

Automobile Engine Fault Diagnosis Using Machine Learning Method

21900416 AN GYEON HEAL

21900727 JIN GA RAM

School of Mechanical & Control Engineering

10.01.2024

Industrial AI & Automation
Prof. Young-Keun Kim

CONTENTS

Engine Fault Diagnosis

Background



Problem Statement



Baseline Survey



Progress Scheme



Additional Study



BACKGROUND

Engine Fault Diagnosis

Necessity

- Vehicle engine malfunctions increase emissions of harmful pollutants like CO, HC, and CO₂, contributing to environmental problems.
- Engine malfunctions lead to high repair costs, unexpected downtime, and potential financial and safety risks.
- Engines may experience mechanical or electronic failures (e.g., sensor issues, pressure problems, injector defects), complicating the repair process.

Technical Challenges

- Traditional Engine Fault Diagnosis Methods Rely on Manual Inspections and Specialized Tools (e.g., fuel pressure measurements, OBD-2 scanners).
- It's Labor-Intensive, Expensive, and Require Expert Knowledge, It's Vulnerable to Human Error.

Solution

- By Applying Machine Learning, Variables from Engine Testing Experiments can be used for Classification to Identify Types of Engine Faults.
- Classification will Help Quickly Diagnose and Resolve Engine Faults, Enabling Fast and Efficient Repairs.

PROBLEM STATEMENT

Engine Fault Diagnosis

Objective

By **Developing** a Model that Outperforms those Implemented in Existing Journals,
 Accurately **Identify** types of Engine faults, Providing Greater **Speed** and **Precision** in Engine Repairs.

Detailed Objectives

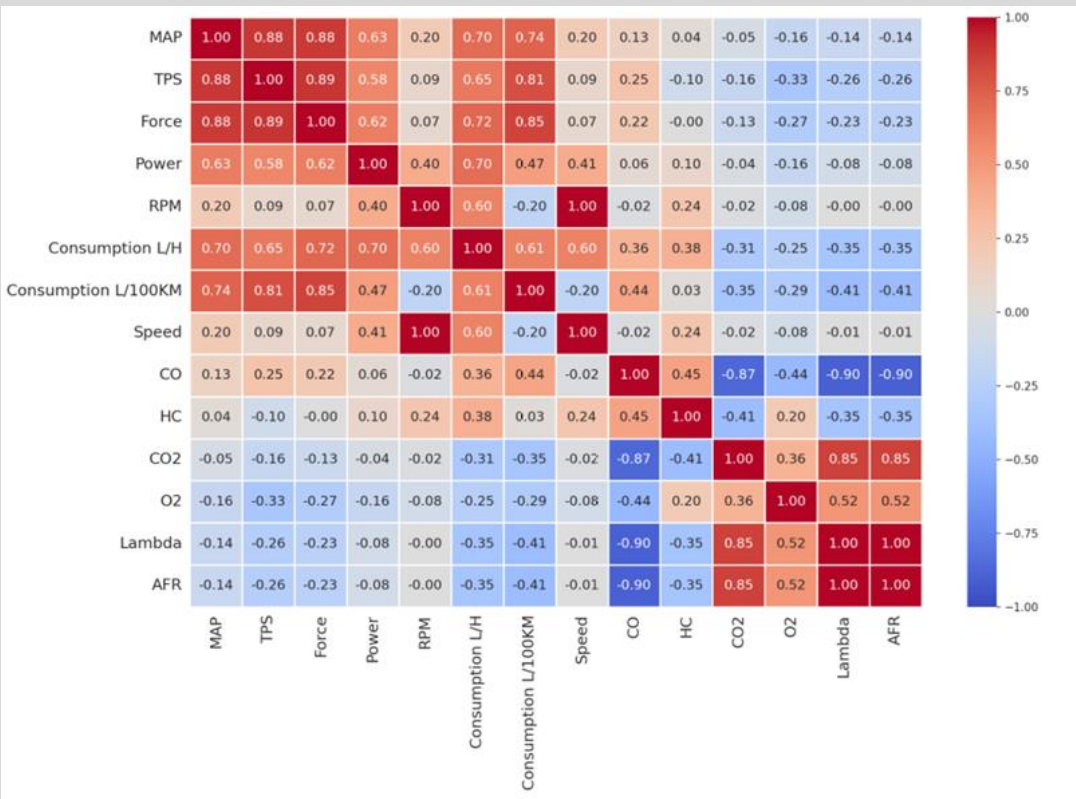
1. Classify Datasets by Machine Learning Model Using “**EngineFaultDB**” Datasets.
2. Understand and Analysis Datasets through **Feature Analysis**.
3. Improve Model Performance through **Feature Extraction** and **Feature Reduction/Selection**
4. Compare the **Model Performance** with Baseline Journal.

BASELINE SURVEY-JOURNAL

Engine Fault Diagnosis

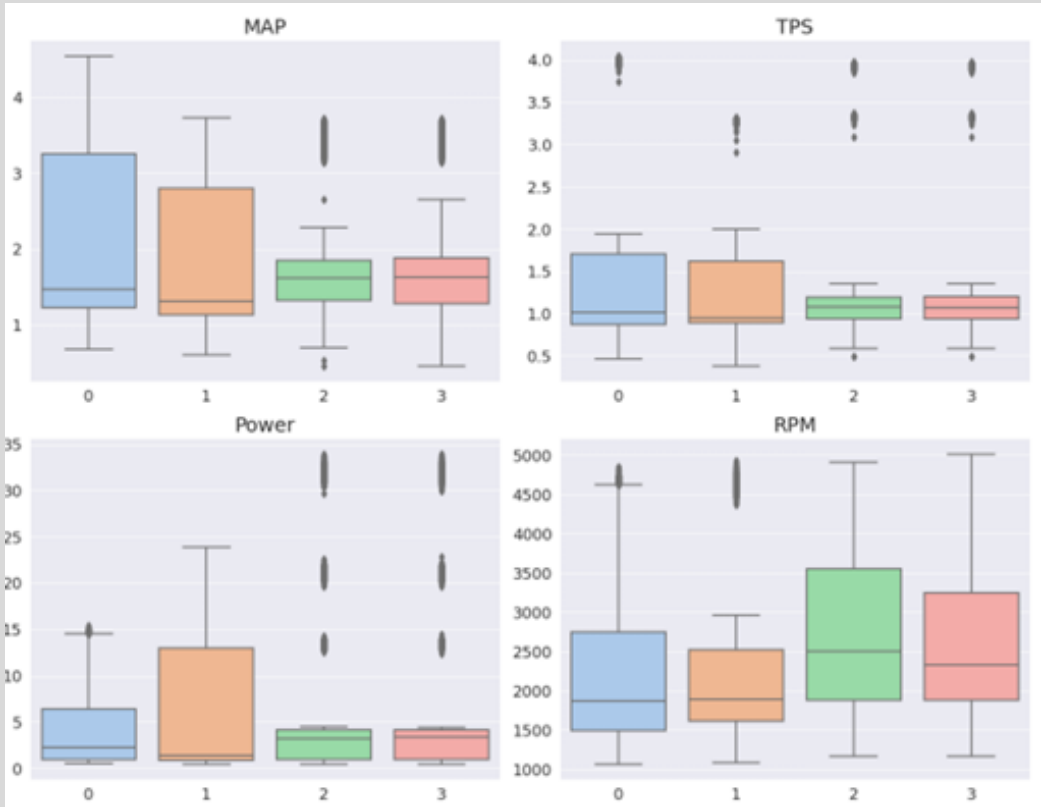
Journal (IEEE): [LINK](#)

1. Correlation Analysis



- Correlation :**
How variables influence each other
- Multicollinearity :**
Reducing highly correlated variables

2. Box Plot Analysis



- Data Distribution Summary**
- Identifying Outlier**
- Skewness Check**
- Range and Variability Check**

3. Data Preprocessing & Classification

$$X_{\text{scaled}} = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$$

- Each Datasets has a different range → **Scaling is necessary.**
- Min-Max Scaler**

- Logistic Regression**
- Decision Tree**
- Random Forest**
- SVM**
- KNN**
- Naive Bayes**
- Feed-Forward Neural Network**

4. Performance Confusion Metrics

- Accuracy**

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

- Precision**

$$\text{Precision} = \frac{TP}{TP + FP}$$

- Recall**

$$\text{Recall} = \frac{TP}{TP + FN}$$

- F1-Score**

$$\text{F1 - Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

Classifier	Accuracy	Precision	Recall	F1-score
LR	0.576	0.574	0.576	0.574
DT	0.750	0.750	0.750	0.750
RF	0.748	0.748	0.748	0.748
SVC	0.747	0.768	0.747	0.715
KNN	0.751	0.751	0.751	0.751
NB	0.394	0.370	0.394	0.353
Neural Net.	0.749	0.748	0.749	0.748

<Performance – Testsets>

Performance of KNN = 0.751

BASELINE SURVEY-DATASETS

Engine Fault Diagnosis

Datasets Name

- EngineFaultDB (Supervised Learning)

Achieved by

- Mary Vergara, Diego Rivera, Francklin Rivas-Echeverría

Test Engine

- C14NE Spark Ignition Engine: Petrol Engine (Gasoline)

Specification	Detail
Maximum power	83.7 HP @ 6000 RPM
Torque	113.56 N.m @ 3000 RPM
Displacement	1388 cc
Injection system	Multipoint
Fuel consumption	6.8 l/100 km
Valve configuration	SOHC

<Test Engine Configuration>

Engine Data Collection Method



<Engine Data Collection / Gas Analyzer Device>

- Gas Analyzer
- USB 6008 Data Acquisition Card (DAQ)

Datasets (Github) : [LINK](#)

INPUT

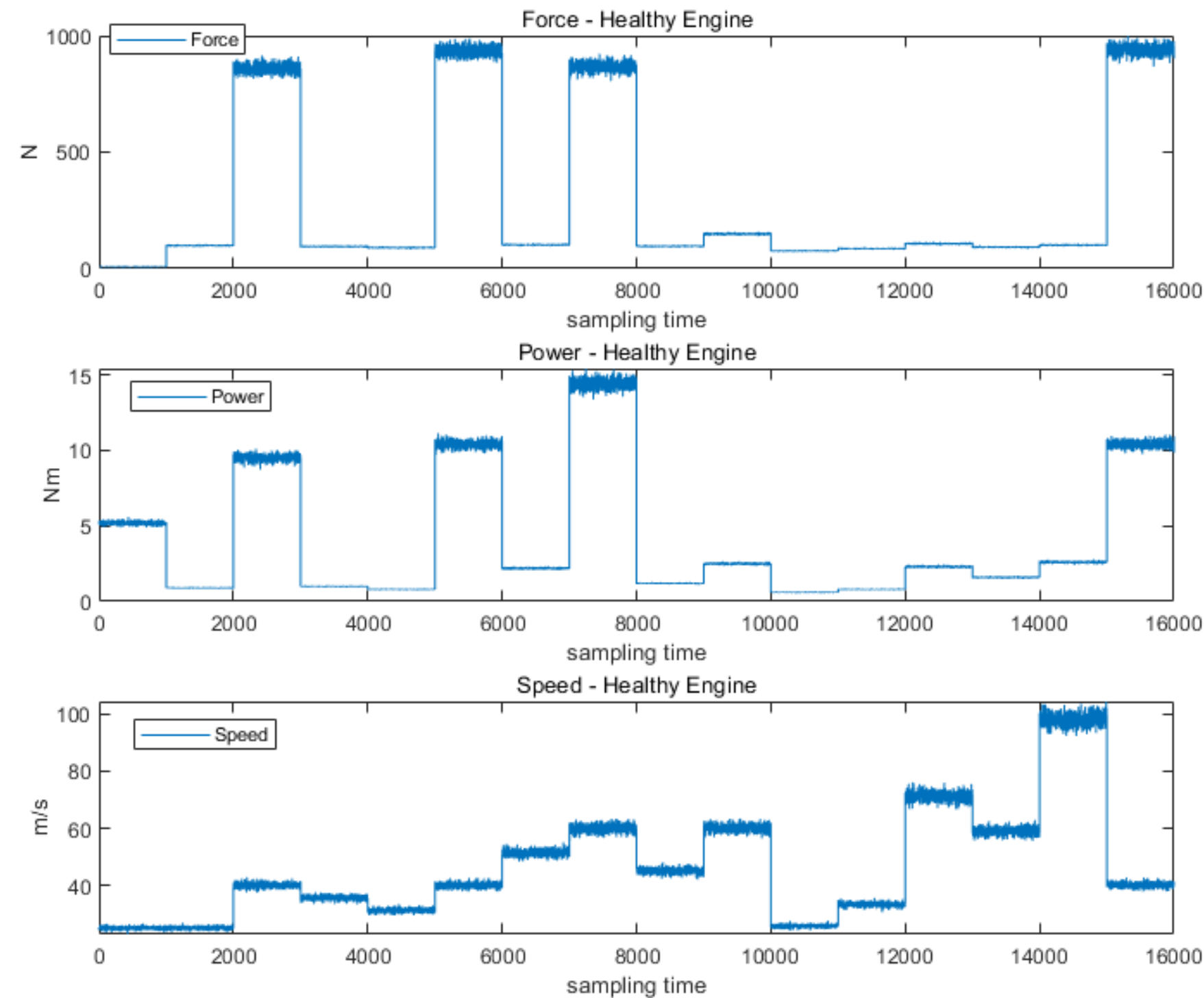
- Manifold Absolute Pressure (MAP)**
 - Pressure inside manifold [kPa]
- Throttle Position Sensor**
 - Position of Throttle: about fuel injection, ignition time, etc. [%]
- Force**
 - Engine torque/rotational force [N]
- Power**
 - Energy transferred in engine [kW]
- Revolutions Per Minute (RPM)**
 - The times crankshaft rotates per minute
- Fuel consumption L/H**
 - Engine's fuel consumption rate
- Fuel consumption L/100KM**
 - Engine's fuel efficiency by distance
- Speed**
 - Vehicle's travel speed [km/h]
- Carbon monoxide (CO)**
 - CO concentration in the exhaust gases [%]
- Hydrocarbons (HC)**
 - Hydrocarbons concentration [%]
- Carbon dioxide (CO2)**
 - CO2 concentration : combustion efficiency [%]
- Oxygen (O2)**
 - Oxygen concentration : insights about combustion [%]
- Lambda**
 - Air-fuel equivalence ratio
- Air-Fuel Ration (AFR)**
 - Ratio of the air fuel in the combustion chambers

OUTPUT

- Fault type 0: Normal (16,000 entries)
- Fault type 1: Rich mixture - High Pressure, Incorrect Sensor, etc. (10,988 entries)
- Fault type 2: Lean mixture – Low Pressure, Incorrect Sensor, etc. (15,000 entries)
- Fault type 3: Low Voltage – Worn Spark, Defective Coil, etc. (14,001 entries)

BASELINE SURVEY – DATASETS

Engine Fault Diagnosis



Asunto: Consulta sobre el tiempo de muestreo en el conjunto de datos EngineFaultDB



안전철학부생 <21900416@handong.ac.kr>
vmaryjose, nrívera, frivas6, theleothomasramos, 숨은 참조: 진가람학부생에게 ▾

10월 7일 (월) 오후 11:33 (34분 전) ☆ ↶ ⋮

Estimada Mary Vergara, Diego Rivera, Francklin Rivas-Echeverría y Leo Ramos,

Mi nombre es An Gyeonhil y soy estudiante de pregrado en la Facultad de Ingeniería Mecánica y Control en la Universidad Global de Handong, Corea del Sur. Recientemente leí su artículo titulado "EngineFaultDB: A Novel Dataset for Automotive Engine Fault Classification and Baseline Results" y lo encontré increíblemente interesante, especialmente en lo relacionado con el proceso de recolección de datos y el uso de aprendizaje automático.

He encontrado el conjunto de datos en el GitHub de Leo Ramos y estoy planeando usarlo para un proyecto de aprendizaje automático en uno de mis cursos. Me gustaría amablemente preguntar si podría proporcionarme más información sobre el **tiempo de muestreo** utilizado en el conjunto de datos EngineFaultDB.

Muchas gracias de antemano por su tiempo y asistencia.

Atentamente,
An Gyeonhil

Dear Mary Vergara, Diego Rivera, Francklin Rivas-Echeverría, and Leo Ramos,

My name is An Gyeonhil, an undergraduate student at the School of Mechanical and Control Engineering, Handong Global University, South Korea. I recently came across your paper titled "EngineFaultDB: A Novel Dataset for Automotive Engine Fault Classification and Baseline Results" and found it incredibly insightful, especially regarding the data collection process and the use of machine learning.

I located the dataset on Leo Ramos's GitHub and I am planning to use it for a machine learning project in one of my courses. I would like to kindly ask if you could provide more information regarding the **sampling time** used in the EngineFaultDB dataset.

Thank you very much in advance for your time and assistance.

Best regards,
An Gyeonhil

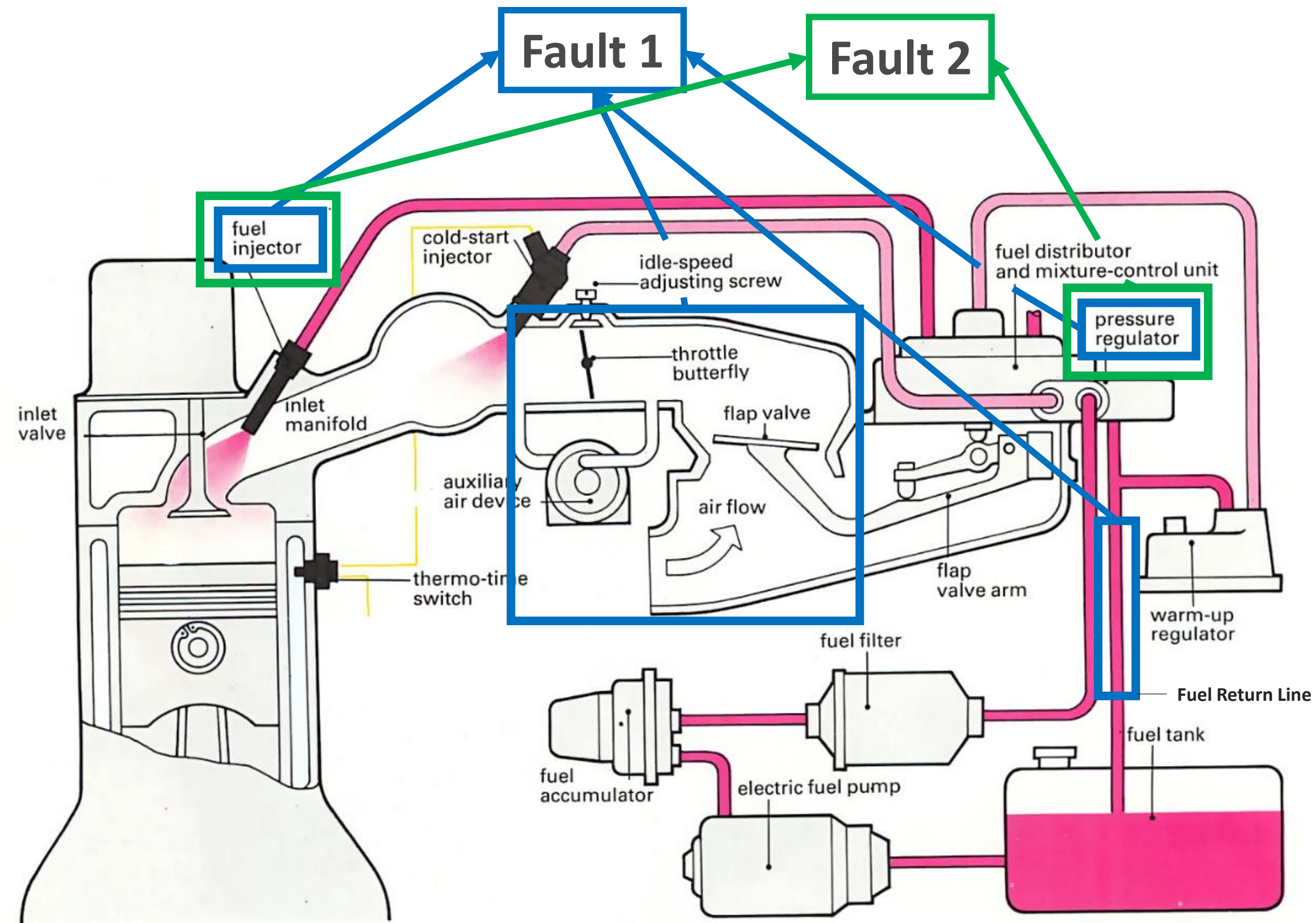
Experiment Period (T) = 1,000

Experiment under Various Conditions, per 1,000 periods

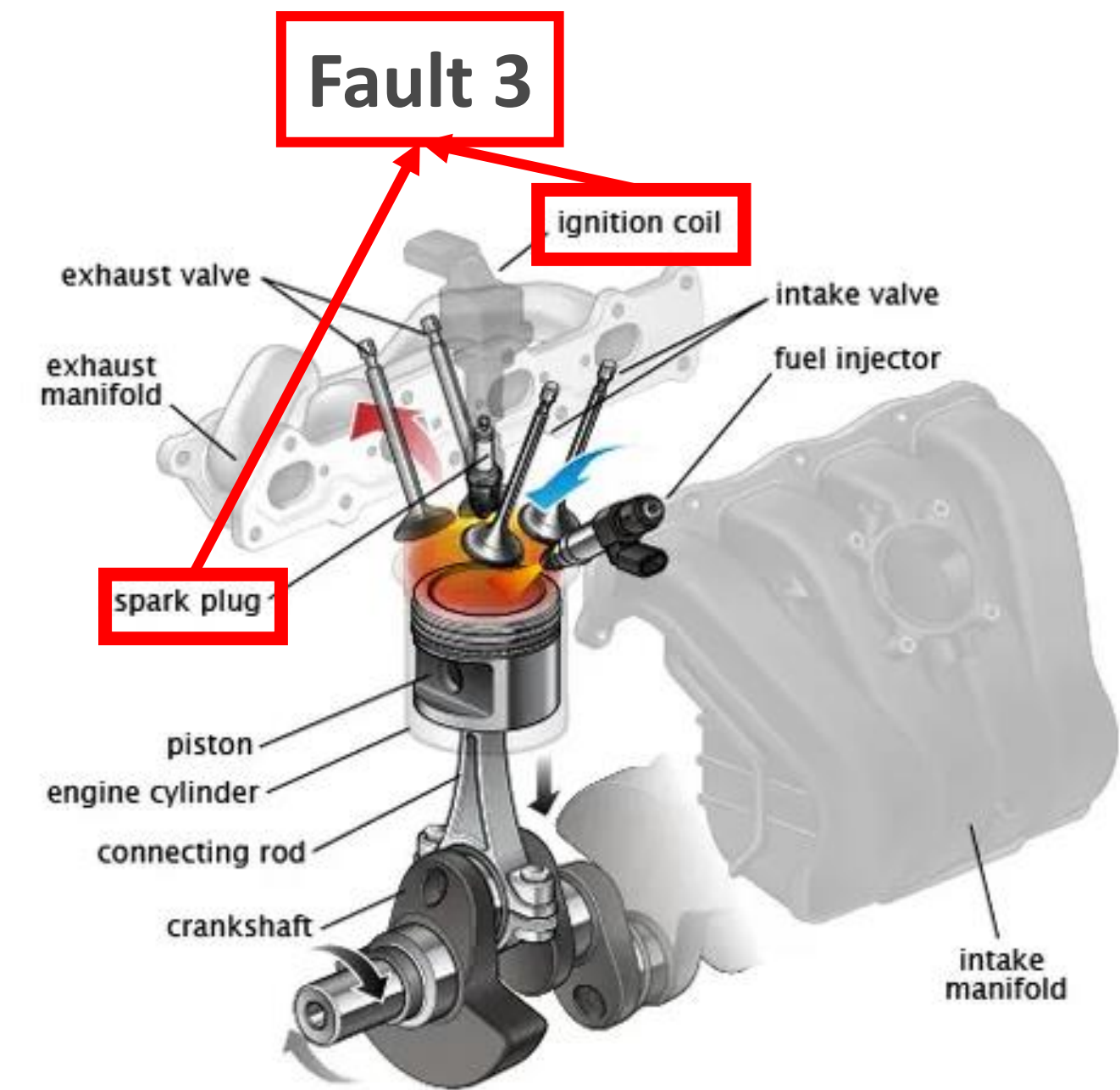
There's **NO Information about Sampling Time**

BASELINE SURVEY – FAULTS

Engine Fault Diagnosis



< Lucas Mechanical Fuel Injection System (Petrol Engine) >



< Ignition System >

Fault 1 Rich mixture

Incorrect sensor performance
High fuel pressure
Defective injector
Faulty pressure regulator
Clogged air filter
Clogged fuel return line

Fault 2 Lean mixture

Incorrect sensor performance
Low fuel pressure
Defective injector
Faulty pressure regulator

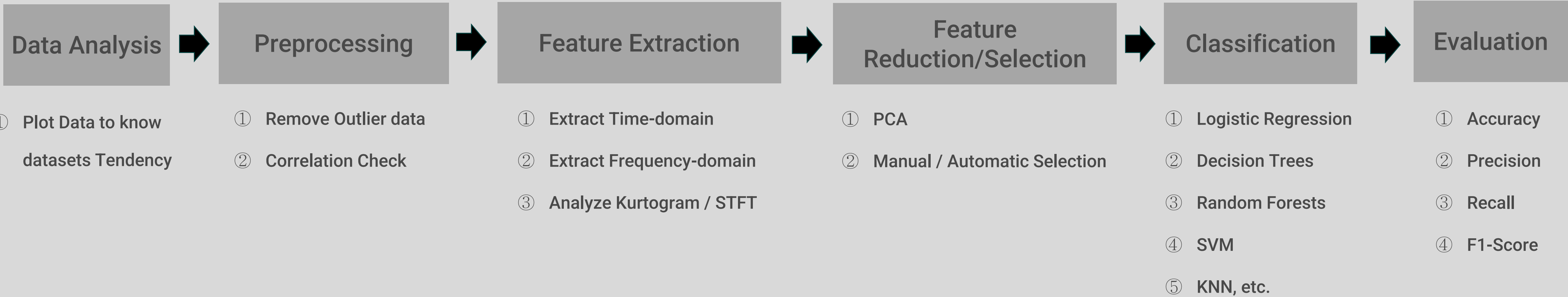
Fault 3 Low voltage

Worn spark plugs
Faulty ignition cables
Defective coil
Faulty sensor wiring

PROJECT PROGRESS SCHEME

Engine Fault Diagnosis

Method



Expected Outcome

Extracting/Reducing/
Selecting Features

- ① Derive Meaningful Information
- ② Migrates Overfitting
- ③ Enhance Computational Efficiency

GOAL

Higher Performance than Baseline Journal (≥ 0.751)

Schedule

	1 st Week	2 nd Week	3 rd Week
Data Analysis			
Preprocessing			
Feature Extraction			
Feature Reduction/Selection			
Classification			
Evaluation			
Writing Report			

Role Distribution

AN GYEON HEAL

- Analyze about Baseline Journal
- Search Engine Structure & Vulnerable Part or Process
- Research about Diagnosis Process

JIN GA RAM

- Search Applicable Way or Field
- Search Additional Datasets
- Research about Additional ML Method that can Improve Model Performance

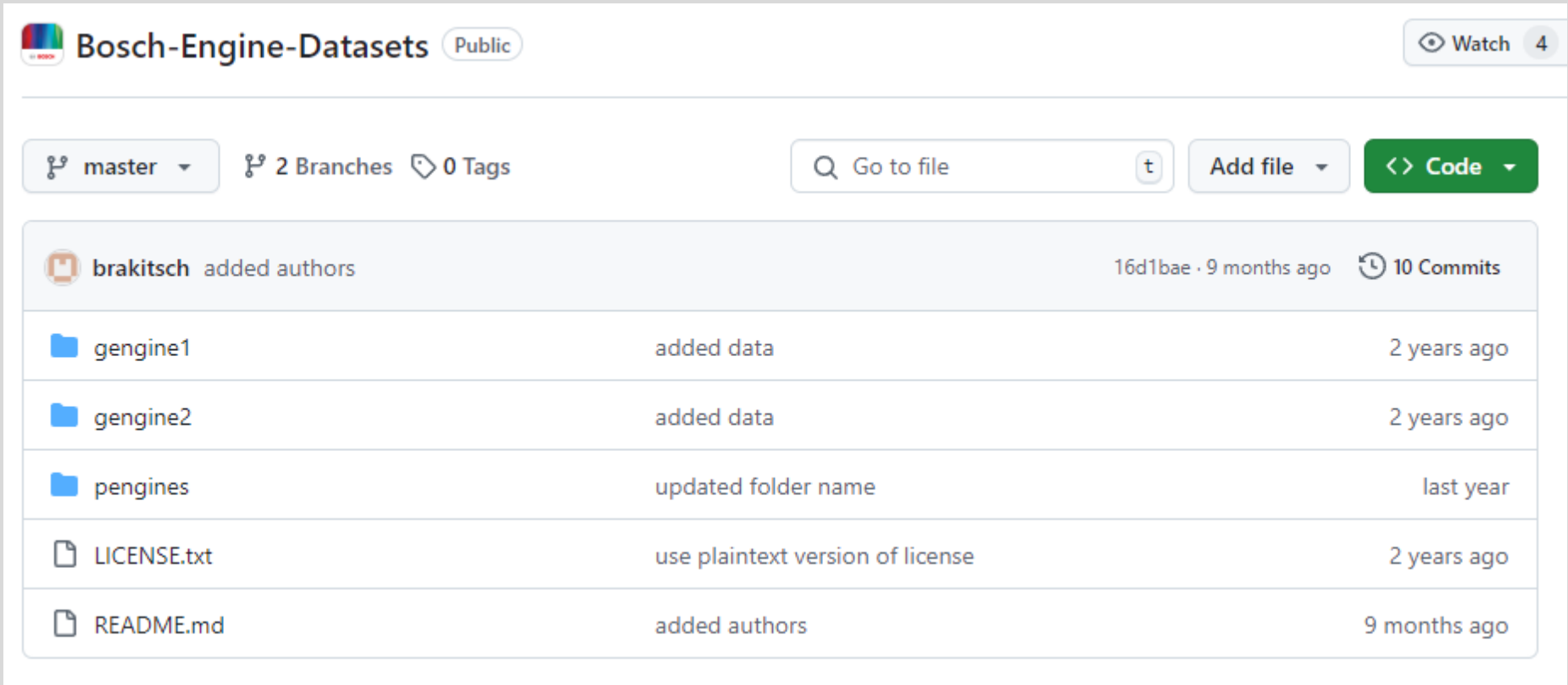
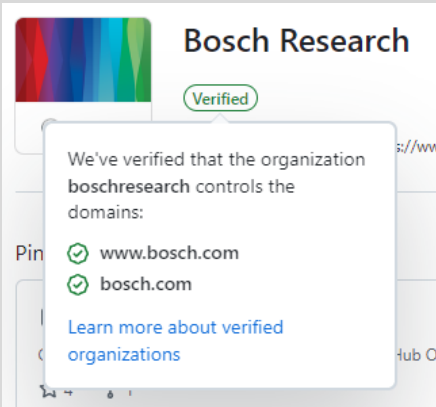
ADDITIONAL RESEARCH

Engine Fault Diagnosis

Further Study

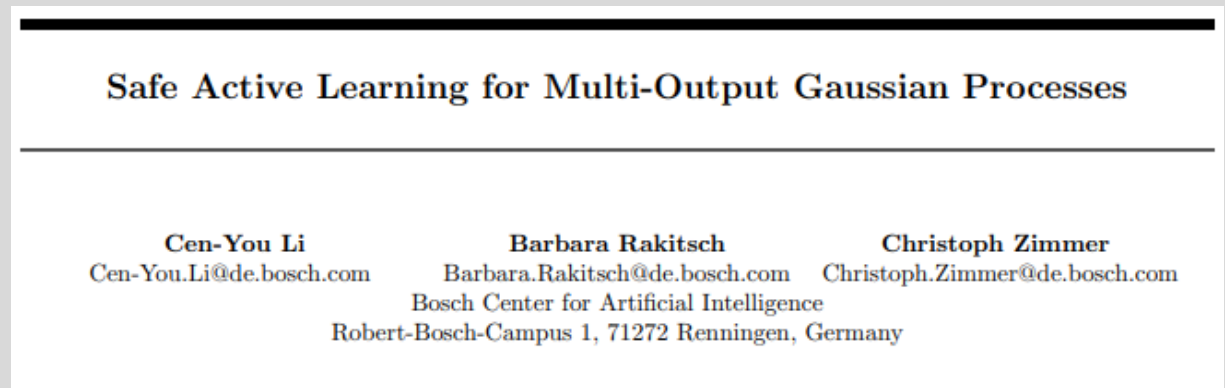
Apply the Model to the Open BOSCH Engine Dataset Available on BOSCH’s Official GitHub.

Datasets Information



[5] Bosch Research. (n.d.). Bosch Engine Datasets. GitHub. <https://github.com/boschresearch/Bosch-Engine-Datasets>

Related Survey/Research



❖ Uses “Multiple Output Gaussian Process” (Safe Active Learning)

❖ RMSE Performance Comparison: Active Learning MOGP (≤ 0.4)

❖ Datasets : BOSCH-Engine-Datasets

❖ Supervised Training : Labeling Output (HC, NOx, O2, etc.)

[6] Li, C. Y., Rakitsch, B., & Zimmer, C. (2022, May). Safe active learning for multi-output gaussian processes. In International Conference on Artificial Intelligence and Statistics (pp. 4512-4551). PMLR.

Plan

❖ Train ML Model which Labeled by (HC, NOx, O2, etc.) and Compare Performance

OR

❖ Diagnose Engine Fault Using BOSCH-Engine-Datasets

REFERENCES

Engine Fault Diagnosis

Baseline Survey – Journal & Datasets



Multidisciplinary | Rapid Review | Open Access Journal

Received 10 October 2023, accepted 2 November 2023, date of publication 8 November 2023, date of current version 14 November 2023.

Digital Object Identifier 10.1109/ACCESS.2023.3331316



RESEARCH ARTICLE

EngineFaultDB: A Novel Dataset for Automotive Engine Fault Classification and Baseline Results

MARY VERGARA¹, LEO RAMOS^{2,3}, (Student Member, IEEE),
NÉSTOR DIEGO RIVERA-CAMPOVERDE⁴, (Member, IEEE),
AND FRANCKLIN RIVAS-ECHEVERRÍA^{3,5}, (Senior Member, IEEE)

¹Higher School of Engineering, Science and Technology, Valencian International University, 46002 Valencia, Spain
²Computer Vision Center, Universitat Autònoma de Barcelona, Bellaterra, 08193 Barcelona, Spain
³Kauel Inc., Houston, TX 77027, USA
⁴Grupo de Investigación en Ingeniería del Transporte (GIIT), Universidad Politécnica Salesiana, Cuenca 010105, Ecuador
⁵Escuela de Ingeniería, Pontificia Universidad Católica del Ecuador Sede Ibarra, Ibarra 100112, Ecuador

Corresponding authors: Mary Vergara (maryjosefina.vergara@professoruniversidadviu.com) and Leo Ramos (leo.ramos@kauel.com)

This work was supported by the Grupo de Investigación en Ingeniería del Transporte (GIIT) at the Universidad Politécnica Salesiana, Ecuador.

[1] Vergara, M., Ramos, L., Rivera-Campoverde, N. D., & Rivas-Echeverría, F. (2023). Enginefaultdb: a novel dataset for automotive engine fault classification and baseline results. *IEEE Access*, 11, 126155-126171.

[2] Thomas, L. (2024). EngineFaultDB Dataset. GitHub. <https://github.com/Leo-Thomas/EngineFaultDB>

Additional Research

Safe Active Learning for Multi-Output Gaussian Processes

Cen-You Li
Cen-You.Li@de.bosch.com

Barbara Rakitsch
Barbara.Rakitsch@de.bosch.com
Bosch Center for Artificial Intelligence
Robert-Bosch-Campus 1, 71272 Renningen, Germany

Christoph Zimmer
Christoph.Zimmer@de.bosch.com

Abstract

Multi-output regression problems are commonly encountered in science and engineering. In particular, multi-output Gaussian processes have been emerged as a promising tool for modeling these complex systems since they can exploit the inherent correlations and provide reliable uncertainty estimates. In many applications, however, acquiring the data is expensive and safety concerns might arise (e.g. robotics, engineering). We propose a safe active learning approach for multi-output Gaussian process regression. This approach queries the most informative data or output taking the relatedness between the regressors and safety constraints into account. We prove the effectiveness of our approach by providing theoretical analysis and by demonstrating empirical results on simulated datasets and on a real-world engineering dataset. On all datasets, our approach shows improved convergence compared to its competitors.

of a machine are not supposed to crash any objects. A system should avoid generating high pressure, high temperature, or explosion. Safe learning addresses this by incorporating and learning safety constraints (Sui et al., 2015). Schreiter et al. (2015) and Zimmer et al. (2018) combine safety considerations with AL so that the data selection is done only in the determined safe domain.

These works, however, rarely considered multi-output (MO) regression problems, despite them commonly encountered in science, engineering and medicine (Xu et al., 2019; Zhang and Yang, 2021; Liu et al., 2018). In such problems, it is possible to consider individual tasks or outputs independently, but the plausibly shared mechanisms are ignored, and the performances or data efficiency might be deteriorated. Zhang et al. (2016) dealt with AL on MO models but focused on efficient computation of AL with large datasets and safe exploration was not addressed.

We consider safe AL for MO regression models that exploit the correlations. In particular, we focus on problems in which different output components may not be synchronously observed (e.g. due to different measuring cost or difficulty). MO Gaussian processes (GPs) are natural candidates for these problems (Bonilla et al., 2008; Álvarez and Lawrence, 2011; Álvarez et al., 2012; van der Wilk et al., 2020), due to their capability of capturing the correlations among different outputs and of quantifying the uncertainty.

In our work, we consider as main model the Linear Model of Coregionalization (LMC, Journel and Huijbregts (1976)), in which each output is modeled as a weighted sum of shared latent functions. Each

1 Introduction

Active learning (AL) selects the most informative data sequentially according to previous measurements and an acquisition function (Krause et al., 2008; Houlisby et al., 2011; Zhang et al., 2016). The objective is to optimize a model without labeling unnecessary data. The problem setup is closely related to Bayesian op-

[5] Li, C. Y., Rakitsch, B., & Zimmer, C. (2022, May). Safe active learning for multi-output gaussian processes. In International Conference on Artificial Intelligence and Statistics (pp. 4512-4551). PMLR.

[6] Bosch Research. (n.d.). Bosch Engine Datasets. GitHub. <https://github.com/boschresearch/Bosch-Engine-Datasets>