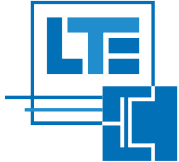


*Presented at the Session of July 2016*

---

**Technische Universität München**



**LEHRSTUHL FÜR TECHNISCHE ELEKTRONIK**

Technische Universität München - Arcisstraße 21 - 80333 München

Tel.: 089/289-22929 - Fax: 089/289-22938 - Email: lte@ei.tum.de

Prof. Dr. Doris Schmitt-Landsiedel



## **Master's Thesis Report**

Topic:

**LONG-TERM TRAVEL TIME PREDICTION**

Carried out by:

**Syed Murtaza Hassan**

Host Organization:

**NEC Laboratories Europe**

**University Advisor:** Prof. Dr. rer. nat. Doris Schmitt-Landsiedel

**Company Supervisor:** Dr. Luis Moreira-Matias

## *Abstract*

Long-term travel time predictions are crucial for tactical and operational public transport planning in schedule design and resource allocation tasks. Similarly to any regression task, its success considerably depend on an adequate feature selection framework. In this project, we approach the myopia of the State-of-the-Art method RReliefF on mining relevant inter-relationships of the feature space relevant for reducing the entropy around the target variable on regression tasks. A comparative study was conducted using baseline regression methods and LASSO as a valid alternative to RReliefF. Experimental results obtained on a real-world case study uncovered the bias/variance reduction obtained by each approach, pointing out promising ideas on this research line.

**Keywords :** travel time prediction, machine learning, regression, feature selection

*To my parents, and siblings*

# Acknowledgements

I would like to express my gratitude and respect to my supervisor, Dr. Luis Moreira-Matias, for the tremendous opportunity he gave me to work within his team. It was a great learning experience to be part of such a highly regarded team of research scientists. I'm extremely thankful for his guidance and and continuous help. I'm also grateful to my colleague Jihed Khiari who helped me immensely throughout.

I would also like to take the opportunity to thank my teachers at the Technische Universität München for considerably contributing to my professional and academic development as well as my supervisor Prof. Dr. rer. nat. Doris Schmitt-Landsiedel.

# Achievement

My research paper "Feature Selection Issues in Long-term Travel Time Prediction" got published in IDA 2016 - The 15th International Symposium on Intelligent Data Analysis (an ERA A-rank conference on Computer Science/Data Science fields).

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Motivation . . . . .	2
1.2	Tasks . . . . .	3
<b>2</b>	<b>Overview</b>	<b>5</b>
2.1	Definitions . . . . .	5
2.1.1	Data Sources . . . . .	5
2.1.2	Machine learning . . . . .	6
2.1.3	Supervised Learning . . . . .	6
2.1.3.1	Regression . . . . .	7
2.1.3.2	Bias – Variance . . . . .	7
2.1.3.3	Bias . . . . .	7
2.1.3.4	Variance . . . . .	7
2.1.4	Unsupervised Learning . . . . .	8
2.1.4.1	Clustering . . . . .	8
2.1.5	Feature Selection . . . . .	8
2.2	General Context of project . . . . .	9
<b>3</b>	<b>Case Study</b>	<b>11</b>
3.1	Available Data . . . . .	11
3.2	Pre-processing . . . . .	11
3.3	Feature Generation . . . . .	12

<b>4</b>	<b>Methodology</b>	<b>14</b>
4.1	Feature Selection Methods . . . . .	15
4.1.1	RReliefF . . . . .	15
4.1.2	LASSO . . . . .	15
4.2	Regression Methods . . . . .	16
4.2.1	Random Forests . . . . .	16
4.2.2	Support Vector Machines . . . . .	16
4.2.3	Projection Pursuit Regression . . . . .	17
<b>5</b>	<b>Implementation</b>	<b>18</b>
5.1	Environment . . . . .	18
5.1.1	Hardware Environment . . . . .	18
5.1.2	Software Environment . . . . .	18
5.2	Experiments . . . . .	19
5.2.1	Hyperparameter Tuning . . . . .	19
<b>6</b>	<b>Results and Discussion</b>	<b>21</b>
6.1	Results . . . . .	21
6.2	Discussion . . . . .	21
<b>7</b>	<b>Conclusion and Future Work</b>	<b>25</b>

# List of Figures

1	Clustering results of vehicle ids for route A2. . . . .	13
2	LASSO results for all routes. A vertical red dashed line is drawn at the best $\log \lambda$ value. This serves as cut.off point. . . . .	23
3	RReliefF results for all routes. A horizontal red line is drawn at $y=0.01$ . . . . .	23
4	RReliefF and LASSO comparative analysis (y-axis) using RMSE (scaled on RRelief side) along different scheduled departure times (x-axis). Bars denote the sample size on each timespan (scaled on LASSO side). . . . .	24



# List of Tables

1	Statistics per Route. The values are as mean $\pm$ s.d.. Times in seconds. . .	12
2	Optimal Hyperparameters setting. . . . .	22
3	SVR results for initial, LASSO and RrF-RReliefF feature sets. . . . .	24
4	RF and PPR results for initial, LASSO and RrF-RReliefF feature sets. . . .	24

# Introduction

## 1.1 Motivation

The rapid urbanization and population growth has led to wider transportation networks, which increases the need for efficient schedule planning. The guarantee of having a safe, fast, comfortable and affordable way to traversing a city from one point to another is an asset that local governments and agencies want to achieve. However, it is common to note a lack of adherence to the schedule which is inconvenient for both the passengers and the operators. That's why, it is more relevant than ever for public transportation agencies to improve their reliability.

One of the most common research problems in transportation is travel time prediction (TTP). The literature on this topic is extensive and covers different application domains such as fleet management, monitoring, control, mass transit and individual navigation [1]. Hereby, we focus on public transport in general and buses in particular. It is possible to distinguish short and long-term travel time prediction problem based on the prediction horizon (e.g. threshold of 2-3 hours). Operational tasks (e.g. timetable design) or resource allocation (e.g. vehicle and crew scheduling) requires long-term TTP.

Regression is one of the most common approach for travel time prediction. Regression models comprise of a large number of techniques to estimate the relationship between a set of predictors and a dependent variable. All of the previous studies have

estimated parametric regression models using the least squared methods.

In transportation science, it is known that the main determinants of bus running times are route length, passenger activity at stops and the number of traffic signals (e.g. [2,3]). Other studies also added driver response to the deviation from the schedule as an explanatory variable [4,5]. However, all of those have estimated linear regression models to identify the impact of potential explanatory variables on bus running times. Consequently, the resulting models often have very limited predictive power.

Attaining better bus travel time predictions can have significant consequences for passenger delays, operator's performance fines and the efficiency of its resource allocation. The inherently complex and uncertain operational environment in which urban bus service operate call for the development of more sophisticated models that can capture non-linear relations between system variables. To the authors' best knowledge, the literature to handle this specific issue is scarce. Mendes-Moreira *et al.* [6] compared Random Forests (RF), Support Vector Machine Regression (SVR) and Progression Pursuit Regression (PPR). On the other hand, the well-known RRelief [7] was proposed to do an adequate feature selection for each route. As many other methods from the RELIEF\*-family, RRelief is an instance-based learning method which leverages on the concept of neighborhood to define features that can (or cannot) contribute significantly to the entropy reduction on estimating the target variable  $Y$ . Consequently, as many other instance-based methods (e.g.  $k$ -nearest neighbors), it is highly dependent on an adequate setting of a distance metric that serves this specific purpose (which can easily vary from problem to problem). Moreover, it also has limitations on evaluating inter-relationships among the feature set  $X$  which can lead to this effect.

## 1.2 Tasks

As the main contribution, we apply feature selection, on the three baseline regressors (SVR, PPR and RF), using the L1 regularization method Least Absolute Shrinkage and Selection Operator (LASSO) and we show that it achieves a better fit for this problem than the state-of-the-art method RRelief.

The project tasks are organized as follows:

- An overview of the project.
- Description of case study and related data sources.
- Feature selection methods' elaboration and limitations of RReliefF.
- Description of baseline regressors (SVR, PPR, RF).
- Implementation of the project; the software and hardware requirement and the experiments performed.
- Result report and a comprehensive discussion demonstrating LASSO performs better than RReliefF.
- Conclusion with final remarks and future research direction.

## Overview

In this chapter, first we will introduce some useful definitions and concepts that shed light on the project and then we will give a general context of the project.

### 2.1 Definitions

We define in this section the key concepts necessary to put our project in its context.

#### 2.1.1 Data Sources

The main data sources in the context of public transportation are AVL and APC.

**Automatic Vehicle Location (AVL)** data determines the exact location of vehicles in real-time based on GPS. This data is transmitted to a data server where it is stored for later preprocessing and mining. Given the availability and accuracy of this data, many transport operators have equipped their fleets with on-board transmitters to collect it. The AVL datasets reference the trips made by the vehicle typically contain information about the actual arrival/departure times from the different stops of the route, but also the scheduled arrival/departure times, and the dwell times. Thus, the round-trip times and the link travel times are directly deduced and stored in the dataset.

**Automatic Passenger Counting (APC)** data determines for every (stop,trip) the number of boarding/alighting passengers. Consequently, it expresses the passenger

demand on a particular route which is valuable source of information for operators. However, the literature shows most operators rely on AVL data for the planning and evaluation of their service. That's why, mining APC data alongside AVL data, presents an opportunity to gain more insight and have a greater impact on the schedule planning.

### 2.1.2 Machine learning

**Machine Learning** [8] is a subfield of Computer Science devoted to the development of algorithms able to learn dependences as well as to generalize behavioral patterns from data in a fully autonomous way. An explanatory **model** is then built as result of the application of such ML algorithms to data. This data is known as **training data**. Machine learning algorithms find patterns in data and make predictions or decisions which can later be used for analyzing new data.

The two main branches of ML are Supervised and Unsupervised Learning

### 2.1.3 Supervised Learning

**Supervised Learning** [9] Algorithms learn by giving a set of inputs and their respective output/label. For example a data set can have input data points with their associated outputs like if input1 and input2 =1, then output = spam, similarly if input1 and input2 = 2, then output = not-spam. The models are trained using actual output. Preparation of model is done using a training process where predictions are made and corrected (when those predictions are wrong). The training process goes on until the model achieves a desired level of accuracy on the training data.

Supervised learning algorithms are utilized for predicting future outcomes. As in our case of travel time prediction where we provide data (historical data) for training the models.

Regression is an example of supervised learning.

### 2.1.3.1 Regression

**Regression** [10] is a method of finding relationship between dependent variables (input) and independent variables (output) or in other words regression is an approach describing relationship between set of **predictors (features)** and outcome. Regression helps us to understand how the dependent variable (outcome) is changing with the independent variable (input). The target (outcome) is a function of input variables. This function is called the regression function or **regressor**. Regression is used for prediction. Using the known values of independent variables (input), regression is used for predicting the output. Narrowly speaking, regression is used for estimation of continuous output variables. In our project, travel time is a continuous variable and we will apply regression methods to predict travel time using large data from input variables.

### 2.1.3.2 Bias – Variance

In machine learning models, there are two main sources of prediction errors i-e error due to **bias**, error due to **variance** [11]. Understanding these two is critical in understanding the performance of machine learning algorithms. The model leading to a high variance issue is said to have an **overfitting** problem and model leading to high bias is said to have an **underfitting** problem. A good model is the one which decreases both sources of error.

### 2.1.3.3 Bias

Error due to bias occurs when the learning model is too simple failing to find relationship between predictors and output. One solution of reducing the bias error is to increase model complexity (add features).

### 2.1.3.4 Variance

Error due to variance occurs when the learning model is too complex and it sort of memorizes the data given in training set but this memorization leads to bad results on

test set. One solution of reducing the variance error is to decrease model complexity (reduce features).

## 2.1.4 Unsupervised Learning

**Unsupervised Learning** [12] algorithms learn by giving a set of inputs without giving any outcome/label. Algorithms learn by figuring out structures in data and form. For example algorithm forms groups/clusters/classes of data by finding structure in data.

Classification is an example of unsupervised learning.

### 2.1.4.1 Clustering

**Clustering** [13] is a technique that segments data into different groups. These groups are known as clusters. Data objects in one cluster are more similar each other than to data objects in other cluster.

## 2.1.5 Feature Selection

Automatic **Feature Selection** [14] is a subfield of study focused on developing algorithms capable of defining adequate feature spaces for supervised learning problems. The idea is to find the feature subset that guarantees solutions (i.e. models) close to the global minima of our generalization error by defining which features to use and which to drop on a particular regression/classification problem. There are mainly two types of feature selection algorithms: (i) **filters**, where the induction model is not take into account to select an adequate feature subset and (ii) **wrappers**, where the feature subset selection process takes into account the induction model (typically through an encapsulated optimization framework). In this project, we are focused on discussing issues around this topic (i), as well as its impact in the context of long-term TTP tasks.



## 2.2 General Context of project

Nowadays, most public transportation companies use Automatic Vehicle Location (AVL)/Automatic Passenger Counting (APC) systems to track the services provided by each vehicle. This data can be used to improve the Schedule Planning. This project aims to develop a Machine Learning framework and features extraction system for long term travel time prediction.

A traditional approach to TTP is regression analysis. It comprises a large number of techniques to estimate the relationship between a set of predictors (i.e features) and a dependent variable:

$$\hat{f} : x_i, \theta \rightarrow \mathbb{R} \text{ such that } \hat{f}(x, \theta) = f(x_i) = y_i, \forall x_i \in X, y_i \in Y \quad (2.1)$$

where  $f(x_i)$  denotes the true unknown function which is generating the samples' target variable and  $\hat{f}(x_i, \theta) = \hat{y}_i$  be an approximation dependent on the feature vector  $x_i$  and an unknown parameter vector  $\theta \in \mathbb{R}^n$  (given by some induction model  $M$ ). Notoriously, this approximation will be as good as the adequacy of  $M$  to the dependence structure of  $f$  as well as the relevancy of the input feature space  $X$ . If it has a low number of features, it may not explain the variance of  $Y$ , thus leading  $M$  to biased models. Conversely, for a large set of features, we may be using features with a low predictive power. In consequence,  $M$  may output very complex models which lead to optimal fits on the input dataset (i.e. *local minima*) but a considerably lower ones when tested in any generic inference task. These phenomenons are known as *underfitting* and *overfitting*, respectively as we discussed earlier too.

This project is focused on studying the effects of RReliefF *myopia* to unrealistic distance functions and/or interrelationship on the feature set relevant for predicting the target variable value. To do it so, we propose an the Least Absolute Shrinkage and Selection Operator (LASSO) as a simple and yet valid alternative to RReliefF for this particular domain. The idea is to leverage on the priority that LASSO gives on the bias error reduction - in contrast to RReliefF. Consequently, our contributions are twofold: (1) a practical demonstration on RReliefF limitations through the study of its impact on particular application area; (2) the introduction of LASSO as a valid alternative to

this problem due to the high number of relevant interactions among different predictors/features that can reduce bias error. Experimental results of applying the same baseline predictors to a particular real-world case study uncovered the potential of our novel approach.

## Case Study

### 3.1 Available Data

Our case study is a large urban bus operator in Sweden. We collected data collected on four high-frequency (maximum planned headway of 11 minutes between 7:00-19:00) routes A1/A2/B1/B2, i.e. two bus lines A/B. Line A connects residential areas to a public transport interchange hub as well as major shopping areas. B connects the southern parts of the city to the city center, traversing through an interchange, major hospitals as well as a logistic center. This study covers six months between August 2011 and January 2013. The study period transcends over both summer schedule (June 19-December 14) and winter schedule (December 15-June 18).

### 3.2 Pre-processing

Before applying our framework, it was necessary to perform a number of preprocessing tasks on the AVL data. The APC data was used as is. In fact, the AVL data presents some missing entries. For instance, some trips in the original datasets do not cover all the stops of the routes. As part of the preprocessing step, a trip pruning was performed by removing trips where more than 80% of link travel times were missing. In addition, we performed data imputation on the remaining samples by following the interpolation procedure suggested in [15]. The dwell times were also pruned by using

Table 1 – Statistics per Route. The values are as mean  $\pm$  s.d.. Times in seconds.

	<b>NTrips</b>	<b>Stops</b>	<b>Daily Trips</b>	<b>Round Trip Times</b>
<b>A1</b>	17953	33	134 $\pm$ 27	3017 $\pm$ 425
<b>A2</b>	16353	33	133 $\pm$ 30	2755 $\pm$ 480
<b>B1</b>	16280	25	137 $\pm$ 23	2607 $\pm$ 465
<b>B2</b>	16353	25	134 $\pm$ 22	2746 $\pm$ 448

the 99% percentile to remove erroneous measurements. Table 1 presents an overview of the resulting dataset, detailed per route. It contains the (i) total number of trips (NT), (ii) number of stops and (iii) Round Trip Times (RTT).

### 3.3 Feature Generation

The original features are schedule departure time, daytype and vehicle ID. Unlike RF, SVR and PPR do not support categorical values. Therefore, it is required to generate new features based on the original ones. For the type of day, we use one-hot encoding which generates 7 numerical features corresponding to the day type. Vehicle ids associated with less than 0.5% of total number of trips were grouped into a single cluster. The remaining vehicle ids were clustered using a clustering technique described in the experiments section. This procedure resulted in four additional features.

Fig. 1 illustrates the clustering results for route A2. It illustrates the clustering plot (top-left) and the kernel density estimations for the vehicle ids within each cluster. We used the Bayesian Information Criterion (BIC) to determine the best number of clusters  $k = 3$  from the interval  $K = [2 : 20]$ . We note that the three clusters are characterized by slightly different p.d.f.. This justifies mapping the ids into three distinct features. Since driver rosters are typically assigned to individual vehicles throughout their shift, vehicle travel times reflect driving style as well as the propagation of delays from one trip to later ones.

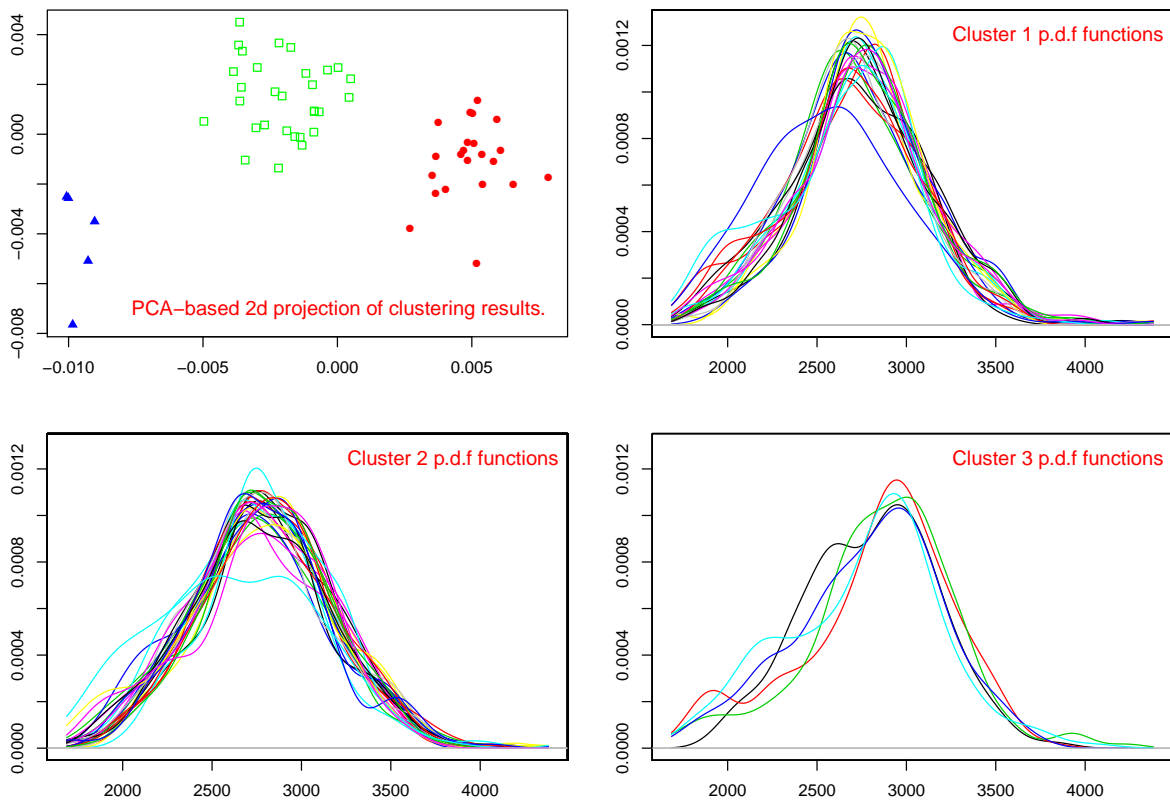


Figure 1 – Clustering results of vehicle ids for route A2.

# Methodology

Feature selection consists in eliminating redundant or non-informative features. Applying feature selection can not only lead to more interpretable models but also attain better results. Redundant features can negatively affect the predictions of models that do not automatically implement feature selection. This is also relevant for our TTP framework, where we seek to determine the best set of features for generating predictions. The state-of-the-art method for this domain (proposed by Mendes-Moreira *et al.* [6] is RreliefF [7]. Initially proposed for binary classification problems, this instance-based learning algorithm is able to determine features relevance on determining the target variable value. It can handle interdependences on the feature space, missing data and/or different type of functional forms for the dependences. However, its success depends largely on an adequate definition of a distance metric. Moreover, it is focused on reducing variance-type error, neglecting the inter-relationships that can potentially reduce the bias-type one.

Hereby, we compare RReliefF to LASSO as filter feature selection method to highlight why the first is not adequate for this task on long-term TTP problems. This section elaborates formally on the two methods as well on the three baseline regressors used to do it so.

## 4.1 Feature Selection Methods

### 4.1.1 RReliefF

was introduced by Kira and Kendell [16]. Its key idea is to rank features based on how well they separate classes. Given a randomly selected instance  $R_i$ , this is achieved by searching for its two nearest neighbors, one from the same class called nearest hit  $H$  and the other from the different class called nearest miss  $M$ . Since a good feature separates different classes, it should have a small distance to  $H$  and a large distance to  $M$ . The estimate of feature  $A$  quality  $W[A]$  is adjusted accordingly. The whole process is repeated for  $m$  iterations, where  $m$  is a user defined parameter. Finally, features that have a higher value than a given threshold  $\phi$  are selected. Similarly, the ReliefF algorithm [17] deals with classification problems with more than two classes, by considering  $k$  hits and misses rather than two. In regression problems, the predicted value is continuous so we cannot determine if two instances are part of the same class or not. To solve this issue, Robnik-Šikonja and Kononenko [7] introduced RreliefF: a probability measure modeled with the relative distance between the predicted values of the two instances. Similarly to ReliefF, a random instance  $R_i$  and its  $k$  nearest instances are selected in order to iteratively calculate the weights of input variables based on an user-defined distance metric.

### 4.1.2 LASSO

is a shrinkage and selection method for linear regression introduced by Tibshirani [18]. Similarly to other shrinkage methods, it aims to improve the least-squares estimator by adding constraints on the value of coefficients noted  $b$ . Given  $p$  vectors  $x$  of size  $N$ , where  $N$  is the number of samples and  $p$  is the number of features and an outcome  $y$ , the LASSO estimate is defined by

$$\hat{b}^{lasso} = \underset{b}{\operatorname{argmin}} \sum_{i=1}^N (y_i - b_0 - \sum_{j=1}^p x_{i,j} b_j)^2 \quad (4.1)$$

subject to

$$\sum_{j=1}^p |b_j| \leq t, t \geq 0 \quad (4.2)$$

The equivalent *Lagrangian form* is

$$\hat{b}^{lasso} = \underset{b}{\operatorname{argmin}} \left\{ \frac{1}{2} \sum_{i=1}^N (y_i - b_0 - \sum_{j=1}^p x_{i,j} b_j)^2 + \lambda \sum_{j=1}^p |b_j| \right\} \quad (4.3)$$

The  $L_1$ -norm penalty of LASSO  $\sum_{j=1}^p |b_j|$  allows a stronger constraint on the coefficients and causes some coefficients to be shrunk exactly to zero. The tuning parameter  $\lambda$  controls the **strength** of the penalty. As it increases, more coefficients are set to zero and hence, less variables are selected.  $\lambda$  is typically set by using a cross-validation search technique over a grid of admissible values.

## 4.2 Regression Methods

RreliefF and LASSO were tested as filter-type feature selection methods to cope with three baseline regressors: RF, SVR and PPR.

### 4.2.1 Random Forests

is an ensemble method based on classification and regression trees (CART [19]) that was introduced by Leo Breiman in 2001 [20]. The trees are grown by randomly choosing a set of candidate predictors at every node for a sample of the data and then producing the split by choosing the best splitter available. RF combines this with a random selection of samples to train the trees which is referred to as bootstrap aggregating or bagging. RF's hyperparameters are (i) the number of randomly selected predictors to choose from at each split *mtry* and the number of grown trees *ntree*.

### 4.2.2 Support Vector Machines

were introduced by Cortes *et al.* in 1995 [21]. They are primarily binary classifiers that perform their task by constructing hyperplanes in a multidimensional space able to separate instances either linear or non-linearly. In  $\epsilon$ -SVM, these hyperplanes are constructed in a way to ensure the largest minimum distance to the training examples. This distance ( $\epsilon$ ) is denominated as *margin*. SVMs can be adapted for regression with



a quantitative response by sequentially optimizing an error function where we seek to maximize the geometrical distance  $\frac{1}{\|w\|}$  which is equivalent to minimizing  $\frac{1}{2}\|w\|^2$ . To allow examples to be in the margin or to be misclassified, slack variables  $\xi_i \geq 0$  are introduced. The optimization problem becomes:

$$\arg \min_{w,b} \frac{\|w\|^2}{2} + C \times \sum_{i=1}^n \xi_i \quad (4.4)$$

where  $C > 0$  is a constant that sets the relative importance of maximizing the margin and minimizing the amount of slack. Kernels are typically used in SVMs to map the data points into higher dimensional feature space. Typical kernel include polynomial and radial basis functions. The choice of the kernel depends on the problem and different functions may depend on different hyperparameters.

### 4.2.3 Projection Pursuit Regression

is an additive model that consists of linear combinations of non-linear transformations of linear combinations of explanatory variables (so-called *ridge functions*) [22]. It firstly projects the data matrix of explanatory variables in the optimal direction before applying smoothing functions to those. If *maxterms* (i.e. the number of linear combinations) is sufficiently large, PPR can be considered a universal approximator with considerable similarities to the so-called feed forward neural networks. However and similarly to the latter, complexity constraints need to be formulated to avoid overfitting. The algorithm starts by adding *maxterms* ridge functions. Then, it removes iteratively the least important term until *nterms* terms remain, which is the number of terms in the final model. Both *maxterms* and *nterms* are hyperparameters that need to be tuned beforehand. *optlevel* is a third hyperparameter which controls how thoroughly the models are refitted during this process. To smooth the ridge functions, we use by default Friedman's 'super smoother' *supsmu* which requires to fit the bass/span control.

# Implementation

In this section we will explain the practical framework used for conducting our project.

## 5.1 Environment

### 5.1.1 Hardware Environment

Our experiments were conducted using a desktop machine with the following characteristics.

- CPU Intel Core i3 3.2 GHz
- 8 GB RAM
- Hard Drive 800 GB
- System Type 64 bit

### 5.1.2 Software Environment

The experiments were conducted using the R Software [23] using Integrated Development Environment (IDE) RStudio.

## 5.2 Experiments

Data was divided into two sets: a training set and a test set (i.e. 70%/30%). Statistical independence was assumed to be in place among the routes. Consequently, we ended up having a total of 4 data sets. Vehicle ids were categorized into four groups: one containing all vehicle ids having less than 0.5% of the total number of trips and 3 obtained through a three-step clustering procedure. First, kernel density estimation was used to generate the p.d.f. for every unique vehicle id. Second, these p.d.f. were clustered by a Gaussian Mixture Model trained using the Expectation-Maximization algorithm. Finally, the Bayesian Information Criterion was used to select the best model.

Package `FSelector` [24] was used for RReliefF. The value used for `neighbour.count` (the number of nearest examples) in [6] was 10. For robustness reasons, we used `neighbour.count = 50` with  $m = 100$  iterations. For illustrative purposes on this particular issue, we used 0.1% of total data set as sample size. Similar results were found for a sample size of 0.5% and 1.0% of total data set length. A minimum weight threshold was set as  $\phi = 0.01$ . The default distance metric of `FSelector`'s implementation of RReliefF was used.

We used `glmnet` [25] procedures for fitting LASSO. The best  $\lambda$  was selected using cross validation.

### 5.2.1 Hyperparameter Tuning

Package `caret` [26] was utilized for hyperparameter tuning of RF, SVR and PPR. The two methods used in our experiments for hyperparameter optimization are (i) Grid Search (e.g. [6]) and (ii) Random Search [27]. (i) Grid Search exhaustively considers all the parameter combinations specified in a grid of parameter values. Hence, a high computational effort is required for large grids. A valid alternative introduced by Bergstra and Bengio [27] is Random Search. It consists on conducting independent draws from a uniform density using the same configuration space as the one defined by a regular grid. This approach only evaluates a random subsample of grid points - set to 60 in our case - and presents similar results to the grid one on an efficient man-

ner [27].

PPR has five different hyperparameters: `nterms`, `max.terms`, `optlevel`, `bass` and `span` (the two latter for *supsmu*). Random Search was used for tuning `nterms` Package `kernlab` [28] was used for SVR. SVR has six different hyperparameters: `kernel`, `C` (for all kernels), `epsilon` (for all kernels), `sigma` (for Radial kernel), `scale` and `degree` (only for polynomial kernel). Random Search was used for tuning `C`, `sigma`, `scale` and `degree`. Finally, Package `randomForest` was used for RF. Grid search was used for tuning both hyperparameters, as well as the ones non explicitly mentioned above.

The three above-mentioned base learners were evaluated based on the three resulting feature spaces: 1) feature set with all features (12 features), and as well as the ones given by 2) LASSO and 3) RReliefF. The obtained results were compared using two metrics of interest: RMSE and MAE.

## Results and Discussion

In this section we will explain the results and discussion using figures and tables .

### 6.1 Results

The optimal hyperparameter values for the three distinct setups are displayed in Table 2 for RF, PPR and SVR. Fig. 3 shows the results of RReliefF for each of the routes. x-axis is the feature set. y-axis is the weight; boxplots. It is evident that only the departure time has a predictive power accordingly with RReliefF. We therefore select departure time as the only feature from RReliefF method for each route. Fig. 2 shows the results of LASSO plots for each of the routes. x-axis are different  $\log(\lambda)$  values while y-axis are the coefficients. Features after the cut-off are selected to be the most suitable ones.

Finally, the evaluation of SVR, PPR and RF for the three feature sets for each of the routes are presented in Table 3 and Table 4 respectively.

### 6.2 Discussion

The tables clearly show that LASSO performs better than RReliefF on this particular task. RF is the algorithm that benefits less of the feature selection process since this task is inherent of its own modelling process. Fig. 3 illustrates the myopia of RReliefF

Table 2 – Optimal Hyperparameters setting.

		PPR			SVR		RF	
		nterms	max.terms	optlevel	$\sigma$	C	mtry	ntree
<b>LASSO</b>	A1	3	3	1	1179.65	909.04	1	500
	A2	3	3	1	997.91	335.15	1	500
	B1	3	3	1	1369.19	0.509	1	500
	B2	3	3	1	1008.36	72.25	1	700
<b>RreliefF</b>	A1	7	7	7	65.42	2.84	3	700
	A2	8	8	3	329.67	3.30	3	900
	B1	6	6	3	17.87	22.80	3	900
	B2	7	7	3	71.46	0.076	3	900
<b>ALL</b>	A1	8	8	3	0.15	909.04	6	900
	A2	9	9	3	0.19	318.79	6	500
	B1	11	11	3	0.17	312.95	6	900
	B2	5	5	3	0.24	72.25	6	900

on identifying some of the daytypes as relevant for reducing the bias-error around the target variable. As result, underfitted models (using only scheduled departure time) produce bad results - especially for PPR and SVR. These effects are depicted in Fig. 4, where the deficiency of the models output by either PPR and SVR during the peak hours when fed by RReliefF feature subspaces is highlighted. This effect happens because the daytype variables do not have a particular effect on the variance-error reduction - but mainly only on the bias one. In the authors' opinion, these results illustrate that RReliefF is not the best technique to handle the feature selection task on this particular problem.

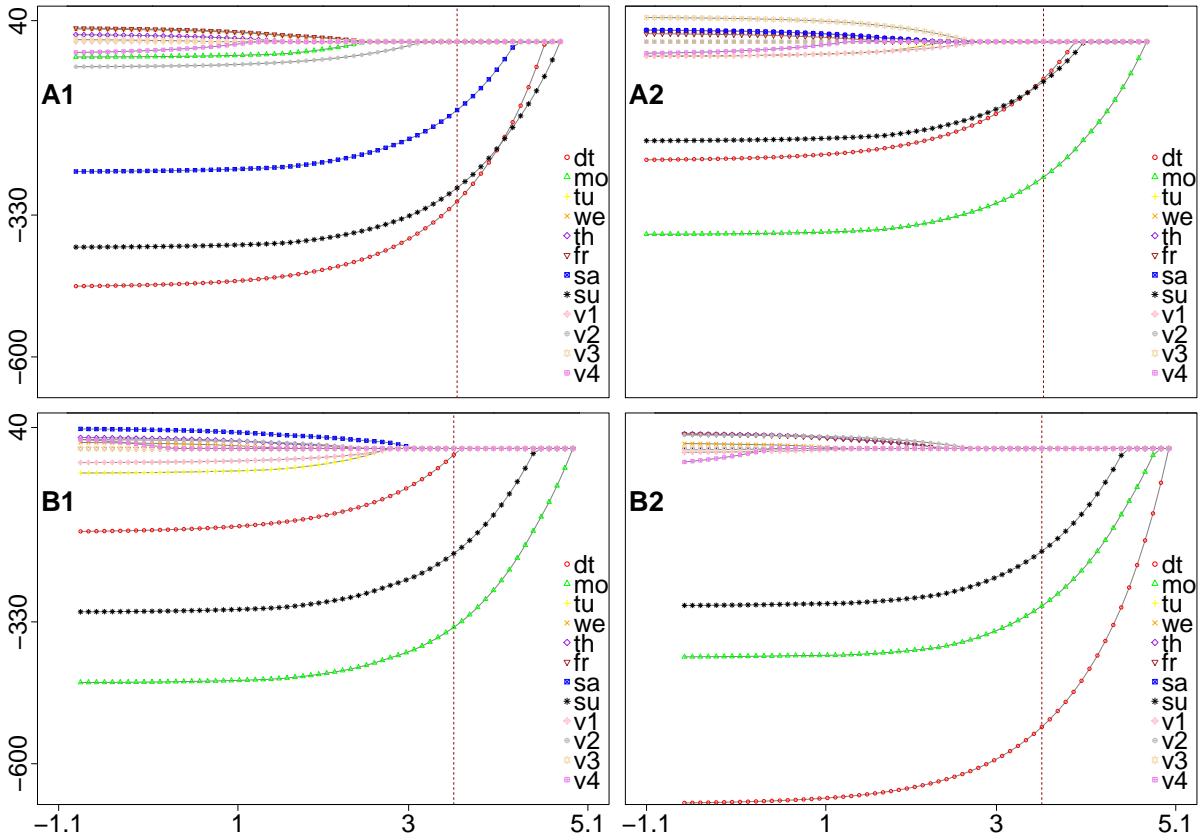


Figure 2 – LASSO results for all routes. A vertical red dashed line is drawn at the best  $\log \lambda$  value. This serves as cut.off point.

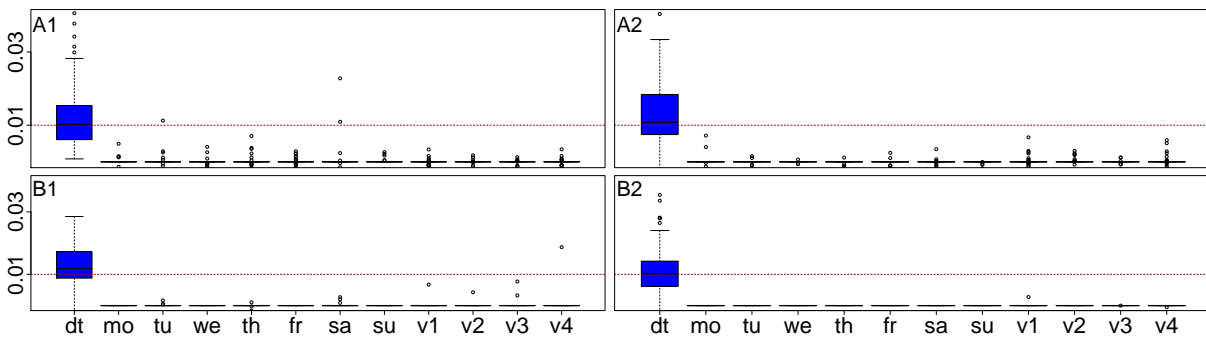


Figure 3 – RRelieF results for all routes. A horizontal red line is drawn at  $y=0.01$ .

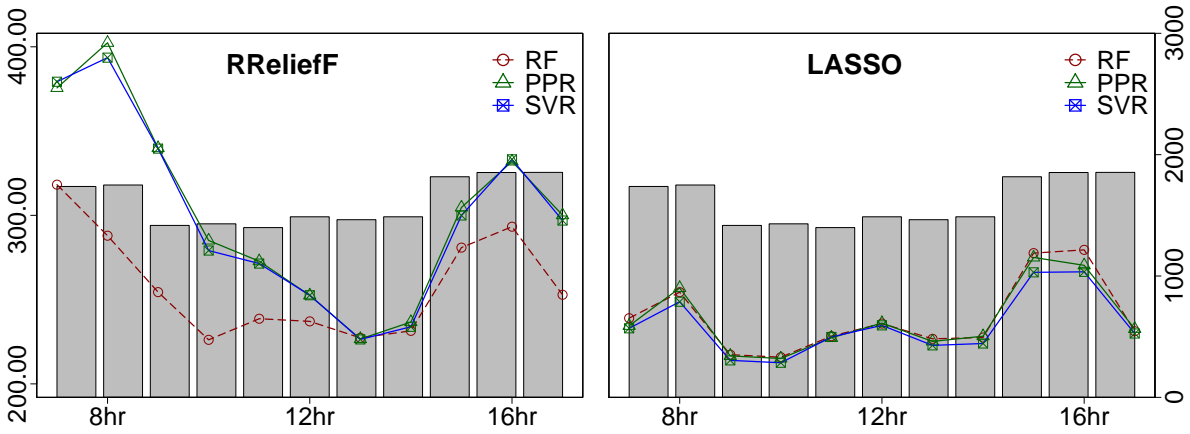


Figure 4 – RRelief and LASSO comparative analysis (y-axis) using RMSE (scaled on RRelief side) along different scheduled departure times (x-axis). Bars denote the sample size on each timespan (scaled on LASSO side).

Table 3 – SVR results for initial, LASSO and RrF-RReliefF feature sets.

Route	RMSE	MAE	RMSE	MAE	RMSE	MAE
	RrF	RrF	LASSO	LASSO	ALL	ALL
A1	293.294	224.455	<b>228.192</b>	<b>182.554</b>	244.888	196.851
A2	260.567	192.483	<b>196.843</b>	<b>154.453</b>	228.977	180.843
B1	309.361	224.084	<b>244.650</b>	<b>180.188</b>	281.480	205.387
B2	311.037	231.711	<b>255.853</b>	<b>204.383</b>	268.029	211.830
ALL	293.564	218.183	<b>231.384</b>	<b>180.394</b>	255.843	198.477

Table 4 – RF and PPR results for initial, LASSO and RrF-RReliefF feature sets.

Route	RF						PPR					
	RMSE RrF	MAE RrF	RMSE LASSO	MAE LASSO	RMSE ALL	MAE ALL	RMSE RrF	MAE RrF	RMSE LASSO	MAE LASSO	RMSE ALL	MAE ALL
A1	240.20	190.18	<b>227.66</b>	<b>181.44</b>	232.72	187.96	311.66	242.49	<b>231.52</b>	<b>186.83</b>	232.94	188.16
A2	235.65	180.76	<b>199.84</b>	<b>158.91</b>	203.11	161.73	263.04	195.19	<b>197.05</b>	<b>158.04</b>	198.80	159.45
B1	266.32	198.87	249.78	189.70	<b>248.95</b>	<b>188.71</b>	311.68	226.25	<b>247.69</b>	<b>188.21</b>	248.33	189.21
B2	283.59	223.68	266.62	215.42	<b>264.51</b>	<b>213.91</b>	316.31	233.50	264.21	213.77	<b>261.05</b>	<b>211.29</b>
ALL	256.44	198.37	<b>235.97</b>	<b>186.37</b>	237.32	188.08	300.67	224.36	<b>235.13</b>	<b>186.71</b>	235.28	187.03



## Conclusion and Future Work

Feature selection is a relevant task in any real-world data mining project. Long-term TTP for public transport planning and/or operational purposes is not an exception. Hereby, we discussed the limitations of RReliefF - the state-of-the-art for this problem. A comprehensive comparison with LASSO was conducted using a real-world case study from a bus operator in Sweden. The obtained results illustrated how dependent RReliefF is on an adequate distance metric that gives different relevance for distinct features - thus leading to a proper normalization of the RReliefF output weights and/or different selection thresholds for each feature accordingly to each one's contribution on model's bias reduction. As future work, we intend to explore further supervised filters for dimensionality reduction purposes on this task - such as Auto-Encoders.

# Bibliography

- [1] Moreira-Matias, L., Mendes-Moreira, J., de Sousa, J.F., Gama, J.: Improving mass transit operations by using avl-based systems: A survey. *IEEE Transactions on Intelligent Transportation Systems* **16**(4) (2015) 1636–1653
- [2] Mishalani, R., McCord, M., Forman, S.: Schedule-based and autoregressive bus running time modeling in the presence of driver-bus heterogeneity. In: *Computer-aided Systems in Public Transport*. Springer (2008) 301–317
- [3] Berkow, M., El-Geneidy, A., Bertini, R., Crout, D.: Beyond generating transit performance measures. *Transportation Research Record: Journal of the Transportation Research Board* **2111**(1) (2009) 158–168
- [4] El-Geneidy, A., Horning, J., Krizek, K.: Analyzing transit service reliability using detailed data from automatic vehicular locator systems. *Journal of Advanced Transportation* **45**(1) (2011) 66–79
- [5] Mazloumi, E., Rose, G., Currie, G., Sarvi, M.: An integrated framework to predict bus travel time and its variability using traffic flow data. *Journal of intelligent Transportation systems* **15**(2) (2011) 75–90
- [6] Mendes-Moreira, J., Jorge, A., de Sousa, J., Soares, C.: Comparing state-of-the-art regression methods for long term travel time prediction. *Intelligent Data Analysis* **16**(3) (2012) 427–449

- [7] Robnik-Šikonja, M., Kononenko, I.: An adaptation of relief for attribute estimation in regression. In: Machine Learning: Proceedings of the Fourteenth International Conference (ICML 97). (1997) 296–304
- [8] : Machine Learning. <http://whatis.techtarget.com/definition/machine-learning> Accessed: 2016-07-03.
- [9] : Supervised Learning. <https://azure.microsoft.com/en-us/documentation/articles/machine-learning-algorithm-choice/> Accessed: 2016-07-03.
- [10] : Regression. <http://www.analyticsvidhya.com/blog/2015/08/comprehensive-guide-regression/> Accessed: 2016-07-04.
- [11] : Bias and Variance. <http://scott.fortmann-roe.com/docs/BiasVariance.html> Accessed: 2016-07-04.
- [12] : Unsupervised Learning. <https://azure.microsoft.com/en-us/documentation/articles/machine-learning-algorithm-choice/> Accessed: 2016-07-03.
- [13] : Clustering. [http://home.deib.polimi.it/matteucc/Clustering/tutorial\\_html/](http://home.deib.polimi.it/matteucc/Clustering/tutorial_html/) Accessed: 2016-07-04.
- [14] Guyon, I., Elisseeff, A.: An introduction to variable and feature selection. The Journal of Machine Learning Research **3** (2003) 1157–1182
- [15] Mendes-Moreira, J., Moreira-Matias, L., Gama, J., de Sousa, J.: Validating the coverage of bus schedules: A machine learning approach. Information Sciences **293** (2015) 299–313
- [16] Kira, K., Rendell, L.A.: A practical approach to feature selection. In: Proceedings of the ninth international workshop on Machine learning. (1992) 249–256
- [17] Kononenko, I.: Estimating attributes: analysis and extensions of relief. In: Machine Learning: ECML-94, Springer (1994) 171–182

- [18] Tibshirani, R.: Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)* (1996) 267–288
- [19] Breiman, L., Friedman, J., Stone, C.J., Olshen, R.A.: *Classification and regression trees*. CRC press (1984)
- [20] Breiman, L.: Random forests. *Machine learning* **45**(1) (2001) 5–32
- [21] Cortes, C., Vapnik, V.: Support-vector networks. *Machine learning* **20**(3) (1995) 273–297
- [22] Friedman, J., Stuetzle, W.: Projection pursuit regression. *Journal of the American statistical Association* **76**(376) (1981) 817–823
- [23] R Core Team: *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. (2012)
- [24] Romanski, P.: Fselector: selecting attributes. *r package version 0.19* (2009)
- [25] Friedman, J., Hastie, T., Tibshirani, R.: Regularization paths for generalized linear models via coordinate descent. *Journal of statistical software* **33**(1) (2010) 1
- [26] Kuhn, M.: Caret package. *Journal of Statistical Software* **28**(5) (2008)
- [27] Bergstra, J., Bengio, Y.: Random search for hyper-parameter optimization. *The Journal of Machine Learning Research* **13**(1) (2012) 281–305
- [28] Zeileis, A., Hornik, K., Smola, A., Karatzoglou, A.: kernlab-an s4 package for kernel methods in r. *Journal of statistical software* **11**(9) (2004) 1–20