

VŠB – Technical University of Ostrava  
Faculty of Electrical Engineering and Computer Science  
Department of Computer Science



---

# **Using Soft Computing Methods to Identify Operational Technical statuses in Smart Home within IoT**

---

**This page will be replaced by official thesis assignment in printed  
and electronic version of the submitted thesis.**

I hereby declare that this bachelor's thesis was written by myself. I have quoted all the references I have drawn upon.

Ostrava, April 2020

.....

I hereby agree to the publishing of the bachelor's thesis as per s.26, ss. 9 of the Study and Examination Regulations for Master's Degree Programmes at VŠB – Technical University of Ostrava.

Ostrava, April, 2019

.....

## **Abstract**

This bachelor project aims to focus on the occupancy monitoring of Smart Home Care (SHC) to optimize functions that SHC provides and optimize energy waste. The primary cause of this project, being focused on this topic, is that, so far, most of the software developed for SHC is based in room occupancy. To know the presence of persons without being invasive, we must use indirect methods to predict it. The project is based on the prediction of the CO<sub>2</sub> waveform based on the use of sensors that measure the indoor temperature and the relative humidity. Support Vector Machine (SVM) is employed to predict the presence of a person with high accuracy.

**Keywords:** IoT, Smart Home, prediction of room occupancy, presence of person monitoring

# Contents

<b>List of symbols and abbreviations</b>	<b>8</b>
<b>List of Figures</b>	<b>9</b>
<b>List of Tables</b>	<b>10</b>
<b>Code Listings</b>	<b>11</b>
<b>1 Introduction</b>	<b>12</b>
1.1 Definitions . . . . .	12
1.2 Project description . . . . .	13
<b>2 Project development</b>	<b>15</b>
2.1 Analysis of requirements . . . . .	15
2.2 Data preprocessing and algorithm design . . . . .	15
2.3 Coding . . . . .	15
2.4 Execution and verification . . . . .	15
2.5 Depuration . . . . .	15
2.6 Benchmarking and experimentation . . . . .	15
2.7 Documentation . . . . .	15
2.8 Outline . . . . .	16
<b>3 IoT</b>	<b>17</b>
<b>4 Support Vector Machine (SVM)</b>	<b>18</b>
4.1 What is SVM? . . . . .	18
4.2 Vector length . . . . .	19
4.3 Direction of a vector . . . . .	19
4.4 Dot product . . . . .	19
4.5 Hyperplane . . . . .	20
4.6 How SVM classifies and regresses? . . . . .	20
4.7 Optimal Hyperplane . . . . .	21
4.8 Hard and Soft Margin . . . . .	21
4.9 Non linearly separable data . . . . .	22
<b>5 Implementation</b>	<b>24</b>
5.1 Hardware limitations . . . . .	24
5.2 Python . . . . .	24
5.3 Enviroment . . . . .	26

5.4	Data preprocessing . . . . .	26
5.5	General ensemble vs. One model per room . . . . .	30
5.6	Code . . . . .	32
<b>6</b>	<b>Discussion</b>	<b>38</b>
6.1	Applying shuffle . . . . .	38
6.2	Applying Savitzky–Golay filter . . . . .	38
6.3	Implementations . . . . .	40
<b>7</b>	<b>Conclusions</b>	<b>47</b>
	<b>References</b>	<b>48</b>
	<b>Appendix</b>	<b>50</b>
<b>A</b>	<b>Code implementation</b>	<b>51</b>
A.1	Ensemble voting . . . . .	51
A.2	Ensemble bagging . . . . .	55
A.3	One model per room . . . . .	57
A.4	Wrapper . . . . .	59

## List of symbols and abbreviations

SH	– Smart Homes
SHC	– Smart Home Care
IB	– Intelligent Buildings
IoT	– Internet of Things
HVAC	– Heating, ventilation, air conditioning
SVM	– Support Vector Machine
SC	– Soft Computing
HC	– Hard Computing
IM	– Indirect Methods
DM	– Direct Methods



## List of Figures

1	Flow of the different project Phases . . . . .	16
2	Linearly separable data . . . . .	18
3	Graphical representation for dot product [26] . . . . .	19
4	Data linearly separated by a hyperplane $w \cdot x + b + \epsilon = 0$ [27] . . . . .	20
5	Living room CO2 signal(4 hours) compared with living room CO2 signal with applied Savitzky–Golay filter, with a window length of 21 and polynomial order of 1 . . . . .	28
6	Bagging approach to regression. Weak learners are trained with subsamples of the data. Then the final regression prediction is made with some averaging method. [38] . . . . .	30
7	Flowchart of the steps followed in code . . . . .	32
8	Outlier visualization from Living Room file: CO2 (ppm) vs H(%) . . . . .	33
9	Scattered data after the removal of the outliers from Living Room file: CO2 (ppm) vs H(%)	34
10	Scattered data after the filter Savitzky–Golay is applied to the data from Living Room file: CO2 (ppm) vs H(%) . . . . .	35
11	Predictions made by the best voting ensemble model. Train and test size: 1440 instances (one day). R2 score: 0.76, MAE: 0.09, MSE: 0.01, RMSE: 0.13 . . . . .	41
12	Predictions made by the best bagging ensemble model. Train and test size: 1440 instances (one day). R2 score: 0.43, MAE: 0.12, MSE: 0.02, RMSE: 0.16 . . . . .	42
13	Predictions made by the best voting ensemble model working without bathroom and kitchen rooms. Train and test size: 1440 instances (one day). R2 score: 0.83, MAE: 0.04, MSE: 0.01, RMSE: 0.10 . . . . .	44
14	Predictions made by the best bagging ensemble model working without bathroom and kitchen rooms. Train and test size: 1440 instances (one day). R2 score: 0.47, MAE: 0.13, MSE: 0.03, RMSE: 0.17 . . . . .	45
15	Predictions made by creating one model per room. A total of five models where created. Train and test size: 1440 instances (one day). R2 score: 0.93, MAE: 0.01, MSE: 0.003, RMSE: 0.0599 . . . . .	46

## List of Tables

1	SVR results with one model per room and non shuffled data. Used two test and train sets per file (1440 instances that equal one day, and 7200 instances that equals five days) . . .	38
2	SVR results trained and tested with one day, and five days of the data (1440 instances, and 7200 instances). The data was not filtered with the Savitzky–Golay filter. . . . .	39
3	SVR results with one model per room, shuffled data and filtered with Savitzky–Golay filter. Used two test and train sets per file (1440 instances that equal one day, and 7200 instances that equals five days) . . . . .	39
4	Voting ensemble results. The results from this ensemble have an average R2 score of 0.60	40
5	Bagging ensemble results. The results from this ensemble have an average R2 score of 0.41 . . . . .	41
6	Voting ensemble results working without bathroom and kitchen rooms. The results from this ensemble have an average R2 score of 0.78 . . . . .	43
7	Bagging ensemble results working without bathroom and kitchen rooms. The results from this ensemble have an average R2 score of 0.46 . . . . .	43

**Code Listings**

1	Python script for the creation of an ensemble using voting ensemble approach explained in section 5.5 . . . . .	51
2	Python script for the creation of an ensemble using bagging ensemble approach explained in section 5.5 . . . . .	55
3	Python script for the creation of an ensemble using one model per room approach explained in section 5.5 . . . . .	57
4	Python script with the common functions for all the scripts coded . . . . .	59

# 1 Introduction

In this section, the project and the problem to solve will be introduced. Also, the motivation and the actual situation of Smart Homes will be explained.

## 1.1 Definitions

In the following we define and explain the terms used in this project.

### 1.1.1 Smart Homes

Smart Homes (SH) are buildings where technical functions like heating, ventilation, air conditioning (HVAC), and lighting are controlled and automatized [1]. The idea of SH is to allow and let the automatic systems make decisions based on the environment [2], behavior, and occupancy of the house.

### 1.1.2 Intelligent Buildings

The concept of Intelligent Buildings (IB) is hard to define because there is not a single definition of it. In [3], Shengwei Wang described and categorized the different types of IB. He formalized three categories.

- **Performance-based definitions:** “Defines an IB as a building created to give its users the most efficient environment, at the same time, the building utilizes and manages resources efficiently and minimizes the life costs of hardware and facilities.” [3]
- **Services- based definitions:** “Defines an IB is a building with the service functions of communication, office automation, and building automation and is convenient for intelligent activities. Services to users are emphasized.” [3]
- **System-based definitions:** “States that IB provide building automation, office automation and communication network systems, and an optimal composition integrates the structure, system, service and management, providing the building with high efficiency, comfort, convenience and safety to users.” [3]

However, generically, in this project, we will describe IB as buildings that use infrastructure to create more efficient and productive buildings. By the automatization and remote control of multiple tasks as HVAC, lighting, and others, it is allowed to increase the performance of the occupants.

### 1.1.3 Soft Computing and Hard Computing

Soft computing (SC) is a set of methods that allow solving complex real-life problems reducing costs with high consistency and robustness [4]. Usually, it deals with problems that cannot be solved by mathematical and exact methods (Hard Computing) because these methods will take an extremely long time

to obtain results.

SC is an alternative to Hard Computing (HC) to solve problems where HC will consume much time to provide suitable solutions [5](17). In SC, the solving of problems is imprecise, tolerant, and uncertain. However, most of the solutions provided can be usable. There are different methods to get solutions to problems. Nevertheless, it is crucial to know the “No free lunch” theorem [6], which means that there is no universal method to solve every problem, meaning that each method will fit better for some determined kind of problems. Besides, the selection of the method to get solutions will depend on the problem to solve. Using SC is a good option as there are large quantities of data to process and compared to HC the time to train the models will be shorter.

#### **1.1.4 Indirect and direct methods**

We define Indirect Methods (IM) as the methods which can predict the occupancy without being invasive. These methods are user friendly, and they do not involve cameras, microphones, or other devices that invade the privacy of the user. The difference between IM and Direct Methods (DM) is straightforward. DM can use all of these invasive devices being easier to predict the occupancy, but for private environments, IM are way better.

### **1.2 Project description**

The SH market made 51,75 billion Euros in 2018, and it is expected to grow twice as big for 2023. Europe represents 14% of this market and is expected to be 30,9% in 2023. Knowing these statistical data, we can assume that the SH is the future for buildings, making them remarkably important with a new level of technology development [7] [8]. That is to say, the more optimal and improved energy efficiency [9] [10], the more benefits. These benefits are both economic and ecologic, making SH more acceptable with what the customer gets and its reliability.

One of the motivations is that in SH, some ways to detect the occupancy are based on direct methods. These methods do not respect the privacy and intimacy, making people feel uncomfortable [11] [12]. Furthermore, these methods raise way too much the costs and usually are difficult to adapt to future changes [10] [13] [14]. The use of mobile phones for detection is a great choice but assumes that every person carries their mobile phone everywhere [15]. Nevertheless, working with CO<sub>2</sub>, relative humidity, and indoor temperature sensors allow making SH systems more affordable and less invasive. Another motivation for this project is the care of the elderly or with disabilities [16]. The detection and recognition of the elderly or with disabilities behavior allows them to live independently, as these IoT systems can be remotely controlled and they can detect critical situations before they happen [14] [17].

The detection of occupants is a hard task, and it can be performed in many ways, and this is because human presence changes its environment, and this can be measured with different sensors [2] [11]. Thus,

as before was said, these systems are expensive and do not have good scalability. This task can be focused on the detection of occupants based on the prediction of CO<sub>2</sub> concentration [18] [19] [20], as it is an excellent way to detect the presence, cheap and not invasive.

The use of statistical models to obtain results in the problem of prediction of human presence will have more accurate results. Methods such as Linear Discriminant Analysis (LDA), Classification and Regression Trees (CART), and Random Forest (RF) have accuracies from 95 to 99% [21]. In the same way, clustering methods (K-means, Self-Organizing Maps), to find patterns in the data to predict the occupancy, will obtain positive results [22]. Even further, the use of models like Multiple Non-linear Regression (MNLN) and Deep neural network (DNN) throws small percentages of error (Between 2,37% and 10,34%) [23].

Our work aims to determine the efficiency of the SVM method, believing that it is innovative and not commonly used in this field [24]. SVM works for classification and regression problems based on the kernel trick, which transforms the data allowing it to find optimal solutions.

## **2 Project development**

In this section, we strip the project in different phases, and we summarize all the chapters in order to get a better comprehension of it.

### **2.1 Analysis of requirements**

In this phase, the problem will be analyzed and investigated so that the algorithm (SVM) can be developed accordingly to obtain better results. The objective of this phase is to understand and work with the data gathering systems from VSB University for a better understanding of the available data: where they come from, what relationships exist, among others.

### **2.2 Data preprocessing and algorithm design**

The algorithm design will be developed out generically, to solve the main problem, and refine the details later, thus being able to obtain optimal results. The data provided will be transformed to be optimal for processing, selecting essential characteristics and attributes, and dropping those that do not provide information.

### **2.3 Coding**

The design of the algorithm will be coded by adding all the details that are necessary for its optimal operation. To accomplish the coding of the algorithm, the use of the laptop, Dell XPS 13 9360 will be done.

### **2.4 Execution and verification**

Once we finish the codification of the algorithm, it is essential to carry out the tests with the available data and see that works and operates optimally.

### **2.5 Depuration**

In the case of finding errors in the execution or miss functionality, we will fix them, and we will repeat the previous step. It would be a cyclical process until the execution is correct and optimal.

### **2.6 Benchmarking and experimentation**

With the developed algorithm, we will work with the available data, proposing different environments, and seeing different results for later comparison.

### **2.7 Documentation**

This step is part of all the other phases. However, in the end, we will work with the results obtained to develop the final discussion and conclusions of the project.

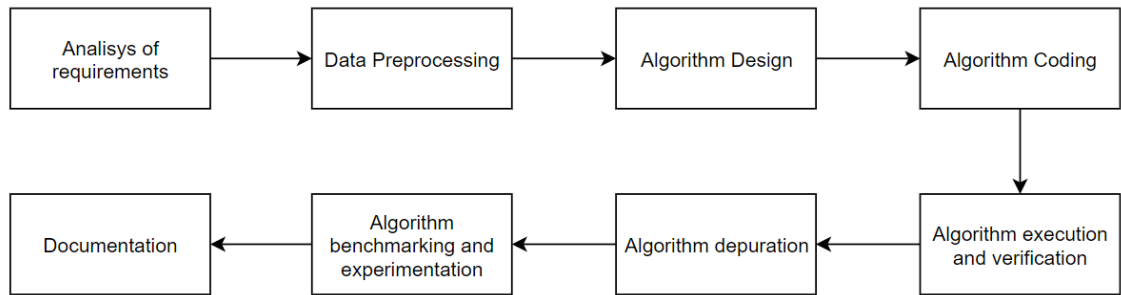


Figure 1: Flow of the different project Phases

## 2.8 Outline

This thesis is organized as follows:

- In **Section 3** iot implementation
- In **Section 4**, theoretical foundations of SVM (Kernels, hyperplanes, margins) are introduced.
- In **Section 5.4**, we explain the steps for data preprocessing.
- In **Section 5**, we give a detailed explanation of the implementation and the faced challenges. Furthermore, Python and libraries used are explained.
- In **Section ??**, we show the process of experimentation made with different configurations of the parameters.
- In **Section 6**, we resume the project and make some conclusions of it, although we explain future work that could be done.
- In **Section 7**, we resume the project and make some conclusions of it, although we explain future work that could be done.



### **3 IoT**

## 4 Support Vector Machine (SVM)

In this section, the theoretical foundations of SVM are explained as all of the parts that form it. Approaches to handle regression and classification with SVM are introduced.

### 4.1 What is SVM?

SVM forms part of Supervised Learning, a part of statistical learning which starts from a series of examples to create a predictor system that tries to predict new values. SVM [25] are learning algorithms that analyze data on classification or regression problems and outliers detection, being one of the most potent methods supervised learning models. SVM methods try to find the optimal line, plane, or hyperplane, which linearly separates the data by maximizing the margin.

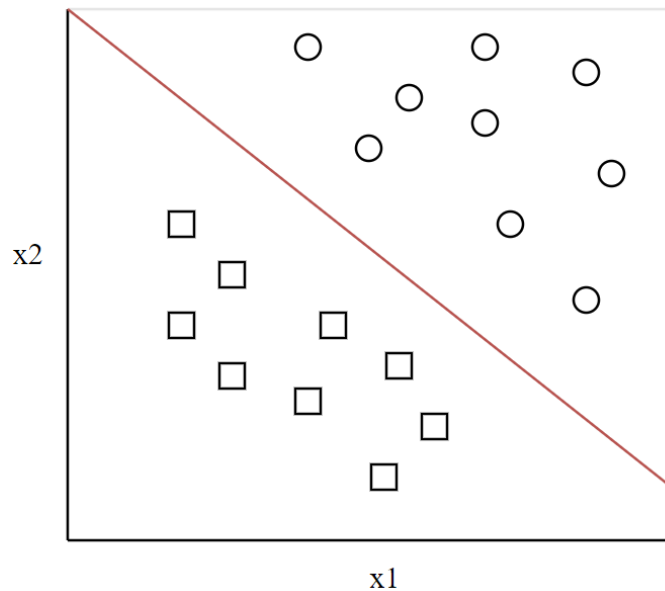


Figure 2: Linearly separable data

As it is possible to see in figure 2, there is a set of two-dimensional points that are linearly separable. SVM works excellent and fast with this kind of data. However, SVM also operates with non-linearly separable data.

## 4.2 Vector length

The length of a vector is also known as the norm of the vector. It is calculated by the euclidean formula, which tells how far are from the origin. Being  $v = (v_1, v_2, v_3, \dots, v_n)$  his norm is [26]:

$$\|v\| = \sqrt{v_1^2 + v_2^2 + v_3^2 + \dots + v_n^2}$$

## 4.3 Direction of a vector

Vector direction ( $w$ ) is calculated as the division of every element of the vector by the norm of itself. Being  $v = (v_1, v_2, v_3, \dots, v_n)$  his direction is [26]:

$$w = \left( \frac{v_1}{\|v\|}, \frac{v_2}{\|v\|}, \frac{v_3}{\|v\|}, \dots, \frac{v_n}{\|v\|} \right)$$

## 4.4 Dot product

The dot product between vectors is a scalar quantity. It tells how two vectors are related.

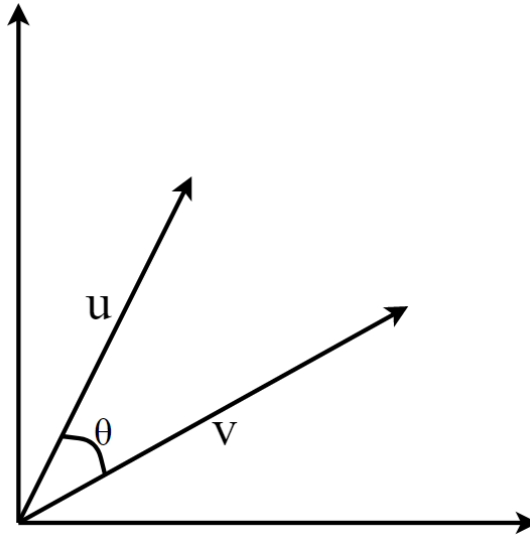


Figure 3: Graphical representation for dot product [26]

Mathematically is defined as [26]:

$$u \cdot v = \|u\| \|v\| \cos(\theta)$$

For a two n-dimensional vector, the dot product is computed [26] as:

$$u \cdot v = \sum_{i=1}^n u_i v_i$$

## 4.5 Hyperplane

A hyperplane is a codimension-1 vector subspace from a vector space [26]. In a more informal way to define a hyperplane, we say that it's a line that divides the data into two different sets. In figure 2, we can see how a hyperplane splits the dataset between squares and circles. The equation for the hyperplane [26] is:

$$w \cdot x + b + \varepsilon = 0$$

Where  $\mathbf{w}$  is a vector normal to the hyperplane ( $\mathbf{w}$  is orthogonal to the hyperplane), and  $\mathbf{b}$  is an offset. For any vector  $\mathbf{x}$ , we can compute the next equation  $w \cdot x + b + \varepsilon = y$ . If  $y = 0$ , then  $x$  is on the hyperplane.

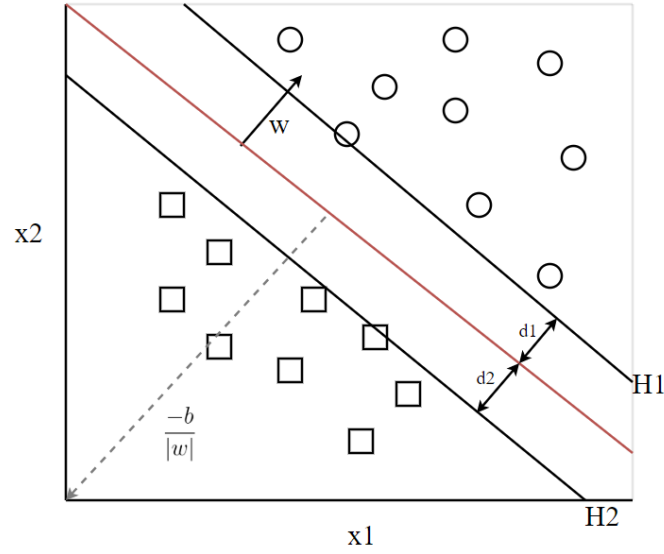


Figure 4: Data linearly separated by a hyperplane  $w \cdot x + b + \varepsilon = 0$  [27]

H1 and H2 are the Support Vectors and they are mathematically defined as:

$$H1 : w \cdot x + b = +1$$

$$H2 : w \cdot x + b = -1$$

The sum of d1 and d2 is the margin, in there will never be data, this distance is known as epsilon. SVM attempts to find the maximum margin. For SVM, it is vital to choose the optimal hyperplane.

## 4.6 How SVM classifies and regresses?

To know how to classify with SVM, all  $\mathbf{x}$  of the data above or on the hyperplane are classified as +1. If they are under the hyperplane, they are classified as -1. The formal definition of the constraints are [25]:

$$prediction = \begin{cases} y_i - wx_i - b \leq \varepsilon \\ wx_i + b - y_i \leq \varepsilon \end{cases}$$

## 4.7 Optimal Hyperplane

The optimal hyperplane [28] is the one with the most significant margin. In section 4.5, we explained H1 and H2, and they are the hyperplanes that go through the Support Vectors. The distance between H2 and the origin is calculated as  $\frac{(-1-b)}{\|w\|}$  and the distance from H1 to the origin is  $\frac{(1-b)}{\|w\|}$ . The margin [25] ( $M$ ) between H1 and H2 is given by:

$$M = \frac{(1-b)}{\|w\|} - \frac{(-1-b)}{\|w\|} = \frac{2}{\|w\|}$$

The margin then is defined by  $\frac{1}{\|w\|}$  because  $\frac{2}{\|w\|}$  is twice the margin. SVM attempts to find the optimal hyperplane, which is the one that maximizes the margin, so the bigger  $\frac{1}{\|w\|}$  it is, the better. This problem can be changed to a minimization problem by trying to find the minimum inverse of the margin [25].

$$\max \frac{1}{\|w\|} = \min \frac{1}{\frac{1}{\|w\|}} = \min \|w\|$$

And this minimization problem is equivalent [25] to:

$$\min \|w\| \equiv \min \frac{\|w\|^2}{2}$$

*such that  $y_i(w \cdot x + b) - 1 \geq 0$ , for  $i = 1 \dots l$*

This equivalent problem is known as a convex quadratic optimization problem. This is the main optimization problem of SVM.

## 4.8 Hard and Soft Margin

SVM, working with linearly separable data, can be formulated in various ways, with Hard Margin or with Soft Margin. The difference between them is that Hard Margin SVM does not work with non-linearly separable data and does not accept outliers. Outliers (noise in data) is an obstacle for real-life problems since these, by definition, are not usually linearly separable, and we can find noise in them. For Hard Margin SVM, it is necessary to satisfy all the constraints in  $y_i(w \cdot x_i + b) \geq 1$ , and the outliers make this impossible, making the optimization problem unsolvable.

### 4.8.1 Soft Margin

Soft Margin is a more flexible model. It consists of adding a slack variable  $\zeta_i$  to the constraints of the problem. This slack [25] variable is where the sample is located in relation to the hyperplane and the margin.

$$(w \cdot x_i + b) \geq 1 - \zeta_i, \text{ for } y_i = +1$$

$$(w \cdot x_i + b) \leq -1 + \zeta_i, \text{ for } y_i = -1$$

If we combine both of the equation above, we get the next equation [25]:

$$y_i(w \cdot x_i + b) - 1 + \zeta_i \geq 0, \text{ for } y_i = +1, -1$$

Where  $\zeta_i \geq 0$  for  $i = 1 \dots l$ .

For a better regularization,  $C$  parameter is added. This parameter controls the importance of the slack variable. For  $C = 0$ , the algorithm does not allow any misclassification. For  $C > 0$  means that  $C$  samples can be misclassified. The more significant  $C$  is, the higher the margin will be. If  $C = \infty$ , then the model will misclassify most of the samples. In conclusion, Soft Margin allows margin violations. This margin violation means choosing a hyperplane, which will allow some data to stay between the margin area or to be misclassified.

## 4.9 Non linearly separable data

In this project, the problem is a real-life problem, which means more complex data. This data is not linearly separable, so we need to work with a version of SVM, which generalizes better and works with different shapes of the data. If the data is not linearly separable in the input space, then transformations to the data are applied. These transformations map the data from the input space into a higher dimensional feature space. The purpose is that after these transformations to a higher-dimensional feature space, the data is now linearly separable.

As is explained before in section 4.7, having linearly separable data, the aim is to obtain the optimal hyperplane which separates the data between the different classes. This hyperplane will be a hyperplane in the higher-dimensional space. These transformations of the data is named as “Kernel trick”.

### 4.9.1 Kernel Trick

The definition of the Kernel Trick [29] is a function that takes as input vectors that are in the original space and returns the dot product of these vectors in a higher-dimensional feature space. This function is named “kernel function” and has to satisfy Mercer’s constraints [30].

The formal for the Kernel Trick [31] definition is:

$$K(x_i, x_j) = (\phi(x_i), \phi(x_j)) \quad (1)$$

It is possible to find multiple kernel functions [32] that transform the input data into a higher-dimensional feature space. Some of the kernels are:

- Polynomial [31]

$$K(x_i, x_j) = (x_i \cdot x_j + a)^p \quad (2)$$

- RBF [31]

$$K(x_i, x_j) = e^{-\gamma(x_i - x_j)^2} \quad (3)$$

- Sigmoid [31]

$$K(x_i, x_j) = \tanh(\eta x_i \cdot x_j + v) \quad (4)$$

#### 4.9.2 Parameters

Working with SVM for regression (SVR), there are two parameters to work with, which are C and Epsilon. However, it depends on the selected kernel. In this project, the chosen kernel is RBF because it's a well-studied method, and it works well in practice. This kernel is explained in the previous section, which adds another parameter to estimate. This new parameter is Gamma.

**4.9.2.1 C** This parameter is the error penalty. It punishes the errors made when the model misclassifies instances. If it is a small value means that the punishment done to the mistake made is small. If it's a high value, it will punish with more significant penalties to the errors. When C is large, it can lead to worse generalization models. Ergo, it's crucial to choose a good C value. [33]

**4.9.2.2 Epsilon** For SVR, this parameter will tell how smooth and how good will the model generalize. This is because it affects the number of support vectors selected. [34]

**4.9.2.3 Gamma** This parameter defines how far the influence of a single training example reaches. Low values mean that the area of influence is significant, and larger values indicate that the area of influence will be small. [33]

## 5 Implementation

This section covers the details of the implementation and obstacles faced in the process of transforming the theoretical concepts into operative python code.

### 5.1 Hardware limitations

For the development of this project, we used as main computer a DELL XPS 13 9365 for computing tasks, research, and model evaluations. This computer is an ultrabook made for long battery duration, being the CPU an i5-7200U. The problem we faced here is that for tasks of cross-validation, parameter estimation, and model fit, it would take long periods of time. This time could be reduced by using dedicated services like cloud computing, GPU processing, or a computer with better specs.

Still, in order to speed up the project, the tool Google Colaboratory has been used. This tool offers cloud computing an environment similar to Jupyter notebook. The specs that the free plan have are the next ones:

- Tesla K80 GPU
- 1 core 2 threads @ 2.2GHz
- 13GB RAM
- 33GB Free Space
- Max 12 hours of code execution

As specified before, Google Colaboratory offers a free GPU. Still, the library used for the model creation and training is scikit-learn, which does not intend to be used as a deep learning framework.

### 5.2 Python

We selected Python as the language for the development and implementation of this project for the following reasons:

- **Readable and easily maintainable code:** Python allows us to develop clear and concise code, without the need for complications. It is simple and easy to use.
- **Compatible with most systems and platforms:** Python is an interpreted programming language that allows us to execute the same code on multiple platforms without recompiling. Therefore, it is not necessary to recompile the code after making any changes.
- **Libraries:** Standard Python libraries allows us to choose from a wide range of modules according to our needs. Each module also allows us to add functionalities to the project without writing additional code. One of the benefits of Python libraries for data analysis is that it offers access to a wide variety of libraries related to data analysis and data science.



- **Many Frameworks and open source tools:** As an open-source programming language, Python helps us eliminate development costs from this project. It is possible to use several open-source Frameworks, libraries, and tools to reduce development time without increasing the costs.
- **Scalability:** It helps with the scalability of projects providing great flexibility and different ways of addressing various problems. We can find Python in multiple industries, allowing faster application development.

### 5.2.1 Libraries

The Python libraries selected for the realization of this project are the following.

**5.2.1.1 Pandas** Pandas is a Python library that allows us to work with different data formats, thus helping to have better data management, whether to perform analysis of them, preprocessing, or other tasks in which it is required to work with large amounts of data. It is an open-source library, and it is kept updated by its great use among the programmer community in Python.

**5.2.1.2 Numpy** This open-source Python library provides the project facilities to perform mathematical calculations. For example, a great facility to work with vectors, since it incorporates a vector object with which we can operate as if it were a scalar so that this tool will be necessary for the development of SVM.

**5.2.1.3 Matplotlib** We are going to use this open-source library for graphics. These graphics will help the visualization of the results, apart from being very useful when comparing and being more visual, it will facilitate the understanding of the results. This module is similar to the functions offered by MATLAB, being very configurable and easy to use.

**5.2.1.4 Scikit-learn** Scikit-learn is an open-source Python library that facilitates and provides us a large number of supervised and unsupervised learning algorithms. With the use of this library, we would already have incorporated all the previously explained libraries (Pandas, NumPy, Matplotlib) and many others (IPython, SciPy, Sympy) that will not necessarily be used in this project but may be helpful. Some of the models offered by this library are:

- **Clustering:** Methods such as KMeans, among others, to perform clustering tasks.
- **Datasets:** Generation of datasets with certain properties to test algorithms.
- **Feature Extraction:** Define attributes in different types of data such as text or image.
- **Feature Selection:** Allows the identification of attributes to create new models.
- **Cross-validation:** View the performance of the different algorithms developed with data that we have not previously worked.

### 5.3 Enviroment

The environment used for the coding phase was Jupyter Notebook combined with Visual Studio. Additionally, as explained in section 5.1, the environment Google Colaboratory was used to speed up the project. Moreover, python scripts were used to automatize the work of training and test of the models. We selected Jupyter Notebook as the development environment because it offers the following advantages:

- **Interactive code:** Jupyter Notebook offers some packages to have better code and data visualization.
- **Easy to share and export:** It uses Jason format to save each notebook allowing an easy way to share notebooks.
- **Text and code combination:** In each notebook, it is possible to combine text, code, and visualizations. This combination is beneficial because it allows us to get rid off the ordinal comments, having a better understanding of every part of the code.

### 5.4 Data preprocessing

Data preprocessing is the step where we transform and shape the data, so when we are working with, it's easier to understand, and the algorithm can interpret the data better. The steps for the realization of the data preprocessing are:

- Outlier detection
- Data filtering (Noise reduction)
- Data normalization

#### 5.4.1 Data

The data we are working with are five files with the measures from three different sensors: CO2 (ppm), Humidity(%), and temperature( $^{\circ}\text{C}$ ) from the day 21 December 2018 to the day 31 December 2018. The measures made were temperature indoor, relative humidity indoor, and CO2 concentration indoor. Those features were recorded in five different rooms of SHC for a period of one week utilizing five stations using Loxone temperature/ humidity/ CO2 sensors. The places included: "living room, children's room, bedroom, kitchen, and bathroom." The sensor specifications are:

- Indoor CO2 measurement (Within 0 and 2000 ppm, accuracy  $\pm 50$  ppm)
- Indoor temperature ( $T_i$ ) measurement (Within 0 and  $50^{\circ}\text{C}$ , accuracy at  $\pm 0.3^{\circ}\text{C}$ )
- Indoor relative humidity (rH) (within 0 and 100% RH, accuracy  $\pm 3\%$  RH)

### 5.4.2 Data normalization

Data normalization is not always necessary. If all the data is in the same range, we don't need to scale it. However, as we have different types and varieties of data (percentage, degree, and ppm), we need to normalize them. Normalize the data means to scale it between 0 and 1. And we have different methods to achieve this. The technique that we are going to use is known as min-max normalization. The formula is for data normalization is [35]:

$$X_{norm} = \frac{X - X_{min}}{X_{max} - X_{min}} \quad (5)$$

### 5.4.3 Outlier detection and removal

An outlier can be defined as a data point that is far away from all the other data points. In this problem, we understand outliers as mistakes done in the measurements made by the sensors. To delete these outliers, we decide to use Z-score, which tells how far is a data point from the mean. The formula to calculate the Z-score is [28]:

$$z = \frac{x - \mu}{\sigma} \quad (6)$$

Where  $\mu$  represents the mean, and  $\sigma$  the standard deviation. After we calculated the Z-score, we need to establish a threshold to delete all of those instances that their Z-score is higher or lower than the threshold are considered outliers. We set a threshold of 3, meaning that every point that is more than three times the standard deviation will be considered as an outlier.

#### 5.4.4 Data filtering (Noise reduction)

Savitzky–Golay filter is a noise reduction and smoothing filter that can be applied to the data to the data. By using this filter, it is possible to increase the accuracy of the predictions without distorting the signal. tendency. Applying this filter reduces the noise and deletes most of the outliers from the data signal, allowing important patterns to stand out and the more reliable prediction of the model. As it is shown in figure 5, all the noise from the data is erased. [36]

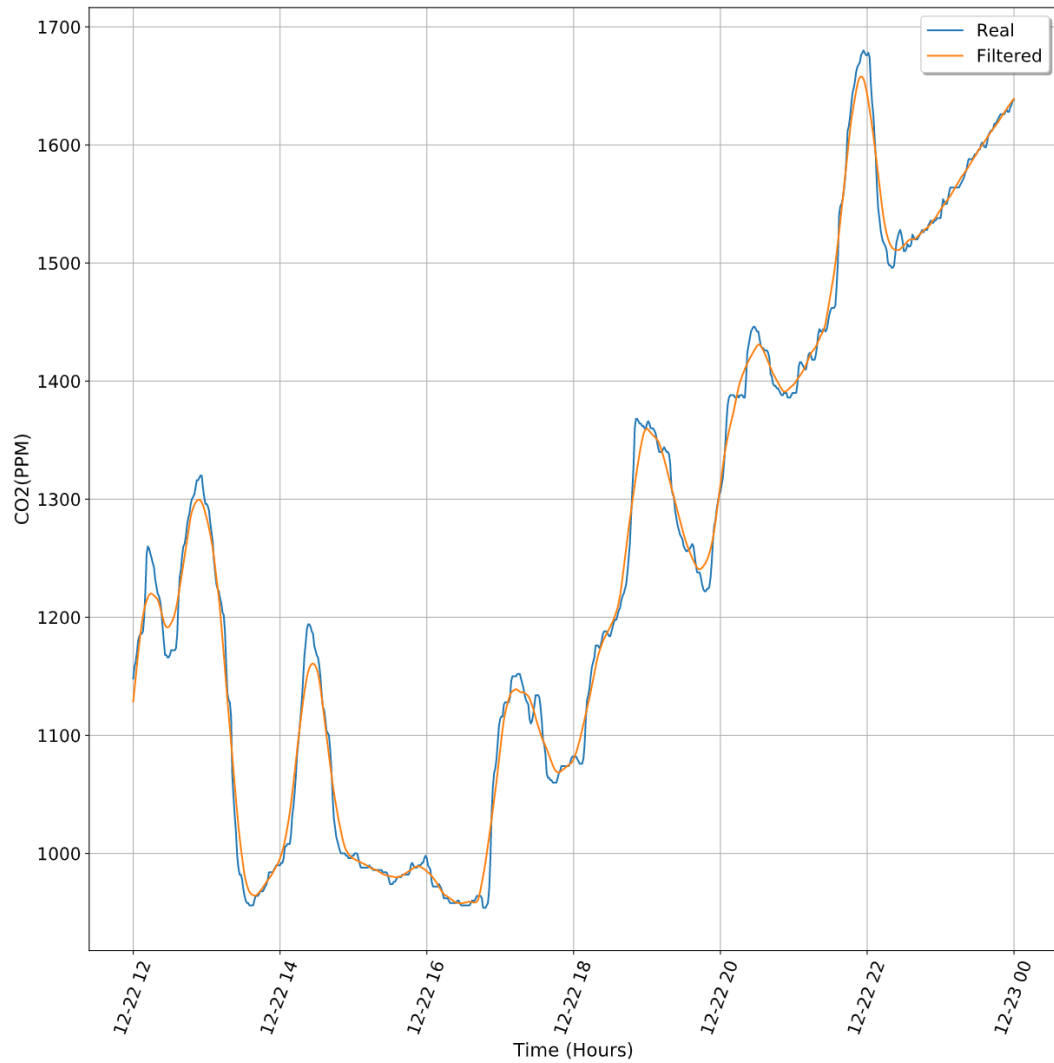


Figure 5: Living room CO2 signal(4 hours) compared with living room CO2 signal with applied Savitzky–Golay filter, with a window length of 21 and polynomial order of 1

### 5.4.5 Metrics for model evaluation

Because we want to be sure that our model is efficient and accurate, we selected three different metrics that are the next ones.

**5.4.5.1 Mean Squared Error** Mean Squared Error (MSE) is one of the most used and known metrics for regression problems. It measures the average of the squared difference between the predicted values and the actual values. Because it is the square of the difference, all the errors, even the small ones, have significant penalties. In conclusion, the lower the score obtained with this metric, the better. The formula to calculate MSE is [37]:

$$MSE = \frac{1}{n} \sum (y_i - \hat{y}_i)^2 \quad (7)$$

**5.4.5.2 Root Mean Squared Error** RMSE is a quadratic score that measures the average of the error. It's the square root of the MSE, which has been explained in the previous section 5.4.5.1. This metric is handy in terms of knowing if there are outliers that might be messing our model.

$$RMSE = \sqrt{MSE} = \sqrt{\frac{1}{n} \sum (y_i - \hat{y}_i)^2} \quad (8)$$

**5.4.5.3 Mean Absolute Error** Mean Absolute Error (MAE) is powerfully robust with outliers and missing data because it does not penalize as hard as MSE. For our problem, it's useful because we are not working with data that has a big group of outliers and missing data. The formula to calculate MAE is [37]:

$$MAE = \frac{1}{n} \sum |y_i - \hat{y}_i| \quad (9)$$

**5.4.5.4 R Squared** R squared (R2), also known as the coefficient of determination, evaluates the performance of our model by comparing our current model with a constant baseline and specifies how much better is our model. It always will be less than 1. The range of the score is from  $-\infty$  to 1, the closest to 1, the better it will be. The formula to calculate R2 is [37]:

$$R2 = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2} \quad (10)$$

## 5.5 General ensemble vs. One model per room

Ensembling is a technique that combines multiple weak models to get a robust and optimal predictor model. Ensemble methods help to minimize factors like noise, improving the stability and performance of the model. There are many different methods to create ensemble models. For our project, two types of ensembles were used: Bagging and Voting regressor. Both are offered in the library of Scikit-learn. Despite that, we also worked on generating a different model per room, having better accuracy than the ensembles tried.

### 5.5.1 Bootstrap Aggregation ensemble (Bagging)

As the name explains, Bagging is composed of two different parts: Bootstrapping and Aggregation, to form one ensemble model.

- **Bootstrapping:** This is a sampling method. Consists of the creation of different random subsamples using replacement. After the selection of these subsamples of the dataset, the learning algorithm that is SVM will proceed to predict.
- **Aggregation:** After the selection of the subsamples, and the multiple model creations, the aggregation method is used to aggregate all the models and form a more efficient predictor.

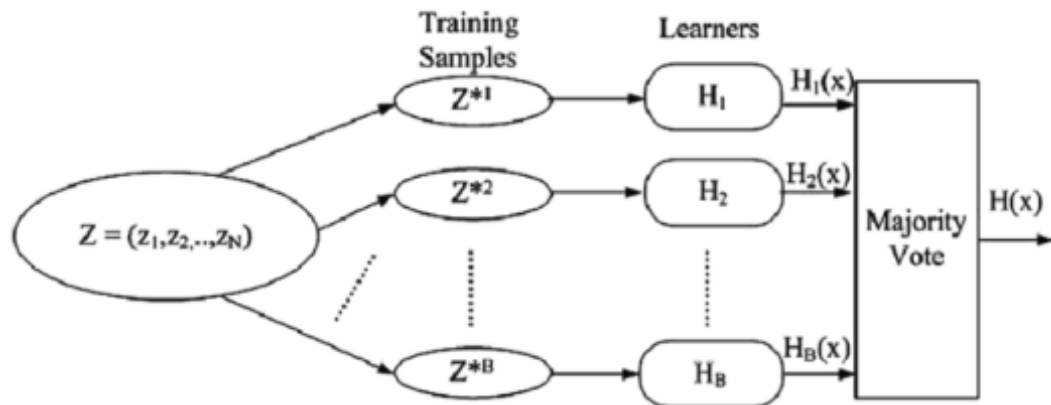


Figure 6: Bagging approach to regression. Weak learners are trained with subsamples of the data. Then the final regression prediction is made with some averaging method. [38]

Bagging can offer multiple advantages as a combination of multiple weak models to produce a strong predictor, helping to reduce the variance, noise, and overfitting. But it also has disadvantages. It has a really high computational cost, and it can be easily underfitted. This ensemble approach is applied to our project by preprocessing the data in a different way.

As explained in section 5.4.1, the data provided consists of 5 files, one for each room (Bedroom, Living room, Kitchen, Bathroom, Kids room). These files are merged into one file in order to train the ensemble bagging model. Each subsample selected can contain measures from different files.

### 5.5.2 Voting ensemble

The voting ensemble is an easy way to create an ensemble model. The idea of voting doesn't differ much from the bagging idea. Multiple heterogeneous weak models are created using subsamples from the dataset. Then the final model makes predictions based on the average of the predictions made by these weak models. The formula to calculate the average value is [28]:

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} \quad (11)$$

This approach differs from Bagging in our project because of the preprocessing of the data, in each ensemble model, the ways to prepare the data to train the weak learners are different. In this ensemble approach, the data is not merged into one file, but a different model is trained for each file. These individual models are merged using the average method to predict the regression.

### 5.5.3 One model per room

Although approaching this problem through the model ensemble is a good option, besides the creation of different models per room, it is also an excellent choice to get higher accuracy in the predictions. The data from the rooms differ much from one to another. For example, the values of temperature in the kitchen vary too much from the living room. This situation happens due to the activities of cooking made in the kitchen. The same problem occurs with the measured humidity in the bathroom and the other rooms.

Furthermore, the creation of a python script to create a predictor model is faster, more efficient, and the models made will have higher accuracy because the data to process and fit is significantly lower. This fact allows that more simple processors could process and predict all of this data. For example, products like Raspberry Pi or Arduino can be used for this task.

## 5.6 Code

This section finally shows the Support Vector Machine code. This is one of the main objectives of the thesis.

### 5.6.1 Code flow

The direction and how the code is executed is shown in figure 7. This code flow is applied for all the approaches to our project. These approaches are explained in section 5.5.

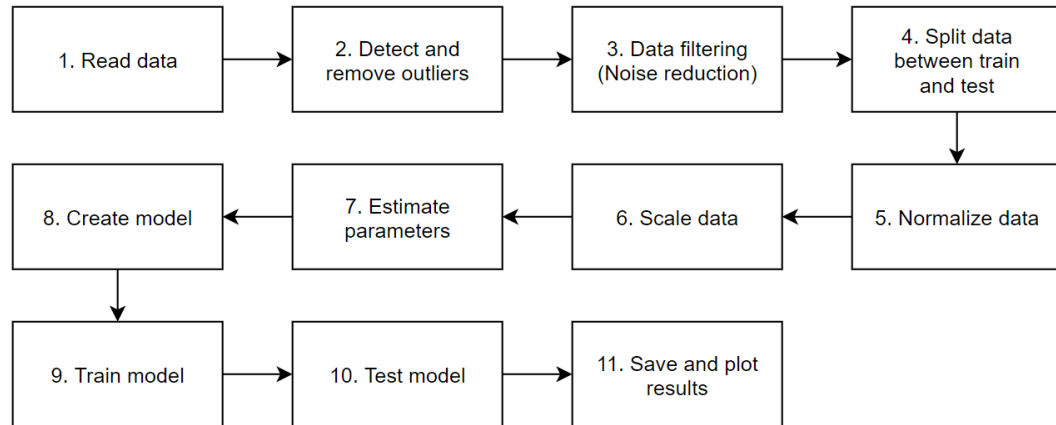


Figure 7: Flowchart of the steps followed in code

### 5.6.2 Code steps

1. **Read data:** Open the csv files and transform the data that they contain.
2. **Detect and remove outliers.**
3. **Data filtering (Noise reduction):** Apply Savitzky–Golay filter.
4. **Split between train and test.**
5. **Normalize data:** Scale all the data into a range from 0 to 1.
6. **Scale data:** Apply a standard scalation to every attribute has the same importance in the algorithm.
7. **Estimate parameters:** Calculate the optimal parameters for the model of each file.
8. **Create model:** With the estimated parameters create the optimal model.
9. **Train model.**
10. **Test model.**
11. **Save and plot results.**



### 5.6.3 Read data

As the first step in our code, we have to read files so we can preprocess them later. For that purpose, we use the library Pandas. This library has the function `read_csv`, which opens a CSV file and loads it as a dataframe. Working with dataframes it is easy and efficient. While the file is loaded into the dataframe, every data is transformed to its type: Float to float, int to int, string to string, etcetera. However, for the attribute date, we have to apply a different process. We use the function contained in the package pandas called `to_datetime`. The code for this part has the same process for every script, though it varies in some details. On line 23, 4 and 9 from files `ensemble_voting.py`, `ensemble_bagging.py` and `svrSplitted_rooms.py` is shown the code for this task.

### 5.6.4 Detect and delete outliers

The outliers are removed from the data with the function `zscore` that Scipy has in their package stats. This function calculates the Z-score for each instance, allowing the posterior removal of those instances that their Z-score is higher than three, considered outliers. In figure 8, is shown the visualization of the scattered data to find possible outliers, that are marked with a black rectangle.

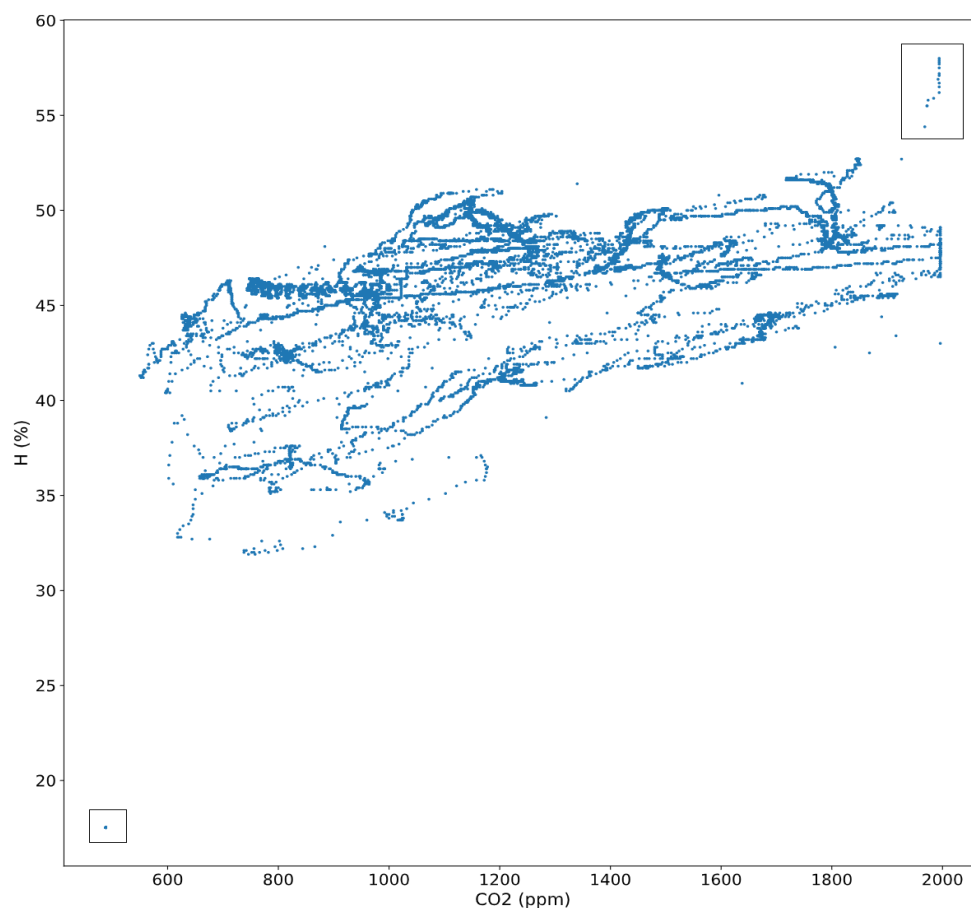


Figure 8: Outlier visualization from Living Room file: CO2 (ppm) vs H(%)

After the detected outliers are removed, the scattered the data looks like shown in figure 9. As it is possible to see all the outliers have been removed.

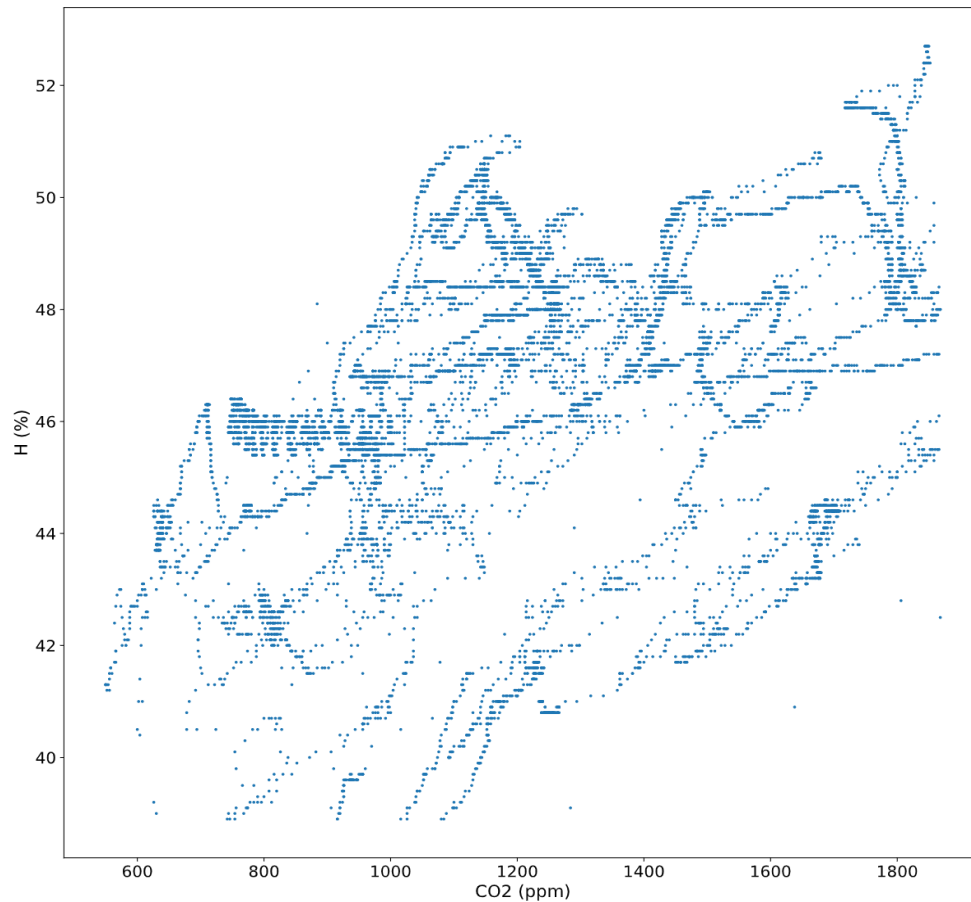


Figure 9: Scattered data after the removal of the outliers from Living Room file: CO2 (ppm) vs H(%)

### 5.6.5 Data filtering (Noise reduction)

Scipy has the package Signal, which has the function `savgol_filter`. This function applies a Savitzky–Golay filter to the data, with the selected window length and polynomial order. We used a window length of 21 and a polynomial order of 1 to reduce all the possible noise that the data signal could have but without distorting the original signal. Additionally, this helps to remove possible outliers that the `zscore` function could not. The result of the applying of this filter is shown in figures 5 and 10.

### 5.6.6 Normalize data

As it is specified in section 5.4.2, the applied formula for normalization is applied to the whole dataframe in only one line. As said in section 5.6.3, working with dataframes is useful and practical. On line 57, 41 and 40 from files `ensemble_voting.py`, `ensemble_bagging.py` and `svr_splitting_rooms.py` is shown the code for this task.

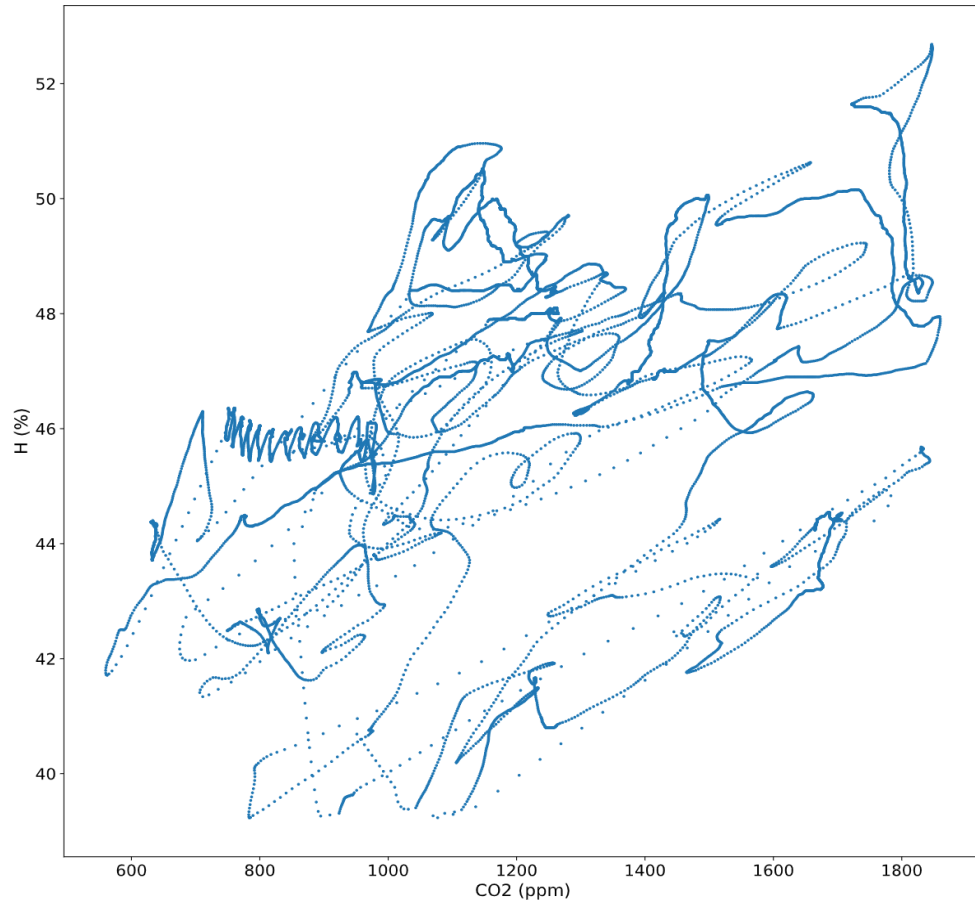


Figure 10: Scattered data after the filter Savitzky–Golay is applied to the data from Living Room file: CO2 (ppm) vs H(%)

### 5.6.7 Scale data

Since we want all the data we are working with, to have the same range, we scale the data. The object `StandardScaler` that Scikit-learn offers scales the data to have the same weight. On lines 108, 48 and 46 from files `ensemble_voting.py`, `ensemble_bagging.py` and `svrSplitted_rooms.py` is shown the code for this task.

### 5.6.8 Split between train and test

Working with the same data to train and test is not right because we might overfit the model. We need new data to test it, so we split it into two sets: Train set and test set.

We selected two different test and train sizes of the dataset to see how the model works in varied circumstances. We set one day (1440 rows) to train and test in the first instance, and five days (7200 rows) to train and five days to test. Because we are working in two different ways to create the ensemble models, two ways to test the ensemble models are made:

1. Two tests for each file: For each file, we will have two test sets, one with a duration of one day and another with a duration of five days.
2. Two global tests: We create two tests from the merged file with random samples with the duration of one day and five days.

To split the data, we use the function `train_test_split`, which will allow us to split between train and test easily. Just with setting the number of instances we want for train and test it will return two objects with train and test subsets with the specified size. On lines 37, 21 and 25 from files `ensemble_voting.py`, `ensemble_bagging.py` and `svrSplitted_rooms.py` is shown the code for this task.

### 5.6.9 Parameter estimation and model training

Because these parameters are problem dependent, so exploration to find the optimal ones, must be done. With the function, `RandomizedSearchCV`, which uses cross-validation, allows the user to enter a parameter grid, and it will randomly select subsets of parameters to test. It will return the best model with the optimal parameters. In the article “A Practical Guide to Support Vector Classification” [39], it says that “We found that trying exponentially growing sequences of  $C$  and  $\gamma$  is a practical method to identify good parameters.”, this is why we use the `logspace` function that the package `NumPy` has. On lines 84 and 57 from files `ensemble_voting.py` and `svrSplitted_rooms.py` is shown the code for this task.

Because of the different ways to process the data, and how the functions from `Scikit-learn` works, we do two different ways to estimate the parameters. For the bagging ensemble approach, the parameters to estimate are not from `SVR`. Instead, we have to estimate the parameters: `max_samples (%)`, `max_features (%)`, and if bootstrap is applied or not. On line 62, the code for this task in file `ensemble_bagging.py` is shown.

Parameter estimation to achieve good results is a real big problem, and this is because it is a problem-dependent. For each problem, there are different optimal parameters. This task after data preprocessing is one of the most significant in the process of data analysis and future prediction.

**5.6.9.1 C** For the estimation of the parameter  $C$ , an exponential range from 1 to 512 has been selected. However, the  $C$  selection is critical because it will influence how good the model will predict. We tried to cover an extensive range for this penalty factor in getting the best possible value.

**5.6.9.2 Epsilon** The values of Epsilon to get good prediction goes on a range from 0 to 1. We decided to use the values: 0.01, 0.05, 0.1, 0.125, 0.2, 0.5, 0.7 to cover most of the range. The best value obtained in most of the trained models is 0.125.

**5.6.9.3 Gamma** For the estimation of the parameter gamma, an exponential range from 1 to 512 has been selected. As we used for the evaluation of the parameter  $C$ . The fact that we had a vast range to

estimate  $C$  and  $\gamma$  made the process of parameter estimation slow, taking up a long time to determine it.

#### **5.6.10 Predict**

With the model created and the parameter estimation finished, we have to predict the CO<sub>2</sub> level from test sets. For this job, the object SVR from Scikit-learn has the function `predict`, which with the test set, will predict all the CO<sub>2</sub> values. On lines 121, 65 and 61 from files `ensemble_voting.py`, `ensemble_bagging.py` and `svr_splitting_rooms.py` is shown the code for this task.

## 6 Discussion

In this section, the experimentation and discussion of the results obtained from the multiple python scripts coded for this project are presented. Some troubles faced, and some improvements discovered are also shown.

### 6.1 Applying shuffle

Once we visualized the data, we noticed that how the data is distributed is not regular and sometimes is day dependent. Because of this, training with specific days does not have sense, and it will lead to bad scores because of the overfitting, so to avoid this situation, we shuffle the data. Obtained results with and without shuffle are shown in table 1 and table 3.

r2	mae	mse	rmse	train size	test size	file
-0.8687	0.2974	0.1314	0.3625	1440	1440	bathroom
-1.0935	0.2493	0.0926	0.3043	7200	7200	bathroom
-1.1107	0.348	0.1785	0.4225	1440	1440	bedroom
-0.26	0.2001	0.0542	0.2328	7200	7200	bedroom
-1.6436	0.3328	0.1434	0.3787	1440	1440	kidsroom
-0.1456	0.1269	0.0483	0.2197	7200	7200	kidsroom
-0.9397	0.2517	0.0921	0.3035	1440	1440	kitchen
-0.392	0.2174	0.0621	0.2493	7200	7200	kitchen
-1.1625	0.3545	0.1549	0.3936	1440	1440	livingroom
-0.3731	0.1323	0.0326	0.1806	7200	7200	livingroom

Table 1: SVR results with one model per room and non shuffled data. Used two test and train sets per file (1440 instances that equal one day, and 7200 instances that equals five days)

When we shuffle the data, SVR has an improvement in the predictive performance and a higher quality of the model. This is because when we shuffle the data, we avoid some patterns while we split the data, having various data in the train and test set. As it is possible to observe, the improvement of the R2 score is enormous, almost 200%.

### 6.2 Applying Savitzky–Golay filter

Applying Savitzky–Golay filter was a great improvement in our problem to predict CO2 levels. This is because of all the noise that was on the dataset. With this filtering, all the noise is erased, smoothed, and we get the signal without distorting it. This filtered signal is shown in figure 5, and the scatter of the filtered data is shown in figure 10.

There is a large difference from the scores obtained without shuffling and filtering the data. The obtained results after both tasks finished are shown in table 3.

r2	mae	mse	rmse	train size	test size	file
0.5619	0.1169	0.0295	0.1716	1440	1440	bathroom
0.7119	0.0726	0.0189	0.1376	7200	7200	bathroom
0.7381	0.0739	0.0164	0.1282	1440	1440	bedroom
0.8119	0.0529	0.0113	0.1065	7200	7200	bedroom
0.8284	0.046	0.0087	0.0932	1440	1440	kidsroom
0.8747	0.0309	0.0065	0.0807	7200	7200	kidsroom
0.5417	0.1211	0.0262	0.1618	1440	1440	kitchen
0.6095	0.0968	0.0221	0.1486	7200	7200	kitchen
0.7593	0.0614	0.0122	0.1104	1440	1440	livingroom
0.8298	0.0419	0.0089	0.0942	7200	7200	livingroom

Table 2: SVR results trained and tested with one day, and five days of the data (1440 instances, and 7200 instances). The data was not filtered with the Savitzky–Golay filter.

r2	mae	mse	rmse	train size	test size	file
0.7748	0.0653	0.0149	0.122	1440	1440	bathroom
0.7597	0.0694	0.0155	0.1246	7200	7200	bathroom
0.8931	0.0339	0.007	0.0838	1440	1440	bedroom
0.8031	0.0586	0.0125	0.1118	7200	7200	bedroom
0.9371	0.0165	0.0036	0.0599	1440	1440	kidsroom
0.9173	0.0219	0.0048	0.0694	7200	7200	kidsroom
0.7014	0.0739	0.0168	0.1298	1440	1440	kitchen
0.7045	0.0722	0.0165	0.1285	7200	7200	kitchen
0.9204	0.0238	0.004	0.0635	1440	1440	livingroom
0.8948	0.0288	0.0055	0.0743	7200	7200	livingroom

Table 3: SVR results with one model per room, shuffled data and filtered with Savitzky–Golay filter. Used two test and train sets per file (1440 instances that equal one day, and 7200 instances that equals five days)

## 6.3 Implementations

We compare three implementations:

- Bagging ensemble: Using the bagging ensemble to train and test, explained in section 5.5.1.
- Voting ensemble: Using the voting ensemble to train and test, explained in section 5.5.2.
- One model per room: Using different models, one per room, to train and test, explained in section 5.5.3.

### 6.3.1 Bagging ensemble vs. Voting ensemble

Since we are working in two different ways to train and test the models (Merging the data of each file and working with it separately), the tables shown are different.

Operating with the voting ensemble made things easier because we can select which room we want to test. This helps set the algorithm in the situation of predicting CO2 values. But, having all the data merged is messy, and we do not know from which room the data selected to train and test comes. That is one of the main reasons why the results from the bagging ensemble are lower than the results from the voting ensemble.

r2	mae	mse	rmse	train size	test size	file
0.4737	0.1186	0.0332	0.1821	1440	1440	bathroom
0.3703	0.1349	0.0388	0.1969	7200	7200	bathroom
0.7635	0.0802	0.0153	0.1238	1440	1440	bedroom
0.7032	0.0919	0.0186	0.1363	7200	7200	bedroom
0.7278	0.0622	0.0146	0.1206	1440	1440	kidsroom
0.6726	0.0731	0.0179	0.1337	7200	7200	kidsroom
0.4759	0.1275	0.0317	0.1781	1440	1440	kitchen
0.4387	0.1355	0.0337	0.1835	7200	7200	kitchen
0.7484	0.0682	0.0127	0.1129	1440	1440	livingroom
0.6882	0.0785	0.0164	0.128	7200	7200	livingroom

Table 4: Voting ensemble results. The results from this ensemble have an average R2 score of 0.60

The results obtained from the voting ensemble are shown in table 4. Although the results from this ensemble are low, this happens because there are two rooms in which their data has more unpredictable patterns. These rooms are the kitchen and the bathroom.



r2	mae	mse	rmse	train size	test size	file
0.4336	0.1268	0.0262	0.1617	1440	1440	merged data
0.4003	0.1335	0.0283	0.1683	7200	7200	merged data

Table 5: Bagging ensemble results. The results from this ensemble have an average R2 score of 0.41

Values predicted for both algorithms are shown in figures 11 and 12. As it is possible to see in both figures, the metrics used to evaluate the models are represented in the predicted values. Since the better the results of the metrics, the closer the predicted values are to the true values.

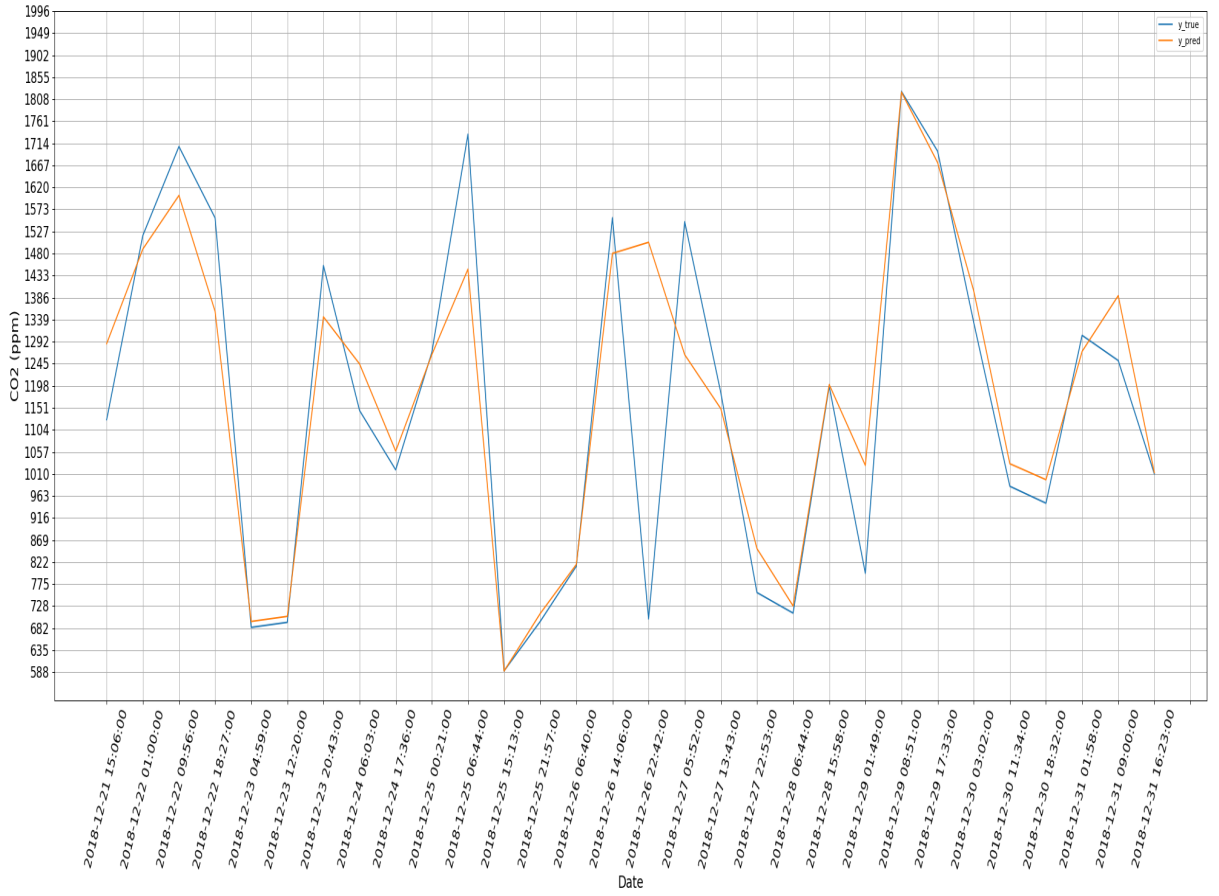


Figure 11: Predictions made by the best voting ensemble model. Train and test size: 1440 instances (one day). R2 score: 0.76, MAE: 0.09, MSE: 0.01, RMSE: 0.13

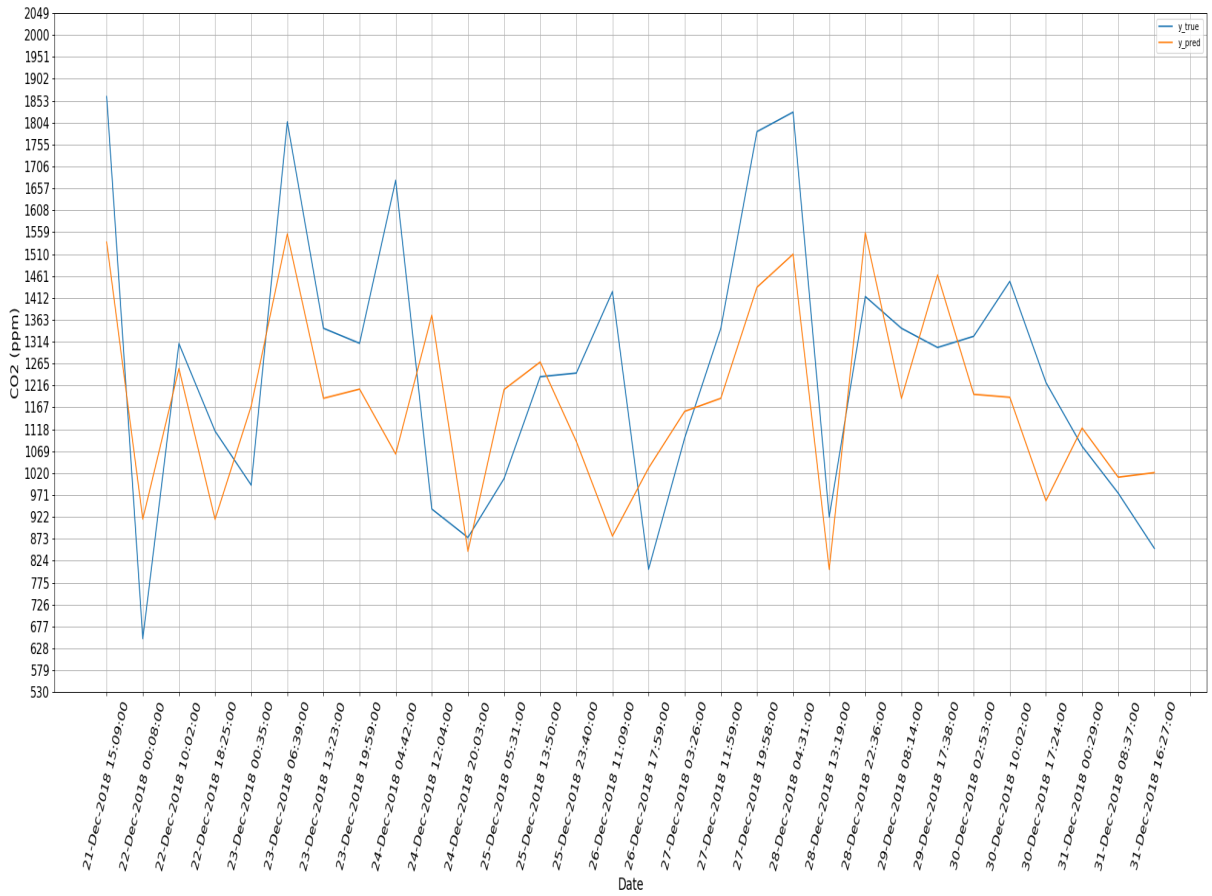


Figure 12: Predictions made by the best bagging ensemble model. Train and test size: 1440 instances (one day). R2 score: 0.43, MAE: 0.12, MSE: 0.02, RMSE: 0.16

### 6.3.2 Working without kitchen and bathroom rooms

Comparing results achieved by the ensemble working with and without the kitchen and bathroom, we found a performance improvement of 18% in the mean of the R2 score. Both rooms have elements (Cooking stove, shower, etcetera), which make drastic changes in the levels of temperature and humidity measures made by the sensors. The same situation happens with the bagging ensemble.

r2	mae	mse	rmse	train size	test size	file
0.8253	0.0661	0.0127	0.1125	1440	1440	bedroom
0.7913	0.0754	0.015	0.1227	7200	7200	bedroom
0.8318	0.0446	0.0105	0.1023	1440	1440	kidsroom
0.8061	0.0522	0.0124	0.1114	7200	7200	kidsroom
0.7863	0.063	0.0127	0.1126	1440	1440	livingroom
0.6967	0.0762	0.0176	0.1328	7200	7200	livingroom

Table 6: Voting ensemble results working without bathroom and kitchen rooms. The results from this ensemble have an average R2 score of 0.78

It is easy to see that the voting ensemble seems more efficient than the bagging ensemble in any situation presented. However, even if the voting ensemble works better, maybe with another data processing to create the bagging ensemble, the performance could improve. Even so, the results obtained by the bagging ensemble have an improvement of 4%, which is not remarkable, but it is still an improvement.

r2	mae	mse	rmse	train_size	test_size	file
0.4754	0.1368	0.0314	0.1772	1440	1440	merged data
0.4577	0.1388	0.0324	0.18	7200	7200	merged data

Table 7: Bagging ensemble results working without bathroom and kitchen rooms. The results from this ensemble have an average R2 score of 0.46

But working without the bathroom and kitchen rooms made the metrics MSE and RMSE higher than working with them. This means that these models have a worse fitting to the data. However, the higher R2 score indicates that most of the data used as the test can be explained, about 83% of the data can be explained.

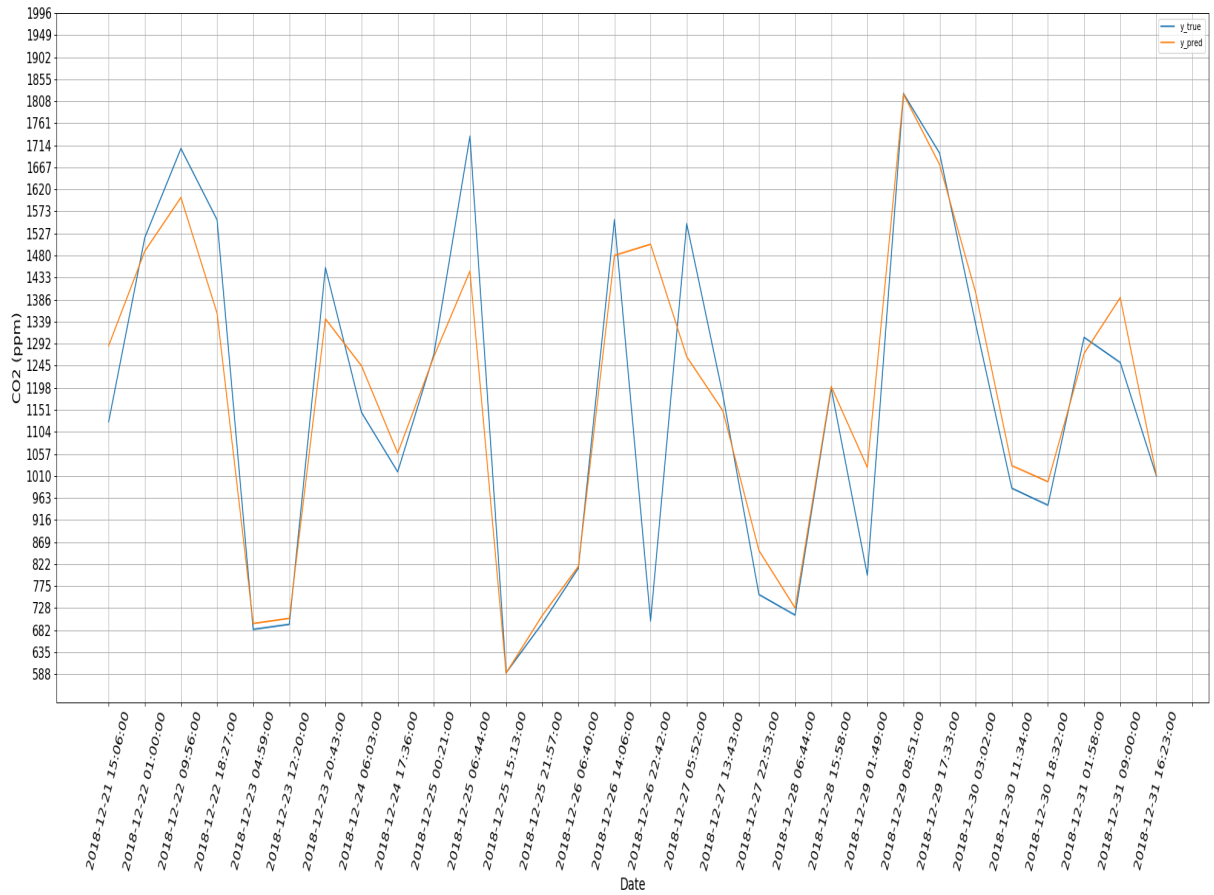


Figure 13: Predictions made by the best voting ensemble model working without bathroom and kitchen rooms. Train and test size: 1440 instances (one day). R2 score: 0.83, MAE: 0.04, MSE: 0.01, RMSE: 0.10

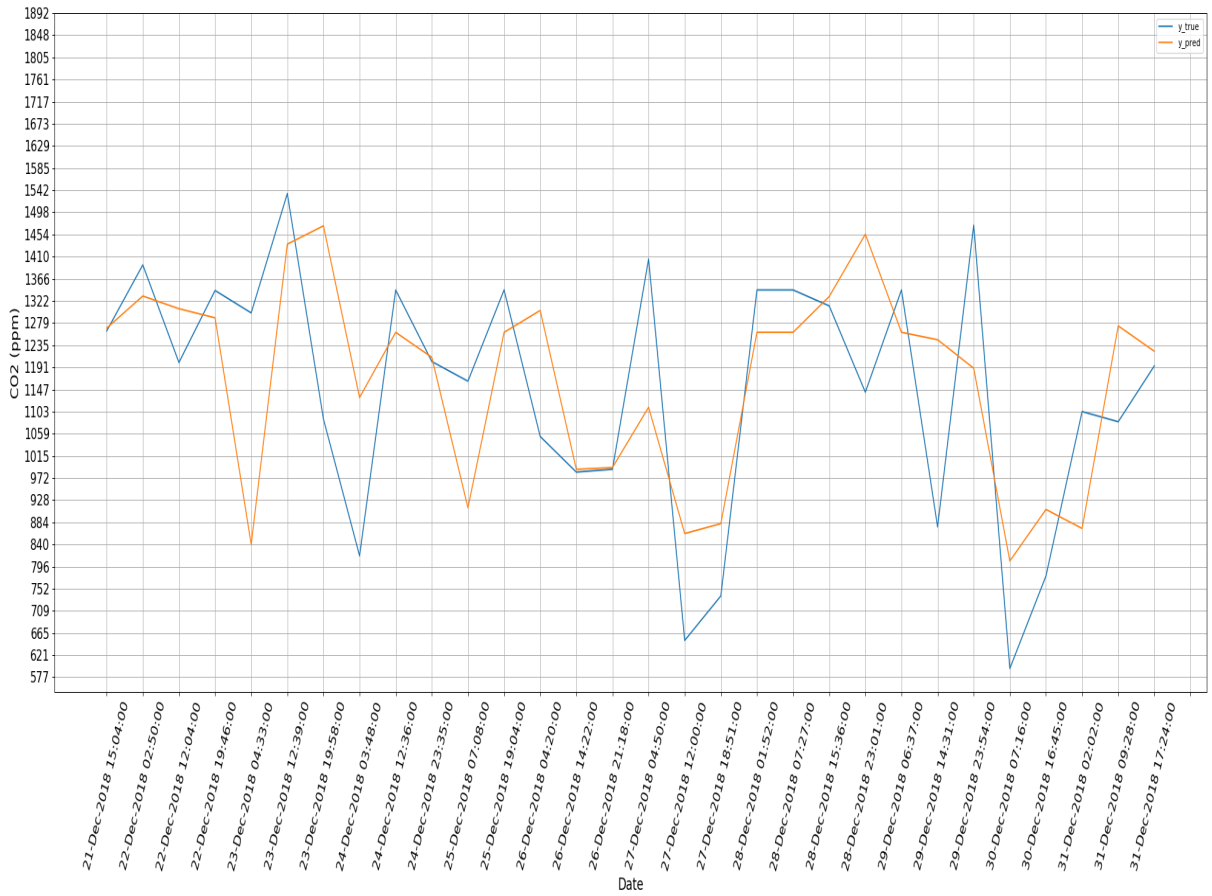


Figure 14: Predictions made by the best bagging ensemble model working without bathroom and kitchen rooms. Train and test size: 1440 instances (one day). R2 score: 0.47, MAE: 0.13, MSE: 0.03, RMSE: 0.17

### 6.3.3 One model per room

Although the ensembles are more robust and stable models, it seems that to work with this problem, the most efficient way is to work with each room individually, creating a model for each one of them. The results of working with a model per room are shown in table 3. These results represent a fairly higher performance improvement, up to 11% improvement in predicting CO2 levels. This happens because apart from the fact that there are elements in each room that alter the measurements taken by the sensors, each room has different sizes and quite different usage patterns. This increases the disparity between rooms data, which means that it is more challenging to create a model that is capable of predicting CO2 values for all rooms. In figure 15, is shown how an R2 score of 0.93 is almost a success in the predictions of patterns used for the test. However, if we look at table 3, when we train with five days, and we predict five days, the results are generally lower, this happens since the model is overfitted.

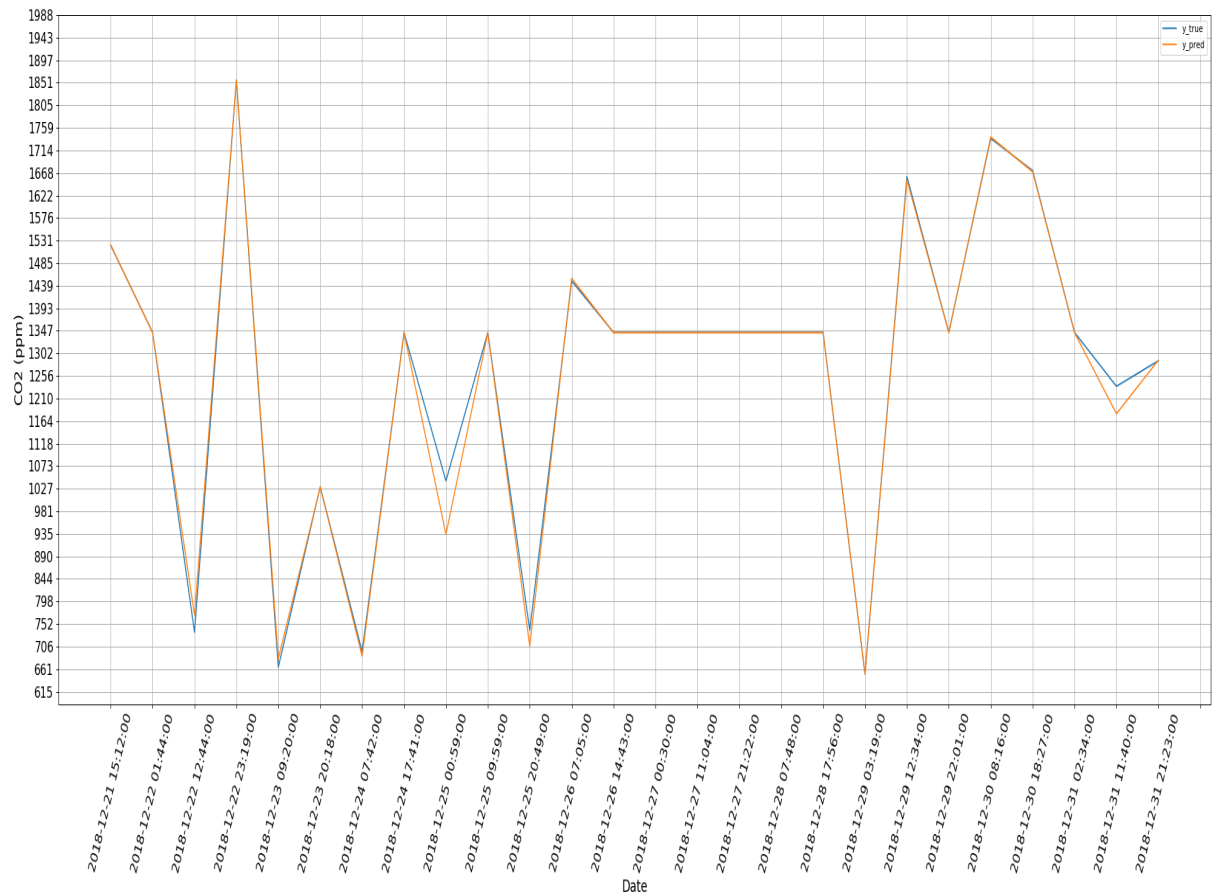


Figure 15: Predictions made by creating one model per room. A total of five models were created. Train and test size: 1440 instances (one day). R2 score: 0.93, MAE: 0.01, MSE: 0.003, RMSE: 0.0599

## **7 Conclusions**

## References

1. LASHKARI, Behzad; CHEN, Yuxiang; MUSILEK, Petr. *Energy Management for Smart Homes-State of the Art*. SWITZERLAND: MDPI, 2019.
2. PENG, Yuzhen; RYSANEK, Adam; NAGY, Zoltán; SCHLÜTER, Arno. Using machine learning techniques for occupancy-prediction-based cooling control in office buildings. *Applied Energy*. 2018, vol. 211, pp. 1343–1358. ISSN 0306-2619. Available from DOI: <https://doi.org/10.1016/j.apenergy.2017.12.002>.
3. WANG, Shengwei. Intelligent buildings and building automation. *Intelligent Buildings and Building Automation*. 2009-01, pp. 1–248. Available from DOI: 10.4324/9780203890813.
4. ZADEH, LOTFI A. Soft Computing and Fuzzy Logic. In: *Fuzzy Sets, Fuzzy Logic, and Fuzzy Systems*, pp. 796–804. Available from DOI: 10.1142/9789814261302\_0042.
5. PRATIHAR, D.K. *Soft Computing*. S.l. Alpha Science International, Ltd, 2007.
6. WOLPER, David H.; MACREADY, William G. *No Free Lunch Theorems for Optimization*. San Jose: IBM Almaden Research Center, 1996.
7. BLUMTRITT, Christoph. *Smart Home Report 2019*. London: Statista, 2019.
8. VADILLO, Laura. The Role of Smart Homes in Intelligent Homecare and Healthcare Environments. In: *book auth.] Health Sciences Elsevier. Ambient Assisted Living and Enhanced Living Environments*. Woburn, United States: Butterworth-Heinemann, 2016.
9. CLARK, Jordan D.; LESS, Brennan D.; DUTTON, Spencer M.; WALKER, Lain S.; SHERMAN, Max H. (eds.). *Efficacy of occupancy-based smart ventilation control strategies in energy-efficient homes in the United States*. California: BUILDING and ENVIRONMENT, 2019.
10. BELTRAN, Alex; CERPA, Alberto E. Optimal HVAC Building Control with Occupancy Prediction. In: *Proceedings of the 1st ACM Conference on Embedded Systems for Energy-Efficient Buildings*. New York: Association for Computing Machinery, 2014.
11. A machine learning approach for indirect human presence detection using IOT devices. In: *2016 Eleventh International Conference on Digital Information Management (ICDIM)*. 2016.
12. RAZAVI, Rouzbeh; GHARIPOUR, Amin; FLEURY, Martin; AKPAN. Occupancy detection of residential buildings using smart meter data: A. *Ikpe Justice. s.l. : Elsevier*. 2019, vol. 183, pp. 0378–7788.
13. SHIH, Huang-Chia. A robust occupancy detection and tracking algorithm for the automatic monitoring and commissioning of a building. *Energy and Buildings*. 2014.
14. HARROU, F.; ZERROUKI, N.; SUN, Y.; HOUACINE, A. An Integrated Vision-Based Approach for Efficient Human Fall Detection in a Home Environment. *IEEE Access*. 2019, vol. 7, pp. 114966–114974. ISSN 2169-3536. Available from DOI: 10.1109/ACCESS.2019.2936320.



15. BARBOUR, E.; DAVILA, C.C.; GUPTA, S.; REINHART, C.; KAUR, J.; GONZÁLEZ (eds.). *Planning for sustainable cities by estimating building occupancy with mobile phones*. Nature Publishing Group, 2019.
16. FAHAD, Labiba Gillani; KHAN, Asifullah; RAJARAJAN, Muttukrishnan sl (eds.). *Activity recognition in smart homes with self verification of assignments*. Neurocomputing, 2015.
17. *Monitoring Activities of Daily Living in Smart Homes: Understanding human behavior*. 2, s.l. IEEE, 2016.
18. CALÌ, Davide; MATTHES, Peter; HUCHTEMANN, Kristian; STREBLOW, Rita; MÜLLER, Dirk. CO2 based occupancy detection algorithm: Experimental analysis and validation for office and residential buildings. *Building and Environment*. 2015, vol. 86, pp. 39–49. ISSN 0360-1323. Available from DOI: <https://doi.org/10.1016/j.buildenv.2014.12.011>.
19. ARIEF-ANG, Irvan; HAMILTON, Margaret; SALIM, Flora. A Scalable Room Occupancy Prediction with Transferable Time Series Decomposition of CO2 Sensor Data. *ACM Transactions on Sensor Networks*. 2018-11, vol. 14. Available from DOI: 10.1145/3217214.
20. BRENNAN, Colin; TAYLOR, Graham; SPACHOS, P.sl (eds.). *Designing Learned CO2-based Occupancy Estimation in Smart Buildings*. IET Wireless Sensor Systems, 2018.
21. CANDANEDO, Luis M.; FELDHEIM, Véronique sl (eds.). *Accurate occupancy detection of an office room from light, temperature, humidity and CO2 measurements using statistical learning models*. Energy and Buildings, 2016.
22. Vázquez and W. Clustering methods for occupancy prediction in smart home control. In: *2011 IEEE International Symposium on Industrial Electronics*. 2011.
23. ANAND, P.; CHEONG, D.; SEKHAR, C.; SANTAMOURIS, M.; KONDEPUDI, S.sl (eds.). *Energy saving estimation for plug and lighting load using occupancy analysis*. Elsevier, 2019.
24. WANG, Wei; CHEN, Jiayu; HONG, Tianzhen (eds.). *Occupancy prediction through machine learning and data fusion of environmental sensing and Wi-Fi sensing in buildings*. 2018.
25. CORTES, Corinna; VAPNIK, Vladimir. Support-Vector Networks. *Mach. Learn.* 1995-09, vol. 20, no. 3, pp. 273–297. ISSN 0885-6125. Available from DOI: 10.1023/A:1022627411411.
26. HEFFERON, Jim. *LINEAR ALGEBRA*. Colchester, Vermont: Mathematics and Statistics Department Saint Michael's College, 2017.
27. PARRELLA, Francesco. Online Support Vector Regression. In: 2007.
28. HASTIE, Trevor; TIBSHIRANI, Robert; FRIEDMAN, Jerome. *The elements of statistical learning: data mining, inference and prediction*. 2nd ed. Springer, 2009. Available also from: <http://www-stat.stanford.edu/~tibs/ElemStatLearn/>.
29. HUH, Myung-Hoe. Kernel-Trick Regression and Classification. *Communications for Statistical Applications and Methods*. 2015-03, vol. 22, pp. 201–207. Available from DOI: 10.5351/CSAM.2015.22.2.201.

30. MINH, Ha Quang; NIYOGI, Partha; YAO, Yuan. Mercer's theorem, feature maps, and smoothing. In: *International Conference on Computational Learning Theory*. 2006, pp. 154–168.
31. HOFMANN, Thomas; SCHÖLKOPF, Bernhard; SMOLA, Alexander J. Kernel methods in machine learning. *The Annals of Statistics*. 2008-07, vol. 36, no. 3, pp. 1171–1220. ISSN 0090-5364. Available from DOI: 10.1214/009053607000000677.
32. GENTON, Marc G. Classes of Kernels for Machine Learning: A Statistics Perspective. *J. Mach. Learn. Res.* 2002-03, vol. 2, pp. 299–312. ISSN 1532-4435.
33. ALPAYDIN, Ethem. Introduction to Machine Learning. In: 2004-01, vol. 56.
34. CELIKYILMAZ, Asli; TRKSEN, I. Burhan. *Modeling Uncertainty with Fuzzy Logic: With Recent Theory and Applications*. 1st. Springer Publishing Company, Incorporated, 2009. ISBN 3540899235.
35. HAN, Jiawei; KAMBER, Micheline; PEI, Jian. *Data mining concepts and techniques, third edition*. Waltham, Mass.: Morgan Kaufmann Publishers, 2012. ISBN 0123814790.
36. SAVITZKY, Abraham.; GOLAY, M. J. E. Smoothing and Differentiation of Data by Simplified Least Squares Procedures. *Analytical Chemistry*. 1964, vol. 36, no. 8, pp. 1627–1639. Available from DOI: 10.1021/ac60214a047.
37. HARRISON, M. *Machine Learning Pocket Reference: Working with Structured Data in Python*. O'Reilly Media, 2019. ISBN 9781492047513.
38. ZHANG, Cha; MA, Yunqian. *Ensemble Machine Learning: Methods and Applications*. Springer Publishing Company, Incorporated, 2012. ISBN 1441993258.
39. HSU, Chih-Wei; CHANG, Chih-Chung; LIN, Chih-Jen. *A Practical Guide to Support Vector Classification*. 2003. Available also from: <http://www.csie.ntu.edu.tw/~cjlin/papers.html>. Technical report. Department of Computer Science, National Taiwan University.

## A Code implementation

### A.1 Ensemble voting

---

```
1 data_path = 'data/'
2 onlyfiles = [f for f in listdir (data_path) if isfile (join(data_path, f))]
3
4 output = []
5
6 test_size = [1440, 7200]
7
8 for i in test_size :
9
10     dict_test_X = {}
11     dict_test_y = {}
12     dict_test_dates = {}
13
14     list_train_X = []
15     list_train_y = []
16
17     estimators = []
18
19     for file in onlyfiles :
20
21         #Read file and change date type
22         filename = os.path.splitext ( file )[0]
23         dataset = pd.read_csv('data/'+filename+'.csv', delimiter=";", parse_dates=True)
24         dataset ['Time'] = pd.to_datetime ( dataset .Time)
25         print ("Working with file : ", filename)
26
27         #Remove outliers
28         dataset = wrapper.remove_outliers ( dataset )
29
30         #Filter dataset
31         dataset = wrapper.filter_data ( dataset , 21, 1)
32         #Get min and max
33         min_CO2=dataset['CO2'].min()
34         max_CO2=dataset['CO2'].max()
35
36         #Split between train and test
37         data_train , data_test = train_test_split ( dataset , train_size=i, test_size=i, random_state=55, shuffle =
            True)
38
39         #Split between X and y
40         X_train = data_train .drop(['CO2', 'Time'], axis=1).copy()
41         y_train = data_train ['CO2']
42
43         X_test = data_test .drop(['CO2', 'Time'], axis=1).copy()
```

```

44     y_test = data_test [['CO2', 'Time']]
45
46     date_test = y_test['Time']
47
48     y_test = y_test['CO2'].to_numpy()
49
50     #Normalize data
51     min_H = X_train['H'].min()
52     max_H = X_train['H'].max()
53
54     min_T = X_train['T'].min()
55     max_T = X_train['T'].max()
56
57     X_train['H'] = (X_train['H']-min_H)/(max_H-min_H)
58     X_train['T'] = (X_train['T']-min_T)/(max_T-min_T)
59
60     X_test['H'] = (X_test['H']-min_H)/(max_H-min_H)
61     X_test['T'] = (X_test['T']-min_T)/(max_T-min_T)
62
63     # Save the train data in order to train the ensemble later
64     list_train_X.append(X_train)
65     list_train_y.append(y_train)
66
67     # Save the test data in order to test the ensemble later
68     dict_test_X[filename] = X_test
69     dict_test_y[filename] = y_test
70     dict_test_dates[filename] = date_test
71
72     #Scale data
73     sc_X = StandardScaler()
74     X_train = sc_X.fit_transform(X_train)
75     X_test = sc_X.transform(X_test)
76
77     # Select tuning parameter range
78     gamma_range = np.logspace(0.0001,9,base=2, num=40)
79     C_range = np.logspace(0.0001,9,base=2, num=40)
80     epsilon_range = [0.01, 0.05, 0.1, 0.125, 0.2, 0.5, 0.7, 0.9]
81     tuned_parameters = [{'kernel': ['rbf'], 'gamma': gamma_range, 'C': C_range, 'epsilon': epsilon_range}]
82
83     #Regressor parameter estimation
84     regressor = RandomizedSearchCV(SVR(), tuned_parameters, scoring='r2', cv=3, verbose=1, n_iter=10)
85     regressor.fit(X_train, y_train)
86
87     estimators.append([filename, regressor.best_estimator_])
88
89     #Predict test values
90     y_true, y_pred = y_test, regressor.best_estimator_.predict(X_test)
91

```

```

102     #Normalize outputs to get normalized scores
103     y_true_norm = [ ((i-min_CO2)/(max_CO2-min_CO2)) for i in y_true ]
104     y_pred_norm = [ ((i-min_CO2)/(max_CO2-min_CO2)) for i in y_pred ]
105
106     #Calculate the score
107     results = regression_results (y_true_norm, y_pred_norm)
108     results [ ' train_size ' ] = i
109     results [ ' test_size ' ] = i
110     results [ ' file ' ] = filename
111     print (pd.DataFrame(results , index=[0]))
112
113     #Concat the train data in order to use for training the ensemble
114     ens_X_train = pd.concat( list_train_X )
115     ens_y_train = pd.concat( list_train_y )
116
117     #Scale train data
118     sc_X = StandardScaler ()
119     ens_X_train = sc_X. fit_transform (ens_X_train)
120
121     # Create ensemble model and fit train data in it
122     print ( ' Fitting data ... ' )
123     ensemble_model = VotingRegressor( estimators = estimators )
124     ensemble_model.fit(ens_X_train, ens_y_train)
125
126     #Test data from the different rooms
127     for X_index, y_index, date_index in zip( dict_test_X , dict_test_y , dict_test_dates ):
128         X_test = sc_X.transform( dict_test_X [X_index])
129         y_test = dict_test_y [y_index]
130         date_test = dict_test_dates [date_index]
131         y_true, y_pred = y_test , ensemble_model.predict(X_test)
132
133         #Normalize outputs to get normalized scores
134         y_true_norm = [ ((i-min_CO2)/(max_CO2-min_CO2)) for i in y_true ]
135         y_pred_norm = [ ((i-min_CO2)/(max_CO2-min_CO2)) for i in y_pred ]
136
137         results = wrapper. regression_results (y_true_norm, y_pred_norm)
138         results [ ' train_size ' ] = i
139         results [ ' test_size ' ] = i
140         results [ ' file ' ] = X_index
141         output.append( results )
142         wrapper.plot_accuracy (y_true, y_pred, int(y_true.size/30), i, i, results [ 'r2' ], X_index, sorted(
143             date_test ))
144
145     results_df = pd.DataFrame.from_dict(output)
146     print ( results_df )
147
148     results_df .to_csv( ' predictions /ensemble_voting.csv' , sep=';')

```

---

Code Listing 1: Python script for the creation of an ensemble using voting ensemble approach explained in section 5.5

## A.2 Ensemble bagging

---

```
1
2 #Read file
3 filename = 'merged'
4 dataset = pd.read_csv('merged_data/'+filename+'.csv', delimiter=";", parse_dates=True)
5 dataset = dataset.drop(['id'], axis=1)
6
7 dataset = wrapper.remove_outliers(dataset)
8
9 #Get min and max
10 min_CO2 = dataset['CO2'].min()
11 max_CO2 = dataset['CO2'].max()
12
13 #Filter dataset data
14 dataset = wrapper.filter_data(dataset, 21, 1)
15
16 output = []
17 tests = [1440, 1440*5]
18 regressor = []
19 for i in range(len(tests)):
20     #Split between train and test
21     data_train, data_test = train_test_split(dataset, train_size=tests[i], test_size=tests[i], random_state=55,
22                                               shuffle=True)
23
24     #Split between X and y
25     X_train = data_train.drop(['CO2', 'Time'], axis=1).copy()
26     y_train = data_train['CO2']
27
28     X_test = data_test.drop(['CO2', 'Time'], axis=1).copy()
29     y_test = data_test[['CO2', 'Time']]
30
31     date_test = y_test['Time']
32
33     y_test = y_test['CO2'].to_numpy()
34
35     #Normalize data
36     min_H = X_train['H'].min()
37     max_H = X_train['H'].max()
38
39     min_T = X_train['T'].min()
40     max_T = X_train['T'].max()
41
42     X_train['H'] = (X_train['H']-min_H)/(max_H-min_H)
43     X_train['T'] = (X_train['T']-min_T)/(max_T-min_T)
44
45     X_test['H'] = (X_test['H']-min_H)/(max_H-min_H)
46     X_test['T'] = (X_test['T']-min_T)/(max_T-min_T)
```

```

46
47 #Scale data
48 sc_X = StandardScaler ()
49 X_train = sc_X.fit_transform (X_train)
50 X_test = sc_X.transform(X_test)
51
52 #If first iteration , create the model and estimate parameters
53 if i ==0:
54     rng = check_random_state(0)
55     grid = ParameterGrid({"max_samples": [0.5, 1.0],
56                           "max_features": [0.5, 1.0],
57                           "bootstrap": [True, False]})
58     for params in grid:
59         regressor = BaggingRegressor(base_estimator=SVR(),
60                                     verbose=10,
61                                     random_state=rng,
62                                     **params).fit (X_train, y_train )
63
64 #Predict the test data with the best values found
65 y_true, y_pred = y_test, regressor . predict (X_test)
66
67 #Normalize outputs to get normalized scores
68 y_true_norm = [ ((i - min_CO2)/(max_CO2 - min_CO2)) for i in y_true ]
69 y_pred_norm = [ ((i - min_CO2)/(max_CO2 - min_CO2)) for i in y_pred ]
70
71 #Save the results for each test
72 results = wrapper. regression_results (y_true_norm, y_pred_norm)
73 results [ ' train_size ' ] = tests [i]
74 results [ ' test_size ' ] = tests [i]
75 results [ ' file ' ] = 'merged'
76 output.append( results )
77 wrapper.plot_accuracy (y_true, y_pred, int(y_pred.size/30), tests [i], tests [i], results [ 'r2' ], 'merged',
78                        sorted( date_test ) )
79
80 results_df =pd.DataFrame.from_dict(output)
81 print( results_df )
82 results_df .to_csv(' predictions /ensemble_bagging.csv', sep=';')

```

---

Code Listing 2: Python script for the creation of an ensemble using bagging ensemble approach explained in section 5.5



### A.3 One model per room

---

```
1 data_path = 'data/'
2 onlyfiles = [f for f in listdir (data_path) if isfile (join (data_path, f))]
3
4 output = []
5 for file in onlyfiles :
6
7     #Read file and change date type
8     filename = os.path.splitext ( file )[0]
9     dataset = pd.read_csv('data/'+filename+'.csv', delimiter=";", parse_dates=True)
10    dataset['Time'] = pd.to_datetime ( dataset .Time)
11
12    dataset = wrapper.remove_outliers ( dataset )
13
14    #Filter dataset
15    dataset = wrapper.filter_data ( dataset , 21, 1)
16
17    test_size = [1440, 1440*5]
18    for i in test_size :
19
20        #Split between X and y
21        X = dataset.drop(['CO2', 'Time'], axis=1)
22        y = dataset [['CO2', 'Time']]
23
24        #Split between train and test
25        X_train, X_test, y_train, y_test = train_test_split (X, y, train_size=i, test_size=i, random_state=55,
26                                                                shuffle = True)
27        X_train = X_train.copy()
28        X_test = X_test.copy()
29        date_test = y_test['Time']
30
31        y_test = y_test['CO2'].to_numpy()
32        y_train = y_train['CO2'].to_numpy()
33
34        #Scale data
35        min_H = X_train['H'].min()
36        max_H = X_train['H'].max()
37
38        min_T = X_train['T'].min()
39        max_T = X_train['T'].max()
40
41        X_train['H'] = (X_train['H']-min_H)/(max_H-min_H)
42        X_train['T'] = (X_train['T']-min_T)/(max_T-min_T)
43
44        X_test['H'] = (X_test['H']-min_H)/(max_H-min_H)
45        X_test['T'] = (X_test['T']-min_T)/(max_T-min_T)
```

```

46     sc_X = StandardScaler ()
47     X_train = sc_X.fit_transform (X_train)
48     X_test = sc_X.transform(X_test)
49
50     # Select tuning parameter range
51     gamma_range = np.logspace(0.0001,9,base=2, num=40)
52     C_range = np.logspace (0.0001,9, base=2, num=40)
53     epsilon_range = [0.01, 0.05, 0.1, 0.125, 0.2, 0.5, 0.7, 0.9]
54     tuned_parameters = [{ 'kernel': [ 'rbf' ], 'gamma': gamma_range, 'C': C_range, 'epsilon': epsilon_range  }]
55
56     #Regressor parameter estimation
57     regressor = RandomizedSearchCV(SVR(), tuned_parameters, scoring='r2', cv=3, verbose=1, n_iter=10)
58     regressor . fit (X_train , y_train )
59
60     #Predict test values
61     y_true , y_pred = y_test , regressor . predict (X_test)
62
63     #Calculate the score
64     results = wrapper. regression_results (y_true , y_pred)
65     results [ ' train_size ' ] = i
66     results [ ' test_size ' ] = i
67     results [ ' file ' ] = filename
68     output.append( results )
69     wrapper.plot_accuracy (y_true , y_pred, int(y_true.size/50), i, i, results [ 'r2' ], filename , sorted (
        date_test ))
70
71     results_df = pd.DataFrame.from_dict(output)
72     print ( results_df )
73     results_df .to_csv(' predictions / results_shuffle .csv', sep=';')

```

---

Code Listing 3: Python script for the creation of an ensemble using one model per room approach explained in section 5.5

## A.4 Wrapper

```
1 def plot_accuracy(y_true, y_pred, xax_, t_s, train_size, score, filename, folder, dates):
2     """Function that plots the predicted and true values on the y ax, and the date in the x ax.
3     Also saves the true and predicted in a csv file .
4
5     Arguments:
6         y_true { list } -- list with true values of CO2
7         y_pred { list } -- list with predicted values of CO2
8         xax_ { int } -- step of the y ax
9         t_s { float } -- test size
10        train_size { float } -- train size
11        score { float } -- R2 score obtained
12        filename { string } -- name of the file which has been tested
13        folder { string } -- name of the folder where the plot is going to be saved
14        dates { list } -- list of date to print in the x ax
15    """
16    xax=xax_
17    fig=plt.figure( figsize=(25, 12), dpi= 80, facecolor='w', edgecolor='k')
18    plt.plot(y_true[0::xax], label='y_true')
19    plt.plot(y_pred[0::xax], label='y_pred')
20    min_y_bound = min(min(y_true), min(y_true))
21    max_y_bound = max(max(y_true), max(y_true))
22    step = abs(max_y_bound-min_y_bound)/30
23
24    plt.yticks(np.arange(min_y_bound, max_y_bound+step, step), fontsize=15)
25    plt.xticks(np.arange(0, y_true[0::xax].size+1,1), dates[0::xax], rotation=70, fontsize=15)
26    plt.xlabel('Date', fontsize=16)
27    plt.ylabel('CO2 (ppm)', fontsize=16)
28    plt.legend()
29    plt.grid()
30    plt.tight_layout()
31
32    plt.savefig("graphs/"+folder+"/"+filename+"_score{score:.3f}_test{test_s:d}_train{train:d}.png".format(
33        test_s=t_s, score = score, train = train_size))
34    fname = "predictions/"+folder+"/"+filename+"_score{score:.3f}_test{test_s:d}_train{train:d}.csv".format(
35        test_s=t_s, score = score, train = train_size)
36    y_true = np.ravel(y_true)
37    y_pred = np.ravel(y_pred)
38    df = pd.DataFrame({"y_true" : y_true, "y_pred" : y_pred})
39    df.to_csv(fname, index=False, sep=";")
40
41 def regression_results(y_true, y_pred):
42     """Function that calculates the different scores used to compare the models.
43
44     Arguments:
```

```

45
46     Returns:
47         dictionary -- dictionary with all the calculated scores
48         """
49         mean_absolute_error=metrics.mean_absolute_error(y_true, y_pred)
50         mse=metrics.mean_squared_error(y_true, y_pred)
51         r2=metrics.r2_score(y_true, y_pred)
52         metrics_dic = {'r2':round(r2,4), 'mae':round(mean_absolute_error,4), 'mse':round(mse,4), 'rmse':round(np.
53             sqrt(mse),4)}
54     return metrics_dic
55
56 def filter_data (data, wl, po):
57     """Function that applies Savitzky Golay filter to the data
58     Arguments:
59         data {dataframe} -- Dataframe with all the data
60         wl {int} -- Windows size to apply the filter
61         po {int} -- Polynomial order to apply the filter
62
63     Returns:
64         dataframe -- Dataframe with all the data filtered
65         """
66     dataset_filtered = data.copy()
67     dataset_filtered ['CO2']= savgol_filter ( dataset_filtered ['CO2'], wl, po)
68     dataset_filtered ['H']= savgol_filter ( dataset_filtered ['H'], wl, po)
69     dataset_filtered ['T']= savgol_filter ( dataset_filtered ['T'], wl, po)
70     return dataset_filtered
71
72 def remove_outliers ( dataset , threshold ):
73     """Function that detect the outliers from the dataset and removes them
74
75     Arguments:
76         dataset {dataframe} -- Dataframe with all the data
77         threshold {float} -- Threshold to consider a data point as outlier
78
79     Returns:
80         dataframe -- dataframe without outliers
81         """
82     z = np.abs( stats . zscore( dataset . iloc[:,1:4]) )
83     return dataset [(z < threshold) . all (axis=1)]

```

---

Code Listing 4: Python script with the common functions for all the scripts coded