

Balancing the Imbalance: Advanced Anomaly Detection for Predictive Maintenance

Ch.Srivatsa Alivelu Mangatayi¹ Vanam Pavan²

Sri Ratna³ G.Bharath⁴ S.Narasimha⁵

¹ Assistant Professor, ACE Engineering College Hyderabad, India

²Student, ACE Engineering College Hyderabad, India ³Student, ACE Engineering College Hyderabad, India ⁴Student, ACE Engineering College Hyderabad, India

⁵Student, ACE Engineering College Hyderabad, India

Email : ¹srivatsajava@gmail.com ²pavanvanam890@gmail.com

³srirathnatirunagari@gmail.com ⁴bharathgurram443@gmail.com

⁵sangishettinarasimha@gmail.com

ABSTRACT:

Predictive maintenance is a crucial application in industrial systems, ensuring timely interventions to prevent failures and reduce downtime. However, the highly imbalanced nature of failure datasets poses a significant challenge to traditional machine learning models. This paper proposes an advanced anomaly detection framework utilizing synthetic data generation, oversampling techniques (SMOTE, ADASYN), and cost-sensitive learning to enhance failure detection accuracy. The system integrates Random Forest and XGBoost classifiers, optimized through hyperparameter tuning, to improve predictive maintenance efficiency. The proposed method demonstrates improved recall and precision, reducing false negatives while maintaining robustness in detecting rare failure events.

I. INTRODUCTION:

In industries reliant on machinery, unexpected failures lead to substantial financial losses and operational delays. Predictive maintenance, leveraging sensor data, aids in forecasting potential failures. However, the rarity of failure cases within datasets skews model performance, favoring normal conditions over anomalies. Addressing this issue requires advanced techniques in data preprocessing and machine learning model training. This paper introduces an innovative approach that enhances anomaly detection through synthetic data generation, resampling techniques, and cost-sensitive learning.

II. OBJECTIVES:

- Develop a robust predictive maintenance model for gas turbine systems.
- Address the issue of imbalanced datasets using SMOTE, ADASYN, and cost-sensitive learning.
- Train and optimize Random Forest and XGBoost models for anomaly detection.
- Deploy a real-time monitoring system for predictive maintenance.

III. PROBLEM STATEMENT

Predictive maintenance is a critical component of industrial asset management, aiming to mitigate unplanned downtimes and optimize operational efficiency. However, the rarity of failure events in sensor datasets leads to a highly imbalanced data distribution, which significantly hampers the effectiveness of traditional anomaly detection models. The system is designed to work with gas turbine Machine models such as GE 7F.05 and Siemens SGT-800. Conventional machine learning approaches often exhibit bias towards the majority class, leading to suboptimal recall rates for failure detection. This challenge necessitates the development of an advanced anomaly detection framework that leverages synthetic data augmentation, cost-sensitive learning, and ensemble-based models to improve predictive accuracy. By integrating techniques such as SMOTE, ADASYN, and hyperparameter-optimized classifiers like Random Forest and XGBoost, this study proposes a robust

methodology to enhance failure prediction capabilities, ensuring proactive maintenance interventions and reducing operational risks in industrial settings.

IV. PROPOSED SYSTEM

The system is designed to work with **gas turbine Machine** models such as **GE 7F.05** and **Siemens SGT-800**, utilizing sensor data to detect anomalies and predict failures.

The proposed system incorporates:

Data Collection: Utilizing real-world sensor datasets such as the NASA CMAPSS dataset and synthetic data generation.

Data Preprocessing: Handling missing values, feature engineering, and normalization.

Handling Imbalance: Implementing SMOTE, ADASYN, and cost-sensitive learning to balance datasets.

Model Training & Optimization: Utilizing Random Forest and XGBoost, fine-tuned via GridSearchCV.

Deployment: Integrating the trained model into a real-time monitoring system.

V. SOFTWARE REQUIREMENTS

- Platform : Jupyter Notebook, VS code
- Technologies :

- **Programming:** Python, JavaScript
- **Libraries:** Pandas, NumPy, Scikit-learn, Imbalanced-learn, XGBoost
- **Frameworks:** Flask (for deployment), Next.js
- **Database:** MongoDB (for storing sensor data)
- **Visualization:** Plotly Dash / Streamlit (for monitoring dashboard)
- Testing :
- Postman

VI. TECHNOLOGY DESCRIPTION

Machine and Model:

- **Gas Turbine Models:** GE 7F.05, Siemens SGT-800
- **Machine Learning Models:** Random Forest,

XGBoost

Machine Learning Techniques:

1. **Random Forest Classifier:** A robust ensemble model trained with class-weight adjustments.
2. **XGBoost:** A gradient boosting model optimized for imbalance handling.
3. **SMOTE & ADASYN:** Synthetic oversampling methods to balance dataset distribution.
4. **Hyperparameter Tuning:** Optimizing models via GridSearchCV.

Deployment Framework:

- A **Flask API** for model inference.
- A **Streamlit dashboard** for real-time anomaly monitoring.
- A **MongoDB database** to store historical sensor data.

VII. ALGORITHM

Step 1: Data Preprocessing

- Load sensor datasets from real-world or synthetic sources.
- Handle missing values using forward-fill and interpolation techniques.
- Feature engineering: Calculate rate-of-change metrics and rolling averages.
- Normalize features using StandardScaler.

Step 2: Handling Imbalanced Data

- Apply **SMOTE** and **ADASYN** to generate synthetic failure cases.
- Implement **cost-sensitive learning** by adjusting class weights in classifiers.

Step 3: Model Training & Evaluation

- Train **Random Forest** and **XGBoost** models on resampled data.
- Evaluate performance using classification report and AUC-ROC scores.
- Optimize hyperparameters via **GridSearchCV**.

Step 4: Deployment & Real-Time Prediction

- Save trained models using **joblib**.
- Develop a **Flask API** for failure prediction.
- Implement **Streamlit-based visualization** for

sensor monitoring.

VIII. OUTPUT SCREENS

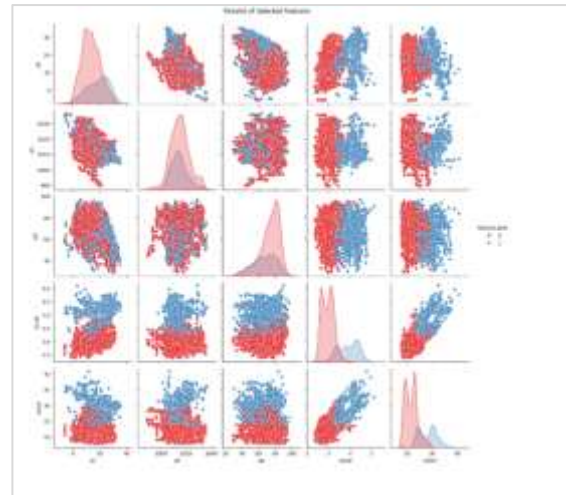
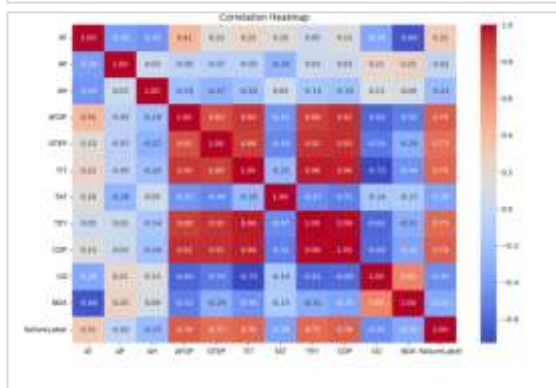
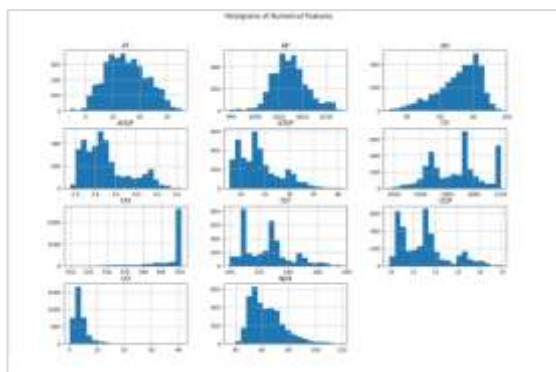
```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

# Load the dataset
file_path = "gh_2013_balanced_hungary.csv"
data = pd.read_csv(file_path)

# Display the first few rows
print(data.head())
```

	AP	AM	AN	ADDP	CTDP	TST	TST	TSY	CDP	
0	38.4018	1089.7	89.765	3.2607	37.284	1872.0	549.89	126.87	11.518	
1	32.4008	1015.3	85.922	2.7219	39.267	1845.7	550.08	189.89	18.367	
2	34.5078	1015.8	85.977	2.5637	31.559	1841.4	549.91	186.89	18.278	
3	34.8108	1018.0	82.784	2.7522	39.488	1831.8	549.98	122.19	18.648	
4	9.4718	1017.5	87.548	2.9802	21.478	1801.0	549.88	123.88	13.848	

```
AP      0
AM      0
AN      0
ADDP    0
CTDP    0
TST     0
TST     0
TSY     0
CDP     0
CN      0
NDR     0
Followed 0
df.shape
(10000, 11)
columns: ['AP', 'AM', 'AN', 'ADDP', 'CTDP', 'TST', 'TSY', 'CDP', 'CN', 'NDR', 'Followed']
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 11 columns):
#   Column      Non-Null Count  Dtype
---  --
0   AP           10000 non-null    float64
1   AM           10000 non-null    float64
2   AN           10000 non-null    float64
3   ADDP         10000 non-null    float64
4   CTDP         10000 non-null    float64
5   TST          10000 non-null    float64
6   TST          10000 non-null    float64
7   TSY          10000 non-null    float64
8   CDP          10000 non-null    float64
9   CN           10000 non-null    float64
10  NDR          10000 non-null    float64
11  Followed     10000 non-null    bool
dtypes: float64(10), bool(1)
memory usage: 335.4 KB
```



```
from imblearn.over_sampling import SMOTE

# Apply SMOTE to balance the dataset
smote = SMOTE(random_state=2)
X_resampled, y_resampled = smote.fit_resample(X_train, y_train)

# Check the class distribution after SMOTE
print('Class distribution after SMOTE:', pd.Series(y_resampled).value_counts())

# Class distribution after SMOTE: Followed
# 0    5000
# 1    5000
# Name: smart, dtype: int64
```

```
Model Classification Report:
precision    recall    F1-score   support
0         0.99      0.98      0.98      5000
1         0.99      0.99      0.99      5000
```

```
from sklearn.metrics import classification_report
from sklearn.metrics import confusion_matrix
from sklearn.metrics import roc_auc_score

# Apply SMOTE to balance the dataset
smote = SMOTE(random_state=2)
X_resampled, y_resampled = smote.fit_resample(X_train, y_train)

# Train the model
model = RandomForestClassifier()
model.fit(X_resampled, y_resampled)

# Evaluate the model
y_pred = model.predict(X_test)
y_prob = model.predict_proba(X_test)[:, 1]

# Print the model's performance
print('Model Performance Report:')
print(classification_report(y_test, y_pred))
print(roc_auc_score(y_test, y_prob))
```



IX. CONCLUSION

This paper presents an advanced anomaly detection framework for predictive maintenance, addressing class imbalance through synthetic oversampling and SMOTE and ADASYN Balancing Techniques. The proposed approach enhances failure detection accuracy, ensuring reliable and timely maintenance interventions. Future work includes integrating IoT-based real-time data collection and XAI Frameworks learning techniques for adaptive anomaly detection.

X. REFERENCES

- [1] Fernández, A., García, S., Galar, M., Prati, R. C., Krawczyk, B., & Herrera, F. (2018). Learning from Imbalanced Data Sets. Springer.

- [2] Buda, M., Maki, A., & Mazurowski, M. A. (2018). A systematic study of the class imbalance problem in convolutional neural networks. *Neural Networks*, 106, 249-259.

- [3] Chen, T., & Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*.

- [4] Saxena, A., & Goebel, K. (2008). Turbofan Engine Degradation Simulation Data Set. NASA Ames Prognostics Data Repository.