



HAL
open science

Online predictive diagnosis of electrical train door systems

Yufei Han, Olivier Francois, Allou Same, Laurent Bouillaut, Latifa Oukhellou, Patrice Aknin, Guillaume Branger

► **To cite this version:**

Yufei Han, Olivier Francois, Allou Same, Laurent Bouillaut, Latifa Oukhellou, et al.. Online predictive diagnosis of electrical train door systems. 10th World Congress on Railway Research (WCRR 2013), Nov 2013, Sydney, Australia. 6p. hal-00863798

HAL Id: hal-00863798

<https://hal.science/hal-00863798>

Submitted on 4 Oct 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Online predictive diagnosis of electrical train door systems

Y. Han¹, O. François¹, A. Same¹, L. Bouillaut¹, L. Oukhellou¹, P. Aknin¹, G. Branger²

¹ IFSTTAR, *French institute of science and technology for transport, development and networks, GRETTIA laboratory, F-77447 Marne la Vallée, France.*

² *Bombardier Transport France SAS, F-59154 Crespin, France.*

Abstract

Considering availability purposes for train transportation, passenger accesses (doors and steps) are often designated as critical systems. To improve global availability of its rolling stock, Bombardier Transportation (BT) aims at reinforcing its maintenance procedure by introducing predictive diagnosis. The SURFER project has been initiated to develop online and in-cars tools to early detect and prevent faults. In this paper, an overview of achieved progress with respect to online predictive diagnosis will be introduced. For this purpose, many signals are recorded using a test bench by BT: electrical motor intensity current, door displacement, binary indicators as door closed and locked. The paper focuses on designing a semi-supervised discriminative probabilistic model that take into account contextual variables (train inclination or constraints due to passengers affluence) to perform a robust predictive diagnosis. The main steps of the proposed method are the followings: the segmentation of the provided signals into opening and closing phases, the extraction of relevant features from opening/closing phases, the setting of the discriminative diagnosis model based on statistical semi-supervised learning. The proposed approach is tested on signals collected from regional trains fleeting around Paris. It allows the earlier detection of anomalies, for instance, those due to maladjustments. The practical implementation of this approach will be detailed together with its preliminary results.

1. Introduction

Rolling stock unavailability could be very expensive for the train builder during warranty period or for operating companies afterwards. Furthermore unavailability could lead to route planning disorganisation and to customer dissatisfaction. The SURFER project aims at deploying technical solution of active monitoring of on-board critical systems that are wisely solicited like customer accesses (the focus of this work). It also aims at optimising maintenance plan considering experiences about the system status evolution (the point developed in the paper entitled “*A rolling stock door system’s dynamic maintenance strategies based on a sensitivity analysis through Bayesian networks*” in WCRR 2013). In these two papers, the door systems are considered as the key point in terms of rolling stock availability.

In this paper, a methodology used to perform an on-board discrimination between normal signals and suspicious signals of the door system due to malfunctions or exterior events as customer malevolence or strong track cant is introduced. If an abnormal event is detected, an alarm could be sent to deported operators that can access to all historical data stored on-board to investigate and give instructions to the train driver if necessary. The proposed approach, which does not require any physical model of the system, is generic and can easily be transposed to other critical transportation systems provided that relevant data are available on these ones. A partial labeling of the signals with faults supplied by the experts is exploited to build the discriminative model. The design of the model also takes into account data variability associated to the normal operating state, which can be attributed to certain exploitation contexts as the train inclination.

The rest of the paper is organized as follows. In Section 2 we focus on extracting features of signals. Sections 3 and 4 introduce the general setting of the dataset and present the methodology. In Section 5, we conclude the paper and discuss the technologies that are employed for on-board diagnosis.

2. Selected indicators with real time transmission and processing constraints

In the case of door systems, the signals of electrical currents consumed during the opening and closing operations, the door displacement signals, the binary switch indicator (door opened/door closed) were found to be the more relevant data for an accurate diagnosis to be performed.

Due to transmission band width limitations, these signals have to be simplified into an indicator vector. The limited size of the corresponding file will allow on-board hard drive writings (during long exploitation periods) together with a 3G transmission to a distant terminal of selected signals upon request. A drastic compression of electrical signals has been done without loss of significant information depending on the requirements of practical applications. Under the bandwidth limitation constraint, the goal of this step is to extract the compact set of features or indicators that contain as much information as possible about the faults to be detected.

In this framework, it has been chosen to summarize the raw signals of electrical current consumed during the openings/closing cycles and their associated door displacement into different phases together with values of important switches. For each cycle of opening and closing, we get a set of 66 indicators including

- the required unlocking power,
- the amount of electrical current consumed during different identified translation phases,
- the time when the first peak of the electrical current occurs,
- the time duration of door translation,
- the instants, positions and speeds while limit switches are going up and down,
- the pressure estimated during locking phase.

During the opening/closing operation, different types of faults can arise from maladjusted door leaves, frictions from the top or the bottom of the mechanism, relaxed belt, lack of lubrication or exterior phenomenon such as customer action on the mechanism.

For the robustness sake of the classification issue, a historicity of degradation can be integrated. However this is not the focus of this paper. In the following section, we present how a discriminative model is constructed without taking into account that degradation evolution.

3. Semi-supervised discriminating model for automatic diagnosis

Now, we suppose that the available data are the extracted feature vectors $\{x_i, i=1,2,\dots,n\}$, where each x_i is a 66 dimensional vector describing physical characteristics of door open/close operations of trains. It is also assumed that each feature vector x_i has an associated binary label y_i equaling to 0 or 1 to indicate whether this one corresponds to a normal operating state or contains mechanical faults. The object of the diagnosis task is to construct a discriminative model based on the pairs $\{(x_i, y_i), i=1,2,\dots,n\}$. The input of the model is one such feature vector and the output is an estimation of signal state justifying whether the signal is normal or not.

In the extracted database, corresponding to 1344 doors distributed over 84 trains, only a very small part (3626 out of totally 1107554 signals) is labeled as signals with faults by the experts. Labels of the rest signals are left unknown. Therefore, construction of the model is actually a semi-supervised learning procedure [4]. Semi-supervised learning (SSL) is a family of algorithms halfway between supervised and unsupervised learning. It makes use of supervision information, but not necessarily for all the data samples. In addition to explicitly labeled samples, SSL also benefits from the data distribution characteristics of unlabeled samples and consider unlabeled data distribution as complementary constraints to enhance the model. In machine learning, SSL is often employed in classification when it is nevertheless unrealistic to collect labels of all training samples explicitly, which fits our task perfectly. In diagnosis of faults, even without explicit labels, the distribution of unlabeled collected signals is still helpful to formulate the learning framework. Firstly, normal signals always present homogeneous appearances. Thus they distribute compactly. In contrast, signals with potential faults have diverse profiles due to different mechanical causes of faults. Therefore, they contain relatively high variances in their distribution. Besides, the extracted signals are intrinsically highly

imbalanced. The door switching system runs smoothly during most times, while falls into faults occasionally. The signals with potential faults thus only compose a small proportion of collected data samples. With this property, the identification of faults can be considered as a procedure of outlier detection.

4. Semi-supervised construction of the discriminative diagnosis model

To attack the above varietal issues of the diagnosis task, we propose to employ a two-step semi-supervised learning to construct the expected discriminative model.

- We first utilize a semi-supervised robust Principal Component Analysis (PCA) to identify signals containing potential faults in the unlabeled data with aid of the partial supervision information. This procedure estimates y_i for each unlabeled signal and therefore complements labels for all the feature vectors.
- Based on the fully labelled signals $\{(x_i, y_i), i=1, 2, \dots, n\}$, the second step is to construct the our discriminative model, the bagging average of logistic regression, including both priory labels provided by the experts and those estimated.

4.1. Semi-supervised robust PCA

PCA is a linear projection of the data conserving most of variances in the data [1]. It has been widely used for representing high-dimensional data in a low-dimensional subspace and remove noises during preprocessing. PCA assumes data follow Gaussian distribution and fit the variance of data using the while noise model. Data samples that are severely biased from the Gaussian noise model are considered as noise or outliers.

$$\min_{(U, \mu)} \sum_{i=1}^n \|x_i - Ux_i - \mu\|^2 \quad (1)$$

where U is the projection matrix and x_i is the projection of x_i under U . The estimated columns of U are the eigenvectors of the covariance matrix of x_i and μ is the average of the data [1]. In the diagnosis task, the Gaussian model embedded in PCA will capture most variances of the normal signal samples. Those data samples distinctively far from the Gaussian centre imply non-homogeneous characteristics, thus identified as signals with potential faults. However, due to the least square intrinsic of PCA, extremely large errors will make the principal components deviated from the majority homogeneous signal samples. Besides, classical PCA model doesn't provide the interface to embed the observed labels. To improve its robustness and concatenate semi-supervision flexibility, we borrow the idea from robust M-estimator [3] to extend the standard PCA framework. Starting from initial parameters $(U^{(0)}, \mu^{(0)}, U^{(0)})$, the resulting robust algorithm iterate the following minimization until the parameters stabilize:

$$\min_{(U, \mu)} \sum_{i=1}^n \rho(\|x_i - Ux_i - \mu\|) \quad (2)$$

where

$$\rho(\|x_i - Ux_i - \mu\|) = \begin{cases} \frac{1}{2} \|x_i - Ux_i - \mu\|^2 & \text{if } \|x_i - Ux_i - \mu\| \leq \tau \\ \tau \|x_i - Ux_i - \mu\| & \text{if } \|x_i - Ux_i - \mu\| > \tau \end{cases} \quad (3)$$

c being the current iteration. Robust PCA learning is an iterative procedure. During the k th iteration, each signal sample x_i is assigned with a scalar $w_i^{(c)}$ valued between 0 and 1 to evaluate its biases from the Gaussian model assumption. For unlabelled signal samples, $w_i^{(c)}$ is defined as the exponential mapping of the reconstruction error between the centred signal and its reconstruction through the PCA projection. Without the explicit label, we assume the signal samples that cannot be well described by the Gaussian noise model as the ones that are mostly likely to contain faults. Compared with the majority homogeneous samples, they leads to higher reconstruction error, thus correspond to lower

$w_i^{(c)}$. With the increasing reconstruction error, the corresponding $w_i^{(c)}$ diminishes to 0 exponentially. The mean value of the signal samples μ and principle components U are updated by calculating the weighted version of the sample mean and covariance matrix of the signals. Through this way, the affections of the large errors to the PCA model are reduced to a large extent. For labelled signal samples with faults, $w_i^{(c)}$ is simply fixed to be 0 in order to explicitly remove faults patterns from the PCA model, which inserts the partially supervision information into the PCA model construction.

In our diagnosis task, the distribution of normal signal samples usually presents a multi-modal characteristic caused by mechanical structure settings which can be attributed to the fact that the learning database of feature vector is related to several doors, each of which may be subject to various contexts (train inclination or constraints due to passengers affluence). In this setting, single-modal PCA assumption doesn't fit well with the multi-modal intrinsic. To refine the model for better data description, we extend the idea of mixture of local PCA [1] into the robust PCA framework. In this setting, normal signal samples are supposed to be divided into m separated local groups defining a partition $P = (P_1, \dots, P_m)$, where each group P_j has its own PCA parameters (U_j, μ_j, U_{jj}) indexed by j . The resulting algorithm starts from the initial parameters $(U_j^{(0)}, \mu_j^{(0)}, U_{jj}^{(0)})$ and partition $P^{(0)} = (P_1^{(0)}, \dots, P_m^{(0)})$ and iterates the following two steps until the parameters and partition stabilize:

- step 1: for $j=1, \dots, m$, compute the parameters $(U_j^{(c)}, \mu_j^{(c)}, U_{jj}^{(c)})$ by solving the problem

$$\min_{U_j, \mu_j, U_{jj}} \left\| \sum_{i \in P_j} w_i (x_i - (U_j \mu_j + U_{jj})) \right\|_2^2 \quad (5)$$

where

$$w_i = \exp \left\{ - \frac{\| x_i - (U_j \mu_j + U_{jj}) \|_2^2}{\theta} \right\} \quad (6)$$

- step 2: compute the partition $P^{(c)} = (P_1^{(c)}, \dots, P_m^{(c)})$ according to the rule

$$x_i \in P_j \iff w_i > \theta \quad (7)$$

Training of the robust PCA algorithm with the multi-modal assumption usually converges after about 50 iterations. Fig.1 indicates the distribution of $\{w_i = \max_j w_{ij}\}$ for all collected feature vectors. As we can see, except labelled signals, 98% of the unlabelled signal samples correspond to $\{w_i\}$ that are larger than 0.3. For the unlabelled samples with w_i less than 0.3, we have strong confidence that they imply occurrence of mechanical faults in the door switch system. In our work, we consider any unlabelled signal with its w_i less than 0.3 as the signals with faults. The corresponding y_i are set to 0 as a result and the rest of y_i are set to 1. Note the threshold on $\{w_i\}$ for justifying labels for the unlabelled data is set in the optimal sense according to empirically experiences and experts' knowledge. Lower or higher thresholds can leads to stricter or looser criterion of fault diagnosis, which depends on practical needs of automatic diagnosis tasks.

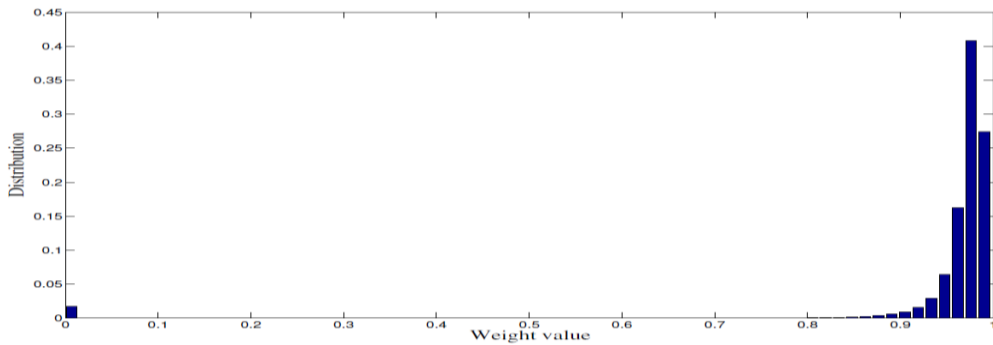


Figure 1: Distribution of weights w_i

4.2. Supervised bagging average of logistic regression

Give a fully labelled signal set $\{(x_i, y_i), i=1, \dots, n\}$ using the proposed robust PCA, a supervised bagging average [4] of discriminative models is constructed through the following iterative procedure. During each iteration, we sample the normal signals (samples with $y_i=1$) with replacement. Each randomly sampled set of normal signals has the same number of samples as the signals with faults (samples with $y_i=0$). A logistic regression model [2] M_j is then constructed based on the selected normal signals and all signals with faults. After T iterations of sampling and training, we can thus obtain T logistic regression models $\{M_j, j=1, 2, \dots, T\}$. With a signal x_i as the input, the model average $(M_1(x_i) + \dots + M_T(x_i))/T$ is used as the final diagnosis decision, where $M_j(x_i)$ is the output of model M_j for input x_i . Since the output of each logistic regression model M_j ranges between 0 and 1, then the average output is well normalized and ranges between 0 and 1. Therefore, the output of the model average acts as a soft label of the input signal. The higher output indicates the stronger confidence of the corresponding input signal to be normal, and vice versa. In our work, we choose the number of training iterations T to be 500, in order to obtain satisfied regression result. Fig.2 illustrates the distribution of the bagging average output of logistic regression models $\{M_j, j=1, 2, \dots, T\}$ on the collected signal set. As we can see, most signals correspond to the output approaching to 1, which is consistent with the fact that most collected signals are normal. A small peak locates around 0 in the figure. The signal samples in this region are highly suspicious to contain potential faults.

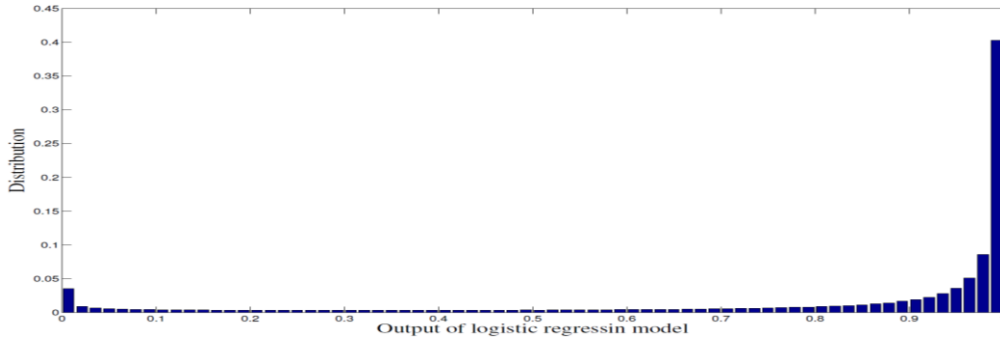


Figure 2: Distribution of logistic regression model output

5. Online predictive diagnosis

Once the parameters of the discriminative model are learned following the previous sections, the predictive diagnosis then consists, for a given door, of tracking the temporal evolution of the score provided by the logistic regression model. Fig.3 shows the temporal variation of the bagging average output for one specific door, lasting for almost one year. Blue curves show directly the variation of the soft diagnosis result of the model average. The green curves are obtained by performing median filter on the blue curves, indicating general temporal tendencies of the discriminative model output. We can observe the transition between the normal state and the faulty state. As we can find, the discriminative model output varies wildly along the temporal scale. This phenomenon is caused by occasional blockage due to passengers for instance. It usually causes a fault-like pattern in indicators, while it neither doesn't last for long time, nor affects the performances of the following door opening/closing operations. Median filter performs a sliding time window scanning the temporal variation of the discriminative model output. For each time window, median filter extracts median value of the model outputs within this model and smoothes the temporal fluctuations using the median value. Therefore, with the proper setting of the sliding window, median filter is able to extract and highlight the real mechanical faults that have distinctive influences in temporal scale. In the online diagnosis application, we first use the outputs of the discriminative model to trace the variation of the signal status. An online median filtering is then used to remove the occasional fluctuations and monitor the temporal tendency of the model output. Once the extracted median value of one specific time window is less than a predefined threshold, the diagnosis system will trigger an alarm for potential faults in the door.

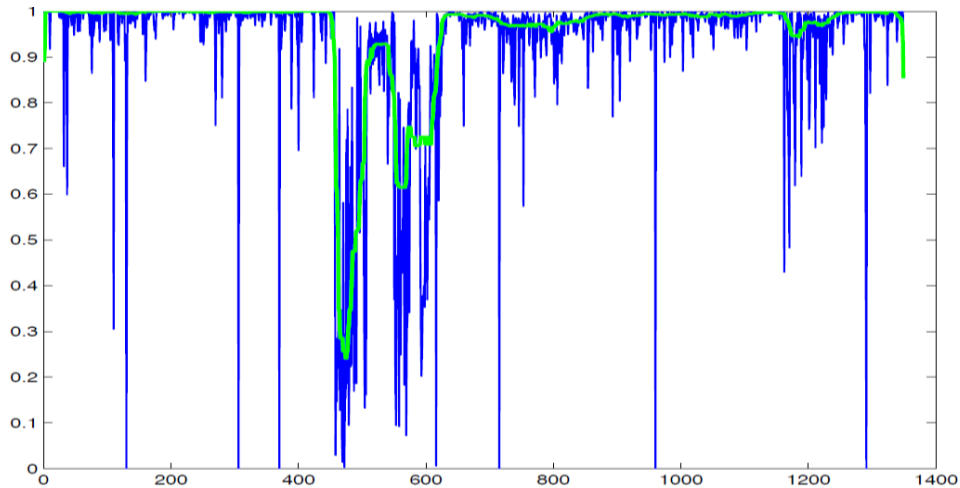


Figure 3: Temporal variation of logistic regression model outputs from February 1st 16:30 to December 31st 18:09.

6. Conclusions and future work

In this paper, we present the detailed construction procedure of a discriminative model for online predictive diagnosis of train door system. This work is composed by extraction of informative indicators from the original electrical signals and the following semi-supervised discriminative model learning based on the partially labelled features. Based on the achievement, our work enables online storing and transmission of the electrical signal features, while conserving as much as possible the needed information. Furthermore, the semi-supervised learning framework provides a flexible way to integrate expert knowledge efficiently and allows interactions between human experiences and system decisions. According to the experimental results, the presented method can track the temporal evolution of the train door system status and manage to early detect the faults of the doors.

Our future work focuses on constructing a temporal dynamic model to describe the variations of the system status. It will integrate the function of early fault detection of the presented method and the probabilistic inference of the system status into a unified framework. With this dynamic model, we could achieve even more accurate monitoring of the system by smoothing the effects of the occasional events and focusing on the distinctive fault patterns.

Acknowledgements

The works presented in this paper are due to the SURFER project between Bombardier Transport and IFSTTAR (French National Institute for Transport, Development and Risk Sciences and Technologies). The authors want to thanks all the people who have collaborated to the critical and laborious step of data extraction to allow the model to be the more efficient as possible.

References

- [1] G. E. Hinton, P. Dayan and M. Revow. Modeling the manifolds of images of handwritten digits. IEEE Transactions on Neural Networks, number 8, pp 65-74, 1997.
- [2] D. W. Hosmer and L. Stanley. Applied Logistic Regression (2nd ed.). Wiley, 2000.
- [3] P. J. Huber. Robust Statistics. Wiley, 1981.
- [4] T. Hastie and J. Friedman. Elements of statistical learning: DataMining, Inference and Prediction. Springer, 2009.