

# Business Scenario

When deciding upon resourcing plans it is important to consider likely fault volumes. Since the orders that fuel the provision work are typically taken 2 weeks in advance, a forward looking forecast of fault volumes would enable a more accurate prediction of the number of provision volumes to take on each week.

You have been provided with the fault volume and rainfall data for a single geographic area for 3 years. Please use this data to create a model to predict the expected daily fault intake based on weather events. Detailed explanation of the fields provided is given below.

At the interview you will be given 10 minutes to present your findings. The presentation should give an overview of your findings and be targeted at a senior stakeholder level. Please provide a copy of your findings for review before the interview.

Things you may want to consider in the model;

- 1 Both short term and long term impact of rainfall
- 2 Impact of weather events on different fault types
- 3 Other things that might impact fault volume

Data provided:

Fault Data: This contains information of individual faults for a single region over a 4 year period.

Fault ID: Fault identifier

Report Date: The date that the fault was reported

Initial MFL: Location of the fault;

CA: Customer appointed (home)

CE: D-side overhead network

EX: Exchange

FU: Frames (exchange)

LN: E-side Underground network

OK: Line has tested OK

OTHER: Other faults

Rainfall: This contains daily rainfall information for the same 4 year period

Observation Date: Date

Rainfall mm: Rainfall (in mm)

Calendar Lookup: This contains a lookup to give calendar information such as day of week and bank holidays

Actual Date: Date

Day of Week: Day of Week

Day Num Cal Week: Numeric value of day within the week

Day Num Cal Month: Numeric value of day within the month

Day Num Cal Year: Numeric value of day within the year

Bank holiday: Flags the bank holiday (Y: bank holiday)

In [1]:

```

1 # Import Libraries
2 from sklearn.model_selection import train_test_split
3 from sklearn.linear_model import LinearRegression
4 import pandas as pd
5 import numpy as np
6 import os
7 import matplotlib.pyplot as plt
8 import matplotlib.dates as mdates
9 import seaborn as sb
10 import datetime
11 from datetime import timedelta,datetime
12 from sklearn import metrics

```

In [2]:

```

1 X = pd.read_excel('RainFallData.xlsx',index_col = 0)
2 Y = pd.read_excel('FaultData.xlsx', index_col = 0)

```

In [3]:

```

1 X['OBS_DATE'] = X.index
2 print(X)

```

OBSERVATION_DATE	RAINFALL_MM	OBS_DATE
2014-01-01	10.6	2014-01-01
2014-01-02	4.2	2014-01-02
2014-01-03	5.8	2014-01-03
2014-01-04	1.8	2014-01-04
2014-01-05	11.0	2014-01-05
2014-01-06	5.6	2014-01-06
2014-01-07	5.9	2014-01-07
2014-01-08	14.8	2014-01-08
2014-01-09	0.8	2014-01-09
2014-01-10	5.2	2014-01-10
2014-01-11	0.2	2014-01-11
2014-01-12	3.8	2014-01-12
2014-01-13	8.1	2014-01-13
2014-01-14	7.6	2014-01-14
2014-01-15	5.3	2014-01-15
2014-01-16	3.3	2014-01-16
2014-01-17	2.0	2014-01-17



In [4]:

```

1 # Rainfall Every 1 to 7 Days
2 X['1_Day_Rain'] = [X['RAINFALL_MM'][max(0,i-1)] for i in range(0,len(X))]
3 X['2_Day_Rain'] = [X['RAINFALL_MM'][max(0,i-2)] for i in range(0,len(X))]
4 X['3_Day_Rain'] = [X['RAINFALL_MM'][max(0,i-3)] for i in range(0,len(X))]
5 X['4_Day_Rain'] = [X['RAINFALL_MM'][max(0,i-4)] for i in range(0,len(X))]
6 X['5_Day_Rain'] = [X['RAINFALL_MM'][max(0,i-5)] for i in range(0,len(X))]
7 X['6_Day_Rain'] = [X['RAINFALL_MM'][max(0,i-6)] for i in range(0,len(X))]
8 X['7_Day_Rain'] = [X['RAINFALL_MM'][max(0,i-7)] for i in range(0,len(X))]
9
10 ## Cummulative Rain for each week
11 X['Cummulative_Rain'] = [sum(X['RAINFALL_MM'][max(0,i-7):i]) for i in range(0,len(X))]
12 X['YMD'] = [X['OBS_DATE'][i].year*10000+X['OBS_DATE'][i].month*100+X['OBS_DATE'][i].day]
13 X

```

Out[4]:

OBSERVATION_DATE	RAINFALL_MM	OBS_DATE	1_Day_Rain	2_Day_Rain	3_Day_Rain	4_
2014-01-01	10.6	2014-01-01	10.6	10.6	10.6	
2014-01-02	4.2	2014-01-02	10.6	10.6	10.6	
2014-01-03	5.8	2014-01-03	4.2	10.6	10.6	
2014-01-04	1.8	2014-01-04	5.8	4.2	10.6	
2014-01-05	11.0	2014-01-05	1.8	5.8	4.2	
2014-01-06	5.6	2014-01-06	11.0	1.8	5.8	
2014-01-07	5.9	2014-01-07	5.6	11.0	1.8	
2014-01-08	14.8	2014-01-08	5.9	5.6	11.0	
2014-01-09	0.8	2014-01-09	14.8	5.9	5.6	
2014-01-10	5.2	2014-01-10	0.8	14.8	5.9	
2014-01-11	0.2	2014-01-11	5.2	0.8	14.8	
2014-01-12	3.8	2014-01-12	0.2	5.2	0.8	
2014-01-13	8.1	2014-01-13	3.8	0.2	5.2	
2014-01-14	7.6	2014-01-14	8.1	3.8	0.2	
2014-01-15	5.3	2014-01-15	7.6	8.1	3.8	
2014-01-16	3.3	2014-01-16	5.3	7.6	8.1	
2014-01-17	2.0	2014-01-17	3.3	5.3	7.6	
2014-01-18	9.2	2014-01-18	2.0	3.3	5.3	
2014-01-19	2.6	2014-01-19	9.2	2.0	3.3	
2014-01-20	0.2	2014-01-20	2.6	9.2	2.0	
2014-01-21	5.6	2014-01-21	0.2	2.6	9.2	
2014-01-22	2.3	2014-01-22	5.6	0.2	2.6	
2014-01-23	1.5	2014-01-23	2.3	5.6	0.2	
2014-01-24	6.4	2014-01-24	1.5	2.3	5.6	

OBSERVATION_DATE	RAINFALL_MM	OBS_DATE	1_Day_Rain	2_Day_Rain	3_Day_Rain	4
2014-01-25	3.2	2014-01-25	6.4	1.5	2.3	
2014-01-26	9.9	2014-01-26	3.2	6.4	1.5	
2014-01-27	10.0	2014-01-27	9.9	3.2	6.4	
2014-01-28	8.0	2014-01-28	10.0	9.9	3.2	
2014-01-29	2.1	2014-01-29	8.0	10.0	9.9	
2014-01-30	0.0	2014-01-30	2.1	8.0	10.0	
...	...	...	...	...	...	
2017-02-28	11.5	2017-02-28	7.7	11.0	11.6	
2017-03-01	14.4	2017-03-01	11.5	7.7	11.0	
2017-03-02	1.9	2017-03-02	14.4	11.5	7.7	
2017-03-03	5.3	2017-03-03	1.9	14.4	11.5	
2017-03-04	3.8	2017-03-04	5.3	1.9	14.4	
2017-03-05	11.4	2017-03-05	3.8	5.3	1.9	
2017-03-06	4.3	2017-03-06	11.4	3.8	5.3	
2017-03-07	7.8	2017-03-07	4.3	11.4	3.8	
2017-03-08	0.1	2017-03-08	7.8	4.3	11.4	
2017-03-09	0.9	2017-03-09	0.1	7.8	4.3	
2017-03-10	0.8	2017-03-10	0.9	0.1	7.8	
2017-03-11	1.7	2017-03-11	0.8	0.9	0.1	
2017-03-12	0.8	2017-03-12	1.7	0.8	0.9	
2017-03-13	0.1	2017-03-13	0.8	1.7	0.8	
2017-03-14	0.0	2017-03-14	0.1	0.8	1.7	
2017-03-15	0.3	2017-03-15	0.0	0.1	0.8	
2017-03-16	1.8	2017-03-16	0.3	0.0	0.1	
2017-03-17	10.9	2017-03-17	1.8	0.3	0.0	
2017-03-18	2.6	2017-03-18	10.9	1.8	0.3	
2017-03-19	5.1	2017-03-19	2.6	10.9	1.8	
2017-03-21	8.4	2017-03-21	5.1	2.6	10.9	
2017-03-22	2.2	2017-03-22	8.4	5.1	2.6	
2017-03-23	0.1	2017-03-23	2.2	8.4	5.1	
2017-03-24	0.0	2017-03-24	0.1	2.2	8.4	
2017-03-25	0.0	2017-03-25	0.0	0.1	2.2	
2017-03-26	0.0	2017-03-26	0.0	0.0	0.1	
2017-03-27	0.0	2017-03-27	0.0	0.0	0.0	
2017-03-28	4.9	2017-03-28	0.0	0.0	0.0	
2017-03-29	9.0	2017-03-29	4.9	0.0	0.0	
2017-03-30	0.0	2017-03-30	9.0	4.9	0.0	

In [5]:

```
1 print(type(X))
2 print(type(Y))
```

```
<class 'pandas.core.frame.DataFrame'>
<class 'pandas.core.frame.DataFrame'>
```

## Calculate Different types of Faults for 29 to 35 days

In [11]:

```
1 Y.head()
2 (Y['REPORT_DATE'][0]).day
```

Out[11]:

1

In [12]:

```
1 Y['YMD'] = [Y['REPORT_DATE'][i].year*10000+Y['REPORT_DATE'][i].month*100+Y['REPORT_DATE'][i].day for i in range(len(Y))]
2 Y.head()
```

Out[12]:

	REPORT_DATE	INITIAL_MFL	YMD	COUNT
Fault ID				
ID000001	2014-01-01	CE	20140101	1
ID000002	2014-01-01	CA	20140101	1
ID000003	2014-01-01	OTHER	20140101	1
ID000004	2014-01-01	OK	20140101	1
ID000005	2014-01-01	LN	20140101	1

In [13]:

```
1 Y['COUNT'] = 1
2 Z = pd.pivot_table(Y,values=['COUNT'],index='YMD',columns='INITIAL_MFL',aggfunc=sum).reindex(index=Z.index,columns=Z.columns)
```

In [14]:

```

1 Z.columns=['YMD', 'FAULTS_CA', 'FAULTS_CE', 'FAULTS_EX', 'FAULTS_FU', 'FAULTS_LN', 'FAULTS_OI'
2 Z.head()

```

Out[14]:

	YMD	FAULTS_CA	FAULTS_CE	FAULTS_EX	FAULTS_FU	FAULTS_LN	FAULTS_OK	FA
0	20140101	15.0	40.0	6.0	4.0	38.0	26.0	
1	20140102	47.0	100.0	3.0	13.0	67.0	83.0	
2	20140103	63.0	87.0	6.0	14.0	77.0	79.0	
3	20140104	24.0	60.0	3.0	6.0	34.0	47.0	
4	20140105	16.0	28.0	2.0	4.0	31.0	41.0	

In [15]:

```

1 for fault in Z.columns[1:]:
2     newfault = fault.replace('FAULTS_', '29to35Faults_')
3     Z[fault]=Z[fault].fillna(0).replace(np.nan,0)
4     Z[newfault]=[sum(Z[fault][min(i+29,len(Z)):min(i+35,len(Z))]) for i in range(0,len(

```

In [16]:

```
1 Z.tail(50)
```

Out[16]:

	YMD	FAULTS_CA	FAULTS_CE	FAULTS_EX	FAULTS_FU	FAULTS_LN	FAULTS_OK	FAULTS_OI
1136	20170210	13.0	89.0	0.0	12.0	37.0	76.0	
1137	20170211	16.0	56.0	0.0	6.0	27.0	47.0	
1138	20170212	15.0	40.0	1.0	3.0	10.0	46.0	
1139	20170213	29.0	84.0	6.0	14.0	36.0	81.0	
1140	20170214	20.0	85.0	2.0	10.0	34.0	75.0	
1141	20170215	26.0	101.0	3.0	13.0	33.0	108.0	
1142	20170216	25.0	77.0	2.0	10.0	25.0	80.0	
1143	20170217	21.0	66.0	1.0	17.0	27.0	72.0	
1144	20170218	10.0	57.0	0.0	13.0	24.0	51.0	

In [17]:

```

1 ## Merge DataFrame
2 df=X.merge(Z,on='YMD')
3 df = df[7:len(df)-35]
4 df = df.reset_index()
5 df

```

Out[17]:

	index	RAINFALL_MM	OBS_DATE	1_Day_Rain	2_Day_Rain	3_Day_Rain	4_Day_Rain	5_Day_Rain	6
0	7	14.8	2014-01-08	5.9	5.6	11.0	1.8	5.8	
1	8	0.8	2014-01-09	14.8	5.9	5.6	11.0	1.8	
2	9	5.2	2014-01-10	0.8	14.8	5.9	5.6	11.0	
3	10	0.2	2014-01-11	5.2	0.8	14.8	5.9	5.6	
4	11	3.8	2014-01-12	0.2	5.2	0.8	14.8	5.9	
5	12	8.1	2014-01-13	3.8	0.2	5.2	0.8	14.8	
6	13	7.6	2014-01-14	8.1	3.8	0.2	5.2	0.8	
7	14	5.3	2014-01-15	7.6	8.1	3.8	0.2	5.2	
8	15	3.3	2014-01-16	5.3	7.6	8.1	3.8	0.2	

In [18]:

```
1 train, test = train_test_split(df, test_size=0.2)
```

In [19]:

```
1 train.head()
```

Out[19]:

FAULTS_LN	FAULTS_OK	FAULTS_OTHER	29to35Faults_CA	29to35Faults_CE	29to35Faults_EX
36.0	80.0	24.0	108.0	289.0	16.0
46.0	126.0	23.0	104.0	295.0	35.0
8.0	40.0	4.0	152.0	419.0	86.0
13.0	34.0	6.0	156.0	422.0	36.0
51.0	95.0	10.0	118.0	341.0	19.0



In [20]:

```
1 test.head()
```

Out[20]:

	index	RAINFALL_MM	OBS_DATE	1_Day_Rain	2_Day_Rain	3_Day_Rain	4_Day_Rain	5_
1117	1124	6.2	2017-01-28	5.6	0.0	0.0	0.5	
890	897	9.5	2016-06-15	12.3	6.8	3.5	4.4	
142	149	0.0	2014-05-30	0.3	1.4	1.2	0.4	
899	906	0.9	2016-06-24	0.2	0.4	0.7	3.8	
905	912	7.3	2016-06-30	14.0	9.3	0.1	6.0	

5 rows × 26 columns

## CA - FAULT MODELS

### Linear Regression

In [24]:

```
1 # Split Train Test Data for CA Faults
2 df_ca = df[['Cummulative_Rain','29to35Faults_CA']]
3 df_ca_train, df_ca_test = train_test_split(df_ca, test_size=0.2,random_state=1801)
```

In [192]:

```
1 df_ca.corr()
```

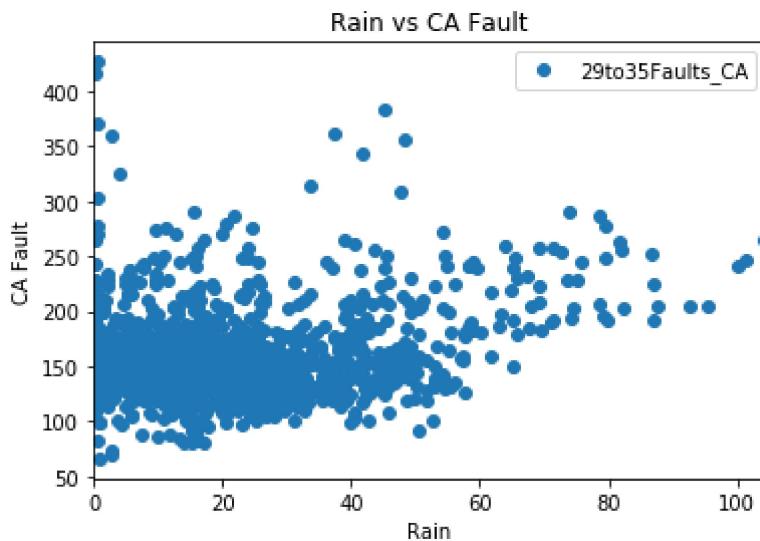
Out[192]:

	Cummulative_Rain	29to35Faults_CA
Cummulative_Rain	1.000000	0.219337
29to35Faults_CA	0.219337	1.000000



In [25]:

```
1 df_ca.plot(x='Cummulative_Rain', y='29to35Faults_CA', style='o')
2 plt.title('Rain vs CA Fault')
3 plt.xlabel('Rain')
4 plt.ylabel('CA Fault')
5 plt.show()
```

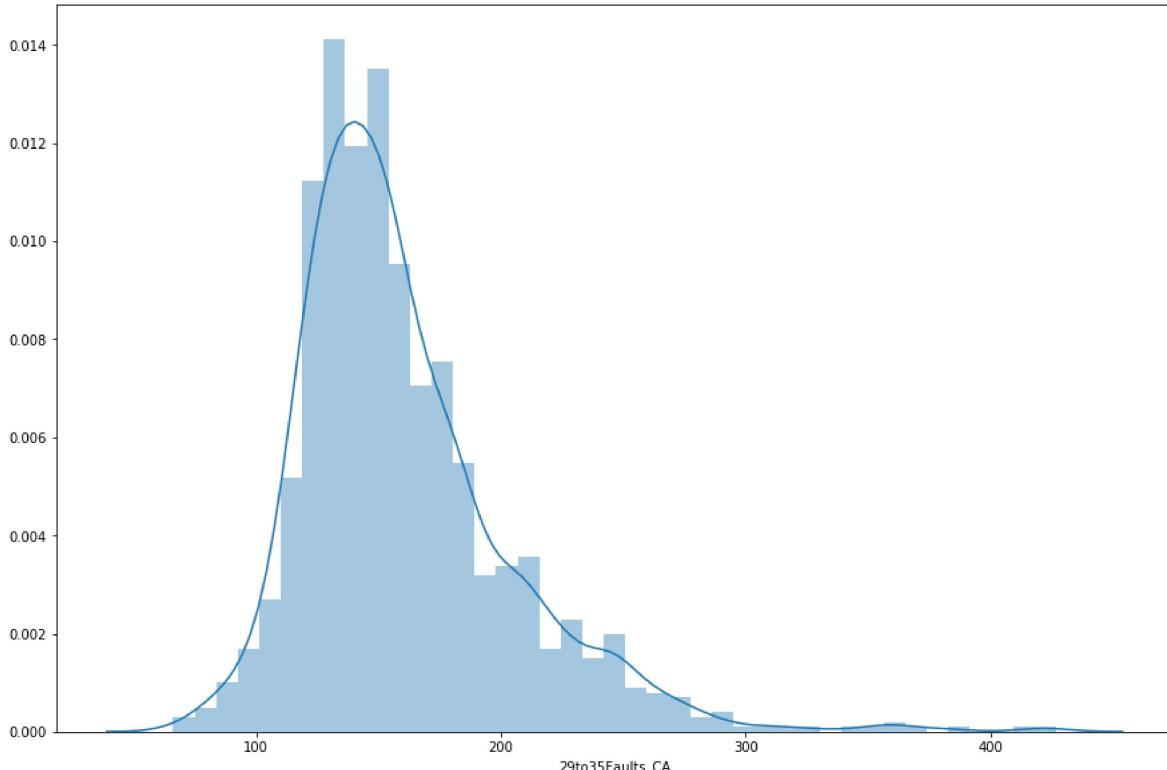


In [26]:

```
1 # Check Average Faults
2 plt.figure(figsize=(15,10))
3 plt.tight_layout()
4 sb.distplot(df_ca['29to35Faults_CA'])
```

Out[26]:

&lt;matplotlib.axes.\_subplots.AxesSubplot at 0x23d96c4d828&gt;





In [27]:

```

1 # Find TRAIN Set Data for 15 to 21 Days for CA Faults
2 df_ca_train_np_rain=df_ca_train['Cummulative_Rain'].to_numpy()
3 df_ca_train_np_fault = df_ca_train['29to35Faults_CA'].to_numpy()
4
5 # Reshape test train data
6 df_ca_train_rain = df_ca_train_np_rain.reshape(len(df_ca_train_np_rain),1)
7 df_ca_train_fault = df_ca_train_np_fault.reshape(len(df_ca_train_np_fault),1)

```



In [28]:

```
1 lr = LinearRegression().fit(df_ca_train_rain,df_ca_train_fault)
```



In [29]:

```

1 print(lr.coef_)
2 print(lr.intercept_)

```

[[0.47427518]]  
[149.49912682]



In [30]:

```

1 # Find TEST Set Data for 15 to 21 Days for CA Faults
2 df_ca_test_np_rain=df_ca_test['Cummulative_Rain'].to_numpy()
3 df_ca_test_np_fault = df_ca_test['29to35Faults_CA'].to_numpy()
4
5 df_ca_test_rain = df_ca_test_np_rain.reshape(len(df_ca_test_np_rain),1)
6 df_ca_test_fault = df_ca_test_np_fault.reshape(len(df_ca_test_np_fault),1)

```



In [31]:

```

1 # Train Model Accuracy
2 lr.score(df_ca_train_rain, df_ca_train_fault)

```

Out[31]:

0.03919502676997954



In [32]:

```

1 # testing model accuracy in test set
2 lr.score(df_ca_test_rain, df_ca_test_fault)

```

Out[32]:

0.08537420789524908

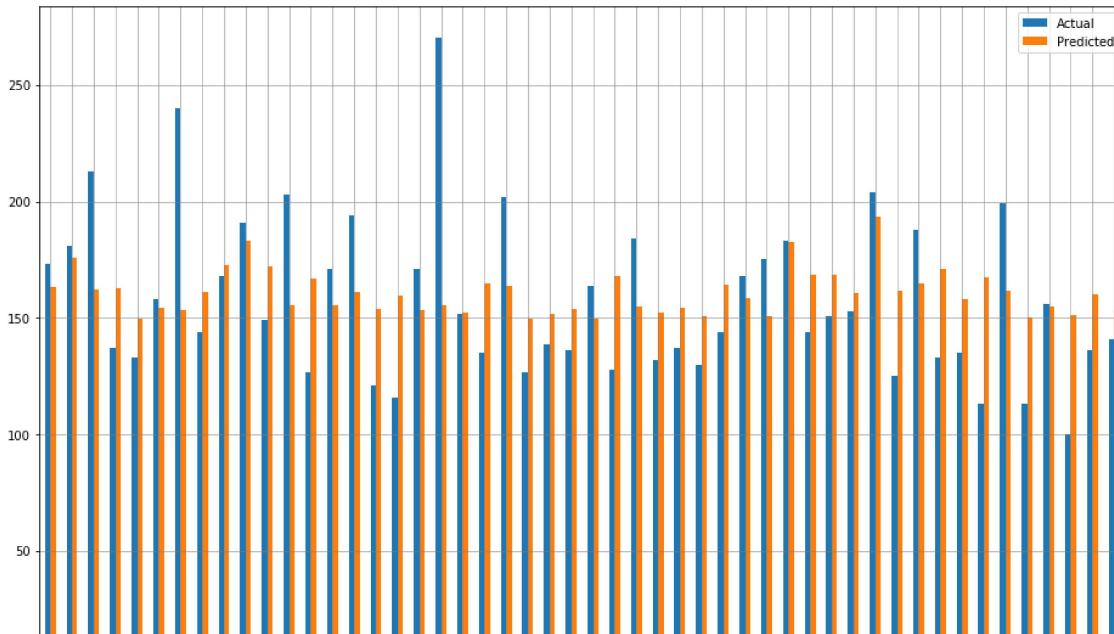


In [33]:

```

1 predict_linear = lr.predict(df_ca_test_rain)
2 actual_pred_linear= pd.DataFrame({'Actual': df_ca_test_fault.flatten(), 'Predicted': p
3 bar_plot_df = actual_pred_linear.head(50)
4 bar_plot_df.plot(kind='bar',figsize=(16,10))
5 plt.grid(which='major', linestyle='-', linewidth='0.5', color='grey')
6 plt.grid(which='minor', linestyle=':', linewidth='0.5', color='black')
7 plt.show()

```

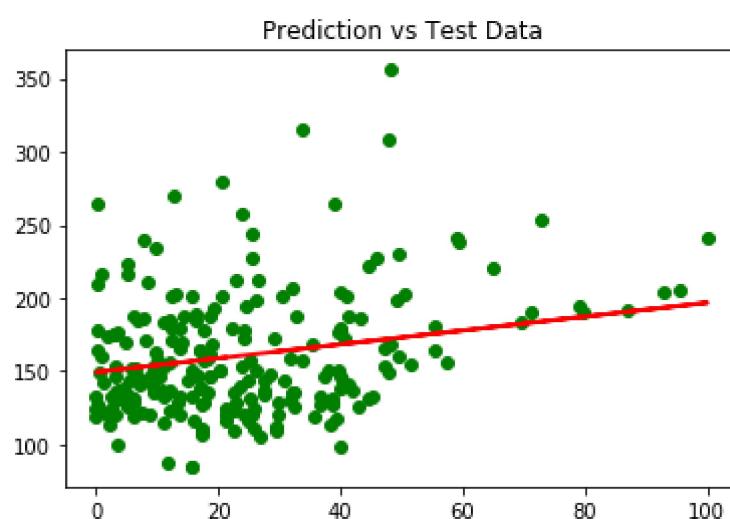


In [34]:

```

1 plt.scatter(df_ca_test_rain, df_ca_test_fault, color='green')
2 plt.plot(df_ca_test_rain, predict_linear, color='red', linewidth=2)
3 plt.title('Prediction vs Test Data')
4 plt.show()

```



In [35]:

```
1 rmspe = (np.sqrt(np.mean(np.square((df_ca_test_fault - predict_linear) / df_ca_test_fa
2 print('Root Mean Squared Error Percentage for CA Faults: ', rmspe, '%'))
```

Root Mean Squared Error Percentage for CA Faults: 24.497530838319857 %

## OK FAULT MODELS

In [36]:

```
1 df_ok = df[['Cummulative_Rain','29to35Faults_OK']]
2 df_ok_train, df_ok_test = train_test_split(df_ok, test_size=0.2, random_state=1801)
```

In [193]:

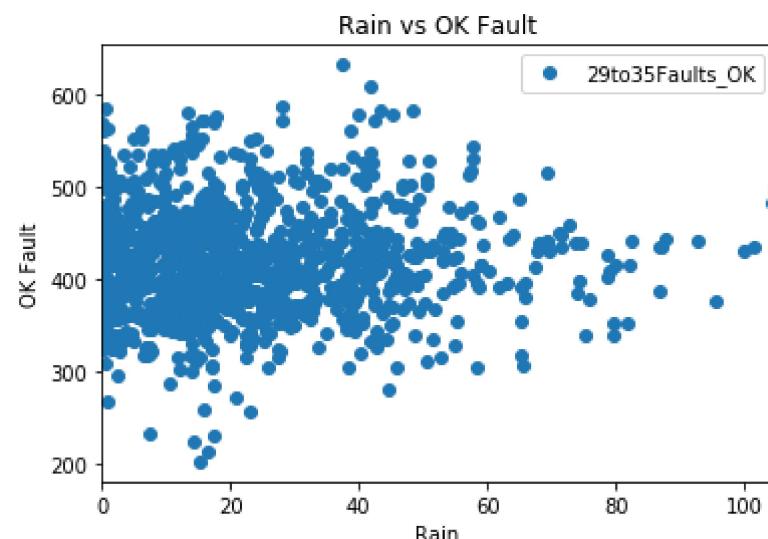
```
1 df_ok.corr()
```

Out[193]:

	Cummulative_Rain	29to35Faults_OK
Cummulative_Rain	1.000000	0.047124
29to35Faults_OK	0.047124	1.000000

In [37]:

```
1 df_ok.plot(x='Cummulative_Rain', y='29to35Faults_OK', style='o')
2 plt.title('Rain vs OK Fault')
3 plt.xlabel('Rain')
4 plt.ylabel('OK Fault')
5 plt.show()
```





In [38]:

```

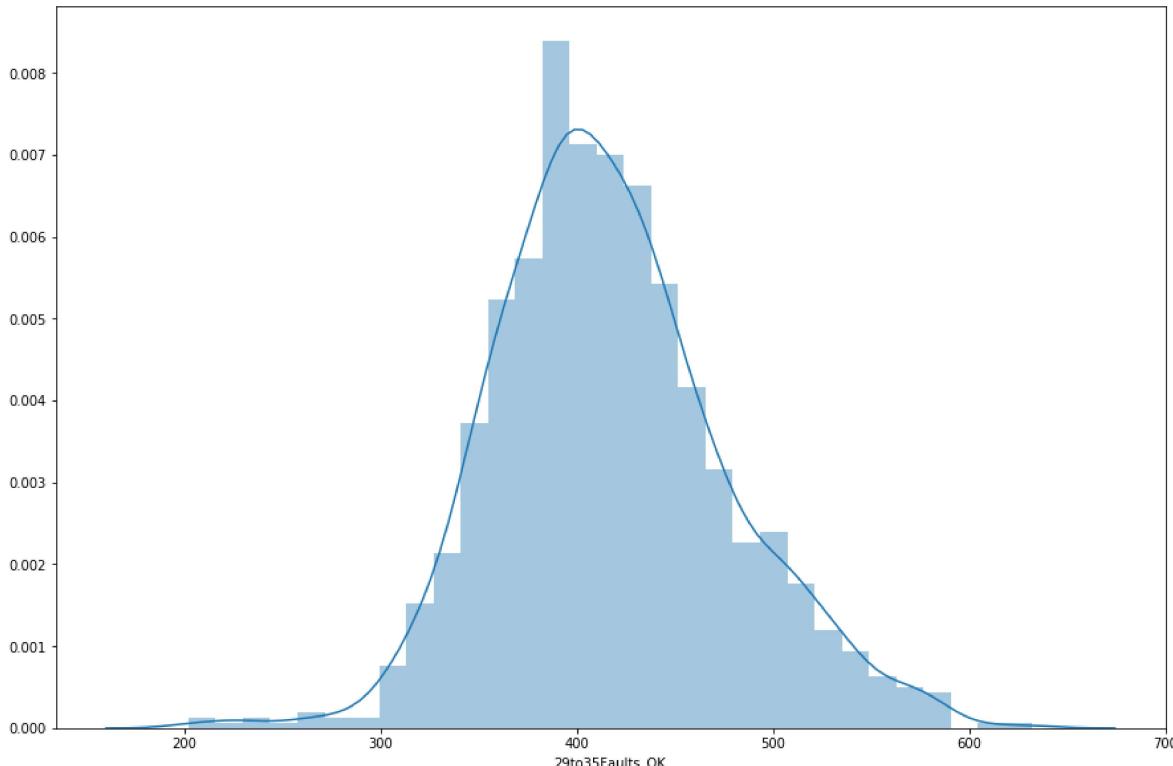
1 # Check Average Faults
2 plt.figure(figsize=(15,10))
3 plt.tight_layout()
4 sb.distplot(df_ok['29to35Faults_OK'])

```



Out[38]:

&lt;matplotlib.axes.\_subplots.AxesSubplot at 0x23d9a295320&gt;



In [39]:

```

1 # Find TRAIN Set Data for 15 to 21 Days for OK Faults
2 df_ok_train_np_rain=df_ok_train['Cummulative_Rain'].to_numpy()
3 df_ok_train_np_fault = df_ok_train['29to35Faults_OK'].to_numpy()
4
5 # Reshape test train data
6 df_ok_train_rain = df_ok_train_np_rain.reshape(len(df_ok_train_np_rain),1)
7 df_ok_train_fault = df_ok_train_np_fault.reshape(len(df_ok_train_np_fault),1)

```



In [40]:

```
1 lr = LinearRegression().fit(df_ok_train_rain,df_ok_train_fault)
```





In [41]:

```
1 # Train Model Accuracy
2 lr.score(df_ok_train_rain, df_ok_train_fault)
```



Out[41]:

0.0008892046714681355



In [42]:

```
1 # Find TEST Set Data for 15 to 21 Days for OK Faults
2 df_ok_test_np_rain=df_ok_test['Cummulative_Rain'].to_numpy()
3 df_ok_test_np_fault = df_ok_test['29to35Faults_OK'].to_numpy()
4
5 df_ok_test_rain = df_ok_test_np_rain.reshape(len(df_ok_test_np_rain),1)
6 df_ok_test_fault = df_ok_test_np_fault.reshape(len(df_ok_test_np_fault),1)
```



In [43]:

```
1 # testing model accuracy in test set
2 lr.score(df_ok_test_rain, df_ok_test_fault)
```



Out[43]:

0.003698449911566648

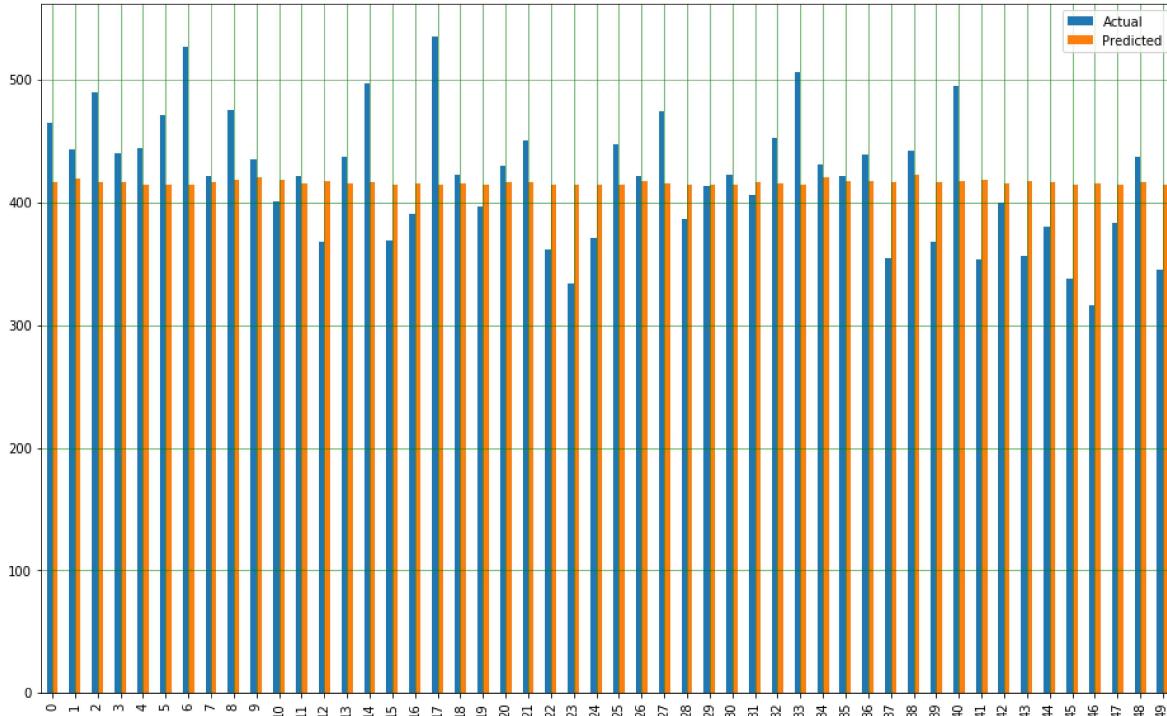


In [44]:

```

1 pred_OK_Fault_linear = lr.predict(df_ok_test_rain)
2 OK_actual_pred_linear= pd.DataFrame({'Actual': df_ok_test_fault.flatten(), 'Predicted'
3 bar_plot_df = OK_actual_pred_linear.head(50)
4 bar_plot_df.plot(kind='bar',figsize=(16,10))
5 plt.grid(which='major', linestyle='-', linewidth='0.5', color='green')
6 plt.grid(which='minor', linestyle=':', linewidth='0.5', color='black')
7 plt.show()

```



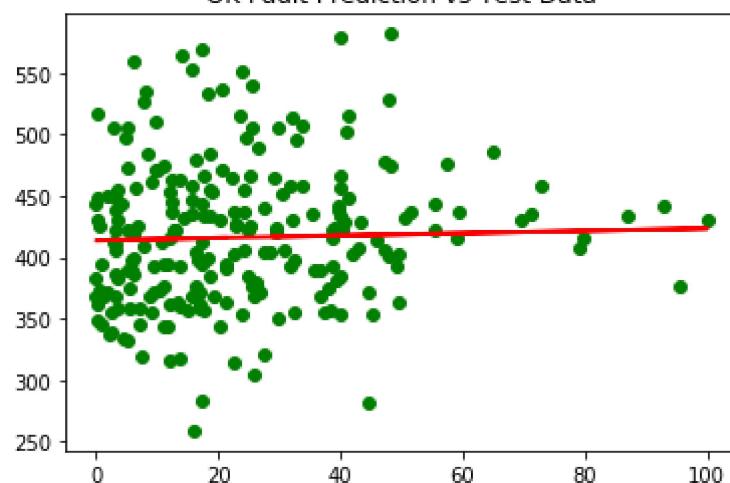
In [45]:

```

1 plt.scatter(df_ok_test_rain, df_ok_test_fault, color='green')
2 plt.plot(df_ca_test_rain, pred_OK_Fault_linear, color='red', linewidth=2)
3 plt.title('OK Fault Prediction vs Test Data')
4 plt.show()

```

OK Fault Prediction vs Test Data



In [46]:

```

1 rmspe_ok = (np.sqrt(np.mean(np.square((df_ok_test_fault - pred_OK_Fault_linear) / df_ok_test_fault))) * 100)
2 print('Root Mean Squared Error Percentage for OK Faults: ', rmspe_ok, '%')

```

Root Mean Squared Error Percentage for OK Faults: 14.240701528516144 %

## CE FAULT MODEL

### Linear Regression on Faults CE

In [47]:

```

1 df_ce = df[['Cummulative_Rain','29to35Faults_CE']]
2 df_ce_train, df_ce_test = train_test_split(df_ce, test_size=0.2, random_state=1801)

```

In [194]:

```
1 df_ce.corr()
```

Out[194]:

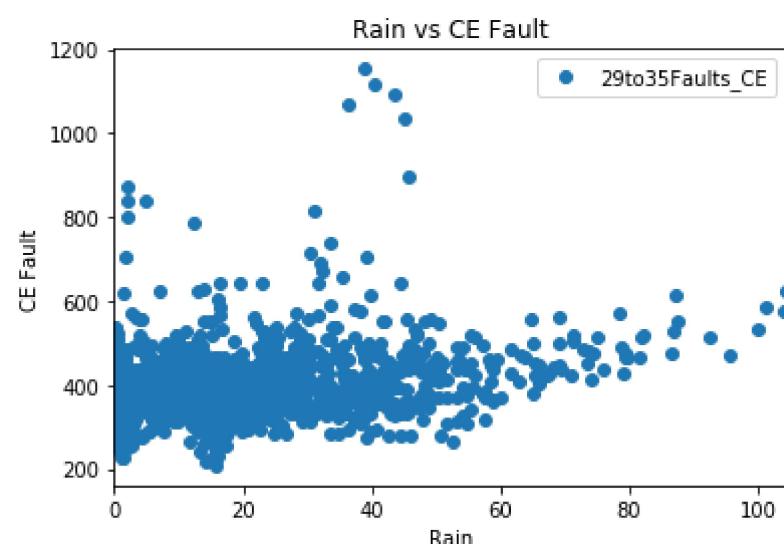
	Cummulative_Rain	29to35Faults_CE
Cummulative_Rain	1.000000	0.221661
29to35Faults_CE	0.221661	1.000000

In [48]:

```

1 df_ce.plot(x='Cummulative_Rain', y='29to35Faults_CE', style='o')
2 plt.title('Rain vs CE Fault')
3 plt.xlabel('Rain')
4 plt.ylabel('CE Fault')
5 plt.show()

```





In [49]:

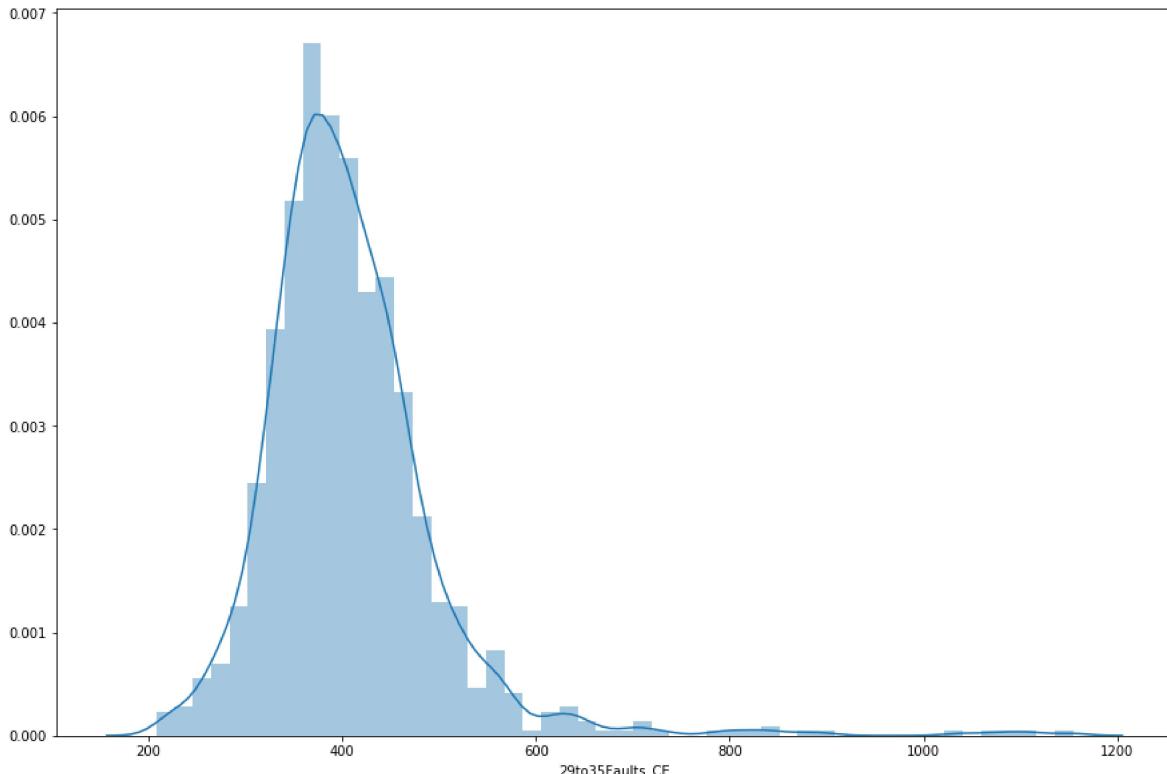
```

1 # Check Average Faults
2 plt.figure(figsize=(15,10))
3 plt.tight_layout()
4 sb.distplot(df_ce['29to35Faults_CE'])

```

Out[49]:

&lt;matplotlib.axes.\_subplots.AxesSubplot at 0x23d97128c18&gt;



In [50]:

```

1 # Find TRAIN Set Data for 15 to 21 Days for CA Faults
2 df_ce_train_np_rain=df_ce_train['Cummulative_Rain'].to_numpy()
3 df_ce_train_np_fault = df_ce_train['29to35Faults_CE'].to_numpy()
4
5 # Reshape test train data
6 df_ce_train_rain = df_ce_train_np_rain.reshape(len(df_ce_train_np_rain),1)
7 df_ce_train_fault = df_ce_train_np_fault.reshape(len(df_ce_train_np_fault),1)

```



In [51]:

```
1 lr = LinearRegression().fit(df_ce_train_rain,df_ce_train_fault)
```





In [52]:

```
1 lr.score(df_ce_train_rain, df_ce_train_fault)
```

Out[52]:

0.04029796464960034



In [53]:

```
1 # Find TEST Set Data for 15 to 21 Days for CA Faults
2 df_ce_test_np_rain=df_ce_test['Cummulative_Rain'].to_numpy()
3 df_ce_test_np_fault = df_ce_test['29to35Faults_CE'].to_numpy()
4
5 df_ce_test_rain = df_ce_test_np_rain.reshape(len(df_ce_test_np_rain),1)
6 df_ce_test_fault = df_ce_test_np_fault.reshape(len(df_ce_test_np_fault),1)
```



In [54]:

```
1 # testing model accuracy in test set
2 lr.score(df_ce_test_rain, df_ce_test_fault)
```

Out[54]:

0.08133819075504167



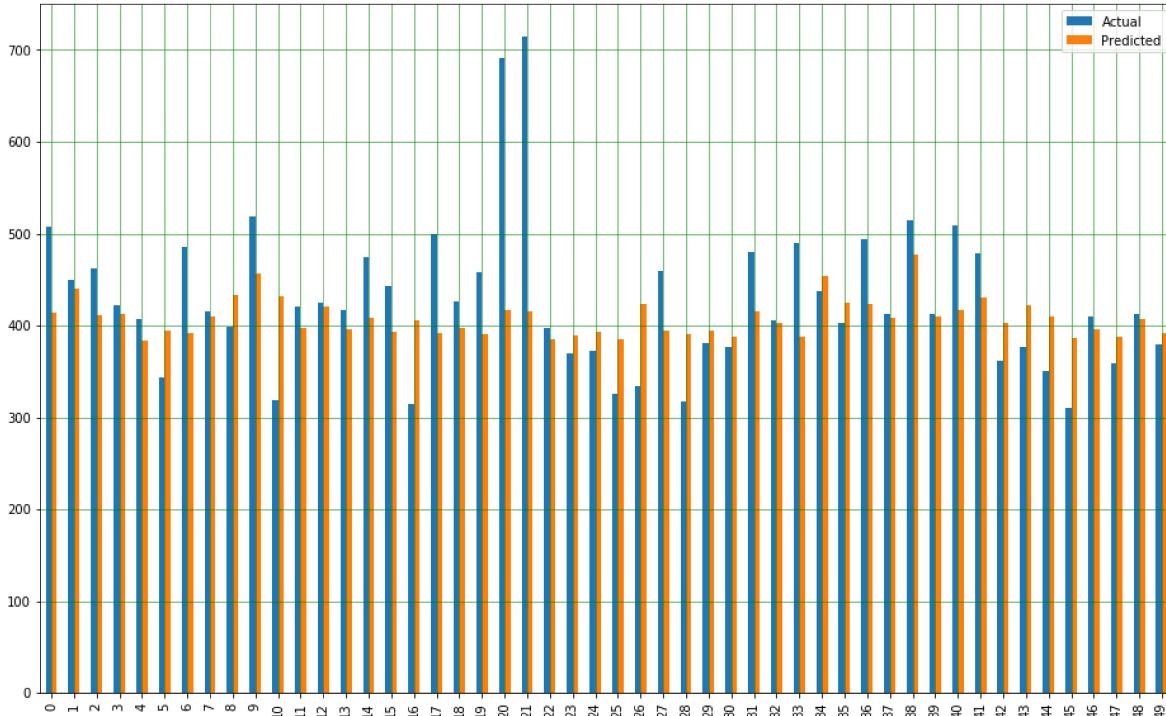


In [55]:

```

1 pred_ce_Fault_linear = lr.predict(df_ce_test_rain)
2 ce_actual_pred_linear= pd.DataFrame({'Actual': df_ce_test_fault.flatten(), 'Predicted'
3 bar_plot_df = ce_actual_pred_linear.head(50)
4 bar_plot_df.plot(kind='bar',figsize=(16,10))
5 plt.grid(which='major', linestyle='-', linewidth='0.5', color='green')
6 plt.grid(which='minor', linestyle=':', linewidth='0.5', color='black')
7 plt.show()

```

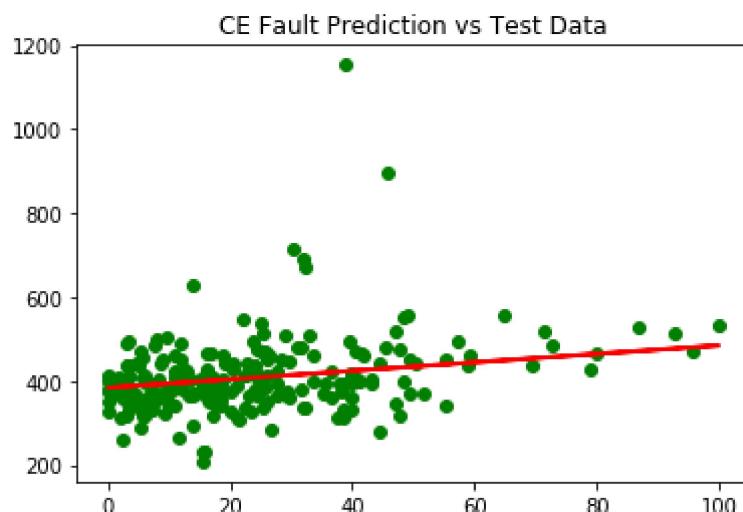


In [56]:

```

1 plt.scatter(df_ce_test_rain, df_ce_test_fault, color='green')
2 plt.plot(df_ce_test_rain, pred_ce_Fault_linear, color='red', linewidth=2)
3 plt.title('CE Fault Prediction vs Test Data')
4 plt.show()

```



In [57]:

```

1 rmspe_ce = (np.sqrt(np.mean(np.square((df_ce_test_fault - pred_ce_Fault_linear) / df_ce
2 print('Root Mean Squared Error Percentage for CE Faults: ', rmspe_ce, '%'))

```

Root Mean Squared Error Percentage for CE Faults: 14.240701528516144 %

## EX FAULT MODEL

### Linear Regression

In [58]:

```

1 df_ex = df[['Cummulative_Rain','29to35Faults_EX']]
2 df_ex_train, df_ex_test = train_test_split(df_ex, test_size=0.2, random_state=1801)

```

In [195]:

```
1 df_ex.corr()
```

Out[195]:

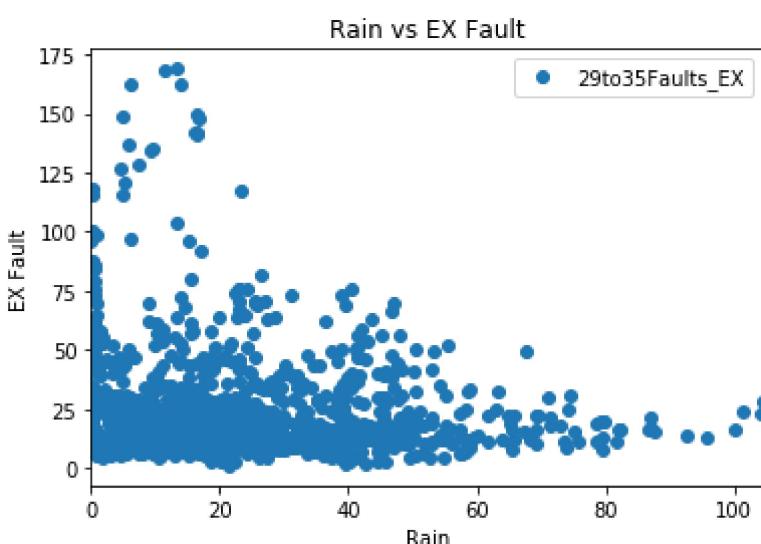
	Cummulative_Rain	29to35Faults_EX
Cummulative_Rain	1.000000	-0.127833
29to35Faults_EX	-0.127833	1.000000

In [59]:

```

1 df_ex.plot(x='Cummulative_Rain', y='29to35Faults_EX', style='o')
2 plt.title('Rain vs EX Fault')
3 plt.xlabel('Rain')
4 plt.ylabel('EX Fault')
5 plt.show()

```





In [60]:

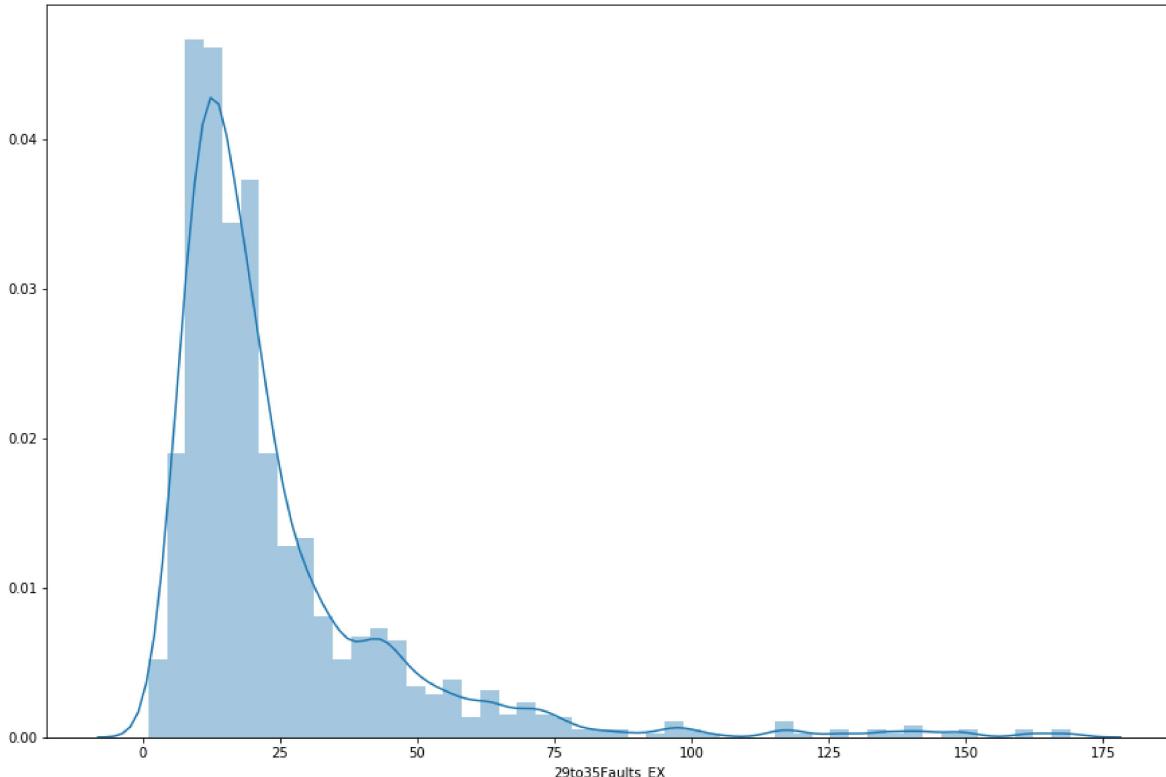
```

1 # Check Average Faults
2 plt.figure(figsize=(15,10))
3 plt.tight_layout()
4 sb.distplot(df_ex['29to35Faults_EX'])

```

Out[60]:

&lt;matplotlib.axes.\_subplots.AxesSubplot at 0x23d97a49f98&gt;



In [61]:

```

1 # Find TRAIN Set Data for 15 to 21 Days for CA Faults
2 df_ex_train_np_rain=df_ex_train['Cummulative_Rain'].to_numpy()
3 df_ex_train_np_fault = df_ex_train['29to35Faults_EX'].to_numpy()
4
5 # Reshape test train data
6 df_ex_train_rain = df_ex_train_np_rain.reshape(len(df_ex_train_np_rain),1)
7 df_ex_train_fault = df_ex_train_np_fault.reshape(len(df_ex_train_np_fault),1)

```

In [62]:

```
1 lr = LinearRegression().fit(df_ex_train_rain,df_ex_train_fault)
```



In [63]:

```
1 lr.score(df_ex_train_rain, df_ex_train_fault)
```

Out[63]:

0.013885844860986651



In [64]:

```
1 # Find TEST Set Data for 15 to 21 Days for CA Faults
2 df_ex_test_np_rain=df_ex_test['Cummulative_Rain'].to_numpy()
3 df_ex_test_np_fault = df_ex_test['29to35Faults_EX'].to_numpy()
4
5 df_ex_test_rain = df_ex_test_np_rain.reshape(len(df_ex_test_np_rain),1)
6 df_ex_test_fault = df_ex_test_np_fault.reshape(len(df_ex_test_np_fault),1)
```



In [65]:

```
1 # testing model accuracy in test set
2 lr.score(df_ex_test_rain, df_ex_test_fault)
```

Out[65]:

0.024408000281454508



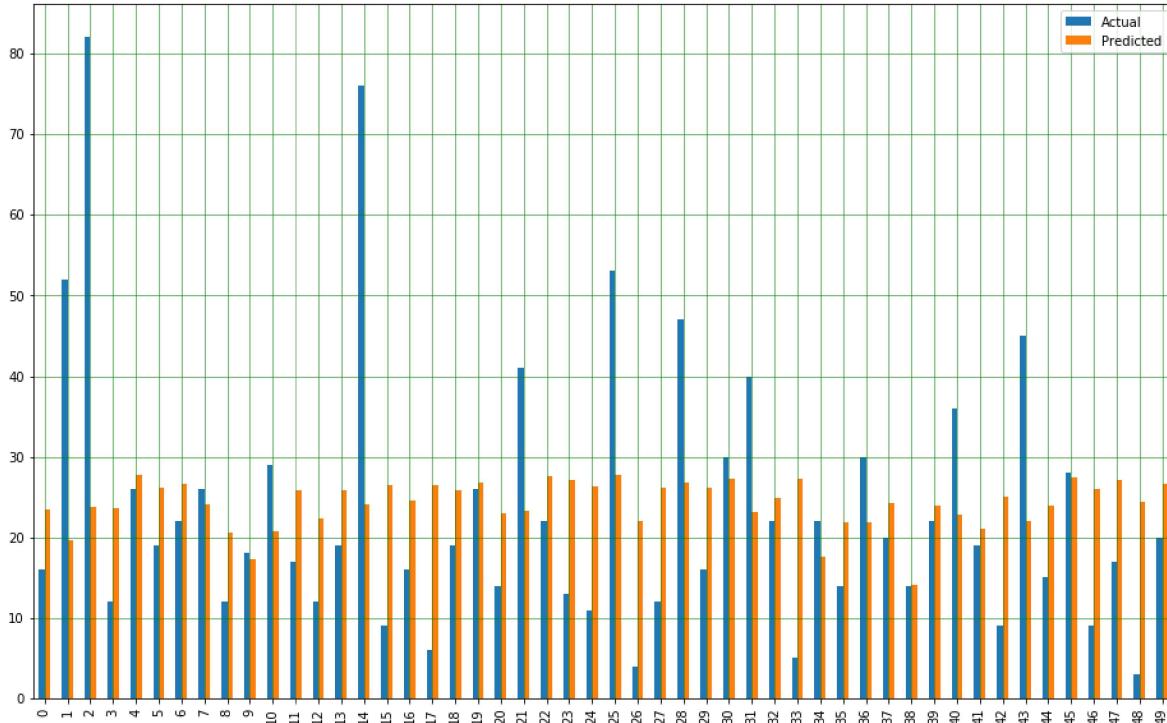


In [66]:

```

1 pred_ex_Fault_linear = lr.predict(df_ex_test_rain)
2 ex_actual_pred_linear= pd.DataFrame({'Actual': df_ex_test_fault.flatten(), 'Predicted'
3 bar_plot_df = ex_actual_pred_linear.head(50)
4 bar_plot_df.plot(kind='bar',figsize=(16,10))
5 plt.grid(which='major', linestyle='-', linewidth='0.5', color='green')
6 plt.grid(which='minor', linestyle=':', linewidth='0.5', color='black')
7 plt.show()

```



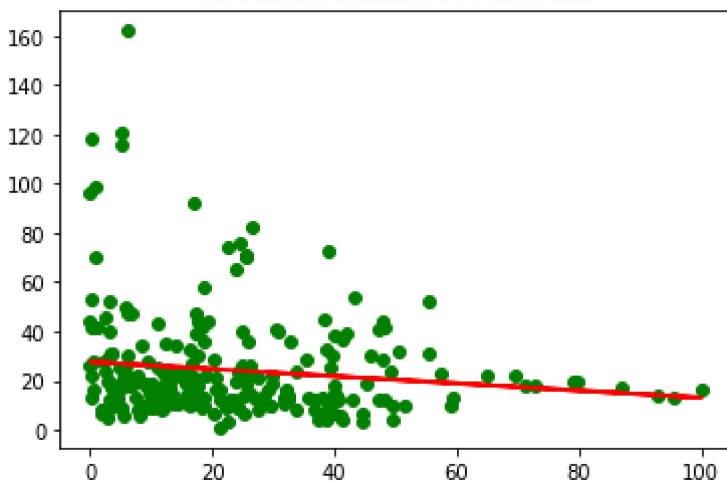
In [67]:

```

1 plt.scatter(df_ex_test_rain, df_ex_test_fault, color='green')
2 plt.plot(df_ex_test_rain, pred_ex_Fault_linear, color='red', linewidth=2)
3 plt.title('EX Fault Prediction vs Test Data')
4 plt.show()

```

EX Fault Prediction vs Test Data



In [68]:

```

1 rmspe_ex = (np.sqrt(np.mean(np.square((df_ex_test_fault - pred_ex_Fault_linear) / df_ex_test_fault))) * 100)
2 print('Root Mean Squared Error Percentage for EX Faults: ', rmspe_ex, '%')

```

Root Mean Squared Error Percentage for EX Faults: 204.0241828291109 %

## FU FAULT MODEL

### Linear Regression

In [69]:

```

1 df_fu = df[['Cummulative_Rain','29to35Faults_FU']]
2 df_fu_train, df_fu_test = train_test_split(df_fu, test_size=0.2, random_state=1801)

```

In [196]:

```
1 df_fu.corr()
```

Out[196]:

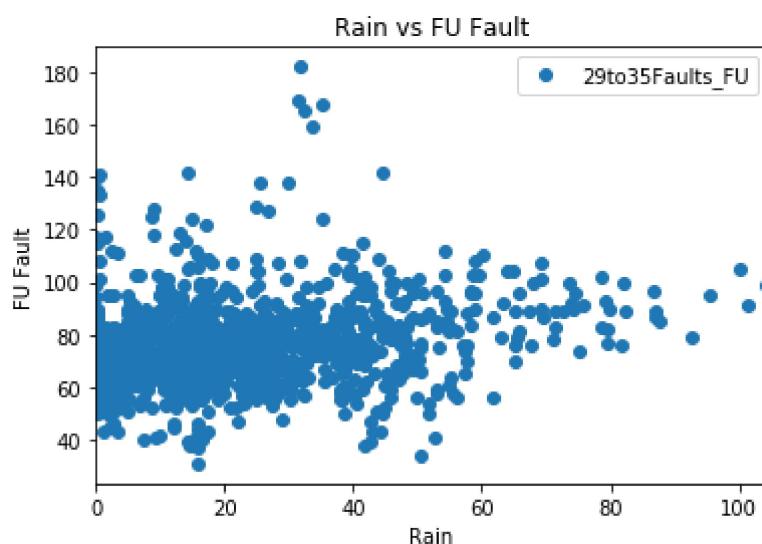
	Cummulative_Rain	29to35Faults_FU
Cummulative_Rain	1.000000	0.200868
29to35Faults_FU	0.200868	1.000000

In [70]:

```

1 df_fu.plot(x='Cummulative_Rain', y='29to35Faults_FU', style='o')
2 plt.title('Rain vs FU Fault')
3 plt.xlabel('Rain')
4 plt.ylabel('FU Fault')
5 plt.show()

```





In [71]:

```

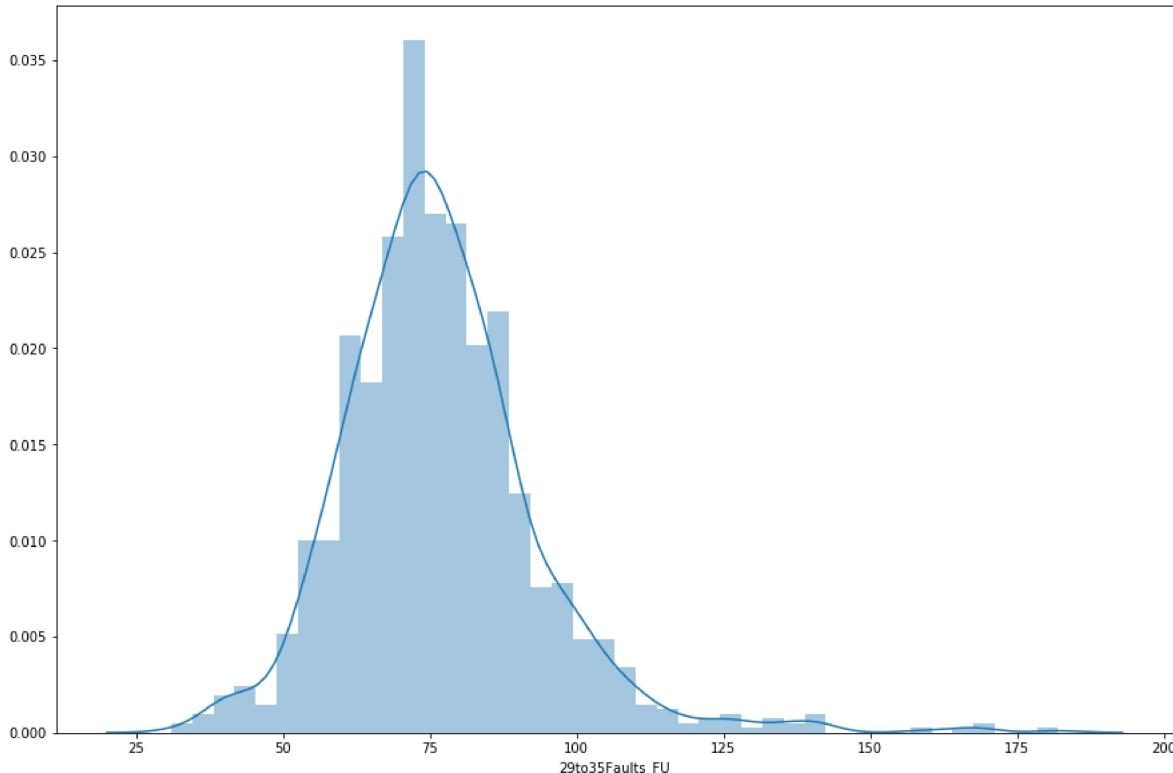
1 # Check Average Faults
2 plt.figure(figsize=(15,10))
3 plt.tight_layout()
4 sb.distplot(df_fu['29to35Faults_FU'])

```



Out[71]:

&lt;matplotlib.axes.\_subplots.AxesSubplot at 0x23d9750af60&gt;



In [72]:

```

1 # Find TRAIN Set Data for 15 to 21 Days for FU Faults
2 df_fu_train_np_rain=df_fu_train['Cummulative_Rain'].to_numpy()
3 df_fu_train_np_fault = df_fu_train['29to35Faults_FU'].to_numpy()
4
5 # Reshape test train data
6 df_fu_train_rain = df_fu_train_np_rain.reshape(len(df_fu_train_np_rain),1)
7 df_fu_train_fault = df_fu_train_np_fault.reshape(len(df_fu_train_np_fault),1)

```



In [73]:

```
1 lr = LinearRegression().fit(df_fu_train_rain,df_fu_train_fault)
```



In [74]:

```
1 lr.score(df_fu_train_rain, df_fu_train_fault)
```

Out[74]:

0.03326805035947833



In [75]:

```
1 # Find TEST Set Data for 15 to 21 Days for FU Faults
2 df_fu_test_np_rain=df_fu_test['Cummulative_Rain'].to_numpy()
3 df_fu_test_np_fault = df_fu_test['29to35Faults_FU'].to_numpy()
4
5 df_fu_test_rain = df_fu_test_np_rain.reshape(len(df_fu_test_np_rain),1)
6 df_fu_test_fault = df_fu_test_np_fault.reshape(len(df_fu_test_np_fault),1)
```



In [76]:

```
1 # testing model accuracy in test set
2 lr.score(df_fu_test_rain, df_fu_test_fault)
```

Out[76]:

0.06734071491453497

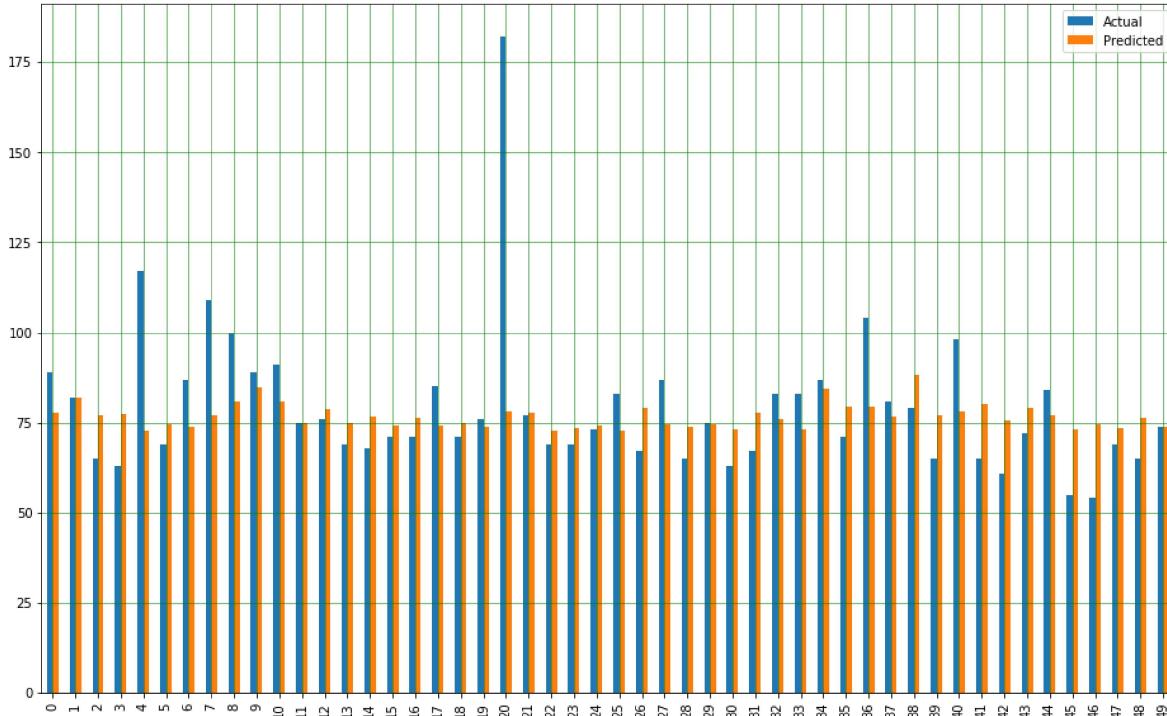


In [77]:

```

1 pred_fu_Fault_linear = lr.predict(df_fu_test_rain)
2 fu_actual_pred_linear= pd.DataFrame({'Actual': df_fu_test_fault.flatten(), 'Predicted'
3 bar_plot_df = fu_actual_pred_linear.head(50)
4 bar_plot_df.plot(kind='bar',figsize=(16,10))
5 plt.grid(which='major', linestyle='-', linewidth='0.5', color='green')
6 plt.grid(which='minor', linestyle=':', linewidth='0.5', color='black')
7 plt.show()

```



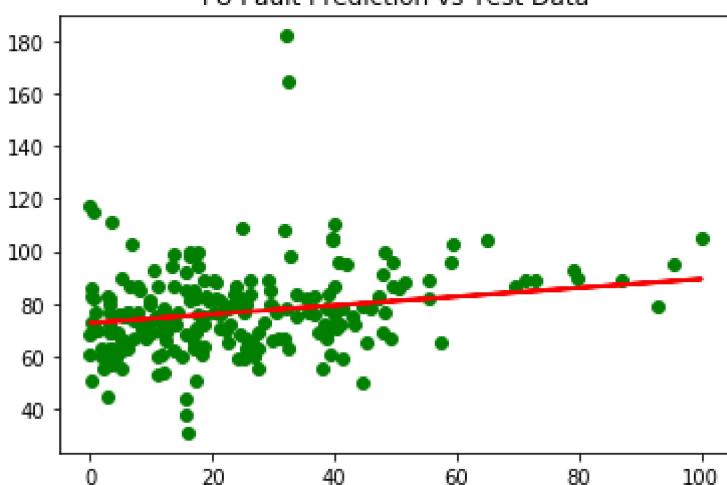
In [78]:

```

1 plt.scatter(df_fu_test_rain, df_fu_test_fault, color='green')
2 plt.plot(df_fu_test_rain, pred_fu_Fault_linear, color='red', linewidth=2)
3 plt.title('FU Fault Prediction vs Test Data')
4 plt.show()

```

FU Fault Prediction vs Test Data



In [79]:

```

1 rmspe_fu = (np.sqrt(np.mean(np.square((df_fu_test_fault - pred_fu_Fault_linear) / df_fu
2 print('Root Mean Squared Error Percentage for FU Faults: ', rmspe_fu, '%'))

```

Root Mean Squared Error Percentage for FU Faults: 21.560791690113394 %

## LN FAULT MODEL

### Linear Regression

In [80]:

```

1 df_ln = df[['Cummulative_Rain','29to35Faults_LN']]
2 df_ln_train, df_ln_test = train_test_split(df_ln, test_size=0.2, random_state=1801)

```

In [197]:

```
1 df_ln.corr()
```

Out[197]:

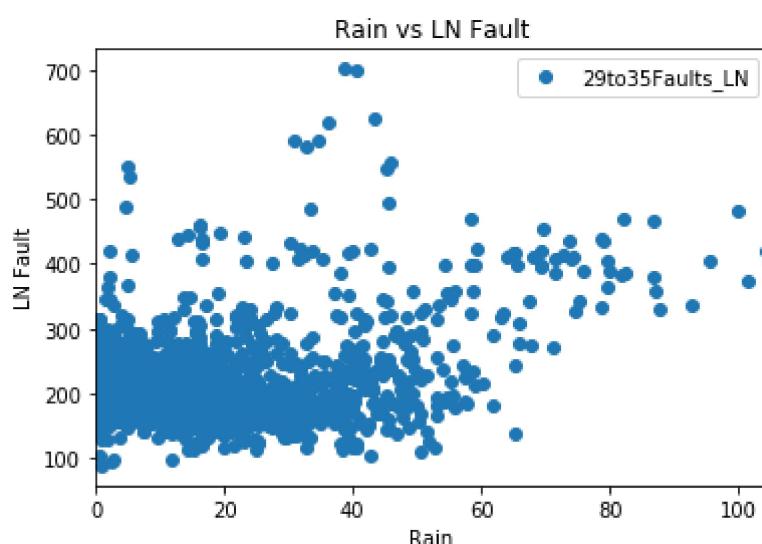
	Cummulative_Rain	29to35Faults_LN
Cummulative_Rain	1.000000	0.302813
29to35Faults_LN	0.302813	1.000000

In [81]:

```

1 df_ln.plot(x='Cummulative_Rain', y='29to35Faults_LN', style='o')
2 plt.title('Rain vs LN Fault')
3 plt.xlabel('Rain')
4 plt.ylabel('LN Fault')
5 plt.show()

```





In [82]:

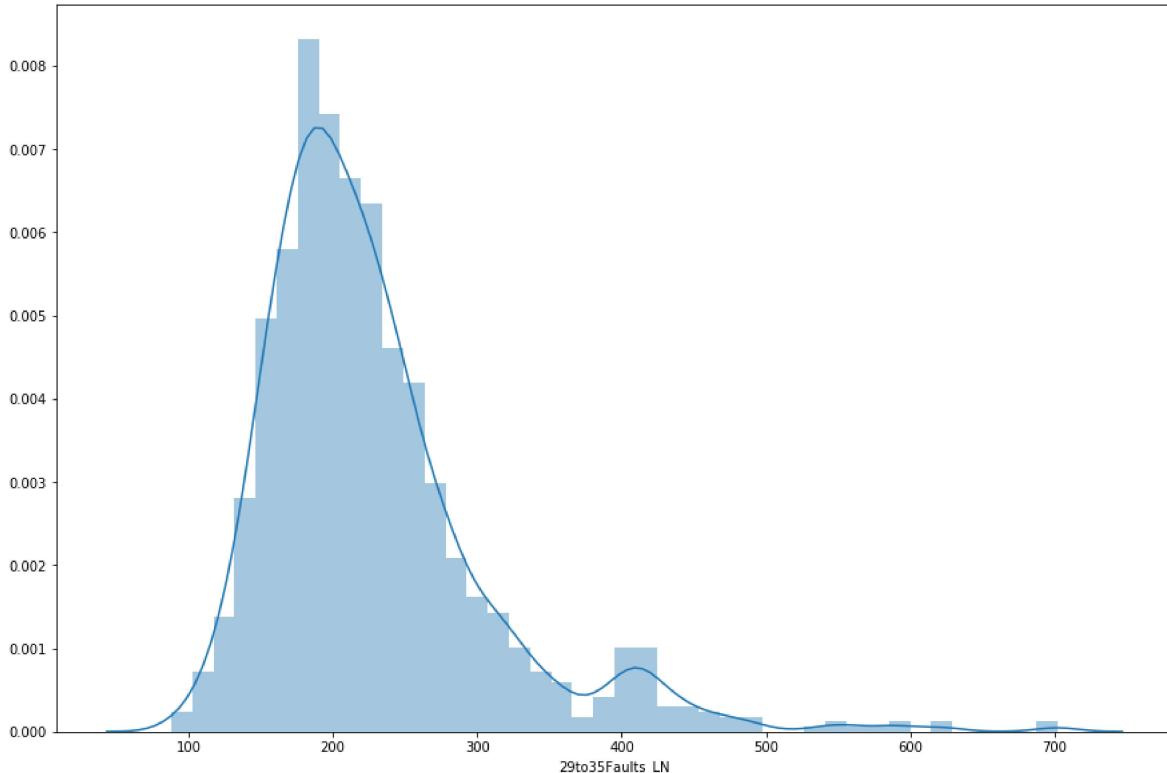
```

1 # Check Average Faults
2 plt.figure(figsize=(15,10))
3 plt.tight_layout()
4 sb.distplot(df_ln['29to35Faults_LN'])

```

Out[82]:

&lt;matplotlib.axes.\_subplots.AxesSubplot at 0x23d96b032b0&gt;



In [83]:

```

1 # Find TRAIN Set Data for 15 to 21 Days for LN Faults
2 df_ln_train_np_rain=df_ln_train['Cummulative_Rain'].to_numpy()
3 df_ln_train_np_fault = df_ln_train['29to35Faults_LN'].to_numpy()
4
5 # Reshape test train data
6 df_ln_train_rain = df_ln_train_np_rain.reshape(len(df_ln_train_np_rain),1)
7 df_ln_train_fault = df_ln_train_np_fault.reshape(len(df_ln_train_np_fault),1)

```

In [84]:

```
1 lr = LinearRegression().fit(df_ln_train_rain,df_ln_train_fault)
```



In [85]:

```
1 lr.score(df_ln_train_rain, df_ln_train_fault)
```

Out[85]:

0.06895246941319677



In [86]:

```
1 # Find TEST Set Data for 15 to 21 Days for FU Faults
2 df_ln_test_np_rain=df_ln_test['Cummulative_Rain'].to_numpy()
3 df_ln_test_np_fault = df_ln_test['29to35Faults_LN'].to_numpy()
4
5 df_ln_test_rain = df_ln_test_np_rain.reshape(len(df_ln_test_np_rain),1)
6 df_ln_test_fault = df_ln_test_np_fault.reshape(len(df_ln_test_np_fault),1)
```



In [87]:

```
1 # testing model accuracy in test set
2 lr.score(df_ln_test_rain, df_ln_test_fault)
```

Out[87]:

0.16658790305590654



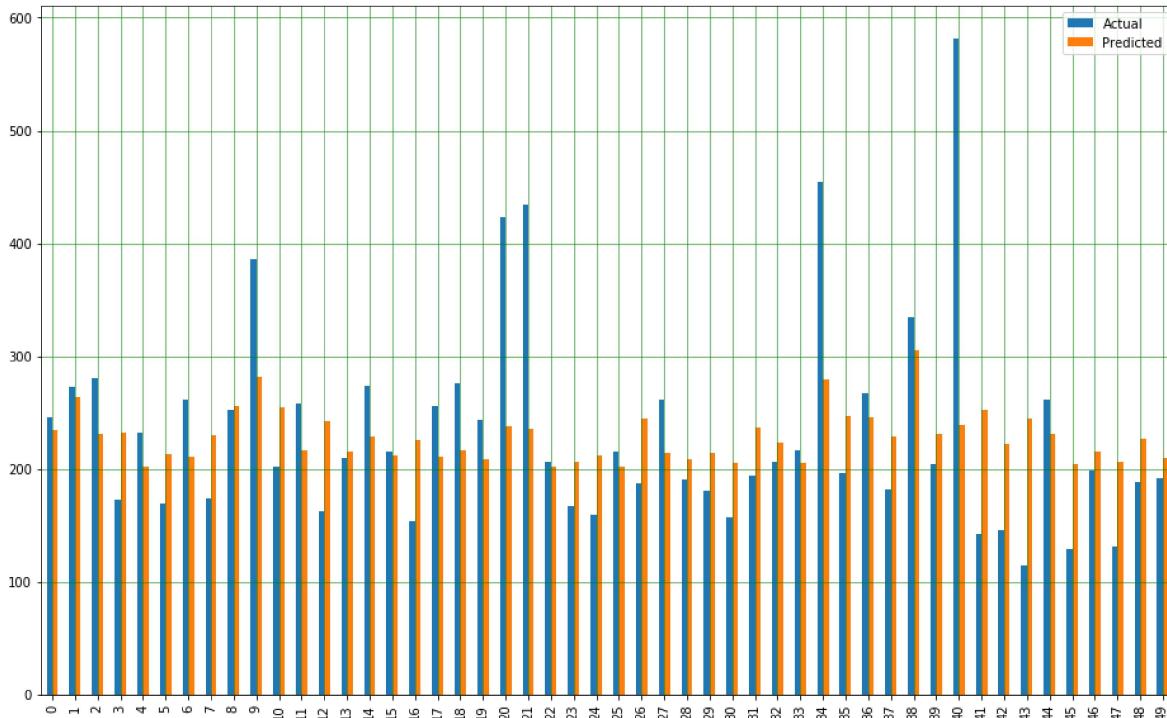


In [88]:

```

1 pred_ln_Fault_linear = lr.predict(df_ln_test_rain)
2 ln_actual_pred_linear= pd.DataFrame({'Actual': df_ln_test_fault.flatten(), 'Predicted'
3 bar_plot_df = ln_actual_pred_linear.head(50)
4 bar_plot_df.plot(kind='bar',figsize=(16,10))
5 plt.grid(which='major', linestyle='-', linewidth='0.5', color='green')
6 plt.grid(which='minor', linestyle=':', linewidth='0.5', color='black')
7 plt.show()

```

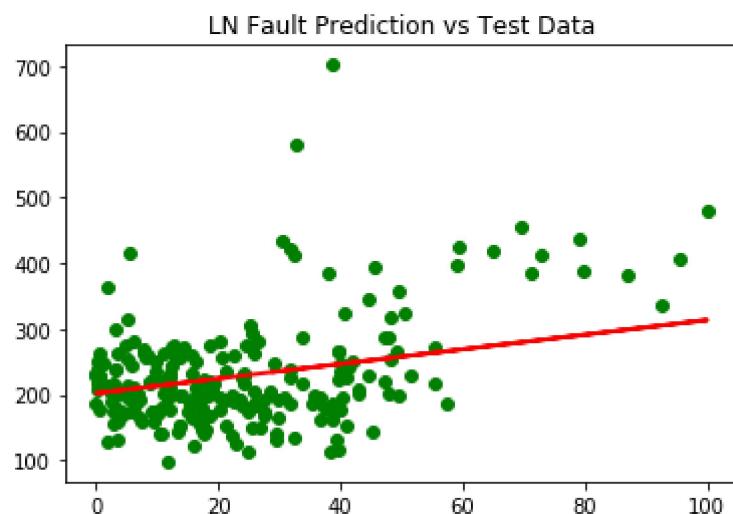


In [89]:

```

1 plt.scatter(df_ln_test_rain, df_ln_test_fault, color='green')
2 plt.plot(df_ln_test_rain, pred_ln_Fault_linear, color='red', linewidth=2)
3 plt.title('LN Fault Prediction vs Test Data')
4 plt.show()

```



In [90]:

```

1 rmspe_ln = (np.sqrt(np.mean(np.square((df_ln_test_fault - pred_ln_Fault_linear) / df_ln
2 print('Root Mean Squared Error Percentage for LN Faults: ', rmspe_ln, '%'))

```

Root Mean Squared Error Percentage for LN Faults: 32.71374764995076 %

## OTHER FAULT MODEL

### Linear Regression

In [91]:

```

1 df_other = df[['Cummulative_Rain', '29to35Faults_OTHER']]
2 df_other_train, df_other_test = train_test_split(df_other, test_size=0.2, random_state=42)

```

In [198]:

```
1 df_other.corr()
```

Out[198]:

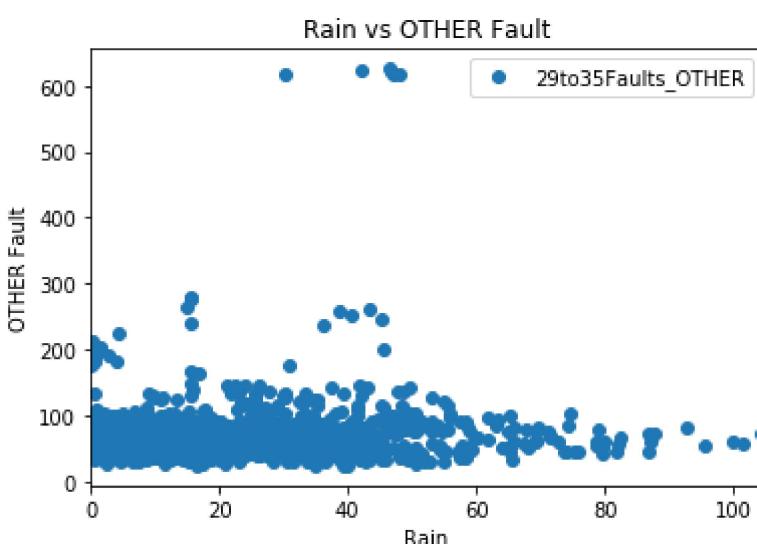
	Cummulative_Rain	29to35Faults_OTHER
Cummulative_Rain	1.000000	0.075699
29to35Faults_OTHER	0.075699	1.000000

In [92]:

```

1 df_other.plot(x='Cummulative_Rain', y='29to35Faults_OTHER', style='o')
2 plt.title('Rain vs OTHER Fault')
3 plt.xlabel('Rain')
4 plt.ylabel('OTHER Fault')
5 plt.show()

```



In [93]:

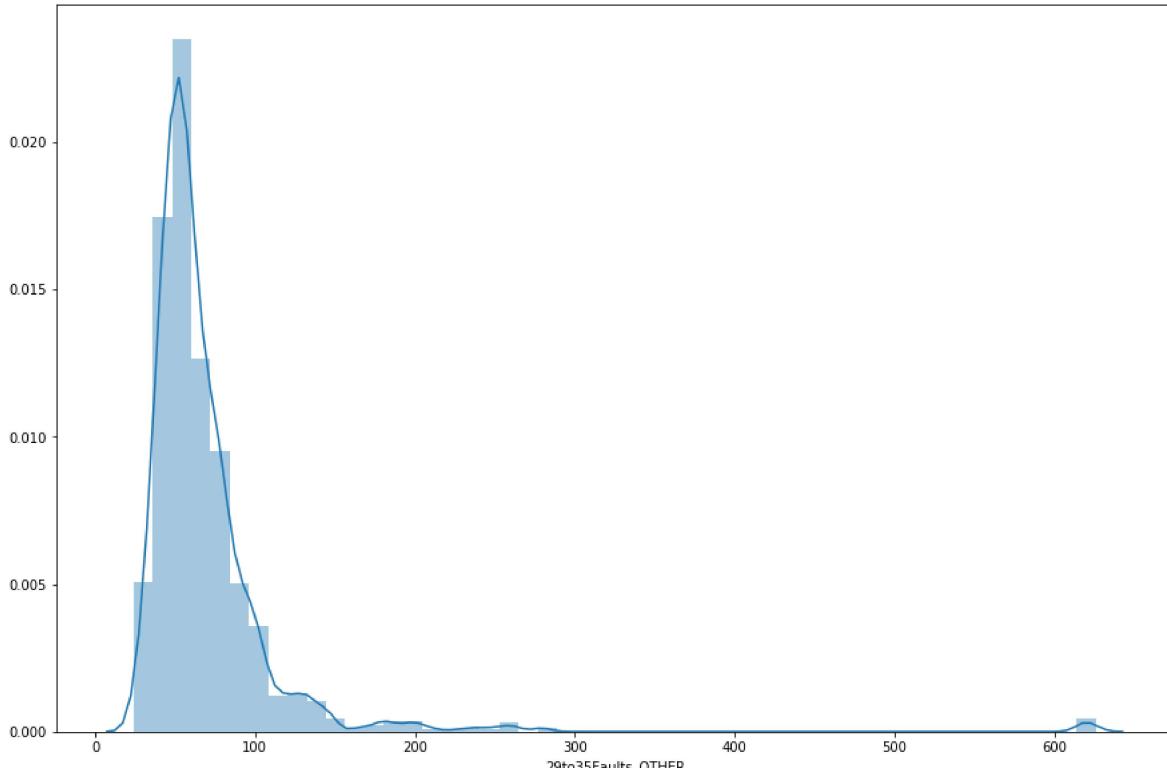
```

1 # Check Average Faults
2 plt.figure(figsize=(15,10))
3 plt.tight_layout()
4 sb.distplot(df_other['29to35Faults_OTHER'])

```

Out[93]:

&lt;matplotlib.axes.\_subplots.AxesSubplot at 0x23d986119b0&gt;



In [94]:

```

1 # Find TRAIN Set Data for 15 to 21 Days for OTHER Faults
2 df_other_train_np_rain=df_other_train['Cummulative_Rain'].to_numpy()
3 df_other_train_np_fault = df_other_train['29to35Faults_OTHER'].to_numpy()
4
5 # Reshape test train data
6 df_other_train_rain = df_other_train_np_rain.reshape(len(df_other_train_np_rain),1)
7 df_other_train_fault = df_other_train_np_fault.reshape(len(df_other_train_np_fault),1)

```

In [95]:

```
1 lr = LinearRegression().fit(df_other_train_rain,df_other_train_fault)
```



In [96]:

```
1 lr.score(df_other_train_rain, df_other_train_fault)
```

Out[96]:

0.003881266621694523



In [97]:

```
1 # Find TEST Set Data for 15 to 21 Days for OTHER Faults
2 df_other_test_np_rain=df_other_test['Cummulative_Rain'].to_numpy()
3 df_other_test_np_fault = df_other_test['29to35Faults_OTHER'].to_numpy()
4
5 df_other_test_rain = df_other_test_np_rain.reshape(len(df_other_test_np_rain),1)
6 df_other_test_fault = df_other_test_np_fault.reshape(len(df_other_test_np_fault),1)
```



In [98]:

```
1 # testing model accuracy in test set
2 lr.score(df_other_test_rain, df_other_test_fault)
```

Out[98]:

0.013061907819305163

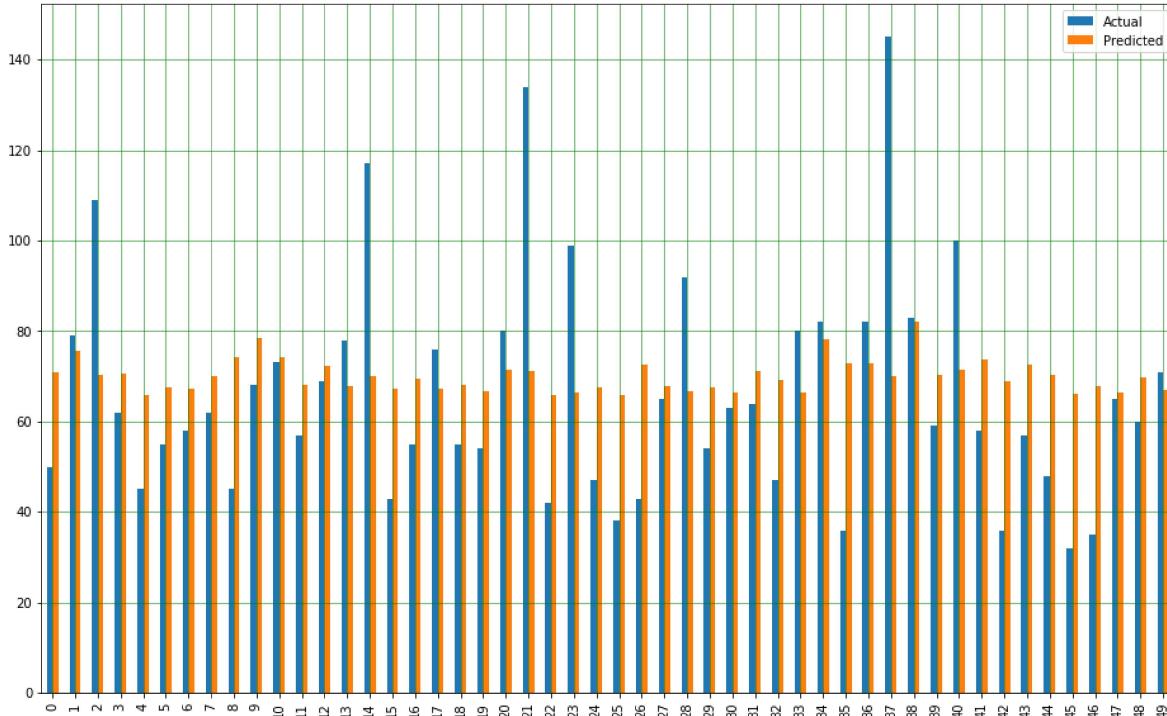


In [99]:

```

1 pred_other_Fault_linear = lr.predict(df_other_test_rain)
2 other_actual_pred_linear= pd.DataFrame({'Actual': df_other_test_fault.flatten(), 'Predicted':pred_other_Fault_linear})
3 bar_plot_df = other_actual_pred_linear.head(50)
4 bar_plot_df.plot(kind='bar',figsize=(16,10))
5 plt.grid(which='major', linestyle='-', linewidth='0.5', color='green')
6 plt.grid(which='minor', linestyle=':', linewidth='0.5', color='black')
7 plt.show()

```



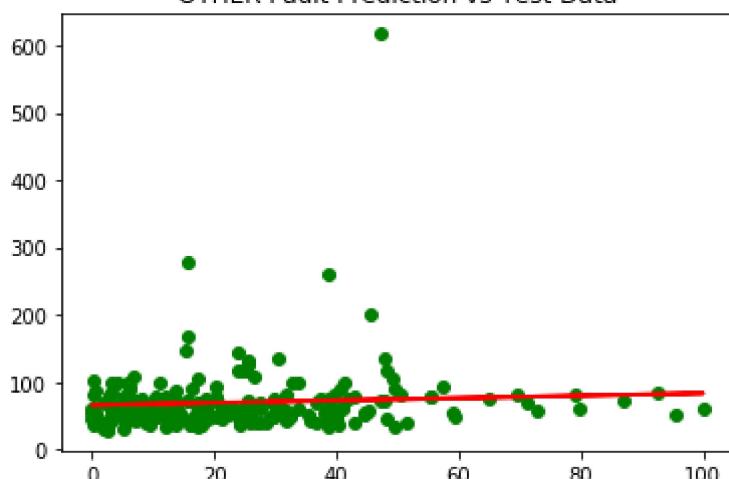
In [100]:

```

1 plt.scatter(df_other_test_rain, df_other_test_fault, color='green')
2 plt.plot(df_other_test_rain, pred_other_Fault_linear, color='red', linewidth=2)
3 plt.title('OTHER Fault Prediction vs Test Data')
4 plt.show()

```

OTHER Fault Prediction vs Test Data





In [101]:

```
1 rmspe_other = (np.sqrt(np.mean(np.square((df_other_test_fault - pred_other_Fault_linear
2 print('Root Mean Squared Error Percentage for Other Faults: ', rmspe_other, '%'))
```

Root Mean Squared Error Percentage for Other Faults: 44.75069998555371 %