



# Gaussian Process Kernel Transfer Enabled Method for Electric Machines Intelligent Faults Detection with Limited Samples

Mrs.S.Savithri

Department of Mca  
Dhanalakshmi srinivasan college  
of engineering and technology

Ms.Aruna.S

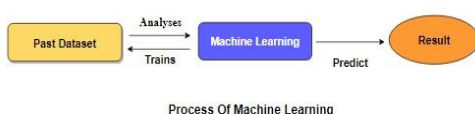
Department of Mca  
Dhanalakshmi srinivasan college  
of engineering and technology

**Abstract:** Electrical fault detection is considered as the process that finds the fault of a dataset that as power and current supply. a short circuit is a fault in which a Live wire touches a Neutral or Earth wire. An open-circuit fault occurs if a circuit is interrupted by cut on any of wires (Phase or Neutral) or blown Fuse. In three-phase systems, a fault may involve one or more phases and ground, or may occur only between phases. In a "ground fault" or "earth fault", current flows into the earth. The prospective short-circuit current of a predictable fault can be calculated for most situations. In power systems, protective devices can detect fault conditions and operate circuit breakers and other devices to limit the loss of service due to a failure. simply identify data for cleaning, before analysis. Electrical fault detection is often applied on both labelled and unlabelled data which is known as unsupervised Electrical fault detection for unlabelled data. In this project the Electrical fault detection is done in according to the dataset and the machine learning algorithms are used on the dataset and the comparison of model is done for better prediction. Machine learning is progressively being used to automate Electrical fault detection.

## I. MACHINE LEARNING

Machine learning is to predict the future from past data. Machine learning (ML) is a type of artificial intelligence (AI) that provides computers with the ability to learn without being explicitly programmed. Machine learning focuses on the development of Computer Programs that can change when exposed to new data and the basics of Machine Learning, implementation of a simple machine learning algorithm using python. Process of training and prediction involves use of specialized algorithms. It feed the training data to an algorithm, and the algorithm uses this training data to give predictions on a new test data. Machine learning can be roughly separated in to three categories. There are supervised learning, unsupervised learning and reinforcement learning. Supervised learning program is both given the input data and the corresponding labeling to learn data has to be labeled by a human being beforehand. Unsupervised learning is no labels. It provided to the learning algorithm. This algorithm has to figure out the clustering of the input data. Finally, Reinforcement learning dynamically interacts with its environment and it receives positive or negative feedback to improve its performance.

Supervised Machine Learning is the majority of practical machine learning uses supervised learning. Supervised learning is where have input variables (X) and an output variable (y) and use an algorithm to learn the mapping function from the input to the output is  $y = f(X)$ . The goal is to approximate the mapping function so well that when you have new input data (X) that you can predict the output variables (y) for that data. Techniques of Supervised Machine Learning algorithms include **logistic regression, multi-class classification, Decision Trees and support vector machines** etc. Supervised learning requires that the data used to train the algorithm is already labeled with correct answers. Supervised learning problems can be further grouped into **Classification** problems. This problem has as goal the construction of a succinct model that can predict the value of the dependent attribute from the attribute variables. The difference between the two tasks is the fact that the dependent attribute is numerical for categorical for classification. A classification model attempts to draw some conclusion from observed values. Given one or more inputs a classification model will try to predict the value of one or more outcomes. A classification problem is when the output variable is a category, such as "red" or "blue".



## II. EXISTING SYSTEM



Traditional Artificial Intelligence (AI) based fault detection approaches need a large amount of data for the model learning. However, in a real-world system, it is very difficult and expensive to obtain massive labeled fault data. In addition, the working conditions of a motor are usually variable, conventional fault diagnosis models with weak generalization ability can only be used for fault detection under constant working condition. The performance of traditional AI based approaches decreases when the working condition changes. A novel GP kernel transfer based few-shot learning method is proposed in this paper for electric machine fault diagnosis under variable working conditions. The diagnostic knowledge learns from limited current signal and vibration signal. The trained model can be transferred to other unseen working conditions without parameters updating and fine-tuning.

#### DISADVANTAGES

1. The performance is not good and its get complicated for electrical Fault.
2. The performance metrics like recall F1 score and comparison of machine learning algorithm is not done.

#### PROPOSED SYSTEM :

#### EXPLORATORY DATA ANALYSIS OF ELECTRICAL FAULT DETECTION :

Multiple datasets from different sources would be combined to form a generalized dataset, and then different machine learning algorithms would be applied to extract patterns and to obtain results with maximum accuracy.

#### DATA WRANGLING

*In this section of the report will load in the data, check for cleanliness, and then trim and clean given dataset for analysis. Make sure that the document steps carefully and justify for cleaning decisions.*

#### DATA COLLECTION

The data set collected for predicting given data is split into Training set and Test set. Generally, 7:3 ratios are applied to split the Training set and Test set. The Data Model which was created using machine learning algorithms are applied on the Training set and based on the test result accuracy, Test set prediction is done.

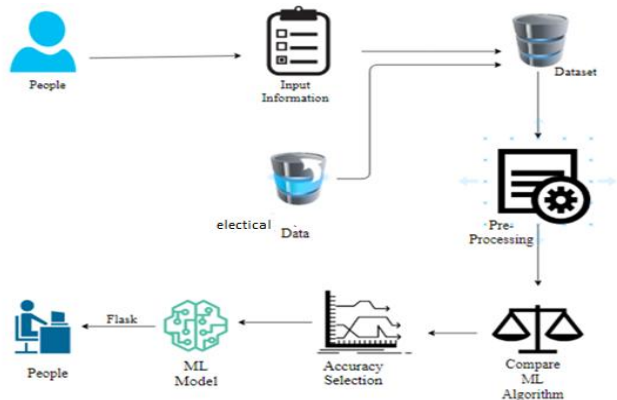
#### BUILDING THE CLASSIFICATION MODEL

The predicting the electrical fault detection, decision tree algorithm prediction model is effective because of the following reasons: It provides better results in classification problem.

- It is strong in preprocessing outliers, irrelevant variables, and a mix of continuous, categorical and discrete variables.

- It produces out of bag estimate error which has proven to be unbiased in many tests and it is relatively easy to tune with.

#### ARCHITECTURE DIAGRAM



#### INTRODUCTION TO ELECTRICAL FAULT DETECTION

#### MODULE DESCRIPTION:

##### Data Pre-processing

Validation techniques in machine learning are used to get the error rate of the Machine Learning (ML) model, which can be considered as close to the true error rate of the dataset. If the data volume is large enough to be representative of the population, you may not need the validation techniques. However, in real-world scenarios, to work with samples of data that may not be a true representative of the population of given dataset. To finding the missing value, duplicate value and description of data type whether it is float variable or integer. The sample of data used to provide an unbiased evaluation of a model fit on the training dataset while tuning model hyper parameters.

The evaluation becomes more biased as skill on the validation dataset is incorporated into the model configuration. The validation set is used to evaluate a given model, but this is for frequent evaluation. It as machine learning engineers use this data to fine-tune the model hyper parameters. Data collection, data analysis, and the process of addressing data content, quality, and structure can add up to a time-consuming to-do list. During the process of data identification, it helps to understand your data and its properties; this knowledge will help you choose which algorithm to use to build your model.

A number of different **data cleaning** tasks using Python's Pandas library and specifically, it focus on probably the biggest data cleaning task, **missing values** and it able to **more quickly clean data**. It wants to **spend less**



**time cleaning data**, and more time exploring and modeling.

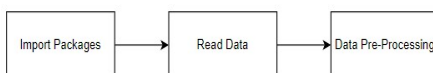
Some of these sources are just simple random mistakes. Other times, there can be a deeper reason why data is missing. It's important to understand these different types of missing data from a statistics point of view. The type of missing data will influence how to deal with filling in the missing values and to detect missing values, and do some basic imputation and detailed statistical approach for dealing with missing data. Before, joint into code, it's important to understand the sources of missing data. Here are some typical reasons why data is missing:

- User forgot to fill in a field.
- Data was lost while transferring manually from a legacy database.
- There was a programming error.
- Users chose not to fill out a field tied to their beliefs about how the results would be used or interpreted.

Variable identification with Uni-variate, Bi-variate and Multi-variate analysis:

- import libraries for access and functional purpose and read the given dataset
- General Properties of Analyzing the given dataset
- Display the given dataset in the form of data frame
- show columns
- shape of the data frame
- To describe the data frame
- Checking data type and information about dataset
- Checking for duplicate data
- Checking Missing values of data frame
- Checking unique values of data frame
- Checking count values of data frame
- Rename and drop the given data frame
- To specify the type of values
- To create extra columns

#### MODULE DIAGRAM



#### GIVEN INPUT EXPECTED OUTPUT

input : data

output : removing noisy data

#### DATA VALIDATION/ CLEANING / PREPARING PROCESS

Importing the library packages with loading given dataset. To analyzing the variable identification by data shape, data type and evaluating the missing values, duplicate values. A validation dataset is a sample of data held back from training your model that is used to give an estimate of model skill while tuning model's and procedures that you can use to make the best use of validation and test datasets when evaluating your models. Data cleaning / preparing by rename the given dataset and drop the column etc. to analyze the uni-variate, bi-variate and multi-variate process. The steps and techniques for data cleaning will vary from dataset to dataset. The primary goal of data cleaning is to detect and remove errors and anomalies to increase the value of data in analytics and decision making.

#### EXPLORATION DATA ANALYSIS OF VISUALIZATION

Data visualization is an important skill in applied statistics and machine learning. Statistics does indeed focus on quantitative descriptions and estimations of data. Data visualization provides an important suite of tools for gaining a qualitative understanding. This can be helpful when exploring and getting to know a dataset and can help with identifying patterns, corrupt data, outliers, and much more. With a little domain knowledge, data visualizations can be used to express and demonstrate key relationships in plots and charts that are more visceral and stakeholders than measures of association or significance. Data visualization and exploratory data analysis are whole fields themselves and it will recommend a deeper dive into some the books mentioned at the end.

Sometimes data does not make sense until it can look at in a visual form, such as with charts and plots. Being able to quickly visualize of data samples and others is an important skill both in applied statistics and in applied machine learning. It will discover the many types of plots that you will need to know when visualizing data in Python and how to use them to better understand your own data.

- How to chart time series data with line plots and categorical quantities with bar charts.
- How to summarize data distributions with histograms and box plots.

#### MODULE DIAGRAM



#### GIVEN INPUT EXPECTED OUTPUT

input : data

output : visualized data





#### ALGORITHM :

- Logistic Regression
- Random Forest
- Decision Tree Classifier
- Naive Bayes

The K-fold cross validation procedure is used to evaluate each algorithm, importantly configured with the same random seed to ensure that the same splits to the training data are performed and that each algorithm is evaluated in precisely the same way. Before that comparing algorithm, Building a Machine Learning Model using install Scikit-Learn libraries. In this library package have to done preprocessing, linear model with logistic regression method, cross validating by KFold method, ensemble with random forest method and tree with decision tree classifier. Additionally, splitting the train set and test set. To predicting the result by comparing accuracy.

#### PREDICTION RESULT BY ACCURACY:

Logistic regression algorithm also uses a linear equation with independent predictors to predict a value. The predicted value can be anywhere between negative infinity to positive infinity. It need the output of the algorithm to be classified variable data. Higher accuracy predicting result is logistic regression model by comparing the best accuracy.

True Positive Rate (TPR) =  $TP / (TP + FN)$

False Positive rate (FPR) =  $FP / (FP + TN)$

**Accuracy:** The Proportion of the total number of predictions that is correct otherwise overall how often the model predicts correctly defaulters and non-defaulters.

#### Accuracy calculation:

Accuracy =  $(TP + TN) / (TP + TN + FP + FN)$

Accuracy is the most intuitive performance measure and it is simply a ratio of correctly predicted observation to the total observations. One may think that, if we have high accuracy then our model is best. Yes, accuracy is a great measure but only when you have symmetric datasets where values of false positive and false negatives are almost same.

**Precision:** The proportion of positive predictions that are actually correct.

Precision =  $TP / (TP + FP)$

Precision is the ratio of correctly predicted positive observations to the total predicted positive observations. The question that this metric answer is of all passengers that labeled as survived, how many actually survived? High precision relates to the low false positive rate. We have got 0.788 precision which is pretty good.

**Recall:** The proportion of positive observed values correctly predicted. (The proportion of actual defaulters that the model will correctly predict)

Recall =  $TP / (TP + FN)$

Recall (Sensitivity) - Recall is the ratio of correctly predicted positive observations to the all observations in actual class - yes.

**F1 Score** is the weighted average of Precision and Recall. Therefore, this score takes both false positives and false negatives into account. Intuitively it is not as easy to understand as accuracy, but F1 is usually more useful than accuracy, especially if you have an uneven class distribution. Accuracy works best if false positives and false negatives have similar cost. If the cost of false positives and false negatives are very different, it's better to look at both Precision and Recall.

General Formula:

F- Measure =  $2TP / (2TP + FP + FN)$

F1-Score Formula:

F1 Score =  $2 * (Recall * Precision) / (Recall + Precision)$

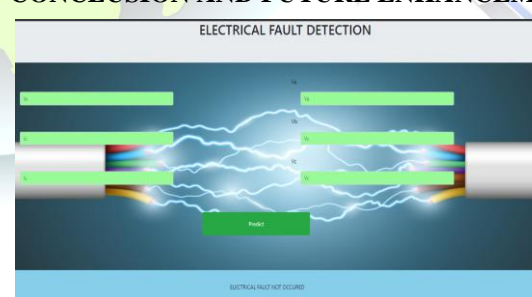
#### ALGORITHM AND TECHNIQUES :

##### Algorithm Explanation

In machine learning and statistics, classification is a supervised learning approach in which the computer program learns from the data input given to it and then uses this learning to classify new observation. This data set may simply be bi-class (like identifying whether the person is male or female or that the mail is spam or non-spam) or it may be multi-class too. Some examples of classification problems are: speech recognition, handwriting recognition, bio metric identification, document classification etc. In Supervised Learning, algorithms learn from labeled data. After understanding the data, the algorithm determines which label should be given to new data based on pattern and associating the patterns to the unlabeled new data.

##### Sample Screens

#### CONCLUSION AND FUTURE ENHANCEMENT



#### CONCLUSION:

The analytical process started from data cleaning and processing, missing value, exploratory analysis and finally model building and evaluation. The best accuracy on public test set is higher accuracy score will be find out. This application can help to find the Prediction of electrical fault or not.

#### REFERENCES



1. U. Singh, M. Rizwan, M. Alaraj, and I. Alsaidan, "A machine learning-based gradient boosting regression approach for wind power production forecasting: a step towards smart grid environments," *Energies*, vol. 14, pp. 1–21, 2021. View at: [Publisher Site](#) | [Google Scholar](#)
2. M. Zhang, C. Shen, N. He et al., "False data injection attacks against smart grid state estimation: construction, detection and defense," *Science China Technological Sciences*, vol. 62, no. 12, pp. 2077–2087, 2019. View at: [Publisher Site](#) | [Google Scholar](#)
3. Y. Wang, X. Ma, L. Zhao, H. Li, and J. Liu, "Analysis of power cable fault diagnosis and electric field detection technology based on computer control system," *Journal of Physics: Conference Series*, vol. 1574, Article ID 012080, 2020. View at: [Publisher Site](#) | [Google Scholar](#)

