

# Explaining Machine Learning Models using Entropic Variable Projection

François Bachoc<sup>1</sup>, Fabrice Gamboa<sup>1,3</sup>, Max Halford<sup>2</sup>, Jean-Michel Loubes<sup>1,3</sup> and Laurent Risser<sup>1,3</sup>

<sup>1</sup>Institut de Mathématiques de Toulouse

<sup>2</sup> Institut de recherche en informatique de Toulouse

<sup>3</sup> Artificial and Natural Intelligence Toulouse Institute (3IA ANITI)

## Abstract

In this paper, we present a new explainability formalism designed to explain how each input variable of a test set impacts the predictions of machine learning models. Hence, we propose a group explainability formalism for trained machine learning decision rules, based on their response to the variability of the input variables distribution. In order to emphasize the impact of each input variable, this formalism uses an information theory framework that quantifies the influence of all input-output observations based on entropic projections. This is thus the first unified and model agnostic formalism enabling data scientists to interpret the dependence between the input variables, their impact on the prediction errors, and their influence on the output predictions. Convergence rates of the entropic projections are provided in the large sample case. Most importantly, we prove that computing an explanation in our framework has a low algorithmic complexity, making it scalable to real-life large datasets. We illustrate our strategy by explaining complex decision rules learned by using XGBoost, Random Forest or Deep Neural Network classifiers on various datasets such as *Adult income*, *MNIST* and *CelebA*. We finally make clear its differences with the explainability strategies *LIME* and *SHAP*, that are based on single observations. Results can be reproduced by using the freely distributed Python toolbox <https://gems-ai.com/>.

## Index Terms

Explainability, black-box decision rules, Kullback-Leibler divergence, Wasserstein distance.

## I. INTRODUCTION

Machine learning algorithms build predictive models which are nowadays used for a large variety of tasks. Over the last decades, the complexity of such algorithms has grown, going from simple and interpretable prediction models based on regression rules to very complex models such as random forests, gradient boosting, and models using deep neural networks. We refer to [15] for a description of these methods. Such models are designed to maximize the accuracy of their predictions at the expense of the interpretability of the decision rule. Little is also known about how the information is processed in order to obtain a prediction, which explains why such models are widely considered as black boxes.

This lack of interpretability gives rise to several issues. When an empirical risk is minimized, the training step may be unstable or highly dependent on the optimization procedure due for example to non-convexity and multimodality. Another subtle, though critical, issue is that the optimal decision rules learned by a machine learning algorithm highly depend on the properties of the learning sample. If a learning sample presents unwanted trends or a bias, then the learned decision rules will reproduce these trends or bias, even if there is no intention of doing so. As a consequence, many users express a lack of trust in these algorithms. The European Parliament even adopted a law called GDPR (General Data Protection Regulation) to protect the citizens from decisions made without the possibility of explaining why they were taken, introducing a right for explanation in the civil code. Hence, building intelligible models is nowadays an important challenge for data scientists.

Different methods have been proposed to make understandable the reasons leading to a prediction, each author using a different notion of explainability and interpretability. We mention early works by [17] for recommender systems, [8] for neural networks and [11] or [26] for generalized additive models. Another generic solution, described in [2] and [6], focused on medical applications. In [24], a discussion was recently opened to refine the discourse on interpretability. Recently, a special attention has also been given to deep neural systems. We refer for instance to [18, 28, 33] and references therein. Clues for real-world applications are given in [14]. In [31] (LIME), the authors recently proposed to locally mimic a black-box model and then to give a feature importance analysis of the variables at the core of the prediction rule. A counterfactual model [37] was also proposed in [13] to explain how the predictions made by a classifier on a query image can be changed by transforming a region of this image. In the same vein, a method called integrated gradients was specifically designed in [34] for the interpretability of single predictions using neural networks. In the *Fair learning* community, counterfactual models are also used to assess whether the predictions of machine learning models are fair [22, 3]. An individual explanation method on images through adversarial examples was also presented in [19]. In [20] the authors finally proposed a strategy to understand black-box models, as we do, but in a parametric setting.

In this work, we specifically present a strategy to generate test samples that present a deviation with respect to an original test distribution. In order to achieve this, we enforce a bias on this distribution. By using an entropic projection method, we then optimally reweight the observations so that their empirical distribution undergoes the chosen bias while still being optimally close to the original test dataset distribution. For instance, the mean value of an input variable of interest can be shifted. We prove the optimality of the reweighting and obtain a theoretical control on the generated distribution with respect to Wasserstein distance. Furthermore, we provide a feasible method that is scalable to large datasets.

This method has several applications in robustness and explainability for algorithms. In this work, we provide a general procedure for global explainability of algorithms. The incorporation of bias in a test dataset makes it possible to quantify the response of a model to a chosen and calibrated stress on the input. It enables to understand how the algorithm reacts to such modifications in distribution. Hence, it explains the role played by the input variables. We can also infer causal properties when looking at changes of predictions caused by changes in the variables, enhancing interpretability of the algorithm. We remark that enforcing a bias to the original distribution is inspired by the fields of sensitivity analysis and computer experiments [23, 32].

We emphasize that contrary to *e.g.* [31], [34] or [27] (SHAP), our method deals with global explainability since it quantifies the global effect of the variables for all the test samples instead of individual observations. We also highlight that our point of view is different from previous works where the importance of each variable was also considered. Sparse models (see [5] for a general introduction on the importance of sparsity) enable to identify important variables. Importance indicators have also been developed in machine learning to detect which variables play a key role in the algorithm. For instance, importance of variables is often computed using feature importance or Gini indices (see in [30] or [15]). Yet such indexes are computed without investigating the particular effects of each variable and without explaining their particular role in the decision process. We also strongly believe that running the algorithm over observations which are created artificially by increasing stepwise the value of a particular variable is not a desirable solution. By doing so, the correlations between variables are indeed not taken into account. Moreover, newly generated observations may be outliers with respect to the learning and test samples.

The paper falls into the following parts: Methodology is explained in Sections II and III. Results are given in Section IV and the discussions are finally developed in Section V. An appendix contains the proofs, different methodological extensions and additional numerical results.

## II. OPTIMAL PERTURBATION OF MACHINE LEARNING DATASETS

We consider a test set  $\{(X_i, Y_i)\}_{i=1, \dots, n}$  (here  $X_i = (X_i^1, \dots, X_i^p)$  are input observations while  $Y_i$  is the true output), on which we consider the outcome of black box decision rules  $f : \mathbb{R}^p \rightarrow \mathbb{R}$ . We consider throughout this paper that  $f$  has been learned by a training set and is fixed. We set  $\hat{Y}_i = f(X_i)$ , the predicted values. Hence we have at hand values  $(X_i, \hat{Y}_i, Y_i)$  for  $i = 1, \dots, n$ .

Our goal is to explain the global behaviour of  $f$ . To achieve this, we propose to study the response of  $f$  to distributional perturbations (or bias) of the input variables. Since  $f$  has been learnt using data following a given distribution, the domain of validity of the algorithm should not deviate too much from this initial distribution. Hence we propose to build perturbed distributions that are as close as possible to the initial distribution using an information theory method. We will show that this amounts to reweighing the observations by proper weights calibrated to incorporate the chosen perturbation on the input variables as explained in Sections II-B and II-C. The methodology to make this problem well posed and to quickly compute the optimal weights is the core of this paper.

### A. General optimal perturbations under moment constraints in machine learning

In order to experiment and to explore the behavior of a predictive model, a natural idea is to study its response to biased inputs. Different ways are possible to create modifications of a probability measure. In this paper, we consider an information theory framework in which we modify the distribution of the original test set  $Q_n = \frac{1}{n} \sum_{i=1}^n \delta_{X_i, \hat{Y}_i, Y_i}$ , by stressing the mean value of a function  $\Phi$  (or simply by stressing the mean value of variables). We will minimize the Kullback-Leibler information (also called mutual entropy) with respect to  $Q_n$ , making the problem well posed.

First, let us recall the definition of the Kullback-Leibler information. Let  $(E, \mathcal{B}(E))$  be a measurable space and  $Q$  a probability measure on  $E$ . If  $P$  is another probability measure on  $(E, \mathcal{B}(E))$ , then the Kullback-Leibler information  $\text{KL}(P, Q)$  is equal to  $\int_E \log \frac{dP}{dQ} dP$ , if  $P \ll Q$  and  $\log \frac{dP}{dQ} \in L^1(P)$ , and  $+\infty$  otherwise.

For a given  $k \geq 1$ , let

$$\Phi : (X, Y, \hat{Y}) \in \mathbb{R}^{p+2} \rightarrow \Phi(X, Y, \hat{Y}) \in \mathbb{R}^k$$

be a measurable function representing the shape of the stress deformation on the whole input. Note that our results are stated for a generic function  $\Phi$  of all variables  $(X, Y, \hat{Y})$ . Of course, this includes the case of functions depending only on  $(X, Y)$ . For  $t \in \mathbb{R}^k$ , we aim at finding a new distribution  $Q_t$  satisfying the constraint

$$\int_{\mathbb{R}^k} \Phi(x) dQ_t(x) = t,$$

and being the closest possible to the initial empirical distribution  $Q_n$  in the sense of Kullback-Leibler divergence, *i.e.* with  $\text{KL}(Q_t, Q_n)$  as small as possible.

We set for two vectors  $x, y \in \mathbb{R}^k$  the scalar product as  $\langle x, y \rangle = x^\top y$ . We next characterize the new distribution  $Q_t$ .

**Theorem II.1.** *Let  $t \in \mathbb{R}^k$  and  $\Phi : \mathbb{R}^{p+2} \rightarrow \mathbb{R}^k$  be measurable. Assume that  $t$  can be written as a convex combination of  $\Phi(X_1, \hat{Y}_1, Y_1), \dots, \Phi(X_n, \hat{Y}_n, Y_n)$ , with positive weights. Assume also that the empirical covariance matrix  $\mathbb{E}_{Q_n}(\Phi \Phi^\top) - \mathbb{E}_{Q_n}(\Phi) \mathbb{E}_{Q_n}(\Phi^\top)$  is invertible.*

Let  $\mathbb{P}_{\Phi,t}$  be the set of all probability measures  $P$  on  $\mathbb{R}^{p+2}$  such that  $\int_{\mathbb{R}^{p+2}} \Phi(x) dP(x) = t$ . For a vector  $\xi \in \mathbb{R}^k$ , let  $Z(\xi) := \frac{1}{n} \sum_{i=1}^n e^{\langle \Phi(X_i, \hat{Y}_i, Y_i), \xi \rangle}$ . Define now  $\xi(t)$  as the unique minimizer of the strictly convex function  $H(\xi) := \log Z(\xi) - \langle \xi, t \rangle$ . Then,

$$Q_t := \operatorname{arginf}_{P \in \mathbb{P}_{\Phi,t}} \operatorname{KL}(P, Q_n) \quad (1)$$

exists and is unique. Furthermore, we have

$$Q_t = \frac{1}{n} \sum_{i=1}^n \lambda_i^{(t)} \delta_{X_i, \hat{Y}_i, Y_i}, \quad (2)$$

with, for  $i = 1, \dots, n$ ,

$$\lambda_i^{(t)} = \exp \left( \langle \xi(t), \Phi(X_i, \hat{Y}_i, Y_i) \rangle - \log Z(\xi(t)) \right). \quad (3)$$

- A particularly appealing aspect of Theorem II.1 is that  $Q_t$  is supported by the same observations as  $Q_n$ , the mean change for  $\Phi$  only leading to different weights for the observations. This desirable property is obtained thorough the choice of the Kullback-Leibler information as a measure of similarity between  $Q_t$  and  $Q_n$ . It also turns the corresponding optimization problem into a favorable problem, both theoretically and numerically. As a consequence, sampling new stressed test sets does not require to create new input-output observations  $(X_i, \hat{Y}_i, Y_i)$  but only to compute the weights  $\lambda_i^{(t)}$ . This can be solved very quickly using (3).
- The optimization problem in Theorem II.1, defining  $\xi(t)$ , is convex and can be tackled very efficiently. The gradient of its objective function is provided in Appendix B. This makes it possible to deal with very large databases without computing new values for new observations. This differs from existing techniques based on perturbed observations as *e.g.* in LIME [31], where the data used for testing are created by changing randomly the labels or by bootstrapping the observations.
- Set

$$t_0 = \int_{\mathbb{R}^k} \Phi(x) dQ_n(x) = \frac{1}{n} \sum_{i=1}^n \Phi(X_i, \hat{Y}_i, Y_i)$$

the empirical mean of  $\Phi$  with the distribution  $Q_n$ . The quantity  $t - t_0$  can be thus understood as the amount of change on the mean of  $\Phi$  induced by changing  $Q_n$  into  $Q_t$ . We remark that, in Theorem II.1, the condition that  $t$  can be written as a convex combination of  $\Phi(X_1, \hat{Y}_1, Y_1), \dots, \Phi(X_n, \hat{Y}_n, Y_n)$ , with positive weights, is almost minimal. Indeed, it is necessary that  $t$  can be written as a convex combination of  $\Phi(X_1, \hat{Y}_1, Y_1), \dots, \Phi(X_n, \hat{Y}_n, Y_n)$  (otherwise the set of distributions that are mutually absolutely continuous to  $Q_n$  and yield expectation  $t$  for  $\Phi$  is empty). For all considered examples in this paper, this condition in Theorem II.1 was not restrictive.

Hereafter, we show how to choose  $\Phi$  specifically, to shed light on the impact of various features of the input variables.

### B. Application to variable importance by stressing the mean

We now apply Theorem II.1 to the special case of perturbing the mean of one of the  $p$  input variables, meaning that  $\Phi$  is valued in  $\mathbb{R}$  (*i.e.*  $k = 1$ ).

**Theorem II.2.** Let  $t \in \mathbb{R}$  and  $j_0 \in \{1, \dots, p\}$ . Assume that  $\min_{i=1}^n X_i^{j_0} < t < \max_{i=1}^n X_i^{j_0}$ .

Let  $\mathbb{P}_{j_0,t}$  be the set of probability measures on  $\mathbb{R}^{p+2}$  such that, when  $(X, \hat{Y}, Y)$  follows a distribution  $P \in \mathbb{P}_{j_0,t}$ , we have  $\mathbb{E}(X^{j_0}) = t$ . For  $\xi \in \mathbb{R}$ , let  $Z(\xi) := \frac{1}{n} \sum_{i=1}^n e^{\xi X_i^{j_0}}$ . Define now  $\xi(t)$  as the unique minimizer of the strictly convex function  $H(\xi) := \log Z(\xi) - \xi t$ . Then,

$$Q_{j_0,t} := \operatorname{arginf}_{P \in \mathbb{P}_{j_0,t}} \operatorname{KL}(P, Q_n)$$

exists and is unique. Furthermore, we have

$$Q_{j_0,t} = \frac{1}{n} \sum_{i=1}^n \lambda_i^{(j_0,t)} \delta_{X_i, \hat{Y}_i, Y_i},$$

with, for  $i = 1, \dots, n$ ,

$$\lambda_i^{(j_0,t)} = \exp \left( \xi(t) X_i^{j_0} - \log Z(\xi(t)) \right).$$

This theorem enables to re-weight the observations of a variable so that its mean increases or decreases. This is then used in Section III to understand the particular role played by each variable.

### C. Stressing several means, variances and covariances

The next corollary enables to stress the means of several variables at the same time.

**Corollary II.3** (perturbing several means). *Let  $1 \leq c \leq p$  and let  $j_1, \dots, j_c$  be two-by-two distinct in  $\{1, \dots, p\}$ . Let  $t_1, \dots, t_c \in \mathbb{R}$ . Assume that there exists a convex combination of  $(X_i^{j_1}, \dots, X_i^{j_c})_{i=1, \dots, n}$  with positive weights that is equal to  $(t_1, \dots, t_c)$ . Assume also that the empirical covariance matrix of  $(X_i^{j_1}, \dots, X_i^{j_c})_{i=1, \dots, n}$  is invertible. Then, there exists a unique distribution  $Q_{t_1, \dots, t_c}$  on  $\mathbb{R}^{p+2}$  such that for  $(X^1, \dots, X^p, \hat{Y}, Y) \sim Q_{t_1, \dots, t_c}$  we have  $\mathbb{E}(X^{j_a}) = t_a$  for  $a = 1, \dots, c$  and such that  $\text{KL}(Q_{t_1, \dots, t_c}, Q_n)$  is minimal. This distribution is obtained by the distribution  $Q_t$  in Theorem II.1, in the case where  $k = c$  and  $\Phi(X^1, \dots, X^p, \hat{Y}, Y) = (X^{j_1}, \dots, X^{j_c})$ .*

The next corollary enables to stress the variance of a variable while preserving its mean  $m_{j_0} = \frac{1}{n} \sum_{i=1}^n X_i^{j_0}$ .

**Corollary II.4** (perturbing the dispersion). *Let  $j_0 \in \{1, \dots, p\}$ . Let  $v \in [0, \infty)$ . Assume that there exists a convex combination of  $(X_i^{j_0}, (X_i^{j_0})^2)_{i=1, \dots, n}$  with positive weights that is equal to  $(m_{j_0}, m_{j_0}^2 + v)$ . Assume also that the empirical covariance matrix of  $(X_i^{j_0}, (X_i^{j_0})^2)_{i=1, \dots, n}$  is invertible. Then, there exists a unique distribution  $Q_{j_0, v}$  on  $\mathbb{R}^{p+2}$  such that for  $(X^1, \dots, X^p, \hat{Y}, Y) \sim Q_{j_0, v}$  we have  $\mathbb{E}(X^{j_0}) = m_{j_0}$  and  $\text{Var}(X^{j_0}) = v$  and such that  $\text{KL}(Q_{j_0, v}, Q_n)$  is minimal. This distribution is obtained by the distribution  $Q_t$  in Theorem II.1, in the case where  $k = 2$ ,  $\Phi(X^1, \dots, X^p, \hat{Y}, Y) = (X^{j_0}, (X^{j_0})^2)$  and  $t = (m_{j_0}, m_{j_0}^2 + v)$ .*

Finally, we next show how to stress the covariance between two variables while preserving their means  $m_{j_1} = \frac{1}{n} \sum_{i=1}^n X_i^{j_1}$  and  $m_{j_2} = \frac{1}{n} \sum_{i=1}^n X_i^{j_2}$ .

**Corollary II.5** (perturbing the covariance). *Let  $j_1, j_2 \in \{1, \dots, p\}$  be distinct. Let  $c \in \mathbb{R}$ . Assume that there exists a convex combination of  $(X_i^{j_1}, X_i^{j_2}, X_i^{j_1} X_i^{j_2})_{i=1, \dots, n}$  with positive weights that is equal to  $(m_{j_1}, m_{j_2}, m_{j_1} m_{j_2} + c)$ . Assume also that the empirical covariance matrix of  $(X_i^{j_1}, X_i^{j_2}, X_i^{j_1} X_i^{j_2})_{i=1, \dots, n}$  is invertible. Then, there exists a unique distribution  $Q_{j_1, j_2, c}$  on  $\mathbb{R}^{p+2}$  such that for  $(X^1, \dots, X^p, \hat{Y}, Y) \sim Q_{j_1, j_2, c}$  we have  $\mathbb{E}(X^{j_1}) = m_{j_1}$ ,  $\mathbb{E}(X^{j_2}) = m_{j_2}$  and  $\text{Cov}(X^{j_1}, X^{j_2}) = c$  and such that  $\text{KL}(Q_{j_1, j_2, c}, Q_n)$  is minimal. This distribution is obtained by the distribution  $Q_t$  in Theorem II.1, in the case where  $k = 3$ ,  $\Phi(X^1, \dots, X^p, \hat{Y}, Y) = (X^{j_1}, X^{j_2}, X^{j_1} X^{j_2})$  and  $t = (m_{j_1}, m_{j_2}, m_{j_1} m_{j_2} + c)$ .*

### D. Asymptotic rate of convergence

While, in this paper, the test set  $(X_i, \hat{Y}_i, Y_i)_{i=1, \dots, n}$  with empirical distribution  $Q_n$  is considered fixed, in Section II-D (and only in Section II-D), we assume that this test set is random, composed of i.i.d. realizations of a distribution  $Q^*$ .

The following proposition proves a statistical rate of convergence of our methodology. We show that the optimally perturbed distribution  $Q_t$  of Theorem II.1 (defined w.r.t.  $Q_n$ ,  $\Phi$  and  $t$ ) converges to the corresponding optimally perturbed distribution  $Q_t^*$ , (defined w.r.t.  $Q^*$ ,  $\Phi$  and  $t$ ). The convergence is measured by the  $L^1$  Wasserstein distance  $\mathcal{W}_1$ , defined by, for two distributions  $P, Q$  on  $\mathbb{R}^{p+2}$  with finite first moments,

$$\mathcal{W}_1(P, Q) = \inf_{U \sim P, V \sim Q} \mathbb{E}(\|U - V\|),$$

where the above infimum is taken over all pairs of (dependent or independent) random variables  $U, V$  such that  $U \sim P, V \sim Q$ . We refer for instance to [29] for the definition and computation of  $\mathcal{W}_1$  and to [1, 7] for recent work related to it.

**Proposition II.6.** *Let  $\Phi : \mathbb{R}^{p+2} \rightarrow \mathbb{R}^k$  and  $t \in \mathbb{R}^k$  be fixed. Assume that the support of  $Q^*$  is bounded and that  $\Phi$  is bounded in absolute value and Lipschitz continuous on the support of  $Q^*$ . Assume also that for  $v \in \mathbb{R}^k$ ,  $b \in \mathbb{R}$ ,  $Q^*(\{x \in \mathbb{R}^{p+2}; \langle v, \Phi(x) \rangle = b\}) = 1$  if and only if  $v = 0$  and  $b = 0$ . Assume finally that there exists a distribution  $Q$ , mutually absolutely continuous to  $Q^*$ , such that  $\int_{\mathbb{R}^{p+2}} \Phi(x) dQ(x) = t$ .*

*Then there exists a unique measure  $Q_t^*$  on  $\mathbb{R}^{p+2}$  such that  $\int_{\mathbb{R}^{p+2}} \Phi(x) dQ_t^*(x) = t$  and  $\text{KL}(Q_t^*, Q^*)$  is minimal. Furthermore, as  $n \rightarrow \infty$ , with  $Q_t$  given in Theorem II.1,*

$$\mathcal{W}_1(Q_t, Q_t^*) = O_p\left(n^{-1/(p+2)}\right).$$

The rate of convergence  $O_p\left(n^{-1/(p+2)}\right)$  is standard for the  $L^1$  Wasserstein distance in dimension  $p+2$ , see [12] for instance. The proof of Proposition II.6 combines techniques from the analysis of Wasserstein distances and from M-estimation with convex objective functions.

For the sake of conciseness, Proposition II.6 is stated under boundedness assumptions and with independent realizations from  $Q^*$ . These conditions could be weakened.

## III. EXPLAINABLE MODELS USING OPTIMALLY PERTURBED DATA SETS

In this section we consider that the transformation  $\Phi : \mathbb{R}^{p+2} \rightarrow \mathbb{R}^k$  and the target multidimensional moment  $t \in \mathbb{R}^k$  have been selected and that the conditions of Theorem II.1 are satisfied. This theorem provides the optimally perturbed distribution  $Q_t$ , given by the weights  $(\lambda_i^{(t)})_{i=1, \dots, n}$ . We now suggest various quantitative properties of  $Q_t$  (that we call quantities of interest), that can quantify the behavior of the studied black box decision rule.

We shall focus on two classical situations encountered in machine learning: binary classification and multi-class classification. The regression case is also explained in Appendix C.

#### A. The case of binary classification

Consider that  $Y_i$  and  $\hat{Y}_i = f(X_i)$  belong to  $\{0, 1\}$  for all  $i = 1 \dots, n$ . This corresponds to the binary classification problem for which the usual loss function is  $\ell(Y, f(X)) = \mathbb{1}_{\{Y \neq f(X)\}}$ . We suggest to use the indicators described hereafter for the perturbed distributions  $Q_t = \frac{1}{n} \sum_{i=1}^n \lambda_i^{(t)} \delta_{(X_i, \hat{Y}_i, Y_i)}$ . Explaining the decision rules may first consist in quantifying the evolution of the error rate as a function of  $t - t_0$ , hence the first indicator is the error rate, *i.e.*

$$\text{ER}_t = \frac{1}{n} \sum_{i=1}^n \lambda_i^{(t)} \mathbb{1}_{\{f(X_i) \neq Y_i\}}.$$

In terms of interpretation, when  $\Phi$  is given as in Theorem II.2, with  $\Phi(X^1, \dots, X^p, \hat{Y}, Y) = X^{j_0}$ ,  $t$  corresponds to the new (stressed) mean of the variable  $X^{j_0}$  while the former (unstressed) mean is  $t_0$ . In this case, plotting  $\text{ER}_t$  as a function of  $t - t_0$  highlights the variables which produce the largest amount of confusion in the error, *i.e.* those for which small or large values provide the most variability among the two predicted classes, hampering the prediction error rate. The False and True Positive Rates may alternatively be represented using

$$\text{FPR}_t = \frac{\frac{1}{n} \sum_{i=1}^n \lambda_i^{(t)} \mathbb{1}_{\{Y_i \neq 1\}} \mathbb{1}_{\{f(X_i)=1\}}}{\frac{1}{n} \sum_{i=1}^n \lambda_i^{(t)} \mathbb{1}_{\{f(X_i)=1\}}}$$

and

$$\text{TPR}_t = \frac{\frac{1}{n} \sum_{i=1}^n \lambda_i^{(t)} \mathbb{1}_{\{f(X_i)=1\}} \mathbb{1}_{\{Y_i=1\}}}{\frac{1}{n} \sum_{i=1}^n \lambda_i^{(t)} \mathbb{1}_{\{Y_i=1\}}}.$$

Again with  $\Phi$  as in Theorem II.2, a ROC curve corresponding to perturbations of the variable  $j_0$  can then be obtained by plotting pairs  $(\text{FPR}_t, \text{TPR}_t)$  for a large number of  $t \in \mathbb{R}$ . We then obtain the evolution of both errors when  $t$  evolves, which allows a sharper analysis of the error evolution (see *e.g.* Appendix F). Finally, the variables influence on the predictions may be quantified by computing the proportion of predicted 1s

$$\text{P1}_t = \frac{1}{n} \sum_{i=1}^n \lambda_i^{(t)} f(X_i)$$

which we suggest to plot similarly as  $\text{ER}_t$ , with  $\Phi$  as in Theorem II.2 (see Figure 1-(top)). The figures representing  $\text{P1}_t$  make it possible to simply understand the particular influence of the variables to obtain a decision  $Y = 1$ , whatever the veracity of the prediction. Importantly, they point out which variables should be positively or negatively modified in order to change a given decision.

#### B. The case of multi-class classification

We now consider the case of a classification into  $q$  different categories, *i.e.* where  $Y_i$  and  $f(X_i)$  belong to  $\{1, \dots, q\}$  for all  $i = 1, \dots, n$ , and where  $q \in \mathbb{N}$  is fixed.

In this case, the strategy described for the binary classification can naturally be extended using

$$\text{Pj}_t = \frac{1}{n} \sum_{i=1}^n \lambda_i^{(t)} \mathbb{1}_{\{f(X_i)=j\}},$$

which denotes the portion of individuals assigned to the class  $j$ .

#### C. Using quantiles to compare multiple mean changes

Consider the case where  $\Phi$  is given as in Theorem II.2, with  $\Phi(X^1, \dots, X^p, \hat{Y}, Y) = X^{j_0}$ , and where we want to plot the quantities of interest of Sections III-A and III-B, as a function of  $t - t_0$ , for all the values of  $j_0 = 1, \dots, p$ . In this case, an issue is that the interpretation of  $t - t_0$  depends on the order of magnitude of the variable  $j_0$ , and thus changes with  $j_0$ .

In order to compare values of  $t - t_0$  across different variables, we suggest a common parametrization of  $t - t_0$  for  $j_0 = 1, \dots, p$ . We consider the empirical quantile function  $q_{j_0}$  associated to the variable  $X^{j_0}$ , so  $q_{j_0}(\rho) = X_{\sigma([n\rho])}^{j_0}$  for  $0 \leq \rho < 1$  and where  $\sigma(\cdot)$  is a function ordering the sample, *i.e.*  $X_{\sigma(0)}^{j_0} \leq X_{\sigma(1)}^{j_0} \leq \dots \leq X_{\sigma(n-1)}^{j_0}$ . Then, the range of the stressed mean  $t$  will be in  $[q_{j_0}(\alpha), q_{j_0}(1 - \alpha)]$ , where  $\alpha \in (0, 1/2)$  (a typical value is 0.05).

We then tune  $t - t_0$  as equal to  $\epsilon_{j_0, \tau}$ , where  $\epsilon_{j_0, \tau} = \tau(t_0 - q_{j_0}(\alpha))$  if  $\tau \in [-1, 0]$ , and  $\epsilon_{j_0, \tau} = \tau(q_{j_0}(1 - \alpha) - t_0)$  if  $\tau \in [0, 1]$ . Parameter  $\tau$  therefore allows to intuitively parametrize the level of stress whatever the distribution of the  $\{X_1^{j_0}, \dots, X_n^{j_0}\}$ .

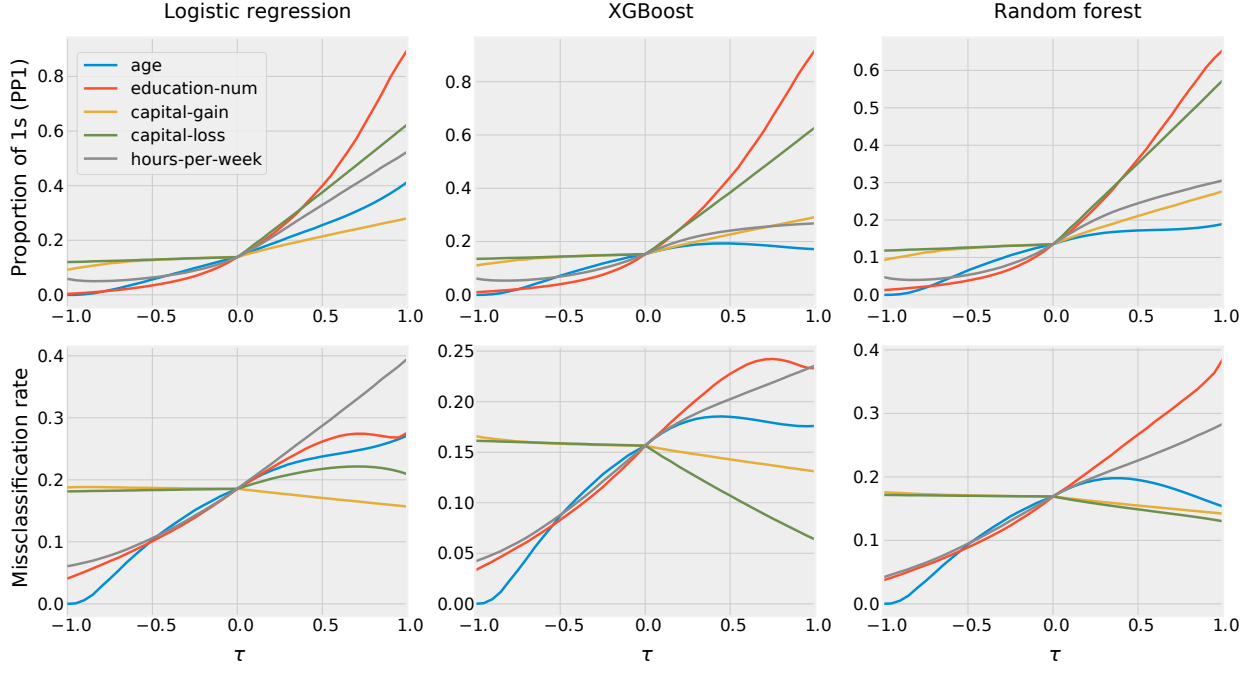


Fig. 1. Results of Section IV-A on the *Adult income* dataset. **(Top - PPIs)** Portion of predicted ones (*i.e.* High Incomes) with respect to the explanatory variable perturbation  $\tau$ . **(Bottom - Mis. Rate)** Classification errors with respect to  $\tau$ . There is no perturbation if  $\tau = 0$ . The larger (resp. the lower)  $\tau$ , the larger (resp. the lower) the values of the selected explanatory variable.

More precisely,  $\tau = 0$  yields no change of mean,  $\tau = -1$  changes the mean from  $t_0$  to the (small) quantile  $q_{j_0}(\alpha)$  and  $\tau = 1$  changes the mean from  $t_0$  to the (large) quantile  $q_{j_0}(1 - \alpha)$ .

For  $j_0 = 1, \dots, p$  and  $\tau \in [-1, 1]$ , we thus naturally suggest to compute

$$ER_{j_0, \tau}, FPR_{j_0, \tau}, TPR_{j_0, \tau}, P1_{j_0, \tau}, Pj_{j_0, \tau} \quad (4)$$

that are defined as  $ER_t, FPR_t, TPR_t, P1_t, Pj_t$  in Sections III-A and III-B, with  $\Phi(X^1, \dots, X^p, \hat{Y}, Y) = X^{j_0}$  and  $t = t_0 + \epsilon_{j_0, \tau}$  as explained above. For a given  $\tau$ , it makes sense to compare the indicators in (4) across  $j_0 = 1, \dots, p$ . For instance, one can plot  $ER_{j_0, \tau}$  as a function of  $\tau$  for  $\tau \in [-1, 1]$  and for each  $j_0 \in \{1, \dots, p\}$ , as shown in Figure 1-(bottom). Remark that  $\tau = 0$  corresponds to the algorithm performance baseline, without any perturbation of the test sample.

## IV. RESULTS

In this section, we illustrate the use of the indices obtained using the entropic projection of samples on two classification cases: In subsection IV-A, the *Adult income* dataset<sup>1</sup> is considered, where  $X$  represents  $n = 32000$  observations of dimension  $p = 14$  and  $Y$  has 2 classes. Results of subsection IV-B are first obtained on the *MNIST* dataset<sup>2</sup>, where  $X$  represents  $n = 60000$  gray level images of  $p = 784$  pixels and  $Y$  has 10 classes, and then obtained on the *Large-scale CelebFaces Attributes (CelebA)* Dataset<sup>3</sup> where  $X$  represents  $n = 200000$  RGB images of 4096 pixels, so  $p = 4096 * 3 = 12288$ , and  $Y$  has 2 classes. Note that the method accuracy is also assessed on synthetic data in Appendix D2 and that the effect of 4 variables on the classification of the well known Iris dataset is shown in Appendix E. Importantly, the Python code to reproduce these experiments is freely available<sup>4</sup>.

### A. Two class classification

In order to illustrate the performance of our procedure, we first consider the *Adult Income* dataset. It is made of  $n = 32000$  observations represented by  $p = 14$  attributes (6 numeric and 9 categorical), each of them describing an adult. We specifically interpret the influence of 5 numeric variables on the categorical output variable representing whether each adult has an income higher ( $Y_i = 1$ ) or lower ( $Y_i = 0$ ) than 50000\$ per year.

<sup>1</sup><https://archive.ics.uci.edu/ml/datasets/adult>

<sup>2</sup><http://yann.lecun.com/exdb/mnist/>

<sup>3</sup><http://mmlab.ie.cuhk.edu.hk/projects/CelebA.html>

<sup>4</sup><https://gems-ai.com/>

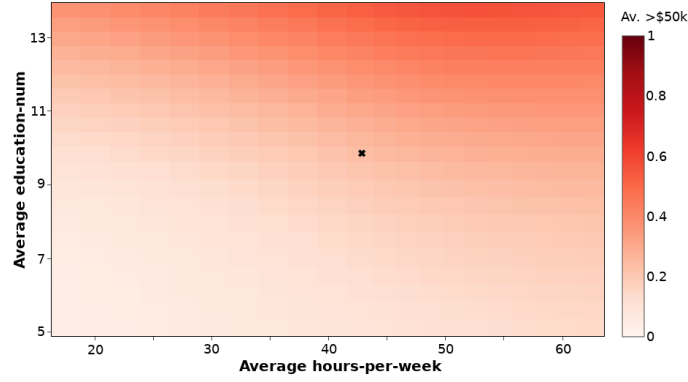


Fig. 2. Result of Section IV-A on the *Adult income* dataset with the *XGboost* classifier. The portion of predicted ones (*i.e.* High Incomes) with respect to the variables *education-num* and *hours-per-week* simultaneously is represented. The black cross represents the average value in the unstressed dataset.

We first trained three different classifiers (Logit Regression, XGboost and Random Forest<sup>5</sup>) on 25600 randomly chosen observations (80% of the whole dataset). We then performed the proposed entropic projection strategy for each learned classifier on a test set made of the 6400 remaining observations. Detailed results are shown in Figure 1. Instead of only quantifying a score for each variable, we display in this figure the evolution of the algorithm confronted at gradually lower or higher values of  $\tau$  (see Section III-C) for each variable. The curves were computed using 21 regularly sampled values of  $\tau$  between  $-1$  and  $1$  with  $\alpha = 0.05$ . For each variable, the weighted observations were therefore stressed so that their mean is contained between the 0.05 and 0.95 quantiles of the original (non-weighted) values distribution in the test set. Note that for a quick and quantified overview of the variables response to a positive (resp. negative) stress, the user can simply interpret the difference of the response for  $\tau = 1$ . and  $\tau = 0$  (resp.  $\tau = 0$  and  $\tau = -1$ .), as illustrated in Section IV-B in the image case.

*a) Influence of the variables in the decision rule:* We present in Figure 1 (*Top*) the role played by each variable in the portion of predicted ones (*i.e.* high incomes) for the *Logit Regression*, *XGboost* and *Random Forest* classifiers. The curves in Figure 1 (*Top*) highlight the role played by the variable *education-num*. The more educated the adult, the higher his/her income will be. The two variables *capital-gain* and *capital-loss* are also testimonial of high incomes since the adults having large incomes can obviously have more money than others on their bank accounts, or may easily contract debts, although the contrary is not true. It is worth pointing out the role played by the age variable which appears clearly in the figure: young adults earn increasingly large incomes with time, which is well captured by the decision rules (left part of the red curves).

We emphasize that these curves enable to intuitively interpret the complex trends captured by *black-box* decision rules. They indeed quantify non-linear effects of the variables and very different behaviors depending on whether the variables increase or decrease.

Contrary to methods that study the influence of the variables by computing information theory criteria between different outputs of the algorithm for changes in the variables (see in Skater [21]), the variable changes we use are plausible since the stressed distributions are as close as possible to the initial distribution. Finally, we work directly on the real *black-box* model and do not approximate it by any surrogate model, as in [31].

It can finally be useful to explain the role played by two variables by stressing them simultaneously. We show for instance in Figure 2 the role played by the variables *education-num* and *hours-per-week* with the *XGboost* classifier. This makes clear the fact that the variable *education-num* has more influence than *hours-per-week* to predict high incomes, as already shown in Figure 1 (*Top*). This additionally shed lights on the fact that the persons with the highest *education-num* values will have slightly decaying predictions if they work more than 52 hours per week, which is not the case for the persons with less than 12 years of academic education.

*b) Influence of the variables in the accuracy of the classifier:* Besides the influence of the variables on the algorithm outcome, it is worth studying their influence on the accuracy or veracity of the model. We then present in Figure 1 (*Bottom*) the evolution of the classification error when each variable is stressed by  $\tau$ . The three sub-figures (one for each prediction model) represent the evolution of the error confronted to the same modification of each variables. The error of the method on the original data is obtained for  $\tau = 0$ . Such curves point out which variables have the strongest influence on prediction errors. Such result may be used to temper the trust in the forecast depending on the values of the variables.

As previously, the curves appear as more informative than single scores: The three models enable to select the same couple of variables that are important for the accuracy of the prediction when they increase *i.e.* education number and numbers of hours worked pro week. The latter makes the prediction task the most difficult when it is increased. Indeed, the persons working a large number of hours per week may not always increase their income, since it relies on different factors. People with high income however usually work a large number of weekly hours. Hence, these two variables play an important role

<sup>5</sup>R command *glm* and packages *xgboost* and *ranger*.

in the prediction and their changes impact the prediction error. In the same flavor, more insight on the error terms could be obtained by dealing with the evolution of the False Positive Rate and True Positive Rate as presented in Appendix F.

### B. Image classification

1) *MNIST dataset*: We now measure the influence of pixel intensities in image recognition tasks. Each pixel intensity is treated as a variable and the stress is used to saturate the intensities towards one side of their spectrum (red if  $\tau = 1$  or blue if  $\tau = -1$ ).

We specifically trained a CNN (Convolutional Neuron Network) on the MNIST dataset using a typical architecture that can be found on the Keras documentation<sup>6</sup>. The CNN was trained on a set of 60000 images whilst the predictions were made on another 10000 images. It achieved a test set accuracy north of 99%. For each of the 784 pixels  $j_0$ , we computed the perturbed distributions in the cases where  $\tau$  equals  $-1$  as well as  $1$ , using the method of Section III-C. The prediction proportions of each of the 10 digits was then computed using the method of Section III-B. The whole process took around 9 seconds to run on a modern laptop (Intel Core i7-8550U CPU @ 1.80GHz, 24GB RAM) running Linux. The results are presented in Figure 3-(top).

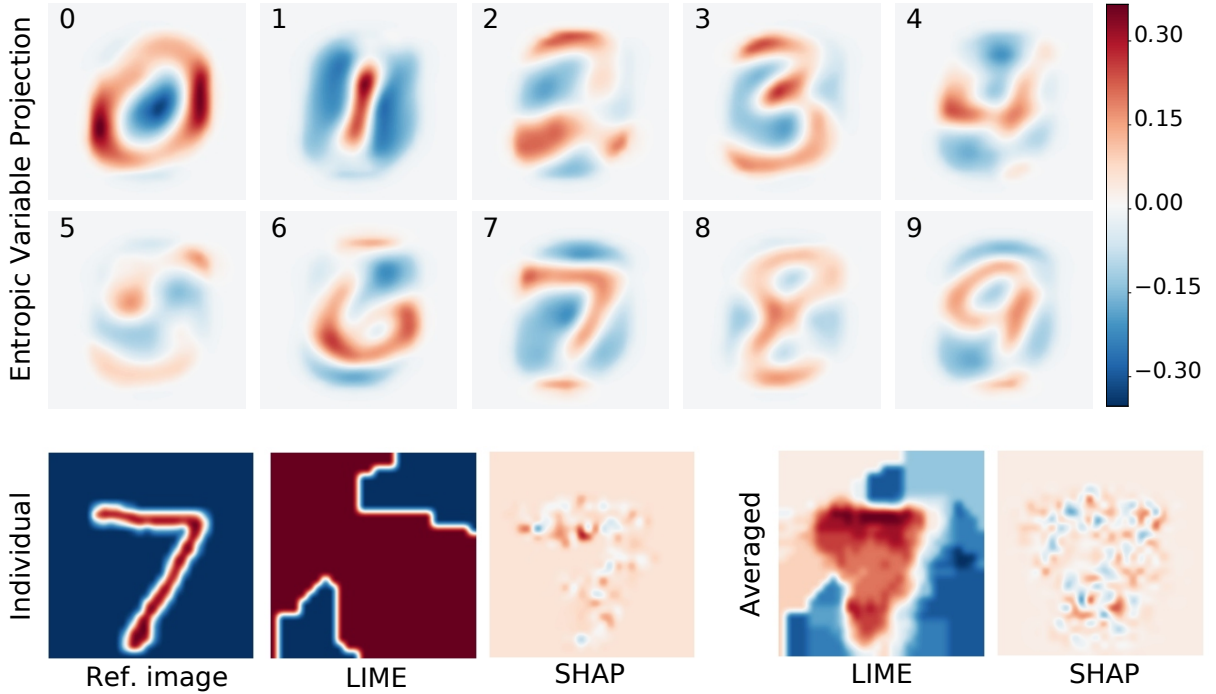


Fig. 3. **(top)** Pixel contributions towards each digit according to our entropic projection method. **(bottom-left)** Pixel contributions to predict seven in an individual image representing a seven, using the LIME and SHAP packages. **(bottom-right)** Average pixel contributions to predict seven in all images of the MNIST test set representing a seven, using the LIME and SHAP packages.

The color of each pixel in Figure 3-(top) represents its contribution towards the prediction of each digit. For example, a value of 0.15 means that the CNN predicts on average this digit 15% more often when the associated pixel is activated ( $\tau = 1$ ) instead of having the background intensity ( $\tau = -1$ ). Although our method is pixel-based, it is still able to uncover regions which the CNN uses to predict each digit. Likewise, redder regions contain pixels that are positively correlated with each digit. Note that the edges of each image don't change color because the corresponding pixels have no impact whatsoever on the predictions. The left part of number 5 has pixels in common with number 6. However, we are able to see that the CNN identifies 6s by using the bottom part of the 6, more so than the top stroke which it uses to recognize 5s. Likewise, according to the CNN, the most distinguishing part of number 9 is the part that links the top ring with the bottom stroke.

We finally emphasize the main difference between our strategy and the two popular interpretability solutions LIME ([31])<sup>7</sup> and SHAP ([27])<sup>8</sup>: we work on whole test sets while these solutions interpret the variables (here pixels) influence when predicting specific labels in *individual observations*. As an illustration, Figure 3-(bottom-left) represents the most influential pixels found with LIME and SHAP to predict a seven in an image of the MNIST test set representing a seven. The results

<sup>6</sup>[https://keras.io/examples/mnist\\_cnn](https://keras.io/examples/mnist_cnn)

<sup>7</sup><https://github.com/marcotcr/lime>

<sup>8</sup><https://github.com/slundberg/shap>



are not straightforward to interpret. To draw similar interpretations as those made on Figure 3-(top), one can represent the average results obtained by using LIME or SHAP over all images of the MNIST test set representing a seven, as represented in Figure 3-(bottom-right). Note that the computations required take about 7 and 10 hours using LIME and SHAP, respectively, which is much longer than when using our method (10 seconds). Compared to our method, averaged results are also less resolved for LIME and harder to interpret for SHAP. Our method can also natively compute other properties of the black-box decision rules with a negligible additional computational cost, as described in Section III.

2) *CelebA dataset*: We now present further results obtained on the *Large-scale CelebFaces Attributes (CelebA)* Dataset [25]. It contains more than 200000 celebrity RGB images, each with 40 binary attribute annotations, such as *Eyeglasses*, *PaleSkin*, *Smiling*, *Young*, *Male* or *Attractive*. It is important to mention here that the images of the CelebA dataset cover large pose variations and background clutter. This makes this dataset far more complex to analyze than the MNIST dataset. In order to explain the decision rules of a state-of-the-art neural-network architecture, we trained a ResNet-18 convolutional Neural Network [16] to predict the attribute *Attractive* based on the CelebA images. The portion of good predictions on the training set was 0.95 and 0.92 for females and males, respectively. This portion was also 0.86 and 0.78 on the test set for females and males, which is good for such a subjective attribute. This suggests that the persons who labeled the data were relatively coherent when choosing who was attractive or not.

We then used almost the same strategy as on the MNIST dataset to explain the decision rules learned by the ResNet-18 neural-network. In order to show the flexibility of its decision rules in two different contexts, we however used the entropic projection strategy on the images representing females and males separately. The results are presented in Figure 4. Note that they were obtained by using all the 200000 images of the dataset and required about three minutes of computations. As shown in Figure 4, it is clear that the trained neural-network takes more into account the hairs for males than for females to make its predictions. It also uses more the contour of the face for females than for males (in particular for the pixels corresponding to the cheeks).

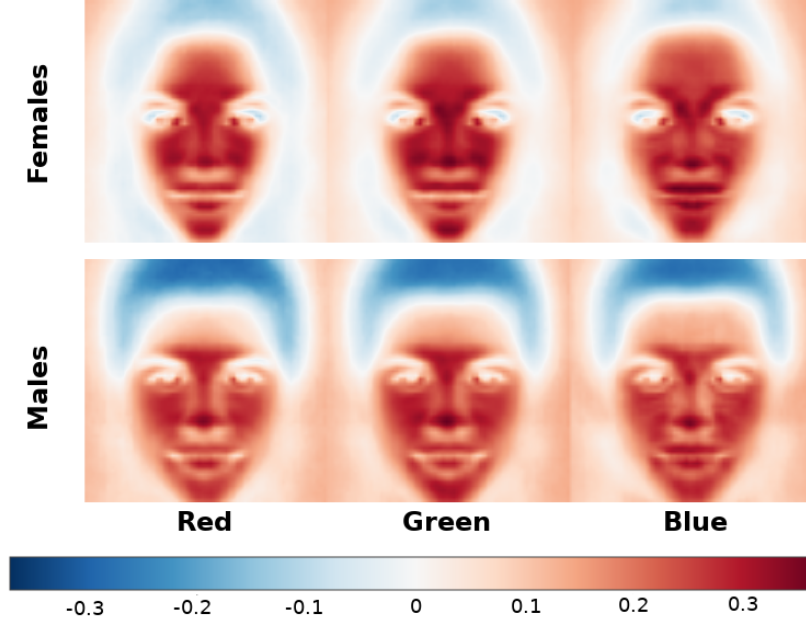


Fig. 4. Pixel contributions to explain the decision rules of a ResNet-18 neural network trained to predict whether someone is considered as attractive or not. The flexibility of the decision rules in different contexts is emphasized by using the entropic projection strategy on images representing females and males separately.

## V. CONCLUSION

Explainability of *black-box* decision rules in the machine learning paradigm has many interpretations and has been tackled in a large variety of contributions. Here, we focused on the analysis of the variables importance and their impact on a decision rule, at the global scale of the entire dataset. We satisfy the important constraint, in the machine learning context, that the test input variables must follow the distribution of the learning sample, as closely as possible. Therefore, evaluating the decision rule at any possible point does not make any sense. To cope with this issue, we have proposed an information theory procedure to bias the original variables without losing the information conveyed by the initial distribution. We proved that this solution amounts in re-weighting the observations of the test sample, leading to very fast computations and the construction of new indices to make clear the role played by each variable.

Remark that our strategy can be seen as a *what if* tool, as counterfactual methods [37]. It indeed explains model decisions by quantifying how their outputs change when the machine learning data are transformed. Nevertheless, existing counterfactual methods substitute individual counterfactual observations for individual baseline observations. In contrast, our strategy substitutes counterfactual data sets, with new variable distribution characteristics, for the original data set.

The first key advantage of this strategy is to preserve as much as possible the distribution of the test set  $(X_i^1, \dots, X_i^p, \hat{Y}_i, Y_i)$ ,  $i = 1, \dots, n$  and thus preserving the correlations of the input variables that have then an impact on the indicators computed by using our procedure. In contrast with other interpretability paradigms such as the PAC learning framework [35], we do not create artificial outliers. Its second key advantage is that, for a given perturbation, the weights are obtained by minimizing a convex function, for which the evaluation cost is  $\mathcal{O}(n)$ . The total cost is then  $\mathcal{O}(np)$  for studying the impact of each of the  $p$  variables (see Appendix D1) and there is no need to generate new data, nor even to compute new predictions from the black box algorithm, which is particularly costly if  $n$  or  $p$  is large. Our procedure therefore scales particularly well to large datasets as e.g. real-world image databases. Finally, the flexibility of the entropic variable projection procedure enables to study the response to various types of stress on the input variables (not only their mean but also their variability, joint correlations, ...) and thus to interpret the decision rules encountered in a wide range of applications encountered in the field of Machine Learning. A package in Python with other examples and industrial use cases is available at <https://gems-ai.com/>.

**Acknowledgements :** This project received funding from the French Investing for the Future PIA3 program within the Artificial and Natural Intelligence Toulouse Institute (ANITI).

## REFERENCES

- [1] F. Bachoc, F. Gamboa, J.-M. Loubes, and N. Venet. A gaussian process regression model for distribution inputs. *IEEE Transactions on Information Theory*, 64(10):6620–6637, 2017.
- [2] D. Baehrens, T. Schroeter, S. Harmeling, M. Kawanabe, K. Hansen, and K.-R. Müller. How to explain individual classification decisions. *Journal of Machine Learning Research*, 11(Jun):1803–1831, 2010.
- [3] E. Black, S. Yeom, and M. Fredrikson. Fliptest: fairness testing via optimal transport. In *FAT\* '20: Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pages 111–121. 2020.
- [4] L. Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- [5] P. Bühlmann and S. Van De Geer. *Statistics for High-Dimensional Data*. Springer, 2011.
- [6] R. Caruana, Y. Lou, J. Gehrke, P. Koch, M. Sturm, and N. Elhadad. Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '15, pages 1721–1730, 2015.
- [7] E. Cazelles, V. Seguy, J. Bigot, M. Cuturi, and N. Papadakis. Geodesic pca versus log-pca of histograms in the Wasserstein space. *SIAM Journal on Scientific Computing*, 40(2):B429–B456, 2018.
- [8] M. Craven and J. W. Shavlik. Extracting tree-structured representations of trained networks. In *Advances in Neural Information Processing Systems 8, NIPS, Denver, CO, November 27-30, 1995*, pages 24–30, 1995.
- [9] I. Csiszár.  $I$ -divergence geometry of probability distributions and minimization problems. *The Annals of Probability*, pages 146–158, 1975.
- [10] I. Csiszár. Sanov property, generalized  $I$ -projection and a conditional limit theorem. *The Annals of Probability*, pages 768–793, 1984.
- [11] M. T. Dzindolet, S. A. Peterson, R. A. Pomranky, L. G. Pierce, and H. P. Beck. The role of trust in automation reliance. *International journal of human-computer studies*, 58(6):697–718, 2003.
- [12] N. Fournier and A. Guillin. On the rate of convergence in Wasserstein distance of the empirical measure. *Probability Theory and Related Fields*, 162(3-4):707–738, 2015.
- [13] Y. Goyal, Z. Wu, J. Ernst, D. Batra, D. Parikh, and S. Lee. Counterfactual visual explanations. In *Proceedings of the 36th International Conference on Machine Learning*, pages 2376–2384, 2019.
- [14] P. Hall, N. Gill, and M. Chan. Practical techniques for interpreting machine learning models: Introductory open source examples using python, h2o, and xgboost, 2018.
- [15] T. Hastie, R. Tibshirani, and J. H. Friedman. *The elements of statistical learning: data mining, inference, and prediction*. New York, NY: Springer, 2009.
- [16] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
- [17] J. L. Herlocker, J. A. Konstan, and J. Riedl. Explaining collaborative filtering recommendations. In *Proceedings of the 2000 ACM conference on Computer supported cooperative work*, pages 241–250. ACM, 2000.
- [18] S. Hooker, D. Erhan, P.-J. Kindermans, and B. Kim. A benchmark for interpretability methods in deep neural networks. In *Advances in Neural Information Processing Systems*, pages 9737–9748, 2019.
- [19] A. Ignatiev, N. Narodytska, and J. Marques-Silva. On relating explanations and adversarial examples. In *Advances in Neural Information Processing Systems 32*, pages 15857–15867. 2019.

- [20] P. W. Koh and P. Liang. Understanding black-box predictions via influence functions. In *Proceedings of the 34th International Conference on Machine Learning (ICML 2017)*, 2017.
- [21] A. Kramer and P. Choudhary. Model Interpretation with Skater. <https://oracle.github.io/Skater/>, 2018. [Online; accessed 28-Jan-2020].
- [22] M. J. Kusner, J. Loftus, C. Russell, and R. Silva. Counterfactual fairness. In *Advances in Neural Information Processing Systems 30*, pages 4066–4076. 2017.
- [23] P. Lemaître, E. Sergienko, A. Arnaud, N. Bousquet, F. Gamboa, and B. Iooss. Density modification-based reliability sensitivity analysis. *Journal of Statistical Computation and Simulation*, 85(6):1200–1223, 2015.
- [24] Z. C. Lipton. The mythos of model interpretability. In *Proceedings of the ICML Workshop on Human Interpretability in Machine Learning (WHI 2016)*, pages 96–100, 2016.
- [25] Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- [26] Y. Lou, R. Caruana, and J. Gehrke. Intelligible models for classification and regression. In *Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '12*, pages 150–158, 2012.
- [27] S. M. Lundberg and S.-I. Lee. A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems (NIPS)*, pages 4765–4774, 2017.
- [28] G. Montavon, W. Samek, and K.-R. Müller. Methods for interpreting and understanding deep neural networks. *Digital Signal Processing*, 2017.
- [29] G. Peyré and M. Cuturi. Computational optimal transport: With applications to data science. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019.
- [30] L. E. Raileanu and K. Stoffel. Theoretical comparison between the gini index and information gain criteria. *Annals of Mathematics and Artificial Intelligence*, 41(1):77–93, 2004.
- [31] M. T. Ribeiro, S. Singh, and C. Guestrin. Why should I trust you?: Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1135–1144. ACM, 2016.
- [32] A. Saltelli, M. Ratto, T. Andres, F. Campolongo, J. Cariboni, D. Gatelli, M. Saisana, and S. Tarantola. *Global sensitivity analysis: the primer*. John Wiley & Sons, 2008.
- [33] R. R. Selvaraju, A. Das, R. Vedantam, M. Cogswell, D. Parikh, and D. Batra. Grad-cam: Why did you say that? *arXiv preprint arXiv:1611.07450*, 2016.
- [34] M. Sundararajan, A. Taly, and Q. Yan. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70, ICML17*, page 33193328, 2017.
- [35] L. Valiant. *Probably Approximately Correct: Nature’s Algorithms for Learning and Prospering in a Complex World*. Basic Books (AZ), 2013.
- [36] C. Villani. *Optimal transport: old and new*, volume 338. Springer Science & Business Media, 2008.
- [37] S. Wachter, B. Mittelstadt, and C. Russell. Counterfactual explanations without opening the black box: Automated decisions and the gdpr. *Harvard journal of law & technology*, 31:841–887, 04 2018.

## APPENDIX

### A. Proofs of the main results

The proofs rely on the following theorem, that is a simplified version of the Theorems in [9] and in [10].

**Theorem A.1.** *Let  $(E, \mathcal{B}(E))$  be a measurable space and  $Q$  a probability measure on  $E$ . Consider  $t \in \mathbb{R}^k$  and a measurable function  $\Phi : E \rightarrow \mathbb{R}^k$ . We assume that, for  $v \in \mathbb{R}^k$ ,  $b \in \mathbb{R}$ ,  $Q(\{x \in E; \langle v, \Phi(x) \rangle = b\}) = 1$  if and only if  $v = 0$  and  $b = 0$ . Let  $\mathbb{P}_{\Phi, t}$  be the set of all probability measures  $P$  on  $(E, \mathcal{B}(E))$  such that  $\int_E \Phi(x) dP(x) = t$ . Assume that  $\mathbb{P}_{\Phi, t}$  contains a probability measure that is mutually absolutely continuous with respect to  $Q$ .*

*For a vector  $\xi \in \mathbb{R}^k$ , let  $Z(\xi) := \int_E e^{\langle \xi, \Phi(x) \rangle} dQ(x)$ . We assume that the set on which  $Z$  is finite is open. Define now  $\xi(t)$  as the unique minimizer of the strictly convex function  $H(\xi) := \log Z(\xi) - \langle \xi, t \rangle$ . Then,*

$$Q_t := \operatorname{arginf}_{P \in \mathbb{P}_{\Phi, t}} \operatorname{KL}(P, Q) \tag{5}$$

*exists and is unique. Furthermore it can be computed as*

$$Q_t = \frac{\exp\langle \xi(t), \Phi \rangle}{Z(\xi(t))} Q.$$

**Proof of Theorem II.1** We will apply Theorem A.1 with  $E = \mathbb{R}^{p+2}$  and  $Q = Q_n$ . Because of the assumption that  $t$  can be written as a convex combination of  $\Phi(X_1, \hat{Y}_1, Y_1), \dots, \Phi(X_n, \hat{Y}_n, Y_n)$ , with positive weights, we have that  $\mathbb{P}_{\Phi, t}$  in Theorem A.1 contains a probability measure that is mutually absolutely continuous with respect to  $Q$ . Furthermore, we have assumed that the empirical covariance matrix of  $\Phi(X_1, \hat{Y}_1, Y_1), \dots, \Phi(X_n, \hat{Y}_n, Y_n)$  is invertible, which means that for any  $v \in \mathbb{R}^k$ ,  $b \in \mathbb{R}$ ,

$\langle v, \Phi(X_1, \hat{Y}_1, Y_1) \rangle, \dots, \langle v, \Phi(X_n, \hat{Y}_n, Y_n) \rangle$  are not all equal to  $b$ . This implies that  $Q_n(\{x \in \mathbb{R}^{p+2}; \langle v, \Phi(x) \rangle = b\}) = 1$  if and only if  $v = 0$  and  $b = 0$ . Hence, all the assumptions of Theorem A.1 are satisfied.

We have, starting from the notation of Theorem A.1,

$$\int_E e^{\langle \xi, \Phi \rangle} dQ(x) = \frac{1}{n} \sum_{i=1}^n e^{\langle \xi, \Phi(X_i, \hat{Y}_i, Y_i) \rangle}$$

and thus the definitions of  $Z(\xi)$  in Theorems A.1 and II.1 indeed coincide. Hence, also the definitions of  $\xi(t)$  in Theorems A.1 and II.1 coincide. Hence, we have

$$\begin{aligned} Q_t &= \frac{\exp\langle \xi(t), \Phi \rangle}{Z(\xi(t))} Q \\ &= \frac{1}{\frac{1}{n} \sum_{i=1}^n e^{\langle \xi(t), \Phi(X_i, \hat{Y}_i, Y_i) \rangle}} \\ &= \frac{1}{n} \sum_{i=1}^n \exp\left(\langle \xi(t), \Phi(X_i, \hat{Y}_i, Y_i) \rangle\right) \delta_{X_i, \hat{Y}_i, Y_i} \\ &= \frac{1}{n} \sum_{i=1}^n \lambda_i^{(t)} \delta_{X_i, \hat{Y}_i, Y_i}. \end{aligned}$$

This concludes the proof.  $\square$

**Proof of Theorem II.2** The proof of this theorem comes from Theorem II.1, by considering  $\Phi : \mathbb{R}^{p+2} \rightarrow \mathbb{R}$  defined by  $\Phi(X^1, \dots, X^p, \hat{Y}, Y) = X^{j_0}$  and by considering the same  $t \in \mathbb{R}$  in Theorems II.1 and II.2. We have assumed that  $\min_{i=1}^n X_i^{j_0} < t < \max_{i=1}^n X_i^{j_0}$ . Hence,  $t$  can be written as a convex combination of  $X_1^{j_0}, \dots, X_n^{j_0}$  with positive weights. Furthermore,  $X_1^{j_0}, \dots, X_n^{j_0}$  are not all equal and thus their empirical variance is non-zero. Hence the conditions of Theorem II.1 hold and the conclusion of this theorem directly provides Theorem II.2.  $\square$

**Proof of Corollary II.3** With the choice of  $\Phi$  of the corollary and with the assumptions there, the conditions of Theorem II.1 hold. Hence the conclusion of this theorem proves the corollary.  $\square$

**Proof of Corollary II.4** With the choice of  $\Phi$  of the corollary and with the assumptions there, the conditions of Theorem II.1 hold. Hence the conclusion of this theorem proves the corollary, since  $(\mathbb{E}(X^{j_0}) = m_{j_0}, \text{Var}(X^{j_0}) = v)$  is equivalent to  $(\mathbb{E}(X^{j_0}) = m_{j_0}, \mathbb{E}((X^{j_0})^2) = m_{j_0}^2 + v)$ .  $\square$

**Proof of Corollary II.5** With the choice of  $\Phi$  of the corollary and with the assumptions there, the conditions of Theorem II.1 hold. Hence the conclusion of this theorem proves the corollary, since  $(\mathbb{E}(X^{j_1}) = m_{j_1}, \mathbb{E}(X^{j_2}) = m_{j_2}, \text{Cov}(X^{j_1}, X^{j_2}) = c)$  is equivalent to  $\mathbb{E}(X^{j_1}) = m_{j_1}, \mathbb{E}(X^{j_2}) = m_{j_2}, \mathbb{E}(X^{j_1} X^{j_2}) = m_{j_1} m_{j_2} + c$ .  $\square$

**Proof of Proposition II.6** The existence and unicity of  $Q_t^*$  follows from Theorem A.1. Also from this theorem,  $Q_t^*$  is of the form

$$dQ_t^*(x) = e^{\langle \xi^*(t), \Phi(x) \rangle - \log(Z^*(\xi^*(t)))} dQ^*(x),$$

where  $\xi^*(t)$  is the minimizer of the strictly convex function

$$\xi \mapsto F(\xi) := \log(Z^*(\xi)) - \langle \xi, t \rangle$$

with

$$Z^*(\xi) = \int_{\mathbb{R}^{p+2}} e^{\langle \xi, \Phi(x) \rangle} dQ^*(x).$$

We also recall from Theorem II.1 that  $Q_t$  is of the form

$$dQ_t(x) = e^{\langle \xi(t), \Phi(x) \rangle - \log(Z(\xi(t)))} dQ_n(x),$$

where  $\xi(t)$  is the minimizer of the strictly convex function

$$\xi \mapsto F_n(\xi) := \log(Z(\xi)) - \langle \xi, t \rangle$$

with

$$Z(\xi) = \int_{\mathbb{R}^{p+2}} e^{\langle \xi, \Phi(x) \rangle} dQ_n(x).$$

Let us first prove that  $n^{1/2}(\xi^*(t) - \xi(t)) = O_p(1)$ . The gradient of  $F$  at  $\xi$  can be computed as

$$\nabla_\xi F(\xi) = \frac{\int_{\mathbb{R}^{p+2}} \Phi(x) e^{\langle \xi, \Phi(x) \rangle} dQ^*(x)}{\int_{\mathbb{R}^{p+2}} e^{\langle \xi, \Phi(x) \rangle} dQ^*(x)} - t.$$

The norm of this gradient is bounded by a constant  $C_{\text{sup},1} < \infty$  from the assumption that  $\Phi$  is bounded on the support of  $Q^*$ .

The Hessian matrix of  $F$  at  $\xi$  is

$$H_\xi F(\xi) = \frac{\int_{\mathbb{R}^{p+2}} \Phi(x) \Phi(x)^\top e^{\langle \xi, \Phi(x) \rangle} dQ^*(x)}{\int_{\mathbb{R}^{p+2}} e^{\langle \xi, \Phi(x) \rangle} dQ^*(x)} - \frac{\int_{\mathbb{R}^{p+2}} \Phi(x) e^{\langle \xi, \Phi(x) \rangle} dQ^*(x)}{\int_{\mathbb{R}^{p+2}} e^{\langle \xi, \Phi(x) \rangle} dQ^*(x)} \frac{\int_{\mathbb{R}^{p+2}} \Phi(x)^\top e^{\langle \xi, \Phi(x) \rangle} dQ^*(x)}{\int_{\mathbb{R}^{p+2}} e^{\langle \xi, \Phi(x) \rangle} dQ^*(x)}.$$

The two last above fractions are  $t$  and  $t^\top$  when  $\xi = \xi^*(t)$ , because  $\nabla_\xi F(\xi^*(t))$  is zero because  $F$  is minimal at  $\xi^*(t)$ . Hence, we obtain

$$H_\xi F(\xi^*(t)) = \frac{\int_{\mathbb{R}^{p+2}} \Phi(x) \Phi(x)^\top e^{\langle \xi^*(t), \Phi(x) \rangle} dQ^*(x)}{\int_{\mathbb{R}^{p+2}} e^{\langle \xi^*(t), \Phi(x) \rangle} dQ^*(x)} - tt^\top.$$

Hence,  $H_\xi F(\xi^*(t))$  is the covariance matrix of  $\Phi(\cdot)$ , under the probability distribution proportional to  $e^{\langle \xi^*(t), \Phi(\cdot) \rangle} dQ^*(\cdot)$  on  $\mathbb{R}^{p+2}$ . Assume that  $H_\xi F(\xi)$  is not invertible. Then there exists an affine subspace  $G$  of  $\mathbb{R}^k$  with dimension strictly less than  $k$  that contains  $\Phi(\cdot)$  almost surely under the probability distribution proportional to  $e^{\langle \xi^*(t), \Phi(\cdot) \rangle} dQ^*(\cdot)$ . Since  $\Phi$  is bounded on the support of  $Q^*$ , this distribution is equivalent to  $dQ^*(\cdot)$  and thus  $\phi(\cdot) \in G$  almost surely under the probability distribution  $dQ^*(\cdot)$ . This is in contradiction with an assumption of Proposition II.6. Hence,  $H_\xi F(\xi^*(t))$  is invertible.

Let  $B(\xi, r)$  be the Euclidean ball with center  $\xi$  and radius  $r$ . One can see that the partial derivatives of  $H_\xi F(\xi)$  are bounded on  $B(\xi^*(t), \delta)$  using that  $\Phi$  is bounded on the support of  $Q^*$ . Hence, there are constants  $C_{\text{inf},2} > 0$  and  $\delta > 0$  such that

$$\inf_{\xi \in B(\xi^*(t), \delta)} \lambda_{\text{inf}}(H_\xi F(\xi^*(t))) \geq C_{\text{inf},2}, \quad (6)$$

with  $\lambda_{\text{inf}}(\cdot)$  the smallest eigenvalue.

Furthermore, we can compute similarly the Hessian matrix of  $F_n$  at  $\xi$ ,

$$H_\xi F_n(\xi) = \frac{\int_{\mathbb{R}^{p+2}} \Phi(x) \Phi(x)^\top e^{\langle \xi, \Phi(x) \rangle} dQ_n(x)}{\int_{\mathbb{R}^{p+2}} e^{\langle \xi, \Phi(x) \rangle} dQ_n(x)} - \frac{\int_{\mathbb{R}^{p+2}} \Phi(x) e^{\langle \xi, \Phi(x) \rangle} dQ_n(x)}{\int_{\mathbb{R}^{p+2}} e^{\langle \xi, \Phi(x) \rangle} dQ_n(x)} \frac{\int_{\mathbb{R}^{p+2}} \Phi(x)^\top e^{\langle \xi, \Phi(x) \rangle} dQ_n(x)}{\int_{\mathbb{R}^{p+2}} e^{\langle \xi, \Phi(x) \rangle} dQ_n(x)}.$$

Hence, one can see that the partial derivatives of  $H_\xi F_n(\xi)$  are bounded uniformly in  $\xi \in B(\xi^*(t), \delta)$ , with a fixed deterministic bound, using again that  $\Phi$  is bounded on the support of  $Q^*$ . Furthermore, for any fixed  $\xi$ , from the law of large number  $H_\xi F_n(\xi) \rightarrow H_\xi F(\xi)$  almost surely. Hence, we can show

$$\sup_{\xi \in B(\xi^*(t), \delta)} |\lambda_{\text{inf}}(H_\xi F_n(\xi)) - \lambda_{\text{inf}}(H_\xi F(\xi))| = o_p(1). \quad (7)$$

Let us consider  $M > 0$  to be fixed later. By convexity, we have

$$\begin{aligned} & \mathbb{P} \left( \|\xi(t) - \xi^*(t)\| \geq \frac{M}{\sqrt{n}} \right) \leq \\ & \mathbb{P} \left( \inf_{\|\xi - \xi^*(t)\| = M/n^{1/2}} (F_n(\xi) - F_n(\xi^*(t))) \leq 0 \right). \end{aligned} \quad (8)$$

Let us bound the probability of this last event. We have, for some random  $\hat{\xi} \in B(\xi^*(t), M/n^{1/2})$  and  $\hat{\xi}$  with  $\|\hat{\xi} - \xi^*(t)\| = M/n^{1/2}$ ,

$$\begin{aligned} & \inf_{\|\xi - \xi^*(t)\| = M/n^{1/2}} (F_n(\xi) - F_n(\xi^*(t))) \\ & = (\nabla_\xi F_n(\xi^*(t)))^\top (\hat{\xi} - \xi^*(t)) + \frac{1}{2} (\hat{\xi} - \xi^*(t))^\top H_\xi F_n(\hat{\xi}) (\hat{\xi} - \xi^*(t)). \end{aligned}$$

Since  $\nabla_\xi F(\xi^*(t)) = 0$  we can show simply  $\|\nabla_\xi F_n(\xi^*(t))\| = O_p(n^{-1/2})$ . Furthermore, from (7) and (6), we obtain

$$\begin{aligned} & \inf_{\|\xi - \xi^*(t)\| = M/n^{1/2}} (F_n(\xi) - F_n(\xi^*(t))) \\ & \geq -O_p \left( \frac{1}{\sqrt{n}} \right) \frac{M}{\sqrt{n}} + \frac{1}{2} \left( \frac{M}{\sqrt{n}} \right)^2 (C_{\text{inf},2} + o_p(1)). \end{aligned}$$

The above  $O_p(1)$  and  $o_p(1)$  can be taken to be independent on  $M$ . Hence, for any  $\eta > 0$  we can take  $M$  large enough such that (8) is smaller than  $\eta$ . Hence we have shown

$$\|\xi^*(t) - \xi(t)\| = O_p\left(\frac{1}{\sqrt{n}}\right). \quad (9)$$

We have

$$\begin{aligned} & Z^*(\xi^*(t)) - Z(\xi(t)) \\ &= \int_{\mathbb{R}^{p+2}} e^{\langle \Phi(x), \xi^*(t) \rangle} dQ^* - \int_{\mathbb{R}^{p+2}} e^{\langle \Phi(x), \xi(t) \rangle} dQ_n \\ &= \int_{\mathbb{R}^{p+2}} e^{\langle \Phi(x), \xi^*(t) \rangle} dQ^* - \int_{\mathbb{R}^{p+2}} e^{\langle \Phi(x), \xi^*(t) \rangle} dQ_n \end{aligned} \quad (10)$$

$$+ \int_{\mathbb{R}^{p+2}} e^{\langle \Phi(x), \xi^*(t) \rangle} dQ_n - \int_{\mathbb{R}^{p+2}} e^{\langle \Phi(x), \xi(t) \rangle} dQ_n. \quad (11)$$

The quantity in (10) is a  $O_p(n^{-1/2})$  by the central limit theorem since  $\Phi$  is bounded on the support of  $Q^*$ . The quantity in (11) is a  $O_p(n^{-1/2})$  from (9) and because  $\Phi$  is bounded on the support of  $Q^*$ . Hence, we have shown

$$Z^*(\xi^*(t)) - Z(\xi(t)) = O_p\left(\frac{1}{\sqrt{n}}\right). \quad (12)$$

Let now  $\mathcal{F} = \{f : \mathbb{R}^{p+2} \rightarrow \mathbb{R}; f \text{ is 1-Lipschitz}, f(0) = 0\}$ . We have

$$\begin{aligned} & \sup_{f \in \mathcal{F}} \left| \int_{\mathbb{R}^{p+2}} f(x) dQ_t^*(x) - \int_{\mathbb{R}^{p+2}} f(x) dQ_t(x) \right| \\ & \leq \frac{1}{Z^*(\xi^*(t))} \sup_{f \in \mathcal{F}} \left| \int_{\mathbb{R}^{p+2}} f(x) e^{\langle \xi^*(t), \Phi(x) \rangle} \right. \\ & \quad \left. (dQ^*(x) - dQ_n(x)) \right| \end{aligned} \quad (13)$$

$$\begin{aligned} & + \frac{1}{Z^*(\xi^*(t))} \sup_{f \in \mathcal{F}} \left| \int_{\mathbb{R}^{p+2}} f(x) e^{\langle \xi^*(t), \Phi(x) \rangle} dQ_n(x) \right. \\ & \quad \left. - \int_{\mathbb{R}^{p+2}} f(x) e^{\langle \xi(t), \Phi(x) \rangle} dQ_n(x) \right| \end{aligned} \quad (14)$$

$$\begin{aligned} & + \left| \frac{1}{Z^*(\xi^*(t))} - \frac{1}{Z(\xi(t))} \right| \\ & \sup_{f \in \mathcal{F}} \left| \int_{\mathbb{R}^{p+2}} f(x) e^{\langle \xi(t), \Phi(x) \rangle} dQ_n(x) \right|. \end{aligned} \quad (15)$$

The term (14) is smaller than a constant times  $\|\xi^*(t) - \xi(t)\|$  because  $\Phi$  is bounded on the support of  $Q^*$  and  $f$  is bounded. Hence this term is a  $O_p(n^{-1/2})$  from (9). The term in (15) is also a  $O_p(n^{-1/2})$  from (12). Let us finally address the term in (13). In this term, the function that is integrated is uniformly bounded, with uniformly bounded Lipschitz norm, because  $\Phi$  is bounded and Lipschitz continuous on the bounded support of  $Q^*$  and since  $f \in \mathcal{F}$ . Also, from for instance Theorem 1 in [12], the  $L^1$  Wasserstein distance between  $Q_n$  and  $Q^*$  satisfies

$$\mathcal{W}_1(Q_n, Q^*) = O_p\left(n^{-1/(p+2)}\right).$$

This implies that the supremum over  $f \in \mathcal{F}$  in (13) is bounded by a constant times  $O_p\left(n^{-1/(p+2)}\right)$ , see for instance [36] for the link between the  $L^1$  Wasserstein distance and differences of expectations of Lipschitz functions. Hence we have proved that as  $n \rightarrow \infty$ ,

$$\sup_{f \in \mathcal{F}} \left| \int_{\mathbb{R}^{p+2}} f(x) dQ_t^*(x) - \int_{\mathbb{R}^{p+2}} f(x) dQ_t(x) \right| = O_p\left(n^{-1/(p+2)}\right).$$

From for instance [36], this implies the conclusion of the proposition.  $\square$

### B. Gradient of the objective function in Theorem II.1

Let us denote  $v_i = (X_i, \hat{Y}_i, Y_i)$  for  $i = 1, \dots, n$ . In Theorem II.1 we want to minimize, over  $\xi \in \mathbb{R}^k$ ,

$$H(\xi) = \log \left( \frac{1}{n} \sum_{i=1}^n \exp(\langle \xi, \Phi(v_i) \rangle) \right) - \langle \xi, t \rangle. \quad (16)$$

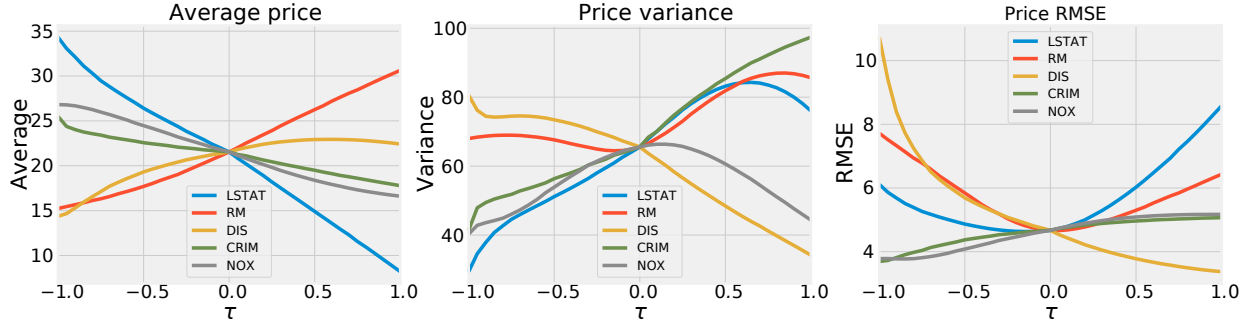


Fig. 5. Results obtained on the *Boston Housing* dataset with Random Forest. The explanatory variable perturbation  $\tau$  has the same signification as in Figure 1.

The gradient of (16) is:

$$\nabla_{\xi} H(\xi) = \frac{\sum_{i=1}^n \Phi(v_i) \exp \langle \xi, \Phi(v_i) \rangle}{\sum_{i=1}^n \exp \langle \xi, \Phi(v_i) \rangle} - t, \quad (17)$$

which makes it possible to compute  $\xi(t)$  in Theorem II.1 using gradient based optimization methods.

### C. Extension to the regression case

1) *Methodology*: As an extension to Section III, we consider now the case of a real valued regression where  $Y_i, f(X_i) \in \mathbb{R}$  for  $i = 1 \dots, n$ . In order to understand the effects of each variable, first we consider the mean criterion

$$M_{i_0, \tau} = \frac{1}{n} \sum_{i=1}^n \lambda_i^{(i_0, \tau)} f(X_i),$$

which will indicate how a change in the variable will modify the output of the learned regression ( $\tau$  is explained in Section III-C and we let  $\lambda_i^{(i_0, \tau)}$ ,  $i = 1, \dots, n$ , be the corresponding weights). Second the variance criterion

$$V_{i_0, \tau} = \frac{1}{n} \sum_{i=1}^n \lambda_i^{(i_0, \tau)} (f(X_i) - M_{i_0, \tau})^2$$

is meant to study the stability of the regression with respect to the perturbation of the variables. Finally the root mean square error (RMSE) criterion

$$\text{RMSE}_{i_0, \tau} = \sqrt{\frac{1}{n} \sum_{i=1}^n \lambda_i^{(i_0, \tau)} (f(X_i) - Y_i)^2}$$

is analogous to the classification error criterion since it enables to detect possibly misleading or confusing variables when learning the regression.

For each  $i_0 \in \{1, \dots, p\}$ , these three criteria can be plotted as a function of  $\tau$  for  $\tau \in [-1, 1]$ .

2) *Application*: We use now our strategy on the Boston Housing dataset<sup>9</sup>. This dataset deals with house prices in Boston. It contains 506 observations with 13 variables that can be used to predict the price of the house to be sold. When considering an optimized Random Forest algorithm, the importance calculated as described in [4] enables to select the 5 most important variables as follows: *lstat* (15227), *rm* (14852), *dis* (2413), *crim* (2144) and *nox* (2042). Remark that the coefficients obtained using a linear model would lead to similar interpretations, with the 5 most important variables as follows: *lstat* (-3.74), *dis* (-3.10), *rm* (2.67), *rad* (2.66), *tax* (-2.07)

As shown in Figure 5, our analysis goes further than these scores. In particular we point out the non linear influence of the variables depending on whether they are high or low. For instance the average number of rooms in a house (variable *rm*) is an important factor that makes the price increase in the case of large houses ( $\tau > 0$ . in Figure 5 (*Average*)). Interestingly, this is far less the case for smaller houses ( $\tau < 0$ . in Figure 5 (*Average*)) since there are other arguments than the number of rooms to keep a high price in this case.

Note that when the number of variables is large, the presence of too many curves may make the graph difficult to understand. In this case, scores that represent average individual evolutions on given ranges of  $\tau$  values for each variables can be computed. Then the highest and lowest scores correspond to the most influential variables on the predictions. For instance, we represent in Table I the evolution of the Average curves in Figure 5 between  $\tau = -0.5$  and  $\tau = 0$ , as well as between  $\tau = 0$  and  $\tau = 0.5$ ,

<sup>9</sup><https://www.kaggle.com/c/boston-housing>

$Mean_0 - Mean_{-0.5}$	$Mean_{0.5} - Mean_0$
black (4.1)	rm (6.80)
rm (3.0)	zn (4.60)
dis (1.7)	chas (2.74)
zn (0.85)	dis (1.64)
...	...
age (-2.78)	rad (-2.99)
indus (-3.2)	indus (-3.05)
ptratio (-3.8)	tax (-3.18)
lstat (-5.1)	lstat (-5.26)

TABLE I

MOST RESPONSIVE VARIABLES TO A POSITIVE OR NEGATIVE STRESS  $\tau$  WHEN ESTIMATING HOUSE PRICES. SCORES ARE SHOWN BETWEEN BRACKETS AND COMPUTED AS THE DIFFERENCE OF THE *Mean* CURVES OF FIGURE 5 FOR (LEFT)  $\tau = -0.5$  AND  $\tau = 0$ , AND (RIGHT)  $\tau = 0$  AND  $\tau = 0.5$ .

which makes clearly understandable which the most influential variables are. It is important to remark that our methodology still allows that the learned decision rules won't be mainly influenced by the same variables depending on whether they increase ( $\tau > 0$ ) or decrease ( $\tau < 0$ ). In Table I, the more influential variables are indeed *rm*, *lstat* and *zn* in the positive direction, while in the negative direction, the variables are *lstat*, *black* and *pratio*. Note that such variables are also cited in studies that relies on LIME [31] or SHAP [27] packages, but the curves we present are more informative and relies on the same distributional input.

#### D. Additional results in the Classification case

1) *Evaluation of the computational burden:* We explained in Section V that our strategy only optimizes, for each of the  $p$  variables, a function which evaluation cost is  $\mathcal{O}(n)$  with no additional outputs predictions out of the *black box* machine learning algorithm. To quantify this, we show in Table II the computational times dedicated to the analysis of synthetic datasets having a different amount of variables  $p$  and observations  $n$ . The variables interpretation was made using 21 values of  $\tau$ , leading to curves as *e.g.* in Figure 1. Computations were run with Python on a standard Intel Core i7 laptop with 24GB memory and no parallelization. It appears that our strategy indeed has a  $\mathcal{O}(np)$  cost, so we then believe it may have a high impact to study the rules learned by black-box machine learning algorithms on large real-life datasets. Remark that when interpreting the influence of the pixel intensities on image test sets, as in Figure 3, only 3 values of  $\tau$  are used. The computations are therefore about 7 times faster. This is coherent with the 10 seconds required on 10000 MNIST images of  $28 \times 28$  pixels in Section IV-B. Note finally that a preliminary implementation of our method in R has lead to very similar results.

$p$	$n$	time (sec)
10	10000	0.76
100	10000	7.79
1000	10000	82.5
10	100000	7.93
10	1000000	86.0

TABLE II

COMPUTATIONAL TIMES REQUIRED ON SYNTHETIC DATASETS, WHERE 21 LEVELS OF STRESS ( $\tau$ ) WERE COMPUTED ON EACH OF THE  $p$  VARIABLES.

2) *Results on simulated data:* In order to further show that our procedure is able to properly recover the characteristics of machine learning algorithms, we again tested it on synthetic data. We have run an experiment with  $p = 5$  variables and  $n = 10^6$  observations, where synthetic data are generated using a logistic regression model, with independent regressors and coefficient vector equal to  $(-4, 2, 0, 2, 4)$ . Figure 6 clearly shows that our method enables to recover the signs and the hierarchy of the coefficients.

#### E. Results on the Iris dataset

As an additional assesment of the method on very well known and simple data, we now consider the *Iris* dataset<sup>10</sup>. This dataset is composed of 150 observations with 4 variables used to predict a label into three categories: *setosa*, *versicolor*, *virginica*. To predict the labels, we used an *Extreme Gradient Boosting* model and a *Random Forest* classifier. Results are show in Figure 7. We first present for both models the Classification error. Then the two other subfigures show the effects of increasing or decreasing the 4 parameters, *i.e.* the width or the length of the sepal or petal. As expected, we recover the well known result that the width of the sepal is the main parameter which enables to differentiate the class *Setosa* while the differentiation between the two other remaining classes is less obvious.

#### F. Other indices: ROC Curves

In the case of two class classification on the Adult Income dataset (Section IV-A), we have shown the evolution of the classification error when the stress parameter  $\tau$  increases. Such results can straightforwardly be extended to True and False

<sup>10</sup><https://archive.ics.uci.edu/ml/datasets/iris>



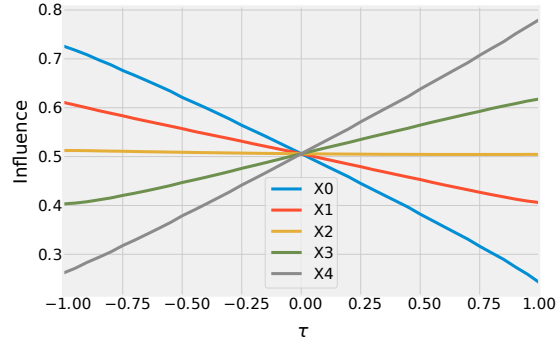


Fig. 6. Proportion of ones found on synthetic data generated using a logistic regression model

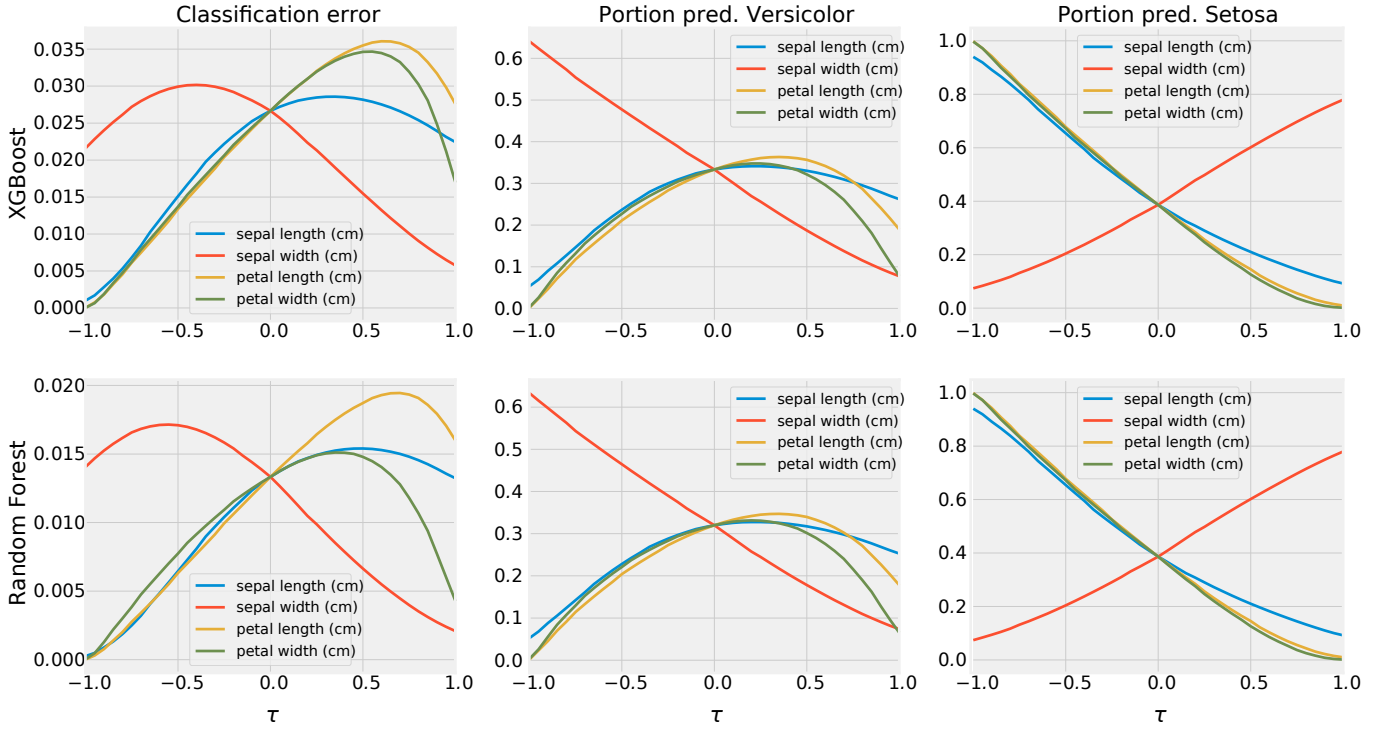


Fig. 7. Evaluation of the classification error and the prediction with respect to the explanatory variable perturbation  $\tau$ , on the *Iris* dataset. The quantity  $\tau$  and the lines have the same signification as in Figure 1. **(Top)** XGBoost Model. The sepal width enables to differentiate the *Setosa* class. **(Bottom)** Random Forest Model. The sepal width again enables to differentiate the *Setosa* class.

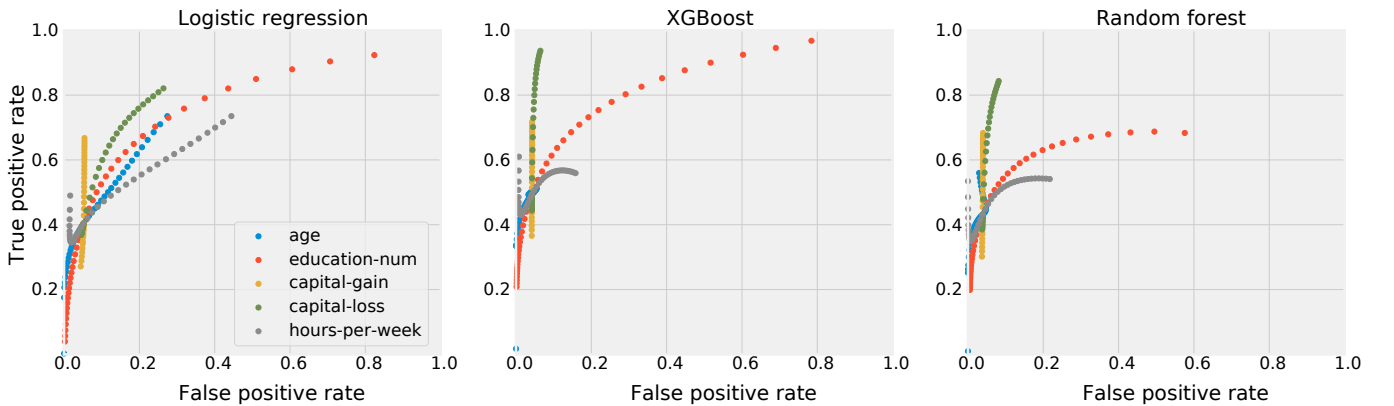


Fig. 8. Evolution of ROC curves in the *Adult income* dataset (Section IV-A). As for the classification errors, we observe that large values of the variable *hoursWeek* make the classification difficult.

Positive Rates, which are commonly represented in ROC curves, that we display in Figure 8. Each point of these curves corresponds to the False Positive Rate and the True Positive Rate, for a sample corresponding to each  $\tau$  and each variable. All curves cross at the same point which corresponds to  $\tau = 0$ . It therefore becomes possible to study the evolution of each criterion.