TALLINN UNIVERSITY OF TECHNOLOGY

School of Information Technologies

Ahmet Caglayan 194234IASM

# Data-Driven Methods for Analysis and Fault Detection of HVAC Systems: Filter Clogging Prediction Example

Master's Thesis

Supervisor:   Eduard Petlenkov

TALLINNA TEHNIKAÜLIKOOL

Infotehnoloogia Teaduskond

Ahmet Caglayan 194234IASM

# Küte ventilatsiooni ja jahutuse (KVJ) seadmete andmepõhise analüüsi ja veatuvastuse meetodid: filtri ummistuse ennustamise näide

Magistritöö

Juhendaja:   Eduard Petlenkov

# Author's declaration of originality

I hereby certify that I am the sole author of this thesis and this thesis has not been presented for examination or submitted for defence anywhere else. All used materials, references to the literature and work of others have been cited.

Author: Ahmet Caglayan

03.01.2022

# Acknowledgement

I would like to express my gratitude to everyone who helped and supported me during the writing of my thesis.

I would like to thank my supervisor for guiding me and taking the time to answer my questions.

I would like to thank my friends who support me any time I need. Gürkan Işık, Kayahan Kaya, Caner Gür, Müge Çelebi Gür. Their support has helped me move forward through difficult times.

Lastly, I would like to express my gratitude to my mother, father, and sister, who have never lost their faith in me.

# Abstract

Today, heating ventilation and air conditioning systems (HVAC) are responsible for a big proportion of total electricity consumption in the world. An overview of the statistics for HVAC energy load for the world is given in the first chapter. Since it creates such a high rate in total energy consumption, many kinds of research are carried out on the subject.

The conventional approach for the maintenance of HVAC systems is applying periodical maintenance and changing particular spare parts in specific periods. This approach causes the unnecessary replacement of many spare parts such as bearings and filters before they complete their efficient lifespan. With the increase in awareness about energy efficiency, the industry has started to look for more efficient methods in the field of maintenance.

Along with the Industry 4.0 innovations, access to system and system components data has been provided in many industrial facilities. With the analysis of these data, it has been possible to develop much more advanced maintenance strategies. The development of data acquisition and data analysis capabilities has paved the way for better digital modelling of system behaviour and better prediction of future behaviour. Thus, predictive maintenance has become the main maintenance strategy for many industrial facilities.

In this thesis, data-driven methods will be examined and applied to the dataset of a shopping mall HVAC system. In the results section, the effectiveness of these data-driven methods will be compared. The filter system of the HVAC system will be modelled by using various machine learning, deep learning and hybrid methods. The goal of this thesis is to provide data-driven analysis to HVAC systems to be operated more efficiently and focus on fault detection with the real data which is retrieved from the system.

This thesis is 39 pages long, including 5 chapters 28 figures, 11 tables. It is written in English.

# List of abbreviations and terms

| | |
|---|---|
| AI | Artificial Intelligence |
| ANN | Artificial Neural Networks |
| CM | Corrective Maintenance |
| DL | Deep Learning |
| DPI | Dots per inch |
| DT | Decision Tree |
| FDD | Fault Detection and Diagnosis |
| HVAC | Heating, Ventilation, and Air Conditioning |
| IA | Department of Computer Systems |
| KNN | K-Nearest Neighbors |
| KPI | Key Performance Indicator |
| LDA | Linear Discriminant Analysis |
| LR | Linear Regression |
| LRN | Layer Recurrent Network |
| LWR | Locally Weighted Regression |
| MAE | Mean Absolute Error |
| MAPE | Mean Absolute Percentage Error |
| ML | Machine Learning |
| MLP | Multi-layer perceptron |
| MSE | Mean Square Error |

| | |
|---|---|
| PdM | Predictive Maintenance |
| PvM | Preventive Maintenance |
| $R^2$ | R Square |
| RF | Random Forest |
| RL | Reinforcement Learning |
| RMSE | Root Mean Square Error |
| RNN | Recurrent Neural Networks |
| SBS | Sequential Backward Selection |
| SFS | Sequential Forward Selection |
| SGD | Stochastic Gradient Descent |
| SR | Stacking Regressor |
| SVM | Support Vector Machine |
| SVR | Support Vector Regression |
| Taltech | Tallinn University of Technology |

# Table of contents

# List of figures

# List of Tables

# 1 Introduction

## 1.1 Problem Statement

The building sector takes the highest share in energy consumption. It is responsible for up to 40% of energy consumption and 36% of greenhouse gas emissions across Europe [1] and 48% to 57% in the USA [2]. HVAC systems are one of the key systems of industrial and residential areas. Commonly, it is used for providing comfortable air conditions (by heating, cooling, moistening, etc.) and waste energy recovery for residential buildings, shopping malls, cold storage for goods, production plants, and so on. Therefore it consumes a significant proportion of energy in the total electricity and gas consumption.

Over the years, maintenance and repair methods applied to HVAC systems have differed. The Most common maintenance types are considered; corrective maintenance (CM), predictive maintenance (PdM), and preventive maintenance (PvM) [3].

Corrective maintenance is the most primitive version of doing maintenance. In CM, maintenance operators intervene in systems based on their own experience or when a malfunction is encountered. Nowadays, CM is considered as a breakdown intervention rather than a maintenance strategy. Many companies accept the CM ratio in maintenance activities as a key performance indicator (KPI) and demand that this ratio be as low as possible.

Preventive maintenance means performing repair and maintenance activities to certain machinery, equipment, and machine parts within the time intervals recommended by the manufacturer, again as determined by the manufacturer, and replacing it when the determined lifespan ends [4]. PvM is still widely used in various industries. PvM has a lot of advantages, it ensures a precise lifespan, reduces malfunctions. In new and existing investments, it makes the cost of spare parts more predictable and makes it easier for companies to create a maintenance budget. It doesn't have an initial setup cost therefore it can be used for uncritical machinery to avoid predictive maintenance initial setup costs.

14

This method, which is considered a more advanced method according to CM, also has various disadvantages. The most important of these is the high cost of labour and spare parts. In the usage and maintenance manuals of industrial systems, the part replacement intervals are generally given by considering the worst conditions. In PM, many parts are actually changed before they expire. Another disadvantage is that the detection of defective parts is slower than PdM. The system will run more inefficiently until the fault is detected, and the faulty part may damage other components of the system and cause the size of the fault to increase. It may cause an increase in the maintenance and repair costs of the system.

PdM has come with the new technologies, protocols, communication systems, standards that come with industry 4.0 maintenance and repair approaches. The need to anticipate the problem before it even occurred gave rise to PdM. PdM aims to detect the instant status of the system and possible failures beforehand by analysing the current data retrieved from the system. Repair and maintenance operations are planned according to these analyses.

PdM is the most advanced maintenance approach and it provides benefits such as real-time system statistics, extends equipment lifespan, and reduces the use of spare parts. PdM needs as much data about the system as possible in order to work efficiently. For this, systems are equipped with various instruments such as temperature and pressure sensors, encoders etc. The data-driven approach is one of the most popular approaches to support PdM. This method is to model the system and train this model continuously with the system data. This model will begin to make better fault predictions according to the amount of data obtained. In this thesis, a filter system of an HVAC system will be modelled by using a real HVAC system of a shopping mall in Tallinn/Estonia.

Today, various control methodologies and algorithms are used to reduce all unnecessary consumption, reuse waste energy and develop more efficient systems. In this study, the system will be modelled with a data-driven approach using ML, DL and hybrid methods, and a method will be developed for the clogging prediction of the HVAC filter system.

## 1.2 Organization of the Thesis

In section 1.2, the methods to be used in solving the problem statement are briefly mentioned. In section 1.3, the literature on the subject has been examined.

In Section 2, the theoretical background of the methods used is explained.

In section 3, the HVAC data of a shopping mall in Tallinn is examined, the methods used to predict the filter clogging failure and the results obtained with these methods are presented.

Finally, in conclusion, section 4, the results obtained from the thesis study are contextualized, and references are included in the 5th section. Appendixes were added to the end of the thesis.

## 1.3 Methodology

In this study, as an example, we will try to develop a methodology for predicting filter clogging failure, which is one of the most common failure types in HVAC systems.

Data-driven fault detection methods were developed based on pattern classification techniques. The data-driven fault detection methods for HVAC systems can be generally categorized into statistic-based methods and artificial intelligence-based methods [5].



Figure 1 Prediction methods [6]

In this paper, an intelligent methodology will be developed for fault detection by using machine learning (ML), deep learning (DL) and hybrid models. Various algorithms will be compared according to their failure prediction performance. There are various studies in the literature on methods of comparing model efficiency [7]. To ensure the accuracy of the results, Mean Absolute Percentage Error (MAPE), Mean Absolute Error (MAE), Mean Square Error (MSE), $R^2$ coefficients will be used as a criterion.

$R^2$ score is used to make performance measurement comparisons of different types of regression models. It takes values between 0 and 1. It shows how close the data are to the fitted regression line, as can be seen from Figure 2 and Figure 3. It can be formulated as follows [8]:

$$R^2(y, \hat{y}) = 1 - \frac{\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}{\sum_{i=1}^{n}(y_i - \bar{y}_i)^2} \tag{1}$$

Where $y_i$ is the true value and $\hat{y}_i$ is the predicted value.



Figure 2 Low $R^2$ coefficient

Figure 3 High $R^2$ coefficient

MSE is the measurement of the average of the squares of the errors. Lower MSE values mean higher prediction accuracy. It can be expressed by the following formula: [9]

$$MSE(y, \widehat{y}) = \frac{1}{n} \sum_{i=1}^{n} (y_i - \widehat{y}_i)^2 \tag{2}$$

Where $y_i$ is the real value and $\widehat{y}_i$ is the predicted value.

The mean absolute error is the sum of the absolute error value, a more direct representation of the sum of the error terms.

$$MAE(y, \widehat{y}) = \frac{1}{n} \sum_{i=1}^{n} |y_i - \widehat{y}_i| \tag{3}$$

MAPE is a measure used in trend estimation to measure the validity and reliability of estimation methods used in constructing time series values.

$$MAPE(y, \widehat{y}) = \frac{100}{n} \sum_{i=1}^{n} \frac{|y_i - \widehat{y}_i|}{y_i} \tag{4}$$

The main motivation of this paper is: Evaluate and compare different types of prediction methods, especially artificial neural networks (ANN) methods and machine learning methods such as Long Short-Term Memory (LSTM) Support Vector Machines (SVM), Linear Regression (LR), and Decision Trees (DT) and so on.

## 1.4 Literature Review

Many studies have been and are being conducted on fault prediction. These studies on the subject divide fault detection methods into 3: the model-based, the rules-based, and the data-driven methods [10].

The model-based approach uses mathematical algorithms to build a model that will be the basis for evaluating differences between actual operation values and expected ones [11] [12]. The model-based approach has some challenges: The data should be pre-processed, and the amount of the data should be high. On the other hand, a rules-based approach refers to applying particular qualitative rules to define the normal operation of a system [13]. The rules-based approach is widely used in industry due to its ease of application.

The data-driven method, which is the subject of this paper, does not require a physical model or expert knowledge of the system, unlike model-based and rules-based methods [14]. With this independence, data-driven methods have shown great advantages while characterizing system operations and developing system models using real system data [15]. It can be seen that one of the biggest advantages of the data-driven approach is that the system modelling can be done very quickly compared to other approaches.

Fast and accurate results and being cost-effective lead to the increasing use of data-driven methods in the industry. In the conducted studies, it was concluded that data-driven models can replace physics-based models with an insignificant loss of predictive accuracy for many applications [16]. In the data-driven method, the relationships between inputs and outputs are determined automatically. The data-driven model consists of 3 parts, input features, output features, and training algorithms [17]. Training data generated with input data train the algorithm and enables consistent predictions

The most used methods in these studies can be listed as Genetic Algorithm (GA), Particle Swarm Optimization (PSO) [18], Ant Colony Optimization (ACO), ML, Evolutionary Programming (EP), ANN, Fuzzy Logic and so on [14] [19].

Sulaiman et al. [20] used HVAC data of 5370 rows and 14 features and normalized all the data for data pre-processing. They applied various optimization algorithms to classify HVAC faults, as an optimization algorithm, they first applied stochastic gradient descent

(SGD) and obtained 94% model accuracy score. With support vector machines (SVM), they achieved a better accuracy score of 97% compared to SGD. The multi-layer perceptron gave the best results between 99.4% among the 3 models.

Schreiber et al. [21] researched on modelling energy systems using the data-driven approach. They used random forest (RF), LR, SVR and ANN for modelling the system and compared R2 and RMSE values for the evaluation of these models. They found that LR performed worse than the others. They also found that SVR and ANN performed better on datasets whose training dataset is not very large.

Shahnazari et al. [22] used a layer recurrent network (LRN), one of the types of recurrent neural networks, for modelling and fault diagnosis in an exemplary HVAC system. As a result of their work, they showed that LRNs can cope with complexities such as nonlinearity, and they were able to improve the results obtained using identification methods such as subspace identification.

Ebrahimifakhar et al. [23] used many classifiers for fault detection and diagnosis in packaged rooftop units, SVM and LR were the best classifiers with an overall accuracy of 96.2% and 93.6%, KNN and linear discriminant analysis (LDA) were found to be the weakest classifiers with the accuracy scores of 83.6% and 76.2% in the same study.

The biggest flaws of the data-driven approach are that it requires a large amount of historical data and that data is pre-processed to be free from all noises. Therefore, the datasets used without pre-processing lead to erroneous predictions [24]. In the next chapter, prediction methods will be explained.

# 2 Theoretical Background

Although the use of traditional methods is intense, there has been an increase in the use of Artificial Intelligence (AI) as a failure prediction method in the research area and industrial facilities. The prediction methods mostly used today are ML, DL and ANN. In this paper, the focus will be on ML and DL.

## 2.1 Machine Learning and Deep Learning

The advances in information technology have enabled us to store, analyse and process large amounts of data over computer networks. Using this data, a model can be defined with various input parameters and training algorithms, and this model is optimized using historical data by executing a computer program [25].

When this definition is taken into consideration, it can be said that an ML algorithm is an algorithm that can learn from data and create a pattern from the same data. Machine learning algorithms work with the principle of training a model with data and making intelligent future predictions according to the pattern created by this model. ML is considered a branch of AI.

Deep Learning (DL), a sub-field of machine learning, is an approach inspired by the structure and functions of the human brain. As in ML, it is tried to predict the outputs with the data set given in DL. Supervised and unsupervised learning methods can be used as in ML.

### 2.1.1 Advantages of ML

While statistical methods usually try to find correlations or relationships between features, ML algorithms try to predict future values of these features with the data given to them [26].

First of all ML Algorithms can handle large amounts of multi-dimensional and complex datasets [27]. It allows for analysing and processing industrial system data in real-time.

ML algorithm is an algorithm that is able to learn from data. ML algorithm is continuously trained according to the data input feed. That makes the algorithm improve its prediction capacity with more data in time. Also, this makes the ML model more flexible than other types of prediction models.
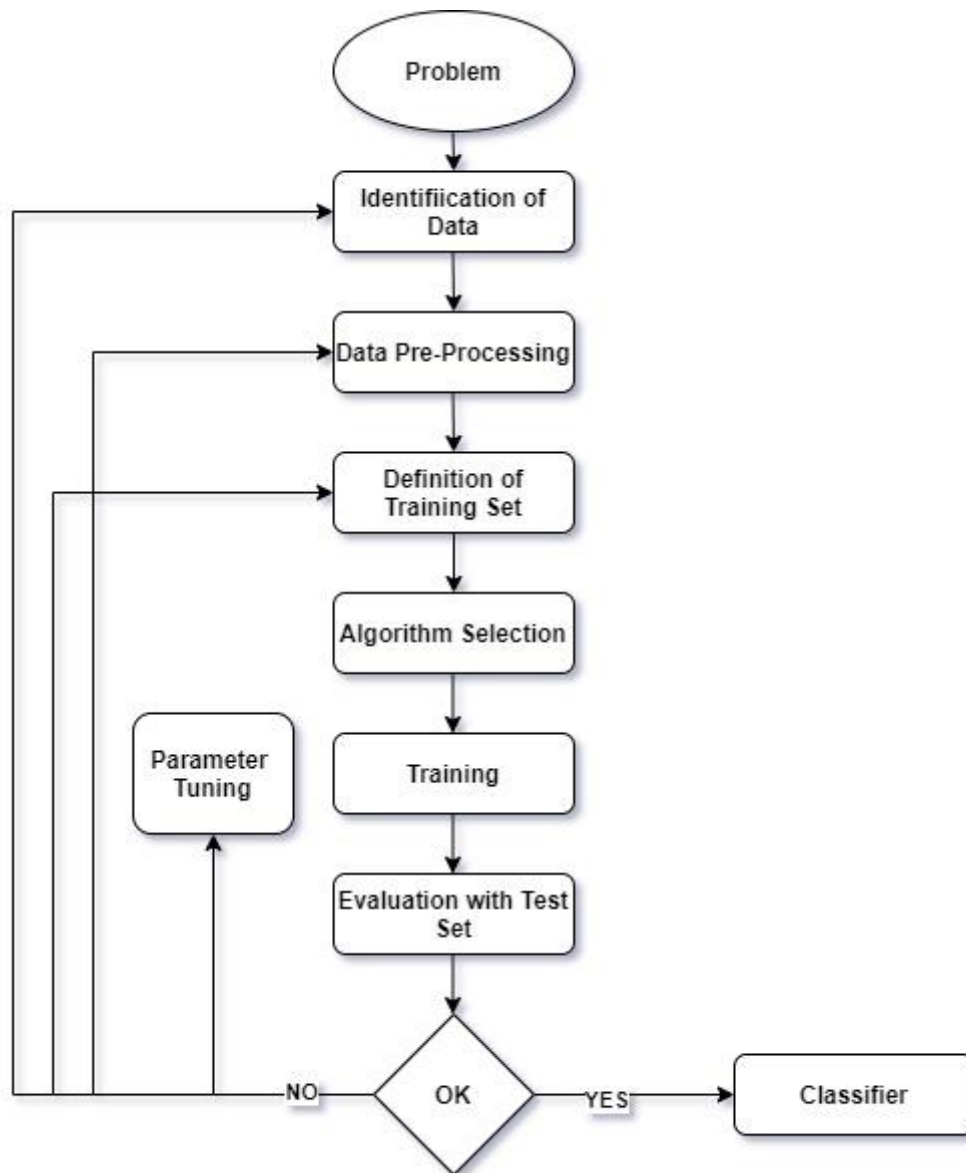
Figure 4 Machine learning process flow [28]

## 2.1.2 Types of ML Algorithms

ML algorithms are of three main types, supervised machine learning unsupervised machine learning, reinforcement learning.

**Supervised Machine Learning**

In supervised machine learning, the input features that will affect the target variable are predetermined and the model to be created is trained with these predetermined features. In the literature, supervised machine learning is mainly used to detect HVAC system faults.

**Unsupervised Machine Learning**

Unlike supervised ML, in unsupervised ML, the model is expected to self-discover the pattern found in the data. In short, it is aimed that the model decides on its own how it will learn. As it is out of the scope of this thesis, unsupervised ML will not be mentioned in the following sections.

**Reinforcement Learning**

Reinforcement learning (RL) is based on the principle that the model observes the results and makes choices based on feedback from the results. As it is out of the scope of this thesis, RL will not be mentioned in the following sections.

## 2.2 Modelling Algorithms

ML and DL algorithms help us to create models that will help us to estimate the output values based on the features. The main goal of training a model is, minimizing the least squared error and finding the best fitting line. The input features are called independent or explanatory variables, and outputs are called dependent variables.

The main purpose of modelling algorithms is to represent a real-life problem digitally. While doing this, our model should have low bias and low variance values for its performance to be high.

The difference between the estimate of our model and the actual value is called bias. High bias is a sign that our model does not fully learn the data and will make erroneous predictions. Figure 5 shows an underfitted model with a high bias.

Figure 5 High Bias

High variance indicates that our model memorizes the data instead of learning it by generalizing it so that it can make meaningful predictions from the data. The model cannot make accurate predictions when it encounters a value that is not existing in the Dataset. Figure 6 shows an overfitted model with high variance.



Figure 6 High variance

To optimize the behaviour of the model, overfitting and underfitting should be avoided. Figure 7 shows an example of an optimized model.



Figure 7 Low bias low variance

### 2.2.1 Linear Regression

One of the most widely used regression types in industry and finance is linear regression. In linear regression, the aim is to define the relationship between independent variables and dependent variables with the help of a coefficient. Linear regression can be formulated with the following expression.

$$y = Ax + B + e \tag{5}$$

With this equation, a linear model can be created. In this model, $A$ is the intercept of the line, B is the slope, and e is the margin of error.

The linear regression aims to find these A and B values, to express the linear relationship that best fits the given data. As can be seen in Figure 8, the best fitting line is found by minimizing the least squared error.



Figure 8 an example of linear regression

### 2.2.2 Ridge, Lasso and Elastic Net Regression

With the increase in the amount of data and the development of analysis methods, the number of features used in the data sets has increased proportionally. An increase in the number of features may lead to an increase in model complexity as well, which may cause the model to memorize rather than learn. Although the training scores of the memorizing model are high, the prediction capacity will be lower. Techniques such as Ridge and Lasso are used to overcome this situation, which we call overfitting.

Ridge regression is a linear regression model. In these regression techniques, the coefficients are again trained from the training data but fitted with a constraint value. Ridge and lasso restrict the model more. The complexity of the more constrained model decreases. The less complex model performs worse on the training data, but the generalization of the model is better. As the alpha value gets larger, the coefficients entering the model are limited.

Lasso is also a linear regression model. Because of the regularization term added to the linear regression using absolute value. Lasso is also called L1 regulation. Lasso's logic is very similar to ridge regression. Lasso significantly reduces the number of features.

Elastic net is a linear regression model that is a combination of ridge and lasso. Elastic net uses both L1 and L2 regularization. By applying the L1 and L2, elastic net reduces the magnitude of the regression coefficients.

**2.2.3 K-Nearest Neighbors**

Today, K-Nearest Neighbors (KNN) is widely used in classification problems. In addition, good results are obtained in the analysis of continuous data and time series. KNN works on the principle of estimating the target value according to the class intensities of the nearest neighbors, which are made up of independent variables. 'K' refers to the number of neighbors to be used in estimation. For example, in the case where K is 1, the target value is assigned to the closest neighbor's class.

In the example seen in Figure 9, it is visualized how the model predicts the classes of randomly placed star objects when the parameter for the number of nearest neighbors is given as 5.

Figure 9 KNN classification example
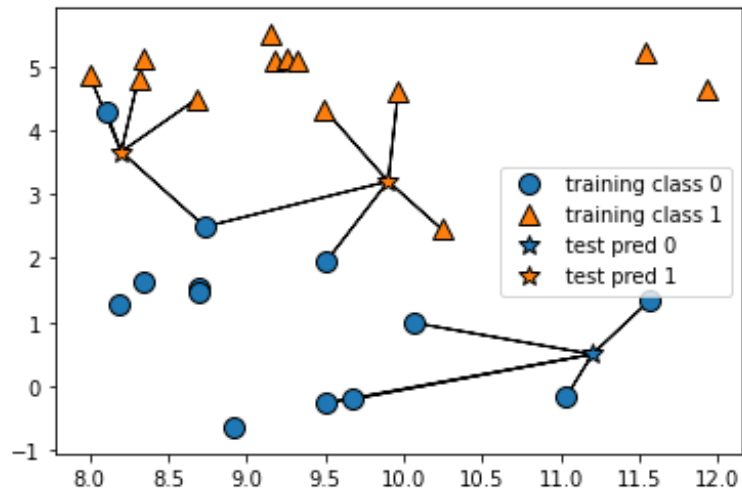
In Figure 10, the working principle of KNN regression can be seen. The determination of the location of 3 randomly selected points according to 5 neighbourhood relations is visualized.



Figure 10 KNN regression example

The model established with a small number of neighbors is relatively complex model, but the accuracy rate of this model is low. The model established with more neighbour

number parameters is simple, but this model may also have performance problems. In KNN models, hyperparameter tuning methods are used to optimize the performance and accuracy of the model.

In KNN, Manhattan, Minkowski and Euclidean equations are generally used to measure distance. The standard Euclidean distance is the most common choice [29]. Euclidean distance function can be expressed as follows:

$$d(x, y) = \sqrt{\sum_{i=1}^{k}(x_i - y_i)^2} \tag{6}$$

### 2.2.4 Decision Trees

Decision tree (DT) algorithms are one of the most used supervised learning methods and it is one of the closest methods to human thinking. As the name suggests, a tree structure is created. Conditions used to form the branches of this tree.

They can be used to solve classification and regression problems and can model nonlinear relationships better than linear models. The advantages of DT algorithms are that they are fast and simple to interpret. It can model very quickly even systems with large datasets.

### 2.2.5 Support Vector Machines

Support vector machines (SVM) are the division and classification of points on a plane or space by a line or plane. SVM was introduced to literature by Vladimir N. Vapnik in 1995 [30]. It is generally suitable for small and medium data modelling. SVM tries to choose the line that separates two different sets of values, providing the furthest possible range. This line can be linear, polynomial, or radial.

In this thesis, support vector regression (SVR) is used for modelling the HVAC system. SVR is the regression implementation of SVM and shares the same principles. Linear and non-linear SVR examples can be seen in Figure 11.

Figure 11 Linear and non-linear SVR [31]

## 2.2.6 Least Angle Regression (LAR)

Least angle regression (LAR) is a regression for high dimensional data which is first introduced by Efron et al [32]. Least Angle Regression is a stylized version of the Stagewise approach that accelerates computations by using a basic mathematical formula [32]. LAR finds the feature closest to the target at every step. It moves in equal angular direction between the features instead of continuing with the same attribute when there are many equally correlated features [33].

## 2.2.7 Long Short Term Memory (LSTM) Networks

LSTMs are a type of recurrent neural network. They can remember the Long-term relationships present in the dataset. They consist of 4 interconnected layers. These are the forget gate, input gate, output gate and cell state. In LSTMs, with the help of forget gates, unimportant inputs are forgotten and important ones are kept. The information from the previous cell and the current information are inserted into the activation function. As a result of this function, a value between 1 and 0 is obtained. Results with 0 are forgotten and results with 1 are transferred to the input gate with the cell state structure. The input gate decides whether to update the input. Finally, the output gate provides the input of the next cell.

Figure 12 Example of an LSTM cell [34]

When using LSTM, it should be kept in mind that LSTM is a prone approach to overfitting. One of the biggest disadvantages compared to other modelling algorithms is that the model training time is longer.

### 2.2.8 Gated Recurrent Units (GRU)

GRU which is one of the RNN, got its name from the gates that regulate the information flow. A body of a GRU cell is similar to LSTM, the difference of GRU, which is a newer algorithm compared to LSTM, LSTM contains 3 gates in the LSTM cell, since there are 2 gates in the GRU cell, reset and update gates.

Figure 13 Example of a GRU cell [35]

Update gate decides which data will be remembered, reset gate decides which data will be forgotten. According to Chung et al. [36] GRU-RNN provides faster progress in terms of both updates and processor time. However, they also stated that their results would not be conclusive in comparing LSTM and GRU.

### 2.2.9 Multi-layer Perceptron

One of the most used DL approaches today is the multi-layer perceptron (MLP). A simple MLP model consists of interconnected layers. Neurons in these layers are linked to all subsequent neurons (Figure 14).

The first layer is the input layer where the inputs are obtained. The hidden layer follows it, the number of layers to be found in the hidden layer changes according to the problem to be solved. The model is initialized by using the number of hidden layers as hyperparameters. Finally, the outputs are obtained from the output layer.



Figure 14 MLP design example

While creating the model, the relationships between the neurons are randomly weighted and the inputs are multiplied by the weights. As these weights increase, the importance of the input value also increases. After this step, the weight inputs are summed and a net input is given to each neuron as below [37].

$$Net = \sum_{i=1}^{n} W_i X_i \qquad (7)$$

Where Z is the number of inputs, $Wi$ is the weighted matrix from layer i to j and $Xi$ is the input of any neuron.

Various activation functions are used to activate neurons ('identity', 'logistic', 'tanh', 'relu'). The output of the current processing element is found by passing the net input

activation function, and this output value is sent to the neurons in the next layer. This process continues until the output layer is output.

Considering the error value between the expected output and the actual one, the error is distributed and the model is optimized by changing the weight values in the next iterations.

## 2.2.10 Hybrid Method

Hybrid methods use the principle of combining several methods to solve a modelling problem. In this study, the stacking regressor method from the sklearn library will be used.

**Stacking Regressor (SR)**

Stacking regressors (SR) is a collective learning method in which several regression models are combined using a meta-regressor. After the regression models are trained, the meta-regressor is fit based on the outputs of the individual regression models in the ensemble [38].

With SR, it is aimed to use the strengths of many regression models in a single model. The SR can be designed with two or more layers. Predictions from the first layer create the inputs of the second layer, and this is repeated for other layers if any. Figure 15 demonstrates a basic 2 layer SR.



Figure 15 Overview of stacking regression [38]

# 3 Case Study

In this section, the dataset of the HVAC system of a shopping center in Tallinn/Estonia, provided by the Tallinn University of Technology (Taltech), will be examined. Based on the given data, it is aimed to develop a methodology for the prediction of filter clogging, which is one of the most frequently encountered failures in HVAC systems.



Figure 16 Case study HVAC system

## 3.1 Understanding the Filter Clogging

Filtration is the process of separating the undesirable parts (dust, metal particles, etc.) from a fluid in a certain process. The purpose of filtration is to operate the existing system in a steady state and to minimize the malfunctions that may occur in the system. When the filters lose their filtration capability over time and become fluid-tight, it is called filter clogging. There are many factors that can affect filter clogging. At the beginning of these, factors such as the amount and sizes of particles in the fluid, fluid pressure, fluid flow,

filter material (fabric, fibre, polyester, etc.) and other specs, ambient temperature and humidity come.



High Pressure      Filteration      Low Pressure

Figure 17 Pressure drop

One of the most basic components of filter units is filter media. Filter media can be produced from many materials and in different pore sizes, depending on the type of fluid. The purpose of the filter media is to keep the unwanted particles in the fluid. These particles accumulate on the filter over time and cause the difference between the pressures measured from the filter inlets and outlets to increase. These particles are called particle/filter cakes [39].

In order to understand filter clogging, the concept of pressure drop must first be understood. The pressure drop value in the filters consists of two components. The first of these is the initial pressure drop value ($\Delta P_0$) that the filter has. The second is the pressure drop caused by the filter cake formed on the filter ($\Delta P_c$). The total pressure drop can be formulated as follows [40].

$$\Delta P = \Delta P_0 + \Delta P_c \tag{8}$$

## 3.2 Understanding the Dataset

Given data covers the dates between 07.09.2019 to 20.05.2020. System data was read and recorded at 15-minute intervals. When the given dataset is examined, it is seen that the system is activated at 6:35 on weekdays and activated at 7:35 on weekends. The system

35

is deactivated at 22:20, 20 minutes after 22:00, which is the closing time of shopping centers in Tallinn. In addition, it has been observed that, exceptionally, on some dates, a 24-hour working order has been adopted.

A total of 139 different data entries are kept in the dataset used. Most of them are numeric data read from sensors. In addition, features derived from this data are used by the system. However, 90 of the features in the dataset consist of null or fixed values. These features can be removed from the dataset as they can be ignored. This will be discussed in detail in Section 3.4.

From Figure 18 it can be seen that there are some radical changes in the pressure drop trend. Many factors can cause this situation. For example, some of the dust accumulated on the filter is scattered by the vibration that occurs during the commissioning of the system, the filter is cleaned and installed in its place during filter maintenance, or the installation of a filter with different micron values or produced from different materials instead of the existing filter can be counted among these factors.

In this study, while examining the pressure drop trend, it was decided to examine a selected part of the existing dataset rather than using it as a whole. In this way, the model to be created will be able to make more accurate predictions. The selected part can be seen in Figure 18 as part 4.

In the first runs of each day, deviations are observed in the pressure drop values until the system starts to work stably. Also, as can be seen in Figure 19, there are many noisy sensor values. These will be examined in Section 3.3 and in the following sections.

Figure 18 Pressure drop on return air filter

## 3.3 Introducing Software Tools

Python 3 is the most frequently used language in subjects such as data analysis, data mining, and data science today. Being open-source code is preferred by many software engineers. Python has paved the way for data scientists to handle many complex tasks very quickly, with many libraries developed for AI, data science, ML, and DL. All these reasons led to the decision to use python in this thesis.

In this thesis, Jupyter Notebook and Spyder are preferred as the software environment. Jupyter notebook and Spyder are the software development environments that have gained popularity today with their support for various programming languages and their open-source code. Compared to other environments, the greatest ease of use is that they divide the code into parts, allowing the desired part to be run at any time.

Many python libraries were used in the stages of ML model creation, data visualization, application of various numerical algorithms to the model and data visualization.

Numpy is one of the basic packages used in scientific calculations. Numpy has various functions for working with Linear algebra, matrices and arrays. Pandas library helps to organize the collected data with various functions and make it suitable for analysis. One of the most important python packages for those working in the field of data science is the scikit-learn package. It is used in this thesis as it offers basic functions such as

preprocessing the data, creating data models, generating training data and test data and measuring their performance as a ready-made package. Matplotlib and seaborn libraries are used to benefit from the quick solutions they offer for data visualization.

In addition, the mlxtend library developed by Sebastian Raschka, a professor of statistics, was used to generate the hybrid models.

## 3.4 Data Pre-processing

As an initial step, features that contain a lot of null or zero values should be eliminated. Secondly, the variables that have low variance or the ones that have no variance should be removed from the model. This may lead to faulty training of the model. In this study case dataset, a lot of columns are encountered which matches this issue. These columns were eliminated. Detailed information about these columns can be found in Appendix 6.

Looking at Figure 19, bars that extend downwards and repeat at fixed intervals are noticeable. Upon careful examination of the dataset, it was found that these were the first data values read in the HVAC system each day. When the system is first commissioned, it must take a while for it to switch to steady-state and for the system operating parameters to reach their working values. Therefore, the initial running values of the system can be considered as faulty values and these can be eliminated to create a better model.



Figure 19 Return air pressure drop data before noise cancellation

When all the above-mentioned erroneous values are cleared, the data graph in Figure 20 is obtained. Data can be further optimized with various noise cancellation methods. Since it is not the subject of this thesis, this subject will not be discussed.



Figure 20 Return air pressure drop data after noise cancellation

## 3.5 Input Feature Selection

In ML and neural network modelling algorithms, input variables are called independent or explanatory variables, and output variables are called dependent variables. As a rule of thumb, input variables are denoted by an uppercase 'X' and output variables are denoted by a lowercase 'y'.

Feature selection process is the process of selecting the 'X' variables which are highly correlated with 'y' variable. The purpose here is to select the input features that will maximize the training and testing capabilities of the created model.

Before starting the input feature selection, there are a few more steps to be followed. Among the remaining variables, those with a high correlation value should be eliminated. And, variables with low correlation with output should not be included in the model. In Table 1, the correlation of all input features with output is indicated. There are three things to note in this table. The first is which input variables have a high correlation with the output variable (pressure drop over return air filter), the second is which input variables have a high correlation coefficient among themselves. Lastly, some features may behave similarly to target variable for particular time interval although they are not related to each other. The code that is used to obtain the results below can be found in Appendix 5.

39

Table 1 Correlation between features and target variable

| All Features | Pressure drop over return air filter_FILTER_RETURN |
|---|---|
| Pressure drop over return air filter_FILTER_RETURN | 1 |
| Total working hours_FAN_SUPPLY | 0.913713982 |
| Total working hours_FAN_RETURN | 0.913703339 |
| Return air temperature__ | 0.755227301 |
| AT_57933 | 0.669879321 |
| Return air pressure diff._HEATRECOVERY_None | 0.582320555 |
| AT_128028 | 0.576962948 |
| AT_58069 | 0.575678682 |
| AT_52520 | 0.549792959 |
| AT_57885 | 0.549766805 |
| AT_52629 | 0.549766447 |
| AT_57952 | 0.549762562 |
| AT_52544 | 0.549762051 |
| AT_52603 | 0.549723155 |
| AT_57850 | 0.549719909 |
| AT_57867 | 0.549706099 |
| AT_52616 | 0.549702231 |
| AT_52537 | 0.549692861 |
| AT_52432 | 0.548850213 |
| AT_52522 | 0.547367693 |
| AT_57833 | 0.547316719 |
| Static pressure_FAN_SUPPLY | 0.483332096 |
| AT_57905 | 0.462984952 |
| Pressure drop over supply air filter_FILTER_SUPPLY | 0.363771216 |
| Fan speed_FAN_SUPPLY | 0.353410335 |
| Air volume_FAN_RETURN | 0.295255954 |
| VFD frequency_FAN_SUPPLY | 0.251672562 |
| Coil temperature_COIL_HEATING | 0.244688324 |
| 10 Rows Difference | 0.221628017 |
| Supply air temperature__ | 0.193857574 |
| Electrical power_FAN_SUPPLY | 0.155599592 |
| A Row Difference | 0.150828664 |
| Supply air temp._HEATRECOVERY_None | 0.02794798 |
| Air volume_FAN_SUPPLY | -0.009004919 |
| Electrical power_FAN_RETURN | -0.150323201 |
| problem | -0.226546891 |
| VFD frequency_FAN_RETURN | -0.351632368 |
| Fan speed_FAN_RETURN | -0.388059602 |
| Valve opening_COIL_HEATING | -0.499017652 |
| BARH_57823 | -0.538280222 |
| BARH_57857 | -0.539249229 |
| BARH_57920 | -0.539283821 |
| BARH_57943 | -0.539349875 |
| BARH_57875 | -0.539355851 |
| BARH_57841 | -0.539465236 |
| BARH_57896 | -0.539544488 |
| Efficiency_HEATRECOVERY_None | -0.571841195 |
| Static pressure_FAN_RETURN | -0.583245838 |
| Rotation speed_HEATRECOVERY_ROTARY | -0.671135855 |

## 3.5.1 Selection of Features with High Correlation with Target

The correlation value changes between -1 and 1 and shows the relationship between the variables. -1 in the negative direction and +1 in the positive direction is an indicator of correlation. In order to increase the prediction performance of the model to be created, the features with the highest correlation coefficient in the negative or positive direction should be used.

As can be seen from Table 1, the features with high correlation value can be extracted as below, correlation heat map can be seen in Figure 21.

Table 2 Highly correlated input features

| Input Variable | Correlation Coefficient with output |
|---|---|
| Total working hours_FAN_SUPPLY | 0.913713982 |
| Return air temperature__ | 0.755227301 |
| Rotation speed_HEATRECOVERY_ROTARY | -0.671135855 |
| AT_57933 | 0.669879321 |
| Static pressure_FAN_RETURN | -0.583245838 |
| Return air pressure diff._HEATRECOVERY_None | 0.582320555 |
| AT_128028 | 0.576962948 |
| Efficiency_HEATRECOVERY_None | -0.571841195 |
| BARH_57896 | -0.539544488 |

Figure 21 Correlation heat map of the features

### 3.5.2 Elimination of Input Features with High Correlation with Each Other

As can be seen from Table 1, there is a high correlation between AT sensors, likewise, a high correlation can be seen between BAR sensors as well. This situation is also visualized using a heat map in Figure 22 and Figure 23. In such cases, instead of training the model with all the features, the effect of these features on the target variable should be investigated and duplicate or highly correlated variables should be eliminated.

Figure 22 Correlation heat map of AT sensors



Figure 23 Correlation heat map of BAR sensors

**Input Feature Selection Using Sequential Forward Selection Method**

Sequential feature selection algorithms are a class of greedy search techniques for condensing a d-dimensional feature space into a k-dimensional feature subspace with k d dimensions [41]. Sequential feature selection algorithms measure the relationships of features with their target value and sort them according to their relationship ratio. The use of features with a high ratio of target value increases the generalization ability of our model and prevents unnecessary noise.

When the effect of AT sensors on the output was examined using the sequential forward selection(SFS) method, it was seen that the performance calculated by taking into account the standard errors of the cross-validation scores increased in the m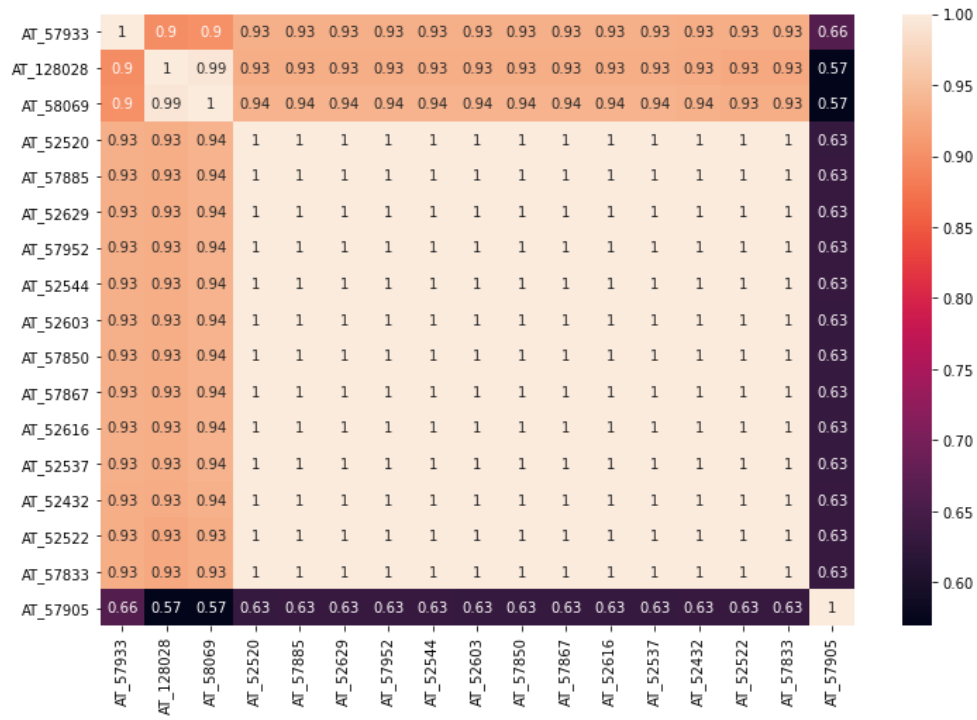odel created by using the AT_57933 and AT_128028 sensors as input variables. However, when other sensor inputs were used, the model performance remained stable and then decreased as seen in Figure 24.



Figure 24 Sequential forward selection of AT sensors

When the same method is applied for BAR sensor inputs, it is observed that there is no effect on the performance of the model as the number of BAR sensor inputs used to train the model increases as can be seen in 24. A similar correlation exists between the Total working hours_FAN_SUPPLY and Total working hours_FAN_RETURN variables. It will be sufficient to use one of these variables while creating the model.

Figure 25 Sequential forward selection of BAR sensors

As a result, using 2 of them instead of 17 AT input variables, 1 instead of 7 BAR input variables and 1 instead of 2 Total working hours variables will be sufficient for training the model.

Finally, when the SFS method is used for the input variables with the highest correlation coefficient with the target variable, it can be seen in Figure 26 that while the first 6 input variables increase the model performance, the variables that will be added to the model later do not have a positive effect on the model performance.



Figure 26 Input feature performance

### 3.5.3 Elimination of Features with High Correlation but Low Impact

Some independent variables may have very similar behaviour to the dependent variable. In order to be sure of its effects on the dependent variable, the same positive correlation relationship should exist in the entire dataset. In order to make a good analysis, the dataset is divided into 4 parts as Figure 27.



Figure 27 Dividing the data set into 4 parts

In the figure above, parts represents following dates:

- 1.Part: 09.07.2019 to 21.10.2019
- 2.Part: 01.11.2019 to 02.12.2019
- 3.Part: 02.12.2019 to 28.01.2020
- 4.Part: 29.01.2020 to 20.05.2020

In Table 3, the correlation value of the previously selected independent variables with the pressure drop variable at different time intervals is seen. It is expected that an independent variable, which is highly correlated with the target, maintains its positive or negative correlation status, while its numeric value does not change significantly.

As can be seen from the tables, the only variable that complies with this rule and has a high correlation with the target variable is the "Total working hours_FAN_SUPPLY"

supply variable. Consequently, this variable will be used to train the model and all other variables will be eliminated.

Table 3 Correlation tables for different time intervals

**1.Part**

| Input Features | Corr. With Target |
|---|---|
| Total working hours_FAN_SUPPLY | 0.601926985 |
| Return air temperature__ | -0.602355745 |
| Rotation speed_HEATRECOVERY_ROTARY | 0.365534728 |
| AT_57933 | -0.525205464 |
| Static pressure_FAN_RETURN | 0.05318627 |
| Return air pressure diff._HEATRECOVERY_None | 0.526997379 |
| AT_128028 | |
| Efficiency_HEATRECOVERY_None | 0.357851103 |
| BARH_57896 | 0.427355477 |

**2.Part**

| Input Features | Corr. With Target |
|---|---|
| Total working hours_FAN_SUPPLY | 0.852034418 |
| Return air temperature__ | -0.356855503 |
| Rotation speed_HEATRECOVERY_ROTARY | 0.38200962 |
| AT_57933 | -0.210384659 |
| Static pressure_FAN_RETURN | 0.02083554 |
| Return air pressure diff._HEATRECOVERY_None | -0.091929588 |
| AT_128028 | |
| Efficiency_HEATRECOVERY_None | 0.407841967 |
| BARH_57896 | -0.013959545 |

**3.Part**

| Input Features | Corr. With Target |
|---|---|
| Total working hours_FAN_SUPPLY | 0.801843231 |
| Return air temperature__ | -0.272353879 |
| Rotation speed_HEATRECOVERY_ROTARY | 0.210643437 |
| AT_57933 | -0.22283112 |
| Static pressure_FAN_RETURN | 0.1515154 |
| Return air pressure diff._HEATRECOVERY_None | -0.159531916 |
| AT_128028 | -0.183241093 |
| Efficiency_HEATRECOVERY_None | 0.328731508 |
| BARH_57896 | -0.158888328 |

**4.Part**

| Input Variable | Corr. With Target |
|---|---|
| Total working hours_FAN_SUPPLY | 0.913713982 |
| Return air temperature__ | 0.755227301 |
| Rotation speed_HEATRECOVERY_ROTARY | -0.671135855 |
| AT_57933 | 0.669879321 |
| Static pressure_FAN_RETURN | -0.583245838 |
| Return air pressure diff._HEATRECOVERY_None | 0.582320555 |
| AT_128028 | 0.576962948 |
| Efficiency_HEATRECOVERY_None | -0.571841195 |
| BARH_57896 | -0.539544488 |

## 3.6 Hyperparameter Tuning

So far, the focus has been on the selection of features to be used in building the model. However, especially with the development of DL approaches, the design of ANN consisting of many layers, the number of layers, the number of neurons, optimization algorithms, and selection of activation functions have gained importance.

The parameters that determine the behavior and scope of the estimator objects used are called hyperparameters. The number of hidden layer sizes in MLP or L1 norm in Lasso, Kernel function in SVM can be shown as an example of hyperparameters. One fit for all approach in the selection of these parameters is not valid for every data set. The performances of hyperparameters also change according to the type of problem, the size

of the data set and many other different conditions. Hyperparameter tuning is the general name given to the selection of parameters that will provide optimum model performance.

There are various methods for hyperparameter selection and hyperparameter tuning. In this study, two of them will be used, grid search and random search methods. The code for Hyperparameter tuning can be found in Appendix 7.

### 3.6.1 Grid Search

Hyperparameters can have many different value ranges. By using intuitive or historical information about the problem, it can determine the range for the values that hyperparameters can take. With grid search, the model is trained with all the specified hyperparameters and the result is observed. The hyperparameter values that give the best results are selected and the model is created with it.

### 3.6.2 Random Search

The random search was brought to the literature by Bergstra and Bengio [42] for the first time. Parameter ranges are determined first in a random search, as in a grid search. Unlike grid search, random hyperparameter groups are created instead of trying all combinations. Then the hyperparameter group that gives the best result is found and the model is trained with it.
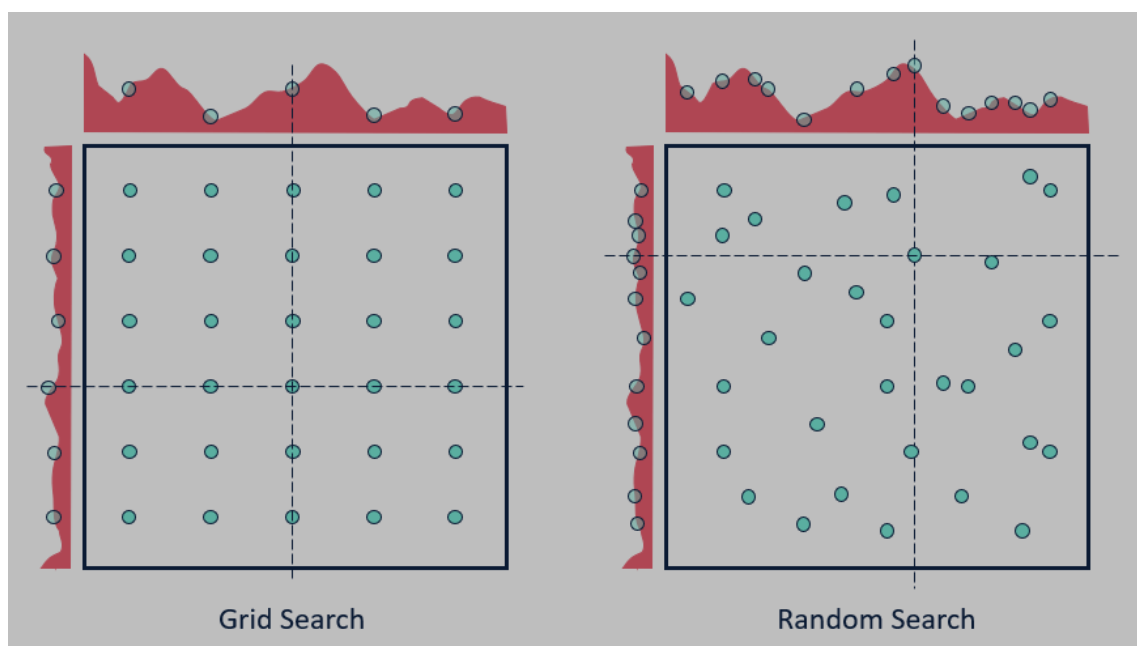


Figure 28 Grid search vs. random search [43]

## 3.7 Creating Models

As mentioned before, in order to obtain a more efficient analysis from the dataset, the existing dataset is divided into 4 parts as in Figure 27. These parts are modelled separately with the algorithms described in Section 2.2. Input feature selection is completed and hyperparameters are tuned in Section 3.5 and Section 3.6 respectively. The codes used to create the model have been added to Appendix 8, 9, 10. The performances of the models are compared in Section 3.8.

## 3.8 Comparison of the Results

### 3.8.1 Data Between 09.07.2019 - 21.10.19

Table 4 Data properties

| Average | Variance | Stdev | Number of Rows |
|---------|----------|-------|----------------|
| 63.58089 | 50.67619 | 7.118721 | 6537 |

Table 5 Comparison of models

| Scores | MSE | MAE | MAPE |
|--------|-----|-----|------|
| Linear | 9.9810 | 2.6410 | 0.0409 |
| Lasso | 10.7393 | 2.7381 | 0.0421 |
| Ridge | 22.0033 | 3.7564 | 0.0574 |
| ElasticNet | 4.4395 | 1.6405 | 0.0224 |
| SVR | 1.6919 | 1.0186 | 0.0162 |
| KNN | 0.6305 | 0.5942 | 0.0095 |
| SGD | 10.1597 | 2.6755 | 0.0414 |
| Decision Tree | 0.7324 | 0.6119 | 0.0097 |
| Lars | 10.4607 | 2.7190 | 0.0418 |
| MLP | 16.3932 | 3.3370 | 0.0504 |
| LSTM | 0.8032 | 0.6655 | 0.0088 |
| GRU | 0.7995 | 0.6599 | 0.0087 |
| Hybrid | 0.6889 | 0.6024 | 0.0096 |

### 3.8.2 Data Between 01.11.19 - 02.12.19

Table 6 Data properties

| Average | Variance | Stdev | Number of Rows |
|---------|----------|-------|----------------|
| 73.87129 | 29.67786 | 5.447739 | 1907 |

Table 7 Comparison of models

| Scores | MSE | MAE | MAPE |
|---|---|---|---|
| Linear | 3.9980 | 1.5528 | 0.0210 |
| Lasso | 4.2569 | 1.6068 | 0.0217 |
| Ridge | 11.6616 | 2.9726 | 0.0404 |
| ElasticNet | 4.4395 | 1.6405 | 0.0224 |
| SVR | 1.8110 | 1.0649 | 0.0145 |
| KNN | 0.7847 | 0.6749 | 0.0092 |
| SGD | 3.9910 | 1.6176 | 0.0221 |
| Decision Tree | 0.9013 | 0.7278 | 0.0100 |
| Lars | 3.7159 | 1.5367 | 0.0209 |
| MLP | 4.1486 | 1.5806 | 0.0216 |
| LSTM | 2.2188 | 1.1697 | 0.0147 |
| GRU | 1.3063 | 0.8912 | 0.0111 |
| Hybrid | 0.8529 | 0.6919 | 0.0093 |

## 3.8.3 Data Between 02.12.19 - 28.01.20

Table 8 Data properties

| Average | Variance | Stdev | Number of Rows |
|---|---|---|---|
| 29.75532 | 3.507557 | 1.872847 | 3447 |

Table 9 Comparison of models

| Scores | MSE | MAE | MAPE |
|---|---|---|---|
| Linear | 1.2433 | 0.8791 | 0.0297 |
| Lasso | 1.1907 | 0.8670 | 0.0293 |
| Ridge | 1.7338 | 1.0353 | 0.0345 |
| ElasticNet | 1.1796 | 0.8582 | 0.0291 |
| SVR | 0.7283 | 0.6383 | 0.0216 |
| KNN | 0.4341 | 0.4961 | 0.0167 |
| SGD | 1.1896 | 0.8576 | 0.0291 |
| Decision Tree | 0.4646 | 0.5033 | 0.0169 |
| Lars | 1.2845 | 0.8948 | 0.0302 |
| MLP | 1.4721 | 0.9223 | 0.0308 |
| LSTM | 0.6199 | 0.6375 | 0.0194 |
| GRU | 0.5845 | 0.6099 | 0.0186 |
| Hybrid | 0.4435 | 0.481 | 0.0162 |

### 3.8.4 Data Between 29.01.20 - 20.05.20

Table 10 Data properties

| Average | Variance | Stdev | Number of Rows |
|---|---|---|---|
| 25.32246 | 4.139173 | 2.034496 | 6152 |

Table 11 Comparison of models

| Scores | MSE | MAE | MAPE |
|---|---|---|---|
| Linear | 0.6781 | 0.6640 | 0.0268 |
| Lasso | 0.6651 | 0.6528 | 0.0262 |
| Ridge | 0.6884 | 0.6696 | 0.0268 |
| ElasticNet | 0.6704 | 0.6578 | 0.0266 |
| SVR | 0.3854 | 0.4633 | 0.0186 |
| KNN | 0.2819 | 0.4005 | 0.0159 |
| SGD | 0.6802 | 0.6634 | 0.0269 |
| Decision Tree | 0.3049 | 0.4242 | 0.0168 |
| Lars | 0.7111 | 0.6800 | 0.0274 |
| MLP | 0.6925 | 0.6767 | 0.0270 |
| LSTM | 0.2072 | 0.3599 | 0.0131 |
| GRU | 0.2463 | 0.3943 | 0.0142 |
| Hybrid | 0.2835 | 0.4056 | 0.0161 |

### 3.8.5 Interpretation of Results

In this study, in the comparison made considering MAE, MSE, MAPE, it was observed that KNN gave the best results for data with different sizes and variances in the estimation of clogging in HVAC system filters. It was seen that KNN, DT and hybrid models give good results, regardless of data size and variance.

The hybrid model was created using KNN and DT, which gave the best results as regressors, and SVR as meta-regressors. The hybrid model gave results close to KNN, which gave the best results in its structure.

In general, the performance of linear models is lower than others. As expected, linear models give worse results as variance increases.

The performance of LSTM and GRU models increased in direct proportion to larger data and inversely proportional to variance. One remarkable point is that as the number of data increases, LSTM and GRU algorithms give very close or better results with KNN and DT.

# 4 Summary

## 4.1 Conclusion

Studies show that HVAC systems still have potential energy savings up to 40% [44]. Considering the spare parts that are changed before completing their life, it is obvious that this ratio occupies a very large place in the world economy. Data-driven methods are easy to implement, are low in cost, and have the potential to provide a great reduction in maintenance costs for HVAC systems.

In this thesis, it is aimed to estimate the return filter clogging time of the HVAC system of a shopping mall by using a data-driven approach. The biggest indicator of filter clogging is the difference in pressures read from the filter inlets and outlets, in other words, pressure drop. An increase in this difference means that the filter cannot fulfill its task and the amount of particles it holds increases. This will cause the filter to clog.

During the study, other system variables related to return filter pressure drop value were examined. Then, models were created with pre-measured data using various ML, DL and hybrid algorithms with associated input features.

In order to get better results from the models created, the existing dataset was divided into 4, and the same models were built for 4 datasets and their performances were observed. Model performance comparison was made in Section 3.8 and the best results were obtained with the KNN method.

## 4.2 Future Work

In future work, extended studies can be done with longer-term collected data and more input features. If data such as the number of people entering and leaving the building where the system is located, indoor and outdoor temperatures, micron values of the filters used in the system, system failures etc. are provided, the prediction models to be created will be more descriptive and more generalizable for the existing system. Thus, more accurate and effective predictions performance can be obtained.

# 5 References

[1]  J. Grözinger, T. Boermans, A. John, F. Wehringer and J. Seehusen, "Overview of member states information on NZEBs," *Background paper Final Report. Cologne, Germany: ECOFYS GmbH,* 2014.

[2]  L. Pérez-Lombard, J. Ortiz and C. Pout, "A review on buildings energy consumption information," *Energy and Buildings,* vol. 40, p. 394–398, 2008.

[3]  P. F. Orrù, A. Zoccheddu, L. Sassu, C. Mattia, R. Cozza and S. Arena, "Machine Learning Approach Using MLP and SVM Algorithms for the Fault Prediction of a Centrifugal Pump in the Oil and Gas Industry," *Sustainability,* vol. 12, p. 4776, 2020.

[4]  S. Vilarinho, I. Lopes and J. A.Oliveira, "Preventive maintenance decisions through maintenance optimization models: a case study," 2017.

[5]  R. Yan, Z. Ma, Y. Zhao and G. Kokogiannakis, "A decision tree based data-driven diagnostic strategy for air handling units," *Energy and Buildings,* vol. 133, p. 37–45, 2016.

[6]  E. Gencalp, "Atm demand prediction with machine learning methods (Master's Thesis)," Galatasaray University, 2020.

[7]  C. J. Willmott and K. Matsuura, "Advantages of the mean absolute error (MAE) over the root mean square error (RMSE) in assessing average model performance," *Climate Research,* vol. 30, p. 79–82, 2005.

[8]  "scikit-learn.org," 01 Nov 2021. [Online]. Available: https://scikit-learn.org/stable/modules/model_evaluation.html#r2-score.

[9]  "wikipedia.org," 1 Nov 2021. [Online]. Available: https://en.wikipedia.org/wiki/Mean_squared_error. [Accessed 22 11 2021].

[10] Z. Du, B. Fan, X. Jin and J. Chi, "Fault detection and diagnosis for buildings and HVAC systems using combined neural networks and subtractive clustering analysis," *Building and Environment,* vol. 73, p. 1–11, 2014.

[11] S. Katipamula and M. Brambley, "Review Article: Methods for Fault Detection, Diagnostics, and Prognostics for Building Systems—A Review, Part II," *HVAC&R Research,* vol. 11, p. 169–187, 2005.

[12] J. M. Vaezi-Nejad, &. W. H. and J. M. House, "An expert rule set for fault detection in air-handling units/discussion.," *Ashrae Transactions,107, 858.,* 2001.

[13] Y. Guo, J. Wang, H. Chen, G. Li, R. Huang, Y. Yuan, T. Ahmad and S. Sun, "An expert rule-based fault diagnosis strategy for variable refrigerant flow air conditioning systems," *Applied Thermal Engineering,* vol. 149, p. 1223–1235, 2019.

[14] M. W. Ahmad, M. Mourshed, B. Yuce and Y. Rezgui, "Computational intelligence techniques for HVAC systems: A review," *Building Simulation,* vol. 9, p. 359–398, 2016.

[15] R. Y. a. G. Rizzoni, "Comparison of model-based vs. data-driven methods for fault detection and isolation in engine idle speed control system," *Annual Conference of the Prognostics and Health Management Society,* 2016.

[16] D. P. Zhou, Q. Hu and C. J. Tomlin, "Quantitative comparison of data-driven and physics-based models for commercial building HVAC systems," in *2017 American Control Conference (ACC)*, 2017.

[17] H. Sha, P. Xu, C. Hu, Z. Li, Y. Chen and Z. Chen, "A simplified HVAC energy prediction method based on degree-day," *Sustainable Cities and Society,* vol. 51, p. 101698, 2019.

[18] H. Alimohammadi, "Parameter estimation in nonlinear dynamic systems using the optimization methods based on the experimental data," 2020.

[19] A. Diez-Olivan, J. D. Ser, D. Galar and B. Sierra, "Data fusion and machine learning for industrial prognosis: Trends and perspectives towards Industry 4.0," *Information Fusion,* vol. 50, p. 92–111, 2019.

[20] N. A. Sulaiman, M. P. Abdullah, H. Abdullah, M. N. Shah and A. M. Yusop, "Fault detection for air conditioning system using machine," *IAES International Journal of Artificial Intelligence,* vol. 9, no. 1, pp. 109-116, 2020.

[21] T. Schreiber, C. Netsch, S. Eschweiler, T. Wang, T. Storek, M. Baranski and D. Müller, "Application of data-driven methods for energy system modelling demonstrated on an adaptive cooling supply system," vol. 230, p. 13, 2021.

[22] H. Shahnazari, P. Mhaskar, J. M. House and T. I. Salsbury, "Modeling and fault diagnosis design for HVAC systems using recurrent," *Computers and Chemical Engineering,* pp. 189-203, 2019.

[23] A. Ebrahimifakhar, A. Kabirikopaei and D. Yuill, "Data-driven fault detection and diagnosis for packaged rooftop unitsusing statistical machine learning classification methods," *Energy & Buildings,* pp. 1-12, 2020.

[24] M. S. Mirnaghi and F. Haghighat, "Fault detection and diagnosis of large-scale HVAC systems in buildings using data-driven methods: A comprehensive review," *Energy and Buildings,* vol. 229, p. 110492, 2020.

[25] E. Alpaydin, Introduction to machine learning, 2010.

[26] M. Paluszek and S. Thomas, "An Overview of Machine Learning," in *MATLAB Machine Learning*, Apress, 2016, p. 3–15.

[27] T. Wuest, D. Weimer, C. I. Thoben and Klaus-Dieter, "Machine learning in manufacturing: advantages, challenges, and applications," 2016.

[28] Y. Zhang, New Advances ın Machıne Learnıng, Rıjeka,Crotia: Intech, 2010.

[29] J. Vanderplas, "scikit-learn.org," 2007-2021. [Online]. Available: https://scikit-learn.org/stable/modules/neighbors.html. [Accessed 08 November 2021].

[30] V. N. Vapnik, The Nature of Statistical Learning Theory, Springer , 1995.

[31] E. Arslan, "Developing Internet of Things and Machine Learning Based Bi-Directional People Counting System With Passive Infrared Sensors (Master's Thesis)," Tallinn University of Technology, 2021.

[32] B. Efron, T. Hastie, I. Johnstone and R. Tibshirani, "Least angle regression," *The Annals of Statistics,* vol. 32(2), pp. 407-499, 2004.

[33] "https://scikit-learn.org/," 2021. [Online]. Available: https://scikit-learn.org/stable/modules/linear_model.html#least-angle-regression. [Accessed 07 12 2021].

[34] S. Varsamopoulos, K. Bertels and C. Almudever, "Designing neural network based decoders for surface codes," 2018.

[35] A. Perambai, "https://medium.com," 09 08 2019. [Online]. Available: https://medium.com/analytics-vidhya/lstm-and-gru-a-step-further-into-the-world-of-gated-rnns-99d07dac6b91. [Accessed 20 12 2021].

[36] J. Chung, C. Gulcehre, K. Cho and Y. Bengio, "Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling," *arXiv:1412.3555,* 2014.

[37] Z. Guo, H. Moayedi, L. K. Foong and M. Bahiraei, "Optimal modification of heating, ventilation, and air conditioning system performances in residential buildings using the integration of metaheuristic optimization and neural computing," *Energy and Buildings,* vol. 214, 2020.

[38] S. Raschka, "github.io," 2020. [Online]. Available: http://rasbt.github.io/mlxtend/user_guide/regressor/StackingRegressor/. [Accessed 2021].

[39] S. Callé, P. Contal, D. Thomas, D. Bémer and D. Leclerc, "Description of the clogging and cleaning cycles of filter media," *Powder Technology,* vol. Volume 123, no. Issue 1, pp. Pages 40-52, 2002.

[40] H. Alimohammadi, K. Vassiljeva, E. Petlenkov, M. Thalfeldt, A. Mikola, T. M. Kull and A. Köse, "Gray Box Time Variant Clogging Behaviour and Pressure Drop Prediction of the Air Filter in the HVAC System," 2021.

[41] S. Raschka, "rasbt.github.io," 2020. [Online]. Available: http://rasbt.github.io/mlxtend/user_guide/feature_selection/SequentialFeatureSele ctor/. [Accessed 5 November 2021].

[42] J. Bergstra and Y. Bengio, "Random Search for Hyper-Parameter Optimization," *Journal of Machine Learning,* pp. 281-305, 2012.

[43] S. Firmin, "https://community.alteryx.com/," [Online]. Available: https://community.alteryx.com/t5/Data-Science/Hyperparameter-Tuning-Black-Magic/ba-p/449289. [Accessed 15 November 2021].

[44] E. S. Cardoso, "Advanced energy management strategies for HVAC systems in smart buildings," 2019.

# Appendix 1 - Non-exclusive licence for reproduction and publication of a graduation thesis[1]

I Ahmet Caglayan

1. Grant Tallinn University of Technology free licence (non-exclusive licence) for my thesis "DATA-DRIVEN METHODS FOR ANALYSIS AND FAULT DETECTION OF HVAC SYSTEMS: FILTER CLOGGING PREDICTION EXAMPLE", supervised by Eduard Petlenkov

    1.1. to be reproduced for the purposes of preservation and electronic publication of the graduation thesis, incl. to be entered in the digital collection of the library of Tallinn University of Technology until expiry of the term of copyright;

    1.2. to be published via the web of Tallinn University of Technology, incl. to be entered in the digital collection of the library of Tallinn University of Technology until expiry of the term of copyright.

2. I am aware that the author also retains the rights specified in clause 1 of the non-exclusive licence.

3. I confirm that granting the non-exclusive licence does not infringe other persons' intellectual property rights, the rights arising from the Personal Data Protection Act or rights arising from other legislation.

03.01.2022

---

# Appendix 2 - Code Repository

Codes used in this study can be found at the below link:

https://github.com/caglayanahmet/Data-Driven-Methods-for-Analysis-and-Fault-Detection-of-HVAC-Systems-Filter-Clogging-Prediction

Appendix – 4, 5, 6, 7, 8, 9, and 10 was added to the code repository.

# Appendix 3 – Graphs

Data between 09.07.2019 - 21.10.19 – (the other graphs can be obtained with the code at appendix 4 at the github repository of this study)

GRU Model