



SINCLAIR



ROBUSTNESS ASSESSMENT OF BLACK-BOX MODELS

QUANTILE-CONSTRAINED WASSERSTEIN PROJECTIONS AND ISOTONIC POLYNOMIAL APPROXIMATIONS

¹EDF Lab Chatou - Département PRISME

²Institut de Mathématiques de Toulouse

³SINCLAIR AI Lab

SINCLAIR : Workshop de l'été

EDF Saclay

Lundi, 20 Juin 2022

Marouane IL IDRISI¹²³, Nicolas BOUSQUET¹³, Fabrice GAMBOA², Bertrand LOOSS¹²³, Jean-Michel LOUBES².

Main question: How does a model's output react to a perturbation of its input ?

Local robustness:

- Does **small perturbations** entail **big prediction changes**?
- For which **subdomain** of the inputs can one be **confident in the predictions** ?

Global robustness:

- How is a **QoI impacted** by an input perturbation ?
- How does the **relative importance of inputs** change w.r.t. a perturbation ?

Long-term goal: Being able to certify the **stability** of a **black-box model's** predictions w.r.t. to a **distributional shift** of its inputs.

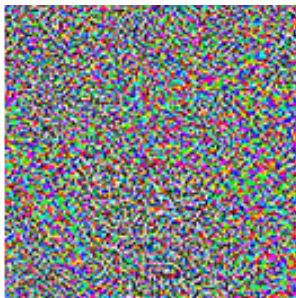
ML setting: Adversarial attacks



"panda"

57.7% confidence

+ ϵ



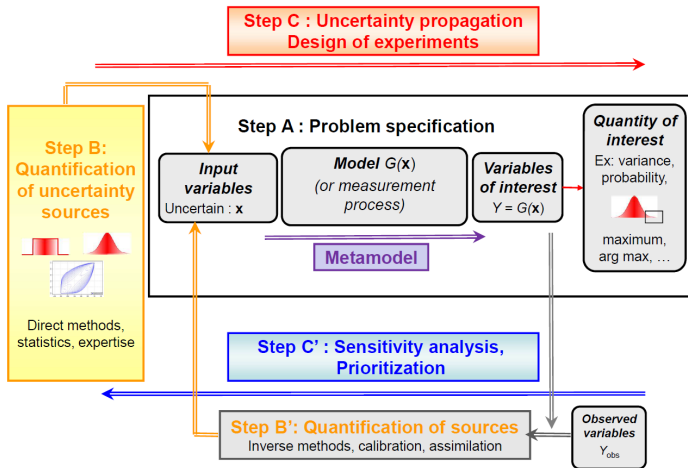
=



"gibbon"

99.3% confidence

SA setting: Input distribution uncertainty



Robustness sensitivity analysis: Analyze variation in the QoI with respect to uncertainty in inputs' distributions (Da Veiga et al. 2021).

This problem can be broken down in **two step**:

1. **Interpretable** and **generic** input perturbation scheme.
2. Global and local robustness to perturbation **diagnostics**.

Main challenge:

- Unify **ML interpretability** and **sensitivity analysis**.

1. Marginal input perturbation methodology
2. Solving the perturbation problem
3. Robustness diagnostics: Acoustic fire extinguisher dataset

1. Marginal input perturbation methodology
2. Solving the perturbation problem
3. Robustness diagnostics: Acoustic fire extinguisher dataset

Marginal input perturbation problem

One seeks to perturb an input such that:

- Some of its interpretable **key statistics are shifted**.
- The perturbation is **minimal**.
- The **dependence structure** between the inputs is **left untouched**.

Marginal input perturbation problem

One seeks to perturb an input such that:

- Some of its interpretable **key statistics are shifted**.
- The perturbation is **minimal**.
- The **dependence structure** between the inputs is **left untouched**.

Formally, let $P \in \mathcal{P}(\mathbb{R})$ be an **initial** probability measure. We seek the solution of the projection problem

$$\begin{aligned} Q = \operatorname{argmin}_{G \in \mathcal{P}(\mathbb{R})} \quad & \mathcal{D}(P, G) \\ \text{s.t.} \quad & G \in \mathcal{C} \end{aligned}$$

where $\mathcal{C} \subseteq \mathcal{P}(\mathbb{R})$ is a **perturbation class**, and \mathcal{D} a discrepancy between probability measures.

P can be modelled as:

- **ML setting:** An empirical measure of an input $P = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$.
- **SA setting:** A marginal probability density function.

Kullback-Leibler divergence and moment perturbation

Several authors proposed perturbations based on the **Kullback-Leibler** (KL) divergence and **generalized moments** perturbations (Csiszar [1975](#)).

- **SA**: Perturbed-law indices (Lemaître et al. [2015](#)).
- **XAI**: Ethik AI (Bachoc et al. [2020](#)).

Kullback-Leibler divergence and moment perturbation

Several authors proposed perturbations based on the **Kullback-Leibler** (KL) divergence and **generalized moments** perturbations (Csiszar 1975).

- **SA**: Perturbed-law indices (Lemaître et al. 2015).
- **XAI**: Ethik AI (Bachoc et al. 2020).

Drawbacks:

- Generalized moments **may not exist**.
- KL divergence is **restrictive** and the result is hard to **interpret**.

Kullback-Leibler divergence and moment perturbation

Several authors proposed perturbations based on the **Kullback-Leibler** (KL) divergence and **generalized moments** perturbations (Csiszar 1975).

- **SA**: Perturbed-law indices (Lemaître et al. 2015).
- **XAI**: Ethik AI (Bachoc et al. 2020).

Drawbacks:

- Generalized moments **may not exist**.
- KL divergence is **restrictive** and the result is hard to **interpret**.

Which key statistics and discrepancy to choose to ensure **genericity** and **interpretability** ?

Kullback-Leibler divergence and moment perturbation

Several authors proposed perturbations based on the **Kullback-Leibler** (KL) divergence and **generalized moments** perturbations (Csiszar 1975).

- **SA**: Perturbed-law indices (Lemaître et al. 2015).
- **XAI**: Ethik AI (Bachoc et al. 2020).

Drawbacks:

- Generalized moments **may not exist**.
- KL divergence is **restrictive** and the result is hard to **interpret**.

Which key statistics and discrepancy to choose to ensure **genericity** and **interpretability** ?

Idea: Quantile-based perturbations and the **Wasserstein distance**.

Why quantiles ?

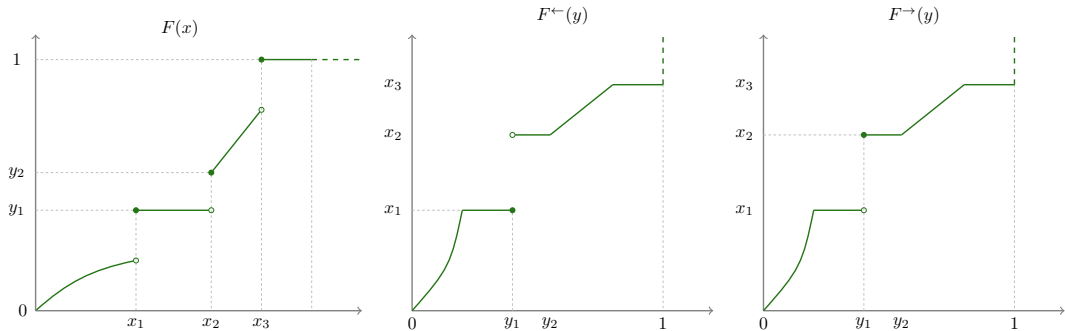
Generalized quantile functions are the generalized inverses (de la Fortelle 2015) of the cdf of random variables.

$$\begin{aligned} F_P^{\leftarrow}(a) &= \sup \{t \in \mathbb{R} \mid F_P(t) < a\} \\ &= \inf \{t \in \mathbb{R} \mid F_P(t) \geq a\}. \end{aligned}$$

$$\begin{aligned} F_P^{\rightarrow}(a) &= \sup \{t \in \mathbb{R} \mid F_P(t) \leq a\} \\ &= \inf \{t \in \mathbb{R} \mid F_P(t) > a\}, \end{aligned}$$

- They **characterize** probability measures (Dufour 1995).
- They exist for any measure in $\mathcal{P}(\mathbb{R})$.
- Quantiles are widely used in many applied fields, and **easy to interpret**.

Generalized quantile functions



The **quantile perturbation class** $\mathcal{Q}_{\mathcal{V}}$ is defined using constraints of the form

$$F_Q^{\leftarrow}(\alpha) \geq b \geq F_Q^{\rightarrow}(\alpha).$$

with $b \in \mathbb{R}$, and leading to the set

$$\mathcal{Q}_{\mathcal{V}} = \{Q \in \mathcal{P}(\mathbb{R}) \mid F_Q^{\leftarrow} \in \mathcal{V}, \quad F_Q^{\leftarrow}(\alpha_i) \geq b_i \geq F_Q^{\rightarrow}(\alpha_i), i = 1, \dots, K\}.$$

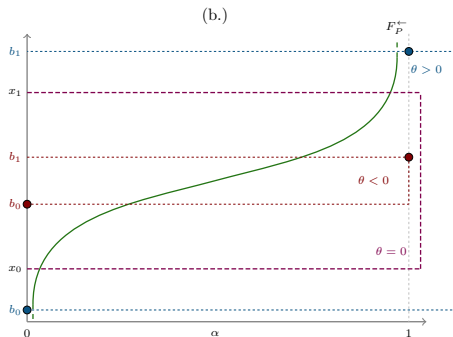
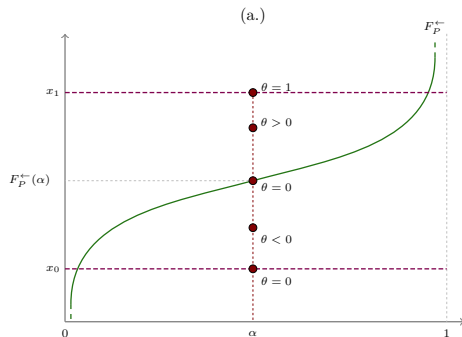
included in $\mathcal{P}(\mathbb{R})$, and where $\mathcal{V} \subseteq \mathcal{F}^{\leftarrow}$ is a **smoothing restriction on the quantile function** characterizing the solution.

Quantile constraints

Perturbations can be driven by an **intensity parameter** $\theta \in [-1, 1]$

- **Quantile shift:** shifting the α -quantile of P between two values.
- **Operating domain dilatation:** widening or narrowing the bounds of the support of P .

Additional **modelling constraints** can also be added (e.g., preservation of empirical quantiles, expert knowledge).



The Wasserstein distance

The p -Wasserstein distance, which naturally comes in the field of **optimal transportation**, can be seen as a **discrepancy between probability measures**.

For two probability measures P and Q in $\mathcal{P}_p(\mathbb{R})$, it simplifies to (Santambrogio 2015)

$$W_p(P, Q) = \left(\int_0^1 |F_P^{-1}(x) - F_Q^{-1}(x)|^p dx \right)^{1/p}$$

Moreover, the 2-Wasserstein distance **metricizes weak convergence** on the set of probability measure with finite 2nd order moments $\mathcal{P}_2(\mathbb{R})$ (Villani 2003).

The Wasserstein distance

The p -Wasserstein distance, which naturally comes in the field of **optimal transportation**, can be seen as a **discrepancy between probability measures**.

For two probability measures P and Q in $\mathcal{P}_p(\mathbb{R})$, it simplifies to (Santambrogio 2015)

$$W_p(P, Q) = \left(\int_0^1 |F_P^{-1}(x) - F_Q^{-1}(x)|^p dx \right)^{1/p}$$

Moreover, the 2-Wasserstein distance **metricizes weak convergence** on the set of probability measure with finite 2nd order moments $\mathcal{P}_2(\mathbb{R})$ (Villani 2003).

It allows to compare the closeness of two probability measures as long as they have a variance.

1. Marginal input perturbation methodology
2. Solving the perturbation problem
3. Robustness diagnostics: Acoustic fire extinguisher dataset

Wasserstein and L^2 projections

The perturbation problem becomes

$$\begin{aligned} Q = \operatorname{argmin}_{G \in \mathcal{P}(\mathbb{R})} \quad & W_2(P, G) \\ \text{s.t.} \quad & G \in \mathcal{Q}_{\mathcal{V}} \end{aligned} \tag{1}$$

Proposition

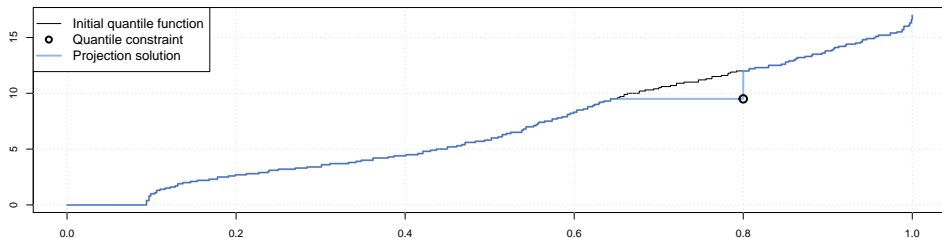
The solution Q of the problem in Eq. (1) is uniquely characterized by its quantile function being the solution

$$\begin{aligned} F_Q^{\leftarrow} = \operatorname{argmin}_{L \in L^2([0,1])} \quad & \int_0^1 (L(x) - F_P^{\rightarrow}(x))^2 \\ \text{s.t.} \quad & L(\alpha_i) \leq b_i \leq L(\alpha_i^+), \quad i = 1, \dots, K, \\ & L \in \mathcal{V} \end{aligned}$$

Solving the perturbation problem

If $\mathcal{V} = \mathcal{F}^{\leftarrow}$, there exists a **unique analytical solution** Q to the problem:

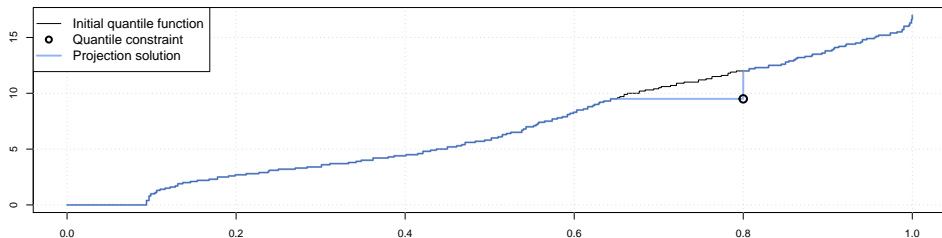
Q is the same as P , except on the intervals between $F_P^{\leftarrow}(\alpha_i)$ and b_i which have no mass, and an atom is added at b_i , taking the initial mass of the interval.



Solving the perturbation problem

If $\mathcal{V} = \mathcal{F}^{\leftarrow}$, there exists a **unique analytical solution** Q to the problem:

Q is the same as P , except on the intervals between $F_P^{\leftarrow}(\alpha_i)$ and b_i which have no mass, and an atom is added at b_i , taking the initial mass of the interval.



How to explicitly add smoothness to the resulting perturbed quantile function ?

Isotonic interpolating piece-wise continuous polynomials

Idea: Using piece-wise continuous polynomials of degree p to ensure continuity.

Partition $[0, 1]$ according into interval $[t_i, t_{i+1}]$, $i = 0, \dots, K$ with $t_0 = 0$, $t_{K+1} = 1$, and $t_i = \alpha_i$ (ordered increasingly), and solve for

$$\begin{aligned} S = \operatorname{argmin}_{G \in \mathbb{R}[x]_{\leq p}} \quad & \int_{t_i}^{t_{i+1}} (F_p^{\rightarrow}(x) - G(x))^2 dx \\ \text{s.t.} \quad & G(t_i) = b_i, G(t_{i+1}) = b_{i+1} \\ & G'(x) \geq 0, \quad \forall x \in [t_0, t_1] \end{aligned} \tag{2}$$

Proposition

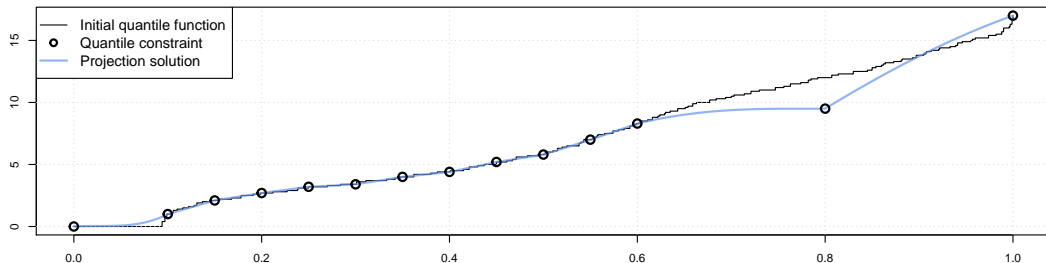
The polynomial solution of Eq. (2) admits as coefficients

$$\begin{aligned} s^* = \operatorname{argmin}_{s \in \mathbb{R}^{p+1}} \quad & s^{\top} M s - 2s^{\top} r \\ \text{s.t.} \quad & s \in \mathcal{K} \end{aligned}$$

where M is the moment matrix of the Lebesgue measure on $[t_i, t_{i+1}]$, r is the moment vector of F_p^{\rightarrow} , and \mathcal{K} is a closed convex subset of \mathbb{R}^{p+1} .

Isotonic interpolation piece-wise continuous polynomials

It is a **Convex Constrained Quadratic Problem** which can be solved numerically using the CVXR solver (Fu, Narasimhan, and Boyd 2020).



The **isotonic interpolating piece-wise continuous polynomial approximation** method can be used for other purposes.

Dependence preservation

ML setting:

Each marginal input $X_i \sim P_i$ can be perturbed using the **monotone perturbation map**

$$T_i = (F_{Q_i}^{\leftarrow} \circ F_{P_i})$$

where the perturbed input is

$$\tilde{X}_i = T_i(X_i) \sim Q_i$$

The empirical copula is preserved, since the ranks of the observation are preserved.

SA setting:

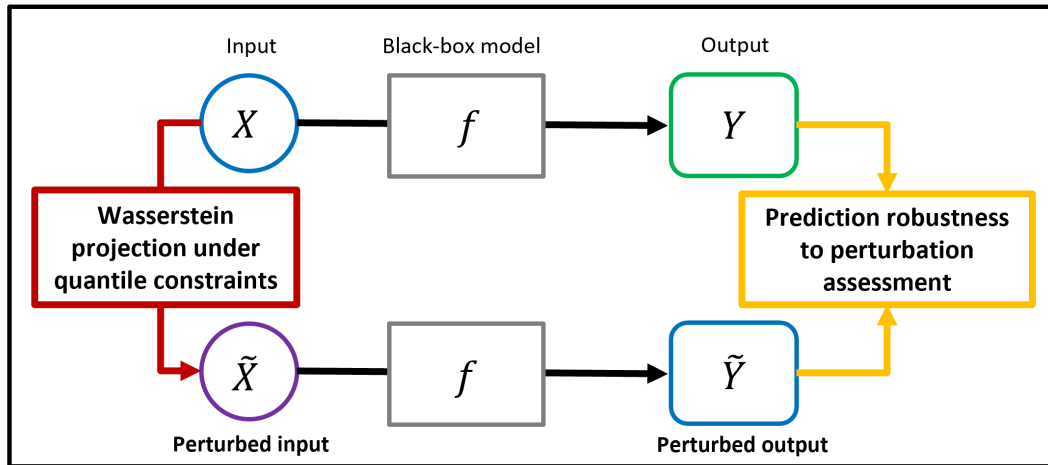
Draw dependent uniform random variables w.r.t. the modeled copula $U = (U_1, \dots, U_d) \sim C$ and work with the perturbed random vector

$$\tilde{X} = (F_{Q_1}^{\leftarrow}(U_1), \dots, F_{Q_d}^{\leftarrow}(U_d))$$

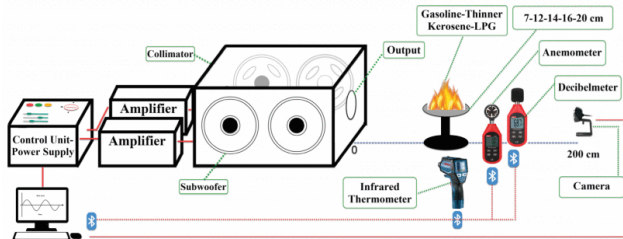
1. Marginal input perturbation methodology
2. Solving the perturbation problem
3. Robustness diagnostics: Acoustic fire extinguisher dataset

Input perturbation and robustness assessment

Methodology



Acoustic Fire Extinguisher



15390 experiments of sound wave fire extinguishing.

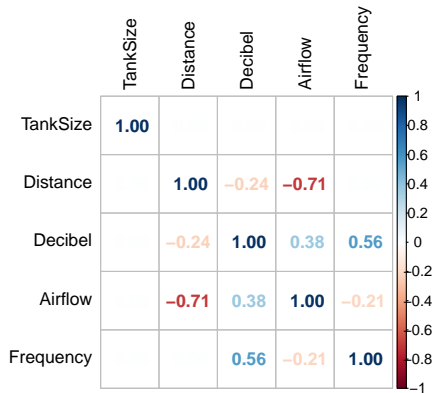
Classification task on 6 variables measured during the experiments.

- Tank Size (L)
- Fuel (Kerosene, Gasoline, Thinner)
- Fire source distance (m)
- Decibel (dB)
- Airflow (m/s)
- Sound frequency (Hz)

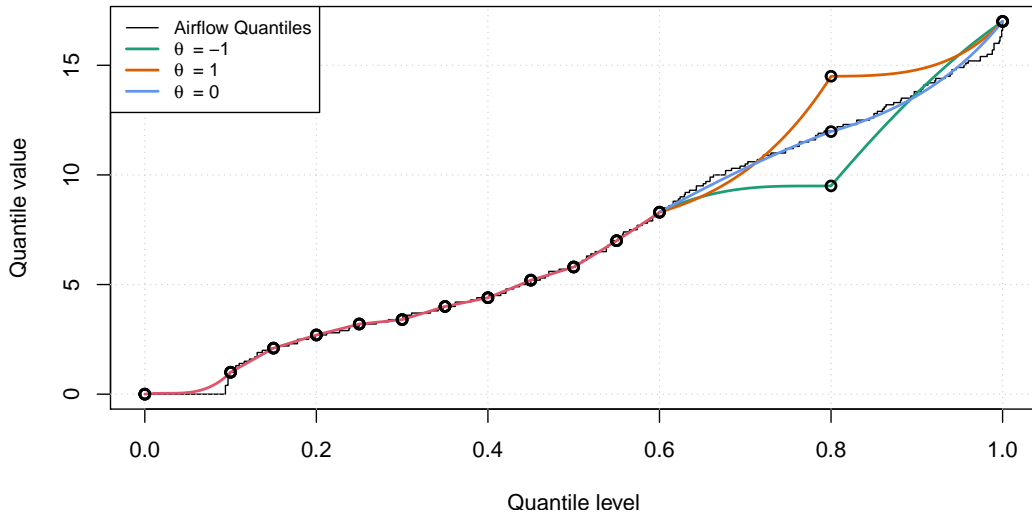
Acoustic Fire Extinguisher

Black-box model: 1-layer neural network (Koklu and Taspinar 2021) trained with an accuracy of 95.15% (validation accuracy of 94.26%).

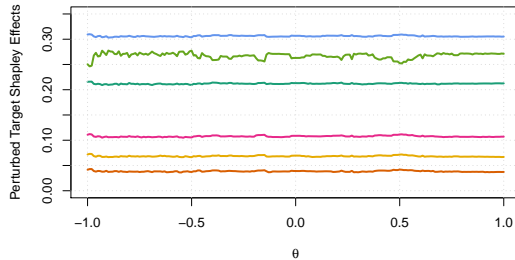
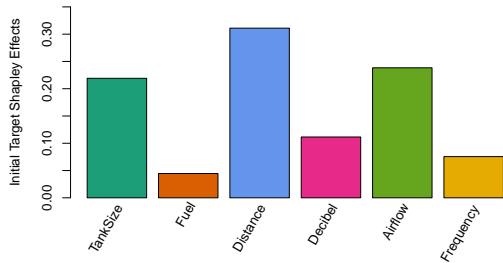
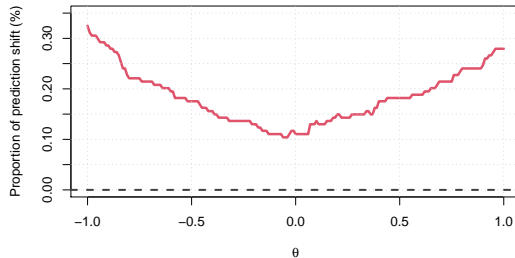
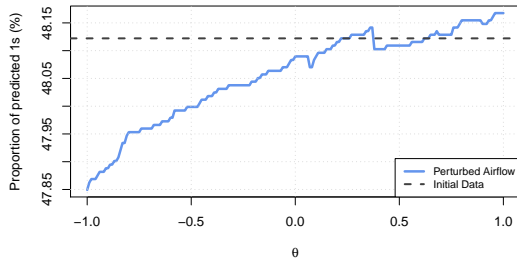
Perturbation scheme: shift of the Airflow 0.8-quantile: initial value at 12, shift between 9.5 ($\theta = -1$) and 14.5 ($\theta = 1$) by polynomial perturbation approximation of degree 9.



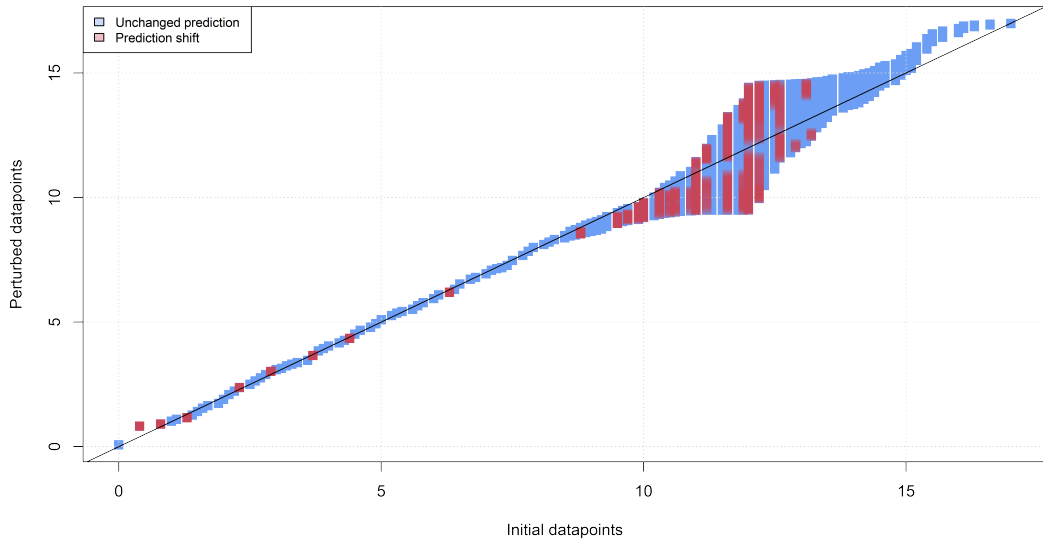
Airflow perturbations



Global robustness



Local robustness



The neural network seem to be **robust globally**:

- Shifting down the 0.8-quantile lead to fewer 1s predicted.
- The more intense the perturbation, the more the predictions change.
- Variable importance order is preserved whatever the perturbation intensity.

However, **locally**, it seems that there are regions where **a small perturbation implies a change in prediction**.

Conclusion & perspectives

Generic and interpretable **marginal perturbation scheme**.

Local and global robustness assessment of black-box numerical (SA) and predictive models (ML).

Perspectives:

- Parallel and efficient computation in \mathbb{R} (soon to be published).
- Optimal degree selection, and derivability of the resulting polynomial.
- Multivariate quantile perturbation, and other discrepancies (Prokhorov).
- More general smoothing spaces.
- Super-quantile constraints.

References i

- Bachoc, F., F. Gamboa, M. Halford, J-M. Loubes, and L. Risser. 2020. "Explaining Machine Learning Models using Entropic Variable Projection" [in en]. ArXiv: 1810.07924, *arXiv:1810.07924 [cs, stat]* (December). <http://arxiv.org/abs/1810.07924>.
- Csiszar, I. 1975. "I-Divergence Geometry of Probability Distributions and Minization problems." *The Annals of Probability* 3 (1): 146–158. <https://doi.org/10.1214/aop/1176996454>. <http://doi.org/10.1214/aop/1176996454>.
- Da Veiga, S., F. Gamboa, B. looss, and C. Prieur. 2021. *Basics and Trends in Sensitivity Analysis. Theory and Practice in R*. SIAM.
- de la Fortelle, A. 2015. "A study on generalized inverses and increasing functions Part I: generalized inverses" [in en], 14. <https://hal-mines-paristech.archives-ouvertes.fr/hal-01255512>.
- Dufour, J-M. 1995. *Distribution and quantile functions* [in en]. https://jeanmariedufour.github.io/ResE/Dufour_1995_C_Distribution_Quantile_W.pdf.
- Fu, A., B. Narasimhan, and S. Boyd. 2020. "CVXR: An R Package for Disciplined Convex Optimization." *Journal of Statistical Software* 94 (14): 1–34. <https://doi.org/10.18637/jss.v094.i14>.
- Koklu, M., and Y. S. Taspinar. 2021. "Determining the Extinguishing Status of Fuel Flames With Sound Wave by Machine Learning Methods." Conference Name: IEEE Access, *IEEE Access* 9:86207–86216. issn: 2169-3536. <https://doi.org/10.1109/ACCESS.2021.3088612>.
- Lemaître, P., E. Sergienko, A. Arnaud, N. Bousquet, F. Gamboa, and B. looss. 2015. "Density modification-based reliability sensitivity analysis." *Journal of Statistical Computation and Simulation* 85 (6): 1200–1223. <https://doi.org/10.1080/00949655.2013.873039>. eprint: <https://doi.org/10.1080/00949655.2013.873039>. <https://doi.org/10.1080/00949655.2013.873039>.

- Santambrogio, F. 2015. *Optimal Transport for Applied Mathematicians*. Vol. 87. Progress in Nonlinear Differential Equations and Their Applications. Cham: Springer International Publishing. ISBN: 978-3-319-20827-5 978-3-319-20828-2.
<https://doi.org/10.1007/978-3-319-20828-2>. <http://link.springer.com/10.1007/978-3-319-20828-2>.
- Villani, C. 2003. *Topics in Optimal Transportation* [in en]. Vol. 58. Graduate Studies in Mathematics. ISSN: 1065-7339. American Mathematical Society, March. ISBN: 978-0-8218-3312-4 978-0-8218-7232-1 978-1-4704-1804-5, accessed June 23, 2021.
<https://doi.org/10.1090/gsm/058>. <http://www.ams.org/gsm/058>.

THANK YOU FOR YOUR ATTENTION!

ANY QUESTIONS?

Projecting without smoothing

Let P be a probability measure in $\mathcal{P}_2(\mathbb{R})$. Let \mathcal{C} be a non-empty perturbation class, defined by a set of quantile constraints \mathcal{Q} . Furthermore, assume, without loss of generality, that, for $i = 1, \dots, K$,

$$\alpha_1 < \dots < \alpha_K, \quad \text{along with,} \quad b_1 < \dots < b_K$$

and let $\beta_i = F_P(b_i)$ for $i = 1, \dots, K$. Denote the following intervals:

$$\begin{aligned} c_1 &= \min(\beta_1, \alpha_1), & c_i &= \min\left[\max(\alpha_{i-1}, \beta_i), \alpha_i\right], i = 2, \dots, K; \\ d_K &= \max(\beta_K, \alpha_K), & d_j &= \max\left[\min(\beta_j, \alpha_{j+1}), \alpha_j\right], j = 1, \dots, K-1. \end{aligned}$$

Furthermore, let $A_i = [c_i, d_i]$ for $i = 1, \dots, K$, $A = \bigcup_{i=1}^K A_i$ and $\bar{A} = [0, 1] \setminus A$.

The solution of the perturbation problem

$$\begin{aligned} Q &= \operatorname{argmin}_{G \in \mathcal{P}_2(\mathbb{R})} W_2(P, G) \\ &\text{s.t. } G \in \mathcal{C} \end{aligned} \tag{3}$$

admits, as a characterizing quantile function :

$$F_Q^{\leftarrow}(y) = \begin{cases} F_P^{\rightarrow}(y) & \text{if } y \in \bar{A} \\ b_i & \text{if } y \in A_i, \quad i = 1, \dots, K \end{cases}$$

Non-negativity of polynomials on closed intervals

Theorem (Non-negativity of polynomials on closed intervals)

Let $t_0, t_1 \in \mathbb{R}$ such that $t_0 < t_1$, and let $p \in \mathbb{N}^*$.

A univariate polynomial S of even degree $d = 2p$ is non-negative on $[t_0, t_1]$ if and only if it can be written as, $\forall x \in [t_0, t_1]$

$$S(x) = Z(x) + (x - t_0)(t_1 - x)W(x)$$

where Z is an SOS polynomial of degree at most equal to d , and W is an SOS polynomial of degree at most equal to $d - 2$.

A univariate polynomial S of odd degree $d = 2p + 1$ is non-negative on $[t_0, t_1]$ if and only if it can be written as, $\forall x \in [t_0, t_1]$

$$S(x) = (x - t_0)Z(x) + (t_1 - x)W(x)$$

where Z, W are SOS polynomials of degree at most equal to d .

SDP representation of SOS polynomials

Let S be an univariate polynomial of even degree $d = 2p$, with coefficients $s = (s_0, \dots, s_d)$, and denote x_p the usual monomial basis of polynomials of degree at most equal to p , i.e., $x_p = (1, x, x^2, \dots, x^{p-1}, x^p)^\top$. S is an SOS polynomial if and only if there exists a $(p \times p)$ symmetric semi definite positive (SDP) matrix

$$\Gamma = [\Gamma_{ij}]_{i,j=1,\dots,p}$$

that satisfies, $\forall x \in \mathbb{R}$,

$$S(x) = x_p^\top \Gamma x_p.$$

Moreover, for $k = 0, \dots, d$, let \mathbb{I}_k^p be the $(p \times p)$ matrix defined by, for $i, j = 1, \dots, p$:

$$[\mathbb{I}_k^p]_{i,j} = \mathbb{1}_{\{i+j=k+2\}}(i,j).$$

If there exists a matrix Γ such that S is SOS, then one has that, for $i = 0, \dots, d$

$$s_i = \langle \mathbb{I}_i^p, \Gamma \rangle_F = \sum_{j+k=i+2} \Gamma_{j,k}$$

where, $\langle \cdot, \cdot \rangle_F$ denotes the Frobenius norm on matrices.

Equivalent optimization formulation

Let $[t_0, t_1] \subset [0, 1]$, and let $s = (s_0, \dots, s_d)^\top \in \mathbb{R}^{d+1}$, M be the symmetric $((d+1) \times (d+1))$ moment matrix of the Lebesgue measure on $[t_0, t_1]$, i.e. for $i, j = 0, \dots, d$,

$$M_{ij} = \int_{t_0}^{t_1} x^{i+j} dx = \frac{(t_1)^{i+j+1} - (t_0)^{i+j+1}}{i+j+1},$$

and denote $r \in \mathbb{R}^{d+1}$ the moment vector of $A(x)$, i.e., for $i = 0, \dots, d$

$$r_i = \int_{t_0}^{t_1} x^i F_P^{\leftarrow}(x) dx$$

Then, the optimization problem can be equivalently solved by finding s as being the solution of the following convex constrained quadratic program,

$$\begin{aligned} s^* &= \operatorname{argmin}_{s \in \mathbb{R}^{d+1}} s^\top M s - 2s^\top r \\ &\text{s.t. } s \in \mathcal{K} \end{aligned}$$

where \mathcal{K} is a closed convex subset of \mathbb{R}^{d+1} .