

Modeling Notebook One

This notebook contains the logic for performing Linear Regression, Ridge Regression, Random Forest Regression, and Random Forest Classification.

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt

1 from google.colab import drive
2 drive.mount('/content/drive')

Mounted at /content/drive

1 !ls

1 cleaned = pd.read_csv('/content/drive/Sharadriya/Data Science 303 Group Project/csv/cleaned_fires_data_with_four_closest_stations_nov21.csv')
2 cleaned = cleaned.replace([np.inf, -np.inf], np.nan)
3 cleaned = cleaned.dropna()

1 display(cleaned.head())
2 display(cleaned.info(verbose=True))
3
4 STATION_LIST = ["BODIE", "BROOKS", "COHASSET", "EEL_RIVER", "HELL_HOLE", "HERNANDEZ"]
5 STATION_LIST += ["HUNTER_MOUNTAIN", "JUANITA_LAKE", "LADDER_BUTTE", "LAS_TABLAS", "LA_BONDA", "OAK_CREEK"]
6 STATION_LIST += ["FAMMERT", "PILOT_HILL", "SCORPION", "SOLDIER_MOUNTAIN", "SQUAW_LAKE", "STAMPEDE", "VAN_BREMNER", "WOLVERTON"]
7 STATION_LIST = set(STATION_LIST)
```

Unnamed: 0	0	RODIE	BROOKS	CONASSET	CONTAINMENT_MONTH	CONF_DATE	CONF_DOY	CONF_TIME	DISCOVERY_DATE	DISCOVERY_DOY	DISCOVERY_MONTH	DISCOVERY_TIME	EEL_RIVER	FIRE_SIZE	FIRE_SIZE_CLASS	FIRE_YEAR	FOD_ID	HELL_HOLE	HERNANDEZ	HUNTER_MOUNTAIN	JUANITA_LAKE	LADDER_BUTTE	LAS_TARLAS	LA_ROMDA	NWCG_REPORTING_AGENCY	NWCG_REPORTING_UNIT_NAME	OAK_CREEK	OWNER_DESCR	PANAMINT	PILOT_HILL	PRIMARY_STATION_10_MONTHS_PRIOR_EXTREME_MAXIMUM_TEMPERATURE_MONTH	PRIMARY_STATION_10_MONTHS_PRIOR_EXTREME_MINIMUM_TEMPERATURE_FOR_MON	
0	0	-0.501752	-0.841998	-1.128433	2	2453403.5	33.0	1730.0	2453403.5	33	2		1300.0	-0.848324	0.10	A	2005	1	-1.000659	0.195610	0.434939	-0.886032	-1.114178	0.448456	-0.394781	FS	Plumas National Forest	0.297877	USFS	0.545114	-0.914818	-0.339638	-0.500K
1	1	-1.237103	-0.936929	-0.802799	5	2453137.5	133.0	1530.0	2453137.5	133	5		845.0	-0.597768	0.25	A	2004	2	-1.471873	-0.562875	-0.243079	-0.487813	-0.669480	-0.281374	-0.869095	FS	Eldorado National Forest	-0.451833	USFS	-0.057172	-1.280961	1.107121	1.248K
2	2	-1.070104	-1.059637	-0.891098	6	2453156.5	152.0	2024.0	2453156.5	152	6		1921.0	-0.695749	0.10	A	2004	3	-1.410920	-0.543260	-0.109734	-0.533929	-0.714978	-0.262071	-0.950563	FS	Eldorado National Forest	-0.317057	STATE OR PRIVATE	0.059122	-1.392229	0.448904	1.179K
3	3	-1.588084	-0.752356	-0.586044	7	2453189.5	185.0	1400.0	2453184.5	180	6		1600.0	-0.401895	0.10	A	2004	4	-1.213089	-0.755069	-0.584467	-0.322165	-0.486222	-0.477873	-0.836997	FS	Eldorado National Forest	-0.795000	USFS	-0.341413	-1.070171	0.054933	-1.171K
4	4	-1.578881	-0.759880	-0.590604	7	2453189.5	185.0	1200.0	2453184.5	180	6		1600.0	-0.407106	0.10	A	2004	5	-1.217813	-0.758004	-0.557924	-0.324379	-0.488628	-0.480494	-0.845502	FS	Eldorado National Forest	-0.788869	USFS	-0.335760	-1.077584	0.054933	-1.171K

```
5 rows x 255 columns
<class 'pandas.core.frame.DataFrame'>
Int64Index: 87106 entries, 0 to 87105
Data columns (total 255 columns):
#      Column                                Dtype
---  -
0      Unnamed: 0                             int64
1      RODIE                                  float64
2      BROOKS                                float64
3      CONASSET                               float64
4      CONTAINMENT_MONTH                     int64
5      CONF_DATE                             float64
6      CONF_DOY                              float64
7      CONF_TIME                             float64
8      DISCOVERY_DATE                        float64
9      DISCOVERY_DOY                        int64
10     DISCOVERY_MONTH                       int64
11     DISCOVERY_TIME                         float64
12     EEL_RIVER                             float64
13     FIRE_SIZE                             float64
14     FIRE_SIZE_CLASS                       object
15     FIRE_YEAR                             int64
16     FOD_ID                                int64
17     HELL_HOLE                             float64
18     HERNANDEZ                            float64
19     HUNTER_MOUNTAIN                       float64
20     JUANITA_LAKE                          float64
21     LADDER_BUTTE                          float64
22     LAS_TARLAS                            float64
23     LA_ROMDA                              float64
24     NWCG_REPORTING_AGENCY                  object
25     NWCG_REPORTING_UNIT_NAME              object
26     OAK_CREEK                             float64
27     OWNER_DESCR                           object
28     PANAMINT                              float64
29     PILOT_HILL                            float64
30     PRIMARY_STATION_10_MONTHS_PRIOR_EXTREME_MAXIMUM_TEMPERATURE_MONTH float64
31     PRIMARY_STATION_10_MONTHS_PRIOR_EXTREME_MINIMUM_TEMPERATURE_FOR_MONTH float64
32     PRIMARY_STATION_10_MONTHS_PRIOR_HEATING_DEGREE_DAYS_TO_DATE float64
33     PRIMARY_STATION_10_MONTHS_PRIOR_NUM_COOLING_DEGREE_DAY float64
34     PRIMARY_STATION_10_MONTHS_PRIOR_NUM_COOLING_DEGREE_DAY_CUMULATIVE float64
35     PRIMARY_STATION_10_MONTHS_PRIOR_NUM_DAYS_WHERE_AVG_TEMP_BELOW_65_FAHRENHEIT float64
36     PRIMARY_STATION_10_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_70_FAHRENHEIT float64
37     PRIMARY_STATION_10_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_90_FAHRENHEIT float64
38     PRIMARY_STATION_10_MONTHS_PRIOR_NUM_DAYS_WITH_MIN_TEMP_BELOW_0_FAHRENHEIT float64
39     PRIMARY_STATION_10_MONTHS_PRIOR_NUM_DAYS_WITH_MIN_TEMP_BELOW_32_FAHRENHEIT float64
40     PRIMARY_STATION_10_MONTHS_PRIOR_TEMPERATURE_AVERAGE float64
41     PRIMARY_STATION_10_MONTHS_PRIOR_TEMPERATURE_MAX float64
42     PRIMARY_STATION_10_MONTHS_PRIOR_TEMPERATURE_MIN float64
43     PRIMARY_STATION_11_MONTHS_PRIOR_EXTREME_MAXIMUM_TEMPERATURE_MONTH float64
44     PRIMARY_STATION_11_MONTHS_PRIOR_EXTREME_MINIMUM_TEMPERATURE_FOR_MONTH float64
45     PRIMARY_STATION_11_MONTHS_PRIOR_HEATING_DEGREE_DAYS_TO_DATE float64
46     PRIMARY_STATION_11_MONTHS_PRIOR_NUM_COOLING_DEGREE_DAY float64
47     PRIMARY_STATION_11_MONTHS_PRIOR_NUM_COOLING_DEGREE_DAY_CUMULATIVE float64
48     PRIMARY_STATION_11_MONTHS_PRIOR_NUM_DAYS_WHERE_AVG_TEMP_BELOW_65_FAHRENHEIT float64
49     PRIMARY_STATION_11_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_70_FAHRENHEIT float64
50     PRIMARY_STATION_11_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_90_FAHRENHEIT float64
51     PRIMARY_STATION_11_MONTHS_PRIOR_NUM_DAYS_WITH_MIN_TEMP_BELOW_0_FAHRENHEIT float64
52     PRIMARY_STATION_11_MONTHS_PRIOR_NUM_DAYS_WITH_MIN_TEMP_BELOW_32_FAHRENHEIT float64
53     PRIMARY_STATION_11_MONTHS_PRIOR_TEMPERATURE_AVERAGE float64
54     PRIMARY_STATION_11_MONTHS_PRIOR_TEMPERATURE_MAX float64
55     PRIMARY_STATION_11_MONTHS_PRIOR_TEMPERATURE_MIN float64
56     PRIMARY_STATION_12_MONTHS_PRIOR_EXTREME_MAXIMUM_TEMPERATURE_MONTH float64
57     PRIMARY_STATION_12_MONTHS_PRIOR_EXTREME_MINIMUM_TEMPERATURE_FOR_MONTH float64
58     PRIMARY_STATION_12_MONTHS_PRIOR_HEATING_DEGREE_DAYS_TO_DATE float64
59     PRIMARY_STATION_12_MONTHS_PRIOR_NUM_COOLING_DEGREE_DAY float64
60     PRIMARY_STATION_12_MONTHS_PRIOR_NUM_COOLING_DEGREE_DAY_CUMULATIVE float64
61     PRIMARY_STATION_12_MONTHS_PRIOR_NUM_DAYS_WHERE_AVG_TEMP_BELOW_65_FAHRENHEIT float64
62     PRIMARY_STATION_12_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_70_FAHRENHEIT float64
63     PRIMARY_STATION_12_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_90_FAHRENHEIT float64
64     PRIMARY_STATION_12_MONTHS_PRIOR_NUM_DAYS_WITH_MIN_TEMP_BELOW_0_FAHRENHEIT float64
65     PRIMARY_STATION_12_MONTHS_PRIOR_NUM_DAYS_WITH_MIN_TEMP_BELOW_32_FAHRENHEIT float64
66     PRIMARY_STATION_12_MONTHS_PRIOR_TEMPERATURE_AVERAGE float64
67     PRIMARY_STATION_12_MONTHS_PRIOR_TEMPERATURE_MAX float64
68     PRIMARY_STATION_12_MONTHS_PRIOR_TEMPERATURE_MIN float64
69     PRIMARY_STATION_13_MONTHS_PRIOR_EXTREME_MAXIMUM_TEMPERATURE_MONTH float64
70     PRIMARY_STATION_13_MONTHS_PRIOR_EXTREME_MINIMUM_TEMPERATURE_FOR_MONTH float64
71     PRIMARY_STATION_13_MONTHS_PRIOR_HEATING_DEGREE_DAYS_TO_DATE float64
72     PRIMARY_STATION_13_MONTHS_PRIOR_NUM_COOLING_DEGREE_DAY float64
73     PRIMARY_STATION_13_MONTHS_PRIOR_NUM_COOLING_DEGREE_DAY_CUMULATIVE float64
74     PRIMARY_STATION_13_MONTHS_PRIOR_NUM_DAYS_WHERE_AVG_TEMP_BELOW_65_FAHRENHEIT float64
75     PRIMARY_STATION_13_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_70_FAHRENHEIT float64
76     PRIMARY_STATION_13_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_90_FAHRENHEIT float64
77     PRIMARY_STATION_13_MONTHS_PRIOR_NUM_DAYS_WITH_MIN_TEMP_BELOW_0_FAHRENHEIT float64
78     PRIMARY_STATION_13_MONTHS_PRIOR_NUM_DAYS_WITH_MIN_TEMP_BELOW_32_FAHRENHEIT float64
79     PRIMARY_STATION_13_MONTHS_PRIOR_TEMPERATURE_AVERAGE float64
80     PRIMARY_STATION_13_MONTHS_PRIOR_TEMPERATURE_MAX float64
81     PRIMARY_STATION_13_MONTHS_PRIOR_TEMPERATURE_MIN float64
82     PRIMARY_STATION_14_MONTHS_PRIOR_EXTREME_MAXIMUM_TEMPERATURE_MONTH float64
83     PRIMARY_STATION_14_MONTHS_PRIOR_EXTREME_MINIMUM_TEMPERATURE_FOR_MONTH float64
84     PRIMARY_STATION_14_MONTHS_PRIOR_HEATING_DEGREE_DAYS_TO_DATE float64
85     PRIMARY_STATION_14_MONTHS_PRIOR_NUM_COOLING_DEGREE_DAY float64
86     PRIMARY_STATION_14_MONTHS_PRIOR_NUM_COOLING_DEGREE_DAY_CUMULATIVE float64
87     PRIMARY_STATION_14_MONTHS_PRIOR_NUM_DAYS_WHERE_AVG_TEMP_BELOW_65_FAHRENHEIT float64
88     PRIMARY_STATION_14_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_70_FAHRENHEIT float64
89     PRIMARY_STATION_14_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_90_FAHRENHEIT float64
90     PRIMARY_STATION_14_MONTHS_PRIOR_NUM_DAYS_WITH_MIN_TEMP_BELOW_0_FAHRENHEIT float64
91     PRIMARY_STATION_14_MONTHS_PRIOR_NUM_DAYS_WITH_MIN_TEMP_BELOW_32_FAHRENHEIT float64
92     PRIMARY_STATION_14_MONTHS_PRIOR_TEMPERATURE_AVERAGE float64
93     PRIMARY_STATION_14_MONTHS_PRIOR_TEMPERATURE_MAX float64
94     PRIMARY_STATION_14_MONTHS_PRIOR_TEMPERATURE_MIN float64
95     PRIMARY_STATION_15_MONTHS_PRIOR_EXTREME_MAXIMUM_TEMPERATURE_MONTH float64
96     PRIMARY_STATION_15_MONTHS_PRIOR_EXTREME_MINIMUM_TEMPERATURE_FOR_MONTH float64
97     PRIMARY_STATION_15_MONTHS_PRIOR_HEATING_DEGREE_DAYS_TO_DATE float64
98     PRIMARY_STATION_15_MONTHS_PRIOR_NUM_COOLING_DEGREE_DAY float64
99     PRIMARY_STATION_15_MONTHS_PRIOR_NUM_COOLING_DEGREE_DAY_CUMULATIVE float64
100    PRIMARY_STATION_15_MONTHS_PRIOR_NUM_DAYS_WHERE_AVG_TEMP_BELOW_65_FAHRENHEIT float64
101    PRIMARY_STATION_15_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_70_FAHRENHEIT float64
102    PRIMARY_STATION_15_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_90_FAHRENHEIT float64
103    PRIMARY_STATION_15_MONTHS_PRIOR_NUM_DAYS_WITH_MIN_TEMP_BELOW_0_FAHRENHEIT float64
104    PRIMARY_STATION_15_MONTHS_PRIOR_NUM_DAYS_WITH_MIN_TEMP_BELOW_32_FAHRENHEIT float64
105    PRIMARY_STATION_15_MONTHS_PRIOR_TEMPERATURE_AVERAGE float64
106    PRIMARY_STATION_15_MONTHS_PRIOR_TEMPERATURE_MAX float64
107    PRIMARY_STATION_15_MONTHS_PRIOR_TEMPERATURE_MIN float64
108    PRIMARY_STATION_16_MONTHS_PRIOR_EXTREME_MAXIMUM_TEMPERATURE_MONTH float64
109    PRIMARY_STATION_16_MONTHS_PRIOR_EXTREME_MINIMUM_TEMPERATURE_FOR_MONTH float64
110    PRIMARY_STATION_16_MONTHS_PRIOR_HEATING_DEGREE_DAYS_TO_DATE float64
111    PRIMARY_STATION_16_MONTHS_PRIOR_NUM_COOLING_DEGREE_DAY float64
112    PRIMARY_STATION_16_MONTHS_PRIOR_NUM_COOLING_DEGREE_DAY_CUMULATIVE float64
113    PRIMARY_STATION_16_MONTHS_PRIOR_NUM_DAYS_WHERE_AVG_TEMP_BELOW_65_FAHRENHEIT float64
114    PRIMARY_STATION_16_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_70_FAHRENHEIT float64
115    PRIMARY_STATION_16_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_90_FAHRENHEIT float64
116    PRIMARY_STATION_16_MONTHS_PRIOR_NUM_DAYS_WITH_MIN_TEMP_BELOW_0_FAHRENHEIT float64
117    PRIMARY_STATION_16_MONTHS_PRIOR_NUM_DAYS_WITH_MIN_TEMP_BELOW_32_FAHRENHEIT float64
118    PRIMARY_STATION_16_MONTHS_PRIOR_TEMPERATURE_AVERAGE float64
119    PRIMARY_STATION_16_MONTHS_PRIOR_TEMPERATURE_MAX float64
120    PRIMARY_STATION_16_MONTHS_PRIOR_TEMPERATURE_MIN float64
121    PRIMARY_STATION_17_MONTHS_PRIOR_EXTREME_MAXIMUM_TEMPERATURE_MONTH float64
122    PRIMARY_STATION_17_MONTHS_PRIOR_EXTREME_MINIMUM_TEMPERATURE_FOR_MONTH float64
123    PRIMARY_STATION_17_MONTHS_PRIOR_HEATING_DEGREE_DAYS_TO_DATE float64
124    PRIMARY_STATION_17_MONTHS_PRIOR_NUM_COOLING_DEGREE_DAY float64
125    PRIMARY_STATION_17_MONTHS_PRIOR_NUM_COOLING_DEGREE_DAY_CUMULATIVE float64
126    PRIMARY_STATION_17_MONTHS_PRIOR_NUM_DAYS_WHERE_AVG_TEMP_BELOW_65_FAHRENHEIT float64
127    PRIMARY_STATION_17_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_70_FAHRENHEIT float64
128    PRIMARY_STATION_17_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_90_FAHRENHEIT float64
129    PRIMARY_STATION_17_MONTHS_PRIOR_NUM_DAYS_WITH_MIN_TEMP_BELOW_0_FAHRENHEIT float64
```



```
51 PRIMARY_STATION_11_MONTHS_PRIOR_TEMPERATURE_MIN float64
52 PRIMARY_STATION_11_MONTHS_PRIOR_TEMPERATURE_MAX float64

1 cols = numerical.columns
2 closest = []
3
4 for col in cols:
5     if (col[16] == "CLOSEST_STATION_" or col in STATION_LIST or col == "S_LATITUDE" or col == "S_LONGITUDE"):
6         closest.append(col)

1 primary_station_data_cols = []
2 for col in cols:
3     if ("PRIMARY_STATION_" in col or col in STATION_LIST or col == "S_LATITUDE" or col == "S_LONGITUDE"):
4         primary_station_data_cols.append(col)
5
6 primary_station_df = pd.DataFrame()
7 for col in primary_station_data_cols:
8     primary_station_df[col] = numerical[col]
9
10 primary_station_df["FIRE_SIZE"] = numerical["FIRE_SIZE"]
11 primary_station_df["DURATION"] = numerical["DURATION"]
12

1 NUM_FEATURES_TO_KEEP = 300 + 2
2 correlation = primary_station_df.corr(method='pearson')
3 highest_correlation = correlation.nlargest(NUM_FEATURES_TO_KEEP, 'FIRE_SIZE').index
4 print(f"TOP {NUM_FEATURES_TO_KEEP} FEATURES WITH HIGHEST CORRELATION TO FIRE SIZE")
5 print(list(highest_correlation))
6
7 print(f"TOP 10 FEATURES WITH HIGHEST CORRELATION TO FIRE SIZE")
8 display(highest_correlation[2:12])
9 display(correlation.nlargest(NUM_FEATURES_TO_KEEP, 'FIRE_SIZE')[2:12][['FIRE_SIZE']])
10
11 fire_size_prediction_df = primary_station_df[highest_correlation]
12 del fire_size_prediction_df["FIRE_SIZE"]
13 del fire_size_prediction_df["DURATION"]

TOP 302 FEATURES WITH HIGHEST CORRELATION TO FIRE SIZE
['FIRE_SIZE', 'DURATION', 'PRIMARY_STATION_13_MONTHS_PRIOR_EXTREME_MAXIMUM_TEMPERATURE_MONTH', 'PRIMARY_STATION_1_MONTHS_PRIOR_EXTREME_MAXIMUM_TEMPERATURE_MONTH', 'PRIMARY_STATION_13_MONTHS_PRIOR_TEMPERATURE_MAX', 'PRIMARY_STATION_1_MONTHS_PRIOR_TEMPERATURE_MAX', 'PRIMARY_STATION_12_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_70_FAHRENHEIT', 'PRIMARY_STATION_12_MONTHS_PRIOR_NUM_DAYS_WI
TOP 10 FEATURES WITH HIGHEST CORRELATION TO FIRE SIZE
Index(['PRIMARY_STATION_13_MONTHS_PRIOR_EXTREME_MAXIMUM_TEMPERATURE_MONTH',
      'PRIMARY_STATION_1_MONTHS_PRIOR_EXTREME_MAXIMUM_TEMPERATURE_MONTH',
      'PRIMARY_STATION_13_MONTHS_PRIOR_TEMPERATURE_MAX',
      'PRIMARY_STATION_1_MONTHS_PRIOR_TEMPERATURE_MAX',
      'PRIMARY_STATION_12_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_70_FAHRENHEIT',
      'PRIMARY_STATION_12_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_90_FAHRENHEIT',
      'PRIMARY_STATION_13_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_90_FAHRENHEIT',
      'PRIMARY_STATION_1_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_90_FAHRENHEIT',
      'PRIMARY_STATION_12_MONTHS_PRIOR_TEMPERATURE_MAX',
      'PRIMARY_STATION_12_MONTHS_PRIOR_EXTREME_MAXIMUM_TEMPERATURE_MONTH'],
      dtype='object')
PRIMARY_STATION_13_MONTHS_PRIOR_EXTREME_MAXIMUM_TEMPERATURE_MONTH    0.022363
PRIMARY_STATION_1_MONTHS_PRIOR_EXTREME_MAXIMUM_TEMPERATURE_MONTH    0.022363
PRIMARY_STATION_13_MONTHS_PRIOR_TEMPERATURE_MAX                     0.020201
PRIMARY_STATION_1_MONTHS_PRIOR_TEMPERATURE_MAX                     0.020201
PRIMARY_STATION_12_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_70_FAHRENHEIT 0.019643
PRIMARY_STATION_12_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_90_FAHRENHEIT 0.019516
PRIMARY_STATION_13_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_90_FAHRENHEIT 0.019503
PRIMARY_STATION_1_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_90_FAHRENHEIT 0.019503
PRIMARY_STATION_12_MONTHS_PRIOR_TEMPERATURE_MAX                     0.019122
PRIMARY_STATION_12_MONTHS_PRIOR_EXTREME_MAXIMUM_TEMPERATURE_MONTH    0.018295
Name: FIRE_SIZE, dtype: float64

1 display(fire_size_prediction_df)

PRIMARY_STATION_13_MONTHS_PRIOR_EXTREME_MAXIMUM_TEMPERATURE_MONTH  PRIMARY_STATION_1_MONTHS_PRIOR_EXTREME_MAXIMUM_TEMPERATURE_MONTH  PRIMARY_STATION_13_MONTHS_PRIOR_TEMPERATURE_MAX  PRIMARY_STATION_1_MONTHS_PRIOR_TEMPERATURE_MAX  PRIMARY_STATION_12_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_70_FAHRENHEIT  PRIMARY_STATION_12_MONTHS_PRIOR_NUM_DAYS_WITH_MAX_TEMP_ABOVE_90_FAHRENHEIT
0 -0.931493 -0.931493 -1.345868 -1.345868 -1.159138 -0.641117
1 -0.997255 -0.997255 -1.277944 -1.277944 0.053157 -0.143038
2 0.515265 0.515265 -0.092363 -0.092363 0.942172 -0.641117
3 -0.471161 -0.471161 -0.283785 -0.283785 0.780533 -0.641117
4 -0.471161 -0.471161 -0.283785 -0.283785 0.780533 -0.641117
... ...
87101 0.712551 0.712551 0.722723 0.722723 1.265451 2.265798
87102 1.370168 1.370168 1.395787 1.395787 0.942172 1.412221
87103 1.698977 1.698977 1.735407 1.735407 1.348270 2.265798
87104 0.317980 0.317980 -0.129413 -0.129413 -0.916679 -0.641117
87105 1.370168 1.370168 1.395787 1.395787 0.942172 1.412221
87106 rows x 231 columns

1 NUM_FEATURES_TO_KEEP = 300 + 2
2 correlation = primary_station_df.corr(method='pearson')
3 highest_correlation = correlation.nlargest(NUM_FEATURES_TO_KEEP, 'DURATION').index
4 print(f"TOP {NUM_FEATURES_TO_KEEP} FEATURES WITH HIGHEST CORRELATION TO FIRE DURATION")
5 display(highest_correlation)
6
7 print(f"TOP 10 FEATURES WITH HIGHEST CORRELATION TO FIRE DURATION")
8 display(highest_correlation[2:12])
9 display(correlation.nlargest(NUM_FEATURES_TO_KEEP, 'DURATION')[2:12][["DURATION"]])
10
11 fire_duration_prediction_df = primary_station_df[highest_correlation]
12 del fire_duration_prediction_df["FIRE_SIZE"]
13 del fire_duration_prediction_df["DURATION"]
```

```
TOP 302 FEATURES WITH HIGHEST CORRELATION TO FIRE DURATION
Index('DURATION', 'FIRE_SIZE',
      'PRIMARY_STATION_5_MONTHS_PRIOR_NUM_DAYS_WHERE_AVG_TEMP_BELOW_45_FARENHEIT',
      'PRIMARY_STATION_6_MONTHS_PRIOR_NUM_DAYS_WITH_MIN_TEMP_BELOW_32_FARENHEIT',
      'PRIMARY_STATION_7_MONTHS_PRIOR_NUM_DAYS_WHERE_AVG_TEMP_BELOW_45_FARENHEIT',
      'PRIMARY_STATION_ELEVATION',
      'PRIMARY_STATION_5_MONTHS_PRIOR_NUM_DAYS_WHERE_AVG_TEMP_BELOW_32_FARENHEIT',
      'PRIMARY_STATION_5_MONTHS_PRIOR_NUM_DAYS_WHERE_AVG_TEMP_BELOW_45_FARENHEIT',
      'PRIMARY_STATION_7_MONTHS_PRIOR_NUM_DAYS_WITH_MIN_TEMP_BELOW_32_FARENHEIT',
      'PRIMARY_STATION_5_MONTHS_PRIOR_NUM_DAYS_WHERE_AVG_TEMP_BELOW_32_FARENHEIT',
      'PRIMARY_STATION_5_MONTHS_PRIOR_NUM_DAYS_WHERE_AVG_TEMP_BELOW_45_FARENHEIT')
1 display(fire_duration_prediction_df)
```

	PRIMARY_STATION_6_MONTHS_PRIOR_NUM_DAYS_WHERE_AVO_TEMP_BELOW_05_FAHRENHEIT	PRIMARY_STATION_6_MONTHS_PRIOR_NUM_DAYS_WITH_MIN_TEMP_BELOW_32_FAHRENHEIT	PRIMARY_STATION_7_MONTHS_PRIOR_NUM_DAYS_WHERE_AVO_TEMP_BELOW_05_FAHRENHEIT	PRIMARY_STATION_ELEVATION	PRIMARY_STATION_9_MONTHS_PRIOR_NUM_DAYS_WITH_MIN_TEMP_BELOW_32_FAHRENHEIT	PRIMARY_STATION_5_MONTHS_PRIOR_NUM_DAYS_WHERE_AVO_TEMP_BELOW_05_FAHRENHEIT
0	-0.869653	-0.935826	-0.995396	1.192259	0.394730	
1	1.227107	1.291149	-0.737623	0.525056	1.291149	
2	1.506884	1.291149	1.227107	0.525056	1.243994	
3	2.613420	1.676955	2.229905	1.808604	1.740558	
4	2.613420	1.676955	2.229905	1.808604	1.740558	
...	...	...	...	...	...	
87101	0.548096	0.394730	0.133145	-1.008744	0.837331	
87102	-0.036608	-0.028466	0.218021	0.180489	0.394730	
87103	-0.055469	-0.232580	0.425497	-1.008744	-0.232580	
87104	-0.391831	-0.935826	-0.844505	-1.266280	-0.505570	
87105	-0.036608	-0.028466	0.218021	0.180489	0.394730	
87106 rows x 231 columns						

```

1 closest_df = pd.DataFrame()
2
3 for col in closest:
4     closest_df[col] = numerical[col]
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036

```

	BOIDE	BRONKOS	COMASSET	KEEL_RIVER	KEEL_HOLE	HERNANDEZ	HUNTER_MOUNTAIN	JUANITA_LAKE	LADDER_BUTTE	LAS_TABLAS	LA_MONDA	OAK_CREEK	PANAMINT	PILOT_HILL	SCORPION	SOLDIER_MOUNTAIN	SQUAW_LAKE	STAMPEDE	VAR_BREXDER	NOLVERTON	S_LATITUDE	S_LONGITUDE
0	-0.501752	-1.84198	-1.128433	-0.848324	-1.000659	0.195610	0.434939	-0.886032	-1.114718	0.448456	-0.394781	0.297877	0.545114	-0.911488	-0.846454	-1.044163	0.687200	-1.139269	-0.918627	0.410736	0.788242	0.330313
1	-0.123703	-0.936929	-0.802799	-0.597768	-1.471873	-0.562875	-0.478713	-0.669480	-0.281374	-0.689095	-0.451833	-0.175217	-1.280961	-0.473194	-0.612901	-0.382626	-1.287797	-0.515304	-0.311271	0.672172	0.388850	
2	-1.070104	-1.059637	-0.891098	-0.695749	-1.410920	-0.543260	-1.067934	-0.533329	-0.714978	-0.262071	-0.950563	-0.317057	0.059122	-1.392229	-0.537461	-0.663661	-0.382626	-1.233892	-0.557916	-0.204882	0.677589	0.356621
3	-1.588084	-0.752356	-0.586044	-0.401895	-1.213089	-0.755069	-0.564467	-0.322165	-0.486222	-0.477673	-0.354703	-0.795000	-0.331763	-1.070171	-0.296636	-0.431126	-0.342827	-1.125369	-0.346639	-0.618217	0.632919	0.436654
4	-1.576981	-0.759850	-0.590004	-0.407106	-1.217813	-0.758604	-0.557924	-0.324379	-0.488626	-0.480404	-0.845502	-0.788869	-0.335760	-1.077584	-0.299859	-0.433722	-1.146406	-1.126335	-0.348787	-0.614275	0.632919	0.437474

### Model 1: Linear Regression

```

1 from sklearn.model_selection import train_test_split
2 from sklearn.metrics import mean_absolute_error, mean_absolute_percentage_error, median_absolute_error, r2_score
3 from sklearn.metrics import explained_variance_score
4 from sklearn import svm
5 from sklearn.preprocessing import PowerTransformer
6 from sklearn.ensemble import RandomForestRegressor
7
8 # def log_transform(column):
9     """
10     Transforms a feature so that it is scaled logarithmically. Useful for correcting floating point errors
11     """
12
13 # def min_max_transform(column):
14     """
15     Transforms a feature using min-max scaling
16     """
17     pass
18
19 def run_full_linear_regression_with_accuracy(features, response, n_highest = 0):
20     print(f"Response Minimum: {response.min()}")
21     print(f"Response Maximum: {response.max()}")
22
23     model = LinearRegression()
24     # Split into test and training set
25     X_train, X_test, y_train, y_test = train_test_split(features, response, test_size=0.2, random_state=2020)
26     model.fit(X_train, y_train)
27
28     training_predictions = model.predict(X_train)
29     test_predictions = model.predict(X_test)
30     print(f"Linear Regression Coefficients: {model.coef_},\nLinear Regression Intercept: {model.intercept_}\n")
31
32     print("Trainings: mean absolute error: ", mean_absolute_error(y_train, training_predictions))
33     print("Test mean absolute error: ", mean_absolute_error(y_test, test_predictions))
34
35     print("Training average error rate: ", median_absolute_error(y_train, training_predictions))
36     print("Test average error rate: ", median_absolute_error(y_test, test_predictions))
37
38     print("Explained variance in training set is: ", explained_variance_score(y_train, training_predictions))
39     print("Explained variance in test set is: ", explained_variance_score(y_test, test_predictions))
40
41     print("R2 Score training set: ", r2_score(y_train, training_predictions))
42     print("R2 Score testing set: ", r2_score(y_test, test_predictions))
43
44     fig = plt.figure()
45     ax1 = fig.add_subplot(1)
46     ax1.set_xlabel("True Value from Test Set")
47     ax1.set_ylabel("Prediction from Test Set")
48     ax1.set_title("True Value vs Predicted Value for Test Set: Linear Regression")
49     ax1.scatter(test_predictions, y_test)
50     plt.show()
51
52 # Residual plot against predictor
53 fig = plt.figure()
54 ax1 = fig.add_subplot(1)
55 ax1.set_xlabel("Test Set Prediction")
56 ax1.set_ylabel("Residuals from Test Set")
57 ax1.set_title("Residual Graph: Linear Regression")
58 ax1.scatter(test_predictions, y_test - test_predictions)
59 plt.show()
60
61
62
63 def run_full_ridge_regression_with_accuracy(features, response, n_highest = 0):
64     print(f"Response Minimum: {response.min()}")
65     print(f"Response Maximum: {response.max()}")
66
67     model = Ridge()
68     # Split into test and training set
69     X_train, X_test, y_train, y_test = train_test_split(features, response, test_size=0.2, random_state=2020)
70     model.fit(X_train, y_train)
71
72     training_predictions = model.predict(X_train)
73     test_predictions = model.predict(X_test)
74     print(f'Ridge Regression Coefficients: {model.coef_},\nRidge Regression Intercept: {model.intercept_}\n')

```

```
77 print("Training: mean absolute error: ", mean_absolute_error(y_train, training_predictions))
78 print("Test mean absolute error: ", mean_absolute_error(y_test, test_predictions))
79
80 print("Training average error rate: ", median_absolute_error(y_train, training_predictions))
81 print("Test average error rate: ", median_absolute_error(y_test, test_predictions))
82
83 print("Explained variance in training set is: ", explained_variance_score(y_train, training_predictions))
84 print("Explained variance in test set is: ", explained_variance_score(y_test, test_predictions))
85
86 print("R2 Score training set: ", r2_score(y_train, training_predictions))
87 print("R2 Score testing set: ", r2_score(y_test, test_predictions))
88
89 fig = plt.figure()
90 ax1 = fig.add_subplot()
91 ax1.set_xlabel("True Value from Test Set")
92 ax1.set_ylabel("Prediction from Test Set")
93 ax1.set_title("True Value vs Predicted Value for Test Set: Ridge Regression")
94 ax1.scatter(y_test, test_predictions)
95 plt.show()
96
97 # Residual plot against predictor
98 fig = plt.figure()
99 ax1 = fig.add_subplot()
100 ax1.set_xlabel("Test Set Prediction")
101 ax1.set_ylabel('Residuals from Test Set')
102 ax1.set_title('Residual Graph: Ridge Regression')
103 ax1.scatter(test_predictions, y_test - test_predictions)
104 plt.show()
105
106 def run_full_svm_regression_with_accuracy(features, response, n_highest = 0):
107     print(f"Response Minimum: {response.min()}")
108     print(f"Response Maximum: {response.max()}")
109
110     model = svm.SVR()
111     # Split into test and training set
112     X_train, X_test, y_train, y_test = train_test_split(features, response, test_size=0.2, random_state=2020)
113     model.fit(X_train, y_train)
114     training_predictions = model.predict(X_train)
115     test_predictions = model.predict(X_test)
116
117     print("Training: mean absolute error: ", mean_absolute_error(y_train, training_predictions))
118     print("Test mean absolute error: ", mean_absolute_error(y_test, test_predictions))
119
120     print("Training average error rate: ", median_absolute_error(y_train, training_predictions))
121     print("Test average error rate: ", median_absolute_error(y_test, test_predictions))
122
123     print("Explained variance in training set is: ", explained_variance_score(y_train, training_predictions))
124     print("Explained variance in test set is: ", explained_variance_score(y_test, test_predictions))
125
126     fig = plt.figure()
127     ax1 = fig.add_subplot()
128     ax1.set_xlabel("True Value from Test Set")
129     ax1.set_ylabel("Prediction from Test Set")
130     ax1.set_title("True Value vs Predicted Value for Test Set: SVM Regression")
131     ax1.scatter(y_test, test_predictions)
132     plt.show()
133
134 # Residual plot against predictor
135 fig = plt.figure()
136 ax1 = fig.add_subplot()
137 ax1.set_xlabel("Test Set Prediction")
138 ax1.set_ylabel('Residuals from Test Set')
139 ax1.set_title('Residual Graph: SVM Regression')
140 ax1.scatter(test_predictions, y_test - test_predictions)
141 plt.show()
142
143 def run_full_random_forest_regression_with_accuracy(features, response, n_highest = 0):
144     print(f"Response Minimum: {response.min()}")
145     print(f"Response Maximum: {response.max()}")
146
147     model = RandomForestRegressor(n_estimators=100)
148     # Split into test and training set
149     X_train, X_test, y_train, y_test = train_test_split(features, response, test_size=0.2, random_state=2020)
150     model.fit(X_train, y_train)
151     training_predictions = model.predict(X_train)
152     test_predictions = model.predict(X_test)
153
154     print("Training: mean absolute error: ", mean_absolute_error(y_train, training_predictions))
155     print("Test mean absolute error: ", mean_absolute_error(y_test, test_predictions))
156
157     print("Training average error rate: ", median_absolute_error(y_train, training_predictions))
158     print("Test average error rate: ", median_absolute_error(y_test, test_predictions))
159
160     print("Explained variance in training set is: ", explained_variance_score(y_train, training_predictions))
161     print("Explained variance in test set is: ", explained_variance_score(y_test, test_predictions))
162
163     fig = plt.figure()
164     ax1 = fig.add_subplot()
165     ax1.set_xlabel("True Value from Test Set")
166     ax1.set_ylabel("Prediction from Test Set")
167     ax1.set_title("True Value vs Predicted Value for Test Set: Random Forest Regression")
168     ax1.scatter(y_test, test_predictions)
169     plt.show()
170
171 # Residual plot against predictor
172 fig = plt.figure()
173 ax1 = fig.add_subplot()
174 ax1.set_xlabel("Test Set Prediction")
175 ax1.set_ylabel('Residuals from Test Set')
176 ax1.set_title('Residual Graph: Random Forest')
177 ax1.scatter(test_predictions, y_test - test_predictions)
178 plt.show()
```

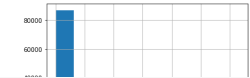
1 SUBSET = -1

## Fire Size Prediction using Linear Regression and Ridge Regression

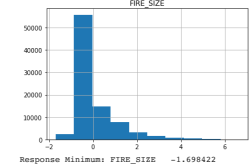
Here, we predict the size of a fire using linear regression

```
1 FIRE_SIZE = numerical["FIRE_SIZE"]
2 FIRE_SIZE.hist()
```

<matplotlib.axes\_subplots.AxesSubplot at 0x7f5d93f86d0>



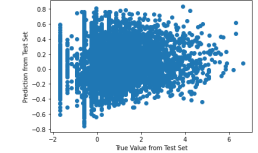
```
1 log_fire_size = pd.DataFrame(np.log(FIRE_SIZE))
2 log_fire_size = log_fire_size.replace([np.inf, -np.inf], np.nan)
3 log_fire_size = log_fire_size.fillna(0)
4 log_fire_size = (log_fire_size - log_fire_size.mean()) / log_fire_size.std()
5 log_fire_size.hist()
6 plt.show()
7
8 # run_full_linear_regression_with_accuracy(fire_size_prediction_df[SUBSET], log_fire_size[SUBSET])
9 run_full_ridge_regression_with_accuracy(fire_size_prediction_df[SUBSET], log_fire_size[SUBSET])
```



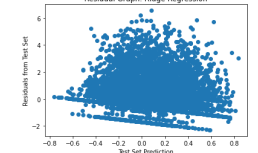
Response Minimum: FIRE\_SIZE -1.698422  
dtype: float64  
Response Maximum: FIRE\_SIZE 6.638909  
dtype: float64  
Ridge Regression Coefficients: [[ 4.33594698e-02 4.33594698e-02 1.03918029e-01 1.03918029e-01  
-3.55139592e-02 1.84372343e-02 -2.37730649e-03 -2.37730649e-03  
2.95197670e-01 -5.50421501e-03 -1.64275460e-02 -1.64275460e-02  
2.09743454e-02 2.09743454e-02 -2.73514574e-03 -2.73514574e-03  
-2.32751510e-01 -2.32751510e-01 5.47128811e-02 1.03293300e-01  
1.03293300e-01 -3.10742784e-02 -9.37498845e-02 -1.59790510e-01  
-8.63986988e-03 -2.02195149e-03 -2.02195149e-03 1.43489531e-01  
1.43489531e-01 4.85216745e-02 -7.13044955e-03 -7.13044955e-03  
-1.53354316e-01 -1.53354316e-01 -1.89397625e-03 -1.89397625e-03  
1.25153780e-01 1.25153780e-01 -3.62207874e-02 -3.62207874e-02  
3.47649977e-02 3.47649977e-02 -2.95336422e-01 -2.95336422e-01  
-7.49515622e-03 -7.49515622e-03 -1.50014523e-02 -1.50014523e-02  
4.71990775e-02 1.57950987e-02 1.57950987e-02 3.95389368e-02  
3.95389368e-02 -2.94773022e-01 2.17586250e-00 -6.20496107e-01  
8.02527227e-03 8.02527227e-03 1.92007303e-02 1.92007303e-02  
-1.21936264e+00 -1.25427278e-01 3.96608944e-02 3.96608944e-02  
2.29678022e-02 2.29678022e-02 8.13479503e-02 8.13479503e-02  
-2.46745638e-02 -2.46745638e-02 1.22420091e-02 -2.83519348e-02  
-1.29021438e+00 -4.49127615e-02 1.43418219e+00 4.66073026e-02  
-1.64123986e+00 6.40358408e-03 6.40358408e-03 1.08638129e+00  
4.82500336e-01 -1.10822859e-02 3.25933961e-01 1.96258308e-01  
3.96177084e-02 -5.15759323e-02 -5.15759323e-02 -4.70122684e-03  
-4.70122683e-03 1.17971539e-02 1.17971539e-02 -1.06617077e-02  
-1.06617077e-02 -2.67885150e-02 1.70111535e-03 5.34686540e-02  
-1.99321352e-02 -1.99321352e-02 9.48966135e-02 1.97705340e-02  
1.08665133e+00 -3.54806340e-02 -3.58292019e-02 -2.88611820e-02  
-6.10815315e-02 2.84367189e-02 2.84367189e-02 -1.93343131e-02  
-1.93343131e-02 2.27118183e-02 -6.85355516e-04 -6.85355516e-04  
1.43410472e-02 -7.64684468e-02 -3.35278327e-02 4.95580216e-02  
3.39636156e+00 1.31892920e-01 -2.77583729e-02 -1.97469144e-02  
-7.52814716e-03 1.16471480e-02 -7.92131710e-02 -3.61492458e-02  
6.83011830e-03 -6.34140872e-02 -2.91079150e-03 -3.86656895e-02  
8.66979561e-03 1.46205124e-01 -5.60352074e-02 9.68871324e-01  
-9.80762352e-02 3.17596654e-02 -4.23048647e-02 1.02074229e-03  
-4.34179526e-02 -7.49110514e-03 -7.49110514e-03 2.83423971e-02  
5.44445929e-02 7.70584044e-02 -5.35758603e-01 2.60590370e-01  
-3.45819863e-02 -6.60734163e-03 -1.77984248e-02 -1.77984248e-02  
6.42804897e-02 2.33412963e-03 -4.30588901e-01 -7.29125279e-01  
-79227273e-02 -7.9227273e-02 4.29888154e-02 4.11278431e-01  
-1.21747881e-01 1.33964358e-04 -4.82635951e-01 -8.23258453e-01  
2.90044843e-03 1.04770930e-02 1.04770930e-02 3.16474061e-01  
-2.91602393e-01 -4.48059033e-04 -4.48059033e-04 -2.11956074e-02  
1.14757576e-02 1.14757576e-02 2.13991992e-01 4.25428791e-01  
9.23604515e-03 1.41532828e-02 1.41532828e-02 -6.99421468e-05  
-7.03892942e-03 2.11925108e-02 5.42289730e-04 1.32431866e-02  
1.32431866e-02 2.52032969e-02 6.70447767e-03 -1.57652797e-02  
3.29509529e-02 -1.74455543e-02 -4.99090504e-04 -5.62255052e-04  
-5.62255052e-04 -9.27821656e-04 -9.27821656e-04 2.52165762e-02  
2.52165762e-02 -6.20699629e-02 4.77016420e-02 -6.73406909e-03  
-4.34379424e-01 4.80174177e-01 2.66107395e-02 2.66107395e-02  
1.06272342e+00 -2.49460239e-01 1.76320789e-02 1.83135442e-02  
2.58548703e-02 -1.71273569e-02 2.77324148e-02 6.47875586e-03  
3.91878129e-02 2.19223807e-02 -7.91176431e-02 -3.99393870e-01  
-5.57587801e-02 -5.57587801e-02 -2.79974557e-02 -3.79393033e-02  
5.09759789e-02 -1.27408458e-04 -1.27408458e-04 5.09244746e-05  
5.09244746e-05 3.14949159e-02 -9.60952526e-02 -9.60952526e-02  
1.17174890e-01 -4.72507300e-02 6.10308156e-04 6.10308156e-04  
-4.65120694e-02 -9.95586537e-02 5.44285772e-03]]\Ridge Regression Intercept: [-1.24172231]

Training: mean absolute error: 0.7025343078227769  
Test mean absolute error: 0.700932062396989  
Training average error rate: 0.540773590564652  
Test average error rate: 0.5394821140280315  
Explained variance in training set is: 0.04953417446802133  
Explained variance in test set is: 0.04486034715927376  
R2 score training set: 0.04953417446802111  
R2 score testing set: 0.04482933788042309

True Value vs Predicted Value for Test Set: Ridge Regression



Residual Graph: Ridge Regression

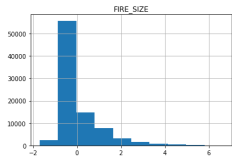


LINEAR REGRESSION and RIDGE REGRESSION for Summer Months

```
1 # LINEAR REGRESSION FOR SUMMER MONTHS
2 log_fire_size = pd.DataFrame(np.log(FIRE_SIZE))
3 log_fire_size = log_fire_size.replace([np.inf, -np.inf], np.nan)
4 log_fire_size = log_fire_size.fillna(0)
5 log_fire_size = (log_fire_size - log_fire_size.mean()) / log_fire_size.std()
6 log_fire_size.hist()
```

```
6 log_fire_size.hist()
7 plt.show()
8
9 #x = 11
10 # x = 6
11 # x = 1
12 summer_fires = ((cleaned['DISCOVERY_MONTH'] >= 5) & (cleaned['DISCOVERY_MONTH'] <= 8))
13 summer_fires_log_size = log_fire_size[summer_fires]
14 summer_fires_df = fire_size_prediction_df[summer_fires]
15 print(len(summer_fires_log_size))
16
17 run_full_linear_regression_with_accuracy(summer_fires_df, summer_fires_log_size)
18 run_full_ridge_regression_with_accuracy(summer_fires_df, summer_fires_log_size)
```





```
56358
Response Minimum: FIRE_SIZE      -1.698422
dtype: float64
Response Maximum: FIRE_SIZE      6.638909
dtype: float64
```

```

Linear Regression Coefficients | 8.00199778e-002 | 8.00199778e-002 | -3.2844170e-01 | -3.2844170e-01
1.41337747e-01 | 1.41337747e-01 | 5.52823668e-02 | 5.52823668e-02
2.41337777e-01 | 2.41337777e-01 | -3.52823668e-02 | -3.52823668e-02
1.31532826e-02 | 1.31532826e-02 | 1.39738464e-02 | 1.39738464e-02
1.54445904e-01 | 1.54445904e-01 | 1.00873978e-02 | 1.00873978e-02
3.81038848e-01 | 3.81038848e-01 | 4.43013866e-03 | 4.43013866e-03
-2.26870866e-01 | -1.00925044e-01 | -1.00925044e-02 | 2.72601478e-01
-2.76404786e-01 | 5.76593256e-02 | 2.68739778e-02 | 2.68739778e-02
-0.62754866e-01 | -0.62754866e-01 | -9.29554666e-03 | -9.29554666e-03
-1.81505739e-01 | -1.81505739e-01 | -5.06344638e-02 | -5.06344638e-02
-2.59006212e-02 | 2.45006212e-02 | -3.16595666e-02 | -3.16595666e-02
-6.29006146e-02 | -6.29006146e-02 | -2.28369312e-02 | -2.28369312e-02
7.68524820e-02 | -3.18647736e-02 | -3.18647736e-02 | 5.67124896e-02
-2.54181896e-02 | 2.59283006e-02 | 3.12549666e-02 | 3.12549666e-02
-2.54576646e-02 | -2.54576646e-02 | 2.53790085e-02 | 2.53790085e-02
-1.66451761e-00 | -1.37235684e-01 | 2.36031586e-01 | 2.36031586e-01
-3.84460506e-02 | -3.84460506e-02 | 9.63327666e-02 | 9.63327666e-02
-6.34594506e-02 | -6.34594506e-01 | 1.03222484e-01 | 1.35825334e-01
-1.1247115e-00 | -5.51723516e-02 | 3.00681080e-02 | -6.72800444e-02
-2.48771516e-02 | -2.62674416e-02 | -2.62674416e-02 | -2.62674416e-02
-1.47333336e-01 | 1.03189586e-02 | 8.45828310e-02 | 6.79477596e-01
-1.27333336e-01 | 1.37333336e-01 | 1.61333336e-01 | 1.61333336e-01
-2.72558116e-02 | 9.8274248e-02 | 9.8274248e-02 | 4.94310456e-01
-3.81034596e-01 | 9.80773254e-02 | 4.33896322e-02 | 5.31314336e-02
-3.82366836e-02 | -3.26226836e-02 | 9.43327666e-02 | 9.43327666e-02
-1.12136566e-02 | -1.01223346e-02 | -9.66773591e-02 | 1.71567424e-01
-1.34124366e-02 | -1.59059380e-01 | -1.59059380e-02 | 2.01640517e-01
-3.26150517e-02 | 8.52943606e-02 | -6.53259666e-02 | -6.53259666e-02
3.22863026e-02 | 4.39602032e-01 | 6.79531216e-02 | 1.90947523e-02
3.83528126e-02 | 3.96192681e-02 | 4.97551585e-02 | 1.40784166e-01
-2.76151886e-02 | 8.52751539e-04 | -1.08925044e-02 | -0.01232516e-01
-9.67176577e-02 | 3.20908686e-02 | 8.12967428e-02 | 2.68813210e-01
4.71847746e-02 | 1.97801167e-01 | -2.55622399e-01 | 1.93576476e-01
-2.70985766e-02 | 4.86031366e-02 | -5.48829666e-02 | -2.58499666e-02
-1.95058316e-02 | -1.22176516e-02 | -1.21765756e-02 | 8.31952402e-02
5.92685677e-02 | 9.25738919e-02 | -5.60675678e-02 | 8.28368096e-02
-5.92685677e-02 | 9.25738919e-02 | -5.60675678e-02 | 8.28368096e-02
-9.58014046e-02 | 4.05485425e-02 | -9.11472646e-01 | 1.39425207e-02
-8.55240696e-02 | 1.88402736e-02 | -9.51316599e-02 | 1.19757189e+00
-1.92034319e-02 | 2.78518406e-03 | -8.30139778e-02 | -8.30139778e-02
6.92812087e-02 | 4.72251516e-02 | 4.72251516e-02 | 5.91649896e-01
6.73891394e-02 | 1.99015596e-02 | 1.99015596e-02 | -2.07267444e-02
-1.27004866e-02 | 4.57378086e-02 | 1.10575056e-02 | -9.92628386e-02
-4.82238886e-02 | 4.40189639e-02 | 2.40189639e-02 | -1.37934464e-02
5.99381806e-02 | 7.36813666e-02 | -4.66710926e-02 | -4.66710926e-02
-3.34384866e-02 | -2.50728846e-02 | 7.93279216e-02 | 1.26793802e-02
-2.54964884e-02 | -1.04667414e-02 | -2.27994121e-02 | 7.04337040e-02
7.64637056e-02 | -7.87998125e-02 | -7.79981250e-02 | 4.40282816e-02
-4.70284726e-02 | -2.70284726e-02 | 4.52226216e-02 | 4.40282816e-02
-2.34305746e-02 | 5.88052546e-02 | 8.13661380e-02 | 8.13661380e-02
-1.02218010e-02 | -1.02298586e-02 | 5.93251576e-02 | 5.93251576e-02
-6.40087036e-02 | -2.28402472e-02 | 3.52315157e-02 | -2.78219205e-02
1.28643356e-02 | 4.79626846e-02 | -8.52599844e-02 | 5.15494773e-01
-1.02218006e-02 | 4.74183066e-02 | 6.88139778e-02 | 6.88139778e-02
3.86930044e-02 | 4.74778196e-02 | -6.74778196e-02 | -1.24272235e-02
-6.74772256e-02 | 8.24631158e-02 | 9.69245284e-02 | -6.99245284e-02
-9.69245284e-02 | 8.24631158e-02 | 9.69245284e-02 | -9.69245284e-02
-6.74772256e-02 | -1.12939610e-01 | 1.93471477e-01 | 1.93471477e-01

```

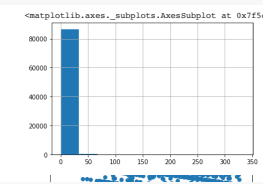
```
Training: mean absolute error: 0.7235662080929671
Test mean absolute error: 0.7327221148900808
Training average error rate: 0.5581052030804253
Test average error rate: 0.5597196652795539
Explained variance in training set is: 0.058073205766394076
Explained variance in test set is: 0.05112624256973641
R2 Score training set: 0.05807320576639419
R2 Score testing set: 0.0510658338813299
True Value vs Predicted Value for Test Set: Linear Regression
```

## Fire Duration Prediction using Linear Regression and Ridge Regression

Here, we predict the duration of a fire using linear regression

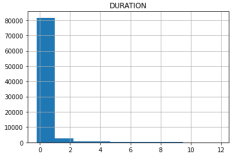
[illegible]

```
1 FIRE_DURATION = numerical["DURATION"]
2 FIRE_DURATION.hist()
```



```
1 log_fire_duration = pd.DataFrame(np.log(FIRE_DURATION))
2 log_fire_duration = log_fire_duration.replace(np.inf, -np.inf, np.nan)
3 log_fire_duration = log_fire_duration.fillna(0)
4 log_fire_duration = (log_fire_duration - log_fire_duration.mean()) / log_fire_duration.std()
5 log_fire_duration.hist()
6 plt.show()
7
8 # run full linear regression with accuracy(fire_duration_prediction_df[SUBSET], log_fire_duration[SUBSET])
9
10 run_full_ridge_regression_with_accuracy(fire_duration_prediction_df[SUBSET], log_fire_duration[SUBSET])
```

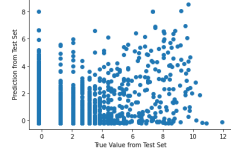
```
/usr/local/lib/python3.7/dist-packages/pandas/core/series.py:726: RuntimeWarning: divide by zero encountered in log
result = getattr(ufunc, method)(*inputs, **kwargs)
```



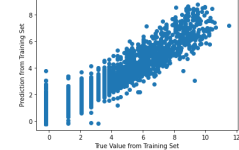
```
Response Minimum: DURATION      -0.214565
dtype: float64
Response Maximum: DURATION      11.92416
dtype: float64
Ridge Regression Coefficients: [[ 1.02358088e-01  1.17706857e-02  2.77658530e-02 -4.89525767e-02
  1.95109615e-02  2.60871570e-02  1.17886907e-03 -3.39702365e-02
 -3.39702365e-02  1.14883335e-02  1.14883335e-02 -1.79400296e-03
 -1.79400295e-03  1.37626317e-03  1.37626317e-03  1.01626380e-01
  0.01626380e-01  1.85498678e-02  1.54935616e-01 -4.37212136e-02
 -3.45442451e-02 -3.45442451e-02 -5.29869195e-02 -5.92115148e-02
 -5.92115149e-02 -7.51207610e-02  2.46944372e-02  2.46944372e-02
 2.76259866e-02  2.95448724e-03 -6.11598490e-02  1.49373314e-02
 1.33683946e-03 -1.20237337e-02  5.05777989e-01 -1.00976037e-01
 2.66687688e-02  2.66687688e-02 -1.06019517e-02 -1.06019517e-02
 -1.72863067e-02 -1.91313449e-02 -2.77578655e-02 -2.77578655e-02
 -2.06476701e-03 -2.06476701e-03 -5.46713631e-02  1.03590442e-02
 1.77977628e-02  2.13905891e-02 -2.90507249e-02  1.11060190e-02
 2.38643166e-02  1.08445871e-02  4.04342440e-02 -3.09167089e-03
 -3.09167089e-03  7.05168606e-03  7.03837774e-03  7.03837774e-03
 6.99616520e-03  6.99616520e-03 -2.45759277e-01  3.09045845e-02]
```

```
1 run_full_random_forest_regression_with_accuracy(fire_duration_prediction_df[1:SUBSET], log_fire_duration[1:SUBSET])
```

```
Response Minimum: DURATION      -0.214565
dtype: float64
Response Maximum: DURATION      11.92416
dtype: float64
/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:154: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().
Training mean absolute error: 0.1177982426916068
Test mean absolute error: 0.32392140126236035
Training average error rate: 0.00953956228133569
Test average error rate: 0.04141452036099269
Explained variance in training set is: 0.8867181072640216
Explained variance in test set is: 0.22640951091631878
True Value vs Predicted Value for Test Set: Random Forest Regression
```



True Value vs Predicted Value for Training Set: Random Forest Regression



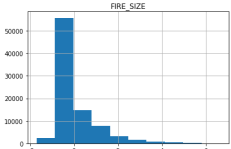
```
-----
ValueError                                Traceback (most recent call last)
<ipython-input-30-f164410162a> in <module>()
----> 1 run_full_random_forest_regression_with_accuracy(fire_duration_prediction_df[1:SUBSET], log_fire_duration[1:SUBSET])
```

```
-----
3 frames
/usr/local/lib/python3.7/dist-packages/pandas/core/ops/_init_.py in to_series(right)
464     if len(left.columns) != len(right):
465         raise ValueError(
--> 466             msg.format(req_len=len(left.columns), given_len=len(right))
467         )
468     right = left._constructor sliced(right, index=left.columns)
```

```
ValueError: Unable to coerce to Series, length must be 1: given 17421
```

## Random Forest Regression FIRE\_SIZE

```
1 |
2 |
3 |
1 log_fire_size = pd.DataFrame(np.log(FIRE_SIZE))
2 log_fire_size = log_fire_size.replace([np.inf, -np.inf], np.nan)
3 log_fire_size = log_fire_size.fillna(0)
4 log_fire_size = (log_fire_size - log_fire_size.mean()) / log_fire_size.std()
5 log_fire_size.hist()
6 plt.show()
7
8 # run_full_gvm_regression_with_accuracy(fire_size_prediction_df[1:SUBSET], log_fire_size[1:SUBSET])
9 run_full_random_forest_regression_with_accuracy(fire_size_prediction_df[1:1], log_fire_size[1:1])
10
```

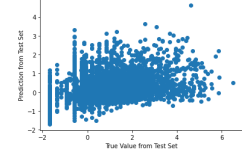


```
Response Minimum: FIRE_SIZE      -1.698422
dtype: float64
Response Maximum: FIRE_SIZE       6.638909
dtype: float64
/usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:154: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples,), for example using ravel().
```

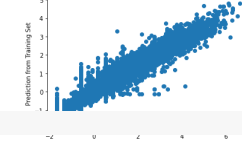
Double-click (or enter) to edit

```
1 # Random forests for summer months
2 run_full_random_forest_regression_with_accuracy(summer_fires_df, summer_fires_log_size)
```

Response Minimum: FIRE\_SIZE -1.698422  
dtype: float64  
Response Maximum: FIRE\_SIZE 6.638909  
dtype: float64  
/usr/local/lib/python3.7/dist-packages/ipykernel\_launcher.py:154: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n\_samples,), for example using ravel().  
Training mean absolute error: 0.2544606226266703  
Test mean absolute error: 0.4681863760662583  
Training average error rate: 0.1714524020813536  
Test average error rate: 0.4603202742169611  
Explained variance in training set is: 0.4693374232241285  
Explained variance in test set is: 0.14721311827405048  
True Value vs Predicted Value for Test Set: Random Forest Regression



True Value vs Predicted Value for Training Set: Random Forest Regression



1

Model 2: Random Forest Classifier

```
1 # Model 2:
2 stat = numerical['STAT_CAUSE_CODE']
3 stat.value_counts()
4
5 cleaned_stat1 = numerical[numerical. STAT_CAUSE_CODE != 13.0]
6 cleaned_stat = cleaned_stat1[cleaned_stat1. STAT_CAUSE_CODE != 9.0]
7
8 col_index = cleaned_stat.columns.get_loc('STAT_CAUSE_CODE')
9 #print(col_index)
10 cleaned_stat.head()
11
12 # use either a random forest or neural network to predict the
```

Unnamed: 0	BODIE	BROOKS	CONASSET	CONTAINMENT_MONTH	CONF_DATE	CONF_DOY	CONF_TIME	DISCOVERY_DATE	DISCOVERY_DOY	DISCOVERY_MONTH	DISCOVERY_TIME	EEL_RIVER	FIRE_SIZE	FIRE_YEAR	FOD_ID	HELL_HOLE	HERRANDEZ	HUNTER_MOUNTAIN	JUANITA_LAKE	LADDER_RUTTE	LAS_TARLAS	LA_MONDA	OAK_CREEK	PANAMINT	PILOT_HILL	PRIMARY_STATION_10_MONTHS_PRIOR_EXTREME_MAXIMUM_TEMPERATURE_MONTH	PRIMA
1	1	-1.237103	-0.939629	-0.802799	5	2453137.5	133.0	1530.0	2453137.5	133	5	845.0	-0.597768	0.25	2004	2	-1.471873	-0.562875	-0.243079	-0.487813	-0.669480	-0.281374	-0.869095	-0.451833	-0.057172	-1.280961	1.107121
2	2	-1.070104	-1.059637	-0.891098	6	2453156.5	152.0	2024.0	2453156.5	152	6	1921.0	-0.695749	0.10	2004	3	-1.410920	-0.543260	-0.109734	-0.533929	-0.714976	-0.262071	-0.950563	-0.317057	0.059122	-1.392229	0.449504
3	3	-1.588084	-0.752356	-0.586044	7	2453189.5	185.0	1400.0	2453184.5	180	6	1600.0	-0.401895	0.10	2004	4	-1.213089	-0.755069	-0.564467	-0.322165	-0.486222	-0.477673	-0.838997	-0.795000	-0.341413	-1.070171	0.054933
4	4	-1.578681	-0.759860	-0.590604	7	2453189.5	185.0	1200.0	2453184.5	180	6	1600.0	-0.407106	0.10	2004	5	-1.217813	-0.758604	-0.557924	-0.324379	-0.488626	-0.480494	-0.845502	-0.788869	-0.335760	-1.077584	0.054933
5	5	-1.475062	-0.828420	-0.650928	7	2453187.5	183.0	1600.0	2453186.5	182	7	1800.0	-0.466098	0.10	2004	6	-1.287362	-0.734812	-0.468782	-0.365956	-0.534467	-0.452476	-0.887998	-0.696533	-0.257419	-1.153409	-0.142352

5 rows x 246 columns

```
1 from sklearn.ensemble import RandomForestClassifier
2 # Instantiate model with 135 decision trees (figured out in the next few blocks of code that this is optimal)
3 forest = RandomForestClassifier(n_estimators = 135, oob_score = True)
```

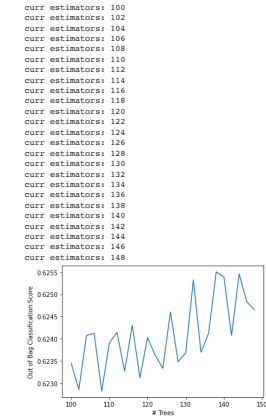
```
1 x_train_mod2 = np.array(cleaned_stat.iloc[:, col_index+1:])
2 y_mod2 = cleaned_stat.iloc[:, col_index]
3 y_train_mod2 = np.array(y_mod2)
4
5 forest.fit(x_train_mod2, y_train_mod2.ravel())
6 print(forest.score(x_train_mod2, y_train_mod2))
7 print(forest.oob_score_)
```

1.0  
0.6247818786422151

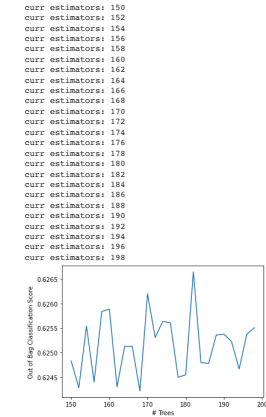
```
1 oob_list = []
2
3 for i in range(20, 100, 2):
4     print('curr estimators:', i)
5     forest = RandomForestClassifier(n_estimators = i, oob_score = True)
6     forest.fit(x_train_mod2, y_train_mod2.ravel())
7     oob_list.append(forest.oob_score_)
8 plt.plot(list(range(20, 100, 2)), oob_list)
9 plt.xlabel('# Trees')
10 plt.ylabel('Out of Bag Classification Score')
11 plt.show()
```

```
curr_estimators: 20
/usr/local/lib/python3.7/dist-packages/sklearn/ensemble/_forest.py:554: UserWarning: Some inputs do not have OOB scores. This probably means too few trees were used to compute any reliable OOB estimates.
  UserWarning,
curr_estimators: 22
/usr/local/lib/python3.7/dist-packages/sklearn/ensemble/_forest.py:554: UserWarning: Some inputs do not have OOB scores. This probably means too few trees were used to compute any reliable OOB estimates.
  UserWarning,
curr_estimators: 24
/usr/local/lib/python3.7/dist-packages/sklearn/ensemble/_forest.py:554: UserWarning: Some inputs do not have OOB scores. This probably means too few trees were used to compute any reliable OOB estimates.
  UserWarning,
curr_estimators: 26
/usr/local/lib/python3.7/dist-packages/sklearn/ensemble/_forest.py:554: UserWarning: Some inputs do not have OOB scores. This probably means too few trees were used to compute any reliable OOB estimates.
  UserWarning,
curr_estimators: 28
curr_estimators: 30
curr_estimators: 32
curr_estimators: 34
curr_estimators: 36
-
```

```
1 oob_list2 = []
2
3 for i in range(100, 150, 2):
4     print("curr_estimators:", i)
5     forest = RandomForestClassifier(n_estimators = i, oob_score = True)
6     forest.fit(x_train_mod2, y_train_mod2.ravel())
7     oob_list2.append(forest.oob_score_)
8 plt.plot(list(range(100, 150, 2)), oob_list2)
9 plt.xlabel('# Trees')
10 plt.ylabel('Out of Bag Classification Score')
11 plt.show()
```



```
1 oob_list3 = []
2
3 for i in range(150, 200, 2):
4     print("curr_estimators:", i)
5     forest = RandomForestClassifier(n_estimators = i, oob_score = True)
6     forest.fit(x_train_mod2, y_train_mod2.ravel())
7     oob_list3.append(forest.oob_score_)
8 plt.plot(list(range(150, 200, 2)), oob_list3)
9 plt.xlabel('# Trees')
10 plt.ylabel('Out of Bag Classification Score')
11 plt.show()
```



```
1 forest = RandomForestClassifier(n_estimators = 182, oob_score = True)
2 forest.fit(x_train_mod2, y_train_mod2.ravel())
3 val_list = forest.feature_importances_
4 idx_list = np.argsort(val_list)[::-1]
5
6 print('From high to low:')
7 for idx in idx_list:
8     print('Feature %d: %f' % (idx, val_list[idx]))

From high to low:
Feature 0: 0.215777
Feature 4: 0.187571
Feature 3: 0.174209
Feature 1: 0.166072
Feature 2: 0.162683
Feature 5: 0.050670
Feature 6: 0.043017
```

```
1 x_train_mod3 = np.array(cleaned_stat.iloc[:, col_index+1:])
```

```
2 y_mod3 = cleaned_stat.iloc[:, [col_index]]
3 y_train_mod3 = np.array(y_mod3)
4
5 forest.fit(x_train_mod3, y_train_mod3.ravel())
6 print(forest.score(x_train_mod3, y_train_mod3))
7 print(forest.oob_score_)
```

```
1.0
0.6264445395581603
```

```
1 forest = RandomForestClassifier(n_estimators = 182, oob_score = True)
2
3 x_train_mod2 = np.array(cleaned_stat.iloc[:, col_index+1:])
4 y_mod2 = cleaned_stat.iloc[:, [col_index]]
5 y_train_mod2 = np.array(y_mod2)
6
7 X_train, X_test, y_train, y_test = train_test_split(cleaned_stat, y_train_mod2, test_size=0.2, random_state=2020)
8
9 forest.fit(X_train, y_train.ravel())
10 print(forest.score(X_test, y_test))
11 print(forest.oob_score_)
```

```
12
13
/usr/local/lib/python3.7/dist-packages/sklearn/base.py:446: UserWarning: X does not have valid feature names, but RandomForestClassifier was fitted with feature names
  "X does not have valid feature names, but"
0.8585185185185186
0.8519837023623343
```