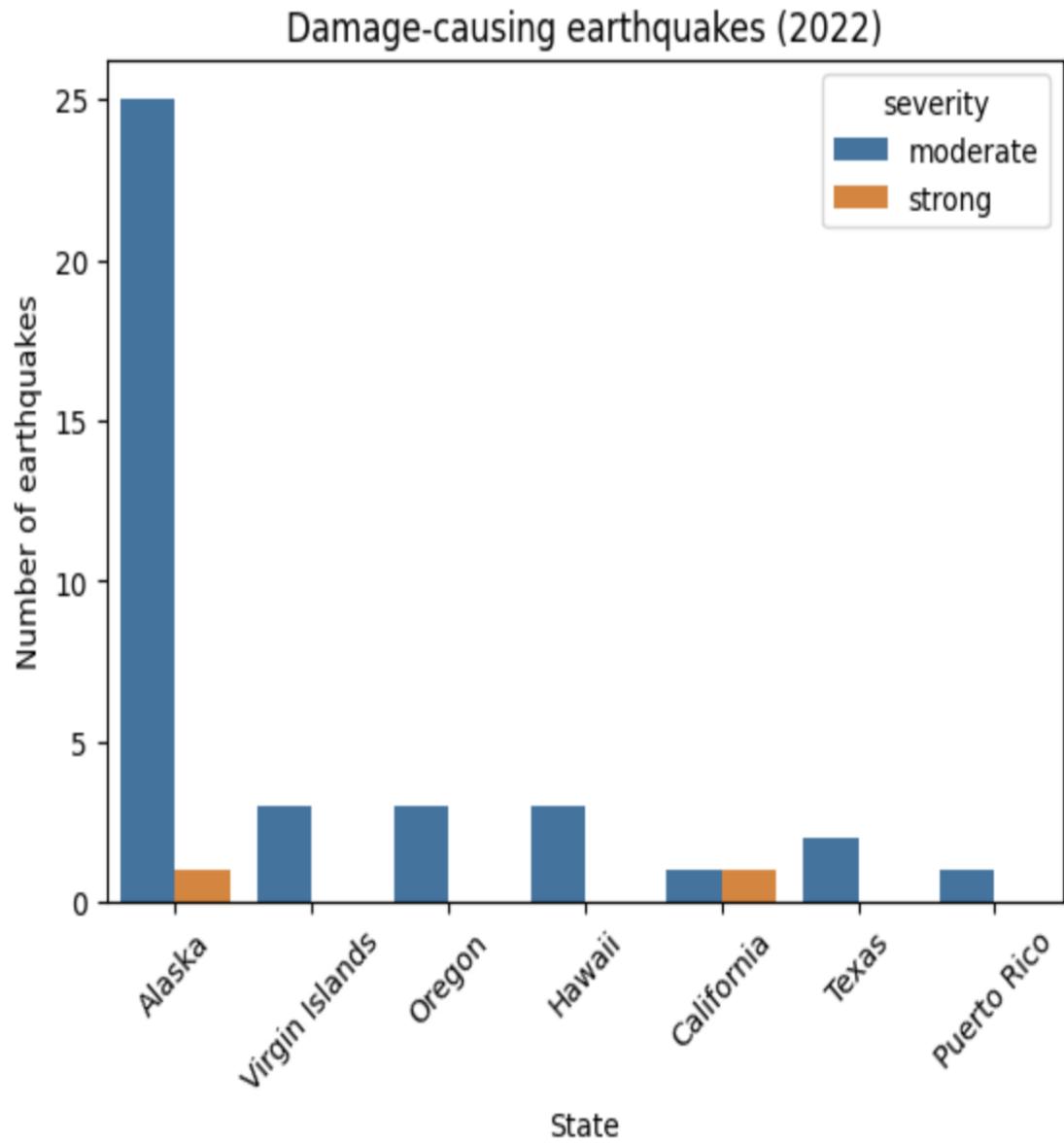
## States with highest magnitude of earthquakes



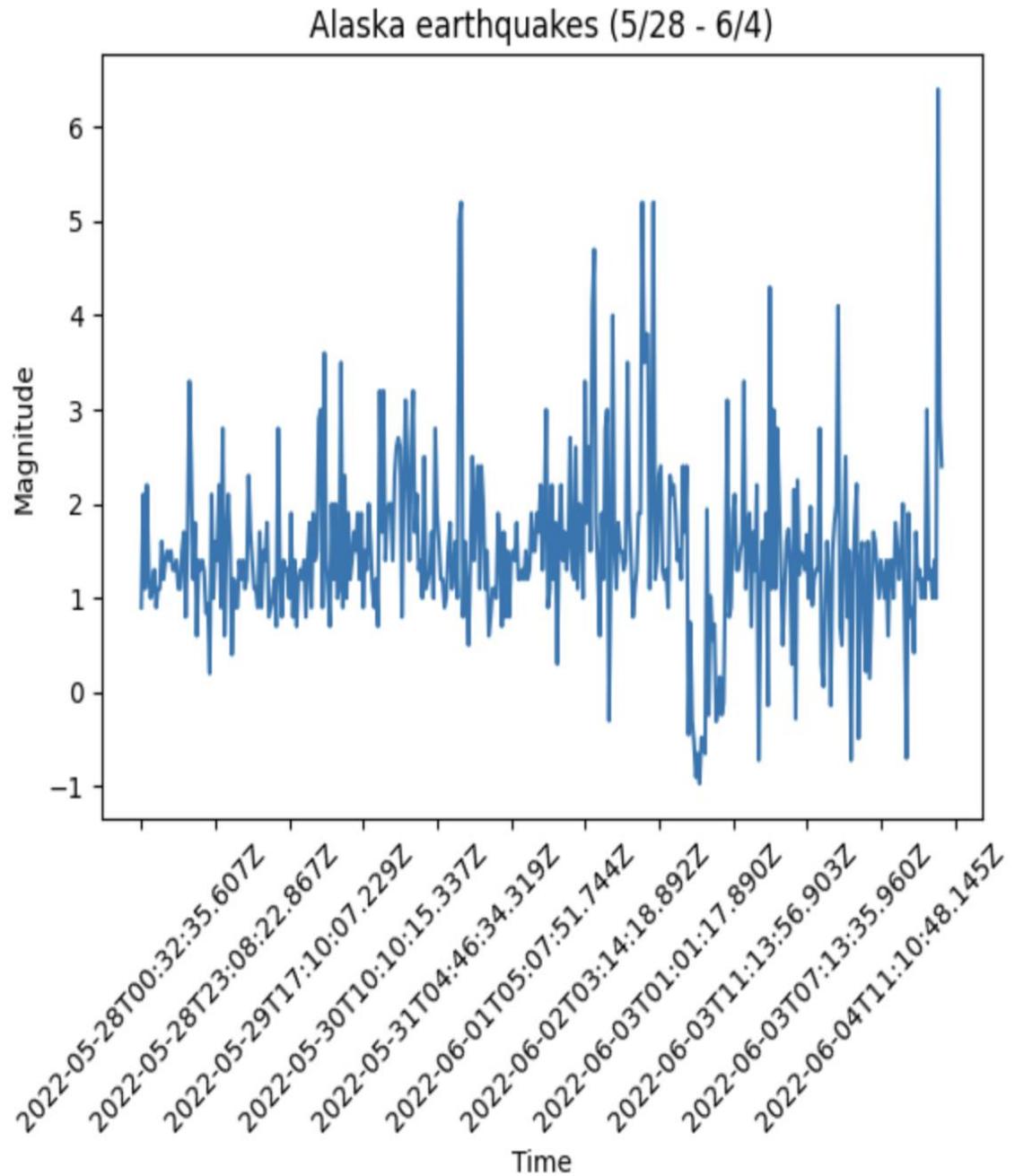
## States with less number of earthquakes



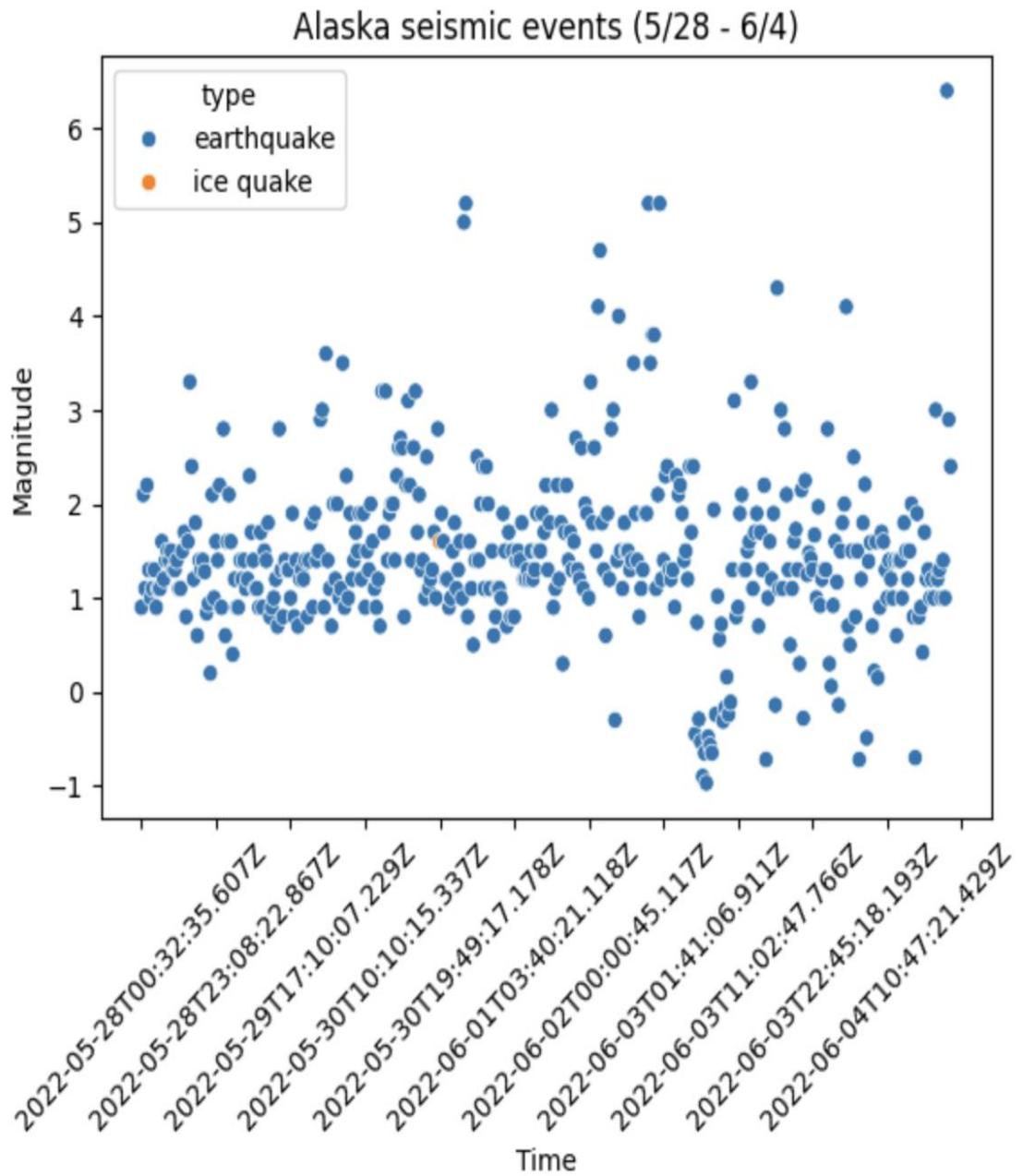
Severity of earthquakes compared

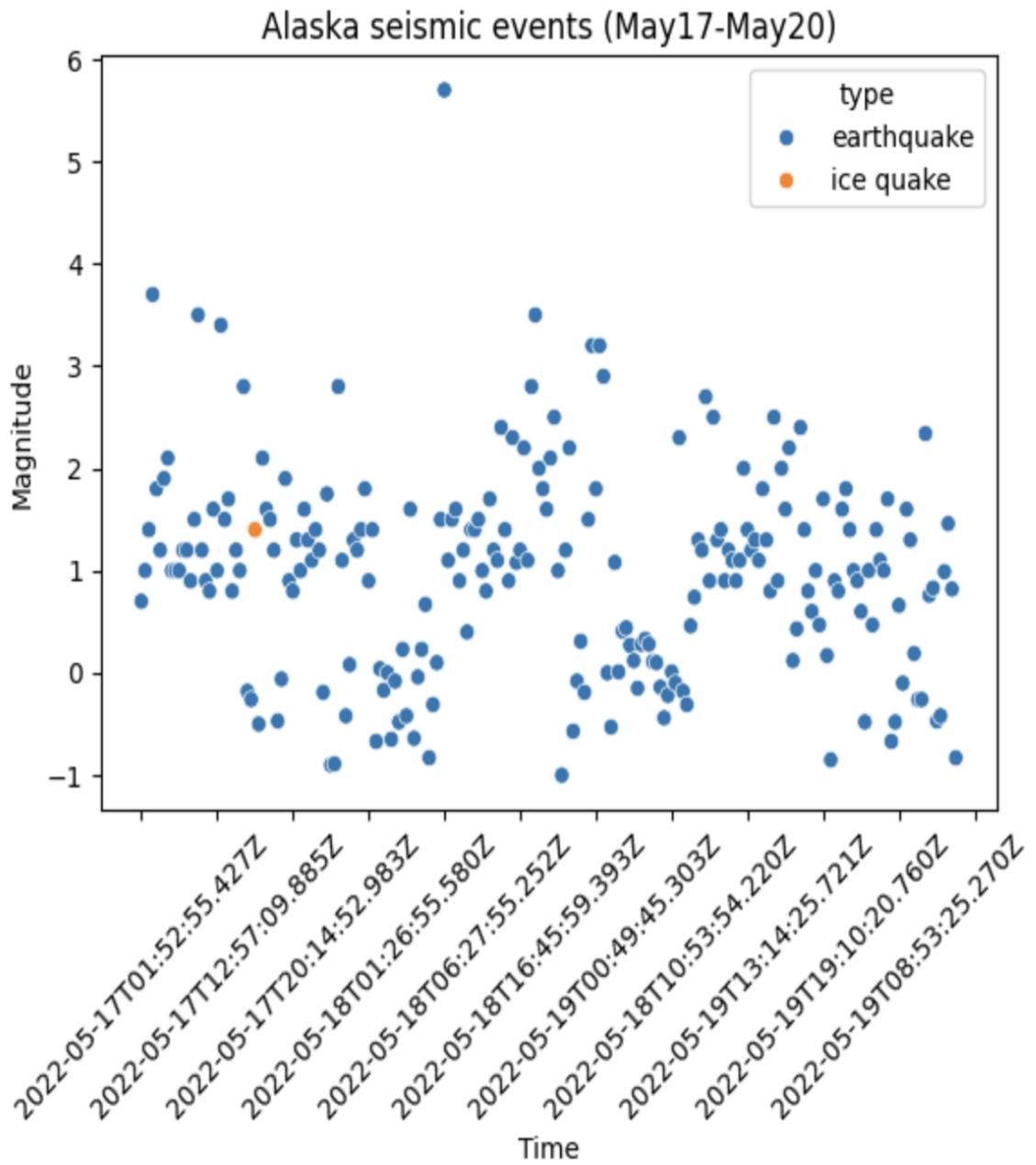


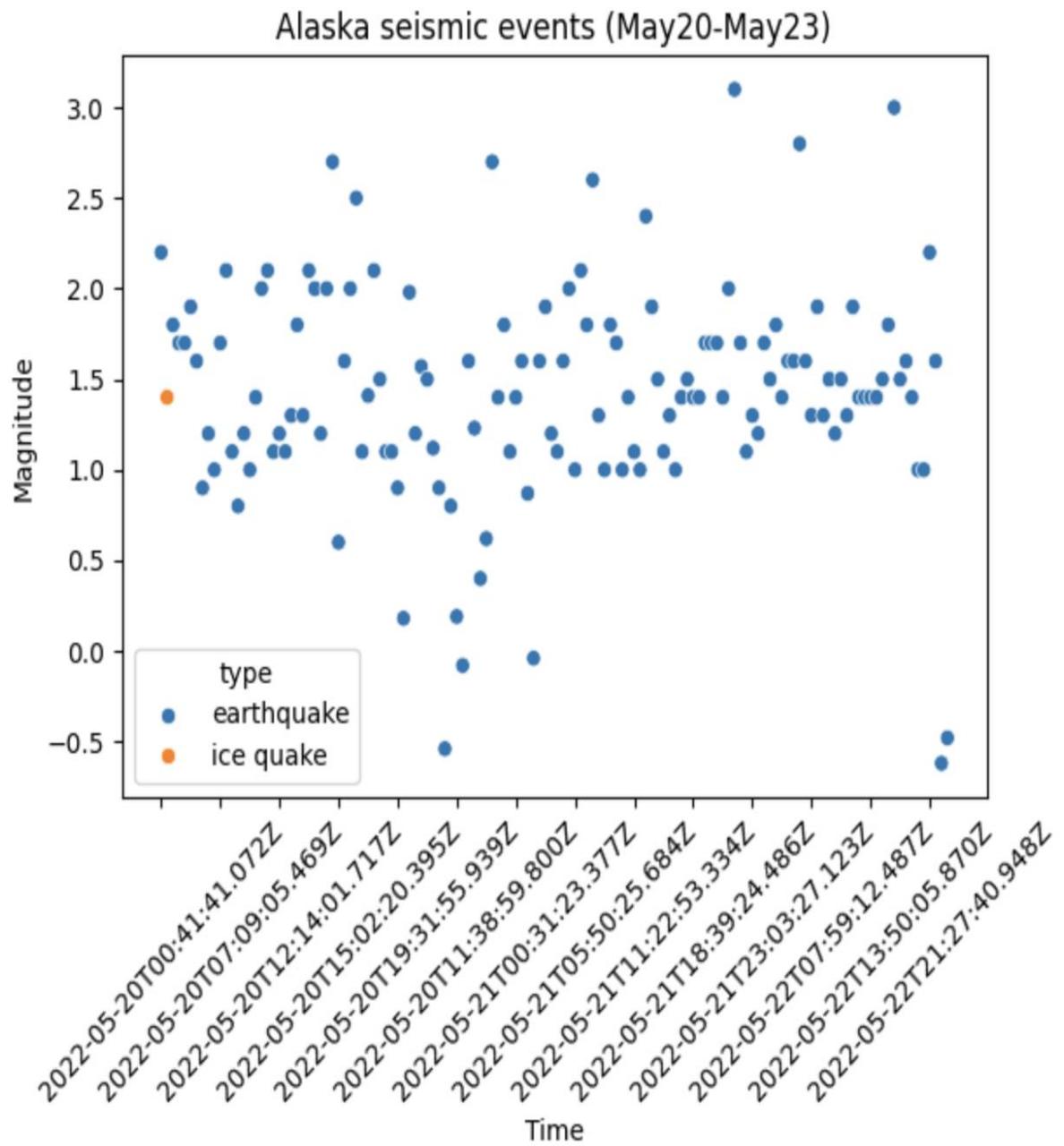
## Earthquakes in Alaska earthquakes

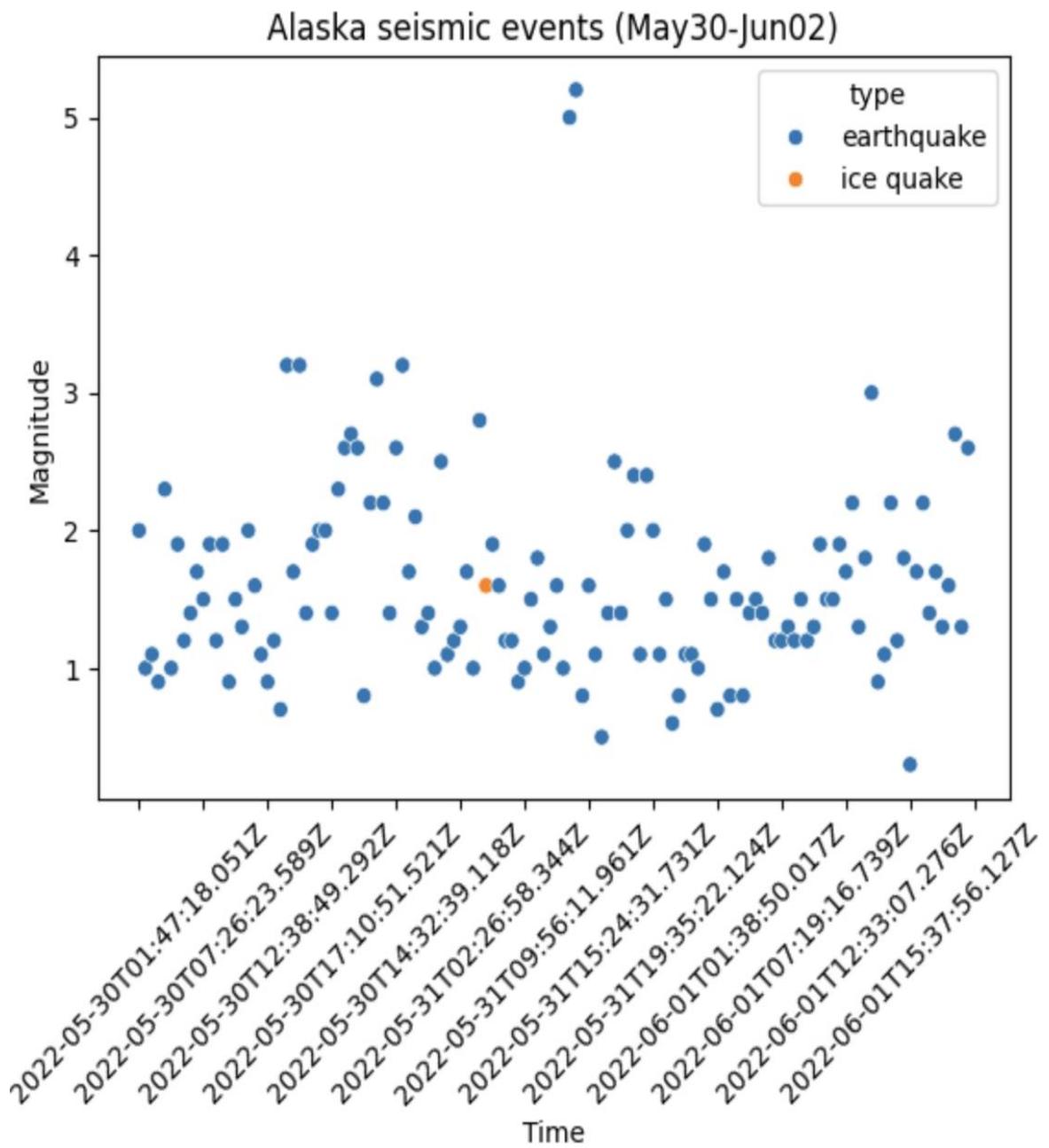


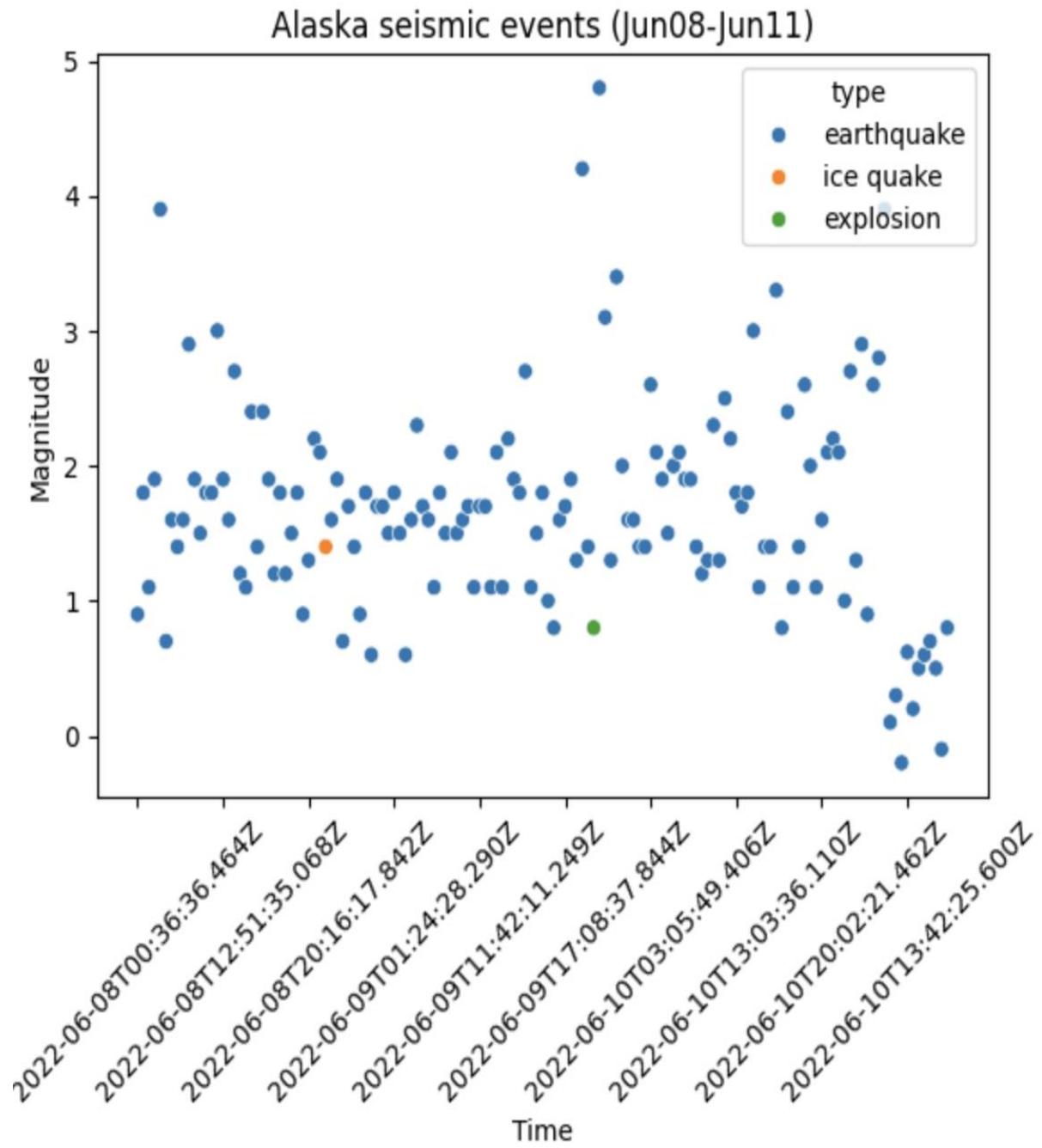
Average and other metrics of California earthquake Alaska seismic events

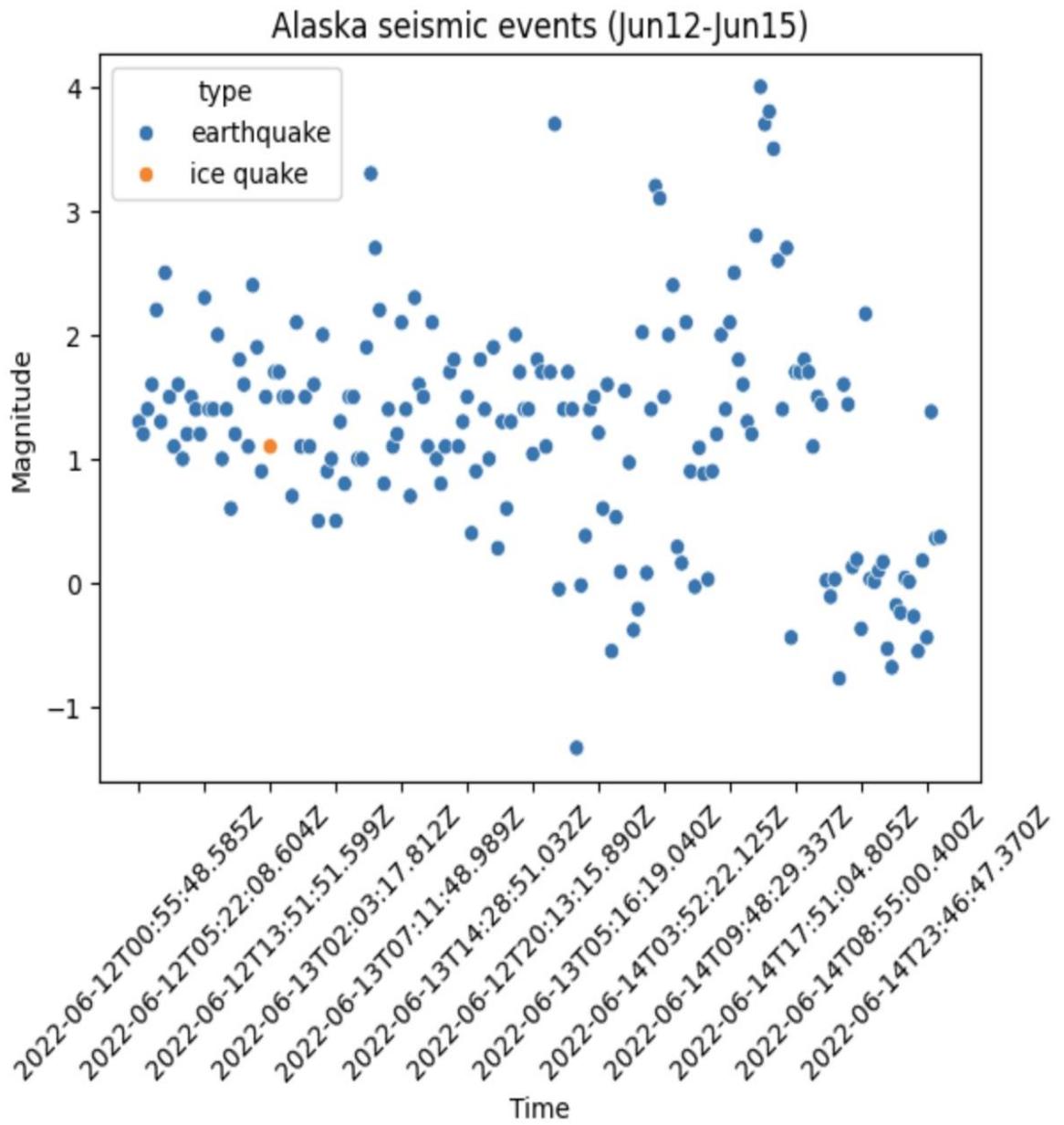


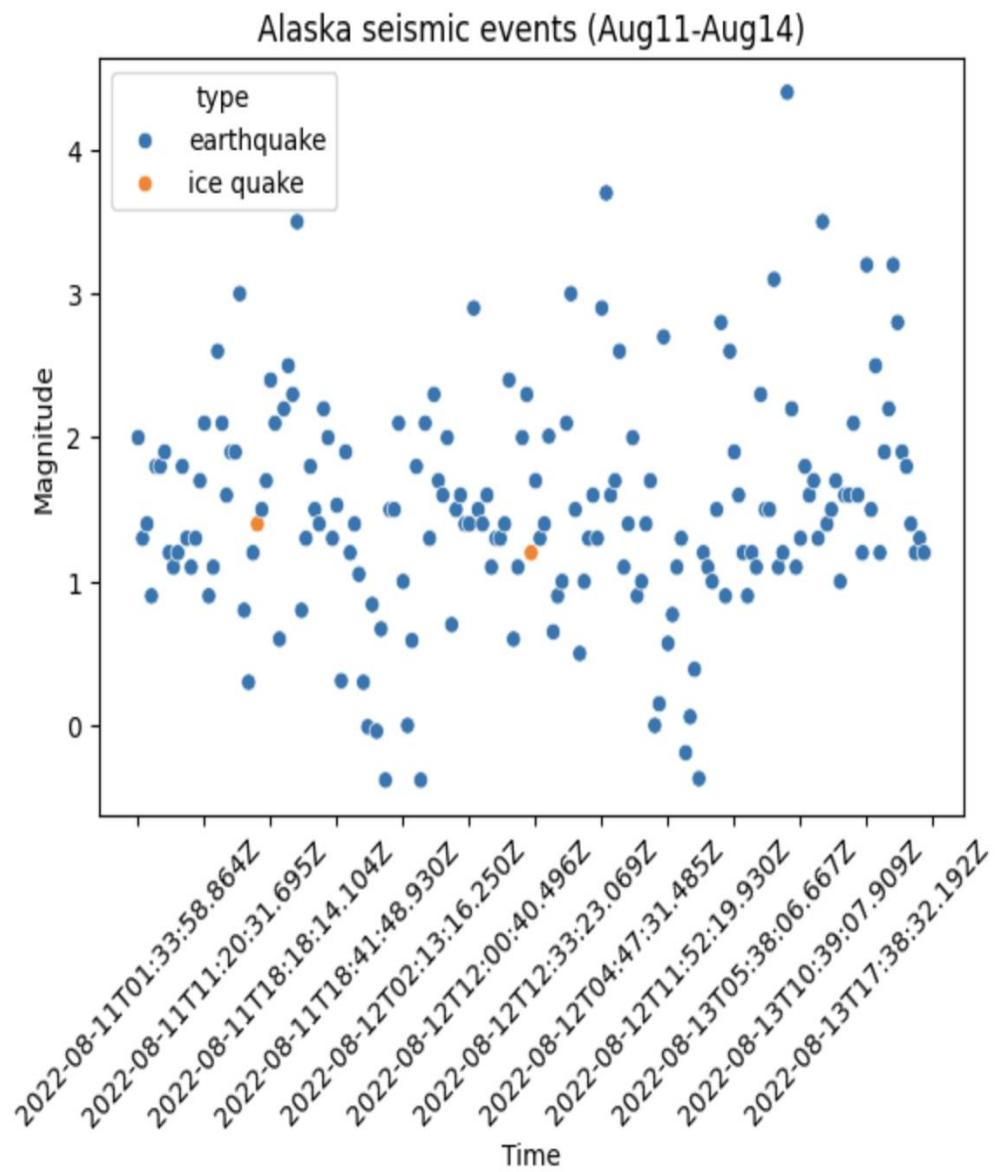


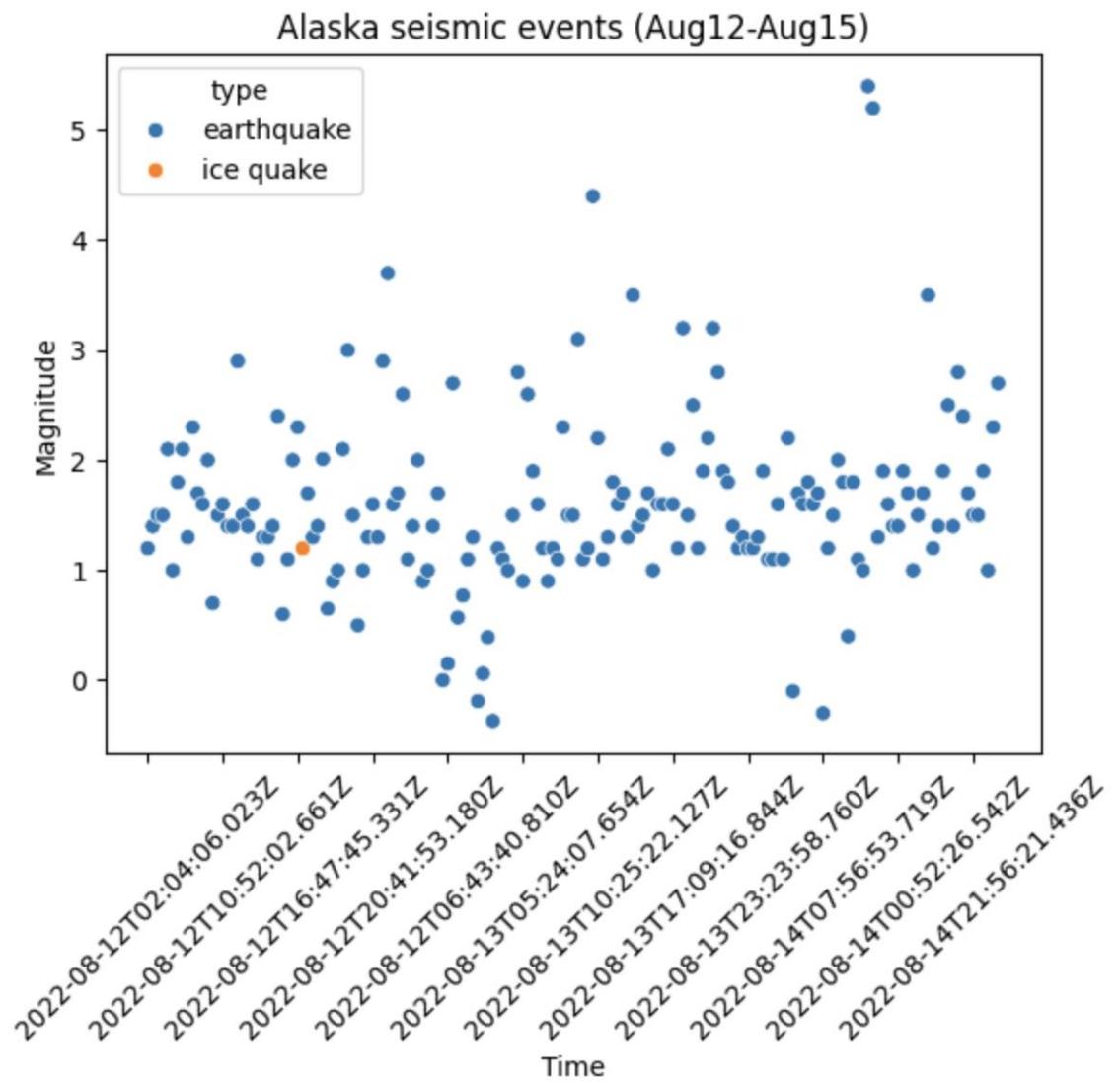


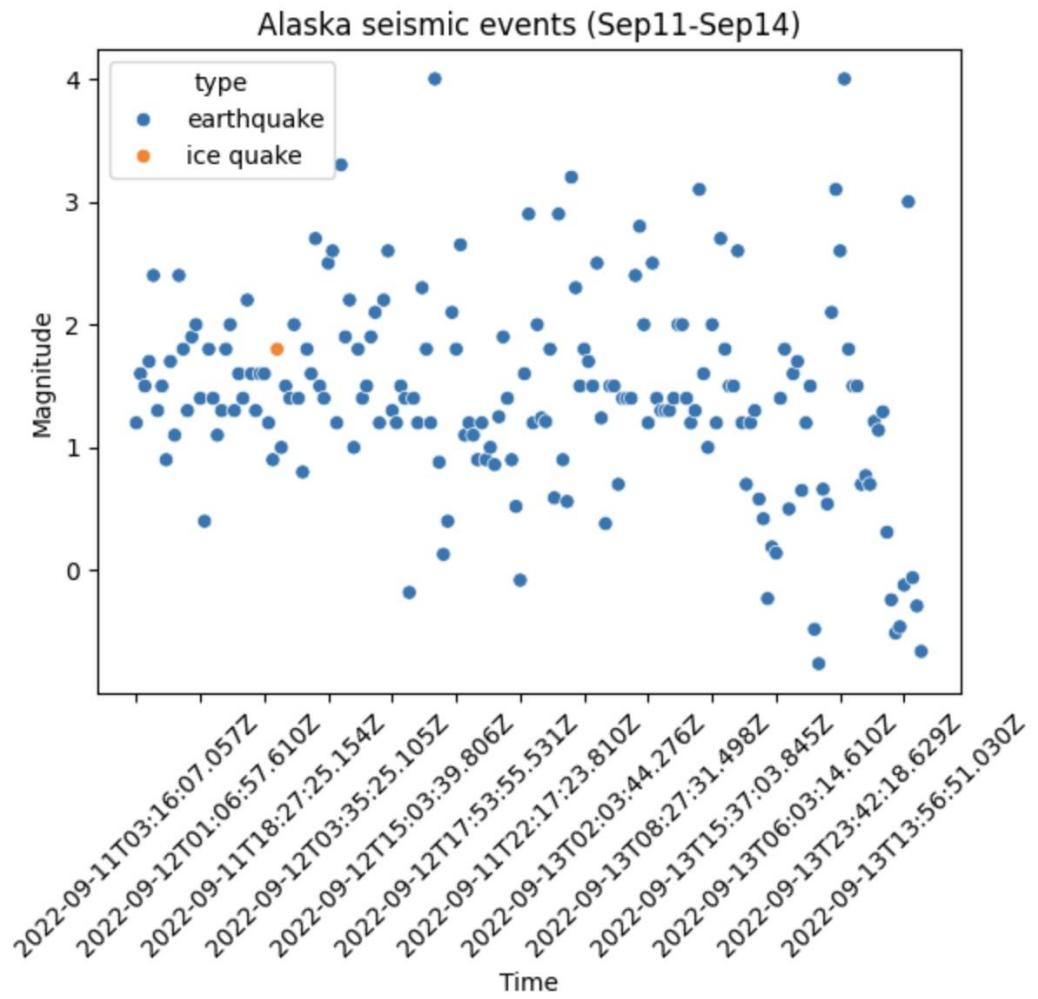


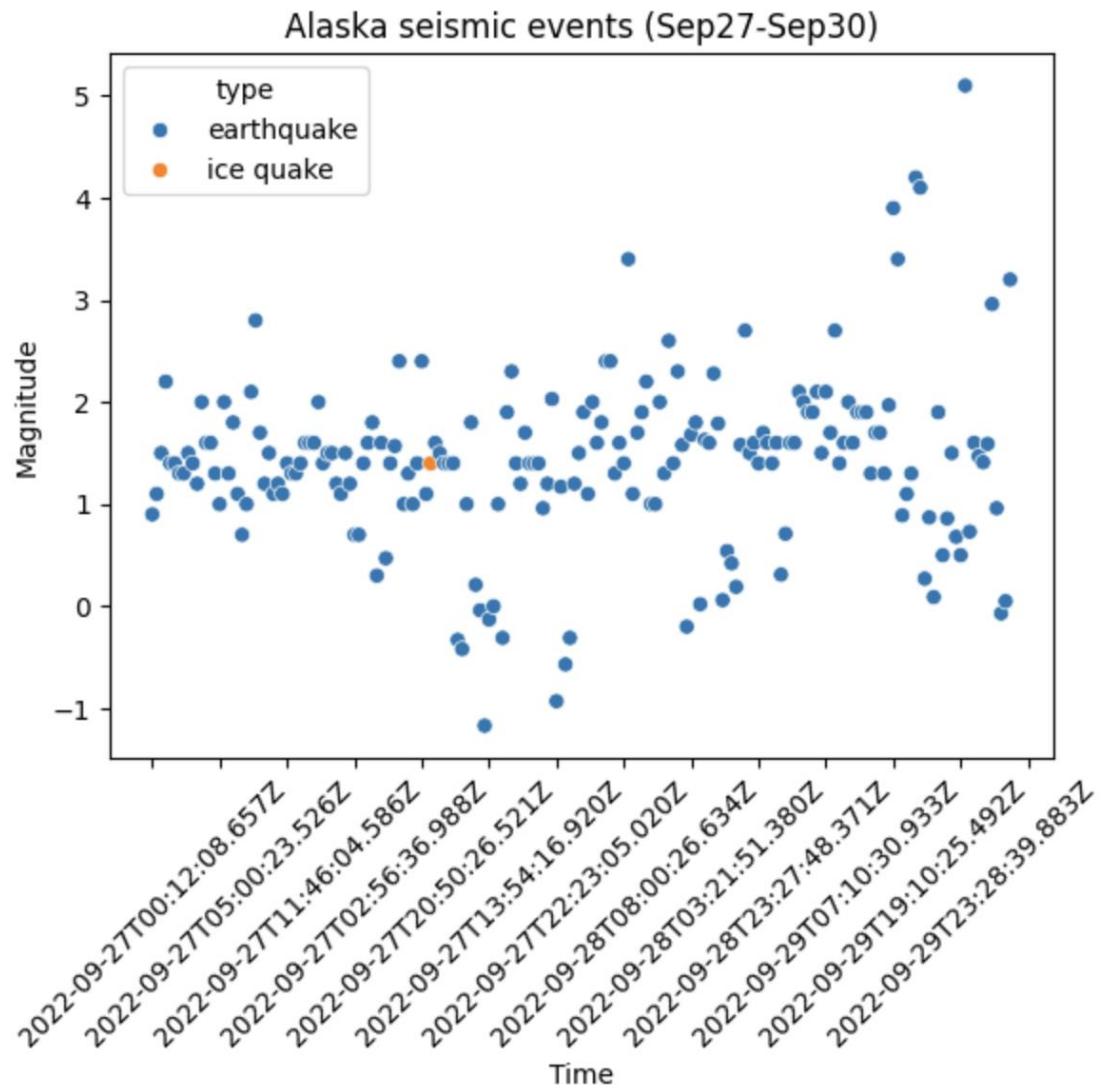


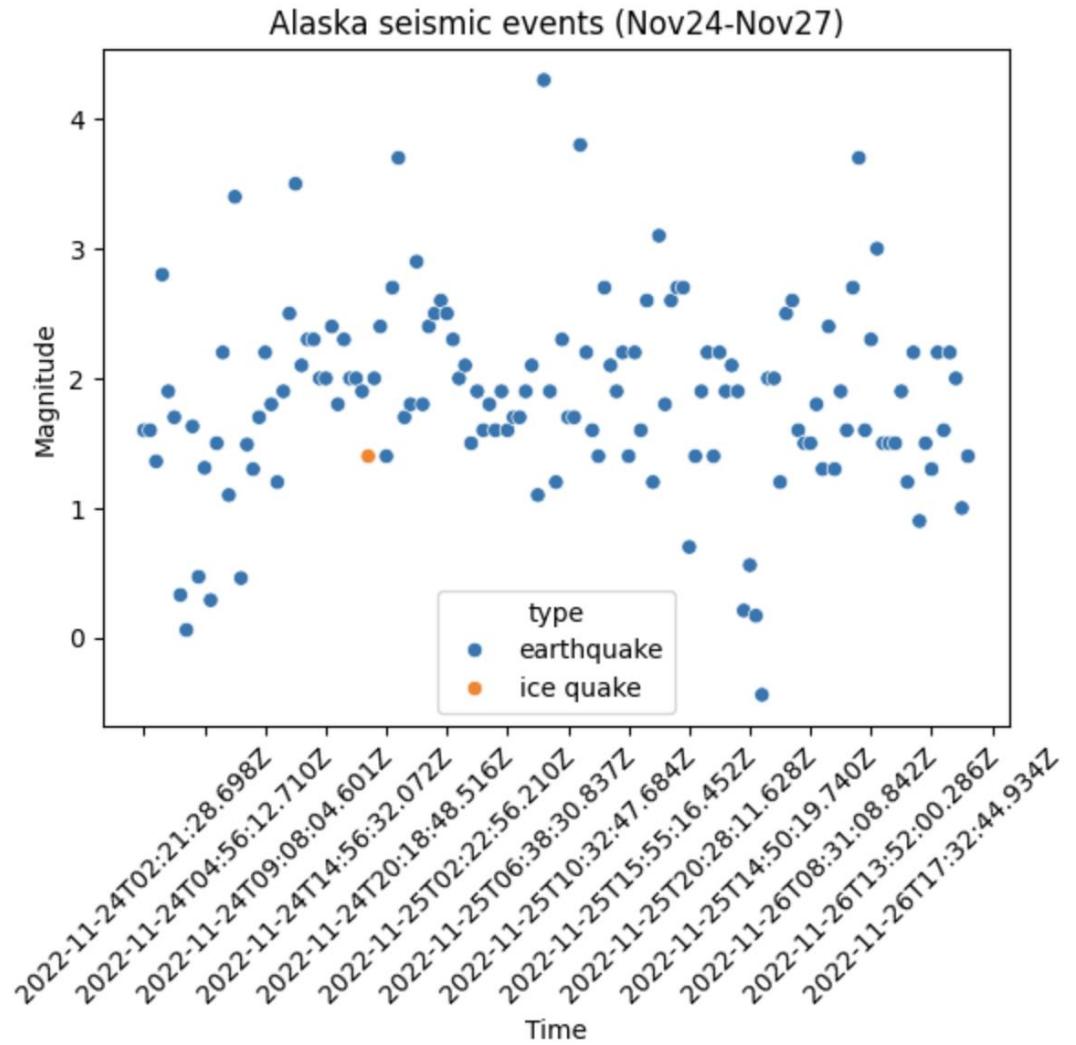


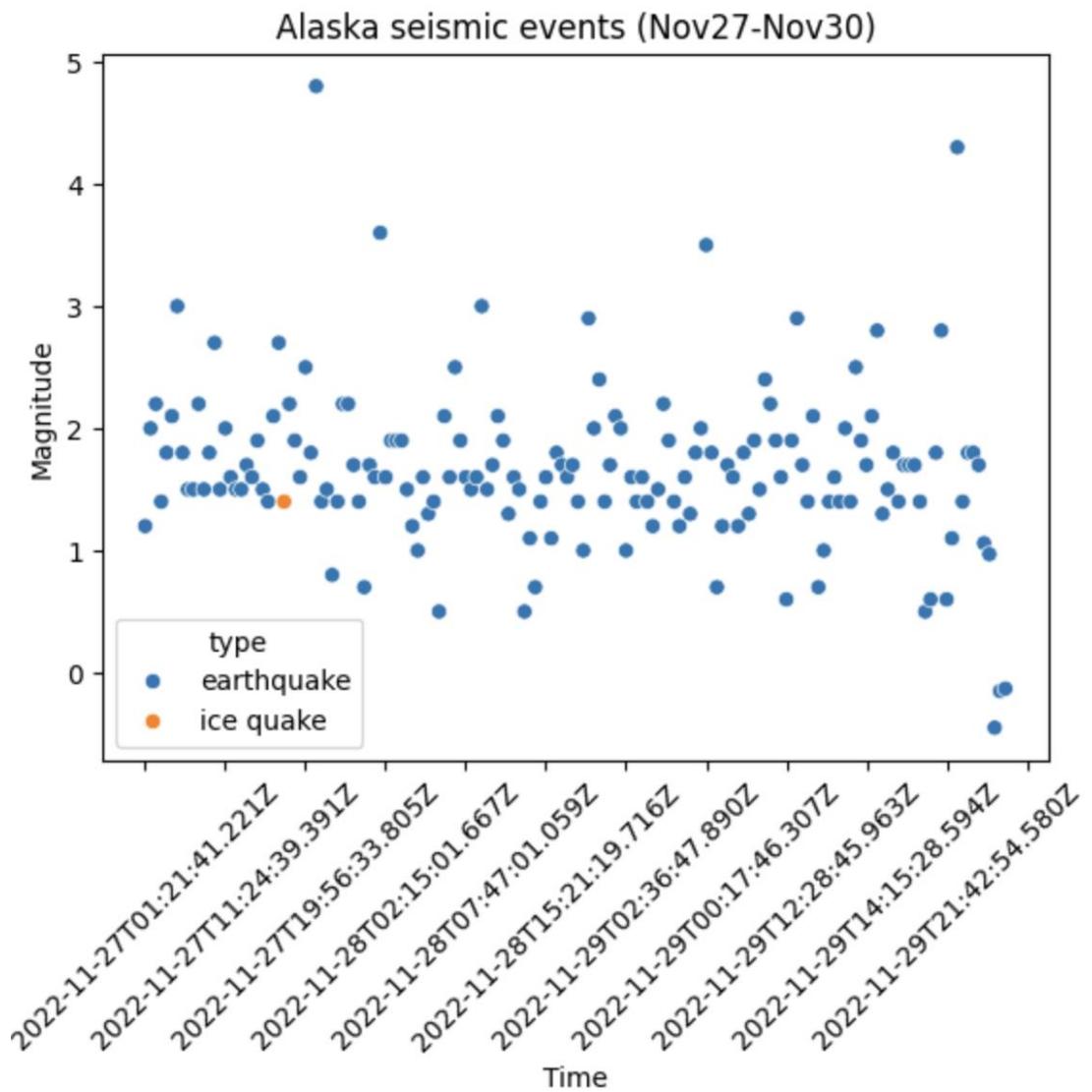












Saved png files into d3 from jupyter hub

Screenshot of the AWS S3 console showing the contents of the 'sainathfinal' bucket.

The browser address bar shows: us-east-2.console.aws.amazon.com/s3/buckets/sainathfinal?region=us-east-2&bucketType=general&tab=objects

The page title is: sainathfinal info

The navigation path is: Amazon S3 > Buckets > sainathfinal

The main tabs are: Objects (selected), Properties, Permissions, Metrics, Management, Access Points.

**Actions:** Copy S3 URI, Copy URL, Download, Open, Delete, Actions (dropdown), Create folder, Upload.

**Objects (13) Info:**

Name	Type	Last modified	Size	Storage class
Alaska earthquakes.png	png	May 6, 2024, 11:32:58 (UTC-05:00)	2.3 KB	Standard
California earthquakes.png	png	May 6, 2024, 11:32:58 (UTC-05:00)	2.3 KB	Standard
California quarry blasts.png	png	May 6, 2024, 11:32:59 (UTC-05:00)	2.3 KB	Standard
Damage-causing earthquakes (2022).png	png	May 6, 2024, 11:32:59 (UTC-05:00)	2.3 KB	Standard
Earthquake count (2022).png	png	May 6, 2024, 11:33:00 (UTC-05:00)	2.3 KB	Standard
earthquakeCountsbylocandtime_.png	png	May 6, 2024, 11:33:01 (UTC-05:00)	2.3 KB	Standard
earthquakeCountsbylocandtime_.csv	csv	May 6, 2024, 02:57:51 (UTC-05:00)	34.6 KB	Standard
usgs_main.csv	csv	May 6, 2024, 02:57:58 (UTC-05:00)	14.2 MB	Standard

CloudShell Feedback © 2024, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

Screenshot of the AWS S3 console showing the contents of the 'sainathfinal' bucket after some objects have been deleted.

The browser address bar shows: us-east-2.console.aws.amazon.com/s3/buckets/sainathfinal?region=us-east-2&bucketType=general&tab=objects

The page title is: sainathfinal info

The navigation path is: Amazon S3 > Buckets > sainathfinal

The main tabs are: Objects (selected), Properties, Permissions, Metrics, Management, Access Points.

**Actions:** Copy S3 URI, Copy URL, Download, Open, Delete, Actions (dropdown), Create folder, Upload.

**Objects (13) Info:**

Name	Type	Last modified	Size	Storage class
earthquakeCountsbylocandtime_.png	png	May 6, 2024, 11:33:01 (UTC-05:00)	2.3 KB	Standard
earthquakefrequencyovertime.png	png	May 6, 2024, 11:33:01 (UTC-05:00)	2.3 KB	Standard
earthquakelocations.png	png	May 6, 2024, 11:33:02 (UTC-05:00)	2.3 KB	Standard
Seismic events - Log scale (2022).png	png	May 6, 2024, 11:33:00 (UTC-05:00)	2.3 KB	Standard
seismic events.png	png	May 6, 2024, 11:33:02 (UTC-05:00)	2.3 KB	Standard
States with moderate earthquakes (2022).png	png	May 6, 2024, 11:33:00 (UTC-05:00)	2.3 KB	Standard
usgs_current.csv	csv	May 6, 2024, 02:57:51 (UTC-05:00)	34.6 KB	Standard
usgs_main.csv	csv	May 6, 2024, 02:57:58 (UTC-05:00)	14.2 MB	Standard

CloudShell Feedback © 2024, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

## Conclusion:

This project successfully demonstrated how to use cloud technologies to handle and analyze massive amounts of seismic data to forecast earthquake features using PySpark and AWS SageMaker. This project not only improved our capacity to anticipate natural disasters but also showed the scalability and real-time data processing advantages of cloud computing. Among the things that have been acquired are the value of properly preparing data and the advantages of incorporating real-time data analysis for quick reaction. Subsequent research efforts could investigate the integration of more diverse data sources and sophisticated machine learning models to strengthen prediction precision and broaden the project's relevance to alternative disaster situations.

## References:

1. Jordan, M. I., & Mitchell, T. M. (2015). Machine Learning: Trends, Perspectives, and Prospects. *Science*, 349(6245), 255-260.
2. Zaharia, M., Chowdhury, M., Franklin, M. J., Shenker, S., & Stoica, I. (2016). Apache Spark: A Unified Engine for Big Data Processing. *Communications of the ACM*, 59(11), 56-65.
3. Liaw, A., & Wiener, M. (2002). Classification and Regression by randomForest. *R News*, 2(3), 18-22. Retrieved from <https://CRAN.R-project.org/doc/Rnews/>
4. Stein, S., & Wysession, M. (2003). An Introduction to Seismology, Earthquakes, and Earth Structure. Blackwell Publishing.