

ML ASSIGNMENT VOLCANO ON DATASET
SE22MAID002
PARAM SHIV ASHISH
AI AND DS

INITIAL DATA THAT WAS GIVEN:

CONTAINS 10 VOLCANO FOLDERS ASSUMING EACH FOLDER TO CORRESPOND WITH A VOLCANO IN REAL LIFE. CERTAIN NUMBER OF OBSERVATIONS IN EACH FOLDER

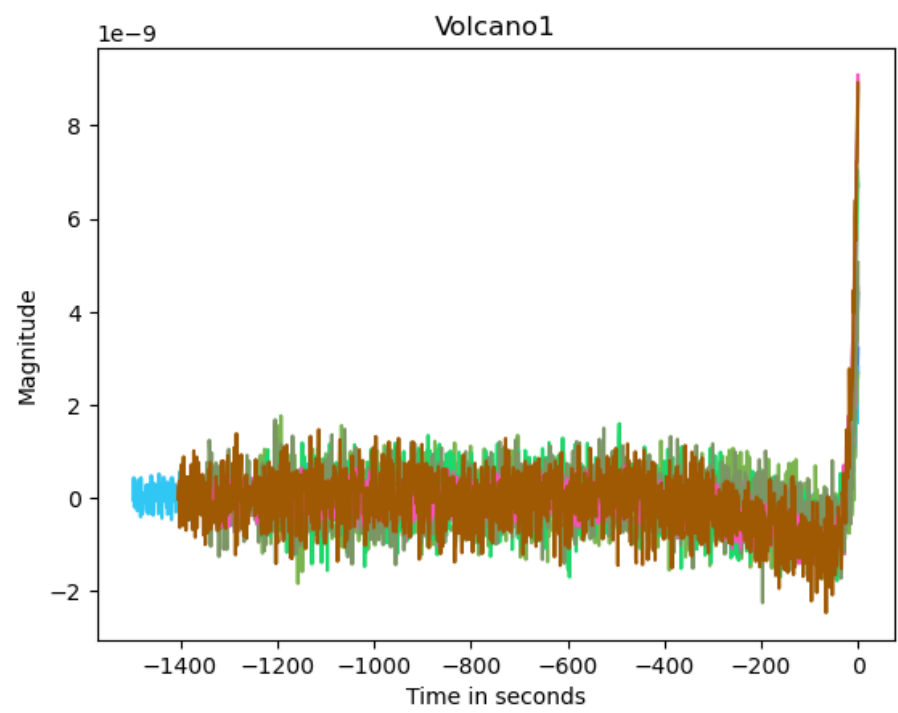
```
Volcano_Dataset/  
├─ Volcano1/  
│   ├── Observation1.txt  
│   ├── Observation2.txt  
│   ├── Observation3.txt  
│   ├── Observation4.txt  
│   └── Observation5.txt  
├─ Volcano2/  
├─ Volcano3/  
├─ Volcano4/  
├─ Volcano5/  
├─ Volcano6/  
├─ Volcano7/  
├─ Volcano8/  
├─ Volcano9/  
└─ Volcano10/
```

AN OBSERVATION TEXT FILE FORMAT INSIDE VOLCANO SUB-FOLDERS

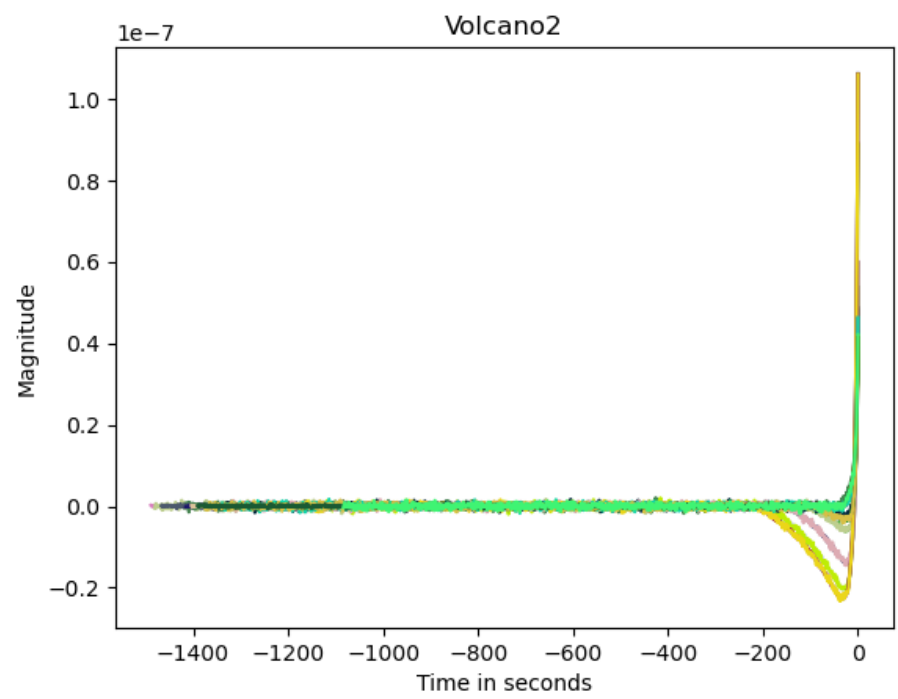
```
observation1 - Notepad  
File Edit View  
v,0.25  
Patm,100000  
g,9.81  
r,500  
g,10^11.12  
rho,2200  
mu,10^3.51  
rc,18  
  
M,10^6.18  
sigma,0.1  
  
tilt_erupt,3.1271nrad  
  
-1052,-1051,-1050,-1049,-1048,-1047,-1046,-1045,-1044,-1043,-1042,-1041,-1040,-1039,-1038,-1037,-1036,-1035,-1034,-1033,-1032,-1031,-1030,-1029,-1028,-1027,-1026,-1025,-1024,-1023,-1022,-1  
63,-862,-861,-860,-859,-858,-857,-856,-855,-854,-853,-852,-851,-850,-849,-848,-847,-846,-845,-844,-843,-842,-841,-840,-839,-838,-837,-836,-835,-834,-833,-832,-831,-830,-829,-828,-827,-826,  
63,-662,-661,-660,-659,-658,-657,-656,-655,-654,-653,-652,-651,-650,-649,-648,-647,-646,-645,-644,-643,-642,-641,-640,-639,-638,-637,-636,-635,-634,-633,-632,-631,-630,-629,-628,-627,-626,  
63,-462,-461,-460,-459,-458,-457,-456,-455,-454,-453,-452,-451,-450,-449,-448,-447,-446,-445,-444,-443,-442,-441,-440,-439,-438,-437,-436,-435,-434,-433,-432,-431,-430,-429,-428,-427,-426,  
63,-262,-261,-260,-259,-258,-257,-256,-255,-254,-253,-252,-251,-250,-249,-248,-247,-246,-245,-244,-243,-242,-241,-240,-239,-238,-237,-236,-235,-234,-233,-232,-231,-230,-229,-228,-227,  
55,-54,-53,-52,-51,-50,-49,-48,-47,-46,-45,-44,-43,-42,-41,-40,-39,-38,-37,-36,-35,-34,-33,-32,-31,-30,-29,-28,-27,-26,-25,-24,-23,-22,-21,-20,-19,-18,-17,-16,-15,-14,-13,-12,-11,-10,-9,-8  
-4.87775135323027e-12,-5.41769735416295e-11,6.28908620278579e-11,8.04923683778754e-11,1.69726718818826e-10,3.35612409306871e-11,1.09193792889304e-10,6.71105189072505e-11,-7.23953342024172e  
981403490877e-11,1.50386285802289e-11,-4.88638956367922e-11,1.00756789074999e-10,-2.59215277860629e-12,1.81193112384341e-11,2.65138095115512e-11,6.94407438681043e-11,1.27911947946209e-10,  
1.81159590808081e-11,5.79703236648474e-11,1.455419432506e-10,1.9027385471466e-10,-7.80755463793587e-11,5.0824458096685e-11,-6.62964695213727e-11,1.41859788888385e-10,-3.49583584602322e-11,  
-8.2255782976458e-11,1.87955134546798e-10,-4.3868245432955e-11,-9.06283973769824e-11,8.0607203040075e-11,-8.35233941666028e-11,1.73489651362758e-10,1.93396226594379e-10,2.20755631244144e-  
-11,-4.53466430629462e-11,3.4443276299607e-11,3.3515780881549e-12,1.8010450850799e-10,9.45821881931583e-11,-7.72554054708933e-11,8.72152415630946e-11,-1.67588086555805e-10,-1.1299349327492  
659169e-11,2.58022982032766e-12,4.76125470429707e-12,1.45512183512001e-10,2.33882021660498e-11,-1.45611336391863e-11,-6.18572296955692e-11,3.78078478142433e-11,2.66997663566903e-11,-1.5895  
791519926e-11,-9.37872611127857e-11,-1.90298121595972e-10,-1.17378812859422e-11,1.90094848386733e-11,7.50042963312224e-11,3.05584684677696e-11,1.00652508669274e-10,3.64946117357078e-11,-6.  
02927609e-12,-2.20304106352417e-10,8.08722144569076e-13,3.13601032182347e-11,-6.42073690227012e-11,1.92700890475349e-10,6.67107682144224e-11,-1.46278549647451e-11,-1.62784549647451e-11,2.81671269414047e-11,-9.3  
1,4.71300362914284e-11,-1.18738708694492e-10,4.89498476798457e-11,6.77696969163277e-12,-6.56262283520482e-12,-1.17311922711348e-10,-2.01623648620347e-10,-2.31782476055435e-10,-9.0014337122  
892806106e-11,-1.47101538421233e-10,1.00609374569994e-10,1.86446008535573e-10,-9.22522590347e-11,1.85276875538499e-10,-9.65532396592069e-11,-7.38742801680741e-11,1.07543148109884e-10,9.813  
515972783e-11,-2.03555883906059e-10,9.70768976732961e-12,3.73717614045361e-11,1.73648517393497e-10,-7.93444734045739e-11,-1.4931616156170e-11,2.9191391770847e-11,8.44614888470958e-11,4.28  
1.41000742754188e-11,3.9975296114847e-11,-6.45126144065970e-11,7.7245372417289e-12,1.47791198003590e-10,2.20526840801279e-10,1.88505834306851e-12,-4.13586005648521e-11,2.20765245415151e-1  
e-11,1.00055897060974e-10,-1.36307098771387e-10,2.75084416827775e-10,-2.83433629384502e-11,-2.0862442697494e-11,-1.6939736639226e-10,-7.82245586726088e-11,8.61244627315034e-12,1.1595733  
49567e-11,-1.470279830261e-10,1.30166550329382e-10,-6.8202172676437e-11,1.32618873946544e-11,-5.54454360899082e-11,-4.77719731855578e-11,-1.7945311218152e-10,-4.89342071114515e-11,1.432714  
437500256861e-10,9.48694383970794e-11,4.93663479603107e-12,-2.99657853260974e-11,-5.8913562776004e-11,6.38864333193813e-11,1.20826853513494e-12,1.577970808400609e-10,-4.55795015921346e-12,5  
005872103e-11,-2.35960959027552e-11,-2.31733785500245e-13,2.29328370920997e-11,9.44913336784882e-11,-5.59803511863652e-11,5.78350353202466e-11,-2.23109140777429e-11,1.68754285030379e-11,4.  
2.89019976144294e-11,-2.48529386346755e-11,1.57988914788439e-10,5.03416663814696e-11,3.59450122485969e-12,-4.87494535520932e-11,2.00290643230339e-12,2.07437712303234e-10,2.71321466746080e  
2.66302742369761e-11,1.97070985068995e-10,-2.6789193978734e-11,-1.89461796685613e-10,6.95432255809419e-11,1.52739037450282e-10,6.98132770680462e-11,2.83230320783311e-11,1.58799768230504e-1  
2,5808766762e-11,-1.77850781488905e-10,-7.008710773750514e-11,-6.43047570436707e-11,-1.34852600681027e-10,-1.43614000118601e-10,5.51400111565663e-11,-1.46587476035206e-10,-1.43076517035170e-1
```

THE FIRST SERIES OF VALUES IN SECONDS STARTING IN THE PAST, INDICATED WITH THE
NEGATIVE SIGN UP TO 0 WHICH CAN BE CONSIDERED AS CURRENT TIME
THE SECONDS SERIES OF VALUES ARE THE CORRESPONDING MAGNITUDES AT THOSE TIME
STAMPS WHICH IS THE HIGHEST AT 0TH SECOND

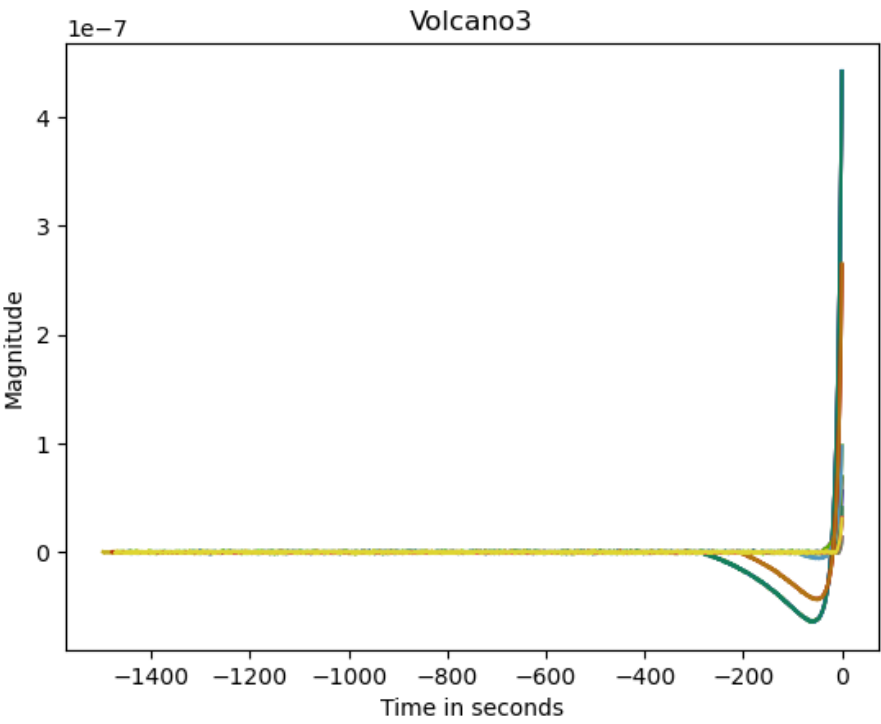
PLOTTING INITIAL DATA THAT WAS GIVEN
VALCANO - 1
OBERVATIONS - 12



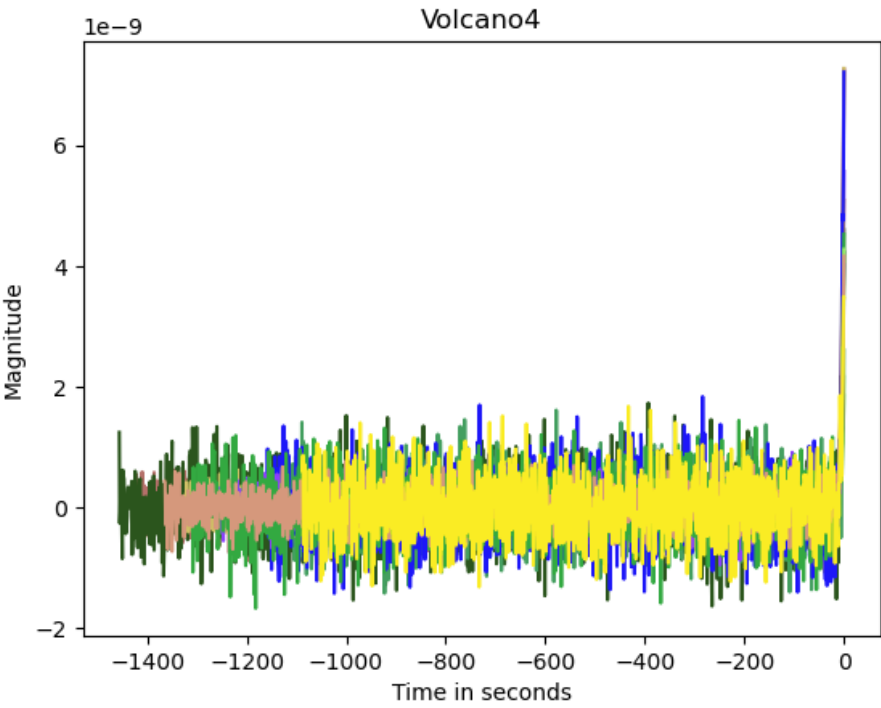
VALCANO - 2
OBERVATIONS - 27



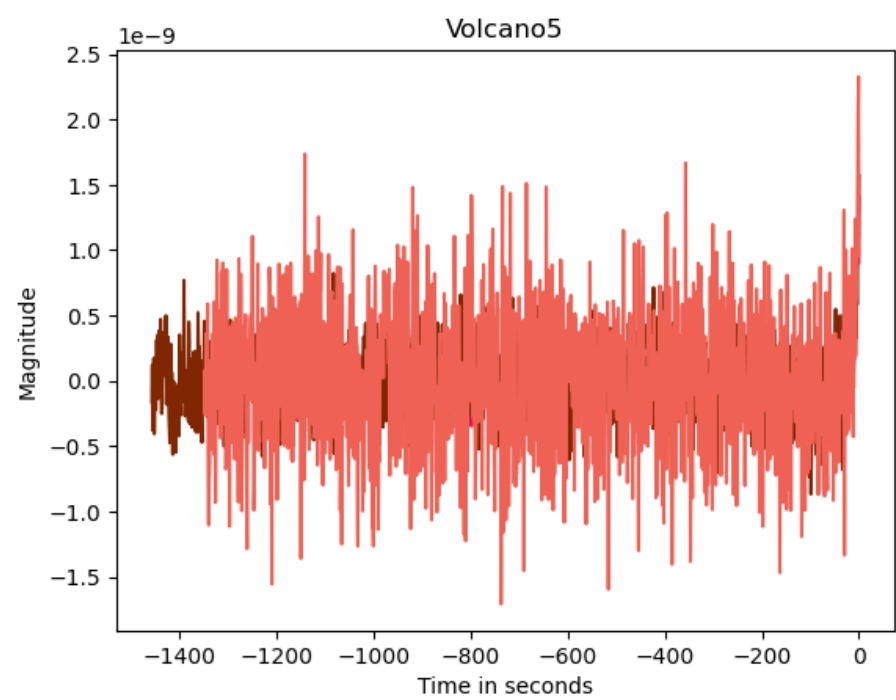
VALCANO - 3
OBERVATIONS - 30



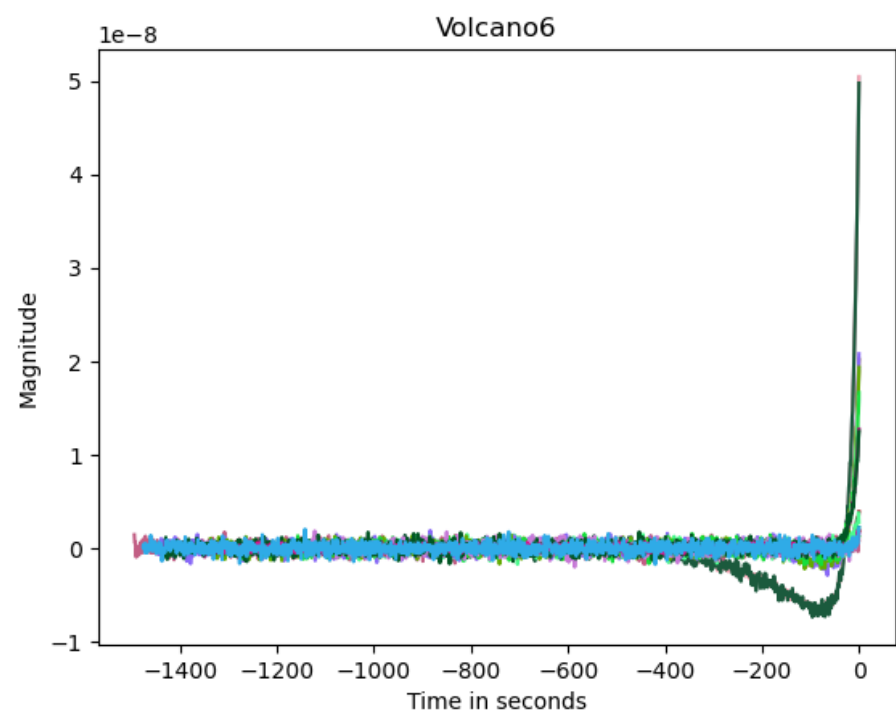
VALCANO - 4
OBERVATIONS - 15



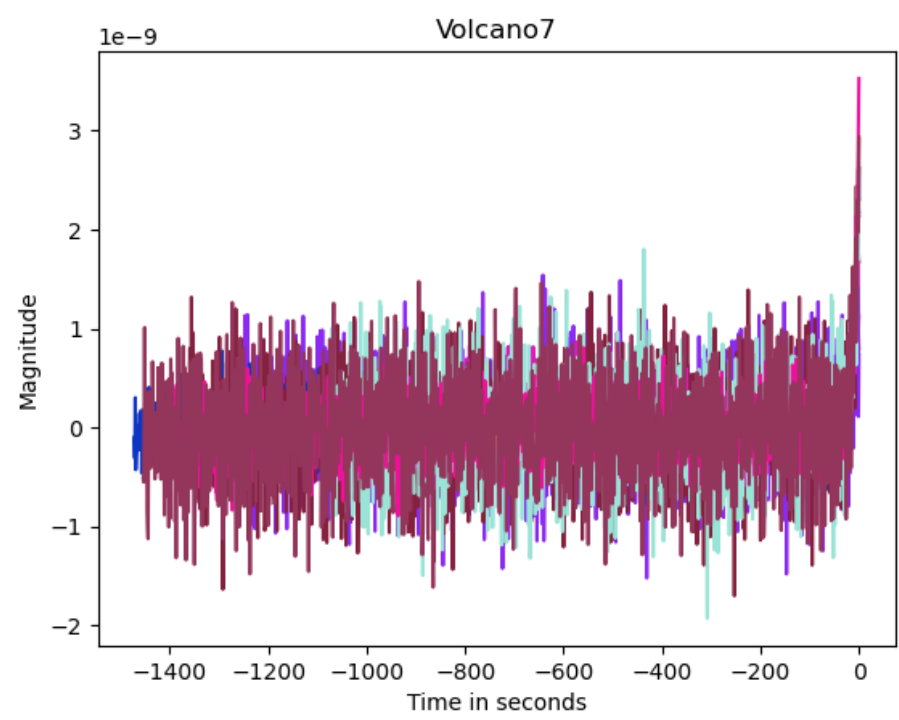
VALCANO - 5
OBERVATIONS - 3



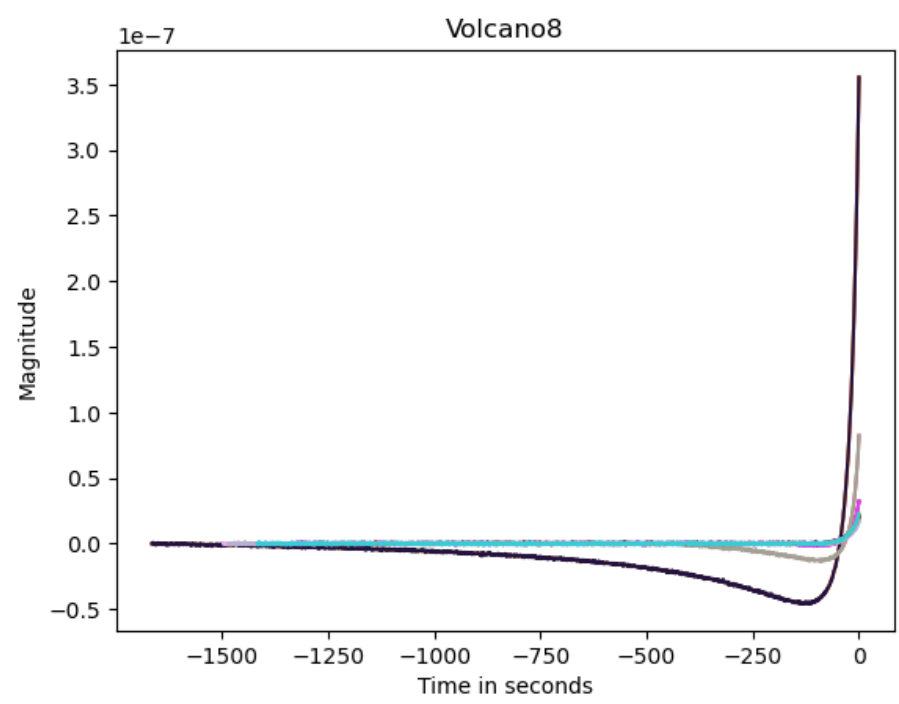
VALCANO - 6
OBERVATIONS - 27



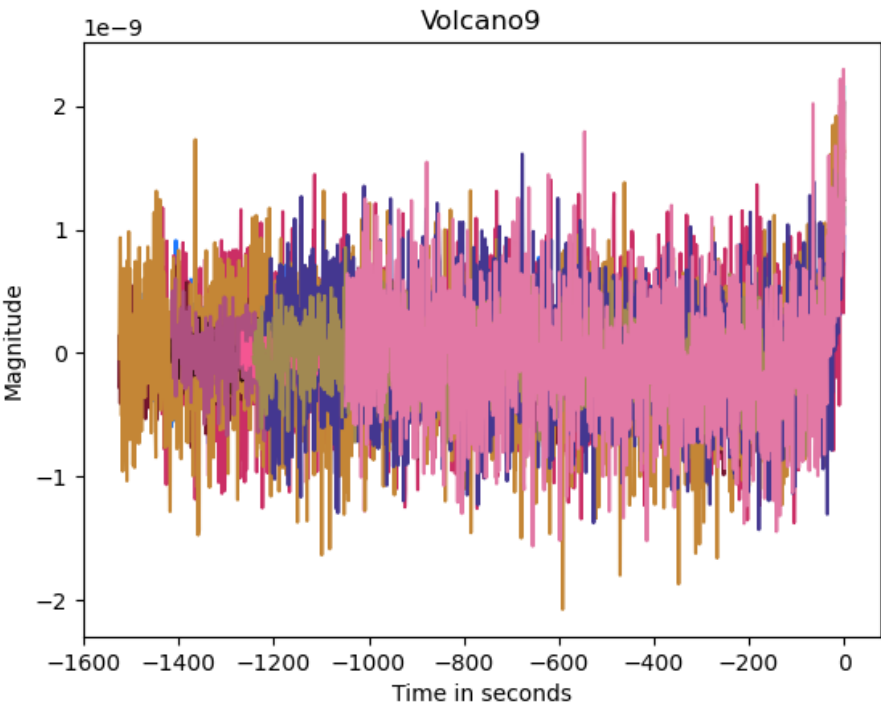
VALCANO - 7
OBERVATIONS - 12



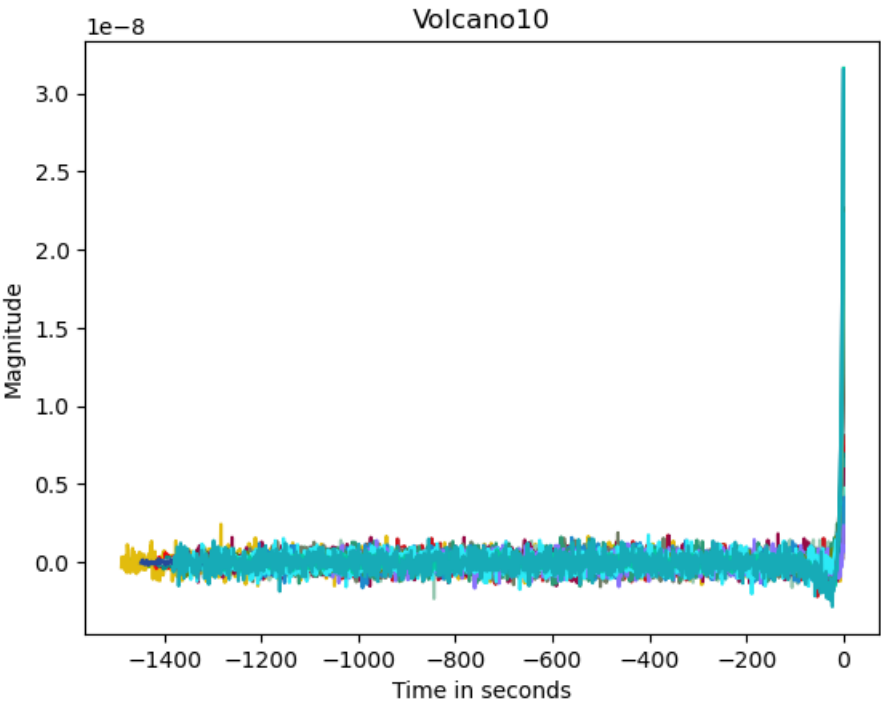
VALCANO - 8
OBERVATIONS - 21



VALCANO - 9
OBERVATIONS - 12



VALCANO - 10
OBERVATIONS - 30



FEATURE ENGINEERING FOR DIRECT FORECASTING

The direct forecasting strategy uses an ML model for each forecasting step. More specifically, *each model is trained using as target the time series shifted of the desired number of time periods into the future*. Now we need to predict the magnitude that sensor reads at the time of eruption which is corresponding to the 0 seconds in our data. In this case, each timestamp in the target time series is chosen as magnitude at 0th second. The feature values are then taken as a consecutive chain of time series magnitudes, here I have considered consecutive chain of length 300. In this way, we create a model trained to predict future value of magnitude. The same procedure is repeated for all forecasting steps.

The direct method outlined above does not generate error accumulation at each forecast, but it has larger computational cost making it not suitable for large forecasting horizons.

Moreover, it cannot model statistical relationships among predictions since the models used for each time step are independent.

IMPLEMENTING A SIMPLE DIRECT FORECASTING FOR VOLCANO DATASET

TIME	MAGNITUDE
-1000	1.23E-20
-999	1.45E-20
-998	1.60E-20
-997	1.87E-20
.	.
.	.
.	.
.	.
.	.
.	.
4	4.12E-20
3	4.22E-20
2	4.45E-20
1	4.56E-20
0	4.79E-20

THE ABOVE TIME SERIES CAN BE TRANSFORMED TO X FEATURES AND Y OUTPUTS AS SHOWN BELOW

X1	X2	X3	.	.	.	X298	X299	X300	Y1	Y2
1.23E-20	1.45E-20	1.60E-20	.	.	.	2.12E-20	2.33E-20	2.60E-20	4.79E-20	700
1.45E-20	1.60E-20	1.87E-20	.	.	.	2.33E-20	2.60E-20	2.95E-20	4.79E-20	701

THE SAMPLES IN THIS STRUCTURE CAN THEN BE USED TO MAKE LINEAR REGRESSION MODELS.

TO ACHIEVE THIS IN PYTHON I HAVE WRITTEN A FUNCTION WHICH EXTRACTS ALL THE SAMPLES FROM THE CONSIDERED OBSERVATIONS.

(ONLY CONSIDERING HALF OF THE OBSERVATIONS AS THE DATASET IS TOO LARGE)

```
import os
import math
def FeatureEngineering(folder,mainDF):
    print(f"VALCANO - {folder}")
    initial_count = 0
    dir = (f"C:\\Users\\crims\\Desktop\\ML Assignment\\Volcano_Dataset\\Volcano{folder}")
    for path in os.listdir(dir):
        if os.path.isfile(os.path.join(dir, path)):
            initial_count += 1
    # Only considering half of the observations in each volcano folder as this is compute intensive and
    # can easily be scaled by considering all the observations of a give volcano
    print(f"Observations Considered - {math.ceil(initial_count/2)}")
    for x in range(1,math.ceil(initial_count/2)):
        # change the path as per your folder path where this project folder is located
        f = open(f"C:\\Users\\crims\\Desktop\\ML Assignment\\Volcano_Dataset\\Volcano{folder}\\observation{x}.txt")
        lines = f.read().splitlines()
        main = [l for l in lines if l != " "]
        tempP = main[-1].split(",")
        p = [float(str(e)) for e in tempP]
        for idx,x in enumerate(p):
            if idx + 300 < len(p):
                instance = [p[i] for i in range(idx,idx+300)]
                instance.append(p[-1])
                instance.append(len(p)-(idx+300))
                mainDF.loc[len(mainDF.index)] = instance
    return math.ceil(initial_count/2)
```

THE ABOVE CODE CONVERTS THE TIME SERIES INTO 300 FEATURES AND 2 TARGETS

THIS RETURNS A DATASET WITH SAMPLES CORRESPONDING TO A SINGLE VOLCANO

SIMPLE FUNCTION FOR METRIC ANALYSIS

```
29 def metric(y_test,y_pred):
30     mse = mean_squared_error(y_test.to_numpy(),y_pred,squared=False,multioutput= 'raw_values')
31     variance_score = explained_variance_score(y_test.to_numpy(), y_pred,multioutput = 'raw_values')
32     r2 = r2_score(y_test.to_numpy(),y_pred, multioutput ="raw_values")
33     mae = mean_absolute_error(y_test.to_numpy(),y_pred, multioutput ="raw_values")
34     return mse,variance_score,r2,mae
```


TRAINING THE LINEAR REGRESSOR MODEL USING PYTHON

TRAINING A MODEL FOR EVERY VOLCANO WITH 2 TARGETS Y1 AND Y2

Q1 outputs and Q2 outputs are both taken in this model

Data set Format

X1 X2 X3 X300 // Y1 Y2

where Xs are the consecutive sensor values forming a chain of length 300 sensor values per sample

Y1 is the sensor value at time = 0

Y2 is the time in seconds which is remaining until the Eruption

```
1 # Q1 outputs and Q2 outputs are both taken in this model
2 # Data set Format
3 # X1 X2 X3 . . . . X300 // Y1 Y2
4 # where Xs are the consecutive sensor values forming a chain of length 300 sensor values per sample
5 # Y1 is the sensor value at time = 0
6 # Y2 is the time in seconds which is remaining until the Eruption
7 from sklearn.model_selection import train_test_split
8 from sklearn.model_selection import GridSearchCV
9 from sklearn.linear_model import LinearRegression
10 from sklearn.model_selection import KFold, cross_val_score
11 from IPython.display import display
12 import pandas as pd
13 modelPerVolcano = []
14 for f in range(1,10+1):
15     print(f"Building Model for Volcano {f}")
16     mainDF = pd.DataFrame(columns=[x for x in range(0,302)])
17 # Actually building our dataset
18 o = FeatureEngineering(f,mainDF)
19 # just shuffling the sample, not really required
20 mainDF = mainDF.sample(frac=1)
21 inputDF = mainDF.iloc[:, :-2]
22 outputDF = mainDF.iloc[:, 300:302]
23 display(mainDF.head(5))
24 x_train, x_test, y_train, y_test = train_test_split(inputDF, outputDF, test_size=.2, random_state = 0)
25 kf = KFold(n_splits = 5)
26 # We can make the model better with different combinations of params, I had left it empty for faster model training
27 params = {}
28 model = LinearRegression()
29 gv = GridSearchCV(model, param_grid=params, scoring = 'neg_mean_absolute_error', cv = kf)
30 gv.fit(x_train, y_train)
31 model = gv.best_estimator_
32 y_pred = model.predict(x_test)
33 # evaluating a few metrics
34 mse, variance_score, r2, mae = metric(y_test, y_pred)
35 modelPerVolcano.append({"volcano": f, 'model': model, 'observationsConsidered' : o, 'neg_mean_absolute' : gv.best_score,
36 # saving models
37 import joblib
38 for i, m in enumerate(modelPerVolcano):
39     model = m["model"]
40     filename = "Q1Q2Models\\volcano"+str(m['volcano'])+".joblib"
41     print(f"Saving Model {i}")
42     joblib.dump(model, filename)
```

```
volcano : 1
model : LinearRegression()
observationsConsidered : 6
neg_mean_absolute : -97.57778057982007
mse : [1.94814466e-09 2.41239771e+02]
variance_score : [-0.03123205  0.30518383]
r2 : [-0.03335851  0.30409217]
mae : [1.75670110e-09 1.93904039e+02]
```

```
volcano : 2
model : LinearRegression()
observationsConsidered : 14
neg_mean_absolute : -122.3613398242079
mse : [1.65646222e-08 2.94893507e+02]
```

```
variance_score : [0.02208243 0.09224182]
r2 : [0.02190775 0.09111543]
mae : [1.30414009e-08 2.44953506e+02]

volcano : 3
model : LinearRegression()
observationsConsidered : 15
neg_mean_absolute : -109.48190136211397
mse : [1.91013411e-07 2.74222673e+02]
variance_score : [0.09854558 0.17476368]
r2 : [0.09832661 0.17456918]
mae : [1.75295522e-07 2.21066750e+02]

volcano : 4
model : LinearRegression()
observationsConsidered : 8
neg_mean_absolute : -126.69547319192353
mse : [1.77968090e-09 2.89567855e+02]
variance_score : [-0.01270171 -0.01879843]
r2 : [-0.01278633 -0.02246905]
mae : [1.56530128e-09 2.44116181e+02]

volcano : 5
model : LinearRegression()
observationsConsidered : 2
neg_mean_absolute : -131.1719472777568
mse : [ 0.          299.01479396]
variance_score : [ 1.          -0.27760749]
r2 : [ 1.          -0.27831909]
mae : [ 0.          245.13714396]

volcano : 6
model : LinearRegression()
observationsConsidered : 14
neg_mean_absolute : -103.06275144527224
mse : [1.56685145e-08 2.54802790e+02]
variance_score : [0.22492813 0.15519286]
r2 : [0.22491849 0.15518569]
mae : [1.26258082e-08 2.01940545e+02]

volcano : 7
model : LinearRegression()
observationsConsidered : 6
neg_mean_absolute : -131.1936882886096
mse : [8.27542064e-10 3.06893522e+02]
variance_score : [-0.11238432 0.04025424]
r2 : [-0.11290911 0.03821124]
mae : [7.81701636e-10 2.59889492e+02]

volcano : 8
model : LinearRegression()
observationsConsidered : 11
neg_mean_absolute : -122.94225051795902
mse : [1.11655615e-07 3.09627219e+02]
variance_score : [0.49734878 0.21298945]
r2 : [0.49714622 0.21287286]
mae : [9.46410480e-08 2.44541169e+02]

volcano : 9
model : LinearRegression()
observationsConsidered : 6
neg_mean_absolute : -124.77685061692526
mse : [3.94753926e-10 2.97057141e+02]
variance_score : [0.33155095 0.21787659]
```

```
r2 : [0.33131857 0.21662227]
mae : [3.32538428e-10 2.39427482e+02]
```

TRAINING A MODEL FOR EVERY VOLCANO WITH 1 TARGETS Y1

```
# Training with Q1 outputs only
# given a time series predict the magnitude at 0 seconds
# Data set Format
# X1 X2 X3 . . . X300 // Y1
# where Xs are the consecutive sensor values forming a chain of
length 300 sensor values per sample
# Y1 is the sensor value at time = 0
```

```
8 from sklearn.model_selection import train_test_split
9 from sklearn.model_selection import GridSearchCV
10 from sklearn.linear_model import LinearRegression
11 from sklearn.model_selection import KFold,cross_val_score
12 from IPython.display import display
13 import pandas as pd
14 modelPerVolcano = []
15 for f in range(1,10+1):
16     print(f"Building Model for Volcano {f}")
17     mainDF = pd.DataFrame(columns=[x for x in range(0,302)])
18     # Actually building our dataset
19     o = FeatureEngineering(f,mainDF)
20     # just shuffling the sample, not really required
21     mainDF = mainDF.sample(frac=1)
22     inputDF = mainDF.iloc[:, :-2]
23     outputDF = mainDF.iloc[:, 300:301]
24     display(mainDF.head())
25     x_train, x_test, y_train, y_test = train_test_split(inputDF, outputDF, test_size=.2, random_state = 0)
26     kf = KFold(n_splits = 5)
27     # We can make the model better with different combinations of params, I had left it empty for faster model training
28     params = {}
29     model = LinearRegression()
30     gv = GridSearchCV(model, param_grid=params, scoring = 'neg_mean_absolute_error', cv = kf)
31     gv.fit(x_train, y_train)
32     model = gv.best_estimator_
33     y_pred = model.predict(x_test)
34     # evaluating a few metrics
35     mse, variance_score, r2, mae = metric(y_test, y_pred)
36     modelPerVolcano.append({'volcano': f, 'model': model, 'observationsConsidered' : o, 'neg_mean_absolute' : gv.best_score})
37 for m in modelPerVolcano:
38     for key, value in m.items():
39         print(f"{key} : {value}")
40     print()
41 #saving models
42 import joblib
43 for i, m in enumerate(modelPerVolcano):
44     model = m["model"]
45     filename = "Q1Models\\volcano"+str(m['volcano'])+"Q"+str(1)+".joblib"
46     print(f"Saving Model {i+1} For Q1")
47     joblib.dump(model, filename)
```

```
volcano : 1
model : LinearRegression()
observationsConsidered : 6
neg_mean_absolute : -1.8028479870665053e-09
mse : [1.94245444e-09]
variance_score : [-0.01232022]
r2 : [-0.01234496]
mae : [1.75011137e-09]
```

```
volcano : 2
model : LinearRegression()
observationsConsidered : 14
neg_mean_absolute : -1.3071602514930183e-08
mse : [1.64608213e-08]
variance_score : [0.05216863]
r2 : [0.05213251]
mae : [1.298671e-08]
```

```
volcano : 3
```

```
model : LinearRegression()  
observationsConsidered : 15  
neg_mean_absolute : -1.7513592680965734e-07  
mse : [1.89849988e-07]  
variance_score : [0.10932453]  
r2 : [0.10918578]  
mae : [1.73319429e-07]
```

```
volcano : 4  
model : LinearRegression()  
observationsConsidered : 8  
neg_mean_absolute : -1.583499832939313e-09  
mse : [1.83162029e-09]  
variance_score : [-0.01034224]  
r2 : [-0.0112746]  
mae : [1.6140701e-09]
```

```
volcano : 5  
model : LinearRegression()  
observationsConsidered : 2  
neg_mean_absolute : -1.2407709188295415e-25  
mse : [0.]  
variance_score : [1.]  
r2 : [1.]  
mae : [0.]
```

```
volcano : 6  
model : LinearRegression()  
observationsConsidered : 14  
neg_mean_absolute : -1.2813019692186025e-08  
mse : [1.55221247e-08]  
variance_score : [0.23431414]  
r2 : [0.23258764]  
mae : [1.26316773e-08]
```

```
volcano : 7  
model : LinearRegression()  
observationsConsidered : 6  
neg_mean_absolute : -7.87277572674364e-10  
mse : [8.38164832e-10]  
variance_score : [-0.09451914]  
r2 : [-0.09703457]  
mae : [7.90873505e-10]
```

```
volcano : 8  
model : LinearRegression()  
observationsConsidered : 11  
neg_mean_absolute : -9.535096813979395e-08  
mse : [1.12595105e-07]  
variance_score : [0.49269383]  
r2 : [0.49269058]  
mae : [9.51329867e-08]
```

```
volcano : 9  
model : LinearRegression()  
observationsConsidered : 6  
neg_mean_absolute : -3.359623036601528e-10  
mse : [4.03828146e-10]  
variance_score : [0.29691795]  
r2 : [0.29691761]  
mae : [3.4255245e-10]
```

TRAINING A MODEL FOR EVERY VOLCANO WITH 1 TARGETS Y2

Training with Q2 outputs only
given a time series predict the time remaining for eruption
Data set Format
X1 X2 X3 . . . X300 // Y2
where Xs are the consecutive sensor values forming a chain of
length 300 sensor values per sample
Y2 is the time in seconds which is remaining until the Eruption

```
8 from sklearn.model_selection import train_test_split
9 from sklearn.model_selection import GridSearchCV
10 from sklearn.linear_model import LinearRegression
11 from sklearn.model_selection import KFold, cross_val_score
12 from IPython.display import display
13 import pandas as pd
14 modelPerVolcano = []
15 for f in range(1,10+1):
16     print(f"Building Model for Volcano {f}")
17     mainDF = pd.DataFrame(columns=[x for x in range(0,302)])
18 # Actually building our dataset
19 o = FeatureEngineering(f,mainDF)
20 # just shuffling the sample, not really required
21 mainDF = mainDF.sample(frac=1)
22 inputDF = mainDF.iloc[:, :-2]
23 outputDF = mainDF.iloc[:, 301:302]
24 display(mainDF.head())
25 x_train, x_test, y_train, y_test = train_test_split(inputDF, outputDF, test_size=.2, random_state = 0)
26 kf = KFold(n_splits = 5)
27 # We can make the model better with different combinations of params, I had left it empty for faster model training
28 params = {}
29 model = LinearRegression()
30 gv = GridSearchCV(model, param_grid=params, scoring = 'neg_mean_absolute_error', cv = kf)
31 gv.fit(x_train, y_train)
32 model = gv.best_estimator_
33 y_pred = model.predict(x_test)
34 # evaluating a few metrics
35 mse, variance_score, r2, mae = metric(y_test, y_pred)
36 modelPerVolcano.append({"volcano": f, 'model': model, 'observationsConsidered' : o, 'neg_mean_absolute' : gv.best_score_})
37 for m in modelPerVolcano:
38     for key, value in m.items():
39         print(f"{key} : {value}")
40     print()
41 #saving models
42 import joblib
43 for i, m in enumerate(modelPerVolcano):
44     model = m["model"]
45     filename = "Q2Models\\volcano"+str(m["volcano"])+str(2)+".joblib"
46     print(f"Saving Model {i+1} For Q2")
47     joblib.dump(model, filename)
```

volcano : 1
model : LinearRegression()
observationsConsidered : 6
neg_mean_absolute : -192.13152926381218
mse : [245.29704525]
variance_score : [0.2775991]
r2 : [0.27684691]
mae : [194.51415443]

volcano : 2
model : LinearRegression()
observationsConsidered : 14
neg_mean_absolute : -243.97487348983728
mse : [296.72994315]
variance_score : [0.11189373]
r2 : [0.11171043]
mae : [245.73568167]

volcano : 3

```
model : LinearRegression()  
observationsConsidered : 15  
neg_mean_absolute : -219.83992863526638  
mse : [272.50377888]  
variance_score : [0.17144118]  
r2 : [0.1713098]  
mae : [218.87593194]
```

```
volcano : 4  
model : LinearRegression()  
observationsConsidered : 8  
neg_mean_absolute : -251.2468354658128  
mse : [300.69253749]  
variance_score : [-0.07836757]  
r2 : [-0.07915243]  
mae : [252.13749916]
```

```
volcano : 5  
model : LinearRegression()  
observationsConsidered : 2  
neg_mean_absolute : -268.5240725210555  
mse : [281.32834511]  
variance_score : [-0.19920282]  
r2 : [-0.2077473]  
mae : [227.66155176]
```

```
volcano : 6  
model : LinearRegression()  
observationsConsidered : 14  
neg_mean_absolute : -206.25666672037778  
mse : [256.13651777]  
variance_score : [0.14015635]  
r2 : [0.13751768]  
mae : [204.51730837]
```

```
volcano : 7  
model : LinearRegression()  
observationsConsidered : 6  
neg_mean_absolute : -262.720835407909  
mse : [305.27324033]  
variance_score : [0.02341456]  
r2 : [0.02100751]  
mae : [257.43734246]
```

```
volcano : 8  
model : LinearRegression()  
observationsConsidered : 11  
neg_mean_absolute : -247.18729343527562  
mse : [302.62960384]  
variance_score : [0.20145598]  
r2 : [0.20134074]  
mae : [239.34060316]
```

```
volcano : 9  
model : LinearRegression()  
observationsConsidered : 6  
neg_mean_absolute : -249.80174768663977  
mse : [297.88335983]  
variance_score : [0.24708372]  
r2 : [0.24681474]  
mae : [244.54250077]
```

LOAD THE SAVED MODELS SO WE DON'T HAVE TO TRAIN THEM AGAIN
WE CAN JUST CHANGE TO THE VALUES INSIDE LOAD() FUNCTION TO LOAD A MODEL FOR A DIFFERENT VOLCANO

LOAD MODEL FOR Q1 AND Q2 PREDICTIONS

```
1 #Load Q1Q2 Models
2 #change Volcano folders to get predictions of different volcano
3 model1 = joblib.load('Q1Q2Models\\volcano'+str(1)+".joblib")
4 mainDF = pd.DataFrame(columns=[x for x in range(0,302)])
5 FeatureEngineering(1,mainDF)
6 mainDF = mainDF.sample(frac=1)
7 inputDF = mainDF.iloc[:, :-2]
8 outputDF = mainDF.iloc[:, 300:302]
9 x_train, x_test, y_train, y_test = train_test_split(inputDF, outputDF, test_size=.2, random_state = 0)
10 y_pred = model1.predict(x_test)
11 for idx, x in enumerate(y_test.to_numpy()):
12     print(f"Q1 {x[0]} : {y_pred[idx][0]} | Q2 {x[1]} : {y_pred[idx][1]}")
13
```

LOAD MODEL FOR Q1 PREDICTIONS

```
1 ### Load Q1 Models
2 #change Volcano folders to get predictions of different volcano
3 model1 = joblib.load('Q1Models\\volcano'+str(1)+"Q"+str(1)+".joblib")
4 mainDF = pd.DataFrame(columns=[x for x in range(0,302)])
5 FeatureEngineering(1,mainDF)
6 mainDF = mainDF.sample(frac=1)
7 inputDF = mainDF.iloc[:, :-2]
8 outputDF = mainDF.iloc[:, 300:302]
9 x_train, x_test, y_train, y_test = train_test_split(inputDF, outputDF, test_size=.2, random_state = 0)
10 y_pred = model1.predict(x_test)
11 for idx, x in enumerate(y_test.to_numpy()):
12     print(f"{x[0]} : {y_pred[idx][0]}")
```

LOAD MODEL FOR Q1 AND Q2 PREDICTIONS

```
1 #Load Q2 Models
2 #change Volcano folders to get predictions of different volcano
3 model1 = joblib.load('Q2Models\\volcano'+str(1)+"Q"+str(2)+".joblib")
4 mainDF = pd.DataFrame(columns=[x for x in range(0,302)])
5 FeatureEngineering(1,mainDF)
6 mainDF = mainDF.sample(frac=1)
7 inputDF = mainDF.iloc[:, :-2]
8 outputDF = mainDF.iloc[:, 301:302]
9 x_train, x_test, y_train, y_test = train_test_split(inputDF, outputDF, test_size=.2, random_state = 0)
10 y_pred = model1.predict(x_test)
11 for idx, x in enumerate(y_test.to_numpy()):
12     print(f"{x[0]} : {y_pred[idx][0]}")
```