**For this dataset, there are 3 insights I found after exploring some features and also target variables:**

1. Imbalance data can be found between binary classes on target variable. After calculation, there are 8404 labelled "0" and 5443 labelled "1" target variable for total of 13847 observations.

2. Using *mutual_info_classif()* method on scikit-learn to evaluating the Information gain of each variable in the context of the target variable or in other words, measures the dependency between the features and the target variable.

   CHP2Temp1(Deg C) 0.5509484000930598

   CHP2Temp2(Deg C) 0.550259620067963

   CHP2Vib1(mm/s) 0.4625149128689099

   CHP2Vib2(mm/s) 0.45978536822018623
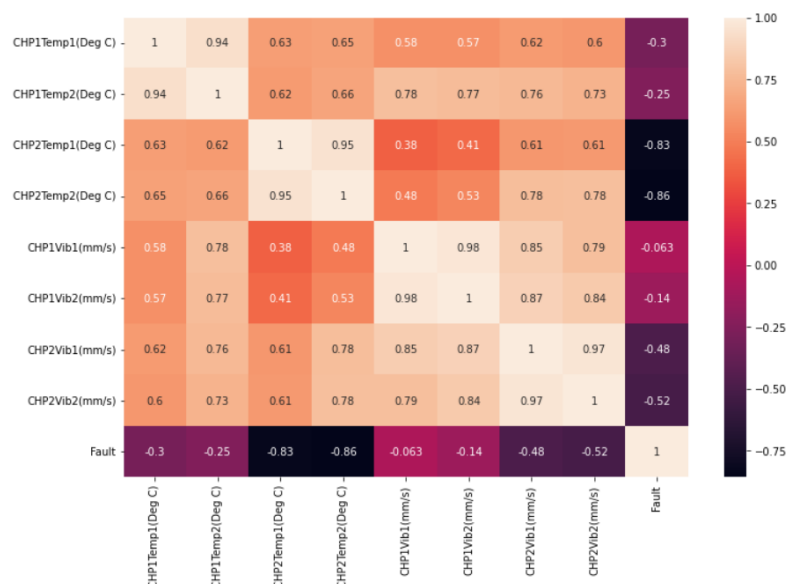
   CHP1Vib2(mm/s) 0.19567828846226898

   CHP1Vib1(mm/s) 0.19285954457513266

   CHP1Temp1(Deg C) 0.1784938851598714

   CHP1Temp2(Deg C) 0.07303022517117852

   As we can see, using this method, target variable ("Fault" label) has higher dependency to Chiller pump CHP2 parameters than CHP1

3. Using Pearson Correlation,



   We can see that both sensor with same parameter on the same Chiller Pump has very high correlation. Also, Temperature sensor on both Chiller pump has high correlation, same case found on Vibration sensor.