

Article

Exploiting Data Analytics and Deep Learning Systems to Support Pavement Maintenance Decisions

Ronald Roberts * , Laura Inzerillo  and Gaetano Di Mino 

DIING—Department of Engineering, University of Palermo, Viale delle Scienze, 90128 Palermo, Italy; laura.inzerillo@unipa.it (L.I.); gaetano.dimino@unipa.it (G.D.M.)

* Correspondence: ronaldanthony.roberts@unipa.it; Tel.: +39-328-442-8206

Abstract: Road networks are critical infrastructures within any region and it is imperative to maintain their conditions for safe and effective movement of goods and services. Road Management, therefore, plays a key role to ensure consistent efficient operation. However, significant resources are required to perform necessary maintenance activities to achieve and maintain high levels of service. Pavement maintenance can typically be very expensive and decisions are needed concerning planning and prioritizing interventions. Data are key towards enabling adequate maintenance planning but in many instances, there is limited available information especially in small or under-resourced urban road authorities. This study develops a roadmap to help these authorities by using flexible data analysis and deep learning computational systems to highlight important factors within road networks, which are used to construct models that can help predict future intervention timelines. A case study in Palermo, Italy was successfully developed to demonstrate how the techniques could be applied to perform appropriate feature selection and prediction models based on limited data sources. The workflow provides a pathway towards more effective pavement maintenance management practices using techniques that can be readily adapted based on different environments. This takes another step towards automating these practices within the pavement management system.

Keywords: pavement management systems; pavement maintenance decisions; road asset databases; data mining; feature importance; deep learning



Citation: Roberts, R.; Inzerillo, L.; Di Mino, G. Exploiting Data Analytics and Deep Learning Systems to Support Pavement Maintenance Decisions. *Appl. Sci.* **2021**, *11*, 2458. <https://doi.org/10.3390/app11062458>

Academic Editor:
Sakdirat Kaewunruen

Received: 3 February 2021
Accepted: 4 March 2021
Published: 10 March 2021

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



Copyright: © 2021 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

1.1. The Need for Information to Support Pavement Management Decisions

Roads are a critical component of any society from a small town scale to the highest levels as they allow for the essential movement of goods and services [1] with road transport being the widely used transport mode for both passenger and freight transport in Europe [2]. In many countries worldwide the public road network is the biggest publicly owned asset [3]. To this end, these networks must be kept in a suitable condition making the road agencies' job of maintaining them a critical component towards the development of any region. Road agencies have the responsibility to develop programs that guide rehabilitation and maintenance practices and they have to enact critical decisions in terms of which areas should be prioritized and when interventions should be made. Further complicating their decisions are steadily reducing budgets for these activities [4,5] that will likely see greater depletions given the current global economic pandemic related situation [6].

The decisions of these road agencies are generally based on the implementation of a pavement management system (PMS) which falls under the analysis stage of the pavement management process and can be executed at either a network or project level [7]. The PMS is the most typical analytical method and it attempts to optimize the use of financial resources based on the particular needs of the road network [8]. The PMS is a combined toolbox, which uses key road data to help decision-makers create and make the best decisions to allow for optimum system conditions over time [9]. Three critical basic requirements for the adequate application of a PMS [10] are:

1. It serves different types of users within the organization.
2. It allows for good decision making concerning the decided programs and projects and allows for the timely execution of projects.
3. It makes good use of existing technologies and new ones once they are available.

The system, while effective, relies heavily on the input data, to ensure appropriate strategies are established. The data used as inputs can generally be categorized as inventory road data and pavement condition data, which are used within the system using various functions and models to produce outputs about the overall network or project assessments and proposed maintenance strategies [9,11]. This is visualized in Figure 1.

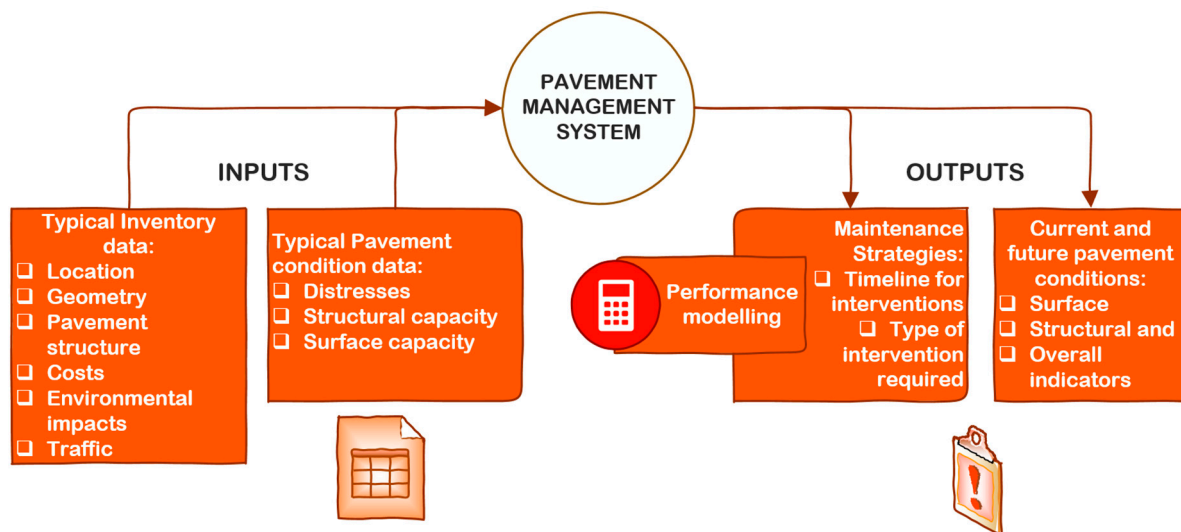


Figure 1. How data are utilized in the PMS.

However, acquiring this data can be problematic due to typically high costs, time, resources and high-level team expertise required to do so [12,13]. In some authorities, the collection of pavement condition data is done manually in over 90% of the occurrences [14] which therefore allocates a tremendous amount of subjectivity into the process. This can lead to the planning of ineffective strategies. As a result, road asset databases need to be kept up to date and effective. There are significant needs in transportation systems, now and in the future to develop smart infrastructure, that can carry out ‘self-reporting and self-managing infrastructure’ through the application of automated information flows [15]. It should also be noted that the mere presence of a road asset database does not necessarily lead to effective solutions and decision-making. The effective management of these databases should be guided by their relevance, reliability and affordability [16]. These factors remain important regardless of the development of new technologies and methods and should be integral to the design and introduction of any new technique or model. Data collection should also be done at the lowest level of detail satisfactory enough to make effective decisions [17]. Therefore, when designing a database or collection system, the networks’ characteristics should be integral to the process along with the ability of the authority to maintain the system both financially and technically.

1.2. Relationship of Factors and Features that Contribute to Pavement Maintenance Interventions

The selection of potential sites for pavement maintenance is usually dependent on the pavement age and condition, frequency of maintenance interventions, traffic rates, accessibility and safety [18]. The preference between these sites can then be decided based on performance and distress surveys. In practice, there are a multitude of different road performance indicators that can be employed during the decision-making processes [3]. The use of these indicators is generally dependent on the road agency, their needs and

the available data such as distress and surface texture [9]. Many PMS's will utilize top-down approaches wherein funding levels required to meet long-term performance goals are estimated but in these scenarios, further steps are still required to prioritize specific segments for maintenance interventions [19]. However, it has been found that issues have been identified within the methodologies utilized for the PMS in effectively considering multi-attribute condition information for modelling processes and with uncertainties linked primarily to pavement deterioration despite the establishment of the importance of environmental and economic conditions [20].

Performance indicators are useful and important because they help to efficiently allocate resources amongst options based on the availability of resources [21]. It has been suggested that frameworks for performance indicators for road assets should be divided across two levels—general performance indicators which can provide an overview and can be ascertained from public statistics and detailed objective indicators of service quality and institutional effectiveness [21]. These overview characteristics are important as they can provide a summarized view of the network's situation. Concerning the overview performance indicators, there are typical feature groups that should be assessed to ascertain the summary of the road asset [21,22] which include:

1. Network parameters—geometric configuration and dimensions of roads
2. Asset values
3. Road users—types of users and trip purposes
4. Demography and economic circumstances—population and land area
5. The density of network and roads
6. Use of roads—travel by class
7. Safety—accidents and fatalities

With regard to these indicators, it is possible to retrieve a lot of this data from local sources such as censuses, therefore making them a viable resource outlet for obtaining data for modelling scenarios for small or under-resourced cities. The data sourcing can be complicated based on specific circumstances but it is a necessary step to adequately understand particular situations and predict future trends.

1.3. How Have Pavement Management Decisions and Approaches Been Supported by Data Analytics?

Once data have been identified in a system, the next important decision is the type of modelling and estimation to be done to support maintenance decisions by predicting future conditions to understand and manage maintenance and rehabilitation strategies. Models are typically grouped into four categories [9] as seen in Table 1. The application of them is heavily based on the available resources. There have been numerous models produced and researched with the majority of them being specific to a combination of variable or conditions, highlighting the need for a workflow that can work with different variable combinations [23].

Table 1. Overview of typical types of pavement performance models.

Approach	Description
Deterministic models	Models that produce a single dependent value such as pavement condition from a combination of different variables concerning the characteristics of the pavement and network
Probabilistic models	Models that predict a range of values for the dependent variable with probabilities of changes based on different conditions and timelines
Bayesian models	Models that combine objective and subjective data and having their variables described in a probabilistic distribution manner
Subjective/Expert-based models	Similar to deterministic ones but with structures based on opinion and not historical data

Considering and using these model types, there are several areas within the PMS decision-making processes where Artificial intelligence-based systems can help support decisions. These include estimation of pavement condition, assessment of maintenance needs, identification and selection of maintenance actions and prioritization of maintenance programs [24]. Studies have tried to develop algorithms using techniques such as decision trees to predict Pavement Condition Index (PCI) [25] values (a common index used to represent road conditions), and therefore future pavement conditions of road sections using input data such as historical PCI values, weather, historical maintenance data and traffic [26]. Additional studies have focused on the International Roughness Index (IRI), which is an index that considers the roughness of the profile of the road [27]. Using this index, random forests regressions have been used to predict values of IRI to determine pavement roughness of sections utilizing traffic information, previous IRI values and pavement distresses [28]. Artificial Neural networks (ANN) have also been constructed using the IRI as an input to identify maintenance strategies with reliance on large national datasets [29,30]. Other studies have tried to predict PCI values using neural networks. Planning of road interventions utilizing Fuzzy based networks have also been considered [18]. In other cases, ANNs were utilized to predict an exact pavement maintenance decision based on inputs of distress data, functional class, traffic and pavement structure [31]. Another method of analysis is the use of Genetic Algorithms, which were shown as an effective tool for predicting maintenance programs and have been combined within neural networks to bolster their effectiveness [32–34]. These previous studies all represent attempts at helping plan maintenance interventions but have commonly relied on large databases. Additionally, whilst interoperability issues in machine learning have been raised about the difficulty in understanding relationships between model outputs and inputs [35], feature selection has been successfully used to filter out redundant data, improve accuracies and help produce more explainable models [36].

1.4. Development of Strategies for Pavement Maintenance in Agencies with Limited Data and Resources

Data Analytics has been considered as an effective way of handling pavement condition data and relating predicting pavement conditions [26]. However, in many instances, many large databases do not provide particularly useful information to allow for efficient decisions to be made with many of the databases having subjective data [37]. There have been attempts at creating datasets that will help boost deep learning activities in the field of civil engineering where datasets of different infrastructures have been set up [38]. However, these studies have tended to focus on image-based datasets. Image-based analysis can be very effective and recent studies have shown how low-cost models can be developed to detect distresses within a network using smartphones [39]. However, even within these systems, the question of where to direct the image survey and how to understand the feature characteristics of the network still needs to be addressed for them to be practically applied.

There is one large database that includes pavement features and factors called the Long-Term Pavement Performance (LTPP) database [40] which is an enormous database developed and maintained by the Federal Highway Administration in the U.S. It has allowed for the study and development of many different algorithms and decision-based studies for predicting different future conditions of pavements and developing maintenance strategies, with reviews showing that the majority of research using ANN for predictions is done using the LTPP database [41]. If similar databases could be set up for all agencies, it would be a great help in their planning processes. However, in most cases, road agencies do not have access or funds to create such a large database and the associated algorithms and models developed for road maintenance. This means the systems and predictions cannot be applied to their cities and Whilst it can be reasoned that more data can produce higher accuracies in models, the costs of acquiring the additional data must be amply considered by the road authority. Therefore, there should be a balance between available budgets and information pursuit [42]. The use of feature selection has been considered useful to help increase the accuracy of models by utilizing fewer feature characteristics but those of the

most significant importance [42]. To this end, machine learning approaches have been used to obtain the most important features within a dataset with conclusions that these approaches could potentially be effective for cases where there is limited data available.

2. Aim of the Study

Given the constraints faced by road agencies worldwide, this study was developed to create a workflow for agencies to carry out efficient network information mining to understand the important features in their networks and when interventions should be done. The inputs to this process include historical data and information on the precise environment with the output being an understanding of which are the important features that contribute to maintenance decisions in this particular environment and a model to predict future maintenance activities. To enable this, a case study was developed in Palermo, Italy utilizing historical information on road interventions, network characteristics and census information to provide a succinct but effective overview of important features as identified by previous studies [21,22] from low-cost and easily available sources. An innovative aspect of this study is that it considers a small database, which does not have the many descriptive and technical features of large-scale databases such as the LTPP one. To this end, the workflow highlights important factors that contribute to road maintenance schedules without relying on expensive data collection systems. The study aims to carry out an effective feature selection with limited data and use these features to create a model for future predictions. A model is constructed using open source and flexible computing libraries to predict when interventions should be done within the network. The output of this process would allow for the optimization of available information at a low-cost and with limited resources, which helps to advance the pavement management planning and implementation process in limited data circumstances.

3. Methodology

Given the state of the research field and the need to create a pipeline that could be augmented to different agencies based on data availabilities, a general workflow was created that could be adapted based on the particular circumstances of a road agency. This workflow is shown in Figure 2. The workflow highlights the sources of data, the type of analysis needed at each stage and the outcome, which is a model to identify times for interventions for various roads within the network.

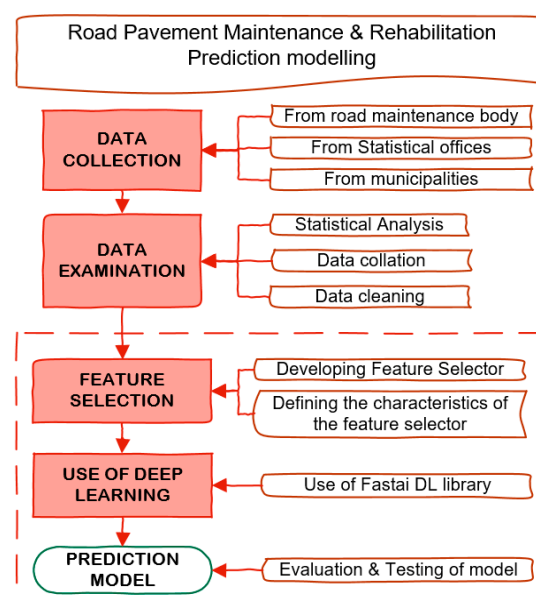


Figure 2. Workflow for study.

Within this methodological framework the use of well-known Artificial Intelligence (AI) tools, including deep learning systems, are particularly well suited due to the intrinsic nature and complexity of the issues addressed. Furthermore, it has been shown by others that AI tools are particularly important to use in pavement management situations because of uncertain and/or incomplete information that appears in the field [43] which is often the case in smaller or under-resourced authorities.

3.1. Data Collection

The data collected was obtained using several open datasets available from government authorities, both locally and nationally namely the statistical offices in Italy and Palermo, the local municipality and the assigned road maintenance company. It was very important given the context of the study to have available data. Data characteristics are further described in Section 4.

3.2. Feature Selection Workflow

Once the data were obtained, the first step was a visual overview data analysis to establish relationships between features. To do this, a pandas data frame [44] was set up to explore the dataset along with the use of matplotlib [45] and seaborn [46] for visualization through graphical relationship plots. This process allows for easy visualization of the statistics of the dataset and any necessary data cleaning before processing can begin. They are also open source packages and therefore can easily be utilized for further works with no large costs to agencies. A descriptive analysis of the dataset was performed to highlight distributions and identify relationships between the various features. The data were also preprocessed to ensure there were no missing values and inconsistencies.

3.2.1. Developing Feature Selector

Once the data were preprocessed, the next step was preparing a feature selector tool adapted to the particular dataset but whose configuration could be modified based on the type of data typically available. Feature selection involves using only appropriate features that explain the dependent variable and the process helps produce good learners and models [47]. In this case, the dependent variable is the year when road intervention took place for a given road in the network. A feature selector is a tool that can be used for producing strong datasets for machine learning purposes. For the study, a feature selector tool was developed considering a selector developed in python utilizing the LightGBM library [48]. For the feature selector, the tool's backbone model was redesigned to identify feature importance using a gradient boosting machine from the CatBoost (Category Boosting) [49] library instead. This construction was done based on the strength of CatBoost to handle categorical features and the high presence of these features within this study and typical of this type of data. For the feature selector's backbone, several ensemble deep learning gradient boosting algorithms were considered. Gradient boosting algorithms are considered powerful tools that work by building ensemble tree predictors carrying out gradient descent within a functional space [49]. The algorithms sequentially create base models. The accuracy of these models versus others are considered higher given the production of multiple models sequentially and the emphasis on training cases that are more difficult to evaluate, making the mistake more evident. In the process, the examples that were considered harder to estimate in the earlier base models are forced to appear more frequently in the training data. The subsequent base models are all aimed at correcting the prior mistakes in earlier models. The boosting process utilizes weak base models that are easier to predict and combines them to get one highly accurate model. The process fits subsequent models that allow for minimization of particular loss function objectives and errors averaged across the training data [50]. The pseudocode for Friedman's gradient boosting algorithm is subsequently shown in Algorithm 1 [51].

Algorithm 1: Friedman's Gradient Boost algorithm**Inputs:**

- input data $(x, y) \ N_i=1$
- number of iterations, M
- choice of the loss-function $\Psi(y, f)$
- choice of the base learner model $h(x, \theta)$

Algorithm:

1. initialize \hat{f}_0 with a constant
2. for t = 1 to M do
3. compute the negative gradient $g_t(x)$
4. fit a new base-learner function $h(x, \theta_t)$
5. find the best gradient descent step-size ρ_f

$$\rho_t = \operatorname{argmin}_{\rho} \sum_{i=1}^N \Psi \left[y_i, \hat{f}_{t-1}(x_i) + \rho h(x_i, \theta_t) \right]$$

6. update the function estimate:
 $\hat{f}_t \leftarrow \hat{f}_{t-1} + \rho_t h(x, \theta_t)$
7. **end for**

Gradient boosting methods offer a high level of customizable flexibility to any specific data-driven job [52] making it appealing for this study. There are several available gradient boosting models that are applied for prediction and regression problems. These include XGBoost [53], AdaBoost [54] and LightGBM [55] with the recent introduction of the CatBoost [49] libraries which provide a new approach, especially for categorical features. As a result of this, the CatBoost library was chosen given its ability to handle categorical features and previous studies identifying a higher performance of this algorithm versus the others previously mentioned [56,57]. This is particularly important to the model given the number of categorical features present in the dataset. CatBoost is different from the other gradient boosting algorithms in the following ways:

- The library handles categorical features during training as opposed to during the preprocessing time. It also utilizes the entire dataset for training. For each training example, the library carries out a random permutation of the dataset and calculates an average label value for the particular example with the same category value placed before the one provided by the permutation [58].
- The library allows unbiased boosting with categorical features and feature combinations where all the categorical features can be combined as a new feature. [49]
- The use of a fast scorer allows utilizing oblivious trees as base predictors [58].

The CatBoost library [49,59] introduced two vital advances: ordered boosting, which is a permutation driven alternative to the classic choice and a state-of-the-art algorithm to process categorical features. Most gradient boosting algorithms use encoding for categorical features but CatBoost utilizes an innovative strategy for this. An ordering principle is utilized which feeds training examples sequentially which allows the target values to rely on the historically observed values. For this to be possible, CatBoost utilizes a random artificial time permutation for the training examples for the gradient boosting process. Based on these factors, the CatBoost library was exploited.

3.2.2. Defining the Characteristics of the Feature Selector

Once the backbone model was decided, the hyperparameters of the model were needed. Within the selector, 12 training runs were utilized to lower variances with the CatBoost model itself being run for 500 iterations, and early stopping using a validation set was used to avoid overfitting of the data. Within the implementation, CatBoost gives indices of categorical columns to allow them to be encoded as one-hot encoding utilizing `one_hot_max_size`. This is based on the number of unique values available for the

features considered. This is important for this study as the process can be calibrated based on the number of unique values within the dataset instead of using an assumptive value. An 'one_hot_max_size' of 25 was utilized given the presence of data for 25 neighborhoods (as is the case in Palermo) which have related categorical features. The use of the `cat_features` parameter allows the user to pass the column indices for the categorical features through the model for preprocessing and use in the one-hot encoding. The model utilized within the setup was defined as follows using the `CatBoostRegressor` model setup: `model = catboost.CatBoostRegressor(eval_metric = 'RMSE', one_hot_max_size = 25, depth = 10, iterations = 500, l2_leaf_reg = 9, learning_rate = 0.05)`.

Within the definition of the model, the `eval_metric` used was the Root Mean Square (RMSE) with an 'l2_leaf_reg' value of 9 used (regularization term used to regularize the objective function and minimize both loss and complexity of the model) and a depth of 10 (depth of tree used) and a learning rate of 0.05. These parameters were used based on iterations of models and were found to be the most appropriate. Early stopping of the model was also carried out with a test set of 20% to validate and ensure overfitting was not done. During the run of the model, the feature importance levels was recorded using the 'feature_importances' function of the `Catboost` library. These were recorded and sorted in descending order to showcase the features with the most impact on the model. Finally, the values were normalized based on the total contribution to the model result and these results were graphically represented.

3.3. Use of Deep Learning for Tabular Sets—Use of FastAI Deep Learning Library

Once the important features were identified, the next step was developing a model that could attempt to predict the year for intervention when roads would reach a similar degradation level as those identified in the dataset, based on the available characteristics. To accomplish this, a deep learning model was developed for the application on a tabular dataset. Deep Learning is a machine learning technique based on "learn by example" principle. A Deep Learning model commonly uses large sets of labeled data which is processed through neural network architectures built-up with many layers. By exploiting the input data, through synaptic connections between adjacent layers with feed-forward propagation, the output data are computed. The output quality depends on the recognition accuracy, which in turn depends on the consistency of both the neural structure and process data. Therefore, it is clear that the choice of the input data (features) and the neural network architecture, including the typical training and testing parameters, are important to achieve decision-making targets. Whilst, a significant amount of research has been focused on the utilization of deep learning for image analysis, there is a lot of merit to leverage the power of deep learning for tabular dataset analysis as well. New base architectures have shown a lot of worth in this type of analysis over traditional classification and grouping methods [60] and by considering a deep learning model, this workflow can leverage both the resources of the gradient boosting for feature selection and deep learning for the model development. Deep learning models also have shown advantages over traditional approaches when datasets have high cardinality categorical variables. This is as a result of the model using embeddings for these categorical features [61] and this is important given the presence of these variables in the dataset.

For this study, the open source `FastAI` framework [60] was utilized. This library was chosen because of its ease of use, therefore, making it easily interpretable and easily replicated for authorities and users. The library was designed around two main goals: to make its implementation and use approachable and productive whilst maintaining flexible configurability. This is important as there is not a significant level of training needed to create the models and therefore explaining its implementation for road professionals would not be difficult. `FastAI` also offers excellent support for tabular datasets with built-in loaders to handle this type of data. Within the `FastAI` library, there is high-level API support for tabular datasets which allows the creation of models based on characteristics of the dataset available. Within the tabular data setup, there are also provisions to denote

categorical elements, which is an added advantage, considering the dataset at hand. The tabular API is made up of the components shown in Figure 3.

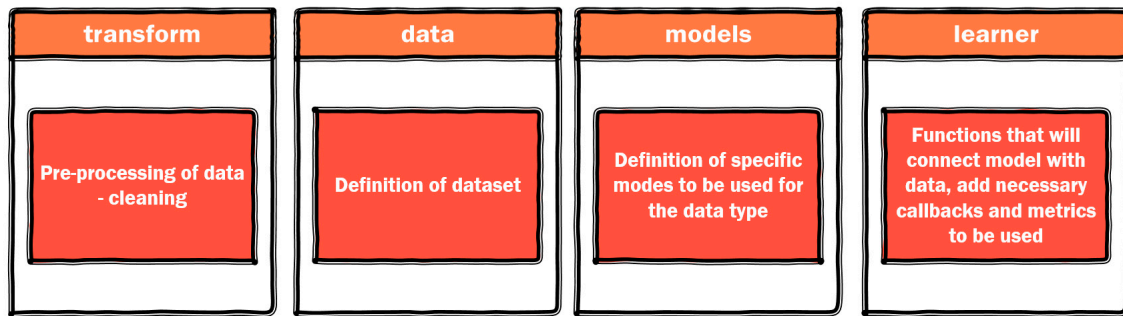


Figure 3. FastAI tabular module structure.

The library also has model-fitting methods within its structure namely the use of 'lr_find', which allows the establishment of a good learning rate to be used rather than going through an iterative process at the training level. Within the model, the training was regularized utilizing dropouts for dense layers and embeddings (0.05). These were features applied to avoid overfitting of the model. A random seed value of 47 was utilized for the beginning of the training of the model. Three preprocessors were also utilized—FillMissing, Categorify and Normalize; FillMissing searches for missing values, Categorify finds the categorical variables and Normalize carries out a normalization for continuous variables within the dataset. The dependent variable ('dep_var') for this study was the year and the features used were grouped based on them being either categorical (cat_vars) or continuous variables (cont_vars). The model structure is shown in Table 2.

Table 2. The layer structure of the model utilized.

Layer	Output Shape	Number of Parameters
Embedding	3	12
Embedding	3	9
Embedding	5	40
Embedding	10	260
Embedding	5	40
Dropout	26	0
BatchNorm 1d	9	18
Linear	50	1800
ReLU	50	0
BatchNorm 1d	50	100
Dropout	50	0
Linear	10	510
ReLU	10	0
BatchNorm 1d	10	20
Dropout	10	0
Linear	1	11
ReLU	1	0
BatchNorm 1d	1	2
Dropout	1	0
Linear	1	2

Total Parameters: 2824, Total trainable parameters: 2824
 Optimized with 'torch.optim.adam.Adam', betas = (0.9, 0.99)
 Using true weight decay, Loss function: Flattened Loss

In the model, three dense layers were utilized with different dropout values for each of the dense layers. The activation function used for the layers is also depicted in the table. Within the model summary, it can be seen there are embedding layers for each categorical

variable and there are dropout layers provided to help reduce biases and variances as previously indicated. Additionally, all the parameters used for the model were trainable and there utilized. The data were split between training, testing and validation sets with 80% allocated for training and 20% for training and validation. It is critical to note that only the most important features as set out by the feature importance tool in Section 3.2 were utilized in the model. The model setup code parameters developed and used is further given in Appendix A. Once these parameters were set up, the learning rate finder was employed to suggest an appropriate learning rate and using this learning rate (lr) the model was trained with 80 epochs (training steps) and using a weight decay factor (wd) of 0.2 to help regularization and to prevent overfitting. All of these factors implemented using the 'learn.fit_one_cycle' function within the library. The final output of the model is a prediction of the year, which should be designated as the next point of maintenance intervention for a given road within the database. Any user using the model at the road agency would have to provide the same information for that road as used in this model configuration. However, as the data used are from easily available sources, input data pose an easier task than other methods or techniques.

3.4. Assessment of the Accuracy of the Model

For the metric evaluation of the model produced, the Root Mean Square Error (RMSE) metric was utilized. This is a common metric used for the evaluation of regression-based models and is estimated by the function given in Equation (1).

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (Y_{i,m} - Y_{i,e})^2} \quad (1)$$

where $Y_{i,m}$ = measured value, $Y_{i,e}$ = estimated value, n = number of observations.

To further explain the model's results accuracy, the RMSE was normalized based on the mean of the dependent variable, using a metric called Normalized Root Mean Square Error (NMRSE) to allow for a comparison across different scales and one that is more interpretable to users. These metrics are typical and representative of metrics utilized when assessing the accuracy of machine learning forecasting models [62,63]. This was generated using the function in Equation (2) below.

$$NMRSE = \frac{\sqrt{\frac{1}{n} \sum_{i=1}^n (Y_{i,m} - Y_{i,e})^2}}{\bar{Y}_{i,m}} \quad (2)$$

The evaluation was done utilizing the validation set during the model to establish the model's accuracy. Once this was completed, the study then tried to establish a pipeline on how the model could be used in a practical scenario combining other low-cost surveys methods. This pipeline and the results are further discussed in Sections 4–6.

4. Description of Case Study—Palermo, Italy

For the study, it is important to understand the particulars of the area. Palermo is located on the north-western coast of Sicily, Italy, covering an area of 158.88 km² with 663,401 inhabitants [64]. It is a part of southern Italy where municipalities have been shown to have a lower efficiency when considering the difference between the assessed spending needs given conditions and the actual spending [65]. This is an important factor as the efficiency level is associated with productivity levels. This is coupled with lower spending on information technology (IT) systems and services within the public sector [66] which limits the capacity of agencies to use highly technical systems. As a result of this, there is a great opportunity to use Palermo as a case study given the limited resources these authorities face and therefore the need to implement systems that do not require significant IT investment and resources.

The municipality's road network consists of approximately 3800 road axes with a total surface area equal to approx. 9 million m². This is a significantly large network and therefore having a system in place that could pinpoint particular areas of interest for intervention would be very useful. Data on the commute from the last population census carried out by ISTAT (Italian National Institute of Statistics) [67] show that there are approximately 278,954 individuals who commute every morning, for work and study purposes, of which 96.3% of the movements are carried out within the municipality. It is also important to consider that car sharing is very low with the average coefficient of vehicle occupancy being 1.3 and the highest mode of transport for these trips being private cars (48.2%) [68]. This is an important consideration as it means that the number of workers within the city will have a significant impact on the traffic levels and in turn the roads with higher debilitating conditions. The area was previously divided into 25 neighborhoods but this was subsequently amended in 2009 to a division of eight circumscriptions [68] as shown in Figure 4. There are overlaps between the divisions and, as the data collected have factors from both divisions, both are considered. Circumscription no.1 is particularly important as it houses the historical district of the city and the ZTL (Limited Traffic Zone) area, which has a traffic congestion charge attached to it and additionally is where a lot of tourism is concentrated given the presence of historical buildings, museums and artifacts. Circumscription no. 8 is also important given its central location and proximity to the port.

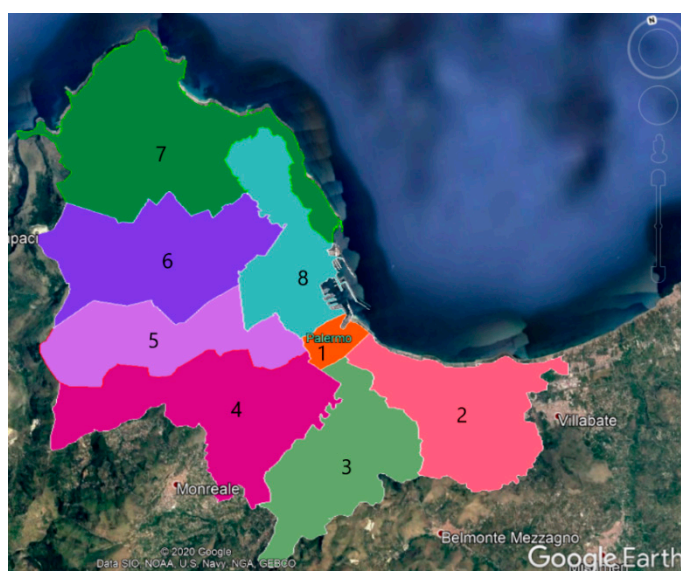


Figure 4. Circumscription divisions within Palermo (screenshot taken from Google Earth [69]).

The circumscription and neighborhood divisions are very important as population and industrial data were collected and aggregated according to these divisions. It is critical for any model development to understand the regional and geographical divisions within a region to adequately collate the data collected and, therefore, this should be carried out in any study of this nature. For traffic analysis, the area is divided into 200 zones ('PUT' traffic zones) across its area. This is particularly important as commute and commercial activity are also recorded based on the interactions between these 'traffic zones'.

Characteristics of Available Data

As previously mentioned, the data collected were obtained from government authorities, both locally and nationally. The data compiled from these sources resulted in a dataset of 1099 maintenance occurrences over 10 years with the recordings made in six time period groups across the ten years. The local company legally required to carry out these maintenance activities is currently Risorse Ambiente Palermo S.p.A. (RAP) [70]. The data sources are summarized in Table 3. It must be noted whilst the data were obtained

from several open source systems, it did involve multiple stakeholders and the research team had to collate the items across matching points for the final dataset. However, this is an unavoidable task to ensure that the characteristics of the area are understood and additionally it is a step that would only be needed to be done at the beginning of the process in a practical implementation scenario.

Table 3. Summary of Data Sources.

Data Group	Source	Data Obtained
Census data	ISTAT [67]	Population of Circumscription,
Traffic and commuting data	ISTAT, Municipality of Palermo [67,71,72]	Employment, commuting statistics, industrial activities, commercial activities
Economic activity data	Palermo Urban Transport Plan [68]	PUT zone activity rate, traffic rate, workers, arrivals and departures, Industrial activities in the zone
Historical maintenance and road network data	RAP and Municipality of Palermo [73]	Year of maintenance, road lengths and area

According to the legal requirements of the company, it has to carry out a minimum of 1.8% surface area intervention on the road network per year. For the planning of interventions, the company employs the application of filtering criterion that produces a list of roads to be maintained that are in the worst deterioration state and the process favors those roads considered critical with higher traffic intensity. To create this list of roads with grave severities, road inspectors carry out laborious monitoring activities throughout the year to update road conditions. They manually observe conditions, using video and physical surveys [74], and generate reports on the road conditions. Whilst this can be an effective method to collect data, it can also be expensive considering the number of surveys and personnel involved. Therefore, having a model that could assist in predicting which roads are likely to need rehabilitation would reduce costs. Additionally should physical surveys not be possible, the model could be utilized to produce a non-subjective indication of where interventions should be done.

The instances of intervention were recorded for six time intervals over the period and therefore within the data, the year group is recorded with numerical values with a range of 1–6 where the interventions occurring at year group 6, represent interventions at 10 years beyond the first recording. The recorded six groups are equally spaced out across the ten years meaning that each group represents one-sixth of 10 years, i.e., approximately 20 months. The roads within the list are those having a high severity rating and therefore requiring full surface depth repair over the area specified in the data. The characteristics of the roads on which these interventions were carried out were compiled and this represented the list of features to be analyzed. Unlike large databases like the LTPP database, where detailed technical information on the road condition and monitoring, such as PCI and IRI values, are available, the dataset focused on physical road characteristics, movements of cars and people and economic activities carried out on and near the roads. This issue is evident in many small countries and cities where there is a lack of funds to have extensive databases and therefore the challenge exists to produce reputable data analyses with limited data. The features produced are given in Table 4 and do provide a good representation of the specific situation considering the typical overview feature groups previously discussed in Section 1.2, allowing for an accurate portrayal of the area. The features produced were linked to the roads identified in the historical data by careful data matching.

Table 4. Description of features within the dataset.

Data Group.	Description of Data
neighborhood	neighborhood where the road is located
circ	Circumscription (circ.) where the road is located
circ_pop	population of circ. where the road is located
street_category	road category classification (labelled 1 or 2; where 1 represents the higher trafficked option)
length	road length (measured in meters)
area	road area covered (measured in sq. meters)
zone	zone where road is located (1-historic center, 2-main city, 3-peripheries)
year	year group of maintenance
pop_den	density of population in circ.
public_buildings	The presence or lack of public buildings near the road (labelled 1 for yes and 0 for no)
commercial_activities	The presence or lack of commercial activities along the road (labelled 1 for yes and 0 for no)
traffic_rate	rate of activities in traffic zone where the road is located
tz_pop	population of traffic zone where the road is located
tz_pop_den	population density of traffic zone where the road is located
tz_workers	number of workers in the traffic zone where the road is located
unemployment	percentage of unemployment in circ. where road is located
industrial_jobs	percentage of industrial jobs in circ.
circ_road_den	road density within circ. where road is located
t_arrivals	number of trip arrivals in the traffic zone where the road is located
tz_departures	number of departures in the traffic zone where the road is located
tz_perdays_rt	number of total trips made within the traffic zone where the road is located

The task was, therefore, to utilize these factors to provide analyses on road maintenance and predict the next instances of interventions. Consequently, using this workflow could provide an easier pathway for the road authority to plan interventions. Within the dataset, it is worth mentioning that a significant number of the features are categorical as they relate to the circumscription or the neighborhood in which the feature occurs.

5. Results and Discussion

5.1. Examination of Data

The first step was the uploading and examination of the data utilizing a pandas data frame. The first analytical step was a descriptive analysis of the features of the dataset. Within the preprocessing step, it was shown that there were no missing data entries, which is important to ensure there are no model inconsistencies later in the process. As previously mentioned there were 21 feature categories present in the dataset each carrying 1099 values corresponding to the maintenance activities carried out over the previous 10 years (2001–2019). This timeline thereby represents a significant period for analysis to ensure credibility in the process. To understand the distributions within the data, histograms were generated for the features to analyze the distributions of the features across the period. Of particular interest were the distributions of the features corresponding to the physical road characteristics (length and area) and the distributions for the occurrences of activities within the different zones, circumscriptions and neighborhoods over the 10 years. The length and area are directly related to each other with their distributions showing a direct similarity. The distribution for the length feature is shown in Figure 5. This showed that the majority of interventions were of a road length of less than 500 m.

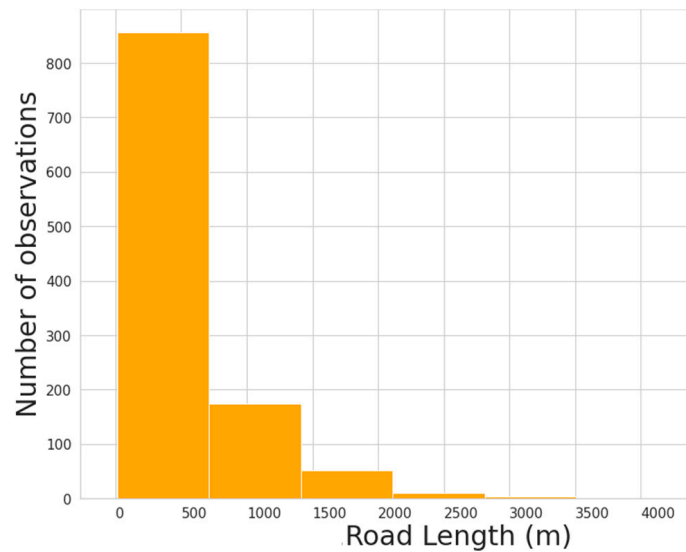


Figure 5. Distribution of feature—Road Length.

Subsequently, the distributions corresponding to the geographical divisions of Palermo were examined. The distribution corresponding to the zoning is shown in Figure 6 (where the zones are numbered as described in Table 4 with 1 representing the historical center, 2—the main city and 3 the peripheries). This highlights that most of the interventions did not take place within the historical center and instead were concentrated outside of the center and the central districts of the city. This was expected given the traffic levels within the center and a higher level of foot traffic that occurs within the historical district due to restrictions and congestion charging being in place. Furthermore, zone 3 is larger and therefore has more roads than the other two zones. This provides another key to understanding necessary intervention scheduling in the municipality.

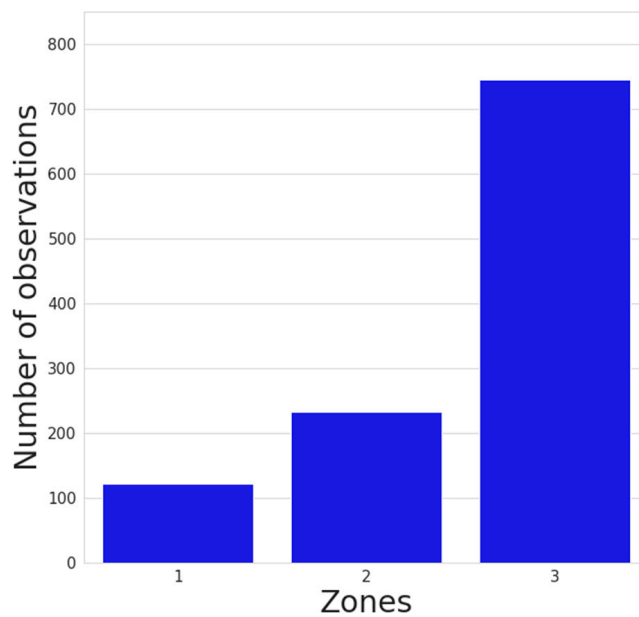


Figure 6. Distribution of interventions across zones.

Concerning the distributions of the district divisions (Figure 7), it was noted that the majority of interventions were recorded in circumscription No. 8 which is near the seaport as previously noted in Section 4. The most prevalent neighborhood was no.10 which is the

‘Politeama’ area, which has many stores and is highly trafficked therefore validating this statistic. This provides an insight into the priority locations of intervention. Otherwise, the number of interventions across the neighborhoods is fairly distributed. Whilst these are interesting results, it was also critical to identify correlations between interventions within these districts and the time of intervention to understand the sequences of interventions concerning not only their location but also time.

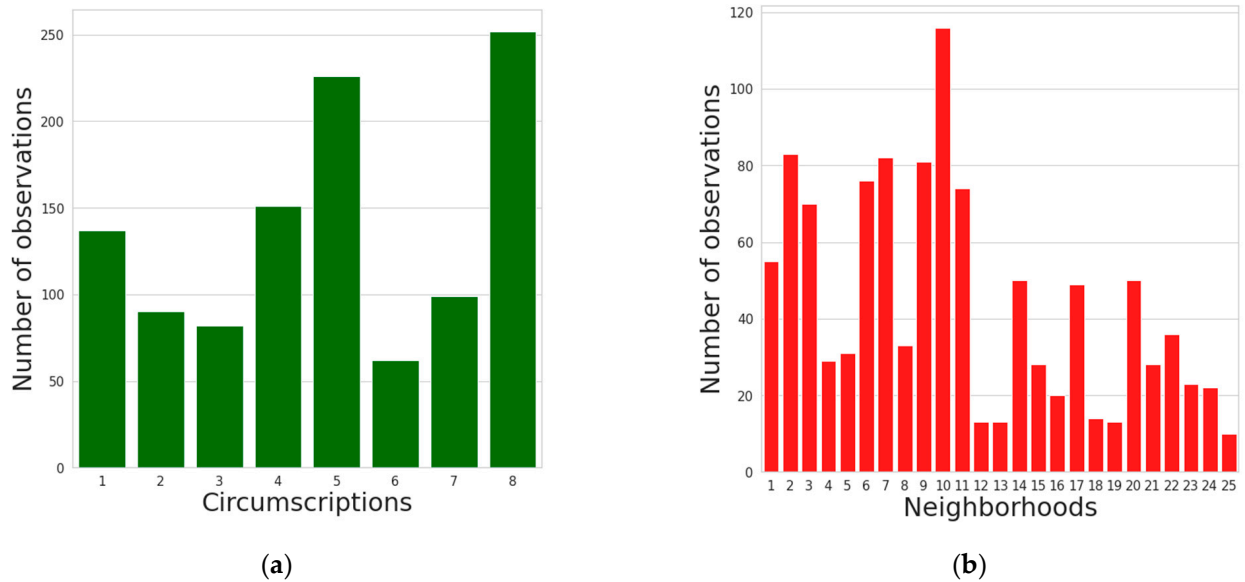


Figure 7. Distributions across (a) circumscriptions (left) and (b) neighborhoods (right).

To analyze the interventions by time, boxplots were used, within the pandas data frame, to relate the divisions to the year group of intervention. These boxplots are given in Figures 8 and 9.

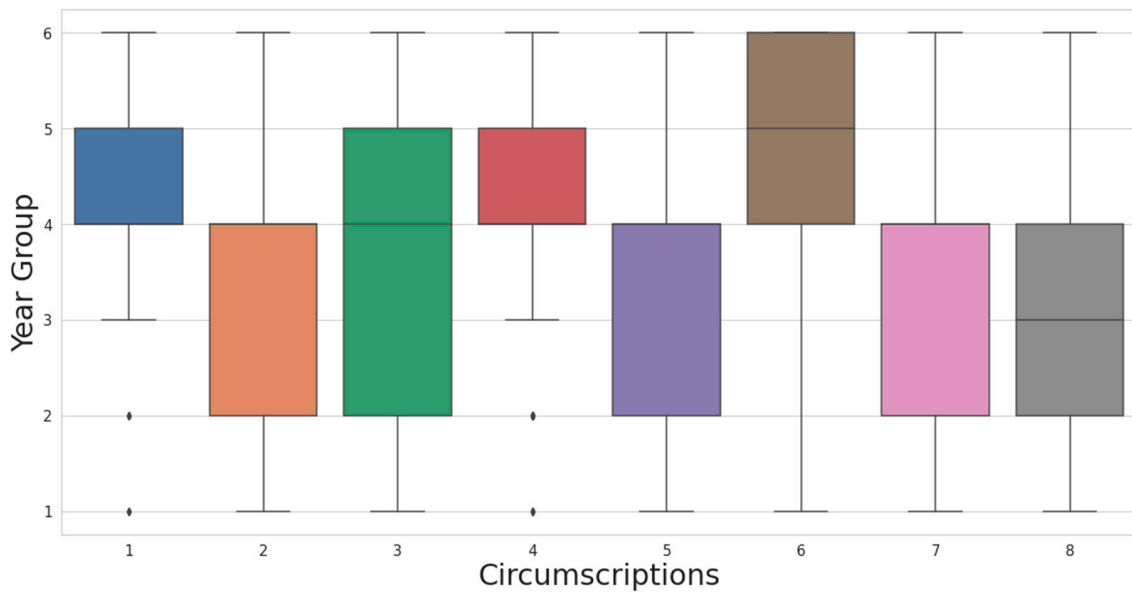


Figure 8. Boxplot of maintenance intervention in each circumscription with respect to time.

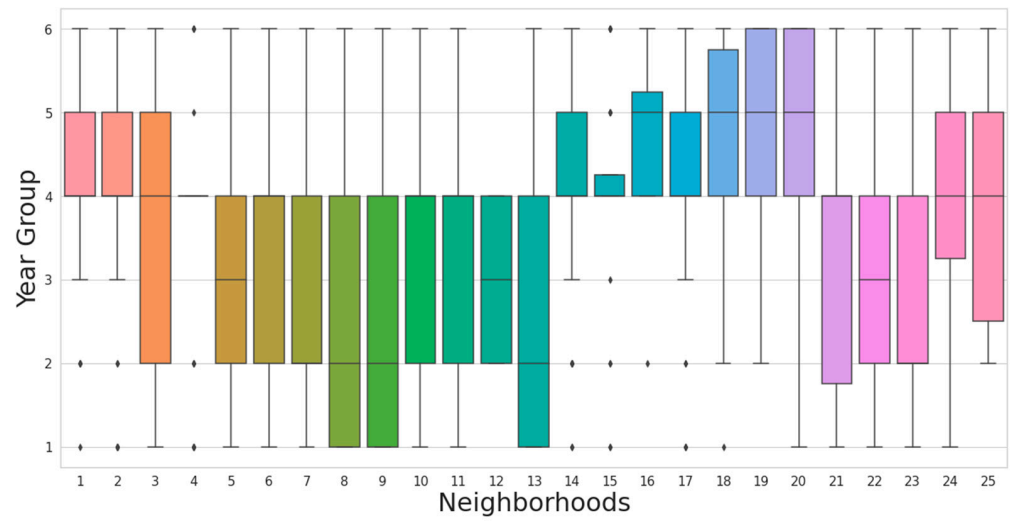


Figure 9. Boxplot of maintenance intervention in the neighbourhood with respect to time.

The plots allowed for analysis on when the interventions were taking place. For the circumscriptions, the interventions were split in time with the majority of interventions taking place within the first 4 year groups for circumscriptions 1, 5, 7 and 8. On the other hand, activities were concentrated in the later years for circumscriptions 2, 4 and 6; with circumscription 3 being the only one with a statistically significant overlap of activities across the period. This is important because it again gives an insight into time intervals that can be set up for monitoring activities within various districts based on the historical activities. It essentially demonstrates timelines of when neighborhoods and circumscriptions need attention and thus intervention.

The next important data analytic was the year of the intervention. A violin plot was constructed to highlight when in time the activities were being carried out. This is shown in Figure 10. Within the plot, it can be seen that the majority of interventions occurred in the middle at year group 4, which represents interventions at year 7. This gives the study another important discussion point as it showcases that the interventions are concentrated around this time and thus a larger number of resources could be allocated in the future for this period. The violin plot is used to easily highlight the peaks in the distribution.

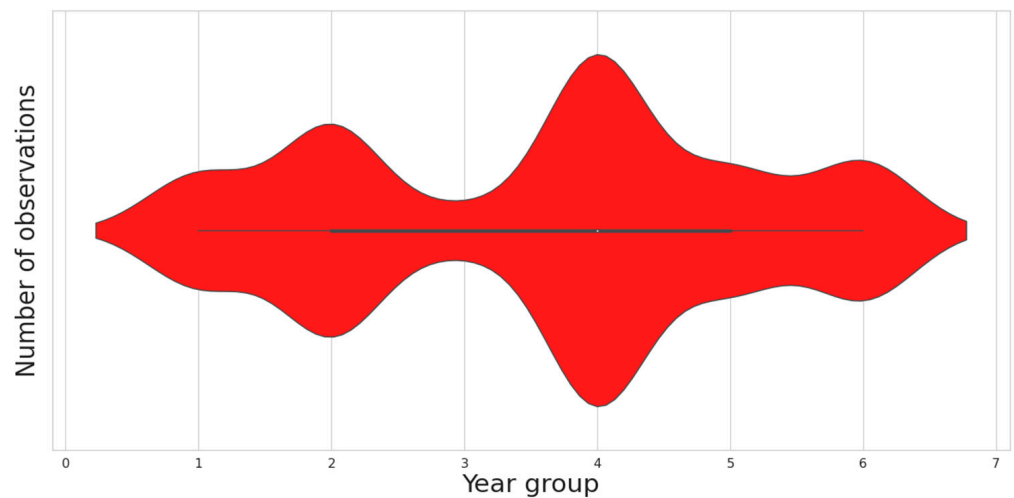


Figure 10. Violin plot of the distribution of interventions in time.

For the other features within the dataset, another interesting feature to highlight was the presence of commercial activities. This is important as commercial activities are likely to lead to higher traffic and general activity and therefore the roads in these areas are likely to suffer from more damage quicker. A boxplot of this feature versus time was plotted and is displayed in Figure 11. This indicates that the areas without commercial activities required maintenance at a later stage than those with the businesses. This was expected but it is important to confirm to validate the use of this data based on typical trends in traffic and road maintenance on pavement deterioration.

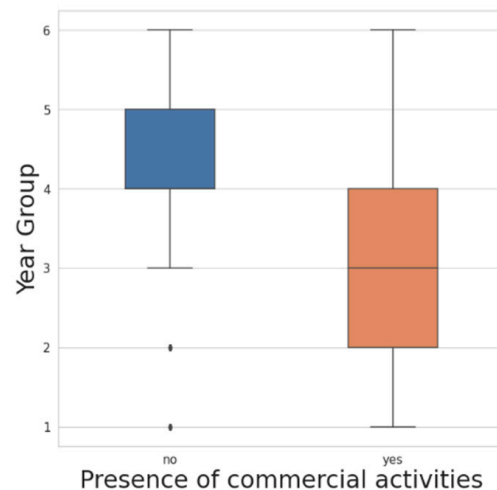


Figure 11. Distribution of maintenance interventions by the presence of commercial activities.

Once the initial examination of the data were completed, the next phase was the feature selection as it is more difficult to understand which of the other features should be utilized in a modelling framework using only a statistical examination of the distributions of features and this is where the use of the deep learning becomes key.

5.2. Feature Selection Results

With the first examination of the features complete, the feature selection algorithm as constructed in Section 3.2 was applied. The model was implemented in python and a normalized plot of feature importance was generated. Early stopping was applied using the validation set of the data to prevent overfitting. At iteration 408, the best result was achieved according to the model. The top 15 important ranked normalized features are given in Figure 12.

Within the figure, it can be seen that the most important feature is the zone and this can be further validated based on the data examined in Section 5.1 concerning the distribution of interventions within the zones. The next important feature is that of commercial activity and this again can be expected given the propensity of interventions being needed earlier when there are commercial activities as highlighted in Figure 11. The next feature on the list is that of the population density within the traffic zone showing how important it is to understand the density of people living near a particular road for the maintenance activities. Following these are the physical features of area and length, which definitely should be considered to be important as they represent information on the actual road's dimensions. It is important however to note that whilst the physical characteristics are important they are not the most important features, which is an important concluding result. The next series of features related to the traffic in the zone and on the type of persons living near a particular road.

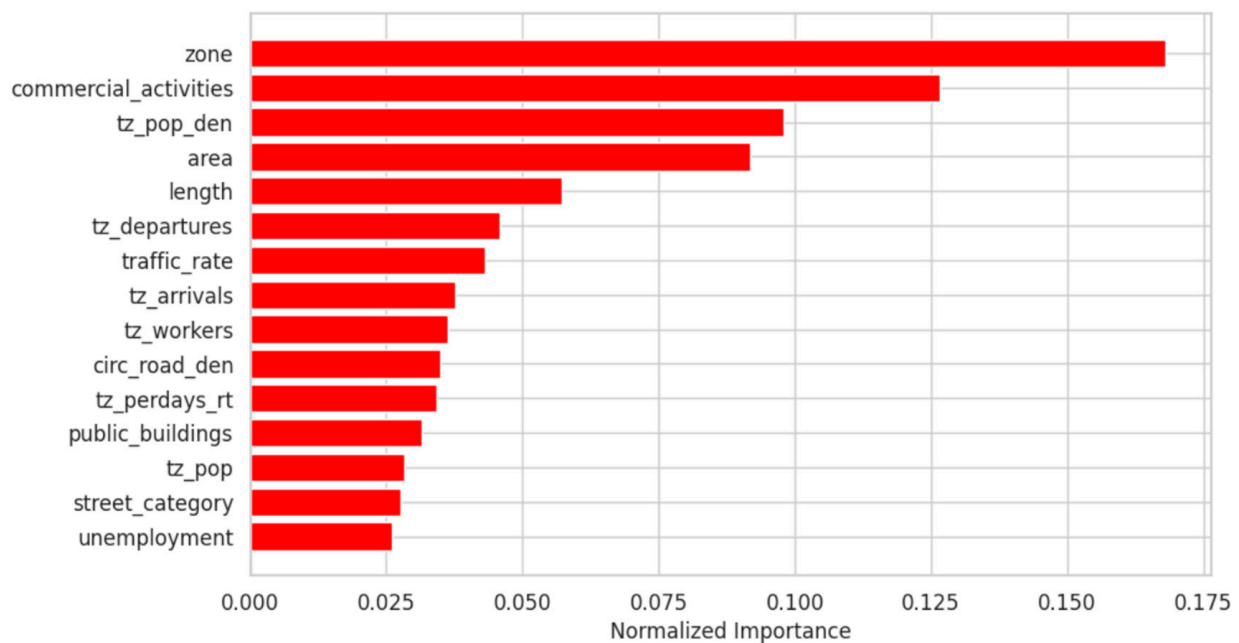


Figure 12. Ranked feature importance within the dataset.

It should be noted that these results are about the particular situation within the municipality under question and the same features could result in a different relationship in a different region. As such, it is critical to utilize results from the particular municipality and not rely on feature representations of other cities or municipalities. Therefore, what is important is the modelling process and also the ability to verify whether the results produced by the model can be explained considering the situation in real life within the area. The feature selection workflow is key to the work as it provides an ability to filter out the most important variables in attempting to make a prediction on future occurrences of the dependent variable. This feature selection can also be used on any database with more variables when available. This includes those that were not available in this study but could be available in other agencies or networks. This means the workflow is not isolated to this particular area. Once these features were computed, the next phase of the study was the use of these particular features within a model to predict maintenance interventions.

5.3. Fastai Model Results

The model was set up using the hyperparameters shown in Section 3.3 and it was run with 80 epochs (model steps) used to reduce variances. For the model, the first step was to utilize the learner within the FastAI library to identify the appropriate learning rate for the training. This was performed within the setup and was based on models carried out to determine appropriate learning rates for machine learning models [75]. This is visualized by the graph in Figure 13, showing the point corresponding to this learning rate to be utilized. The FastAI library has a built-in function that allows for programming a suggested value which can then be utilized directly within the model for more effective training. This process was done with the suggested training rate being highlighted as seen in Figure 13. From this process, a learning rate of approximately 0.02 was used for the model.

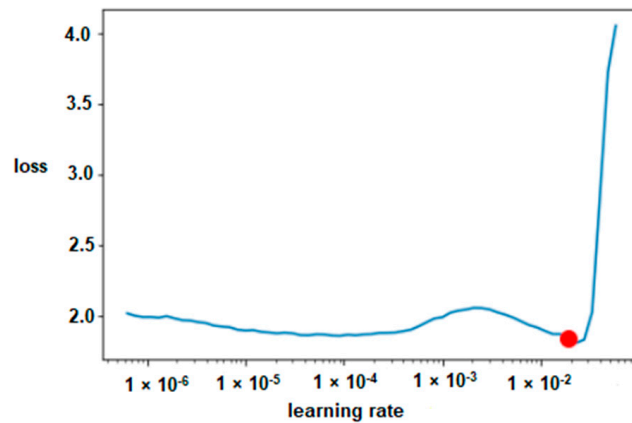


Figure 13. Utilization of Learning rate finder within FastAI library.

This learning rate was subsequently utilized in the training process. During the training, the validation loss was monitored against the training loss to ensure there was no overfitting. Overfitting can typically be assumed to be taking place if the validation loss is significantly higher than the training loss. This was monitored within the setup with the graph shown in Figure 14. Within the figure, it can be seen that the values stabilized by the end of the training run and the validation losses were not significantly higher than the training losses.

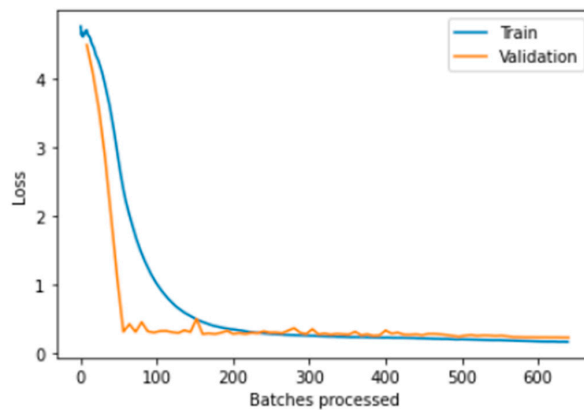


Figure 14. Loss observed during the training process.

The RMSE values were also recorded during the process as shown in Figure 15.

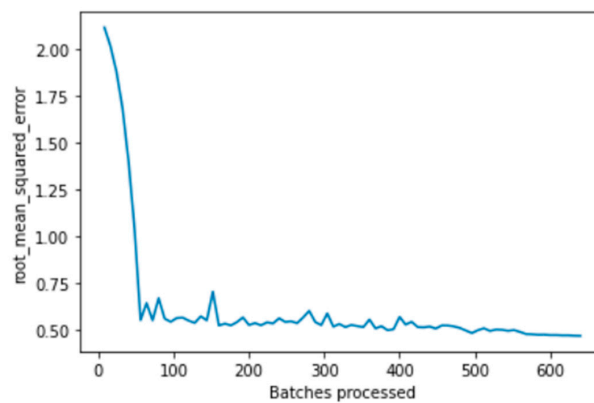


Figure 15. Observed training metric during training.

The values of RMSE were also stabilized by the end of the training. The RMSE, as identified in Section 3.4 was used as the metric for assessment of the accuracy of the model. After training was completed, the RMSE observed at the final epoch was 0.469. Using the values of RMSE, the NRMSE was computed using the average of the input dataset to normalize the value and make it interpretable on different scales and to provide clearer interpretation. This resulted in an NRMSE of 0.128, which therefore equates to an accuracy rate of 0.872, meaning the accuracy of the derived model is 87.2%. For comparison and validation of the feature selection, the model was run using all the features from Table 4 and this model achieved an NRMSE of 0.196 indicating an accuracy of 80% showcasing an increase of 7% in performance when using the feature selection, which could prove key in carrying out the correct maintenance at the correct time. It should also be noted that the increase could be more if more features are eliminated from a potentially larger database, eliminating those features that are not key to the dependent variable. Consequently, the model can accurately predict future maintenance sequences. In the model's application, new roads that have not undergone similar interventions would be put into the model to produce an idea of the timeline for intervention. The validation set used, which the model never saw during the training highlights the accuracy of the model. It must be stated that this is where deep learning has its strengths in its ability to search out deeper connections and patterns that are not easy to see and understand [76] and is why the general research has seen rapid growth and is considered more and more for engineering applications. Whilst it is important to mention that the model still shows an error, it is clear that the model is adequate and useful. This is in line with other studies, which have used different forms of artificial intelligence to predict maintenance intervention timelines and metrics. This is validated in recent studies which have achieved accuracies of: 87% when predicting the infrastructure's condition using a neural network with embedding [77], 86% for predicting road surface milling and overlays interventions [78], 85% when using combinations of gradient boosting trees to predict the PCI [79] and 87% predicting global road performance indicators [29]. Other studies have achieved higher accuracies when using the LTPP database [80] and other studies showing that the majority of research using ANN for predictions is done using this database [41]. This is expected given the size and resources used to create that database but does not assist the authorities without similar resources. Given the level of accuracy achieved from this study using the available data and resources, the result represents a substantial advance for the targeted types of road agencies.

Additionally, in scenarios where there are no models in place because of the lack of data, this would represent a substantial upgrade to their planning processes. With respect to time, the error would represent less than 1.5 months of time of a gap between when the model predicts maintenance should occur and when it does occur. Considering that plans are made on a yearly scale, this error can, therefore, be considered acceptable as opposed to the alternative where the activity is scheduled in a different year or not at all for a particular road without the model. For a graphical representation of the possible scheduling, a test dataset of 220 road instances was considered and a hypothetical prediction of the required maintenance intervention was plotted for the next 10-year period, grouped by the respective circumscriptions and is shown in Figure 16. This is a prediction plot so the prediction points appear between the year groups as shown in the figure.

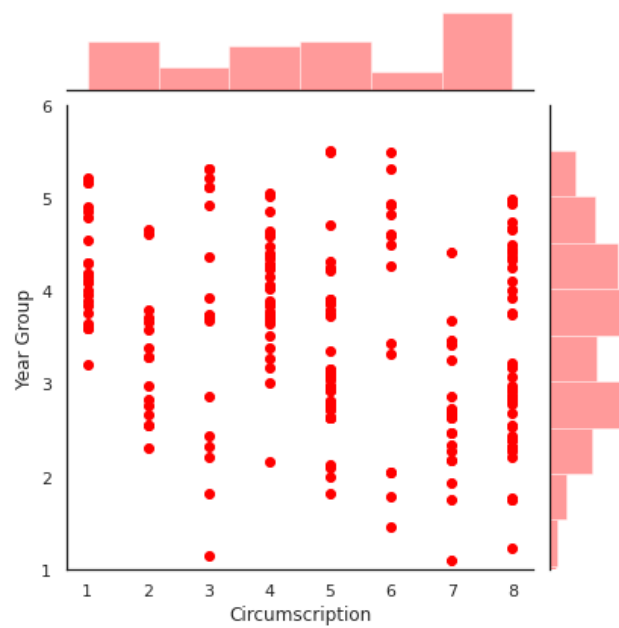


Figure 16. Jointplot of predicted future maintenance interventions grouped by circumscription.

This figure highlights the time points where the model has predicted that maintenance should occur and should be programmed over the next period of years. Within the figure, it is shown which circumscriptions are critical at which points in time. Using this, a plan of activities can easily be drawn up for the next period and budgets can be planned with a better idea of how much interventions are likely required in a particular year. Additionally, to provide a clearer picture of criticalities, the plot was further grouped to generate a heatmap of points in Figure 16 to understand the important periods and circumscriptions. This was generated using the seaborn library and is given in Figure 17.

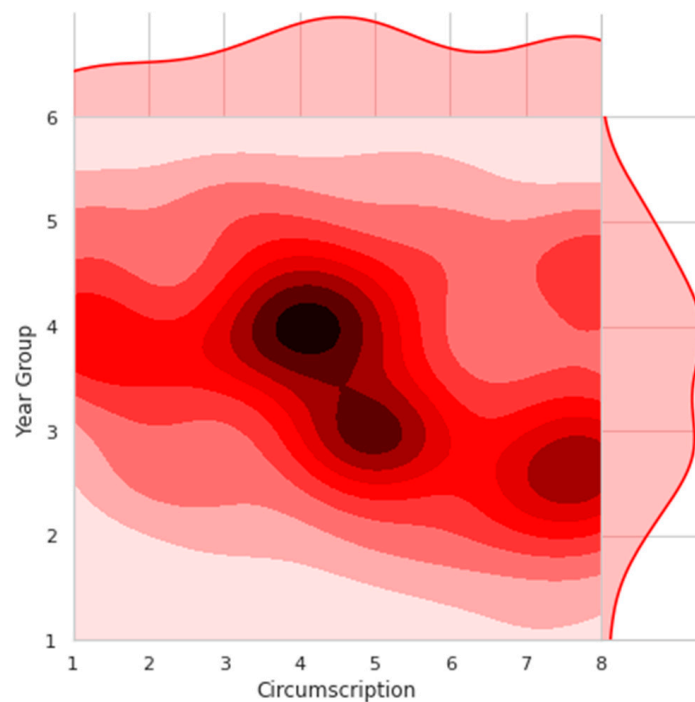


Figure 17. Heatmap of Jointplot of predicted interventions grouped by circumscription.

As expected from the discussions, the plot shows that the period where more action is required are the mid to later year periods but it also importantly highlights the circumstances where more intervention is needed, which is critical for planning and budget allocations. Whilst there is room for improvement in the model, it must be emphasized that given the size of the dataset and the resources available, it is a significantly useful result that could help better plan interventions within the municipality. It is also a pathway towards planning interventions without relying on excessive and expensive databases.

6. Workflow for Practical Implementation with Other Low-Cost Techniques

Given the results of Section 5, a pipeline was developed where the model and workflow utilizing small datasets can practically be integrated along with other low-cost methodologies for the assessment of pavement distresses in a road network. This pipeline is depicted in Figure 18, with the objective being to better plan decisions on maintenance interventions within the network.

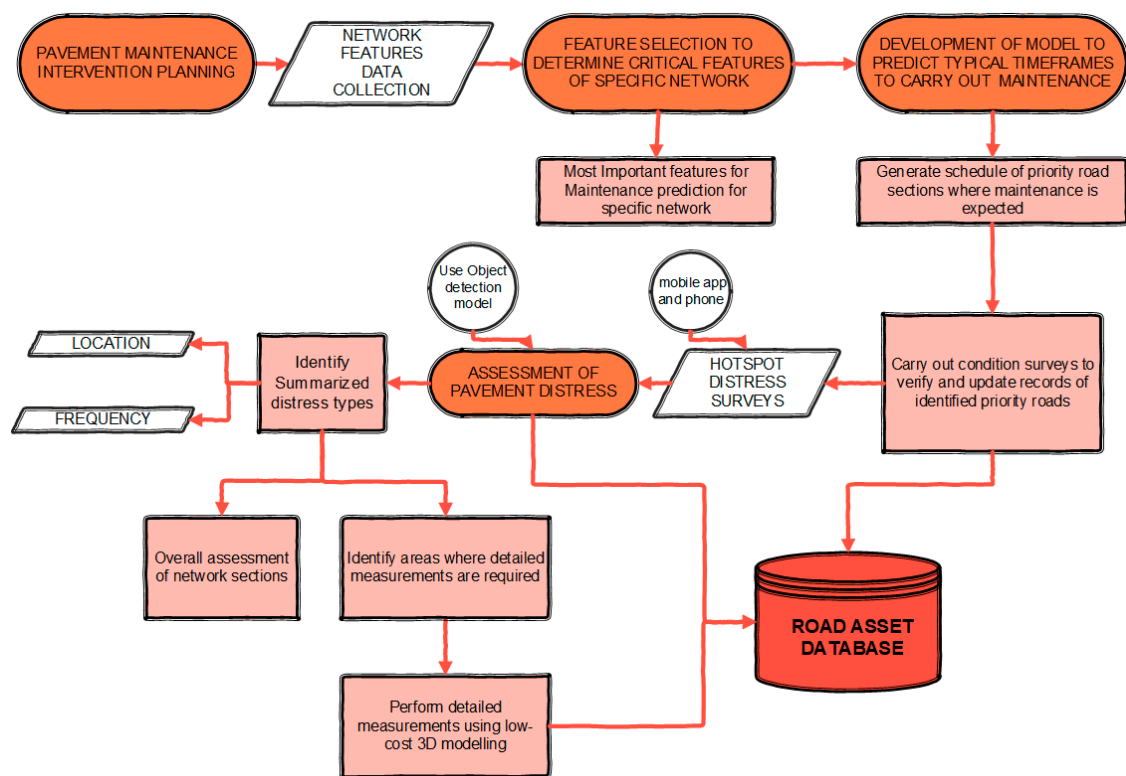


Figure 18. Workflow for the practical implementation of the methodology.

Within the workflow, the process begins with the collection of network data followed by the feature selection and model development as explained and portrayed in this study. Following this, the model can be utilized to generate a schedule of road sections where interventions are likely to be needed over time. The model could be run to predict where interventions would be scheduled over the next ten year period. It must be noted however that all of the features utilized in this study would have to be found for each of the roads in the network before the model is run. Whilst this still requires some data collection, it would be significantly easier to do so than having to physically survey the entire network. The schedule produced as a result of this would indicate the particular geographical areas where interventions are likely needed and the roads that are allocated to a particular year group of intervention. To validate this schedule, conditional surveys can then be done on the sections highlighted by the model. Whilst these surveys can be costly, the pipeline utilizes low-cost surveys using object detection models embedded on smartphones to detect the presence of pavement distresses. This would subsequently produce a hotspot

analysis of the roads covered in the surveys as identified in an earlier study [39]. This survey would yield information about the location and frequency of the distresses and be able to validate the work carried out by the model within this study predicting which roads need intervention. The surveys would be run on the roads as deemed necessary from the results of the model in this study.

Following this, detailed measurements can be determined, where necessary, with low-cost 3D modelling techniques using smartphones and drones to produce metric models of the distressed areas as validated by previous work [81]. It is significant to note that each step of this workflow utilizes a low-cost technique or system to not only acquire the information needed but to process it. This is vital to help road authorities with their under-budgeted scenarios whilst still maintaining a high level of accuracy and adequate planning practices. At each stage of the pipeline, the data and analyses would be fed towards the road asset database to allow for updated records of road assets and continual low-cost network monitoring by the road authority. The workflow and the combination of low-cost measures within it are what offers the greatest contribution to the work.

7. Conclusions

This study was designed to present a flexible framework towards utilizing easily available information on a city and its road network for planning road maintenance interventions. The purpose was to create a workflow that could be replicated by small and under-resourced road authorities who cannot create and/or access large databases but still have to create effective rehabilitation plans. The analyses and models were developed relying on data analytics tools, open source algorithms and deep learning models.

For the workflow's validation, a case study in Palermo, Italy was utilized. Within the study, several open source datasets that are easily available were utilized from the census, road traffic and maintenance history. From the information, feature selection tools were developed to identify which of the features are the most important towards explaining the point in time at which intervention activities need to be done. The feature selection tools were based on gradient boosting algorithms that adequately can handle the presence of categorical variables, which commonly exist in these types of datasets. Once the feature selection tool was developed, it was used to pinpoint the most important features, which were then utilized in a deep learning model to help predict the intervention time for a particular road. The most important features were analyzed and validated based on analyses of feature distributions and the specific surroundings. Once the features were integrated into the deep learning model, a model with high accuracy was achieved. It provides a very good assessment of when interventions should occur. The model is also non-contact and does not require excessive physical surveys. This process, therefore, utilizes the strengths of gradient boosting algorithms and deep learning, leveraging their combined power to handle categorical variables and search out connections on a deeper level than traditional approaches.

Future work will look at different cities and the use of different feature sets to further establish the accuracy of this data exploitation technique. The use of the process in other cities will further validate the workflow along with checks in the following years to verify the accuracy of the next sequence of predicted road interventions. The predicted models in the study can aptly be used to generate schedules of which roads and sections should be prioritized over time. Further to the results of the study, a pipeline was also identified as to how the models and techniques could be used with other low-cost techniques for detecting pavement distresses and planning pavement interventions. The workflow helps exploit several low-cost element techniques, which would be helpful for many agencies given their budgetary allocations. Finally, it is worth pointing out this study exploits both low-cost investigation techniques on road pavement conditions and indirect data analytics, to establish over time and space, the best network maintenance strategy. This would allow more budget allocation towards effective interventions rather than on complex and

expensive systematic investigations. With the techniques combined, they offer another step towards low-cost automation of elements of the PMS.

Author Contributions: Conceptualization, R.R.; methodology, R.R., and G.D.M.; software, R.R.; validation, R.R., L.I. and G.D.M.; formal analysis, R.R. and G.D.M.; investigation, R.R. and G.D.M.; resources, G.D.M.; data curation, G.D.M.; writing—original draft preparation, R.R.; writing—review and editing, R.R., L.I. and G.D.M.; visualization, R.R. and G.D.M. All authors have read and agreed to the published version of the manuscript.

Funding: The research presented in this paper was carried out as part of the SMARTI ETN project under the H2020-MSCA-ETN-2016. This project has received funding from the European Union’s H2020 Programme for research, technological development and demonstration under grant agreement number 721493.

Data Availability Statement: Data used for the study can be found in references given in Table 3.

Acknowledgments: The authors would like to acknowledge and thank Comune di Palermo for assisting with the data collection and explanation of datasets within their portfolio concerning pavement maintenance, network configurations and census information. The authors would also like to thank Roberto Biondo who was integral in collecting the data from the municipality.

Conflicts of Interest: The authors declare no conflict of interest. The funders had no role in the design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript, or in the decision to publish the results.

Appendix A. Fastai Setup

Parameter setups:

```
dep_var = 'year'
cat_vars = ['zone', 'commercial_activities', 'unemployment', 'street_category', 'public_buildings',
'circ_road_den']
cont_vars = ['area', 'length', 'traffic_rate', 'tz_pop', 'tz_pop_den', 'tz_workers', 'tz_arrivals',
'tz_departures', 'tz_perdays_rt']
procs = [FillMissing, Categorify, Normalize]
test = TabularList.from_df(test_df, cat_names = cat_names, cont_names = cont_vars,
procs = procs)
valid =
TabularList.from_df(val_df, cat_names = cat_vars, cont_names = cont_vars, procs = procs)
data
= (TabularList.from_df(train_df, path = '.', cat_names = cat_vars, cont_names = cont_vars,
procs = procs).split_by_rand_pct(valid_pct = 0.2, seed = 47)
.label_from_df(cols = dep_var, label_cls = FloatList, log = True)
.add_test(test)
.databunch())
```

Model setup:

```
learn = tabular_learner(data, layers = [1,10,50], ps = [0.01,0.01,0.1], metrics = rmse, emb_
drop = 0.05, callback_fns = ShowGraph)
Learning rate set up: learn.fit_one_cycle(80, lr, wd = 0.2)
```

References

1. Vandam, T.J.; Harvey, J.T.; Muench, S.T.; Smith, K.D.; Snyder, M.B.; Al-Qadi, I.L.; Ozer, H.; Meijer, J.; Ram, P.V.; Roesier, J.R.; et al. *Towards Sustainable Pavement Systems: A Reference Document FHWA-HIF-15-002*; Federal Highway Administration: Washington, DC, USA, 2015.
2. *Eurostat Energy, Transport and Environment Statistics*, 2019 ed.; European Union: Brussels, Belgium, 2019; ISBN 9789276109716.
3. Karleuša, B.; Dragičević, N.; Deluka-Tibljaš, A. Review of multicriteria-analysis methods application in decision making about transport infrastructure. *J. Croat. Assoc. Civ. Eng.* **2013**, *65*, 619–631. [\[CrossRef\]](#)
4. International Road Federation (IRF). *IRF World Road Statistics 2018 (Data 2011–2016)*; IRF: Brussels, Belgium, 2018.
5. Mbara, T.C.; Nyarirangwe, M.; Mukwashi, T. Challenges of raising road maintenance funds in developing countries: An analysis of road tolling in Zimbabwe. *J. Transp. Supply Chain Manag.* **2010**, *4*, 151–175. [\[CrossRef\]](#)

6. Fernandes, N. Economic effects of coronavirus outbreak (COVID-19) on the world economy. *SSRN Electron. J.* **2020**. [[CrossRef](#)]
7. Inzerillo, L.; Di Mino, G.; Roberts, R. Image-based 3D reconstruction using traditional and UAV datasets for analysis of road pavement distress. *Autom. Constr.* **2018**, *96*, 457–469. [[CrossRef](#)]
8. Peterson, D. *National Cooperative Highway Research Program Synthesis of Highway Practice Pavement Management Practices. No. 135*; Transportation Research Board: Washington, DC, USA, 1987; ISBN 0309044197.
9. American Association of State Highway and Transportation Officials (AASHTO). *Pavement Management Guide*; AASHTO: Washington, DC, USA, 2012.
10. Haas, R.; Hudson, W.R.; Falls, L.C. *Pavement Asset Management*; Wiley: Hoboken, NJ, USA, 2015. [[CrossRef](#)]
11. Roberts, R.; Inzerillo, L.; Di Mino, G. Using UAV based 3D modelling to provide smart monitoring of road pavement conditions. *Information* **2020**, *11*, 568. [[CrossRef](#)]
12. Schnebele, E.; Tanyu, B.F.; Cervone, G.; Waters, N.M. Review of remote sensing methodologies for pavement management and assessment. *Eur. Transp. Res. Rev.* **2015**, *7*, 1–19. [[CrossRef](#)]
13. Amador, L.E.; Magnuson, S. Adjacency modeling for coordination of investments in infrastructure asset management. *Transp. Res. Rec. J. Transp. Res. Board* **2011**, *2246*, 8–15. [[CrossRef](#)]
14. Radopoulou, S.C.; Brilakis, I. Improving road asset condition monitoring. *Transp. Res. Proc.* **2016**, *14*, 3004–3012. [[CrossRef](#)]
15. Mallela, S.S.J.; Lockwood, S. National Cooperative Highway Research Program. In *Transportation Research Board Strategic Issues Facing Transportation, Volume 7: Preservation, Maintenance, and Renewal of Highway Infrastructure*; The National Academies Press: Washington, DC, USA, 2020. [[CrossRef](#)]
16. Paterson, W.D.O.; Scullion, T. *Information Systems for Road Management: Draft Guidelines on System Design and Data Issues*; The World Bank: Washington, DC, USA, 1990.
17. Bennett, C.R.; Chamorro, A.; Chen, C.; De Solminihaç, H.; Flintsch, G.W. *Data Collection Technologies for Road Management*; The World Bank: Washington, DC, USA, 2007.
18. Singh, A.P.; Sharma, A.; Mishra, R.; Wagle, M.; Sarkar, A. Pavement condition assessment using soft computing techniques. *Int. J. Pavement Res. Technol.* **2018**, *11*, 564–581. [[CrossRef](#)]
19. Zimmerman, K.A. *Pavement Management Methodologies to Select Projects and Recommend Preservation Treatments*; Transportation Research Board: Washington, DC, USA, 1995; p. 102.
20. Swei, O.; Gregory, J.; Kirchain, R. Pavement management systems: Opportunities to improve the current frameworks. In *Proceedings of the Transportation Research Board 95th Annual Meeting*; Transportation Research Board: Washington, DC, USA, 2016.
21. Haas, R.; Felio, G.; Lounis, Z.; Falls, L.C. Measurable performance indicators for roads: Canadian and international practice. In *Proceedings of the Annual Conference of Transportation Association of Canada Best Practices in Urban Transportation Planning, Measuring Change*, Vancouver, BC, Canada, 18–21 October 2009.
22. Humplick, F.; Paterson, W. Framework of performance indicators for managing road infrastructure and pavements. In *Proceedings of the 3rd International Conference on Managing Pavements*; National Academy Press: Washington, DC, USA, 1994; pp. 123–133.
23. Gupta, A.; Kumar, P.; Rastogi, R. Critical review of flexible pavement performance models. *KSCE J. Civ. Eng.* **2013**, *18*, 142–148. [[CrossRef](#)]
24. Sundin, S.; Braban-Ledoux, C. Artificial intelligence–Based decision support technologies in pavement management. *Comput. Civ. Infrastruct. Eng.* **2001**, *16*, 143–157. [[CrossRef](#)]
25. American Society for Testing and Materials (ASTM). *ASTM D 6433-18 Standard Practice for Roads and Parking Lots Pavement Condition Index Surveys*; ASTM International: West Conshohocken, PA, USA, 2018. [[CrossRef](#)]
26. Piryonesi, S.M.; El-Diraby, T. *Using Data Analytics for Cost-Effective Prediction of Road Conditions: Case of the Pavement Condition Index*; Federal Highway Administration: McLean, VA, USA, 2018.
27. Paterson, W. International roughness index: Relationship to other measures of roughness and riding quality. *Transp. Res. Rec. J. Transp. Res. Board* **1986**, *1084*, 49–59.
28. Gong, H.; Sun, Y.; Shu, X.; Huang, B. Use of random forests regression for predicting IRI of asphalt pavements. *Constr. Build. Mater.* **2018**, *189*, 890–897. [[CrossRef](#)]
29. Domitrović, J.; Dragovan, H.; Rukavina, T.; Dimter, S. Application of an artificial neural network in pavement management system. *Teh. Vjesn. Tech. Gaz.* **2018**, *25*, 466–473. [[CrossRef](#)]
30. Attoh-Okine, N.O. Analysis of learning rate and momentum term in backpropagation neural network algorithm trained to predict pavement performance. *Adv. Eng. Softw.* **1999**, *30*, 291–302. [[CrossRef](#)]
31. Elbagalati, O.; Elseifi, M.A.; Gaspard, K.; Zhang, Z. Development of an enhanced decision-making tool for pavement management using a neural network pattern-recognition algorithm. *J. Transp. Eng. Part B Pavements* **2018**, *144*, 04018018. [[CrossRef](#)]
32. Di Mino, G.; De Blasiis, M.; Di Noto, F.; Noto, S. An advanced pavement management system based on a genetic algorithm for a motorway network. In *Proceedings of the 3rd Conference on Soft Computing Technology in Civil, Structural and Environmental Engineering*, Cagliari, Italy, 3–6 September 2013. [[CrossRef](#)]
33. Bosurgi, G.; Trifiro, F. A model based on artificial neural networks and genetic algorithms for pavement maintenance management. *Int. J. Pavement Eng.* **2005**, *6*, 201–209. [[CrossRef](#)]
34. Santos, J.; Ferreira, A.; Flintsch, G. An adaptive hybrid genetic algorithm for pavement management. *Int. J. Pavement Eng.* **2017**, *20*, 266–286. [[CrossRef](#)]

35. Samek, W.; Wiegand, T.; Müller, K.R. Explainable artificial intelligence: Understanding, visualizing and interpreting deep learning models. *arXiv* **2017**, arXiv:1708.08296.
36. Cai, J.; Luo, J.; Wang, S.; Yang, S. Feature selection in machine learning: A new perspective. *Neurocomputing* **2018**, *300*, 70–79. [[CrossRef](#)]
37. Pantelias, A.; Flintsch, G.W.; Bryant, J.W.; Chen, C. Asset management data practices for supporting project selection decisions. *Public Work. Manag. Policy* **2008**, *13*, 239–252. [[CrossRef](#)]
38. Al Qurishee, M.; Wu, W.; Atolagbe, B.; Owino, J.; Fomunung, I.; Onyango, M. Creating a dataset to boost civil engineering deep learning research and application. *Engineering* **2020**, *12*, 151–165. [[CrossRef](#)]
39. Roberts, R.; Giancontieri, G.; Inzerillo, L.; Di Mino, G. Towards low-cost pavement condition health monitoring and analysis using deep learning. *Appl. Sci.* **2020**, *10*, 319. [[CrossRef](#)]
40. Federal Highway Administration LTPP InfoPave. Available online: <https://infopave.fhwa.dot.gov/> (accessed on 14 April 2020).
41. Bashar, M.Z.; Torres-Machi, C. Performance of machine learning algorithms in predicting the pavement international roughness index. *Transp. Res. Rec. J. Transp. Res. Board* **2021**. [[CrossRef](#)]
42. Marcelino, P.; Antunes, M.D.L.; Fortunato, E. Comprehensive performance indicators for road pavement condition assessment. *Struct. Infrastruct. Eng.* **2018**, *14*, 1433–1445. [[CrossRef](#)]
43. *The Handbook of Highway Engineering*; Fwa, T.F. (Ed.) CRC Press: Boca Raton, FL, USA, 2006; ISBN 9780849319860.
44. McKinney, W. Data Structures for Statistical Computing in Python. In Proceedings of the 9th Python in Science Conference, Austin, TX, USA, 9–15 July 2010; pp. 56–61. [[CrossRef](#)]
45. Hunter, J. Matplotlib: a 2D Graphics Environment. *Comput. Sci. Eng.* **2007**, *9*, 90–95. [[CrossRef](#)]
46. Waskom, M.; Gelbart, M.; Botvinnik, O.; Ostblom, J.; Hobson, P.; Lukauskas, S.; Gemperline, D.C.; Augspurger, T.; Halchenko, Y.; Warmenhoven, J.; et al. mwaskom/Seaborn: v0.11.1. Available online: [mwaskom/seaborn](http://mwaskom.com/seaborn) (accessed on 30 April 2020). [[CrossRef](#)]
47. Sandru, E.-D.; David, E. Unified feature selection and hyperparameter bayesian optimization for machine learning based regression. In Proceedings of the 2019 International Symposium on Signals, Circuits and Systems (ISSCS), Iasi, Romania, 11–12 July 2019; pp. 1–5. [[CrossRef](#)]
48. Koehrsen, W. *Feature-Selector 1.0.0*; Github online program, 2019. Available online: github.com/Jie-Yuan/FeatureSelector (accessed on 30 April 2020).
49. Prokhorenkova, L.; Gusev, G.; Vorobev, A.; Dorogush, A.V.; Gulin, A. Catboost: Unbiased boosting with categorical features. In *Proceedings of the 32nd Conference on Neural Information Processing Systems, Montreal, Canada, 3–8 December 2018 (NeurIPS 2018)*; Curran Associates Inc.: Montreal, QC, Canada, 2018; pp. 6638–6648.
50. Zhang, Y.; Haghani, A. A gradient boosting method to improve travel time prediction. *Transp. Res. Part C Emerg. Technol.* **2015**, *58*, 308–324. [[CrossRef](#)]
51. Friedman, J.H. Greedy function approximation: A gradient boosting machine. *Ann. Stat.* **2001**, *29*, 1189–1232. [[CrossRef](#)]
52. Natekin, A.; Knoll, A. Gradient boosting machines, a tutorial. *Front. Neurobot.* **2013**, *7*, 21. [[CrossRef](#)]
53. Chen, T.; Guestrin, C. XGBoost: A Scalable tree boosting system. In *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*; Association for Computing Machinery: San Francisco, CA, USA, 2016; pp. 785–794. [[CrossRef](#)]
54. Freund, Y.; Schapire, R.E. A decision-theoretic generalization of on-line learning and an application to boosting BT—Computational learning theory. In *Proceedings of the Second European Conference on Computational Learning Theory*; Springer: Barcelona, Spain, 1995; Volume 904, pp. 23–37. [[CrossRef](#)]
55. 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 of the 31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA, 4–7 December 2017*; Curran Associates Inc.: Long Beach, CA, USA, 2017; Volume 2017, pp. 3147–3155.
56. Al Daoud, E. Comparison between XGBoost, LightGBM and CatBoost using a home credit dataset. *Int. J. Comput. Inf. Eng.* **2019**, *13*, 6–10.
57. Jhaveri, S.; Khedkar, I.; Kantharia, Y.; Jaswal, S. Success Prediction Using Random Forest, CatBoost, XGBoost and AdaBoost for Kickstarter Campaigns. In Proceedings of the 3rd International Conference on Computing Methodologies and Communication, ICCMC, Erode, India, 27–29 March 2019; IEEE: Erode, India, 2019; pp. 1170–1173. [[CrossRef](#)]
58. Huang, G.; Wu, L.; Ma, X.; Zhang, W.; Fan, J.; Yu, X.; Zeng, W.; Zhou, H. Evaluation of CatBoost method for prediction of reference evapotranspiration in humid regions. *J. Hydrol.* **2019**, *574*, 1029–1041. [[CrossRef](#)]
59. Dorogush, A.V.; Ershov, V.; Gulin, A. CatBoost: Gradient boosting with categorical features support. *arXiv* **2018**, arXiv:1810.11363.
60. Howard, J.; Gugger, S. Fastai: A layered API for deep learning. *Information* **2020**, *11*, 108. [[CrossRef](#)]
61. Guo, C.; Berkhahn, F. Entity embeddings of categorical variables. *arXiv* **2016**, arXiv:1604.06737.
62. Chen, D.; Mastin, N. Sigmoidal models for predicting pavement performance conditions. *J. Perform. Constr. Facil.* **2016**, *30*, 04015078. [[CrossRef](#)]
63. Fan, J.; Wu, L.; Zhang, F.; Cai, H.; Zeng, W.; Wang, X.; Zou, H. Empirical and machine learning models for predicting daily global solar radiation from sunshine duration: A review and case study in China. *Renew. Sustain. Energy Rev.* **2019**, *100*, 186–212. [[CrossRef](#)]
64. ISTAT. Istat Italy Resident Population 2020. Available online: <http://dati.istat.it/Index.aspx?QueryId=18460&lang=en> (accessed on 28 April 2020).

65. OECD. *OECD Economic Surveys: Italy 2009*; OECD Publishing: Paris, France, 2019.
66. OECD. *Tax Administration 2017: Comparative Information on OECD and Other Advanced and Emerging Economies*; OECD Publishing: Paris, France, 2018.
67. Istituto Nazionale di Statistica—ISTAT. Permanent Census—Italy. 2011. Available online: <http://dati-censimentopopolazione.istat.it/Index.aspx?lang=en> (accessed on 29 April 2020).
68. *Città di Palermo—Ufficio Traffico ed Authority Piano Generale del Traffico Urbano*; Ufficio Traffico ed Authority: Palermo, Italy, 2010.
69. Google Earth Pro v7.3.2.5776 38°07′18.69″ N, 13°19′42.81″ E, Eye alt 19.55 mi. SIO, NOAA, U.S. Navy, NGA, GEBCO. Available online: <http://www.earth.google.com> (accessed on 28 March 2020).
70. Risorse Ambiente Palermo (RAP). *Carta dei Servizi—Edizione 2019*; Risorse Ambiente Palermo: Palermo, Italy, 2019.
71. *Città di Palermo PANORMUS—Annuario di Statistica del Comune di Palermo 2014*; Comune di Palermo: Palermo, Italy, 2014.
72. *Città di Palermo Servizio Trasporto Pubblico di Massa e Piano Urbano del Traffico Piano Urbano della Mobilità Sostenibile Quadro Conoscitivo*; Comune di Palermo: Palermo, Italy, 2019.
73. Comune di Palermo Portale Open Data. Available online: <https://opendata.comune.palermo.it/opendata-ultimi-dataset.php> (accessed on 28 March 2020).
74. Risorse Ambiente Palermo (RAP). *Piano Industriale 2019–2021*; Risorse Ambiente Palermo: Palermo, Italy, 2019; p. 249.
75. Smith, L.N. Cyclical learning rates for training neural networks. In Proceedings of the 2017 IEEE Winter Conference on Applications of Computer Vision (WACV), Santa Rosa, CA, USA, 24–31 March 2017; pp. 464–472. [CrossRef]
76. Khanafer, M.; Shirmohammadi, S. Applied AI in instrumentation and measurement: The deep learning revolution. *IEEE Instrum. Meas. Mag.* **2020**, *23*, 10–17. [CrossRef]
77. Bukhsh, Z.A.; Stipanovic, I.; Saeed, A.; Doree, A.G. Maintenance intervention predictions using entity-embedding neural networks. *Autom. Constr.* **2020**, *116*, 103202. [CrossRef]
78. Morales, F.J.; Reyes, A.; Caceres, N.; Romero, L.M.; Benitez, F.G.; Morgado, J.; Duarte, E. A machine learning methodology to predict alerts and maintenance interventions in roads. *Road Mater. Pavement Des.* **2020**, 1–22. [CrossRef]
79. Piryonesi, S.M.; El-Diraby, T.E. Data analytics in asset management: Cost-effective prediction of the pavement condition index. *J. Infrastruct. Syst.* **2020**, *26*, 04019036. [CrossRef]
80. Elhadidy, A.A.; El-Badawy, S.M.; Elbeltagi, E.E. A simplified pavement condition index regression model for pavement evaluation. *Int. J. Pavement Eng.* **2019**, 1–10. [CrossRef]
81. Roberts, R.; Inzerillo, L.; Di Mino, G. Exploiting low-cost 3D imagery for the purposes of detecting and analyzing pavement distresses. *Infrastructures* **2020**, *5*, 6. [CrossRef]