

Device Failure Prediction using Machine Learning

Table of Contents

1. Importing the Libraries
 - A. Dataset Overview
 - B. Exploratory Data Analysis
 - C. Data Preprocessing
 - D. Model Building
 1. Logistic Regression
 2. Random Forest Classifier
 3. Extreme Gradient Boosting Classifier
 4. K Nearest Neighbors Classifier
 5. Gradient Boosting Classifier
 - E. Model Saving
 - F. Model Evaluation
 - G. Voting Classifier
 - H. Conclusion

1. Importing the Libraries

```
# Importing the libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from collections import Counter
from sklearn.pipeline import Pipeline
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import StandardScaler, MinMaxScaler
from sklearn.model_selection import train_test_split, StratifiedShuffleSplit, cross_val_score
from sklearn.metrics import accuracy_score, roc_auc_score, confusion_matrix, classification_report, f1_score
from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier, VotingClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import GridSearchCV
from xgboost import XGBClassifier
from sklearn.neighbors import KNeighborsClassifier
from imblearn.over_sampling import SMOTE
from imblearn.under_sampling import RandomUnderSampler
import warnings
warnings.filterwarnings("ignore")
plt.style.use("ggplot")
```

```
# Creating a config class to store all the configurations
class config:
    # training dataset path
    DATASET_PATH = "../data/predict_failure.csv"

    # Input parameters for data preprocessing
    TARGET = "failure"
    TEST_SIZE = 0.2
    SEED = 42
    N_SPLITS = 10

    # Input parameters for data sampling techniques
    SMOTE_SAMPLING_STRATEGY = 0.1
    UNDER_SAMPLING_STRATEGY = 0.75

    # Features to be dropped
    DROP_ATTRIBUTES = ["date", "device", "attribute7"]

    # List of features to be scaled based on Standardization/Normalization
    STANDARD_SCALAR_ATTRIBUTES = ["attribute1", "attribute3", "attribute4", "attribute5", "attribute8", "attribute9"]
    MIN_MAX_SCALAR_ATTRIBUTES = ["attribute2", "attribute6"]
```

2. Dataset Overview

```
# Importing the dataset.
failure_data = pd.read_csv("../data/predict_failure.csv")
failure_data.head()
```

	date	device	failure	attribute1	attribute2	attribute3	attribute4	attribute5	attribute6	attribute7	attribute8	attribute9	
0	2015-01-01	S1F0F80	0	215630672		56	0	52	6	407438	0	0	7
1	2015-01-01	S1F0F68	0	61370680		0	3	0	6	403174	0	0	0
2	2015-01-01	S1F0F69	0	173295968		0	0	0	12	237394	0	0	0
3	2015-01-01	S1F0F85	0	79694024		0	0	0	6	401186	0	0	0
4	2015-01-01	S1F0F828	0	135974080		0	0	0	15	313173	0	0	3

```
# Shape of the dataset.
failure_data.shape

(124494, 13)
```

```
# Checking the dataset info.
failure_data.info()
```

```
Out[3]:
Out[4]:
```

```
Out[5]:
```

```
# Changing the date column type to datetime.
failure_data["date"] = pd.to_datetime(failure_data["date"], format="%Y-%m-%d")
```

```
# Checking for missing values
print(failure_data.isnull().sum()/len(failure_data)*100)
```

	failure	attribute1	attribute2	attribute3	attribute4	attribute5	attribute6	attribute7	attribute8	attribute9
count	124494.000000	124494.000000	124494.000000	124494.000000	124494.000000	124494.000000	124494.000000	124494.000000	124494.000000	124494.000000
mean	0.000000	1.223866e+08	1.59494762	9.940455	1.741120	14.222693	260172.858025	0.292528	0.292528	0.292528
std	0.029167	7.045900e+07	2.171965770	175.747321	22.906507	15.343021	99151.009852	7.436924	7.436924	7.436924
min	0.000000	0.000000e+00	0.000000	0.000000	0.000000	1.000000	8.000000	0.000000	0.000000	0.000000
25%	0.000000	6.12675e+07	0.000000	0.000000	0.000000	8.000000	221452.000000	0.000000	0.000000	0.000000
50%	0.000000	1.227957e+08	0.000000	0.000000	0.000000	10.000000	249799.000000	0.000000	0.000000	0.000000
75%	0.000000	1.833064e+08	0.000000	0.000000	0.000000	12.000000	310266.000000	0.000000	0.000000	0.000000
max	1.000000	2.441405e+08	64968.000000	24929.000000	1666.000000	98.000000	689161.000000	832.000000	832.000000	832.000000

```
# Checking for class imbalance based on target variable failure.
failure_data["failure"].value_counts()
```

failure	count
0	124388
1	106

Name: failure, dtype: int64

3. Exploratory Data Analysis

```
# Checking the unique number of devices in the dataset.
failure_data["device"].unique()
```

Out[9]: 1168

```
# Checking the most frequent devices in the dataset.
failure_data["device"].value_counts().head(10)
```

device	count
S1F0F80	304
W1F0S372	304
S1F0F828	304
S1F0F804	304
W1F0F702	304
W1F0X20L	304
W1F0S87	304
S1F0S38P	304
W1F0S977	304
Name: device, dtype: int64	

```
# Printing the correlation between attributes.
failure_data.corr()
```

	failure	attribute1	attribute2	attribute3	attribute4	attribute5	attribute6	attribute7	attribute8	attribute9
failure	1.000000	0.001984	0.052902	-0.000948	0.067398	0.002270	-0.000550	0.119055	0.119055	0.001622
attribute1	0.001984	1.000000	-0.004248	0.003702	0.001837	-0.003370	-0.001516	0.000151	0.000151	0.001226
attribute2	0.052902	-0.004248	1.000000	-0.002670	0.146593	-0.013999	-0.026350	0.141367	0.141367	-0.002736
attribute3	-0.000948	0.003702	-0.002670	1.000000	0.097452	-0.006696	0.009027	-0.001884	-0.001884	0.532366
attribute4	0.067398	0.001837	0.146593	0.097452	1.000000	0.009773	0.024870	0.045631	0.045631	0.036069
attribute5	0.002270	-0.003370	-0.013999	-0.006696	-0.009773	1.000000	-0.017051	-0.009384	-0.009384	0.005949
attribute6	-0.000550	-0.001516	-0.026350	0.009027	0.024870	-0.017051	1.000000	-0.012207	-0.012207	0.021152
attribute7	0.119055	0.000151	0.141367	-0.001884	0.045631	-0.009384	-0.012207	1.000000	1.000000	0.006861
attribute8	0.119055	0.000151	0.141367	-0.001884	0.045631	-0.009384	-0.012207	1.000000	1.000000	0.006861
attribute9	0.001622	0.001122	-0.002736	0.532366	0.036069	0.005949	0.021152	0.006861	0.006861	1.000000

```
# Visualizing the correlations between attributes using heatmap.
plt.figure(figsize=(12, 8))
sns.heatmap(failure_data.corr(), cmap="magma", annot=True, vmin=0, vmax=0.2, linewidth=1)
plt.title("Heatmap for correlations between attributes")
plt.show()
```



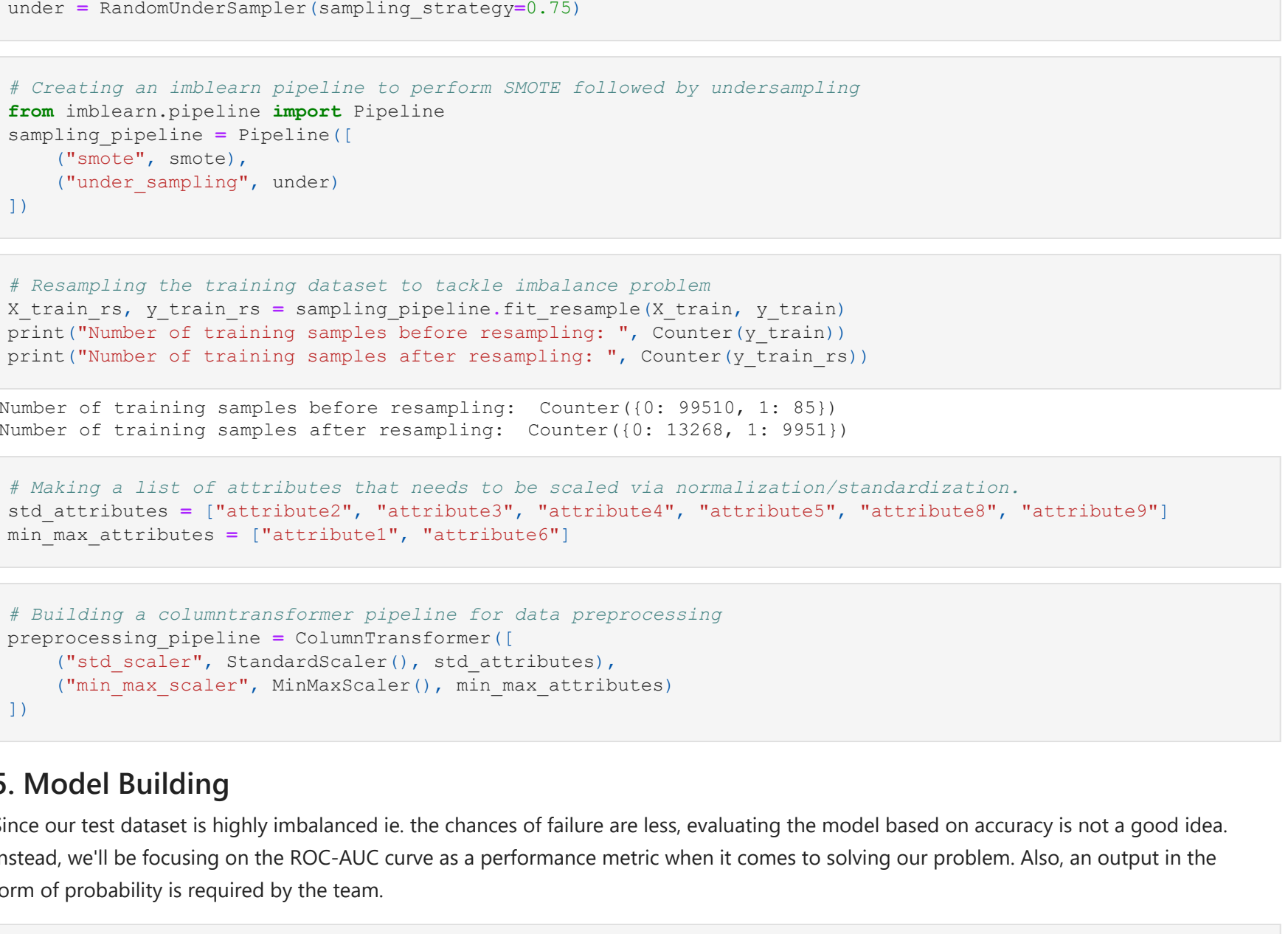
We can see from the correlation chart that attributes 7 and attribute 8 have a correlation of 1 which means they are the same. Also, there's not much relation between the rest of them.

```
# Looks like attribute 7 and attribute 8 are similar. Let's confirm the same.
failure_data["attribute7"].equals(failure_data["attribute8"])
```

Out[13]: True

```
# Creating a separate dataframe for records when the device has failed.
failure_data[failure_data["failure"] == 1]
```

```
# Plotting a distplot to see the distribution of attribute 1 and attribute 6.
fig, axes = plt.subplots(2, 2, figsize=(12, 8), sharey=False, sharex=False)
fig.suptitle("Frequency distribution of other attributes when the device failure status = 1", fontsize=15)
sns.distplot(axes[0, 0], data=failure_data[failure_data["failure"] == 1]["attribute1"], color="green")
axes[0, 0].set_title("Attribute 1")
sns.distplot(axes[0, 1], data=failure_data[failure_data["failure"] == 1]["attribute6"], color="red")
axes[0, 1].set_title("Attribute 6")
sns.distplot(axes[1, 0], data=failure_data[failure_data["failure"] == 1]["attribute7"], color="blue")
axes[1, 0].set_title("Attribute 7")
sns.distplot(axes[1, 1], data=failure_data[failure_data["failure"] == 1]["attribute8"], color="orange")
axes[1, 1].set_title("Attribute 8")
fig.tight_layout()
plt.show()
```



```
# Checking probability distributions for attributes 2, 3, 4, 5, 8 and 9 when failure is 1.
fig, axes = plt.subplots(2, 2, figsize=(15, 10), sharey=False, sharex=False)
fig.suptitle("Frequency distribution of other attributes when the device failure status = 1", fontsize=15)
sns.histplot(axes[0, 0], data=failure_data[failure_data["failure"] == 1], x="attribute2", kde=True, bins=15, color="yellow")
axes[0, 0].set_title("Attribute 2")
sns.histplot(axes[0, 1], data=failure_data[failure_data["failure"] == 1], x="attribute3", kde=True, bins=15, color="violet")
axes[0, 1].set_title("Attribute 3")
sns.histplot(axes[1, 0], data=failure_data[failure_data["failure"] == 1], x="attribute4", kde=True, bins=15, color="red")
axes[1, 0].set_title("Attribute 4")
sns.histplot(axes[1, 1], data=failure_data[failure_data["failure"] == 1], x="attribute5", kde=True, bins=15, color="green")
axes[1, 1].set_title("Attribute 5")
sns.histplot(axes[2, 0], data=failure_data[failure_data["failure"] == 1], x="attribute8", kde=True, bins=15, color="orange")
axes[2, 0].set_title("Attribute 8")
sns.histplot(axes[2, 1], data=failure_data[failure_data["failure"] == 1], x="attribute9", kde=True, bins=15, color="blue")
axes[2, 1].set_title("Attribute 9")
fig.tight_layout()
plt.show()
```



Looks like only the attribute1 and attribute6 from our given dataset are normally distributed. But the rest of the attributes don't have a valid distribution as the data is skewed. But if we take a look at the boxplot of these attributes to check for outliers, there are just too many for them to be removed. They could be useful features for detecting device failures, so we are just going to keep them all. Also, since all the attributes are of different scales, we can fix that by performing Standardization and Normalization.

```
# Let's check out the date column in more detail.
failure_data["date"].nunique()
```

Out[16]: 304

```
# Checking the earliest and latest date available in our dataset.
print("Earliest date: ", failure_data["date"].min())
print("Latest date: ", failure_data["date"].max())
```

Earliest date: 2015-01-01 00:00:00
Latest date: 2015-11-02 00:00:00

```
# Assuming we have records of device from the day it was set up, creating a device uptime (in days) column.
failure_data["device_uptime"] = failure_data.groupby("device")["date"].rank(method="dense")
failure_data["device_uptime"] = failure_data["device_uptime"] - 1
```

```
# Checking the frequency of device failures based on their uptime
plt.figure(figsize=(10, 6))
sns.lineplot(data=failure_data[failure_data["failure"] == 1], x="device_uptime", y="failure", marker="o", lines=plt.show())
```

```
# Looks like there is not much of correlation between the uptime and failure of devices.
# Also, assuming here that devices are of the same type.
# Dropping the non-required columns from our dataset.
failure_data.drop(["attribute7", "device", "device_uptime", "date"], axis=1, inplace=True)
failure_data.head()
```

	failure	attribute1	attribute2	attribute3	attribute4	attribute5	attribute6	attribute8	attribute9
0	0	215630672	56	0	52	6	407438	0	7
1163	0	1650664	56	0	52	6	407438	0	7
2326	0	124017388	56	0	52	6	407438	0	7
3489	0	128073224	56	0	52	6	407439	0	7
4651	0	97393448	56	0	52	6	408114	0	7

4. Data Preprocessing

```
# Creating array of independent variables and target.
X = failure_data.drop("failure", axis=1)
y = failure_data["failure"]
print("Shape of X: ", X.shape)
print("Shape of y: ", y.shape)

Shape of X: (124494, 8)
Shape of y: (124494, 1)
```

```
# Splitting the dataset into training and testing sets.
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42, shuffle=True, stratify=y)
print(Counter(y_train))
print(Counter(y_test))
```

Counter((0: 99510, 1: 85))
Counter((0: 24878, 1: 21))

Our dataset is severely imbalanced and to tackle this issue, we'll be using SMOTE for oversampling followed by random undersampling to get the right amount of data for training our model.

```
# Creating a SMOTE based sampler.
smote = SMOTE(sampling_strategy=0.1)
under = RandomUnderSampler(sampling_strategy=0.75)
```

```
# Creating an imblearn pipeline to perform SMOTE followed by undersampling
from imblearn.pipeline import Pipeline
sampling_pipeline = Pipeline([
    ("smote", smote),
    ("under_sampling", under)
])
```

```
# Resampling the training dataset to tackle imbalance problem
X_train_rs, y_train_rs = sampling_pipeline.fit_resample(X_train, y_train)
print("Number of training samples before resampling: ", Counter(y_train))
print("Number of training samples after resampling: ", Counter(y_train_rs))
```

Number of training samples before resampling: Counter((0: 99510, 1: 85))
Number of training samples after resampling: Counter((0: 13269, 1: 9551))

```
# Making a list of attributes that needs to be scaled via normalization/standardization.
std_attributes = ["attribute2", "attribute3", "attribute4", "attribute5", "attribute6", "attribute8", "attribute9"]
min_max_attributes = ["attribute1", "attribute7"]
```

```
# Building a ColumnTransformer pipeline for data preprocessing
preprocessing_pipeline = ColumnTransformer([
    ("std_scaler", StandardScaler(), std_attributes),
    ("min_max_scaler", MinMaxScaler(), min_max_attributes)
])
```

5. Model Building

Since our test dataset is highly imbalanced, the chances of failure are less, evaluating the model based on accuracy is not a good idea. Instead, we'll be focusing on the ROC-AUC curve as a performance metric when it comes to solving our problem. Also, an output in the form of probability is required by the team.

```
# Create a Stratified K fold split
cv = StratifiedKFold(n_splits=10, random_state=42, shuffle=True)
```

```
# Creating a function to display scores
def display_scores(scores):
    print("ROC-AUC Scores: ", scores.mean())
    print("ROC-AUC Mean: ", scores.mean())
    print("ROC-AUC SD: ", scores.std())
```

```
# A function to print all metrics of results on test set
def print_metrics(prediction):
    print("Accuracy: ", round(accuracy_score(y_test, prediction), 2))
    print("ROC-AUC score: ", round(roc_auc_score(y_test, prediction), 2))
    print("F1-score: ", round(f1_score(y_test, prediction, average="macro"), 2))
    print("Precision: ", round(precision_score(y_test, prediction, average="macro"), 2))
    print("Recall: ", round(recall_score(y_test, prediction, average="macro"), 2))
    print("-----")
    print(confusion_matrix(y_test, prediction))
    print("-----")
    print(classification_report(y_test, prediction))
```

5.1 Logistic Regression

```
# Let's first build a Simple Logistic regression Model
lr_pipeline = Pipeline([
    ("preprocess", preprocessing_pipeline),
    ("lr", LogisticRegression(solver="saga", max_iter=1000, class_weight="balanced"))
])
```

```
scores = cross_val_score(lr_pipeline, X_train_rs, y_train_rs, scoring="roc_auc", cv=cv)

display_scores(scores)
```

ROC-AUC Scores: [0.99973113 0.99896421 0.99951945 0.99949029 0.99950393 0.9996713
0.9997815 0.99972621 0.99959226 0.99958171]
ROC-AUC Mean: 0.9996344742898593
ROC-AUC SD: 0.0005215674852853765

```
# Create a param grid with values for tuning hyperparameters
lr_param_grid = [{"lr_solver": ["newton-cg", "lbfgs", "saga", "sag", "lr_penalty": [None, 'l1', 'l2', 'el1', 'l3', 'l4', 'l5', 'l6', 'l7', 'l8', 'l9', 'l10', 'l11', 'l12', 'l13', 'l14', 'l15', 'l16', 'l17', 'l18', 'l19', 'l20', 'l21', 'l22', 'l23', 'l24', 'l25', 'l26', 'l27', 'l28', 'l29', 'l30', 'l31', 'l32', 'l33', 'l34', 'l35', 'l36', 'l37', 'l38', 'l39', 'l40', 'l41', 'l42', 'l43', 'l44', 'l45', 'l46', 'l47', 'l48', 'l49', 'l50', 'l51', 'l52', 'l53', 'l54', 'l55', 'l56', 'l57', 'l58', 'l59', 'l60', 'l61', 'l62', 'l63', 'l64', 'l65', 'l66', 'l67', 'l68', 'l69', 'l70', 'l71', 'l72', 'l73', 'l74', 'l75', 'l76', 'l77', 'l78', 'l79', 'l80', 'l81', 'l82', 'l83', 'l84', 'l85', 'l86', 'l87', 'l88', 'l89', 'l90', 'l91', 'l92', 'l93', 'l94', 'l95', 'l96', 'l97', 'l98', 'l99', 'l100', 'l101', 'l102', 'l103', 'l104', 'l105', 'l106', 'l107', 'l108', 'l109', 'l110', 'l111', 'l112', 'l113', 'l114', 'l115', 'l116', 'l117', 'l118', 'l119', 'l120', 'l121', 'l122', 'l123', 'l124', 'l125', 'l126', 'l127', 'l128', 'l129', 'l130', 'l131', 'l132', 'l133', 'l134', 'l135', 'l136', 'l137', 'l138', 'l139', 'l140', 'l141', 'l142', 'l143', 'l144', 'l145', 'l146', 'l147', 'l148', 'l149', 'l150', 'l151', 'l152', 'l153', 'l154', 'l155', 'l156', 'l157', 'l158', 'l159', 'l160', 'l161', 'l162', 'l163', 'l164', 'l165', 'l166', 'l167', 'l168', 'l169', 'l170', 'l171', 'l172', 'l173', 'l174', 'l175', 'l176', 'l177', 'l178', 'l179', 'l180', 'l181', 'l182', 'l183', 'l184', 'l185', 'l186', 'l187', 'l188', 'l189', 'l190', 'l191', 'l192', 'l193', 'l194', 'l195', 'l196', 'l197', 'l198', 'l199', 'l200', 'l201', 'l202', 'l203', 'l204', 'l205', 'l206', 'l207', 'l208', 'l209', 'l210', 'l211', 'l212', 'l213', 'l214', 'l215', 'l216', 'l217', 'l218', 'l219', 'l220', 'l221', 'l222', 'l223', 'l224', 'l225', 'l226', 'l227', 'l228', 'l229', 'l230', 'l231', 'l232', 'l233', 'l234', 'l235', 'l236', 'l237', 'l238', 'l239', 'l240', 'l241', 'l242', 'l243', 'l244', 'l245', 'l246', 'l247', 'l248', 'l249', 'l250', 'l251', 'l252', 'l253', 'l254', 'l255', 'l256', 'l257', 'l258', 'l259', 'l260', 'l261', 'l262', 'l263', 'l264', 'l265', 'l266', 'l267', 'l268', 'l269', 'l270', 'l271', 'l272', 'l273', 'l274', 'l275', 'l276', 'l277', 'l278', 'l279', 'l280', 'l281', 'l282', 'l283', 'l284', 'l285', 'l286', 'l287', 'l288', 'l289', 'l290', 'l291', 'l292', 'l293', 'l294', 'l295', 'l296', 'l297', 'l298', 'l299', 'l300', 'l301', 'l302', 'l303', 'l304', 'l305', 'l306', 'l307', 'l308', 'l309', 'l310', 'l311', 'l312', 'l313', 'l314', 'l315', 'l316', 'l317', 'l318', 'l319', 'l320', 'l321', 'l322', 'l323', 'l324', 'l325', 'l326', 'l327', 'l328', 'l329', 'l330', 'l331', 'l332', 'l333', 'l334', 'l335', 'l336', 'l337', 'l338', 'l339', 'l340', 'l341', 'l342', 'l343', 'l344', 'l345', 'l346', 'l347', 'l348', 'l349', 'l350', 'l351', 'l352', 'l353', 'l354', 'l355', 'l356', 'l357', 'l358', 'l359', 'l360', 'l361', 'l362', 'l363', 'l364', 'l365', 'l366', 'l367', 'l368', 'l369', 'l370', 'l371', 'l372', 'l373', 'l374', 'l375', 'l376', 'l377', 'l378', 'l379', 'l380', 'l381', 'l382', 'l383', 'l384', 'l385', 'l386', 'l387', 'l388', 'l389', 'l390', 'l391', 'l392', 'l393', 'l394', 'l395', 'l396', 'l397', 'l398', 'l399', 'l400', 'l401', 'l402', 'l403', 'l404', 'l405', 'l406', 'l407', 'l408', 'l409', 'l410', 'l411', 'l412', 'l413', 'l414', 'l415', 'l416', 'l417', 'l418', 'l419', 'l420', 'l421', 'l422', 'l423', 'l424', 'l425', 'l426', 'l427', 'l428', 'l429', 'l430', 'l431', 'l432', 'l433', 'l434', 'l435', 'l436', 'l437', 'l438', 'l439', 'l440', 'l441', 'l442', 'l443', 'l444', 'l445', 'l446', 'l447', 'l448', 'l449', 'l450', 'l451', 'l452', 'l453', 'l454', 'l455', 'l456', 'l457', 'l458', 'l459', 'l460', 'l461', 'l462', 'l463', 'l464', 'l465', 'l466', 'l467', 'l468', 'l469', 'l470', 'l471', 'l472', 'l473', 'l474', 'l475', 'l476', 'l477', 'l478', 'l479', 'l480', 'l481', 'l482', 'l483', 'l484', 'l485', 'l486', 'l487', 'l488', 'l489', 'l490', 'l491', '
```



```
In [13]: # Create ROC Graph
lr_fpr, lr_tpr, lr_thresholds = roc_curve(y_test, final_lr_model.predict_proba(X_test)[:,:])
rf_fpr, rf_tpr, rf_thresholds = roc_curve(y_test, final_rf_model.predict_proba(X_test)[:,:])
knn_fpr, knn_tpr, knn_thresholds = roc_curve(y_test, knn_pipeline.predict_proba(X_test)[:,:])
xgb_fpr, xgb_tpr, xgb_thresholds = roc_curve(y_test, xgb_pipeline.predict_proba(X_test)[:,:])
gb_fpr, gb_tpr, gb_thresholds = roc_curve(y_test, gb_pipeline.predict_proba(X_test)[:,:])

lr_roc_auc = roc_auc_score(y_test, y_pred_lr)
rf_roc_auc = roc_auc_score(y_test, y_pred_rf)
knn_roc_auc = roc_auc_score(y_test, y_pred_knn)
xgb_roc_auc = roc_auc_score(y_test, y_pred_xgb)
gb_roc_auc = roc_auc_score(y_test, y_pred_gb)

plt.figure(figsize=(8, 6))

# Plot Logistic Regression ROC
plt.plot(lr_fpr, lr_tpr, label='LR (area = %0.2f)' % lr_roc_auc)

# Plot Random Forest ROC
plt.plot(rf_fpr, rf_tpr, label='RF (area = %0.2f)' % rf_roc_auc)

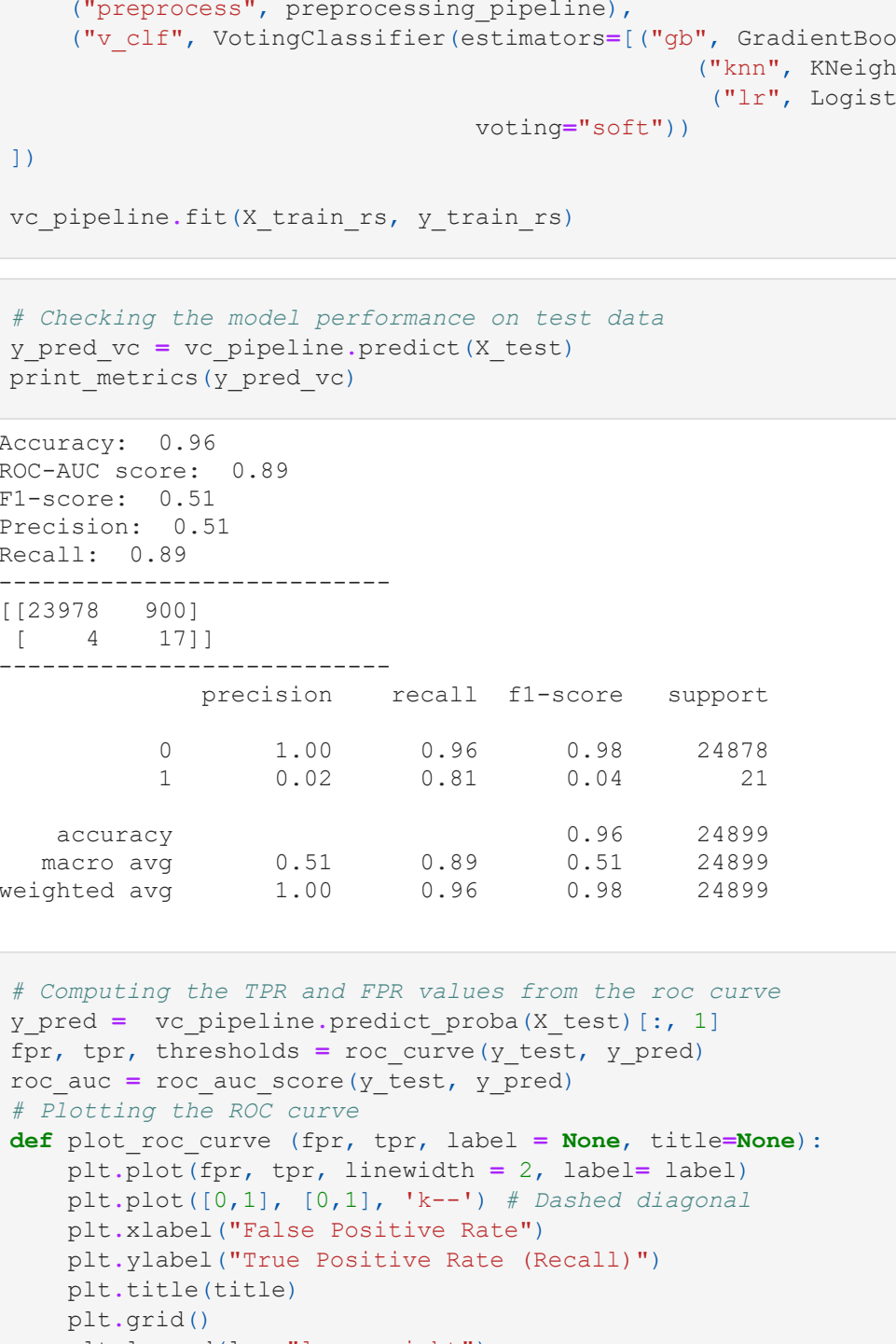
# Plot K Nearest Neighbors ROC
plt.plot(knn_fpr, knn_tpr, label='KNN (area = %0.2f)' % knn_roc_auc)

# Plot GBoost ROC
plt.plot(gb_fpr, gb_tpr, label='GB (area = %0.2f)' % gb_roc_auc)

# Plot XGBoost ROC
plt.plot(xgb_fpr, xgb_tpr, label='XGB (area = %0.2f)' % xgb_roc_auc)

# Plot Base Rate ROC
plt.plot([0,1], [0,1],label='Base')

plt.xlim((0.0, 1.05))
plt.ylim((0.0, 1.05))
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Graph')
plt.legend(loc='lower right')
plt.show()
```



8. Voting Classifier

Since the models trained by us are lacking in some parts (either TP or FP) when it comes to predicting the failure of devices, we will ensemble three of our top performing models using a Voting classifier. This will help us reduce the number of FP and increase the TP.

```
In [1]: # Creating a Voting Classifier consisting of our top 3 performing models i.e GB, LR and KNN Classifiers
vc_pipeline = Pipeline((
    ("preprocess", preprocessing_pipeline),
    ("vc_clf", VotingClassifier(estimators=[ ("gb", GradientBoostingClassifier(n_estimators=500, learning_rate = 0.01,
                                                                                   ("knn", KNeighborsClassifier(n_neighbors=25)),
                                                                                   ("lr", LogisticRegression(solver="saga", max_iter=1000, class_weight="soft"))
    ])

vc_pipeline.fit(X_train_rs, y_train_rs)
```

```
In [234]: # Checking the model performance on test data
y_pred_vc = vc_pipeline.predict(X_test)
print_metrics(y_pred_vc)

Accuracy: 0.96
ROC-AUC score: 0.89
F1-score: 0.51
Precision: 0.51
Recall: 0.89
-----
[[23978  900]
 [ 4      17]]
-----
      precision    recall  f1-score   support

      0         1.00      0.96      0.98      24878
      1         0.02      0.81      0.04         21

 accuracy          0.96      0.96      24899
 macro avg         0.51      0.89      0.51      24899
 weighted avg         1.00      0.96      0.98      24899
```

```
In [241]: # Computing the FPR and FPR values from the roc curve
y_pred = vc_pipeline.predict_proba(X_test)[:,:].ravel()
fpr, tpr, thresholds = roc_curve(y_test, y_pred)
roc_auc = roc_auc_score(y_test, y_pred)
# Plotting the ROC curve
def plot_roc_curve(fpr, tpr, label=None, title=None):
    plt.plot(fpr, tpr, linewidth=2, label=label)
    plt.plot([0,1], [0,1], 'k--') # dashed diagonal
    plt.xlabel("False Positive Rate")
    plt.ylabel("True Positive Rate (Recall)")
    plt.title(title)
    plt.grid()
    plt.legend(loc="lower right")

plot_roc_curve(fpr, tpr, 'VC (area = %0.2f)' % roc_auc, "ROC-AUC for Voting Classifier")
plt.show()
```



```
In [242]: # This looks great. We'll go ahead and save this as our final model
joblib.dump(vc_pipeline, "../models/final_voting_classifier_model.pkl")
```

```
Out[242]: ['../models/final_voting_classifier_model.pkl']
```

9. Conclusion

To conclude, we have implemented an Ensemble of trained Machine Learning models such as Logistic Regression, Gradient Boosting Classifier and K Nearest Neighbors using Voting Classifier to predict the probability of device failures given some input features. The ensemble model achieved an ROC-AUC score of almost 0.90 on the unseen test data.

Thank you!