

Capstone Project 1

(BY MULA SHASHANK)

ANOMA DATA -Automated Anomaly Detection for Predictive Maintenance

Problem Statement:

Many different industries need predictive maintenance solutions to reduce risks and gain actionable insights through processing data from their equipment. Although system failure is a very general issue that can occur in any machine, predicting the failure and taking steps to prevent such failure is most important for any machine or software application. Predictive maintenance evaluates the condition of equipment by performing online monitoring. The goal is to perform maintenance before the equipment degrades or breaks down. This Capstone project is aimed at predicting the machine breakdown by identifying the anomalies in the data. The data we have contains about 18000+ rows collected over few days. The column 'y' contains the binary labels, with 1 denoting there is an anomaly. The rest of the columns are predictors.

Loading the required libraries

In [1]:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
import seaborn as sns
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.preprocessing import StandardScaler, LabelEncoder
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
```

In [2]:

```
# Loading the dataset
data = pd.read_excel('AnomaData.xlsx')
# Displaying the first few rows of the dataset
data.head()
```

Out[2]
5 rows * 62 columns

In[3]

data.columns

Out[3]:

```
Index(['time', 'y', 'x1', 'x2', 'x3', 'x4', 'x5', 'x6', 'x7', 'x8', 'x9',
```

```

        'x10', 'x11', 'x12', 'x13', 'x14', 'x15', 'x16', 'x17', 'x18', 'x19'
    ,
        'x20', 'x21', 'x22', 'x23', 'x24', 'x25', 'x26', 'x27', 'x28', 'x29'
    ,
        'x30', 'x31', 'x32', 'x33', 'x34', 'x35', 'x36', 'x37', 'x38', 'x39'
    ,
        'x40', 'x41', 'x42', 'x43', 'x44', 'x45', 'x46', 'x47', 'x48', 'x49'
    ,
        'x50', 'x51', 'x52', 'x54', 'x55', 'x56', 'x57', 'x58', 'x59', 'x60'
    ,
        'y.1'],
    dtype='object')

```

In [4]:

```
data.drop(['y.1', 'time'], axis=1, inplace=True)
```

Exploratory Data Analysis

In []:

```

from ydata_profiling import ProfileReport
rep=ProfileReport(data)
rep

Summarize dataset:   0%|          | 0/5 [00:00<?, ?it/s]

```

In [7]:

```

# Checking null values
data.isna().sum()

```

Out[7]:

```

y      0
x1      0
x2      0
x3      0
x4      0
x5      0
x6      0
x7      0
x8      0
x9      0
x10     0
x11     0
x12     0
x13     0
x14     0
x15     0
x16     0
x17     0
x18     0
x19     0
x20     0
x21     0
x22     0

```

No null values present

Out[8]:

[illegible]

	y	x 1	x 2	x 3	x 4	x 5	x 6	x 7	x 8	x 9	..	x 5 0	x 5 1	x 5 2	x 5 4	x 5 5	x 5 6	x 5 7	x 5 8	x 5 9	x 6 0
	0	0	0	0	0	0	0	0	0	0		0	0	0	0	0	0	0	0	0	0
	0	0	0	0	0	0	0	0	0	0		0	0	0	0	0	0	0	0	0	0
	0	0	0	0	0	0	0	0	0	0		0	0	0	0	0	0	0	0	0	0
	0	0	0	0	0	0	0	0	0	0		0	0	0	0	0	0	0	0	0	0
m e a n	0.	0.	0.	0.	-	0.	2.	0.	-	-		0.	-	0.	0.	2.	9.	0.	-	-	0.
	0	0	1	5	9.	0	3	0	0.	0.		6	3.	3	1	3	2	0.	0.	0	0
	0	1	5	6	9	0	8	0	0	0	.	0	3	8	7	7	3	0	0	6	0
	6	1	7	9	5	6	7	1	0	0	.	2	5	0	3	9	4	0	1	5	1
	7	8	9	3	8	5	5	6	4	3	.	5	7	5	7	1	9	4	8	2	2
	4	2	8	0	3	1	3	4	2	5		5	3	1	0	5	5	9	6	5	5
	0	4	6	0	4	8	3	7	5	6		3	9	9	8	4	3	3	1	2	8
s t d	0.	0.	4.	5.	1	0.	3	0.	0.	0.		6.	3	6.	3.	6	8	2.	0.	1	0.
	0	7	9	9	3	6	7.	1	0	1		4	8.	2	0	7.	1.	3	0	0.	0
	8	4	3	3	0	3	1	0	7	5	.	5	2	1	2	9	2	2	4	3	0
	1	2	9	7	3	4	0	8	5	6	.	4	5	1	9	4	7	6	8	9	4
	8	8	7	1	3	0	4	8	4	0	.	1	6	5	5	0	4	8	7	0	7
	2	7	6	7	7	5	1	7	6	4		5	7	9	1	9	1	3	3	8	2
	2	5	2	8	1	4	2	0	0	7		6	1	8	6	4	3	8	2	5	1
					2								6								
m i n	0.	-	-	-	-	-	-	-	-	-		-	-	-	-	-	-	-	-	-	-
	0	3.	7.	8.	3	1.	7	0.	0.	0.		2	3	1	8.	2	2	1	0.	0	0.
	0	7	3	1	2.	6	9.	4	4	1	.	3.	6	7.	2	0.	9.	2.	1	0.	0
	0	8	1	9	7	2	4	2	5	2	.	4	5	9	4	5	0	6	4	8	1
	0	7	6	8	8	3	0	9	1	0	.	8	8	3	3	7	3	0	9	1	2
	0	2	5	5	1	9	8	2	1	0	.	9	9	4	4	0	9	4	7	0	2
	0	7	5	0	6	8	4	7	4	8		8	0	7	7	0	5	9	5	2	2
		9	0	9	1	8	0	3	1	7		5	0	0	0	0	0	0	0	0	9
					0								0								
2 5 %	0.	-	-	-	-	-	-	-	-	-		-	2	-	0.	-	-	-	0.	0.	-
	0	4	1	5	1	4	4.	0	0	0	.	3.	9.	3.	4	4	4	1.	0	2	0
	0	0	5	3	1.	6	3	5	5	9	.	8	8	6	8	0	5	5	0	9	0
	0	5	8	7	3	4	4	8	1	9	.	3	4	7	7	5	1	9	5	0	1
	0	6	2	0	7	7	5	5	0	9		6	6	6		0	9	8	4	0	8

	y	x 1	x 2	x 3	x 4	x 5	x 6	x 7	x 8	x 9	.	x 5 0	x 5 1	x 5 2	x 5 4	x 5 5	x 5 6	x 5 7	x 5 8	x 5 9	x 6 0
	0 0	8 1	3 5	5 4	8 3 7 2	8 7	2 6 8	2 0	4 3	6 6		7 8	2 4	8 4	8 0	0 4	1 4	0 4	7 0	2 3	0 5
5 0 %	0 0 0 0 0 0	0 1 2 8 2 4 5	- 0 0 7 5 0 5	- 0 1 9 0 6 8 3	- 1 4 8 8 1 5 8	- 0 1 2 0 7 4 5	1 0 5 2 8 4 3 5	- 0 0 0 9 3 3 8	- 0 0 0 0 9 9 3 7	- 0 0 3 0 5 7		0 0 6 5 6 3	2 9 9 8 4 6 2 4	0 2 9 4 8 6 4 6	0 7 0 2 2 9 9 7	1 7 4 7 1 3 1 7	1 4 3 8 8 0 6 6	0 0 8 5 8 2 8 8	0 0 1 2 8 8 1 1	0 7 3 4 5 9 1 0	0 0 0 7 1 0
7 5 %	0 0 0 0 0 0	0 4 2 1 2 2	2 3 1 9 2 7	3 4 2 1 2 3 4	9 2 1 9 1 3 4	0 3 2 5 1 5 2 4	3 2 7 2 9 7 4	0 0 6 0 5 1 5 6	0 0 3 8 9 8 6 0	0 0 0 1 9 9 0		6 1 4 8 6 1 0	2 9 9 8 4 6 2 4	5 1 0 9 5 4 3 1	2 6 7 5 7 5 1	4 0 9 3 3 6 8 7	6 3 2 9 9 6 8 1	2 2 2 0 1 1 8	0 0 2 0 9 9 1 1	1 2 6 6 5 0 6 7	0 0 0 4 0 8 7
m a x	1 0 0 0 0 0 0	3 0 5 4 1 5 6	1 6 7 4 2 1 5	1 5 9 0 1 1 6	3 3 4 6 9 3 0 8 8	4 2 3 9 3 5	9 6 0 6 0 7 8	1 7 0 5 5 9 0	0 7 8 8 0 2 6	4 0 6 8 0 3 3		1 7 8 2 8 4 7	4 0 1 5 2 3 4 8	1 4 1 8 0 5 8 8	6 6 3 7 2 6 5	2 8 7 2 5 0 1 5	2 5 1 4 7 4 8 5	6 9 2 2 0 4 9	0 0 6 7 2 4 0 9	6 9 8 5 4 6 0	0 0 2 0 5 1 0

8 rows × 60 columns

```
# Separate features and target variables
X = data.drop('y',axis=1)
y = data['y'] # Binary target variable for anomaly detection
```

In [10]:

```
# Checking the skewness of the columns
from scipy.stats import skew
```

```
for col in X.columns:  
    print(col+" --> Skew=",skew(data[col]))
```

```
x1 --> Skew= -1.247625732759167  
x2 --> Skew= 0.24618128774105086  
x3 --> Skew= 0.4055454880638477  
x4 --> Skew= 0.12514637378313767  
x5 --> Skew= 0.6837495582949324  
x6 --> Skew= -0.7783583177473792  
x7 --> Skew= 2.4357311339796475  
x8 --> Skew= -1.6026471752443343  
x9 --> Skew= 9.012488007639806  
x10 --> Skew= 8.943562386353957  
x11 --> Skew= -7.185593183352757  
x12 --> Skew= -8.361297553646873  
x13 --> Skew= -0.12379559435046963  
x14 --> Skew= 10.143169032655202  
x15 --> Skew= 7.49124549847556  
x16 --> Skew= -7.889659645806755  
x17 --> Skew= 0.8767548887172883  
x18 --> Skew= 0.6193348339124728  
x19 --> Skew= -14.642912582263218  
x20 --> Skew= -1.2036063560137629  
x21 --> Skew= 0.25536084410168486  
x22 --> Skew= 0.5803436878881983  
x23 --> Skew= -3.1418220516846533  
x24 --> Skew= 0.554933106032238  
x25 --> Skew= -9.498804786346405  
x26 --> Skew= 0.7062128721051979  
x27 --> Skew= -1.09724465412704  
x28 --> Skew= 0.8690546825826609  
x29 --> Skew= -0.3144094618837163  
x30 --> Skew= -0.38429639300540624  
x31 --> Skew= -0.06350061244827468  
x32 --> Skew= -9.61533865007797  
x33 --> Skew= 0.27416672685782983  
x34 --> Skew= -0.7427083382729173  
x35 --> Skew= -0.05333477661865808  
x36 --> Skew= 4.796638845977901  
x37 --> Skew= -2.01988817570014  
x38 --> Skew= 9.851534132930421  
x39 --> Skew= 3.7708045847451492  
x40 --> Skew= -4.413442472540411  
x41 --> Skew= 0.29427380773473477  
x42 --> Skew= 1.51225912751579  
x43 --> Skew= 7.522535603813891  
x44 --> Skew= -0.751908745142232  
x45 --> Skew= 9.14525391875508  
x46 --> Skew= 2.085440032593785  
x47 --> Skew= -0.9776387034518712  
x48 --> Skew= -1.8096800837955427  
x49 --> Skew= -0.41343459474009875  
x50 --> Skew= 0.23970568808172968  
x51 --> Skew= -10.384495189618546  
x52 --> Skew= -1.3077363575331327  
x54 --> Skew= -0.9314628850541723  
x55 --> Skew= -0.46067152341753304
```

```
x56 --> Skew= 0.5252953172596814
x57 --> Skew= 0.20301850665091772
x58 --> Skew= -2.460119223931329
x59 --> Skew= -9.49916162199559
x60 --> Skew= 0.4793085411594726
```

```
# Applying yeo-johnson transformation as there are both negative and positive skew
```

```
from sklearn.preprocessing import PowerTransformer
pt=PowerTransformer(method='yeo-johnson')
```

In [12]:

```
# Performing train test split first to prevent data leakage later
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2,random_state=42,stratify=y)
```

In [13]:

```
print("Shape of X train: ",X_train.shape)
print("Shape of y train: ",y_train.shape)
print("Shape of X test: ",X_test.shape)
print("Shape of y test: ",y_test.shape)
```

Out[13]:

```
Shape of X train:  (14718, 59)
Shape of y train:  (14718,)
Shape of X test:   (3680, 59)
Shape of y test:   (3680,)
```

```
# Transforming the X train and X test
X_train=pt.fit_transform(X_train)
X_test=pt.transform(X_test)
```

In [15]:

```
# Checking the value counts to see if there is imbalance in target column
y_train.value_counts()
```

Out[15]:

```
Y
0    14619
1         99
Name: count, dtype: int64
```

Dataset is imbalanced. So we will apply oversampling method SMOTE

In [16]:

```
# SMOTE oversampling
from imblearn.over_sampling import SMOTE
smote=SMOTE()
X_train,y_train=smote.fit_resample(X_train,y_train)
```

In [17]:

```
y_train.value_counts()
```

Out[17]:

```
Y
0    14619
1    14619
Name: count, dtype: int64
```

Now the dataset is balanced

```
print("Shape of X train: ",X_train.shape)
print("Shape of y train: ",y_train.shape)

Shape of X train:  (29238, 59)
Shape of y train:  (29238,)
```

Featuring Scaling

In [19]:

```
# Scaling the data using Standard Scaler
sc=StandardScaler()
X_train_sc=sc.fit_transform(X_train)
X_test_sc=sc.transform(X_test)
```

In [20]:

```
train_data_processed=pd.concat([pd.DataFrame(X_train_sc,columns=X.columns),
pd.DataFrame(y_train,columns=['y'])],axis=1)
```

```
train_data_processed.to_csv('/content/train_data_processed.csv')
```

In [21]:

```
test_data_processed=pd.concat([pd.DataFrame(X_test_sc,columns=X.columns),
pd.DataFrame(y_test,columns=['y'])],axis=1)
```

In [22]:

```
test_data_processed.to_csv('/content/test_data_processed.csv')
```

In [23]:

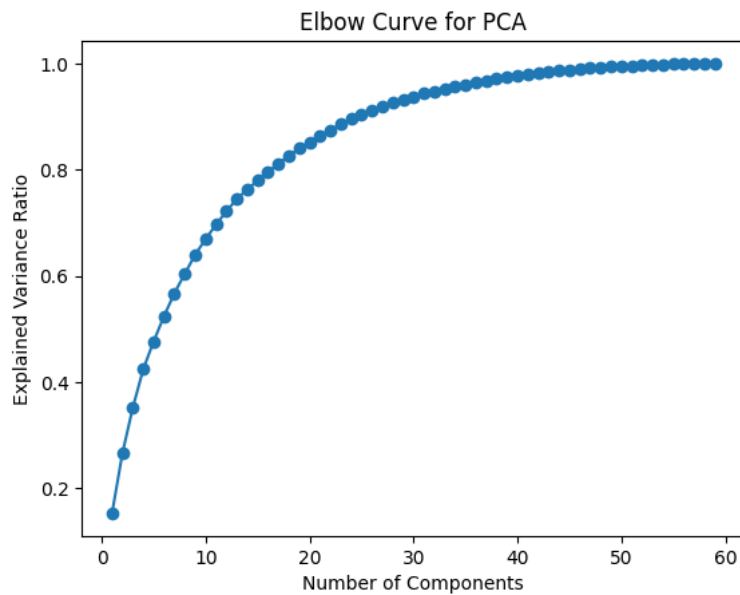
```
# Applying PCA on n number of components to see how many components needed
to
# to achieve more than 90% explained variance ratio
```

```
from sklearn.decomposition import PCA

n_components_range=np.arange(1,X_train_sc.shape[1]+1)
explained_variances=[]

for n_component in n_components_range:
    pca=PCA(n_components=n_component)
    X_train_pca=pca.fit_transform(X_train_sc)
    explained_variance=pca.explained_variance_ratio_
    explained_variances.append(np.sum(pca.explained_variance_ratio_))

plt.plot(n_components_range,explained_variances,marker='o')
plt.xlabel('Number of Components')
plt.ylabel('Explained Variance Ratio')
plt.title('Elbow Curve for PCA')
plt.show()
```

We see that for `n_component=30`, we can reach 90% of the explained variance ratio. So we choose `n_components=30`

In [24]:

```
pca=PCA(n_components=30)
X_train_pca=pca.fit_transform(X_train_sc)
```

In [25]:

```
print("Shape of X train before PCA: ",X_train_sc.shape)
print("Shape of X train after PCA: ",X_train_pca.shape)
```

```
Shape of X train before PCA:  (29238, 59)
Shape of X train after PCA:  (29238, 30)
```

In [26]:

```
X_test_pca=pca.transform(X_test_sc)
```

In [27]:

```
print("Shape of X test before PCA: ",X_test_sc.shape)
print("Shape of X test after PCA: ",X_test_pca.shape)
```

```
Shape of X test before PCA:  (3680, 59)
Shape of X test after PCA:  (3680, 30)
```

In[28]

```
# Predict and evaluate the binary classifier
```

```
def evaluate_model(clf,X_train,X_test):
    clf.fit(X_train,y_train)
    y_pred = clf.predict(X_test)
    print("=====")
    print("Binary Classification - Predictictive Maintenance")
    print("=====")
    print("Accuracy:", accuracy_score(y_test, y_pred))

    print("Classification Report:\n", classification_report(y_test, y_pred,
target_names=['No Anomaly','Anomaly']))
    sns.heatmap(confusion_matrix(y_test, y_pred),annot=True,fmt='g',
        yticklabels=['No Anomaly','Anomaly'],xticklabels=['No
Anomaly','Anomaly'])
```

```
plt.xlabel('Predicted')
plt.ylabel('True')
plt.tight_layout()
plt.show()
```

In[29]

```
from sklearn.linear_model import LogisticRegression
from sklearn.svm import SVC
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score, recall_score, f1_score
```

In [30]:

```
from sklearn.ensemble import AdaBoostClassifier, GradientBoostingClassifier
from xgboost import XGBClassifier
```

In [31]:

```
model_dict={'Support Vector Classifier':SVC(),
            'Random Forest': RandomForestClassifier(),
            'Gradient Boosting Classifier' : GradientBoostingClassifier(),
            'XGBoost':XGBClassifier()
            }
```

In [32]:

```
def model_results(model_dict):
    results=[]
    for modelname,model in model_dict.items():
        model.fit(X_train_pca,y_train)
        acc=accuracy_score(y_test,model.predict(X_test_pca))
        rec=recall_score(y_test,model.predict(X_test_pca))
        f1=f1_score(y_test,model.predict(X_test_pca))
```

```
crossvalscore=np.mean(cross_val_score(model,X_train_pca,y_train,cv=3))
    results.append({'Model Name':modelname,
                    'Accuracy':acc,
                    'Recall':rec,
                    'F1 Score':f1,
                    'Cross Val Score': crossvalscore})
    results_df=pd.DataFrame(results)
    return results_df
```

```
model_results(model_dict)
```

Out[33]:

	Model Name	Accuracy	Recall	F1 Score	Cross Val Score
0	Support Vector Classifier	0.982609	0.36	0.219512	0.990184
1	Random Forest	0.994837	0.32	0.457143	0.998940
2	Gradient Boosting Classifier	0.943750	0.52	0.111588	0.965045

	Model Name	Accuracy	Recall	F1 Score	Cross Val Score
3	XGBoost	0.992935	0.32	0.380952	0.997093

For all the models, recall is extremely low, which means model isn't performing well on test data

In[33]

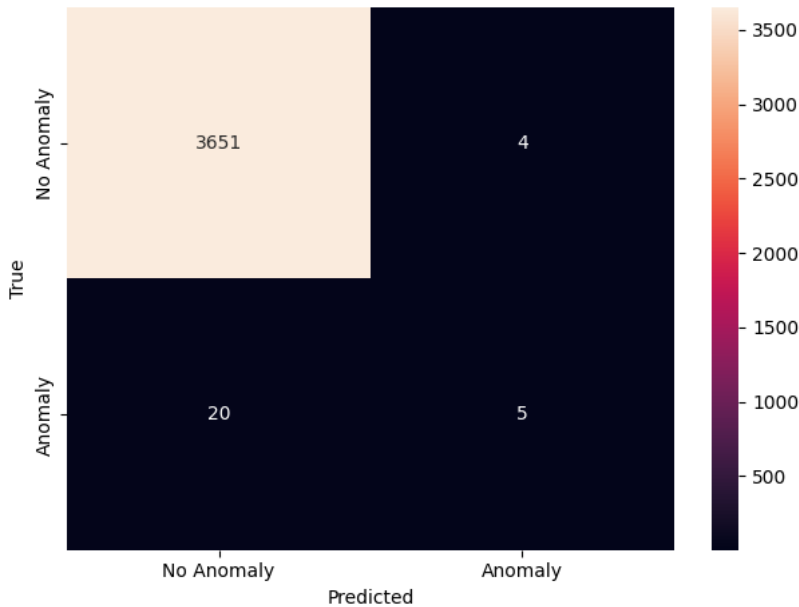
```
rf_clf=RandomForestClassifier()
evaluate_model(rf_clf,X_train_pca,X_test_pca)
```

Out[33]

```
=====
Binary Classification - Predictictive Maintenance
=====
Accuracy: 0.9934782608695653
Classification Report:
              precision    recall  f1-score   support

   No Anomaly       0.99      1.00      1.00      3655
     Anomaly       0.56      0.20      0.29         25

   accuracy              0.99              0.99      3680
  macro avg              0.78              0.60      3680
weighted avg              0.99              0.99      3680
```



In[34]

```
gb_clf=GradientBoostingClassifier()
evaluate_model(gb_clf,X_train_pca,X_test_pca)
```

Out[34]

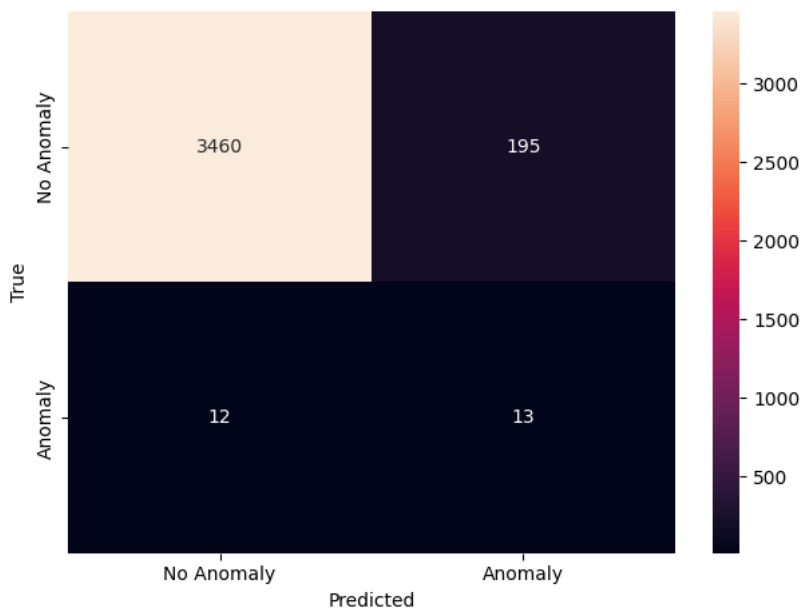
```
GradientBoostingClassifier()
```

```

evaluate_model(gb_clf,X_train_pca,X_test_pca)
=====
Binary Classification - Predictictive Maintenance
=====
Accuracy: 0.94375
Classification Report:

```

	precision	recall	f1-score	support
No Anomaly	1.00	0.95	0.97	3655
Anomaly	0.06	0.52	0.11	25
accuracy			0.94	3680
macro avg	0.53	0.73	0.54	3680
weighted avg	0.99	0.94	0.97	3680



For both the models, we see from the confusion matrix and classification report that the recall is extremely low. So we use a different method instead of PCA.

Using Feature Selection

```

In [35]:
from sklearn.feature_selection import SelectFromModel

In [36]:
# Function to select the best features and show the evaluation metrics for
that best model
def select_from_model(clf):
    sfm=SelectFromModel(clf)
    sfm.fit(X_train_sc, y_train)
    sfm.transform(X_train_sc)
    X_train_sc_df=pd.DataFrame(X_train_sc,columns=X.columns)
    X_train_sc_df=X_train_sc_df[X_train_sc_df.columns[sfm.get_support()]]
    X_test_sc_df=pd.DataFrame(X_test_sc,columns=X.columns)

```

```
X_test_sc_df=X_test_sc_df[X_test_sc_df.columns[sfm.get_support()]]
clf=sfm.estimator_
evaluate_model(clf,X_train_sc_df,X_test_sc_df)
```

In [37]:

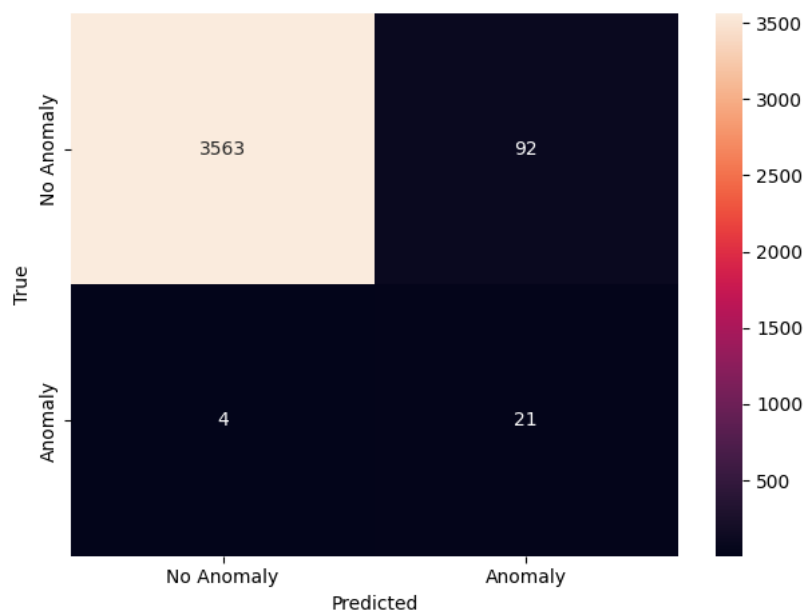
```
gb_clf2=GradientBoostingClassifier()
gb_clf2.fit(X_train_sc,y_train)
select_from_model(gb_clf2)
```

OUT[37]

```
=====
Binary Classification - Predictictive Maintenance
=====
Accuracy: 0.9739130434782609
Classification Report:
              precision    recall  f1-score   support

   No Anomaly         1.00      0.97      0.99      3655
     Anomaly         0.19      0.84      0.30         25

   accuracy                   0.97      3680
  macro avg         0.59      0.91      0.65      3680
 weighted avg         0.99      0.97      0.98      3680
```



In[38]

```
rf_clf2=RandomForestClassifier()
rf_clf2.fit(X_train_sc,y_train)
select_from_model(rf_clf2).
```

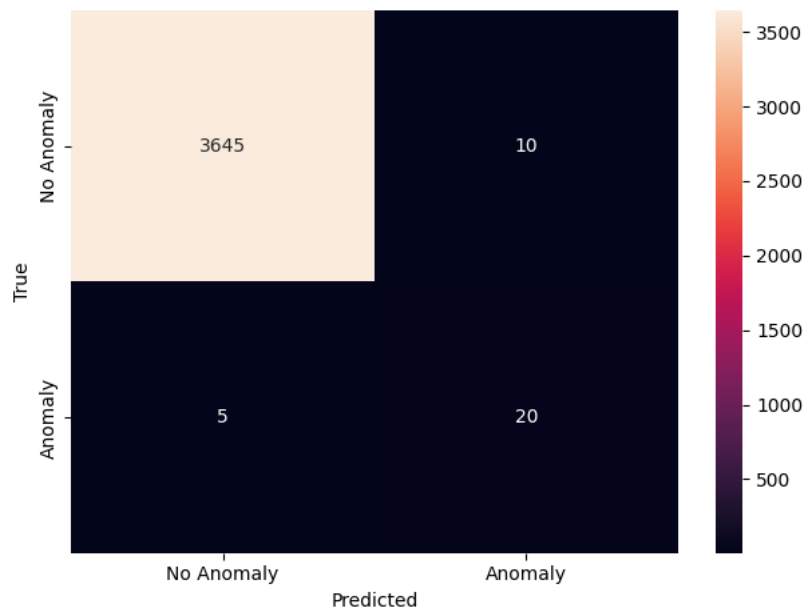
Out [38]

```
=====
Binary Classification - Predictictive Maintenance
=====
Accuracy: 0.9959239130434783
Classification Report:
              precision    recall  f1-score   support

   No Anomaly         1.00      1.00      1.00      3655
     Anomaly         0.67      0.80      0.73         25

   accuracy                   1.00      3680
```

macro avg	0.83	0.90	0.86	3680
weighted avg	1.00	1.00	1.00	3680



In[39]

```
xgb_clf=XGBClassifier()
xgb_clf.fit(X_train_sc,y_train)
select_from_model(xgb_clf).
```

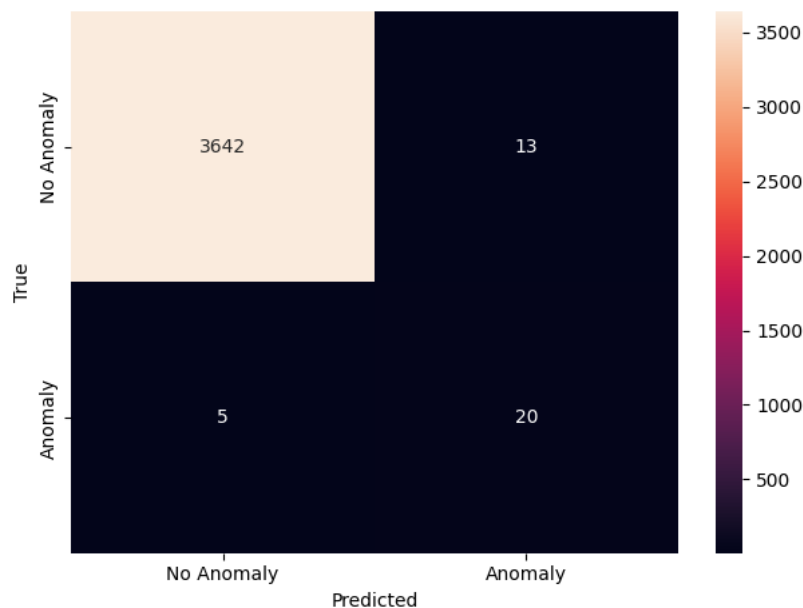
Out[39]

=====
Binary Classification - Predictictive Maintenance
=====

Accuracy: 0.9951086956521739

Classification Report:

	precision	recall	f1-score	support
No Anomaly	1.00	1.00	1.00	3655
Anomaly	0.61	0.80	0.69	25
accuracy			1.00	3680
macro avg	0.80	0.90	0.84	3680
weighted avg	1.00	1.00	1.00	3680



As we can see, the Gradient Boosting Classifier has excellent recall on the 'Anomaly'(1) class, but overall accuracy and F1 score is higher for Random Forest Model. The Random Forest Model is also slightly performing better than the XGBoost model

We see that the test data is highly imbalanced but still model performs very well with feature selection.

Hyperparameter tuning for the Random Forest Model

In[40]

```
# Selecting the best features and building the Random forest model
sfm=SelectFromModel(rf_clf2)
sfm.fit(X_train_sc, y_train)
sfm.transform(X_train_sc)
X_train_sc_df=pd.DataFrame(X_train_sc,columns=X.columns)
X_train_sc_df=X_train_sc_df[X_train_sc_df.columns[sfm.get_support()]]
X_test_sc_df=pd.DataFrame(X_test_sc,columns=X.columns)
X_test_sc_df=X_test_sc_df[X_test_sc_df.columns[sfm.get_support()]]
rf_clf2=sfm.estimator_
```

In [41]:

```
from sklearn.model_selection import GridSearchCV

# Defining the parameter grid
param_grid = {
    'n_estimators': [300,400],
    'max_depth': [None, 10, 20],
    'min_samples_split': [2, 5,10],
    'min_samples_leaf': [1, 2, 4],
}

# Initialize GridSearchCV
```

```

grid_search = GridSearchCV(estimator=rf_clf2, param_grid=param_grid, cv=3,
scoring='accuracy', verbose=2, n_jobs=-1)

# Fit GridSearchCV
grid_search.fit(X_train_sc_df, y_train)

# Get the best parameters
best_params = grid_search.best_params_
print("Best Parameters:", best_params)
Fitting 3 folds for each of 54 candidates, totalling 162 fits
Best Parameters: {'max_depth': None, 'min_samples_leaf': 1, 'min_samples_sp
lit': 2, 'n_estimators': 300}

# The best Random Forest Model
best_rf_clf=grid_search.best_estimator_

evaluate_model(best_rf_clf,X_train_sc_df,X_test_sc_df)

```

In [42]:

In [43]:

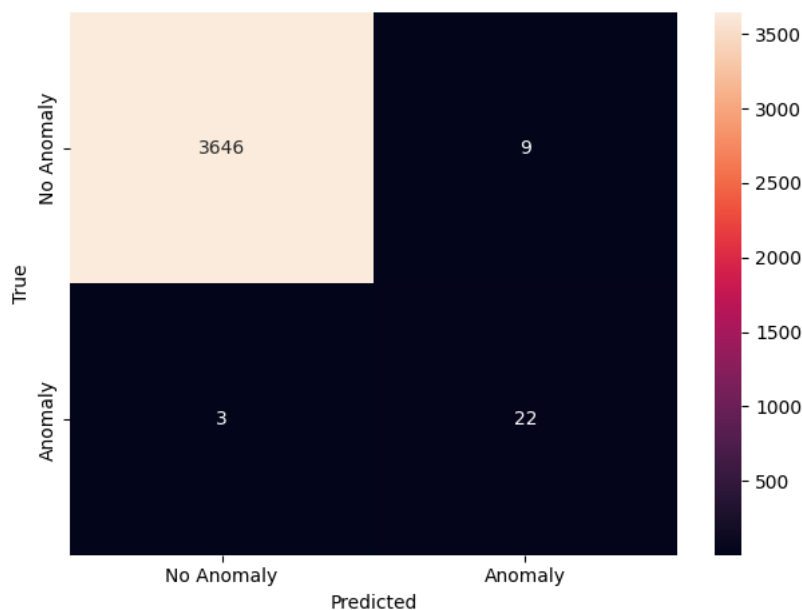
Out[43]

```

=====
Binary Classification - Predictictive Maintenance
=====
Accuracy: 0.9967391304347826
Classification Report:

```

	precision	recall	f1-score	support
No Anomaly	1.00	1.00	1.00	3655
Anomaly	0.71	0.88	0.79	25
accuracy			1.00	3680
macro avg	0.85	0.94	0.89	3680
weighted avg	1.00	1.00	1.00	3680



As we can see that for our tuned Random Forest model, the accuracy, F1 score has increased compared to the untuned model.

Saving the model as a pickle file

```
import pickle
```

In [44]:

```
with open('/content/random_forest_model.pkl', 'wb') as file:
    pickle.dump(best_rf_clf, file)
```

In [45]:

```
from google.colab import files

# Download the file
files.download('/content/random_forest_model.pkl')
```

In [46]:

Using Deep Learning

```
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense, Dropout
```

In [47]:

```
model=Sequential()
model.add(Dense(128,activation='relu',input_shape=(59,)))
model.add(Dropout(0.5))
model.add(Dense(64,activation='relu'))
model.add(Dropout(0.5))
model.add(Dense(1,activation='sigmoid'))

model.compile(optimizer='adam',loss='binary_crossentropy',metrics=['accuracy'])
```

In [48]:

```
model.fit(X_train_sc,y_train,validation_data=(X_test_sc,y_test),batch_size=
32,epochs=20)
```

In [49]:

```
Epoch 1/20
914/914 [=====] - 4s 3ms/step - loss: 0.2843 - acc
uracy: 0.8747 - val_loss: 0.2015 - val_accuracy: 0.9560
Epoch 2/20
914/914 [=====] - 3s 3ms/step - loss: 0.1018 - acc
uracy: 0.9655 - val_loss: 0.3090 - val_accuracy: 0.9878
Epoch 3/20
914/914 [=====] - 4s 4ms/step - loss: 0.0601 - acc
uracy: 0.9811 - val_loss: 0.5282 - val_accuracy: 0.9875
Epoch 4/20
914/914 [=====] - 2s 3ms/step - loss: 0.0442 - acc
uracy: 0.9871 - val_loss: 0.6363 - val_accuracy: 0.9883
```

Epoch 5/20
914/914 [=====] - 2s 2ms/step - loss: 0.0318 - accuracy: 0.9903 - val_loss: 0.8399 - val_accuracy: 0.9913
Epoch 6/20
914/914 [=====] - 2s 3ms/step - loss: 0.0267 - accuracy: 0.9921 - val_loss: 0.9352 - val_accuracy: 0.9908
Epoch 7/20
914/914 [=====] - 3s 3ms/step - loss: 0.0246 - accuracy: 0.9926 - val_loss: 1.0304 - val_accuracy: 0.9918
Epoch 8/20
914/914 [=====] - 4s 4ms/step - loss: 0.0185 - accuracy: 0.9945 - val_loss: 1.1935 - val_accuracy: 0.9867
Epoch 9/20
914/914 [=====] - 3s 3ms/step - loss: 0.0204 - accuracy: 0.9934 - val_loss: 1.1871 - val_accuracy: 0.9927
Epoch 10/20
914/914 [=====] - 3s 3ms/step - loss: 0.0160 - accuracy: 0.9952 - val_loss: 1.3342 - val_accuracy: 0.9932
Epoch 11/20
914/914 [=====] - 3s 3ms/step - loss: 0.0159 - accuracy: 0.9952 - val_loss: 1.5859 - val_accuracy: 0.9935
Epoch 12/20
914/914 [=====] - 3s 3ms/step - loss: 0.0166 - accuracy: 0.9956 - val_loss: 1.5284 - val_accuracy: 0.9927
Epoch 13/20
914/914 [=====] - 4s 4ms/step - loss: 0.0151 - accuracy: 0.9961 - val_loss: 1.6516 - val_accuracy: 0.9932
Epoch 14/20
914/914 [=====] - 3s 3ms/step - loss: 0.0140 - accuracy: 0.9961 - val_loss: 1.6380 - val_accuracy: 0.9940
Epoch 15/20
914/914 [=====] - 2s 3ms/step - loss: 0.0112 - accuracy: 0.9969 - val_loss: 1.8362 - val_accuracy: 0.9937
Epoch 16/20
914/914 [=====] - 2s 3ms/step - loss: 0.0151 - accuracy: 0.9953 - val_loss: 1.8580 - val_accuracy: 0.9924
Epoch 17/20
914/914 [=====] - 2s 3ms/step - loss: 0.0129 - accuracy: 0.9963 - val_loss: 1.9435 - val_accuracy: 0.9929
Epoch 18/20
914/914 [=====] - 4s 4ms/step - loss: 0.0124 - accuracy: 0.9965 - val_loss: 2.0000 - val_accuracy: 0.9940
Epoch 19/20
914/914 [=====] - 2s 2ms/step - loss: 0.0121 - accuracy: 0.9962 - val_loss: 1.9488 - val_accuracy: 0.9929
Epoch 20/20
914/914 [=====] - 2s 3ms/step - loss: 0.0108 - accuracy: 0.9970 - val_loss: 1.8913 - val_accuracy: 0.9935

Out[49]:

<keras.src.callbacks.History at 0x7c809023a6b0>

In [50]:

```
loss,accuracy=model.evaluate(X_test_sc,y_test)
115/115 [=====] - 0s 2ms/step - loss: 2.0330 - acc
uracy: 0.9937
```

In [51]:

```
threshold = 0.5
y_pred_prob=model.predict(X_test_sc)
y_pred_classes = (y_pred_prob >= threshold).astype(int)
115/115 [=====] - 0s 1ms/step
```

In [52]:

```
print("Predicted probabilities:", y_pred_prob[:10].flatten())
print("Predicted classes:", y_pred_classes[:10].flatten())
```

Out[52]

```
Predicted probabilities: [1.2075962e-25  6.9271626e-25  1.4649272e-07  8.66433
89e-08  6.5011552e-09
 1.6045803e-30  0.0000000e+00  2.7664161e-20  2.7629006e-15  2.2056058e-19]
Predicted classes: [0 0 0 0 0 0 0 0 0 0]
```

In [53]:

```
print("=====")
print("Binary Classification - Predictictive Maintenance")
print("=====")
print("Accuracy:", accuracy_score(y_test, y_pred_classes))

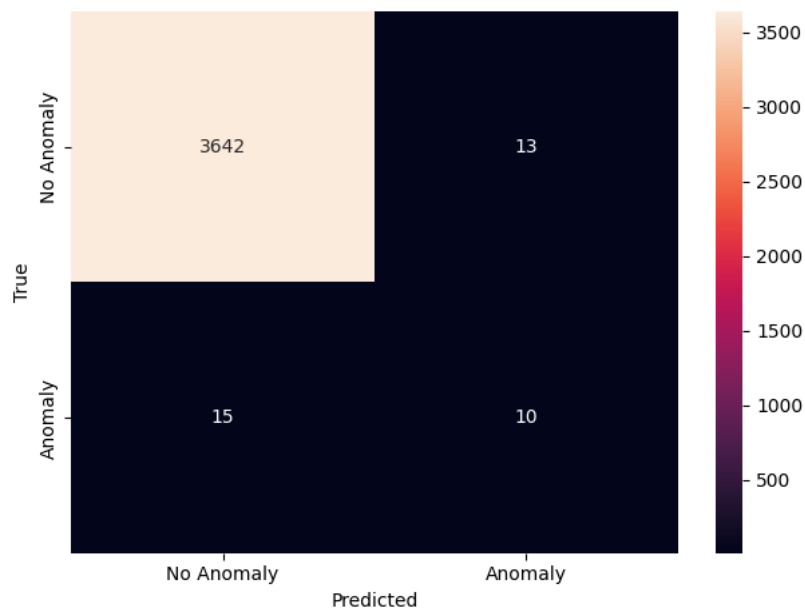
print("Classification Report:\n", classification_report(y_test,
y_pred_classes,
                                                         target_names=['No
Anomaly','Anomaly']))
sns.heatmap(confusion_matrix(y_test, y_pred_classes),annot=True,fmt='g',
            yticklabels=['No Anomaly','Anomaly'],xticklabels=['No
Anomaly','Anomaly'])
plt.xlabel('Predicted')
plt.ylabel('True')
plt.tight_layout()
plt.show()
```

Out[53]

```
=====
Binary Classification - Predictictive Maintenance
=====
Accuracy: 0.9923913043478261
Classification Report:

```

	precision	recall	f1-score	support
No Anomaly	1.00	1.00	1.00	3655
Anomaly	0.43	0.40	0.42	25
accuracy			0.99	3680
macro avg	0.72	0.70	0.71	3680
weighted avg	0.99	0.99	0.99	3680



Deep learning model not performing as good as our Random Forest model.