

## Article

# Comparative Evaluation of Data-Driven Approaches to Develop an Engine Surrogate Model for NO<sub>x</sub> Engine-Out Emissions under Steady-State and Transient Conditions

Alessandro Brusa <sup>1,\*</sup>, Emanuele Giovannardi <sup>1</sup>, Massimo Barichello <sup>2</sup> and Nicolò Cavina <sup>1,\*</sup> <sup>1</sup> Industrial Engineering Department, University of Bologna, 40136 Bologna, Italy<sup>2</sup> Ferrari S.p.A., 41053 Maranello, Italy

\* Correspondence: alessandro.brusa6@unibo.it (A.B.); nicolo.cavina@unibo.it (N.C.)

**Abstract:** In this paper, a methodology based on data-driven models is developed to predict the NO<sub>x</sub> emissions of an internal combustion engine using, as inputs, a set of ECU channels representing the main engine actuations. Several regressors derived from the machine learning and deep learning algorithms are tested and compared in terms of prediction accuracy and computational efficiency to assess the most suitable for the aim of this work. Six Real Driving Emission (RDE) cycles performed at the roll bench were used for the model training, while another two RDE cycles and a steady-state map of NO<sub>x</sub> emissions were used to test the model under dynamic and stationary conditions, respectively. The models considered include Polynomial Regressor (PR), Support Vector Regressor (SVR), Random Forest Regressor (RF), Light Gradient Boosting Regressor (LightGBR) and Feed-Forward Neural Network (ANN). Ensemble methods such as Random Forest and LightGBR proved to have similar performances in terms of prediction accuracy, with LightGBR requiring a much lower training time. Afterwards, LightGBR predictions are compared with experimental NO<sub>x</sub> measurements in steady-state conditions and during two RDE cycles. Coefficient of determination (R<sup>2</sup>), normalized root mean squared error (nRMSE) and mean average percentage error (MAPE) are the main metrics used. The NO<sub>x</sub> emissions predicted by the LightGBR show good coherence with the experimental test set, both with the steady-state NO<sub>x</sub> map (R<sup>2</sup> = 0.91 and MAPE = 6.42%) and with the RDE cycles (R<sup>2</sup> = 0.95 and nRMSE = 0.04).

**Keywords:** data-driven models; machine learning; NO<sub>x</sub> emission; internal combustion engine; surrogate model



**Citation:** Brusa, A.; Giovannardi, E.; Barichello, M.; Cavina, N. Comparative Evaluation of Data-Driven Approaches to Develop an Engine Surrogate Model for NO<sub>x</sub> Engine-Out Emissions under Steady-State and Transient Conditions. *Energies* **2022**, *15*, 8088. <https://doi.org/10.3390/en15218088>

Academic Editor: Stefania Falfari

Received: 3 October 2022

Accepted: 24 October 2022

Published: 31 October 2022

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Copyright:** © 2022 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

These days, due to the increasing number of sensors installed on the engine and the high number of experimental tests, the amount of data available for the car manufacturers is always wider. This, along with the enhancement of computational power of common devices, is affecting the way with which data are managed, processed and analyzed [1].

To extract insights from such a high amount of data, Artificial Intelligence (AI) techniques are also spreading in the automotive fields, especially for applications related to autonomous driving, vehicle control, smart connections, virtual sensing and fault diagnosis [2].

The AI models are computationally cheap and capable of learning the main characteristics of systems based only on experimental measurements, without requiring explicit programming, thus reducing the effort required to model complex physical and chemical phenomena [3,4]. Moreover, the ability to process large amounts of information in a short time and to learn system behavior from experimental data makes AI models interesting for many applications. Therefore, the effort of research on development and implementation of such methodologies in the automotive field is strongly increasing.

In addition, in the last decades, emission regulations have become increasingly demanding, causing a growth in time and costs needed for engine calibration and development.

To reduce the number of physical experiments needed and to limit the cost increase in engine development, AI and machine learning can be exploited for modeling and predicting engine emissions in a virtual environment [5–9]. Differently from the physical and semi-physical models [10–15], a data-driven approach could be helpful, since the processes at the basis of emission formation, such as combustion and turbulence, are quite difficult to model analytically [16] and require much time to run in virtual environments. Despite 0-D models [17–21] being computationally efficient, the analytical formulation of the physical phenomena can be difficult to determine when many independent variables are affecting the output.

Some applications of machine learning aimed at emission modeling are already present in the literature [22–24]. For example, in their review, Shivansh Khurana et al. [25] show that most of the machine learning models proposed in the literature for emission predictions are based on Support Vector Machines (SVMs), ensembles of tree-based models (Random Forest or Gradient Boosted Trees) and Neural Networks (NNs).

The NN could seem the most reliable approach due to its high complexity, but the resulting accuracy strongly depends on the particular application. As an example, in their study, Altuğ and Küçük [26] made a comparison between Elastic-Net, eXtreme Gradient Boosting (XGBoost) and Long Short-Term Memory (LSTM) neural networks for NO<sub>x</sub> prediction, showing that XGBoost outperforms the most complex LSTM recurrent neural network.

Moreover, Papaïouannou [27] proposes a Random Forest algorithm to predict particulate emissions in a GDI engine. It is a simple model, easier to understand and less computationally expensive than deep neural networks.

However, there are many applications where artificial neural networks are used for predicting emissions with satisfying results, such as in [28,29].

Other approaches involve the use of advanced techniques for time series modeling derived from the deep learning (such as NARX, ARIMAX and RegARMA) [30] and heavy preprocessing techniques, such as in the work of Yu et al., where the Long Short-Term Memory neural network is applied to predict the NO<sub>x</sub> processed with Complete Ensemble Empirical Mode Decomposition with the adaptive noise (CEEMDAN) technique [31].

From the literature, there is not a best approach or algorithm uniquely adopted in emission modeling, and it is difficult to state it a priori, since each method has its own advantages and drawbacks.

Therefore, the present work presents a comparison of four different state-of-the-art techniques, namely the Support Vector Regressor (SVR) [32,33], the Random Forest (RF) [34], the Light Gradient Boosting (LightGBM) [35] and the Feed-forward Neural Network (FNN) [36], in order to estimate the NO<sub>x</sub> engine-out emissions. These models are compared with a Polynomial Regressor (PR), chosen as the benchmark, in terms of prediction accuracy and training time, in order to assess the best approach for this specific application. A brief description of these models is given below.

Most of the research in the literature refers to models for predicting the pollutant emissions under steady-state conditions. Nevertheless, one challenging aspect in this field is related to emission prediction under transient operating conditions [37]. Therefore, the present work wants to outline a procedure for the data preprocessing and analysis through machine learning techniques aimed at the development of a data-driven engine surrogate model capable of predicting the emissions not only under steady-state operations but also for dynamic and transient conditions.

Therefore, the models are trained and validated with data coming from Real Driving Emission (RDE) cycles that are well representative of a wide range of operating conditions, including highly dynamic maneuvers. Differently from standard homologation cycles, the RDE tests are performed on roads, and the emissions are measured by means of a PEMS (Portable Emission Measurement System). This means that the RDE cycles do not follow a

specific speed profile; instead, there are many variables such as the weather, the environmental temperature and humidity, as well as the altitude and the traffic conditions [38].

On the other hand, the conventional homologation cycles such as NEDC (New European Driving Cycle) or WLTC (Worldwide harmonized Light Test Cycle) are carried out in laboratory and they are set to follow a defined speed and pedal profiles. Therefore, these are not completely representative of a real driving condition, and they usually underestimate the pollutant emissions with respect to the real usage on the road [39].

For this work, the RDE cycles are performed reproducing a real road condition, and then the speed and load profiles are reproduced with a car on a roll bench. This is carried out mainly for two reasons. First, engine-out emission can be measured only in laboratory, since the PEMS is installed after the tailpipe; secondly, the measurement systems of the laboratory are much more accurate and reliable than the PEMS.

The inputs for the model are selected within the engine control unit (ECU) signals, and the output are obtained from the continuous measurement of the NO<sub>x</sub> emissions.

In Section 2, the experimental campaign, as well as the data preprocessing techniques, are described in detail. In this part, the methodology is applied to real industrial experimental data that require a wide preprocessing activity. First of all, the emission measurements present a delay with respect to the ECU channels that is compensated through an alignment over time, based on the first engine firing. Then, an optimal set of features is defined by combining different feature selection techniques, such as the correlation analysis, the Feature Importance Permutation (FIP) and domain knowledge. Moreover, a novel Sliding Window over time is applied to the input matrix to keep into account the partial history of the inputs and to enhance the performance of the model under transient conditions. In other words, the innovative contribution of the proposed activity consists of the implementation of the preprocessing methodologies and the machine learning algorithms to estimate the NO<sub>x</sub> emissions produced during real driving maneuvers.

In Section 3, the comparison between different data-driven models is presented. Each model is calibrated and tuned through the Randomized Search Cross-Validation, and a summary of their performance is provided. Moreover, a sensibility analysis to the Sliding Window size is reported, showing the impact that it has on the training time and the model accuracy. The LightGBM is selected for its low training time and high accuracy, and it is used to predict the NO<sub>x</sub> emissions of two RDE cycles and under steady-state operating conditions. The model trained on the RDE cycles is also validated calculating the NO<sub>x</sub> emissions for steady-state conditions. In this way, it is possible to highlight the accuracy of this methodology also when it is tested on different operating conditions.

Section 4 summarizes the main conclusions with a particular focus on the future developments of the proposed work.

## 2. Methodology

In this section, the methodologies implemented for the data preprocessing, the feature selection and engineering and the comparison of the models performance are presented.

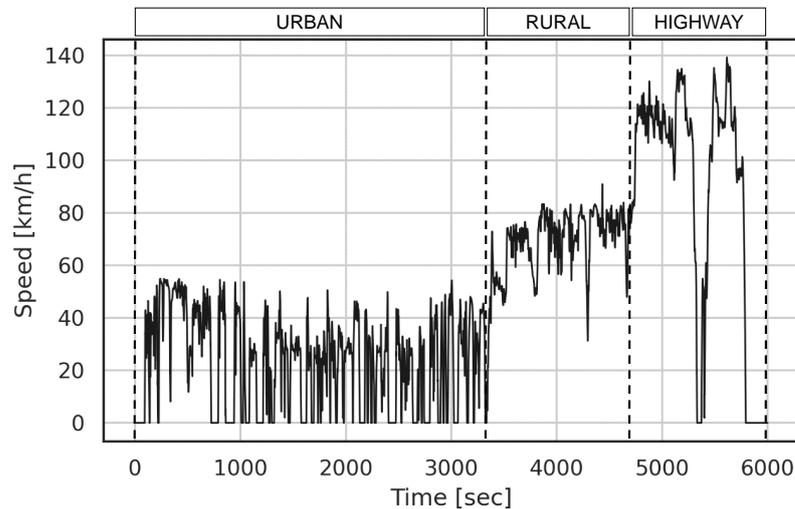
### 2.1. Experimental Dataset

The models are trained and tested using data from eight different experimental RDE cycles and from the engine steady-state emission measurements. The complete dataset is composed of 9.5 h of ECU and emissions recordings for a total of almost 350,000 time samples. More details about the dataset composition are provided in the following sections.

All the RDE cycles are reproduced on a roll bench, where a virtual driver is set to follow the target speed and pedal profile of a real on-road RDE cycles. However, it is possible to reproduce such maneuvers in a controlled laboratory environment, taking advantage of the more accurate and robust tools for the engine parameters and emissions measurements. This choice is also made to get as close as possible to the measurements achievable with the most accurate sensing tools.

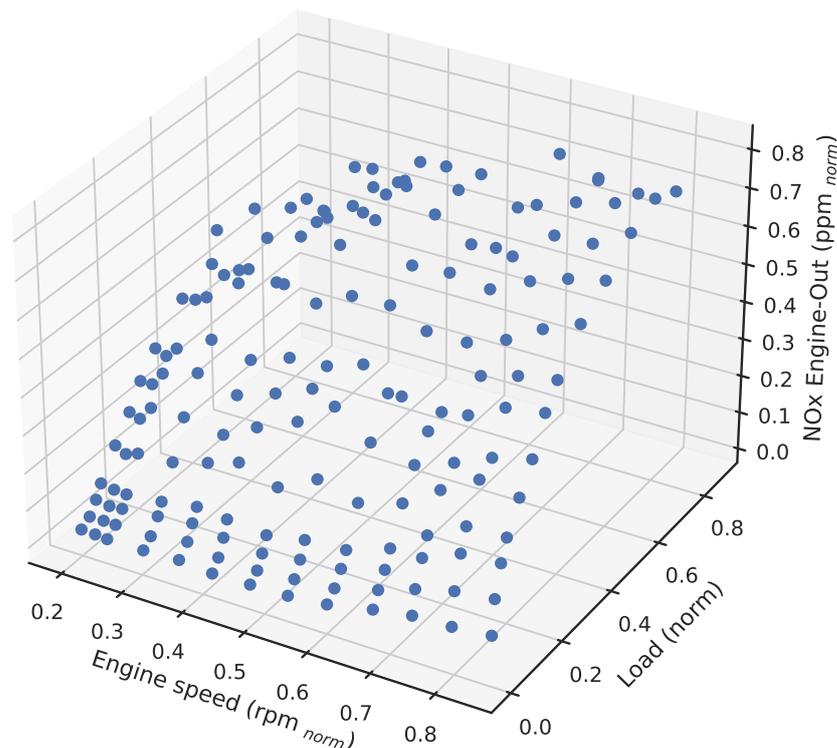
An example of a typical speed profile for an RDE cycle is shown in Figure 1. Three different sections of the cycle are noticeable:

- Urban: Vehicle speed below 60 km/h;
- Rural: Vehicle speed comprises between 60 km/h and 90 km/h;
- Highway: Vehicle speed up to 140 km/h.



**Figure 1.** An example of a Real Driving Emission cycle divided into three main phases: urban, rural and highway paths.

As a reference for steady-state operating conditions, the dataset also includes an experimental map of engine-out NO<sub>x</sub> measurements obtained by testing the engine in many stationary engine points. A 3-D scatter of the NO<sub>x</sub> map normalized along all the axes is reported in Figure 2.



**Figure 2.** NO<sub>x</sub> engine-out concentration under steady-state conditions depending on engine speed and load.

During all the experimental tests, both the NO<sub>x</sub> emissions and the ECU channels are recorded, corresponding to the output and the inputs of the models, respectively. Therefore, a supervised learning approach can be adopted.

The data reported here come from a real industrial test, and the reported data cannot be disclosed. Therefore, for the sake of confidentiality, the whole dataset is normalized in a range between zero and one according to (1):

$$x_{norm} = \frac{x - x_{min}}{x_{max} - x_{min}} \quad (1)$$

For the same reason, the results and plots shown in the following text are normalized with the same equation.

## 2.2. Experimental Setup

The experimental tests were conducted on a laboratory roll bench on a vehicle equipped with a state-of-the-art spark-ignited V12 naturally aspirated engine, whose specifications are reported in Table 1.

The NO<sub>x</sub> emissions are measured by means of a Chemiluminescent Detector analyzer (CLD). This device measures the NO<sub>x</sub> concentration in the exhaust gases, exploiting the fact that the nitric oxide (NO) combines with the ozone (O<sub>3</sub>) to create electronically excited NO<sub>2</sub> molecules, which, returning to the equilibrium state, emit visible radiations with intensity proportional to the concentration of NO in the gas. Being able to measure the irradiated light, the CLD can assess NO<sub>x</sub> concentration in the exhausts. The key features of the CLD analyzer are reported in Table 2.

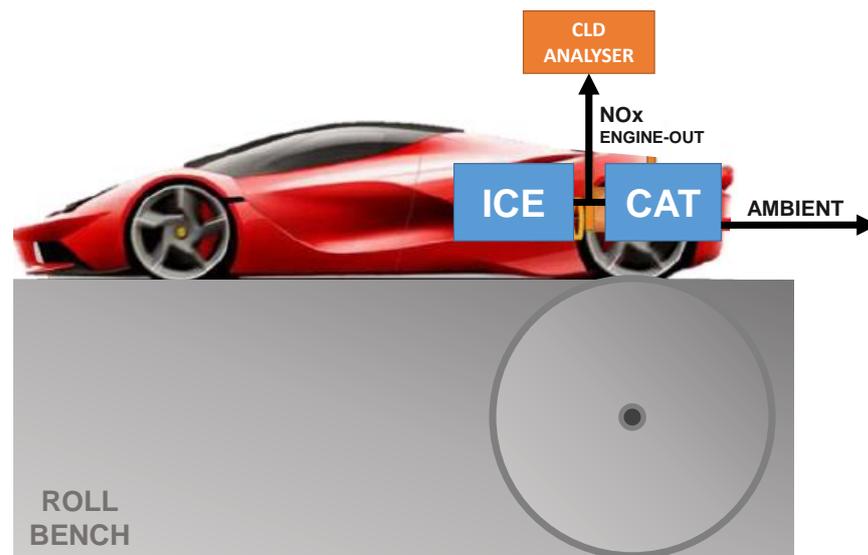
**Table 1.** Engine Specifications.

Engine Specifications	
Engine Type	V12
Displacement Volume [cc]	6495.6 cc
Intake Type	Naturally Aspirated
Combustion System	GDI Spark-ignition
Number of Cylinders [#]	12 (6 per bank)
Valves per Cylinder [#]	4 (2 int + 2 exh)
Bore × Stroke [mm]	94.0 × 78.0

**Table 2.** Chemiluminescent Detector Sensor (CLD) Specifications.

CLD Sensor Specifications	
Acquisition Frequency [Hz]	10
NO <sub>x</sub> Range [ppm]	min 10–max 10,000
Response Dynamics	90% full-scale in 1 s
Temperature Effect	<2% of full-scale per 10K of T

To figure out the structure of the laboratory roll bench equipment, a schematic layout is reported in Figure 3. The CLD analyzer is located between the engine and the after-treatment system of the vehicle to intercept the NO<sub>x</sub> concentration in the exhausts coming out from the engine.



**Figure 3.** Schematic layout of the roll bench with CLD analyzer.

### 2.3. Data Pre-Processing

Starting from the experimental measurements, a raw dataset is firstly generated by including the ECU signals and parameters selected according to the criteria described below, and these are coupled with the NO<sub>x</sub> engine-out values sensed by the CLD analyzer.

Typically, the ECU channels have different sampling frequency depending on their nature and on the quantity they outline. For example, their acquisition frequency can be 10 Hz or 100 Hz, or in some cases, they can be acquired with a frequency proportional to the engine speed. Therefore, to obtain a homogeneous dataset, a resampling of each channel is needed.

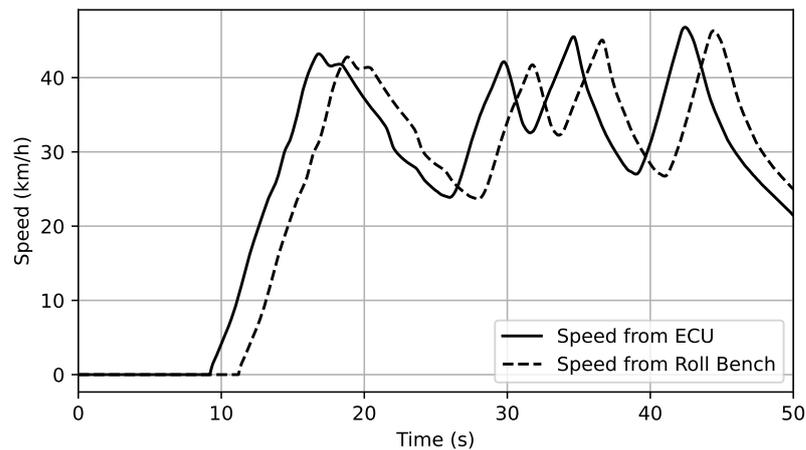
To this end, all the ECU channels are resampled at 10 Hz, which is the same acquisition frequency of the CLD sensor used for measuring NO<sub>x</sub> concentration in the exhausts.

A second issue is related to the presence of two different acquisition systems, namely the ECU, installed on the vehicle, and the CLD sensor, installed on the roll bench. So, the synchronization in time between them is needed in the preprocessing phase. The total delay between them can be split into two main components:

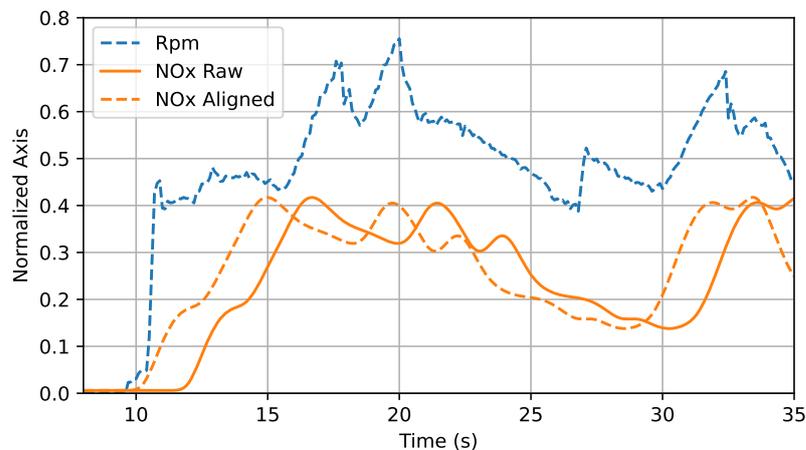
- Different timestamps of the ECU with respect to the CLD analyzer because they are independent measurement tools;
- Delay in the emission measurement due to: (i) the system dynamics, because of the time needed for the exhaust gases to achieve the CLD sensor, and (ii) the sensor dynamics, since it acts as a first-order system that makes the transient measurement smoother.

For the aforementioned reasons, it is understandable how the emission delay has an impact on data synchronization and on data quality. Apart from the systematic delay between the CLD sensor and the ECU, there are components of the delay that depend on the dynamics of the phenomena that are taking place. Indeed, this delay is related to the emission pick-up point position on the exhaust pipe and to the distance to the CLD sensor. Moreover, since the exhaust gas mass flow strongly depends on the engine operating conditions, the time needed to move across the ducts is variable during the test. Nevertheless, as a first approximation, the delay is considered constant, supposing that a rigid shift of the NO<sub>x</sub> signal is sufficient. The delay due to different timestamps between the roll bench and the ECU channels is compensated by means of a reference signal that is acquired both from the roll bench and from the ECU, namely the vehicle speed. Knowing the vehicle speed calculated from the rolls' angular speed and aligning it with the speed coming from the ECU, it is possible to compensate the mentioned delay, as shown in Figure 4.

On the other hand, the delay due to the dynamics cannot be easily compensated since it depends on many variables; firstly, the speed of the exhaust gas flow. Therefore, as a simple and robust solution, the delay was compensated by aligning the first positive gradient of the emissions (NO<sub>x</sub>) with the engine start, as reported in Figure 5. Here, the hypothesis is that the first emission peak occurs immediately after the engine starts and that, as an acceptable approximation, the variable part of the delay, related to the changing exhaust mass flow, can be neglected.



**Figure 4.** Vehicle speed reference to compensate the delay between the roll bench and the ECU.



**Figure 5.** Synchronization of NO<sub>x</sub> signal with engine firing, using engine rotational speed as reference.

Then, a table is generated with the ECU channels synchronized with the NO<sub>x</sub> emission traces. Different physical quantities are in different columns of the dataset, while each row represents a timestep. Moreover, having many RDE cycles, it was possible to concatenate them inside a unique tabular dataset.

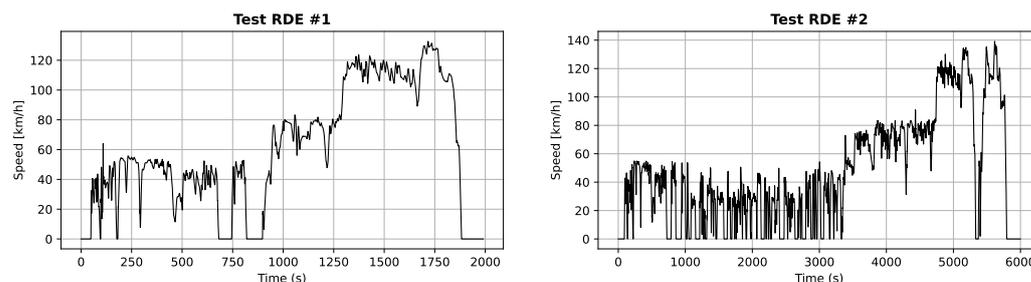
As already mentioned above, for confidentiality reasons, every column of such database is normalized with respect to the maximum value, following the Equation (1).

To develop a data-driven model, the dataset needs to be split into training, validation and testing sets. The training set is the part of the data from which the model learns relations between the input and output. The validation set is a part of the data which is held out from the training process to assess whether the model is incurring overfitting or to perform hyperparameter tuning. Finally, the testing set is the part of the data used after model training and hyperparameter tuning to evaluate prediction capability. During testing, only inputs are provided to the model. Then, by comparing the predicted output with the real one, it is possible to give a score to the model according to different possible metrics.

For this activity, a 10-fold cross-validation was performed instead of a typical validation process, as described in detail later.

The available experimental dataset consists in a set of eight different RDE cycles and the shortest has a duration of 1800 s while the longest of 6000 s.

Two RDE cycles were held out from the dataset for the model testing, reported in Figure 6. These two RDE cycles have different durations, with the first one lasting about 2000 s and the second one lasting 6000 s. Since these experiments are used exclusively for model testing, they are excluded from the training process.



(a) First RDE cycle with duration around 2000 s. (b) Second RDE cycle with duration of 6000 s.

**Figure 6.** Speed profile of two Real Driving Emission cycles within the test set.

Within the remaining six experiments, 90% of the data were used for training and the remaining 10% for the cross-validation of the models. To avoid possible biases due to the time order of the experiments, a random shuffling of the samples was applied.

Table 3 summarizes the split of these sets, specifying the time and the portion of the total dataset.

**Table 3.** Dataset split.

Dataset	Composition	Duration [h]	Proportion [%]
Training	6 RDE cycles	6.25	65
Validation	10% of training set	1	10
Test	2 RDE + Steady-state map	2.25	25
Total	8 RDE + Steady-state map	9.5	100

## 2.4. Model Development

### 2.4.1. Feature Selection and Processing

To develop a suitable dataset, the input features are chosen by means of different techniques of features selection, such as the correlation analysis and the features importance permutation combined with the physical domain knowledge. Afterwards, the selected features are processed, considering that they are describing dynamical phenomena. So, in this work, the introduction of a temporal Sliding Window is assessed, and its impact on the model performances is evaluated. Details and results about the feature selection and feature engineering processes are provided in Section 3.

### 2.4.2. Model Selection and Testing

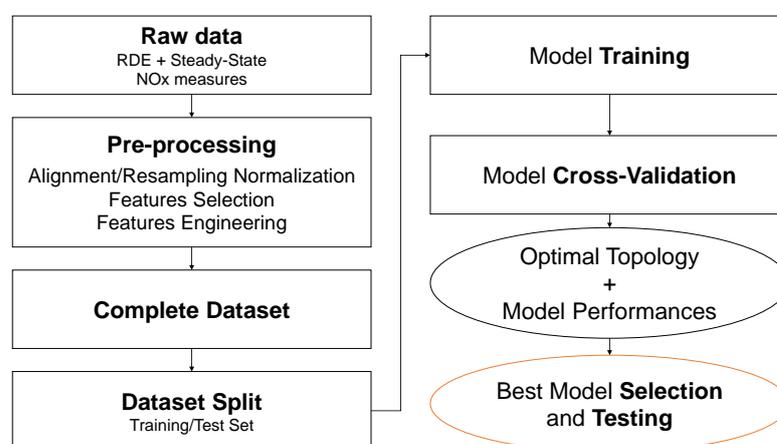
Five data-driven regressors (PR, SVR, RF, LightGBR and FNN) are trained and validated on a set of six RDE cycles. The optimal model is selected as the best compromise between the training time and accuracy, assessed by means of a 10-fold cross-validation.

Then, the best model is applied to the test set composed of two RDE cycles and one steady-state map. With this approach, the capabilities in the NO<sub>x</sub> prediction are assessed both in highly dynamic and steady-state conditions. All the results are discussed in Section 3. It is emphasized that the results are obtained by adopting a single experimental setup, meaning that the vehicle type, the engine type and the emission measurement system

are the same during all the experimental tests. Nevertheless, it is expected that a change in the hardware components that have an impact on the combustion and on the emission production may require a recalibration of the models. In particular, three different cases can be distinguished:

- Change in the ECU control software: The models do not require further actions if training data are still representative of the test data.
- Change in hardware components that does not affect the engine's functional layout: If there is an impact on the engine combustion and performances, the model needs to be retrained on an updated experimental dataset.
- Change in the engine's functional layout: The model needs to be developed from scratch; the feature selection and the model topology should be assessed.

In all the aforementioned cases, the general methodology proposed in this paper for preprocessing, feature selection, feature engineering and model selection is still valid. The algorithm of such methodology is reported in Figure 7.



**Figure 7.** Flow chart of the main methodology steps.

### 3. Results

#### 3.1. Features Selection

The features of the model are selected from the available ECU channels because these signals are also available on-board, and thus can be used for future real-time implementations.

However, a lot of ECU channels are recorded during an experiment, and most of them represent physical quantities that do not affect NOx formation phenomena.

Therefore, the selection of the relevant ECU channels is a critical step of the features selection process, since, on one hand, it is crucial to consider all the relevant inputs that affect the NOx emission, but, on the other hand, it is important having as little redundancy as possible. In some cases, even if some of the inputs are highly correlated with the NOx concentration, they are not providing any additional information to the model. Instead, they are concurring to increase the model complexity and the computational effort.

Moreover, the feature redundancy increases data noisiness, as well as the probability for the model to find inconsistent input–output relationships.

The main steps of the feature selection are summarized in Figure 8. First of all, the ECU channels that were either faulty, incomplete, or not related at all with NOx formation are removed from the input set. As a second selection, the input-to-input Spearman correlation analysis [40] is calculated for each remaining ECU channel to highlight strong correlations between them. Then, a similar correlation analysis is applied to underline existing relationships between the inputs and the output.

However, such approach can lead to incomplete results, since the typical correlation indexes (i.e., the Pearson or Spearman ones) can highlight only the linear or monotonic

relationships that can be applied to a single input, neglecting any possible interaction between multiple inputs.

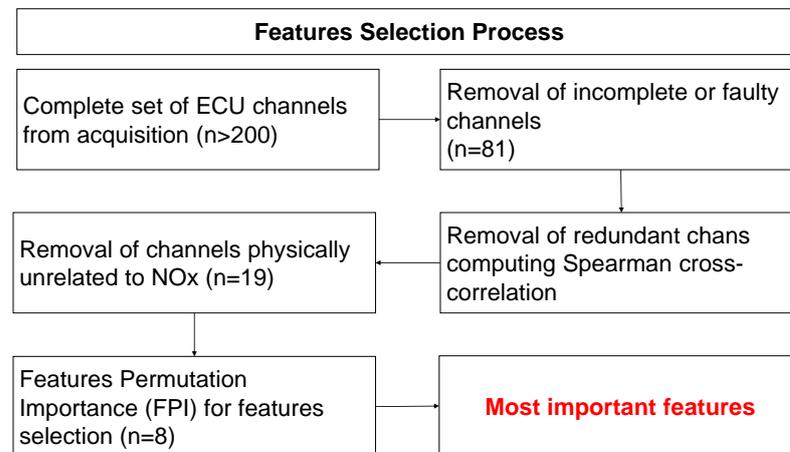


Figure 8. Flow chart of main feature selection steps.

Hence, there are several techniques which take advantage of the model itself to define the key features. One of them is the Feature Importance Permutation (FIP) technique [41].

This method evaluates the drop of model accuracy when a feature is removed from the inputs. The process is iterated removing one features at a time and repeated 30 times to increase the statistical relevance of the results. The features that lead to the highest performance drops can be considered the most important for the model. Since this technique is a greedy process, the domain knowledge can further help to isolate only the features physically related to the outputs, such as the engine actuators, sensors and strategies that are effectively involved in NOx production.

Thereby, only the features considered important from a physical point of view and not redundant are considered. A plot of the results using FIP with a base Random Forest model is shown in Figure 9.

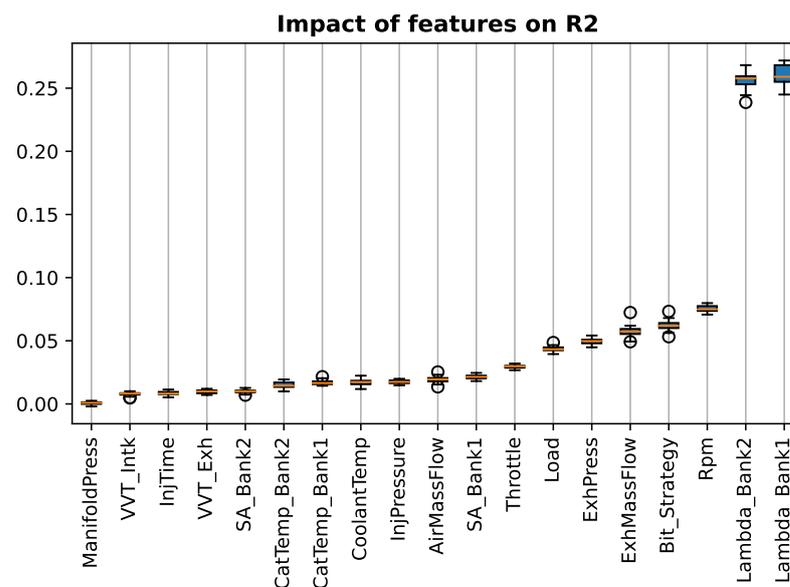


Figure 9. Results of Feature Importance Permutation applied on a subset of ECU channels.

It is clear that the air-to-fuel ratio measured by the lambda probes in both engine banks is the most relevant feature for the model. This finds confirmation from physics, since the NO<sub>x</sub> formation is favored by the lean mixtures.

A further note is needed because some of the inputs are measured independently for the two engine banks. Normally, the signals of each engine bank are aligned with the same variable of the other bank. However, some quantities of the particular engine bank are selected, such as, for instance, the lambda measured in the exhaust line in order to achieve the highest correlation between the inputs and the output of the model.

In the end, the FIP results are used as general guidelines for the selection of the optimal set of features. However, the resulting selection of inputs for the proposed models were conducted by evaluating their impacts on the physical process that affects the emission production (i.e. the combustion). In other words, even if the ranking reported in Figure 9 shows a high importance of the flags associated, for instance, to the activation of the component protection strategy (named Bit-Strategy), such Boolean values are not included in the final set of features because of their poor physical contribution to pollutant emission production. In this way, the training process leads to calibrate a robust model, since the selected features have a physical relationship with the estimated outputs. The complete set of inputs is reported below:

- Lambda (both engine banks);
- Engine speed;
- Engine load;
- Exhaust mass flow;
- Intake valve opening;
- Exhaust valve closing;
- Spark advance.

### 3.2. Feature Processing

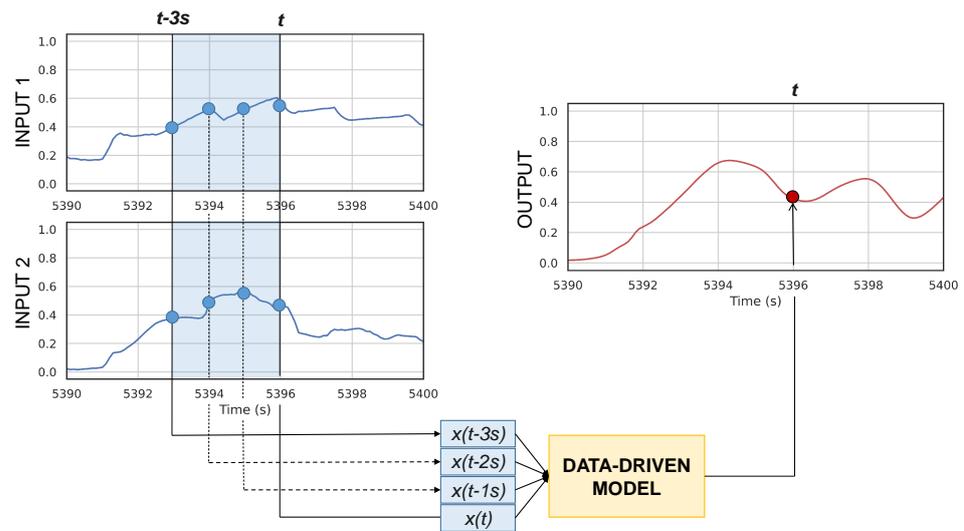
Differently from the typical regression tasks, in this case, the inputs and the output are time series measurements recorded during the RDE driving cycles, and thus under highly dynamic operating conditions. This adds a complication to the task, since the output is also affected by the history of the inputs.

Typical machine learning regressors are not able to keep into account such kind of dynamic behavior, since each sample is considered temporally independent from the others.

To overcome this limitation, a Sliding Window over time is applied to the inputs. With this technique, the inputs are not considered on a single time instant, instead, the partial history of each signal is provided to the model.

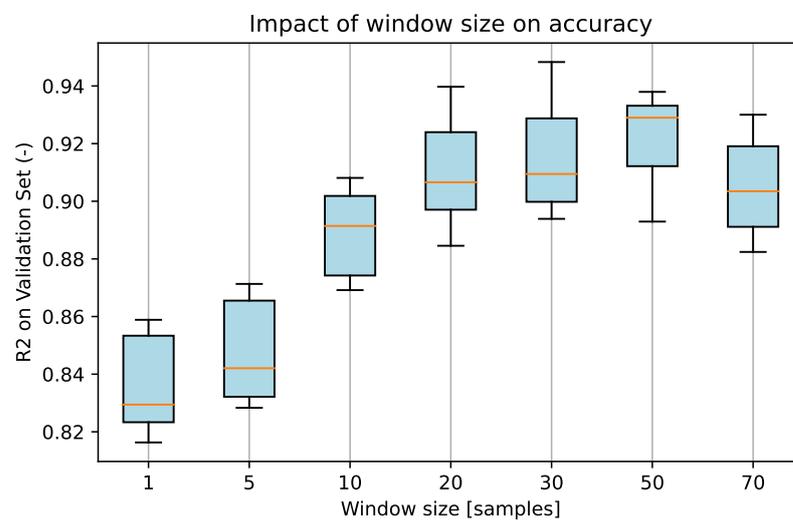
The scheme in Figure 10 is a graphical explanation of the Sliding Window working principle in the case of a 30-sample window size and only two inputs. To predict the output at time  $t$ , each timestep of the inputs within the window from  $t - 3s$  to  $t$  is used to build the input matrix, which is given to the model to infer the output at time  $t$ . Then, to predict the next output sample, the window is shifted one sample forward and the process is repeated.

The window size affects the model performances; thus, a sensitivity analysis is performed to quantify the effect of this hyperparameter on the model output. To this end, the performance of the LightGBR model is tested when Sliding Windows of different sizes are applied. Namely 1-sample (no window), 5-sample, 10-sample, 20-sample, 30-sample, 50-sample and 70-sample windows are the tested configurations.



**Figure 10.** Sliding Window techniques, considering 3 s-sized window.

Figure 11 summarizes the results of this analysis, plotting on the x-axis the window size and on the y-axis the prediction accuracy, represented here by the correlation coefficient between the real and the predicted output. Generally, an increase in the window size leads to an increase in the prediction accuracy until a 50-sample window. This appears to be the optimal size, because for further window enlargements, the prediction accuracy decreases. However, performances are also assessed in terms of training time, and, since, the increase in window size corresponds to an increase in computational time, a 30-sample Sliding Window is chosen as an optimal trade-off between the prediction accuracy and the computational efficiency.



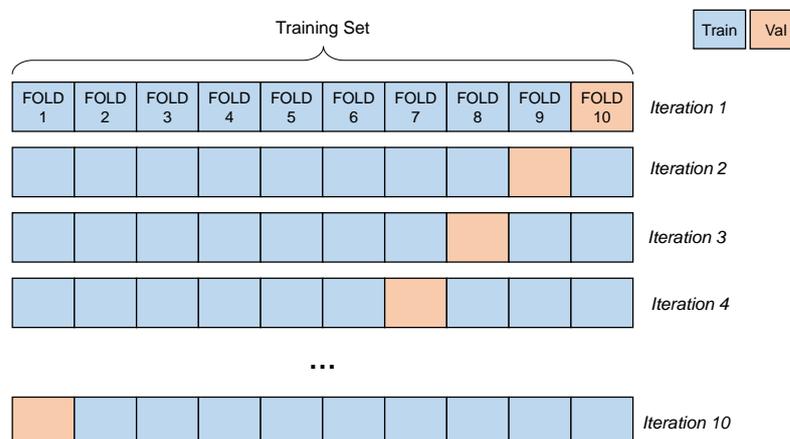
**Figure 11.** Effect of the Sliding Window size on R2 index.

### 3.3. Model Selection

From the literature, even once the modeling problem has been defined, it is not possible to detect a priori an optimal data-driven model that can outperform in the emission prediction task. One of the primary objectives of this study is to investigate the most suitable learning algorithms for NOx emission prediction.

Each learning algorithm has its own advantages and drawbacks and, since it is not possible to define the best model a priori, a specific procedure was applied to a set of machine learning regressors, namely Polynomial Regressor, Support Vector Regressor, Random Forest, Light Gradient Boosting and Feed-forward Neural Network.

The procedure consists in a 10-fold cross-validation of the training set. The training set was split into 10 folds, nine of which are used to train the model, whereas the remaining fold is used to validate the model and to understand the prediction capability of each model on new unseen data. This process is iterated, changing the validation fold each time, until every fold is used, as shown in Figure 12. Moreover, to increase the robustness of the results, this procedure is repeated five times for each model.



**Figure 12.** Iterations of the 10-fold cross-validation procedure. Blue folds are used for training, whereas orange folds are used for validating the model.

Each iteration supplies information about the time needed to train the model and about the model accuracy, assessed on each validation fold.

Moreover, a further step in the model development is the optimal hyperparameters definitions. These have an important effect on the model performance; thus, it is fundamental to choose them properly.

To this purpose, the Grid Search method can be used in combination with the cross-validation procedure. With the Grid Search cross-validation, a grid of possible hyperparameters is defined manually, following general guidelines, then the algorithm itself tests every possible combination of hyperparameters and finds the best one. As easily understandable, this process is time-consuming; therefore, a lighter variant of the Grid Search, namely the Randomized Search [42,43], is used for the same purpose in this work. With this technique, the total number of iterations can be imposed to reduce the computational time needed.

The Randomized Search cross-validation is repeated for each type of regressor with different Sliding Window sizes. This is a complete method to provide a robust summary about model performances, keeping into considerations the many degrees of freedom available during the model definition, namely:

- Type of learning algorithm;
- Set of hyperparameters for each model;
- Size of the Sliding Window.

The resulting output allows to assess the best model according to the following:

- Model accuracy assessed with different metrics (R2 and nRMSE);
- Average training time over several repetitions.

Results obtained using three different Sliding Window sizes are reported in Table 4 (no window), Table 5 (10-sample window) and Table 6 (30-sample window), considering only the best set of hyperparameters for each model. The performances are averaged on the 10 validation folds for all the repetitions and are expressed by means of coefficient for determination (R2), normalized root mean squared error (nRMSE) and training time.

The coefficient of determination defines the proportion of the output data variation that can be explained by the model. On the other hand, the RMSE reports the deviation between the model prediction and the target value. In this case, since the data are normalized,

the resulting RMSE is a non-dimensional index. For this reason, it is referred to as nRMSE. Finally, the training time represents the computational time needed to train the model, expressed in seconds.

**Table 4.** Models Performance without Sliding Window.

Model (No Window)	nRMSE [-]	R2 [-]	Training Time [s]
Polynomial Regressor (PR)	0.080	0.815	0.57
Support Vector Regressor (SVR)	0.076	0.825	22.3
Random Forest Regressor (RF)	0.048	0.874	42.1
Light Gradient Boosting (LGBR)	0.052	0.867	1.52
Neural Network (FNN)	0.055	0.859	64.7

**Table 5.** Models' performance using 10-sample Sliding Window (WS = 1 s).

Model (WS = 1 s)	nRMSE [-]	R2 [-]	Training Time [s]
Polynomial Regressor (PR)	0.072	0.819	27.9
Support Vector Regressor (SVR)	0.061	0.846	32.0
Random Forest Regressor (RF)	0.044	0.890	210
Light Gradient Boosting (LGBR)	0.048	0.888	3.91
Neural Network (FNN)	0.052	0.873	84.3

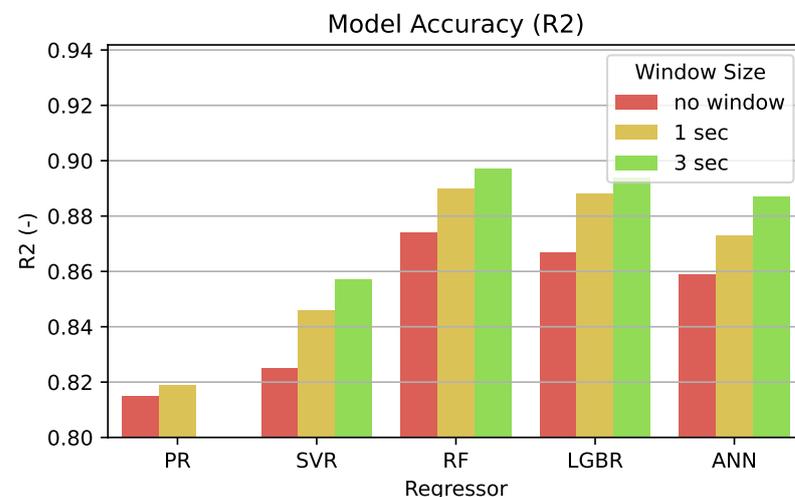
**Table 6.** Models' performance using 30-sample Sliding Window (WS = 3 s).

Model (WS = 3 s)	nRMSE [-]	R2 [-]	Training Time [s]
Polynomial Regressor (PR)	-	-	-
Support Vector Regressor (SVR)	0.059	0.857	72
Random Forest Regressor (RF)	0.042	0.897	657
Light Gradient Boosting (LGBR)	0.043	0.894	11.1
Neural Network (FNN)	0.049	0.887	130

### 3.3.1. Model Accuracy

As shown in Figure 13, regardless of the model considered, the accuracy increases along with the size of the Sliding Window. Here, we reported the results, respectively, for the following model layouts:

- 1-sample window (no window applied);
- 10-sample window (1 s time interval);
- 30-sample window (3 s time interval).



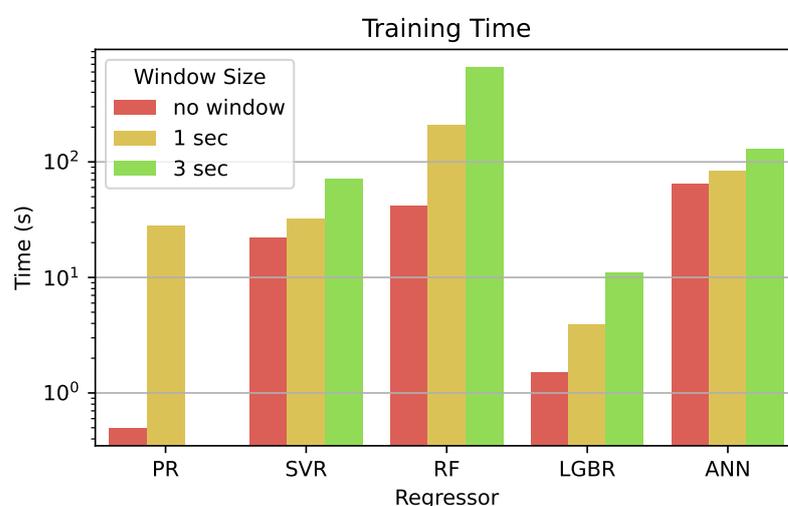
**Figure 13.** Comparison of the accuracy (represented by R2) of different data-driven models depending on Sliding Window size.

The correlation coefficient ( $R^2$ ) and the normalized RMSE (nRMSE) show consistent results, indicating the Random Forest and the LightGBR as the best models in terms of accuracy, followed by the Neural Network. On the other hand, the Support Vector Regressor is only slightly better than the Polynomial Regressor, which represents the baseline.

Within the three best models, in order: the RF, the LightGBR and the FNN, an interesting aspect to highlight is that the LightGBR seems to improve its accuracy more when window size increases with respect to the other two models. For example, considering nRMSE, the relative difference between RF and LightGBR is 8% (in favor of RF) when no window is applied, but it reduces to only 2% when a 30-sample window is used.

### 3.3.2. Computational Time

Figure 14 shows the average time needed for training each model in the log-scale. As expected, the training time shows an increasing trend with the window size adopted.



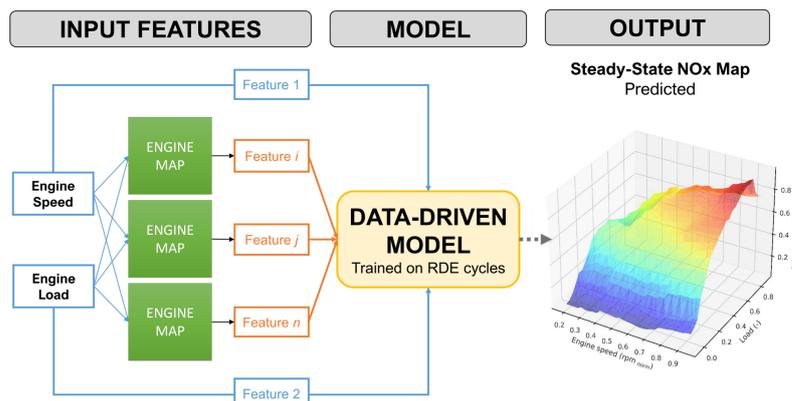
**Figure 14.** Comparison of the training time of different data-driven models depending on Sliding Window size.

The slowest model is the Random Forest, which shows a dramatic increase in time when the window is enlarged. Probably, this is because the hyperparameters are optimized to maximize the model accuracy instead of the training time, preferring a forest with more trees and penalizing the computational efficiency of the model. The FNN is computationally heavy as well, but its training time is less affected by the window size. On the other hand, the LightGBR is considerably faster than all the other models, even when wide windows are used. When using a 30-sample Sliding Window, the training time is around 10 s, and it is about two orders of magnitude lower than the Random Forest, which needs 657 s to train. This is a further experimental verification of how efficient the LightGBR algorithm can be when dealing with large datasets. Therefore, the LightGBR is the optimal model, since it proved to be much faster than any other tested model, with a minimum lack in accuracy with respect to the Random Forest.

### 3.4. Model Testing

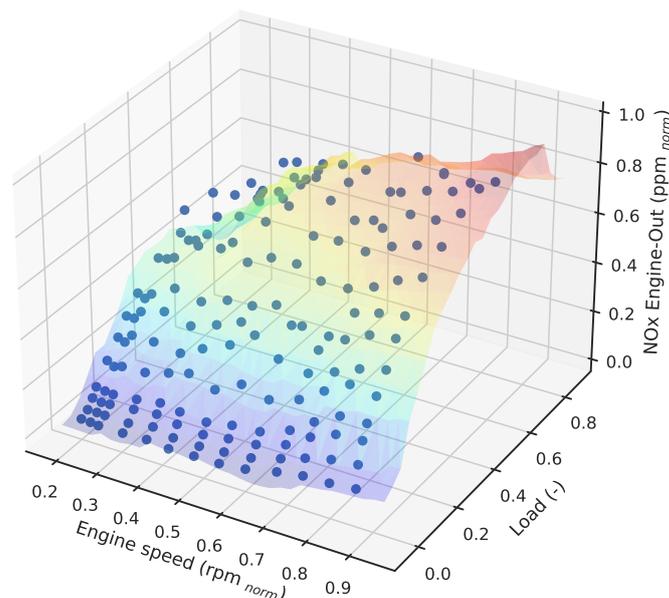
In this section, the LightGBR model with a 30-sample backward Sliding Window is considered for the final performance evaluation on the test set, which was considered yet. To assess the capabilities of the selected data-driven model, two different validation tests are performed. Firstly, the trained model is used to produce a NO<sub>x</sub> map under steady-state conditions. The virtually generated map is compared with an experimental map obtained under specific engine steady-state points, defined by different combinations of constant engine speed and load. Then, the accuracy of the NO<sub>x</sub> estimation on the two RDE cycles excluded from the training dataset is evaluated as second validation test.

To produce a map depending only on engine speed and load, a set of breakpoints is firstly defined. Under steady-state conditions the main control parameters such as the spark advance, the valve phasing and the air–fuel ratio are defined, depending on the engine speed and load, by means of base maps. Therefore, by exploiting these maps, all the model features can be defined depending only on the engine speed and load breakpoints. In this case, the engine speed and load are the independent variables from which the other features are obtained. This is because the majority of the engine actuations and target values are calibrated as a function of the engine speed and load. As a consequence, the NO<sub>x</sub> predicted by the model can also be defined depending only on the engine speed and load, represented in a three-dimensional map, as shown in Figure 15.



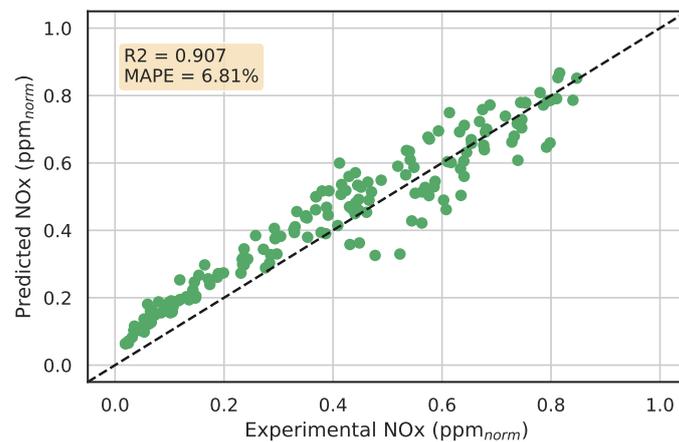
**Figure 15.** Steady-state NO<sub>x</sub> map prediction. Engine speed and load (blue) are the only independent features, whereas the other features are obtained from base maps. Thereby, the final prediction of the model can be represented as a 3-D map depending only on the engine speed and load.

In Figure 16, the experimental scatter plot of steady-state NO<sub>x</sub> measurements is superimposed to the map surface plot obtained from the LightGBM. The model prediction shows a good fit with the experimental scatter points, excluding the region of the low engine speed and high engine load, where the model underestimates the real NO<sub>x</sub> values. However, the general trend is well-represented, showing NO<sub>x</sub> increase corresponding to load and speed increase.



**Figure 16.** Three-dimensional comparison between experimental and modeled NO<sub>x</sub> emissions under steady-state conditions. Experimental values are represented by the scatter points; model prediction is represented by the surface.

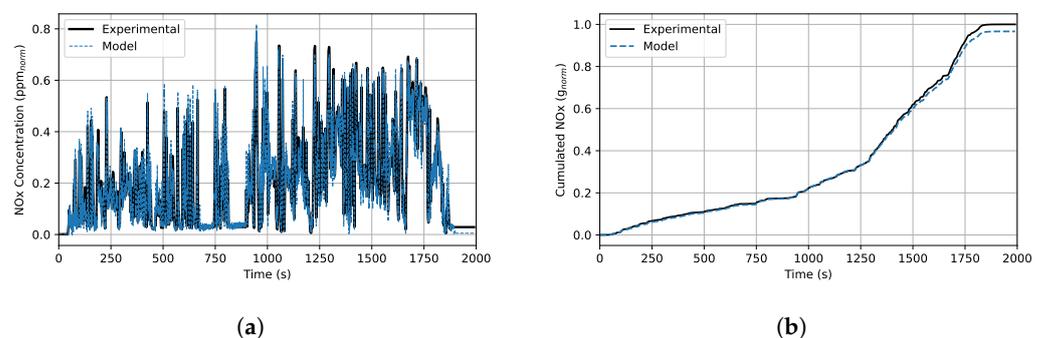
A two-dimensional scatter plot is also shown in Figure 17 to highlight the correlation between experimental and predicted NOx. Apart from the aforementioned region of high load and low speed conditions, the LightGBR slightly overestimates NOx. This can be due to training on highly dynamic conditions. Indeed, the model is trained only on a set of RDE cycles, where steady-state conditions are almost never met. Nevertheless, the general accuracy of the prediction is good, showing an R2 coefficient higher than 0.9 and a MAPE lower than 7%.



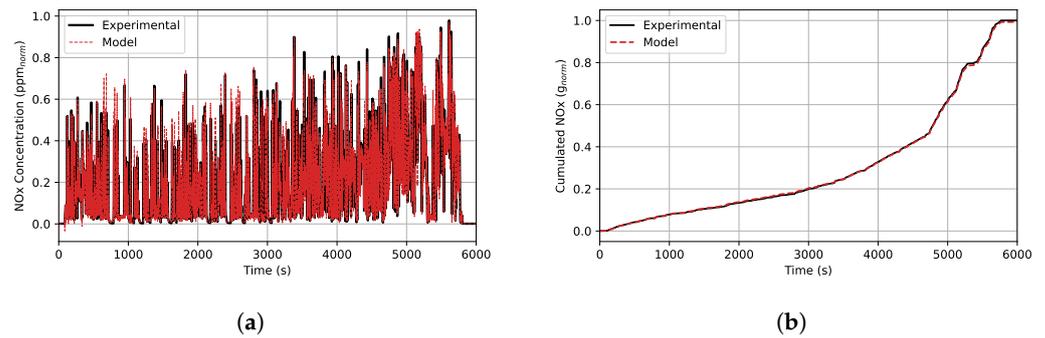
**Figure 17.** Correlation between experimental and predicted NOx under engine steady-state operations.

The second test consists in applying the LightGBR to predict NOx during two RDE cycles, whose speed profiles are reported in Figure 6. This part of the test is helpful for determining prediction accuracy under dynamic conditions. The continuous ECU channels are provided to the model as inputs in the form of time vectors to predict the output. Then, comparison between the real and the predicted outputs is performed to validate the model. Figures 18 and 19 show the results on the first and second RDE cycle of the test set, respectively.

The plots in Figures 18a and 19a show the normalized NOx concentration measured by the CLD sensor during the experiment with a black continuous line and the corresponding model prediction with a colored dashed line. In both cases, the model accurately predicts the trend of the NOx concentration. Indeed, the modelled and the experimental traces are well superposed in the graphs.



**Figure 18.** Comparison between experimental and modeled NOx emissions on the first test RDE cycle. (a) Instantaneous NOx concentration. (b) Cumulated NOx mass. Normalization makes the final experimental value equal to 1.



**Figure 19.** Comparison between experimental and modeled NO<sub>x</sub> emissions on the second test RDE cycle. (a) Instantaneous NO<sub>x</sub> concentration. (b) Cumulated NO<sub>x</sub> mass. Normalization makes the final experimental value equal to 1.

The plots in Figures 18b and 19b show the cumulated mass of NO<sub>x</sub> obtained from NO<sub>x</sub> concentration according to (2) and (3).

$$m_{NO_x} = \int_{t_1}^{t_2} \dot{m}_{NO_x}(t) dt \quad (2)$$

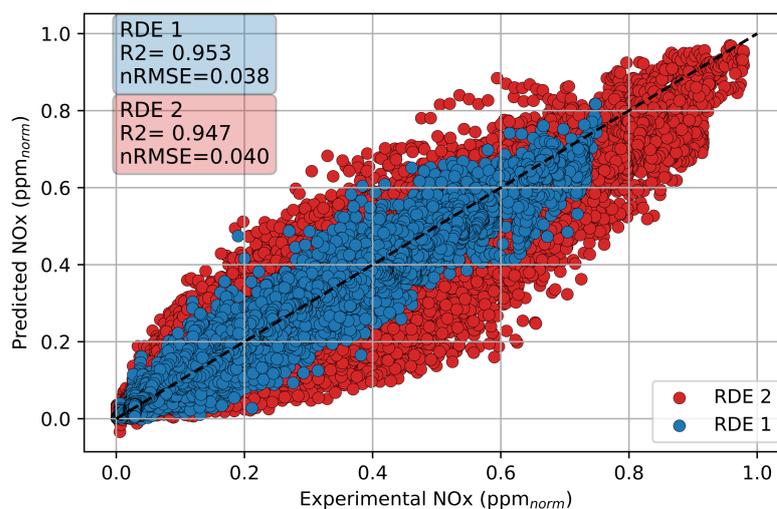
where  $m_{NO_x}$  is the mass of NO<sub>x</sub>,  $t_1$  and  $t_2$  are the time instants when the RDE test starts and finishes, respectively, and  $\dot{m}_{NO_x}$  is the NO<sub>x</sub> mass flow calculated as:

$$\dot{m}_{NO_x}(t) = \rho_{NO_x} k(t) Q_{exh}(t) C_{NO_x}(t) \quad (3)$$

where  $\rho_{NO_x}$  is the density of NO<sub>x</sub>,  $k(t)$  is the dry-to-wet emission correction factor,  $Q_{exh}$  is the exhaust volumetric flow calculated by the ECU and  $C_{NO_x}$  is the NO<sub>x</sub> concentration in the exhausts, either measured experimentally or predicted by the LightGBR. In this case, the normalization is performed separately for each plot, keeping the values in the range between 0 and 1. The relative error between the modelled and the experimental cumulative NO<sub>x</sub> is thereby reported.

The relative absolute error made on the final NO<sub>x</sub> mass is around 4% in the first RDE test and lower than 1% in the second RDE test, confirming the high accuracy of the model.

Finally, a correlation plot is shown in Figure 20. These scatter plots carry out a sample-by-sample comparison between real and predicted NO<sub>x</sub>. The black dashed line represents the ideal correlation line, where all the points of the scatter should lay in case of a perfect model. The correlation is very good for both RDE tests, showing  $R^2 = 0.953$  and  $nRMSE = 0.038$  in the first RDE and  $R^2 = 0.947$  and  $nRMSE = 0.04$  in the second RDE. This shows that the model correctly learned the dynamics behind NO<sub>x</sub> formation. An evident dispersion of samples is visible, especially for the RDE 2 (red scatter). This is due to local delay between experimental and predicted values that for many reasons cannot be always perfect (due to sensor dynamics, flow dynamics, error in the models, measurement noise, etc.). However, the dispersion of the scatter is symmetric with respect to the ideal correlation line, meaning that the general trend is very well-represented by the model.



**Figure 20.** Sample-by-sample correlation scatter. In blue, the result of RDE 1, in red the results of RDE 2.

#### 4. Conclusions

The obtained results confirm that the data-driven approach represents an effective tool for predicting the instantaneous NO<sub>x</sub> engine-out emissions during all the typical maneuvers performed during an RDE cycle. The Sliding Window approach increases the accuracy of all the models, as well as the training time. The result of the sensitivity analysis on window size shows an optimal width of 30 samples, corresponding to 3 s. The cross-validation technique highlighted that the Feed-forward Neural Network ( $R^2 = 0.887$ ), the Random Forest ( $R^2 = 0.897$ ) and the Light Gradient Boosting ( $R^2 = 0.894$ ) are the most accurate models on the validation set, with the use of a 30-sample Sliding Window. In particular, the Light Gradient Boosting Regressor is the best compromise between accuracy ( $R^2 = 0.947$  and  $R^2 = 0.953$  on the RDE test cycles) and training time ( $t = 11$  s) for this application. A good correlation is obtained when the LightGBR is applied to infer the NO<sub>x</sub> emissions under steady-state conditions, showing a coefficient of determination of  $R^2 = 0.947$  and  $MAPE = 6.81\%$ , which can be considered good results. The high accuracy also achieved under steady-state conditions demonstrates the remarkable reliability of the model, even for different applications. The proposed methodology can support the engine development phase to reduce the number of experimental tests at the bench needed for the emission calibration, leading to a potential cost and time reduction.

#### 5. Future Works

The dataset considered in this work is composed of RDE cycles with different aggressiveness and maneuvers. Therefore, the high accuracy of the model prediction shows that the training set conditions are quite representative of the test set, and that the input features selected are sufficiently descriptive of the physical phenomena occurring inside the engine and affecting NO<sub>x</sub> formation. As a further development, this NO<sub>x</sub> model could be tested on other types of homologation cycles to verify its robustness under different working conditions.

Moreover, the methods and the techniques explained in this paper will be extended to other pollutant species, such as CO and HC. Another interesting development will be the introduction of a data-driven model of the after-treatment system to also extend this approach to the prediction of the tailpipe emissions, which are relevant for emission legislation. In particular, the output of the engine-out model presented in this paper can be used as an input of the catalyst data-driven model. Therefore, once the two data-driven models are coupled, they represent a digital twin of the complete exhaust system. The models of the tailpipe emissions can be trained with the experimental on-road measurements to obtain a model as close as possible to the variability of typical real driving conditions.

In this case, the emissions will be measured by means of a Portable Emissions Measurement System (PEMS) directly installed on the vehicle.

Finally, a future application of this method is in the field of virtual sensing, where NO<sub>x</sub> emissions are estimated without the presence of physical sensors. This approach could be used to estimate emissions where physical sensors cannot be installed due to engineering constraints or the cost limitations. For example, an interesting application of such models is for durability tests, where it is not possible to install a PEMS for the emission measurements, where such models can be used to detect the critical maneuvers that produce peaks of the pollutant emissions. This can provide an important contribution to lead towards a more aware and reliable calibration of the ECU control strategies.

**Author Contributions:** Formal analysis, A.B.; Supervision, M.B. and N.C.; Writing—original draft, E.G. All authors have read and agreed to the published version of the manuscript.

**Funding:** This research received no external funding.

**Conflicts of Interest:** The authors declare no conflict of interest.

### Abbreviations

The following abbreviations are used in this manuscript:

CV	Cross-Validation
EFB	Exclusive Feature Bundling
FIP	Feature Importance Permutation
GOSS	Gradient-based One-Side Sampling
LightGBR	Light Gradient Boosting Regressor
NEDC	New European Driving Cycle
PR	Polynomial Regressor
RDE	Real Driving Emission
RF	Random Forest
SVR	Support Vector Regressor
WLTC	Worldwide harmonized Light vehicle Test Cycle

### References

- Luckow, A.; Kennedy, K.; Manhardt, F.; Djerekarov, E.; Vorster, B.; Apon, A. Automotive big data: Applications, workloads and infrastructures. In Proceedings of the 2015 IEEE International Conference on Big Data (Big Data), Santa Clara, CA, USA, 29 October–1 November 2015. [\[CrossRef\]](#)
- Wei, H. Analysis on the Applications of AI in Vehicles and the Expectation for Future. In Proceedings of the 2020 5th International Conference on Information Science, Computer Technology and Transportation (ISCTT), Shenyang, China, 13–15 November 2020. [\[CrossRef\]](#)
- Montáns, F.J.; Chinesta, F.; Gómez-Bombarelli, R.; Kutz, J.N. Data-driven modeling and learning in science and engineering. *C. R. Mécanique* **2019**, *347*, 845–855. [\[CrossRef\]](#)
- Zhou, D.P.; Hu, Q.; Tomlin, C.J. Quantitative comparison of data-driven and physics-based models for commercial building HVAC systems. In Proceedings of the 2017 American Control Conference (ACC), Seattle, WA, USA, 24–26 May 2017. [\[CrossRef\]](#)
- Karri, V.; Ho, T.N. Predictive models for emission of hydrogen powered car using various artificial intelligent tools. *Neural Comput. Appl.* **2009**, *18*, 469–476. [\[CrossRef\]](#)
- Liu, J.; Huang, Q.; Ulishney, C.; Dumitrescu, C.E. Comparison of Random Forest and Neural Network in Modeling the Performance and Emissions of a Natural Gas Spark Ignition Engine. *J. Energy Resour. Technol. Trans. ASME* **2022**, *144*, 032310. [\[CrossRef\]](#)
- Papaioannou, N.; Fang, X.; Leach, F.; Lewis, A.; Akehurst, S.; Turner, J. A Random Forest Algorithmic Approach to Predicting Particulate Emissions from a Highly Boosted GDI Engine. In Proceedings of the SAE Naples ICE 2021 Conference, Naples, Italy, 12–16 September 2021. [\[CrossRef\]](#)
- Moradi, M.H.; Heinz, A.; Wagner, U.; Koch, T. Modeling the emissions of a gasoline engine during high-transient operation using machine learning approaches. *Int. J. Engine Res.* **2022**, *23*, 1708–1716. [\[CrossRef\]](#)
- Huang, Q.; Liu, J.; Ulishney, C.; Dumitrescu, C.E. On the use of artificial neural networks to model the performance and emissions of a heavy-duty natural gas spark ignition engine. *Int. J. Engine Res.* **2022**, *23*, 14680874211034409. [\[CrossRef\]](#)
- Netzer, C.; Franken, T.; Seidel, L.; Lehtiniemi, H.; Mauss, F. Numerical Analysis of the Impact of Water Injection on Combustion and Thermodynamics in a Gasoline Engine Using Detailed Chemistry. *SAE Int. J. Engines* **2018**, *11*, 1151–1166. [\[CrossRef\]](#)

11. Wang, Z.; Wang, J.X.; Shuai, S.J.; Zhang, F. *Numerical Simulation of HCCI Engine with Multi-Stage Gasoline Direct Injection Using 3D-CFD with Detailed Chemistry*; SAE International: Warrendale, PA, USA, 2004. [[CrossRef](#)]
12. Choi, S.; Kolodziej, C.P.; Hoth, A.; Wallner, T. *Development and Validation of a Three Pressure Analysis (TPA) GT-Power Model of the CFR F1/F2 Engine for Estimating Cylinder Conditions*; Argonne National Lab. (ANL): Argonne, IL, USA, 2018. [[CrossRef](#)]
13. Brusa, A.; Cavina, N.; Rojo, N.; Mecagni, J.; Corti, E.; Moro, D.; Cucchi, M.; Silvestri, N. Development and experimental validation of an adaptive, piston-damage-based combustion control system for SI engines: Part 2-implementation of adaptive strategies. *Energies* **2021**, *14*, 5367. [[CrossRef](#)]
14. Riegler, U.G.; Bargende, M. *Direct Coupled 1D/3D-CFD-Computation (GT-Power/Star-CD) of the Flow in the Switch-Over Intake System of an 8-Cylinder SI Engine with External Exhaust Gas Recirculation*; SAE International: Warrendale, PA, USA, 2002. [[CrossRef](#)]
15. Millo, F.; di Lorenzo, G.; Servetto, E.; Capra, A.; Pettiti, M. Analysis of the performance of a turbocharged s.i. engine under transient operating conditions by means of fast running models. *SAE Int. J. Engines* **2013**, *6*, 968–978. [[CrossRef](#)]
16. Han, Z.; Reitz, R.D. Turbulence Modeling of Internal Combustion Engines Using RNG  $k-\epsilon$  Models. *Combust. Sci. Technol.* **1995**, *106*, 267–295. [[CrossRef](#)]
17. Boiarciuc, A.; Floch, A. *Evaluation of a 0D Phenomenological SI Combustion Model*; SAE International: Warrendale, PA, USA, 2011. [[CrossRef](#)]
18. Ravaglioli, V.; Moro, D.; Serra, G.; Ponti, F. *MFB50 On-Board Evaluation Based on a Zero-Dimensional ROHR Model*; SAE International: Warrendale, PA, USA, 2011. [[CrossRef](#)]
19. Cavina, N.; Migliore, F.; Carmignani, L.; Palma, S.D. *Development of a Control-Oriented Engine Model Including Wave Action Effects*; SAE International: Warrendale, PA, USA, 2009. [[CrossRef](#)]
20. Scocozza, G.F.; Silvagni, G.; Brusa, A.; Cavina, N.; Ponti, F.; Ravaglioli, V.; Cesare, M.D.; Panciroli, M.; Benedetti, C. *Development and Validation of a Virtual Sensor for Estimating the Maximum in-Cylinder Pressure of SI and GCI Engines*; SAE International: Warrendale, PA, USA, 2021. [[CrossRef](#)]
21. Ranuzzi, F.; Cavina, N.; Brusa, A.; Cesare, M.D.; Panciroli, M. *Development and Software in the Loop Validation of a Model-Based Water Injection Combustion Controller for a GDI TC Engine*; SAE International: Warrendale, PA, USA, 2019. [[CrossRef](#)]
22. Najafi, G.; Ghobadian, B.; Tavakoli, T.; Buttsworth, D.R.; Yusaf, T.F.; Faizollahnejad, M. Performance and exhaust emissions of a gasoline engine with ethanol blended gasoline fuels using artificial neural network. *Appl. Energy* **2009**, *86*, 630–639. [[CrossRef](#)]
23. Mohammad, A.; Rezaei, R.; Hayduk, C.; Delebinski, T.O.; Shahpouri, S.; Shahbakhti, M. *Hybrid Physical and Machine Learning-Oriented Modeling Approach to Predict Emissions in a Diesel Compression Ignition Engine*; SAE International: Warrendale, PA, USA, 2021. [[CrossRef](#)]
24. Paul, A.; Bhowmik, S.; Panua, R.; Debroy, D. Artificial Neural Network-Based Prediction of Performances-Exhaust Emissions of Diesohol Piloted Dual Fuel Diesel Engine Under Varying Compressed Natural Gas Flowrates. *J. Energy Resour. Technol.* **2018**, *140*, 112201. [[CrossRef](#)]
25. Khurana, S.; Saxena, S.; Jain, S.; Dixit, A. Predictive modeling of engine emissions using machine learning: A review. *Mater. Today: Proc.* **2020**, *38*, 280–284. [[CrossRef](#)]
26. Altug, K.B.; Kucuk, S.E. Predicting Tailpipe NO<sub>x</sub> Emission using Supervised Learning Algorithms. In Proceedings of the 2019 3rd International Symposium on Multidisciplinary Studies and Innovative Technologies (ISMSIT), Ankara, Turkey, 11–13 October 2019. [[CrossRef](#)]
27. Papaioannou, N.; Fang, X.H.; Leach, F.; Davy, M.H. Prediction of NO<sub>x</sub> emissions for a range of engine hardware configurations using artificial neural networks. In *Internal Combustion Engine Division Fall Technical Conference*; American Society of Mechanical Engineers: New York, NY, USA, 2021. [[CrossRef](#)]
28. Alonso, J.M.; Alvarruiz, F.; Desantes, J.M.; Hernández, L.; Hernández, V.; Moltó, G. Combining neural networks and genetic algorithms to predict and reduce diesel engine emissions. *IEEE Trans. Evol. Comput.* **2007**, *11*, 46–55. [[CrossRef](#)]
29. Cornec, C.M.L.; Molden, N.; van Reeuwijk, M.; Stettler, M.E. Modelling of instantaneous emissions from diesel vehicles with a special focus on NO<sub>x</sub>: Insights from machine learning techniques. *Sci. Total Environ.* **2020**, *737*, 139625. [[CrossRef](#)]
30. Ozmen, M.I.; Yilmaz, A.; Baykara, C.; Ozsoysal, O.A. Modelling Fuel Consumption and NO Emission of a Medium Duty Truck Diesel Engine with Comparative Time-Series Methods. *IEEE Access* **2021**, *9*, 81202–81209. [[CrossRef](#)]
31. Wang, X.; Liu, W.; Wang, Y.; Yang, G. A hybrid NO<sub>x</sub> emission prediction model based on CEEMDAN and AM-LSTM. *Fuel* **2022**, *310*, 122486. [[CrossRef](#)]
32. Smola, A.J.; Schölkopf, B. A tutorial on support vector regression. *Stat. Comput.* **2004**, *14*, 199–222. [[CrossRef](#)]
33. Murty, M.N.; Raghava, R. Kernel-based SVM. In *Support Vector Machines and Perceptrons*; Springer: Cham, Switzerland, 2016. [[CrossRef](#)]
34. Breiman, L. Random forests. *Mach. Learn.* **2001**, *45*, 5–32. [[CrossRef](#)]
35. Ke, G.; Meng, Q.; Finley, T.; Wang, T.; Chen, W.; Ma, W.; Ye, Q.; Liu, T.Y. LightGBM: A Highly Efficient Gradient Boosting Decision Tree. In Proceedings of the Advances in Neural Information Processing Systems, Long Beach, CA, USA, 4–9 December 2017.
36. Szandała, T. Review and comparison of commonly used activation functions for deep neural networks. In *Bio-Inspired Neurocomputing*; Springer: Singapore, 2021; Volume 903. [[CrossRef](#)]
37. Fischer, M. Transient NO<sub>x</sub> Estimation using Artificial Neural Networks. *IFAC Proc. Vol.* **2013**, *46*, 101–106. [[CrossRef](#)]
38. EU. Commission Regulation 2017/1154. *Off. J. Eur. Union* **2017**, *1154*.

39. Chindamo, D.; Gadola, M. What is the Most Representative Standard Driving Cycle to Estimate Diesel Emissions of a Light Commercial Vehicle? *IFAC-PapersOnLine* **2018**, *51*, 73–78. [[CrossRef](#)]
40. Zar, J.H. Significance testing of the spearman rank correlation coefficient. *J. Am. Stat. Assoc.* **1972**, *67*, 578–580. [[CrossRef](#)]
41. Altmann, A.; Toloşi, L.; Sander, O.; Lengauer, T. Permutation importance: A corrected feature importance measure. *Bioinformatics* **2010**, *26*, 1340–1347. [[CrossRef](#)]
42. Anggoro, D.A.; Mukti, S.S. Performance Comparison of Grid Search and Random Search Methods for Hyperparameter Tuning in Extreme Gradient Boosting Algorithm to Predict Chronic Kidney Failure. *Int. J. Intell. Eng. Syst.* **2021**, *14*, 198–207. [[CrossRef](#)]
43. Asif, M.A.A.R.; Nishat, M.M.; Faisal, F.; Dip, R.R.; Uday, M.H.; Shikder, M.F.; Ahsan, R. Performance evaluation and comparative analysis of different machine learning algorithms in predicting cardiovascular disease. *Eng. Lett.* **2021**, *29*, 731–741.