



**HAL**  
open science

## Identifying potential significant factors impacting zero-inflated proportions data

Melina Ribaud, Edith Gabriel, Joseph Hughes, Samuel Soubeyrand

► **To cite this version:**

Melina Ribaud, Edith Gabriel, Joseph Hughes, Samuel Soubeyrand. Identifying potential significant factors impacting zero-inflated proportions data. 2023. hal-02936779v4

**HAL Id: hal-02936779**

**<https://hal.science/hal-02936779v4>**

Preprint submitted on 7 Jun 2023

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

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

# Identifying potential significant factors impacting zero-inflated proportion data

Mélina Ribaud<sup>1,\*</sup>, Edith Gabriel<sup>1</sup>, Joseph Hughes<sup>2</sup>, and Samuel Soubeyrand<sup>1,\*</sup>

<sup>1</sup>INRAE, BioSP, 84914 Avignon, France

<sup>2</sup>MRC-University of Glasgow, Centre for Virus Research, Glasgow, Scotland, United Kingdom

\*Corresponding authors: [melina.ribaud@gmail.com](mailto:melina.ribaud@gmail.com) and [samuel.soubeyrand@inrae.fr](mailto:samuel.soubeyrand@inrae.fr)

April 3, 2023

## Abstract

Classical supervised methods like linear regression and decision trees are not completely adapted for identifying impacting factors on a response variable corresponding to zero-inflated proportion data (ZIPD) that are dependent, continuous and bounded. In this article we propose a within-block permutation-based methodology to identify factors (discrete or continuous) that are significantly correlated with ZIPD, we propose a performance indicator quantifying the percentage of correlation explained by the subset of significant factors, and we show how to predict the ranks of the response variables conditionally on the observation of these factors. The methodology is illustrated on simulated data and on two real data sets dealing with epidemiology. In the first data set, ZIPD correspond to probabilities of transmission of Influenza between horses. In the second data set, ZIPD correspond to probabilities that geographic entities (e.g., states and countries) have the same COVID-19 mortality dynamics.

**Keywords:** COVID-19; Equine Influenza; Performance Indicator; Permutation Test; Ranking; Spearman's correlation.

## 1 Introduction

Proportion data are encountered in many fields such as biology, epidemiology and marketing. A common objective is the identification of external factors impacting these data. In

27 marketing, a typical study could be the identification of factors impacting the proportions  
28 of products sold to different age groups of customers. In biology, an analogous study could  
29 be the analysis of proportions of cures with certain drugs by age groups. In epidemiology,  
30 one may be interested in identifying the external factors impacting the transmission of a  
31 virus within a host population when the knowledge about the transmissions (i.e., who in-  
32 fected whom) is uncertain (hence, in this case, the probability of transmission from host A  
33 to host B can be viewed as a proportion datum). In these typical examples, the ‘age groups  
34 and products’, the ‘age groups and drugs’ or the ‘hosts’ can be viewed as the nodes of a  
35 network whose edges are weighted by the aforementioned proportions measuring the links  
36 between age groups and products / drugs or the links between hosts. The edges from the  
37 ‘contributing nodes’ (i.e., the source hosts transmitting the virus, the products or the drugs)  
38 toward a specific ‘target node’ (i.e., a recipient host or an age group) correspond to a vector  
39 of proportions whose sum is equal to one (or eventually lower than one if some contributing  
40 nodes are unobserved). Note that in the epidemiological example, recipient hosts can also  
41 be source hosts and vice versa (i.e., a host can be both a target and a contributing node).  
42 The network vision of data is illustrated by two schematic configurations shown in Figure 1,  
43 which includes notations that are introduced in Section 2, and our general objective in this  
44 paper is to unravel which factor(s) characterizing the pairs of nodes *explain* the network  
45 edges.

46 In this article, we are specifically interested in epidemiological applications. As recently  
47 illustrated with the COVID-19 pandemic, grounding strategies for the management of infec-  
48 tious diseases on accurate knowledge about risk factors is paramount for effectively prevent-  
49 ing a health crisis. Indeed, assessing the influence of social, biological and environmental  
50 factors in the spread of epidemics contributes to identifying mechanisms for controlling the  
51 disease dynamics. The spread of epidemics can be understood by quantifying epidemiologi-  
52 cal links between hosts or, more generally, nodes. Typical examples of epidemiological links  
53 that we have in mind are: probabilities of disease transmission between individuals<sup>1</sup>, and  
54 similarity measures of disease dynamics in several geographic entities<sup>2</sup>. Such measures of  
55 epidemiological links (*i*) have an intrinsic correlation structure and (*ii*) are usually estimated  
56 (i.e., uncertain). These features make the investigation of the relationship between epidemi-  
57 ological links and risk factors challenging because dependencies and overdispersion may lead  
58 to biased results if ignored<sup>3,4</sup>. Here, we focus on epidemiological links defined as proportions  
59 and we aim to provide a statistical methodology for explaining these epidemiological links  
60 (hereafter, the response variable) by multiple potentially impacting factors.

61 Many statistical methods can be used to identify the correlation between factors and a  
62 response variable. Parametric prediction models can identify the set of factors impacting  
63 the response through statistical tests. When the response is normally distributed, or when

64 data are transformed to make it fit a Gaussian distribution<sup>5</sup>, the linear regression model<sup>6</sup>  
65 predicts response values and identifies influencing factors. When the response variable follows  
66 other frequently used distributions (e.g., binomial or Poisson), the generalized linear models<sup>7</sup>  
67 (GLM) can be considered. When one prefers to avoid making a distributional assumption for  
68 the response variable, non-parametric predictive models<sup>6</sup> may be a solution. However, non-  
69 parametric models do not standardly provide direct testing procedure to identify impacting  
70 factors.

71 In the exploration of the link between factors and a response variable, we are particularly  
72 interested in zero-inflated response variables (i.e., random variables with a positive mass at  
73 zero, or, in other words, variables presenting an excess of zeros). In classification problems  
74 (i.e., problems with categorical response variables), excess of any class can be handled by  
75 balancing classes using resampling methods<sup>8</sup>. In regression problems (i.e., problems with  
76 quantitative response variables), zero-inflation is typically handled by defining a model as a  
77 mixture of two processes: the first process generating only zeros, the second process being  
78 governed by a usual distribution such as the zero-inflated Poisson, zero-inflated binomial  
79 and zero-inflated beta distributions<sup>9</sup>. For such zero-inflated models, generally assuming  
80 independence between observations, the influencing factors can be identified with statistical  
81 tests<sup>10</sup>.

82 The above-mentioned parametric models are defined for independent and identically dis-  
83 tributed (i.i.d.) realizations. However, proportion data are not independent since they sum  
84 to a fixed value equal to or lower than one. Such data are often referred to as compositional  
85 data<sup>11</sup>, that have been classified with respect to the nature of the response<sup>12</sup> (proportions  
86 arising from counts *versus* from continuous measurements). Regarding the case of a zero-  
87 and/or one-inflated continuous response, the zero- and/or one-inflated beta regression is a  
88 solution when the proportions work in pairs (e.g., the proportions of males and females for a  
89 given species). When the number of observed categories is greater than two, the Dirichlet's  
90 regression can be used. For instance, an adaptation of the zero-inflated Dirichlet regression  
91 (ZIDR) model was proposed for microbiome compositional data<sup>13</sup>.

92 Statistical tests are generally associated with the parametric approaches mentioned above  
93 for quantifying the significance of a factor (the test generally depends on the type of factors:  
94 discrete *versus* continuous). The statistical test accompanying the linear model can treat all  
95 types of factors. ANOVA can handle discrete factors with more than two levels. The GLM  
96 (including zero-inflated data) and the ZIDR can treat continuous factors as well as discrete  
97 factors with only two levels, even if this restriction is minor since a factor with multiple levels  
98 can be treated as several factors with two levels each.

99 Here we investigate the relationship between zero-inflated, non-Gaussian, correlated pro-  
100 portion data and several factors of any type. In the epidemiological contexts that we are

101 interested in, this objective translates into the investigation of the impact of individual,  
102 environmental, economical, climatic... factors on epidemiological links. The structure of  
103 the data and the objectives generate constraints on the statistical approach to be used.  
104 The response takes values between 0 and 1 (inclusive) and is generally zero-inflated (the  
105 zero-inflation makes classical transformations yielding normally-distributed variables inap-  
106 plicable). Moreover, the realizations are mutually dependent due to the constraint over the  
107 sum of probabilities for a given target. Furthermore, factor values for a given target node  
108 not only depend on the characteristics of this node but also on the characteristics of the  
109 contributing nodes and the target-contributor interaction. Common methods do not match  
110 all of these constraints, as illustrated by Table [1](#).

111 Therefore, we propose a model-free (or more precisely a distribution-free) approach based  
112 on permutation tests aiming (*i*) to identify factors (discrete or continuous; characterizing the  
113 target, the contributor or the target-contributor pair) that are significant, (*ii*) to quantify via  
114 a performance indicator the percentage of correlation explained by the subset of significant  
115 factors, and (*iii*) to predict the ranks of proportions from the significant factors. To take  
116 into account the dependence of proportion data linking the contributing nodes to any target  
117 node, we define a test grounded on the principle of within-block permutations [14](#)[15](#)[16](#). In  
118 this test, permutations are constrained by the dependence structure of data by shuffling  
119 proportion data only within each ‘contributors–target’ block (for any target). In addition,  
120 the test deals with both factor types, continuous and discrete, by adapting the test statistic,  
121 which is based on Spearman’s correlation if the factor is continuous, and on a difference in  
122 mean ranks if it is discrete. The test is applied for each factor separately, but significant  
123 factors are then jointly used to compute a performance indicator quantifying the percentage  
124 of correlation explained by the selected set of factors. The zero-inflation is not explicitly  
125 handled in the test itself (permutation procedures are nevertheless considered as versatile  
126 for variables clumping at zero because they do not require distributional assumptions [17](#)),  
127 but it is accounted for in the performance indicator whose range is robust to the proportion  
128 of zeros (and more generally to the proportion of ties). Supporting Information, Figure S1,  
129 gives a general picture about the workflow of the methodology briefly described above.

130 In what follows, Section [2](#) sets the framework and notations and presents the procedure  
131 based on permutation tests to identify the factors correlated to the response as well as the  
132 performance indicator. The method is then applied to simulations (Section [3](#)) and to real  
133 data dealing with Equine Influenza and COVID-19 epidemics (Sections [4](#) and [5](#)).

Table 1: Comparison of models in their ability to match the constraints considered in this article.

Methods	Response			Factor		Dependency
	Distribution	[0, 1]	Zero	Tests		
	free		inflated	Discrete	Continuous	
Linear regression <sup>6</sup>				✓ <sup>a</sup>	✓	
Beta regression <sup>9</sup>		✓	✓	✓	✓	
Dirichlet regression <sup>18</sup>		✓	✓	✓	✓	✓
Decision tree <sup>19</sup>	✓	✓				✓

<sup>a</sup> ANOVA and ANCOVA

## 2 Identification and quantification of impacting factors

### 2.1 Framework and notations

Hereafter, let  $n_t$  be the number of target nodes,  $n_c$  the number of contributing nodes and  $d$  the number of factors. The response variable  $Z_j^i$  is a random variable measuring the (directed) epidemiological link between target  $i \in \{1, \dots, n_t\}$  and contributor  $j \in \{1, \dots, n_c\}$ , as illustrated in Figure 1. The higher the value, the stronger the link and the weaker the other links. We assume that  $Z_j^i$  is continuous,  $Z_j^i \in [0, 1]$  and, for any target node, the sum over all contributors cannot exceed 1, i.e.:

$$\sum_{j=1}^{n_c} Z_j^i \leq 1, \quad \forall i \in \{1, \dots, n_t\}. \quad (1)$$

In addition, in typical applications we will consider that the distribution of  $Z_j^i$  is zero-inflated.

The factors characterize any target-contributor pair (i.e., the two nodes and their interaction), as illustrated in Figure 1. We denote by  $\mathbb{X} = (\mathbf{x}_1, \dots, \mathbf{x}_d) \in \mathbb{R}^{n_t n_c \times d}$  the set of  $d$  factors ( $d \in \mathbb{N}^*$ ). Thus, any pair  $(i, j)$  is described by  $(x_1^{(i,j)}, \dots, x_d^{(i,j)})$ .

In practice, factors are often grounded on information separately collected from the target node  $i$  and the contributing node  $j$ . In this case,  $x_k^{(i,j)}$  is defined from any application taking as arguments the values of a variable observed from the contributor and target nodes. Intuitive examples include the geographical distance measured from the spatial coordinates of the contributor and the target, the absolute difference between the ages of the contributor and the target, etc. Further examples provided in the application sections lead to discrete, continuous and even categorial factors.

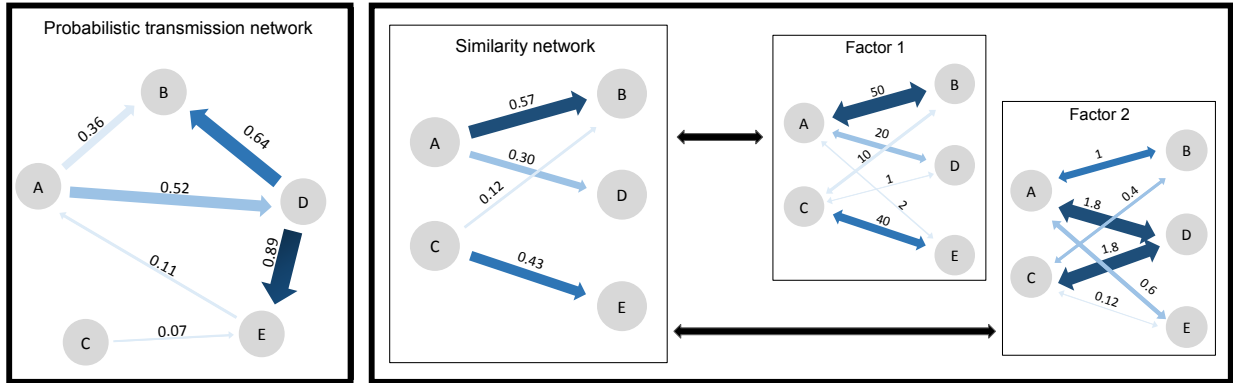


Figure 1: Typical configurations of epidemiological networks of interest and links to factors. Left: Probabilistic transmission network; Arrows indicate probable transmissions (with associated probabilities given by numbers above the arrows); Host C is a contributor (i.e., probable source of infection), host B is a target (i.e., an infection receptor) and hosts A, D and E are both contributors and targets (thus,  $n_t = n_c = 4$ ); E.g., A infected B with probability  $Z_A^B = 0.36$ . Right: Similarity network; Arrows indicate probabilities that targets (namely host units B, D and E;  $n_t = 3$ ) have the same features as (e.g., ‘follow the mortality dynamics of’) contributors (namely hosts units A and C;  $n_c = 2$ ); E.g., B follows the mortality dynamics of A with probability  $Z_A^B = 0.57$ . Notes: The objective of our approach is to assess the link between an epidemiological network and multiple factors that characterize pairs of host units, as illustrated on the right-hand side where Factor 1 seems to better reflect the similarity network than Factor 2 does (the numbers above the arrows for the panels representing Factors 1 and 2 are the values taken by these undirected factors); In the main text, probabilities  $Z_j^i$  are indexed by host indices instead of host names; The sum of inward probabilities for a given target is lower than 1 if some contributors are unobserved.

## 2.2 A within-block permutation-based approach to identify influencing factors

The specific characteristics of our response variable make impossible the use of classical correlation tests such as Spearman’s test<sup>[20]</sup>, because the response has numerous ties (zeros). Solutions to treat ties in ranking problems were proposed<sup>[21]</sup>, but for cases with numerous ties, some hypotheses on the moments have to be satisfied and checking them may be laborious. Furthermore, the response is dependent within each target–contributors block (see Equation (1)) and classical correlation tests do not take into account such a dependence structure. Within-block permutation tests, grounded on the assumption of exchangeability of data within blocks<sup>[14][15][16]</sup>, appear to be a possible alternative to take into account these constraints.

Let  $\mathbf{x}_k \in \mathbb{R}^{n_t n_c}$ ,  $k = 1, \dots, d$ , be the observations of the factor to be tested and let

158  $\mathbf{z} \in \mathbb{R}^{n_t n_c}$  be the observations of the response, whose element  $(i, j)$  denoted by  $z_j^i$  is the  
 159 observed value of  $Z_j^i$ . We adapt the Conditional Monte Carlo (CMC) algorithm<sup>[22]</sup> for block-  
 160 permutation to test the correlation between the response and the factor. We denote by  $T_k$   
 161 the test statistic, which is dependent on the type factor, and by  $\lambda_k(\mathbf{z})$  the p-value. The  
 162 CMC algorithm for block-permutation test adapted to data that we consider in this article  
 163 consists of the following steps:

- 164 1. Compute the statistic  $T_k$  on the original data set  $(\mathbf{x}_k, \mathbf{z})$ ;
- 165 2. Do  $B$  independent repetitions of what follows: randomly permute the response by block  
 166 of type target–contributors, set the new response vector denoted  $\mathbf{z}^b$ ,  $b = 1, \dots, B$ , and  
 167 compute the statistic  $T_k^b$  on the permuted data set  $(\mathbf{x}_k, \mathbf{z}^b)$ ;
- 168 3. Estimate the p-value by  $\hat{\lambda}_k(\mathbf{z}) = \frac{1}{B} \sum_{b=1}^B \mathbf{1}_{\{T_k^b \geq T_k\}}$ .

169 Remarks: The  $T_k$  statistic must be positive to calculate the p-value using the CMC algo-  
 170 rithm. The permutation by block of type target–contributors is carried out by permuting  
 171 the components of the vector  $(z_1^i, \dots, z_{n_c}^i)$  for each target node  $i$  (the permuted response  $\mathbf{z}^b$   
 172 hence satisfies mutual dependencies between its components summarized in the system of  $n_t$   
 173 constraints given by Equation (1), and keeps eventual heterogeneity in the distribution char-  
 174 acteristics of  $(z_1^i, \dots, z_{n_c}^i)$  between targets). Block permutations are required to minimize  
 175 the second species risk of the test (see Supporting Information, Appendix A).

176 For a continuous factor  $\mathbf{x}$ , by omitting the subscript  $k$ , the statistic  $T$  is the non-  
 177 parametric Spearman’s correlation<sup>[23]</sup> between  $\mathbf{x}$  and  $\mathbf{z}$ , say  $r_s(\mathbf{x}, \mathbf{z})$ , i.e. it is defined as  
 178 the Pearson correlation between the rank variables  $r_s(\mathbf{x}, \mathbf{z}) = \rho(R_{\mathbf{x}}, R_{\mathbf{z}})$ , where  $\rho$  is the Pear-  
 179 son correlation,  $R_{\mathbf{x}}$  (resp.  $R_{\mathbf{z}}$ ) is the random vector that gives the ranks of the elements of  
 180  $\mathbf{x}$  (resp.  $\mathbf{z}$ ). Hence, we define the following tests:

181  $H_0$ : “the response and factor ranks are not correlated” versus

- 182 (i)  $H_1$ : “the response and factor ranks are correlated” and the test statistic is  $T = r_s^2(\mathbf{x}, \mathbf{z})$ ;
- 183 (ii)  $H_1$ : “the response and factor ranks are positively correlated” and  $T = r_s(\mathbf{x}, \mathbf{z})$ ;
- 184 (iii)  $H_1$ : “the response and factor ranks are negatively correlated” and  $T = -r_s(\mathbf{x}, \mathbf{z})$ .

For a discrete factor  $\mathbf{x}$  (still omitting the subscript  $k$ ) with  $Q$  levels, the test hypotheses  
 are  $H_0$ : “level-by-level mean ranks are equal” versus  $H_1$ : “mean ranks are different for at  
 least two levels” and the statistic corresponds to the one defined in the H-test<sup>[24]</sup>.

$$T = (n_t n_c - 1) \frac{\sum_{q=1}^Q n_q (\bar{R}_{\mathbf{z}, q} - \bar{R}_{\mathbf{z}})^2}{\sum_{i=1}^{n_t} \sum_{j=1}^{n_c} (R_{z_j^i} - \bar{R}_{\mathbf{z}})^2}, \quad (2)$$



185 where  $R_{z_j^i}$  denote the rank of the element  $(i, j)$  of  $\mathbf{z}$ ,  $\bar{R}_{\mathbf{z}} = \frac{1}{n_t n_c} \sum_{i=1}^{n_t} \sum_{j=1}^{n_c} R_{z_j^i}$ ,  $n_q =$   
186  $\sum_{i=1}^{n_t} \sum_{j=1}^{n_c} \mathbf{1}_q(x^{(i,j)})$ ,  $\bar{R}_{\mathbf{z},q} = \frac{1}{n_q} \sum_{i=1}^{n_t} \sum_{j=1}^{n_c} R_{z_j^i} \mathbf{1}_q(x^{(i,j)})$  and  $\mathbf{1}_q(x^{(i,j)}) = 1$  if  $x^{(i,j)} = q$ ,  
187  $\mathbf{1}_q(x^{(i,j)}) = 0$  otherwise. If the p-value of the factor being considered is less than the signifi-  
188 cance level, *post-hoc* tests can be constructed to test the impact of factor levels. Let  $q$  and  $\tilde{q}$   
189 be two levels, we can then make the following tests  $H_0$ : “there is no difference between the  
190 two mean ranks” versus

191 (i)  $H_1$ : “there is a difference between the two mean ranks” and the statistic is  $T =$   
192  $(\bar{R}_{z,q} - \bar{R}_{z,\tilde{q}})^2$ ;

193 (ii)  $H_1$ : “the mean ranks of level  $q$  is lower than the mean ranks of  $\tilde{q}$ ” and  $T = \bar{R}_{z,\tilde{q}} - \bar{R}_{z,q}$ ;

194 (iii)  $H_1$ : “the mean ranks of level  $q$  is greater than the mean ranks of  $\tilde{q}$ ” and  $T = \bar{R}_{z,q} - \bar{R}_{z,\tilde{q}}$ .

195 Notes: Here, the statistics are the differences in mean ranks<sup>25</sup>. In addition, if the dis-  
196 crete factor has more than two levels, the problem becomes a multiple comparison prob-  
197 lem. A correction can be applied accordingly to control the occurrence of false positives,  
198 e.g., the Bonferroni correction which consists in multiplying the p-values by the number  
199 of comparisons, or the less conservative and sharper improved Bonferroni correction called  
200 Benjamini-Hochberg correction<sup>26,27</sup>. As an illustration, we provide both the Bonferroni and  
201 the Benjamini-Hochberg corrected p-values for the post-hoc tests performed in the applica-  
202 tion dealing with Equine Influenza.

## 203 **2.3 A performance indicator to quantify the monotonous depen-** 204 **dency**

205 To take into account the multivariate aspect of the correlation, we develop a performance  
206 indicator that simultaneously accounts for all discrete and continuous factors previously  
207 identified. The indicator can be viewed as a surrogate for the coefficient of determination  
208 used in linear regression, representing the monotonous relationship between a single linear  
209 combination of all factors and the response. It varies in  $[0, 1]$ ; the closer to 1, the stronger  
210 the correlation between the ranks of the *best* combination of the set of factors  $\mathbb{X}$  (which  
211 corresponds to the term  $M_{\mathbb{X}}\hat{\beta}$  defined below) and the ranks of the response variable  $\mathbf{z}$ . To  
212 ensure that the performance indicator can effectively reach the maximum value 1, it is defined  
213 as the ratio between the Spearman correlation and its actual upper bound. The upper bound  
214 is computed analytically by following the reasoning proposed by Kendall<sup>21</sup> for treating ties  
215 in ranking problems. The reasoning adapted to our problem consists in identifying the  
216 situation where the Spearman correlation is maximum for given  $\mathbf{z}$  and  $M_{\mathbb{X}}\hat{\beta}$  ignoring the  
217 actual pairing between these sets of variables. This situation occurs when the ranks of

218 the  $(i, j)$ -th components of  $M_{\mathbb{X}}\hat{\boldsymbol{\beta}}$  and  $\mathbf{z}$  are equal for any  $(i, j)$  such that  $z_j^i \neq 0$ . Under  
 219 this assumption, the computation of the Spearman correlation simplifies and one obtains  
 220 an explicit expression for the upper bound. In what follows, we derive the expression of  
 221 the performance indicator; details on the computation of the upper bound of the Spearman  
 222 correlation are provided in Supporting Information, Appendix B.

The performance indicator is built from the following function of  $\boldsymbol{\beta}$ :

$$I_{\boldsymbol{\beta}}(\mathbb{X}, \mathbf{z}) = r_s^2(M_{\mathbb{X}}\boldsymbol{\beta}, \mathbf{z})(1 + \Delta_{M_{\mathbb{X}}\boldsymbol{\beta}, \mathbf{z}}), \quad (3)$$

where  $(1 + \Delta_{\mathbf{x}, \mathbf{y}})^{-1}$  is the upper bound of  $r_s^2(\mathbf{x}, \mathbf{z})$  and  $\Delta_{\mathbf{x}, \mathbf{y}} = \frac{\sum_{i \in I_0} (R_{x_i}^2 - R_{y_i}^2)}{(n-1)\hat{\sigma}_{R_{\mathbf{y}}}^2}$  for all  
 $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^n \times \mathbb{R}^n$ ,  $\hat{\sigma}_{R_{\mathbf{y}}}^2$  being the variance of  $R_{\mathbf{y}}$  and  $I_0 = \{i | y_i = 0\}$ , and where the elements  
 of the design matrix  $M_{\mathbb{X}} \in \mathbb{R}^{n \times n_c \times d'}$  are defined by:

$$M_{\mathbb{X}}(\ell, k) = \begin{cases} \frac{x_k^{(i,j)} - \min\{\mathbf{x}_k\}}{\max\{\mathbf{x}_k\} - \min\{\mathbf{x}_k\}}, & \text{if } \mathbf{x}_k \text{ is a continuous factor} \\ \left( \mathbf{1}_q(x_k^{(i,j)}) \right)_{q=1, \dots, Q_k-1}, & \text{if } \mathbf{x}_k \text{ is a discrete factor,} \end{cases}$$

223  $\ell = (i-1)n_c + j$  and  $k$  are the indices of the rows and the columns of the design matrix,  
 224 respectively,  $\min\{\mathbf{x}_k\}$  (resp.  $\max\{\mathbf{x}_k\}$ ) is the minimum (resp. maximum) element of the  
 225 vector  $\mathbf{x}_k$ ,  $d' = \sum_{k=1}^d (Q_k - 1)$ , with  $Q_k = 2$  if  $\mathbf{x}_k$  is a continuous factor and  $Q_k$  is equal to  
 226 the number of levels if  $\mathbf{x}_k$  is a discrete factor.

We then have to estimate the set of parameters  $\boldsymbol{\beta}$  which maximizes the Spearman corre-  
 lation  $r_s^2(M_{\mathbb{X}}\boldsymbol{\beta}, \mathbf{z})$  (to obtain the *best* combination of the set of factors  $\mathbb{X}$  with the form  
 $M_{\mathbb{X}}\boldsymbol{\beta}$ , as evoked above). The values of the components of  $\boldsymbol{\beta}$  associated with the fac-  
 tors identified as insignificant are set to zero, and the optimization is carried out with  
 respect to the remaining subset of parameters (of dimension  $d'' \leq d'$ ) using a genetic  
 algorithm implemented in the R package `rgenoud`<sup>28</sup> (the `genoud` function in this pack-  
 age combines an evolutionary search algorithm with a derivative-based Newton or quasi-  
 Newton method to solve optimization problems). Hence, we estimate  $\boldsymbol{\beta}$  as follows:  $\hat{\boldsymbol{\beta}} =$   
 $\arg \max_{\boldsymbol{\beta} \in \mathbb{R}^{d''}} r_s^2(M_{\mathbb{X}}\boldsymbol{\beta}, \mathbf{z})$ . The solution of this maximization is obviously not unique (if  $\boldsymbol{\beta}_0$  is  
 a solution,  $a\boldsymbol{\beta}_0$  is also a solution for all real value  $a \neq 0$ ), but this is not an issue in the pro-  
 posed framework since only the rank are taken into account and  $I_{\boldsymbol{\beta}_0}(\mathbb{X}, \mathbf{z}) = I_{a\boldsymbol{\beta}_0}(\mathbb{X}, \mathbf{z})$ ,  
 $\forall a \neq 0$ . In practice, since the maximizer can only be identified up to a scale factor,  
 each component of  $\boldsymbol{\beta}$  is constrained within the interval  $[-10, 10]$ , and the genetic algo-  
 rithm is stopped if the value of the objective function  $r_s^2(M_{\mathbb{X}}\boldsymbol{\beta}, \mathbf{z})$  has not increased in the  
 last 50 iterations (with a tolerance level equal to 0.001) or if the maximum number of  
 iterations (set to 200) has been reached. This is specified in the `genoud` function by set-  
 ting the following options: `Domains = matrix(c(-10, 10), byrow=TRUE, nrow=d'', ncol=2)`,

max.generations = 200, hard.generation.limit = TRUE, wait.generations = 50 and solution.tolerance = 0.001. Finally, we calculate the performance indicator by plugging in  $\hat{\beta}$ :

$$I_{\hat{\beta}}(\mathbb{X}, \mathbf{z}) = r_s^2(M_{\mathbb{X}}\hat{\beta}, \mathbf{z})(1 + \Delta_{M_{\mathbb{X}}\hat{\beta}, \mathbf{z}}).$$

## 2.4 Relative importance of factors

The optimal parameter vector  $\hat{\beta}$  must not be directly used for assessing the effect sizes of factors since it is not unique as explained in Section 2.3 (if  $\beta_0$  is a solution of the maximization,  $a\beta_0$  is also a solution for all real value  $a \neq 0$ ). Nevertheless, the components of  $\hat{\beta}$  can be used to compare the relative importance of factors in explaining, through  $M_{\mathbb{X}}\hat{\beta}$ , the links between targets and contributors summarized by  $\mathbf{z}$ . Thus, we define the relative importance of factor  $k \in \{1, \dots, d\}$  (or the level  $q$  of factor  $k$  for discrete factors) with respect to the *average* factor:

$$\tilde{e}_{k,q} = \frac{e_{k,q}}{\frac{1}{d'} \sum_{k'=1}^{d'} \sum_{q=1}^{Q_{k-1}} e_{k',q}}, \quad (4)$$

where  $e_{k,q} = |\hat{\beta}^{(k)}|$  and  $q = 1$  if factor  $k$  is continuous ( $\hat{\beta}^{(k)}$  being the component of  $\hat{\beta}$  corresponding to the continuous factor  $k$ ) and  $e_{k,q} = |\hat{\beta}^{(k,q)}|$  if factor  $k$  is discrete ( $\hat{\beta}^{(k,q)}$  being the component of  $\hat{\beta}$  corresponding to the  $q$ -th level of the discrete factor  $k$ ,  $q \in \{1, \dots, Q_k - 1\}$ ).

Remarks: (a)  $\tilde{e}_k$  is unchanged if one substitutes  $a\hat{\beta}$  for  $\hat{\beta}$  ( $a \neq 0$ ). (b) The operator  $M_{\mathbb{X}}$  in the regression  $M_{\mathbb{X}}\hat{\beta}$  homogenizes the amplitudes of variation of the factors and, therefore, of the coefficients in  $\hat{\beta}$  which are hence comparable. (c) The relative importance of factors measured by Equation (4) has to be understood in terms of factor contributions to the regression  $M_{\mathbb{X}}\hat{\beta}$  explaining  $\mathbf{z}$ . Hence, these contributions might be subject to nonlinear effects of factors or correlation between factors. One may alternatively compute a marginal indicator of the importance of each factor by calculating for example the performance indicator  $I_{\hat{\beta}}(\mathbb{X}, \mathbf{z})$  where  $\mathbb{X}$  is reduced to the factor of interest. Such a proposal however requires additional computation of the optima  $\hat{\beta}$  for every factor considered individually.

## 2.5 Rank prediction

If one ignores the value of  $\mathbf{z}$ , ranks of contributors can be predicted for any target by the ranks  $\hat{R}_{\mathbf{z}}$  of  $M_{\mathbb{X}}\hat{\beta}$ . In other words, the first contributor to a given target is predicted to be the contributor with the largest component of the sub-vector of  $M_{\mathbb{X}}\hat{\beta}$  restricted to the target under focus, the second contributor corresponds the second largest component, and so on. In the applications, we compare this ranking with rank predictions obtained from linear regression and decision tree (for which ranks are computed directly from the predictions

255 of proportions provided by these models) and with the rank-based estimation for linear  
 256 models<sup>[29]</sup>, which uses a distance based on a dispersion function<sup>[30]</sup> instead of the Euclidean  
 257 distance. Note that the four rank-prediction approaches are implemented with the same set  
 258 of factors identified by the tests presented in Section [2.2](#).

259 We use cross-validation to compare the robustness and the quality of the performance  
 260 indicator and the ranking obtained from our multitest-based multivariate analysis (MMA),  
 261 the linear regression model (LM), the rank-based linear regression model (LMRank) and the  
 262 decision tree (Tree). Target hosts are randomly divided into a train sample (80% of targets)  
 263 and a test sample (20% of targets). Tests for factor identification are applied to the global  
 264 sample (union of train and test samples), while indicators are computed separately for each  
 265 sub-sample. This procedure is independently repeated 100 times.

266 We consider two indicators: the performance indicator defined in the previous subsec-  
 267 tion and the contributor ranking indicator (CR). The CR indicator, defined in Supporting  
 268 Information, Appendix C, is the average over the targets of the proportion of the  $N_i$  con-  
 269 tributors with positive transmission probabilities for target  $i$  that are ranked among the top  
 270  $N_i$  contributors by the predictor under consideration (MMA, LM, LMRank or Tree).

## 271 3 Simulation study

272 We carry out a simulation study to investigate the performance of the proposed method.  
 273 All R codes to implement the methods have been incorporated into the package `ZIprop`,  
 274 freely available on R CRAN (<https://cran.r-project.org/package=ZIprop>) and GitLab  
 275 (<https://gitlab.paca.inrae.fr/meribaud/ziprop>).

### 276 3.1 Simulated data

277 We simulate data under the constraints described in Section [2.1](#). The algorithm applied to  
 278 simulate the factors and the response is described below:

- 279 1. Set values for  $n_c > 1$ ,  $n_t > 2$ ,  $m \in [1/n_c, 1]$  (the proportion of non-zero values for the  
 280 responses  $z_j^i$ ),  $d > 1$ ,  $\boldsymbol{\beta} \in \mathbb{R}^d$  and  $\eta > 0$ .
- 281 2. Generate the matrix  $\mathbb{X} = (\mathbf{x}_1, \dots, \mathbf{x}_d) \in \mathbb{R}^{n_t n_c \times d}$  such that the components of the  $n_t n_c$ -  
 282 tuple  $\mathbf{x}_k$ ,  $k \in \{1, \dots, d\}$ , are independently drawn from the continuous (resp. discrete)  
 283 uniform distribution  $U(0, 1)$  (resp.  $U(\{0, 1\})$ ) if  $\mathbf{x}_k$  is a continuous (resp. discrete)  
 284 factor.
- 285 3. Generate the un-constrained response vector  $\tilde{\mathbf{z}} \in \mathbb{R}^{n_t n_c}$  from the Gaussian distribution  
 286  $N(\mathbb{X}\boldsymbol{\beta}, \eta^2\mathbf{I})$  with mean vector  $\mathbb{X}\boldsymbol{\beta}$  and diagonal variance matrix  $\eta^2\mathbf{I}$  where diagonal

287 elements are equal to  $\eta^2$ ; subtract the minimum value of  $\tilde{\mathbf{z}}$  to each element of  $\tilde{\mathbf{z}}$  to get  
 288 only non-negative elements:  $\tilde{\mathbf{z}} \leftarrow \tilde{\mathbf{z}} - \min \tilde{\mathbf{z}}$ ; set to zero the  $n_0 = \lceil (1 - m)n_c n_t \rceil$  lowest  
 289 elements in  $\tilde{\mathbf{z}}$  excluding its maximal elements for each target  $i \in \{1, \dots, n_t\}$ , i.e., the  
 290  $n_0$  lowest elements in  $\tilde{\mathbf{z}} - \{\tilde{z}_j^i : \tilde{z}_j^i = \max_{j'} \{\tilde{z}_{j'}^i\}\}$  ( $\lceil \cdot \rceil$  is the ceiling function).

291 4. Compute the elements  $z_j^i$  of the response vector  $\mathbf{z}$  by scaling  $\tilde{z}_j^i$  for each target  $i \in$   
 292  $\{1, \dots, n_t\}$ :  $z_j^i = \frac{\tilde{z}_j^i}{\sum_{j'=1}^{n_c} \tilde{z}_{j'}^i}$ . Therefore, the response is simulated in such a way that for  
 293 any target  $i$ ,  $\exists j \in \{1, \dots, n_c\}$  such that  $z_j^i > 0$ , and  $\sum_{j=1}^{n_c} z_j^i = 1$ .

We test the effect of each factor and compute the performance indicator setting  $n_c = 20$ ,  
 $n_t = 22$ ,  $d = 20$ , and the proportion of non-zero data  $m \in \{0.1, 0.15, 0.2, 0.25\}$ . The first  
 (resp. last) half of factors are continuous (resp. discrete) and  $\boldsymbol{\beta} \in \mathbb{R}^{d'}$  ( $d' = d$  since discrete  
 factors have only two levels) satisfies:

$$\boldsymbol{\beta} = (\beta_1, \beta_2, \beta_3, -\beta_4, -\beta_5, 0, 0, 0, 0, 0, \beta_1, \beta_2, \beta_3, -\beta_4, -\beta_5, 0, 0, 0, 0, 0) \quad (5)$$

294 where  $\beta_k$ ,  $k = \{1, \dots, 5\}$ , are independently drawn from the following uniform distributions:

$$\beta_1, \beta_5 \sim U(8, 10), \quad \beta_2, \beta_4 \sim U(12, 14) \quad \text{and} \quad \beta_3 \sim U(16, 18). \quad (6)$$

296 The first 10 components of  $\boldsymbol{\beta}$  correspond to the continuous factors, the 10 following com-  
 297 ponents of  $\boldsymbol{\beta}$  correspond to the first level of the discrete factors with two levels (the second  
 298 level having a null effect). In addition, the standard deviation of the noise is set to  $\eta = 10$   
 299 that gives a median performance indicator of 0.69 between  $\mathbf{z}$  and  $\mathbb{X}\boldsymbol{\beta}$  for 100 runs. The main  
 300 characteristics of the simulation setting are given in Table 2.

301 In the simulation algorithm proposed above, the link between  $\mathbb{X}\boldsymbol{\beta}$  and  $\mathbf{z}$  is non-linear due  
 302 to the step where some values in  $\tilde{\mathbf{z}}$  are set to zero (Stage 3 of the algorithm) and the step where  
 303  $\tilde{\mathbf{z}}$  are scaled to obtain a vector of probabilities for each target (Stage 4). To increase the non-  
 304 linearity in an additional simulation study, we modify Stage 3 in the algorithm by generating  
 305  $\tilde{\mathbf{z}}$  in the Gaussian distribution  $N(f(\mathbb{X})\boldsymbol{\beta}, \eta^2\mathbf{I})$ , where  $f$  transforms four of the continuous fac-  
 306 tors included in  $\mathbb{X}$ , two with an expected effect on the response and two without effect given  
 307 the form of  $\boldsymbol{\beta}$  specified above:  $f(\mathbb{X}) = (\mathbf{x}_1^2, \exp(\mathbf{x}_2), \sqrt{\mathbf{x}_3}, \mathbf{x}_4, \mathbf{x}_5, \mathbf{x}_6^2, \exp(\mathbf{x}_7), \sqrt{\mathbf{x}_8}, \mathbf{x}_9, \dots, \mathbf{x}_d)$ .

### 308 3.2 Estimated errors of permutation tests

309 We assess the performance of the two-tailed permutation test for continuous and discrete  
 310 factors and different proportions  $m$  of non-zero data. Figure 2a shows the distribution of  
 311 p-values for each factor for  $m = 0.25$  (we get similar results for the other values of  $m$ , see  
 312 Supporting Information, Figure S2). The factors  $\mathbb{X}_{1:5}$  and  $\mathbb{X}_{11:15}$  are generally identified as  
 313 correlated to the response while  $\mathbb{X}_{6:10}$  and  $\mathbb{X}_{16:20}$  are not. The estimated type I errors of the

Table 2: Main specifications of the simulation study.

Object	Value
Number of Targets ( $n_t$ )	22
Number of Contributors ( $n_c$ )	20
Number of factors ( $d$ )	20
Type of factors	Continuous (10) and discrete (10)
Response variable ( $Z_j^i$ )	Simulated proportion
Number of observations ( $n_t \times n_c$ )	440

314 test at the risk level 0.05 are given in the top part of Table 3 for different values of  $m$ , and  
 315 show that the test is relatively well calibrated. The type II errors (bottom part of Table 3)  
 316 are very small for discrete factors whatever the value of  $m$ . In contrast, they are larger  
 317 for continuous factors (in particular those with relatively small effect) and decrease with  $m$ .  
 318 We carried out the same analysis for the one-tailed permutation tests and we obtained very  
 319 similar results as shown by Supporting Information, Tables S1 to S4.

320 Remark: Type II errors can be relatively large for continuous factors with weak effect (it  
 321 is around 0.6 when 10% of response values are non-zero), clearly showing a potential limit of  
 322 the proposed methodology. However, the large type II errors observed with small  $m$  might  
 323 be an artifact resulting from the simulation scheme. Indeed, by setting to zero the  $n_0$  lowest  
 324 elements in  $\tilde{\mathbf{z}}$  to build  $\mathbf{z}$  (see Stage 3 of the simulation algorithm), one simply deletes a part  
 325 of the information contained in the linear relationship between  $\mathbb{X}\boldsymbol{\beta}$  and the initial value of  $\tilde{\mathbf{z}}$ .  
 326 In other words, one gets  $z_j^i = 0$  for components of  $\mathbb{X}\boldsymbol{\beta}$  in a large range of values (i.e., small  
 327 and intermediate values). In contrast, in real cases,  $z_j^i = 0$  means that  $(i, j)$  is not likely to  
 328 be a target-contributor pair, and if  $\mathbb{X}\boldsymbol{\beta}$  is consistent with  $\mathbf{z}$ , its component corresponding to  
 329  $z_j^i = 0$  should tend to be only small, not intermediate, and should therefore reinforce the power  
 330 of the test. Further investigations are however required to test this artifact assumption.

### 331 3.3 Assessment of the performance indicator

332 For each repetition performed for  $m = 0.25$  (yielding the largest test power for continuous  
 333 variables), the performance indicator is computed for the  $k$  factors with the lowest p-values,  
 334  $k$  varying from 2 to 20. Figure 2b shows the distribution of the performance indicator with  
 335 respect to  $k$ . The indicator increases until it reaches a plateau at the value one (which is  
 336 its maximum value) approximately when  $k = 10$  (vertical line), which corresponds to the  
 337 actual number of factors having a significant effect. The indicator is robust in the sense that

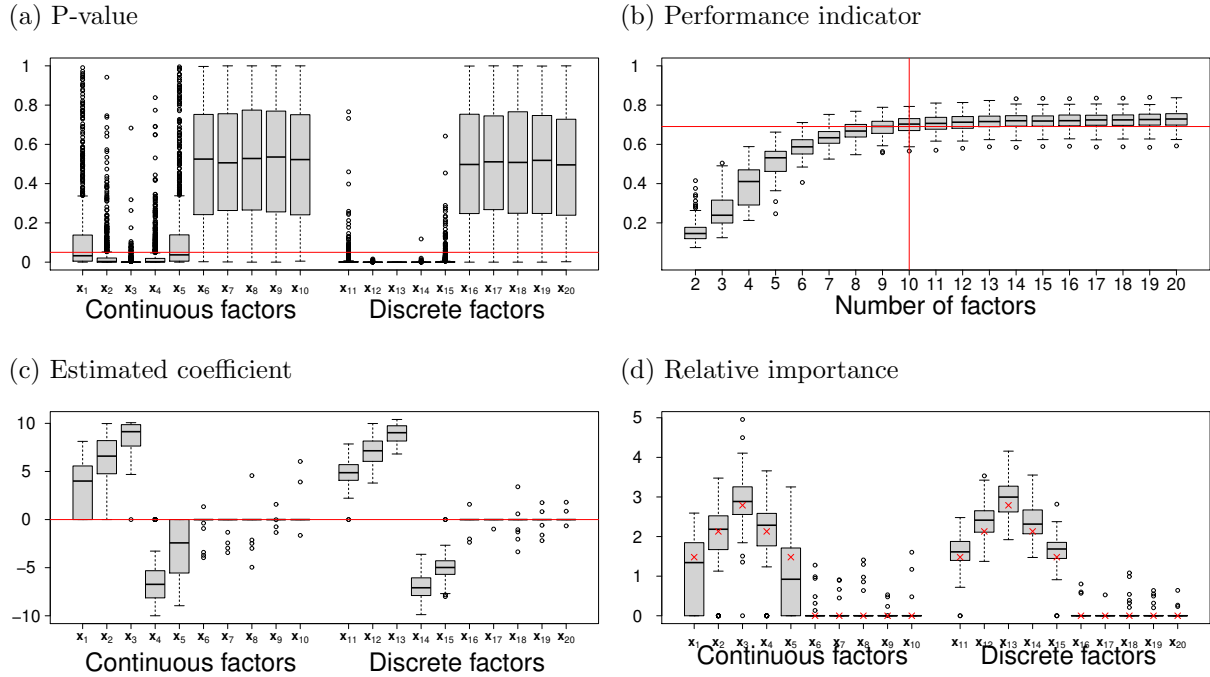


Figure 2: Factor significance and importance in the simulation study. a) P-values of two-tailed permutation tests for each factor (the red line indicates the 0.05 value). b) Distribution of the performance indicator for varying number of included factors (factors are successively incorporated by first including those with lowest p-values). The horizontal red line gives the median value of the performance indicator computed with the true value of  $\beta$  (0.69). c) Distribution of estimated coefficients (i.e., the components of  $\hat{\beta}$ ) for each factor. d) Distribution of the relative importance  $\tilde{e}_k$  of each factor; the red cross gives the expected relative importance given the simulation scheme described in Section 2.4. The distributions are drawn with  $m = 0.25$  and from 1000 repetitions for a) and 100 repetitions for b), c) and d).

338 adding more factors than the actual number of factors with significant effects does not affect  
 339 the performance.

### 340 3.4 Assessment of the relative importance of factors

341 Figure 2c shows that estimated coefficients for non-significant factors are close to zero, while  
 342 they take values that are clearly positive or negative for significant factors in agreement with  
 343 the specifications given by Equations (5)–(6). In addition, we note that the amplitudes of  
 344 the estimated coefficients are not correct, as expected, since the coefficients were constrained  
 345 between -10 and +10 in the optimization process whereas they were simulated between -14  
 346 and +18; see Equations (5)–(6). However, the relative values of coefficients are estimated



Table 3: Estimated type I and type II errors of the two-tailed permutation tests based on 1000 repetitions in the simulation study ( $m$  is the proportion of non-zero values for the response).

<b>Type I</b>		Continuous factors					Discrete factors				
m	$\mathbf{x}_6$	$\mathbf{x}_7$	$\mathbf{x}_8$	$\mathbf{x}_9$	$\mathbf{x}_{10}$	$\mathbf{x}_{16}$	$\mathbf{x}_{17}$	$\mathbf{x}_{18}$	$\mathbf{x}_{19}$	$\mathbf{x}_{20}$	
0.1	0.043	0.044	0.046	0.036	0.047	0.046	0.047	0.049	0.063	0.052	
0.15	0.046	0.049	0.051	0.042	0.040	0.054	0.043	0.054	0.054	0.054	
0.2	0.057	0.043	0.043	0.056	0.039	0.051	0.034	0.055	0.052	0.045	
0.25	0.058	0.043	0.049	0.046	0.050	0.041	0.044	0.055	0.052	0.049	

<b>Type II</b>		Continuous factors					Discrete factors				
m	$\mathbf{x}_1$	$\mathbf{x}_2$	$\mathbf{x}_3$	$\mathbf{x}_4$	$\mathbf{x}_5$	$\mathbf{x}_{11}$	$\mathbf{x}_{12}$	$\mathbf{x}_{13}$	$\mathbf{x}_{14}$	$\mathbf{x}_{15}$	
0.1	0.607	0.293	0.097	0.321	0.616	0.147	0.008	0.000	0.004	0.159	
0.15	0.492	0.214	0.044	0.223	0.522	0.093	0.002	0.000	0.001	0.086	
0.2	0.449	0.176	0.037	0.189	0.460	0.063	0.001	0.000	0.001	0.066	
0.25	0.433	0.147	0.021	0.151	0.447	0.046	0.000	0.000	0.001	0.041	

347 satisfactorily as deduced from Figure 2d, which shows that the relative importance  $\tilde{e}_k$  of  
 348 each factor  $k$  is approximately distributed around its expected value.

349 Very similar results are obtained when the standard deviation of the noise  $\eta$  is decreased  
 350 to  $\eta = 5$  or increased to  $\eta = 15$ , even if the largest standard deviation leads to wider  
 351 distributions of the relative importance (and higher p-values for factors with non-zero effects,  
 352 i.e., larger type II errors); see Supporting Information, Figures S3-S4, which are analogous  
 353 to Figure 2.

354 Similar results are also obtained when we include additional non-linearity with the trans-  
 355 formation function  $f$  (see Supporting Information, Figure S5) the sum  $\sum_{j=1}^{n_c} z_j^i$  is drawn from  
 356 a uniform distribution between 0.5 and 1 (which mimics the non observation of some con-  
 357 tributors), instead of fixing  $\sum_{j=1}^{n_c} z_j^i$  to the value 1 as described in stage 4 of the simulation  
 358 algorithm detailed in Section 3.1 (see Supporting Information, Figure S6).

359 In addition, Supporting Information, Figure S7, shows that the estimates of non-zero  
 360 coefficients are pushed towards the limits of the range, because of the non-identifiability  
 361 evoked at the end of Section 2.3, when we consider wider constraining intervals for the  
 362 optimization of  $\beta$ .



### 3.5 Assessment of the rank prediction with cross-validation

The methodology proposed in this article (MMA) is compared in terms of ranking performance for the simulated data with the three other methods (LM, LMRank and Tree) introduced in Section 2.5. The prediction methods are applied to the set of factors selected by the permutation tests. Figure 3.a shows the good performance of the multitest-based multivariate analysis (MMA). The linear models (LM and LMRank) have relatively similar efficiency, whereas the decision tree (Tree) is clearly less efficient when it is applied to the test samples. Similar results are obtained when one includes the transformation function  $f$  to increase the non-linearity (see Figure 3.b), and when one considers the contributor ranking indicator (CR) instead of the performance indicator except for LMRank that is less efficient than MMA and LM based on this criterion (see Supporting Information, Figure S8).

## 4 Application I: Equine Influenza

We consider an Equine Influenza outbreak in 2003 in race horses from different training yards in Newmarket. Genomic data collected during this outbreak from 48 horses were studied to explore the virus transmissions across the observed horse population<sup>31</sup>. Intra-host sequences were obtained for each horse and these sequences were used to estimate the probabilities of disease transmission between hosts using the BadTrIP software<sup>32;33, chap. 3</sup>. BadTrip was run in BEAST2<sup>34</sup> using two independent MCMC chains of 5 million steps. We used the dates of first positive swabs as epidemiological data in BadTrIP and allowed the horses to be infected for 9 days except for the first horse A01, which was allowed to be infected for 15 days to provide overlap in the infection periods of the first horses in the transmission chain. The estimated transmission probabilities (shown in Supporting Information, Figure S9) are used in the present study as response proportion data; see Table 4. Many of the estimated transmission probabilities are equal to zero, which means that, for each target, BadTrIP identified only a small number of potential contributors. In what follows, we use four discrete factors and one continuous factor computed from the observed variables ‘age’, ‘sex’ and ‘training yard’ described in Table 5.

Some of these factors are missing for some target-contributor pairs (see Table 5 and Supporting Information, Table S5). Hence, the tests for assessing the effect of a given factor on the transmission probability are applied on the subset of complete data for this factor. Permutation tests are applied factor by factor on subsets without missing values (if any). The data set used for this study is available on a public archive repository<sup>35</sup>. Remark: A non-completely random structure of missing values could generate biased results, but we did not observe clear signs of such a structure in this case study. We simply observed that the age and sex variables were systematically missing for one of the yards with several infected

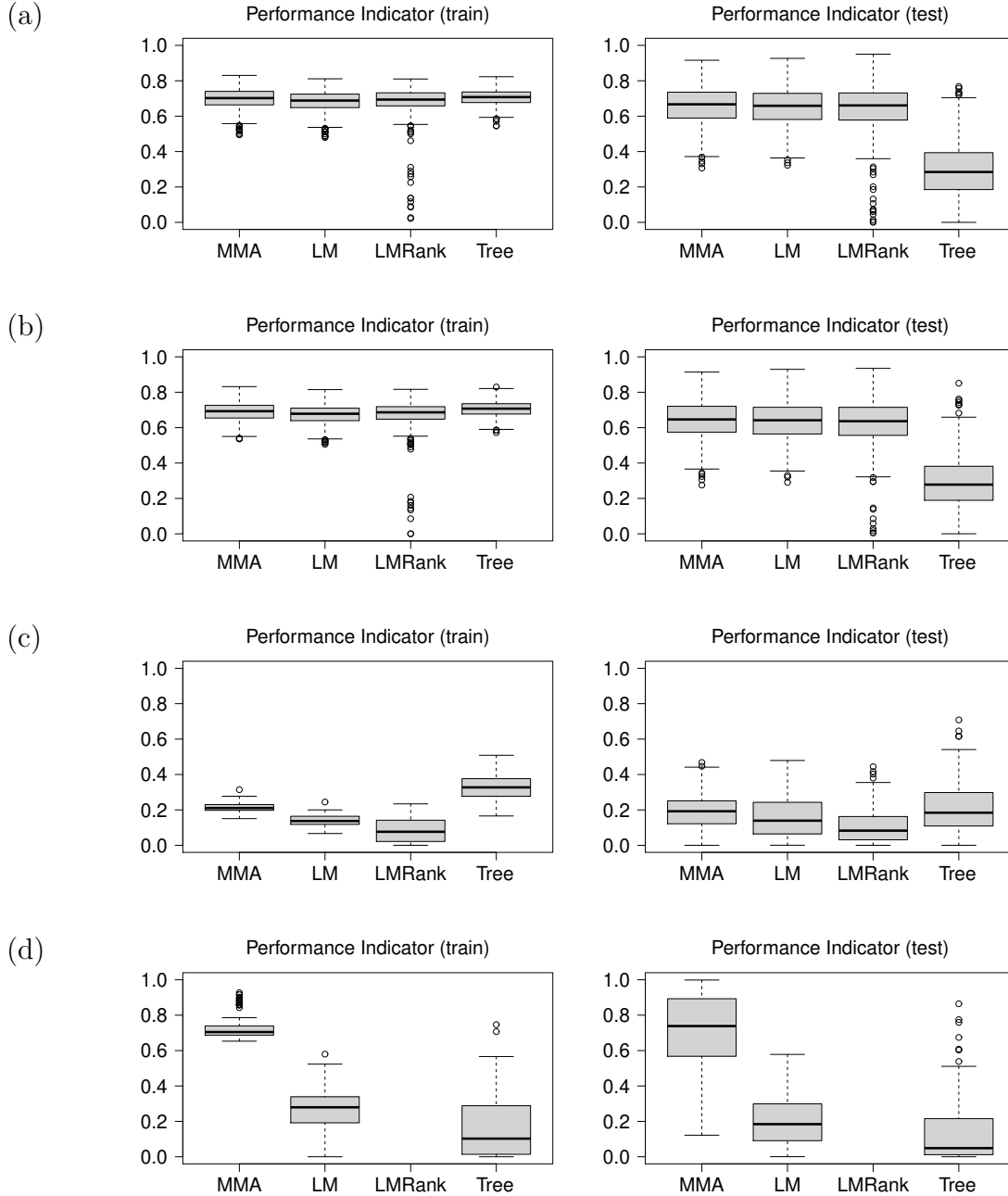


Figure 3: Boxplots of the performance indicator calculated from the training and test samples for MMA, LM, LMRank and Tree: a) in the simulation study without the non-linear transformation  $f$ ; b) in the simulation study with the non-linear transformation  $f$ , c) in the Equine Influenza application; d) in the COVID-19 application (the performance indicator could not be computed for LMRank in this case because this method only predicts null probabilities, which result on a null variance of the ranks and, therefore, an undefined value for the performance indicator).

Table 4: Main specifications of the equine influenza application. Note that the number of observations is  $n_t \times (n_c - 1)$  and not  $n_t \times n_c$  because we do not account for auto-infection (i.e. for pairs where the contributor-horse is also the target-horse).

Object	Value
Targets/Contributors	Horses
Number of Targets/Contributors ( $n_t/n_c$ )	48
Type of factors	Variables related to horses and their environment
Number of factors ( $d$ )	5
Response variable ( $Z_j^i$ )	Transmission probabilities
Number of observations ( $n_t \times (n_c - 1)$ pairs)	2256

Table 5: Explanatory factors for the equine influenza application.

Factor	Description	Number of NA (%)
Same_Yard	1 if the target and contributor are trained in the same yard 0 otherwise	0 (0%)
Same_Sex	1 if the target and contributor have the same sex 0 otherwise	1606 (71%)
Diff_Age	0 if the target and contributor have the same age 1 for a one-year difference 2 for more than one year	1554 (69%)
Dist_Yard	geographic distance (in km) between the training yards of the target and the contributor	0 (0%)
Trans_Sex	“F→F” if a female infected another female “M→F” if a male infected a female “F→M” if a female infected a male “M→M” if a male infected another male	1606 (71%)

398 horses but nothing else in particular characterizes this yard based on the available data.

399 Factors Same\_Yard, Same\_Sex, Dist\_Yard and Trans\_Sex are significantly correlated to

Table 6: Test results for the equine influenza application. Top: Statistic ( $T$ ), p-value (pv) and Spearman’s correlation ( $r_s$ ; for the continuous factor only) associated with the two-sided permutation tests performed for the five factors. Bottom: Statistic ( $T$ ), p-value (pv), Hochberg-corrected p-value (pv\*); for discrete factors with more than 2 levels) associated with the post-hoc permutation tests applied to significant discrete factors. Lines with a significant p-value are highlighted in gray.

Factor	$T$	pv	$r_s$
Same_Yard	0.05	0	
Same_Sex	0.007	0.031	
Diff_Age	0.001	0.8	
Dist_Yard	0.05	0	-0.22
Trans_Sex	0.042	0	

Factor	Factor level	$T$	pv	pv*
Same_Yard	0 - 1	-444	0	
Same_Sex	0 - 1	-24.77	0.04	
Trans_Sex	F→F - F→M	65	0	0.01
	F→F - M→F	111	0	0
	F→F - M→M	80	0	0
	F→M - M→F	46.2	0.01	0.02
	F→M - M→M	15	0.33	0.33
	M→F - M→M	-31.1	0.04	0.08

400 the transmission probability whereas Diff\_Age is not; see Table 6. Among these factors,  
401 Same\_Yard, Same\_Sex and two modalities of Trans\_Sex have the largest and comparable  
402 relative importance with respect to the average factor; see Figure 4. The post-hoc statistic  
403 of Same\_Yard ( $T$  in the bottom part of Table 6) and the Pearson correlation of Dist\_Yard ( $r_s$   
404 in the top part of Table 6) being negative, horses trained in the same yard or in nearby yards  
405 have a higher chance to be linked by a transmission. This is a clearly intuitive result certainly  
406 due to higher contact rate in shared training areas. The statistics of post-hoc univariate tests  
407 for factor Same\_Sex is also negative, which means that the virus better circulates between  
408 horses with the same sex. Moreover, the post-hoc tests on the Trans\_Sex modalities show  
409 that only the difference between “F→M - M→M” (and “M→F - M→M” when one considers  
410 the corrected p-values) are not significant. The results on the p-values and the sign of the  
411 statistics show that transmissions between females are favored compared to all other possible

412 combinations ( $F \rightarrow F$  transmissions have positive probabilities 1.8 times more than expected  
413 under complete randomness; see Table D.1 in Supporting Information, Appendix D.

414 In addition, there is more intersex transmission when females are the sources ( $F \rightarrow M$ ),  
415 than when males are the sources ( $M \rightarrow F$ ). Supporting Information, Appendix D, shows that  
416 the significance of gender-related factors is neither confounded with the effect of the other  
417 available factors nor a consequence of heterogeneous sex frequencies.

418 The performance indicator is calculated on the table containing the four selected factors,  
419 discarding the transmissions containing one or more missing values (NA). The performance  
420 indicator takes the value  $I_{\hat{\beta}}(\mathbb{X}, Z) = 0.21$  using the four selected factors. This relatively low  
421 value, which indicates that there is a moderate correlation between the combination of the  
422 four factors and the transmissions, can actually be viewed as quite large given the fact that  
423 we only consider very basic factors to *predict* the transmissions.

424 To investigate the robustness and the quality of our approach in this case study, we apply  
425 cross-validation and perform the comparison with the three benchmark methods presented in  
426 Section 2.5. For this comparison we use the four factors selected with the permutation tests.  
427 The decision tree seems to outperform the multivariate analysis on the training samples but  
428 this is not confirmed on the test samples; see Figure 3c. Both methods tend to be more  
429 efficient than the two linear models, possibly because of the highly discrete nature of the  
430 factors: three factors among four are discrete and the continuous factor (Dist\_Yard) takes  
431 only 37 different values out of 650 observations. Similar conclusions are drawn when the  
432 ranking performance is measured with the CR indicator; see Supporting Information, Figure  
433 S10. Remark. In an additional simulation study, we fixed all the coefficients of the continuous  
434 factors at zero to see whether the performance of the decision tree is improved when only  
435 discrete factors have an impact on the response variable. No significant improvement of the  
436 decision tree performance was observed (compare Figure 3a and Supporting Information,  
437 Figure S11) and further investigations are required to understand the reason why the decision  
438 tree is relatively efficient in the influenza case study.

## 439 5 Application II: COVID-19

440 A recent article proposed a data-driven method to predict the mortality curve of a target  
441 country with a mixture of the mortality curves of countries that are ahead of time in terms  
442 of mortality rate<sup>2</sup>. The mixture is more exactly formed by the mortality curves of *contribut-*  
443 *ing countries* as well as an additional parametric predictor, and the method is essentially  
444 grounded on the estimation of the mixture probabilities. Real-time predictions based on  
445 this method are available for more than 100 countries via the following web application:  
446 <http://covid19-forecast.biosp.org/>.

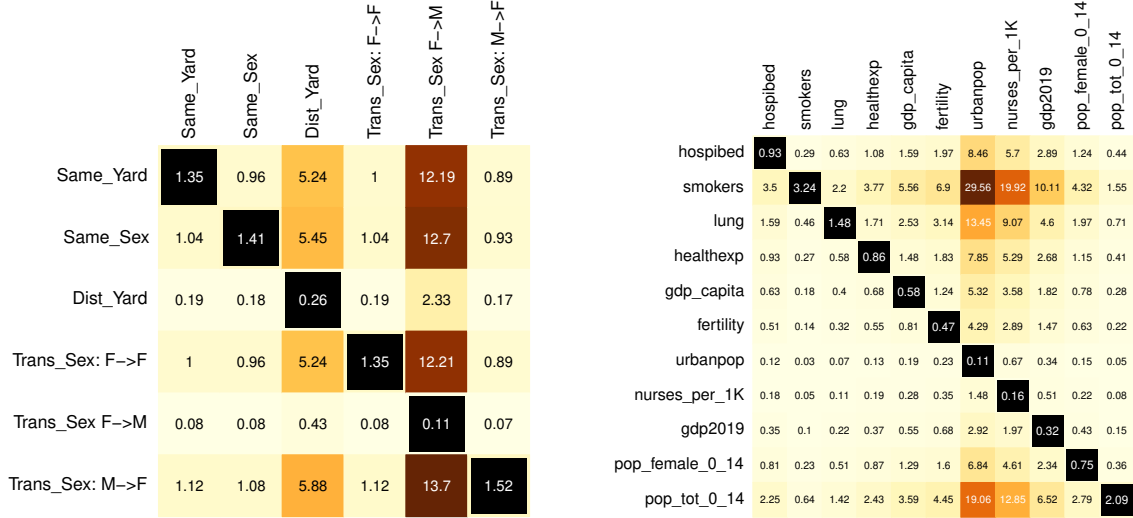


Figure 4: Relative importance of significant factors for the Equine influenza (left) and COVID-19 (right) applications. Diagonal values give the relative importance  $\tilde{e}_k$  of each factor  $k$  with respect to the average factor as defined by Equation (4). Non-diagonal values give the relative importance of any factor with respect to any other factor, i.e.  $\tilde{e}_k/\tilde{e}_{k'} = e_k/e_{k'}$ ; large values are highlighted with brownish colors.

447 Here, we use the estimated mixture probabilities as proportion data. Targets are states  
 448 from the USA and provinces from Canada; contributors are members of the European Eco-  
 449 nomic Area (EEA) and the European Free Trade Association (EFTA). We only consider  
 450 geographic entities with at least 5,000,000 inhabitants (leading to 23 targets and 21 con-  
 451 tributors) and the first epidemic wave by using data up to June 6, 2020. Mortality data  
 452 used to estimate the mixture probabilities were collected from the Johns Hopkins Uni-  
 453 versity Center For Systems Science and Engineering<sup>36</sup> and The Covid Tracking Project  
 454 (<https://covidtracking.com>). The choice of considering Northern American targets and  
 455 European contributors was made because Europe was on average ahead of time in terms of  
 456 mortality rate, at least during the first COVID-19 epidemic wave.

457 To explain the mixture probabilities (i.e., the similarity between targets and contributors  
 458 in terms of mortality dynamics), we consider 29 variables related to economy, demography,  
 459 health, healthcare system and climate; see Table 7 and Supporting Information, Table S6.  
 460 More precisely, our objective is to identify factors negatively correlated with the response, i.e.,  
 461 the lower the distance between two geographic entities with respect to a given variable, the  
 462 higher the mixture probability. Consequently, we use the univariate test (iii) for continuous  
 463 factors, which are computed for each target-contributor pair by  $x_k^{(i,j)} = |x_k^i - x_k^j|$ . The data  
 464 set used for this study is available on a public archive repository<sup>37</sup>.

Table 7: Main specifications of the Covid-19 application.

Object	Value
Targets	States from the USA and provinces from Canada
Number of Targets ( $n_t$ )	23
Contributors	Members of the European Economic Area
Number of contributors ( $n_c$ )	21
Type of factors	Variables related to economy, demography, health, healthcare system and climate
Number of factors ( $d$ )	29
Response variable ( $Z_j^i$ )	Mixture probabilities
Number of missing values	0
Number of observations ( $n_t \times n_c$ )	483

Figure 5 shows the p-values obtained for each factor and the Spearman’s correlation for significant factors. We identified eleven impacting factors whose definitions are provided in Supporting Information, Table S6: hospibed, smokers, lung, healthexp, gdp\_capita, fertility, urbanpop, nurses\_per\_1K, gdp2019, pop\_female\_0\_14, and pop\_tot\_0\_14. The figure shows that Spearman’s correlation is negative for significant factors. This result is consistent with our objective: to identify the significant factors negatively correlated to the response (since we expect that the similarity of the mortality dynamics of two countries decreases when the difference in the factor values for the two countries increases).

Then, we applied the multivariate analysis based on the eleven significant factors. The factors smokers (percentage of smokers within the population), pop\_tot\_0\_14 (percentage of population in the age group 0-14) and lung (death rate for lung diseases per 100,000 people) have the largest relative importance in the linear combination of factors explaining the mortality dynamics similarity; see Figure 4. Hence, small deviations of these factors strongly favor the similarity of mortality curves. In addition, the performance indicator is equal to  $I_{\hat{\beta}}(\mathbb{X}, Z) = 0.73$ , which shows that a high monotonous dependency exists between the mixture probabilities and these factors. This is confirmed by Supporting Information, Figure S12, which illustrates the relative good match between predicted ranks based on selected factors and mixture probabilities interpreted as probabilities of similarity.

Figure 3.d and Supporting Information, Figure S13, show the high performance of the MMA with respect to LM and Tree (we could not provide, for this case study, the performance of LMRank as explained in the caption of Figure 3). The performance indicator is quite variable when it is computed from the test samples. In comparison the CR indicator is

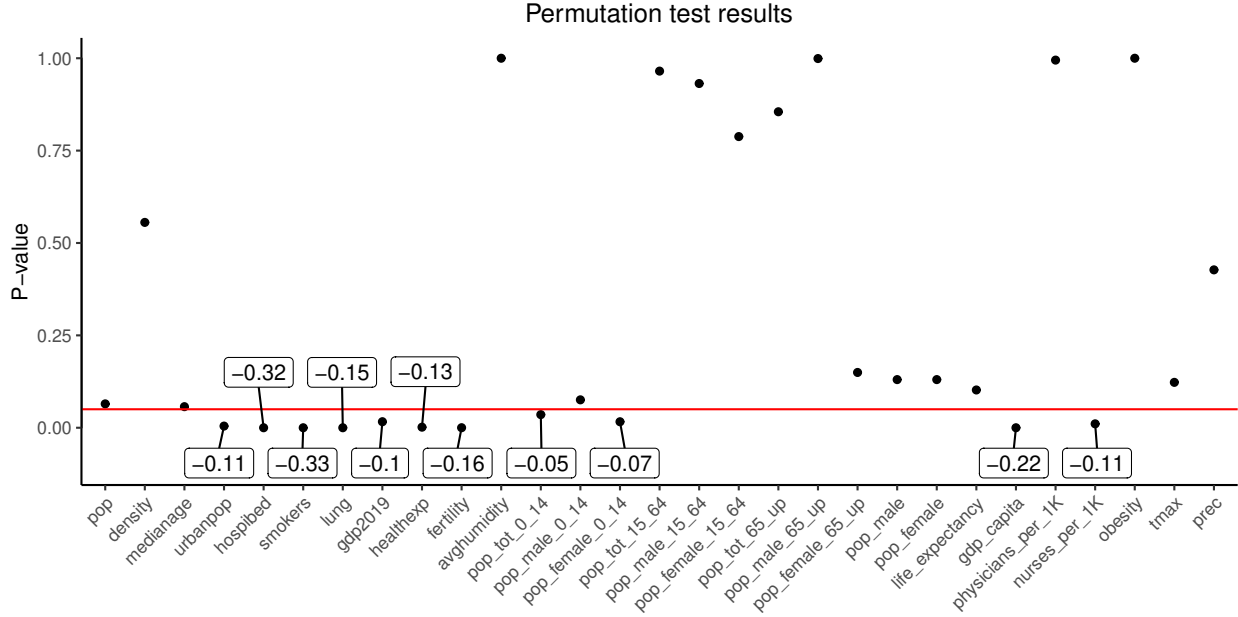


Figure 5: Comparison of the p-values (black dots) for the univariate tests (iii) associated to each factor and the significance level  $\alpha = 0.05$  (horizontal line). The Spearman correlation is given by the framed value for factors with a p-value below the threshold.

487 less variable. However, in both cases, MMA remains globally more efficient than the other  
 488 methods.

## 489 6 Discussion

490 For exploring the relationships between zero-inflated proportion data (ZIPD) and poten-  
 491 tially impacting factors, we developed permutation tests explicitly taking into account the  
 492 dependence structure in the proportions. These tests being uni-dimensional (i.e., each fac-  
 493 tor being treated separately), we propose a posterior analysis simultaneously handling the  
 494 multiple factors. This posterior analysis is grounded on a performance indicator quantifying  
 495 the percentage of correlation explained by the subset of significant factors and enables the  
 496 ranking of proportions based on the observation of the factors. Ranking the proportions as  
 497 we proposed is particularly useful in epidemiology because it especially allows the identi-  
 498 fication of the most likely sources of infection for a given recipient host simply from host  
 499 characteristics.

500 Based on the simulation study, the uni-dimensional tests are well calibrated with respect  
 501 to type I error in the situations under consideration and, overall, type II error is satisfactory



502 and generally lower for discrete factors than for continuous factors. These results are con-  
503 sistent with the generally-observed reliability of permutation tests. The main assumption  
504 of permutation tests (that are distribution-free) is the exchangeability of data under the  
505 null hypothesis. However, this assumption can be violated when the dependence structure  
506 is known and especially in the case of block permutations<sup>16</sup>. This structure must be taken  
507 into account during permutations otherwise the tests lose power (see Supporting Informa-  
508 tion, Appendix A). This is why we test the null hypothesis only on a subset of all possible  
509 permutations.

510 The performance indicator, allowing a simultaneous treatment of all the factors, was  
511 shown to provide a relevant and parsimonious description of the strength of the link between  
512 ZIPD and multiple factors. Indeed, a plateau in terms of performance is reached when all  
513 the significant factors are added. We also observed a relative robustness in the performance  
514 using training and test sets in a cross-validation framework applied to the simulation and  
515 real studies. Cross-validation was also used to compare the performance of our multivariate  
516 analysis with respect to the linear regression model, the rank-based linear regression model  
517 and the decision tree. Contrasted results were obtained across the three studies, but the  
518 multitest-based multivariate analysis appeared as a versatile approach adapted to continuous,  
519 discrete and mixed factors.

520 We have challenged our approach in various settings. However, several particularities of  
521 the data not tested in the simulation study could be considered to deepen the conditions  
522 of validity of the tests and of the performance indicator in particular (for example, partic-  
523 ularities concerning the collinearity or dependence between the factors, the proportion and  
524 the structure of missing values, and the number of modalities of discrete factors). The code  
525 accompanying this article will facilitate the exploration of the limitations of our approach.

526 **Methodological perspectives** Target/contributor pairs containing missing values (NA)  
527 are handled differently for the permutation tests and for the calculation of the indicator (see  
528 application on equine influenza in Section 4). Since the permutation tests are applied factor  
529 by factor, the pairs containing a NA for a factor are removed only for the test associated to this  
530 factor. In contrast, for the calculation of the performance indicator, which simultaneously  
531 handles the factors, any pair containing a NA for any selected factor is ignored. A more subtle  
532 treatment of NA, in particular for the calculation of the indicator, deserves to be explored.

533 If the zero-inflation in the response variable is accounted for in the performance indicator  
534 using a scaling term (see Equation (3)), it is not specifically addressed in the permutation  
535 tests because these tests are valid despite the zero-inflation and, in practice, work relatively  
536 well. Nevertheless, balancing techniques<sup>8</sup> such as those evoked in the introduction might be  
537 adapted to the dependency structure<sup>38</sup> and applied to eventually improve the performance

538 of the approach. Other approaches grounded on distinct treatments of zeros and non-zeros  
539 values and a resulting modification of the test statistics may also be considered to explicitly  
540 handle the zero-inflation in the permutation test<sup>39,40,41</sup>.

541 A non-negligible computational cost in our approach is due to the optimization of  $\beta$  with  
542 a genetic algorithm (the dimension of the optimization domain is equal to the number of  
543 significant factors retained with the uni-dimensional tests). This cost might be reduced with  
544 linear algebra and analysis tools<sup>42</sup>.

545 The method that we propose only provides a point estimate of the performance indicator  
546  $I_{\hat{\beta}}(\mathbb{X}, \mathbf{z})$  (that depends on  $\mathbb{X}$  and  $\hat{\beta}$  via the term  $M_{\mathbb{X}}\hat{\beta}$ ) and point predictors of the ranks  
547 of  $\mathbf{z}$  (i.e.,  $\hat{R}_{\mathbf{z}}$ ) within each target-contributors block using  $M_{\mathbb{X}}\hat{\beta}$  (see Sections 2.3-2.5). The  
548 robustness of these estimates can be approached by the cross-validation technique that we  
549 use in this manuscript. A more advanced evaluation of their uncertainty could be developed  
550 in further study by propagating two sources of uncertainty: (i) the uncertainty resulting  
551 from the test outputs that lead the components of  $\beta$  to be set to zero or to be optimized  
552 in the maximization of  $r_s^2(M_{\mathbb{X}}\beta, \mathbf{z})$ , (ii) the uncertainty associated to the maximization  
553 of  $r_s^2(M_{\mathbb{X}}\beta, \mathbf{z})$  given the components of  $\beta$  set to zero. If a method can be developed to  
554 account for both sources of uncertainty, one would be able to assess the uncertainty of  
555  $M_{\mathbb{X}}\hat{\beta}$  and subsequently assess the uncertainty of the performance indicator estimate and the  
556 rank predictors. Practically, non-parametric bootstrap may allow us to derive an empirical  
557 distribution of  $M_{\mathbb{X}}\hat{\beta}$  and, therefore, empirical distributions of  $I_{\hat{\beta}}(\mathbb{X}, \mathbf{z})$  and  $\hat{R}_{\mathbf{z}}$ .

558 The indicators that we use to evaluate the performance of the method in the simulation  
559 study and the applications, namely the performance indicator  $I$  and the contributor ranking  
560 indicator CR, reflect distinct properties of our approach:  $I$  measures the adequateness of the  
561 prediction of all proportion ranks whereas CR focuses on the rank prediction for proportions  
562 whose actual values are positive. Additional indicators could be envisioned. For instance, one  
563 could consider an indicator evaluating how much the target-contributor pair with the largest  
564 proportion is correctly ranked for any target (in other words, this indicator would assess the  
565 ability of the method(s) to identify, for any target, the most likely target-contributor pair).

566 Interestingly, the uni-dimensional tests can be applied as a first stage for performing an  
567 initial factor selection, whatever the posterior multivariate analysis that is carried out. Here,  
568 in addition to our multivariate analysis, we considered relatively simple tools for the poste-  
569 rior analysis, namely linear regression and decision tree. We could consider more complex  
570 approaches with known ability to handle, e.g., non-linearity and interactions of high order,  
571 and optionally embedding an additional factor-selection stage to solve possible issues gener-  
572 ated by eventual dependence between factors selected with our unit tests. Thus, one could  
573 explore for example the use of neural networks<sup>43</sup>, multivariate adaptive regression splines<sup>44</sup>,  
574 random forest<sup>45</sup> and boosted generalized linear models<sup>46,47</sup>.

575 In the two applications that we considered, the response variable (i.e., the transmission  
576 probability or the similarity probability) is actually estimated from external data and we  
577 only consider a point estimate of the probability for each pair. Thus, we deal with the  
578 first level of uncertainty, namely the fact that the transmission or the similarity is uncertain  
579 and therefore represented by a probability instead of a true/false variable. However, we  
580 do not handle the second level of uncertainty, namely the uncertainty of the probability  
581 estimates. This source of uncertainty could be handled as follows: Suppose that we have at  
582 disposal distributions of probabilities (e.g., posterior distributions obtained from a Bayesian  
583 approach) instead of point estimates, then we could propagate the uncertainty about the  
584 probabilities<sup>48</sup> into our test by (i) sampling the probabilities from their distributions and  
585 (ii) applying the permutation to each sample of probabilities (at stage 2 of the conditional  
586 Monte-Carlo algorithm described in Section 2.2). Moreover, by sampling the probabilities  
587 from their distributions, we could obtain the distribution of the performance indicator given  
588 by Equation (3) instead of a single value. It has however to be noted that if the estimation  
589 of  $\mathbf{z}$  is biased and hence the weights of network edges are misspecified, the method proposed  
590 in this article will certainly miss impacting factors and possibly lead to the identification of  
591 unimportant factors as significant.

592 **Epidemiological perspectives** The study of Equine Influenza data confirmed the obvi-  
593 ous importance of direct contact between hosts for virus transmission: the more frequent  
594 contact between horses (same or nearby yard), the higher the probability of transmission.  
595 The interpretation of sex differentiation in transmission potential is more complex. We have  
596 ruled out an issue of confounding effect between sex and other variables available in the data  
597 set and an issue of sex balancing in the observed population of horses. Behavioral, immuno-  
598 logical, physiological or organizational factors should be explored to unravel the mechanisms  
599 explaining the excess of female-to-female transmissions and transmissions in which a female  
600 is the source of infection. Typically, boys were shown to more likely transmit H1N1 to boys  
601 and girls to girls probably as a result of assortative mixing among playmates<sup>49</sup>. Another po-  
602 tential explanation could stem from the tendency of horse owners to group horses according  
603 to gender, in an attempt to reduce aggressive interactions and the risk of injuries<sup>50</sup>. More  
604 generally, given the low value of the performance indicator (0.21), it would be interesting to  
605 introduce other factors in the analysis of transmissions by considering other equine influenza  
606 data sets, which would allow the results obtained in this article to receive further checks.

607 Beyond this case study, estimating the probabilities that individuals are linked by trans-  
608 mission events during epidemics of infectious diseases or by progeny relationships in popu-  
609 lation dynamics has been the subject of numerous studies tackled, in particular, with joint  
610 models of epidemiological dynamics and evolutionary processes<sup>1,51,52,53,54,55,56</sup> or with phy-

611 logeny, phylogeography and some forms of birth–death processes<sup>[32,57,58,59,60,61,62]</sup>. In many  
612 of these studies, it would be interesting to take matters further by exploring the statistical  
613 relationship between the inferred links (generally corresponding to ZIPD) and factors char-  
614 acterizing the individuals and the environment. Indeed, determining how factors favor the  
615 spread of pathogens or species is crucial to better understand the underlying dynamics<sup>[63]</sup>  
616 and to design adequate control or conservation strategies. In the phylogeography litera-  
617 ture, a framework grounded on randomization was proposed to test hypotheses about the  
618 effect of environmental variables on pathogen spatial spread<sup>[64]</sup>, but this framework requires  
619 the spatio-temporal reconstruction of phylogenetic trees using a software such as BEAST.  
620 In contrast, the method that we propose can be applied to random transmission trees and  
621 random phylogenetic trees, whatever the way these trees are obtained.

622 The study of COVID-19 data allowed us to correlate similarity in COVID-19 mortality  
623 curves with similarity in certain macroscopic factors related to demography, economy, popu-  
624 lation health and healthcare system. The high value of the performance indicator (0.73) and  
625 its relative robustness observed via cross-validation indicate that these factors can be used  
626 to predict with a certain accuracy the similarity between COVID-19 mortality dynamics in  
627 geographic entities of Northern America and Europe. Thus, these factors can be viewed as  
628 characteristics intrinsically measuring the preparedness and/or vulnerability of geographic  
629 entities. The unexplained part of the rank correlation between the mixture probabilities and  
630 the linear combinations of factor differences could be due, in particular, to heterogeneous ini-  
631 tial conditions of the outbreaks and heterogeneous control measures. It would be interesting  
632 to include in the analysis explanatory variables reflecting these components of the epidemics.  
633 Furthermore, the analysis that we performed is based on mortality data up to June 6, 2020,  
634 which approximately correspond to the first wave of the COVID-19 epidemics in Northern  
635 America and Europe. It would be interesting to repeat this analysis across time to assess the  
636 temporal (in-)stability of our findings. Specifically, one could expect a change in significant  
637 factors between the first and the second (or subsequent) epidemic waves, especially because  
638 of heterogeneous levels of immunity after the first wave and heterogeneous impacts of the  
639 first wave on the awareness of populations.

## 640 Acknowledgements

641 This work was funded by an ANR grant (SMITID project; ANR-16-CE35-0006).

## Data availability statement

The data that support the findings of this study are openly available<sup>35,37</sup> in the public archive repository Zenodo at <http://doi.org/10.5281/zenodo.4837560> and <http://doi.org/10.5281/zenodo.4769671>. R codes are available in the package ZIprop deposited in R CRAN (<https://cran.r-project.org/package=ZIprop>) and GitLab (<https://gitlab.paca.inrae.fr/meribaud/ziprop>). The package also includes the data sets used in this study.

## Supplementary Materials

Appendices A to D, Figures S1 to S13 and Tables S1 to S6 cited in the manuscript are available in a single Supporting Information PDF file.

## References

- [1] Alamil M, Hughes J, Berthier K, Desbiez C, Thébaud G, Soubeyrand S. Inferring epidemiological links from deep sequencing data: a statistical learning approach for human, animal and plant diseases. *Philosophical Transactions of the Royal Society B* 2019; 374(1775): 20180258.
- [2] Soubeyrand S, Ribaud M, Baudrot V, Allard D, Pommeret D, Roques L. COVID-19 mortality dynamics: The future modelled as a (mixture of) past(s). *PLoS ONE* 2020; 15(9): e0238410.
- [3] Soubeyrand S, Sache I, Lannou C, Chadœuf J. Residual-based specification of the random-effects distribution for cluster data. *Statistical Methodology* 2006; 3: 464-482.
- [4] Soubeyrand S, Chadœuf J. Residual-based specification of a hidden random field included in a hierarchical spatial model. *Computational Statistics & Data Analysis* 2007; 51: 6404–6422.
- [5] Weisberg S. *Applied linear regression*. 528. John Wiley & Sons . 2005.
- [6] Hastie T, Tibshirani R, Friedman J. *The elements of statistical learning: data mining, inference, and prediction*. Springer Science & Business Media . 2009.
- [7] Nelder JA, Wedderburn RW. Generalized linear models. *Journal of the Royal Statistical Society: Series A* 1972; 135(3): 370–384.

- 670 [8] Estabrooks A, Jo T, Japkowicz N. A multiple resampling method for learning from  
671 imbalanced data sets. *Computational Intelligence* 2004; 20(1): 18–36.
- 672 [9] Stasinopoulos DM, Rigby RA, others . Generalized additive models for location scale  
673 and shape (GAMLSS) in R. *Journal of Statistical Software* 2007; 23(7): 1–46.
- 674 [10] Rigby RA, Stasinopoulos DM. Generalized additive models for location, scale and shape.  
675 *Journal of the Royal Statistical Society: Series C* 2005; 54(3): 507–554.
- 676 [11] Aitchison J. The statistical analysis of compositional data. *Journal of the Royal Statis-*  
677 *tical Society: Series B (Methodological)* 1982; 44(2): 139–160.
- 678 [12] Douma JC, Weedon JT. Analysing continuous proportions in ecology and evolution: A  
679 practical introduction to beta and Dirichlet regression. *Methods in Ecology and Evolu-*  
680 *tion* 2019; 10(9): 1412–1430.
- 681 [13] Tang ZZ, Chen G. Zero-inflated generalized Dirichlet multinomial regression model for  
682 microbiome compositional data analysis. *Biostatistics* 2019; 20(4): 698–713.
- 683 [14] Welch WJ. Construction of permutation tests. *Journal of the American Statistical As-*  
684 *sociation* 1990; 85: 693–698.
- 685 [15] Duffy DE, Quiroz AJ. A permutation-based algorithm for block clustering. *Journal of*  
686 *Classification* 1991; 8: 65–91.
- 687 [16] Winkler AM, Webster MA, Vidaurre D, Nichols TE, Smith SM. Multi-level block per-  
688 mutation. *Neuroimage* 2015; 123: 253–268.
- 689 [17] Keele L, Miratrix L. Randomization inference for outcomes with clumping at zero. *The*  
690 *American Statistician* 2019; 73: 141–150.
- 691 [18] Tsagris M, Stewart C. A Dirichlet regression model for compositional data with zeros.  
692 *Lobachevskii Journal of Mathematics* 2018; 39(3): 398–412.
- 693 [19] Breiman L, Friedman J, Stone C, Olshen R. *Classification and Regression Trees*. The  
694 Wadsworth and Brooks-Cole Statistics-Probability Series Taylor & Francis . 1984.
- 695 [20] Hollander M, Wolfe DA, Chicken E. *Nonparametric statistical methods*. 751. John Wiley  
696 & Sons . 2013.
- 697 [21] Kendall MG. The treatment of ties in ranking problems. *Biometrika* 1945; 33: 239–251.
- 698 [22] Pesarin F, Salmaso L. *Permutation tests for complex data: theory, applications and*  
699 *software*. John Wiley & Sons . 2010.

- 700 [23] Spearman C. The Proof and Measurement of Association between Two Things. *The*  
701 *American Journal of Psychology* 1904; 15(1): 72–101.
- 702 [24] Kruskal WH, Wallis WA. Use of ranks in one-criterion variance analysis. *Journal of the*  
703 *American Statistical Association* 1952; 47(260): 583–621.
- 704 [25] Dunn OJ. Multiple comparisons using rank sums. *Technometrics* 1964; 6(3): 241–252.
- 705 [26] Hochberg Y. A sharper Bonferroni procedure for multiple tests of significance.  
706 *Biometrika* 1988; 75: 800–802.
- 707 [27] Hochberg Y, Benjamini Y. More powerful procedures for multiple significance testing.  
708 *Statistics in Medicine* 1990; 9: 811–818.
- 709 [28] Mebane Jr WR, Sekhon JS. Genetic optimization using derivatives: the rgenoud package  
710 for R. *Journal of Statistical Software* 2011; 42(11): 1–26.
- 711 [29] Kloke JD, McKean JW. Rfit: Rank-based estimation for linear models. *The R Journal*  
712 2012; 4(2): 57–64.
- 713 [30] Jaeckel LA. Estimating regression coefficients by minimizing the dispersion of the resid-  
714 uals. *The Annals of Mathematical Statistics* 1972; 43(5): 1449–1458.
- 715 [31] Hughes J, Allen RC, Baguelin M, et al. Transmission of equine influenza virus dur-  
716 ing an outbreak is characterized by frequent mixed infections and loose transmission  
717 bottlenecks. *PLoS Pathogens* 2012; 8(12): e1003081.
- 718 [32] De Maio N, Worby CJ, Wilson DJ, Stoesser N. Bayesian reconstruction of transmission  
719 within outbreaks using genomic variants. *PLoS Computational Biology* 2018; 14(4):  
720 e1006117.
- 721 [33] Alamil M. *Reconstruction of the transmission of a virus during an epidemic by statistical*  
722 *learning on genomic data*. PhD thesis. Aix-Marseille University, Marseille; 2020.
- 723 [34] Bouckaert R, Heled J, Kühnert D, et al. BEAST 2: A software platform for Bayesian  
724 evolutionary analysis. *PLoS Computational Biology* 2014; 10: e1003537.
- 725 [35] Ribaud M, Hughes J. Equine Influenza dataset. *10.5281/zenodo.4837560* 2021. doi:  
726 [10.5281/zenodo.4837560](https://doi.org/10.5281/zenodo.4837560)
- 727 [36] Dong E, Du H, Gardner L. An interactive web-based dashboard to track COVID-  
728 19 in real time. *The Lancet Infectious Diseases* 2020; 20: 533–534. doi:  
729 [https://doi.org/10.1016/S1473-3099\(20\)30120-1](https://doi.org/10.1016/S1473-3099(20)30120-1)



- 730 [37] Ribaud M, Martinetti D, Soubeyrand S. Data for the comparison of COVID-19 mortality  
731 in European and North American geographic entities. *10.5281/zenodo.4769671* 2021.  
732 doi: [10.5281/zenodo.4769671](https://doi.org/10.5281/zenodo.4769671)
- 733 [38] Lahiri S. *Resampling methods for dependent data*. New York: Springer-Verlag . 2003.
- 734 [39] Hallstrom AP. A modified Wilcoxon test for non-negative distributions with a clump of  
735 zeros. *Statistics in Medicine* 2010; 29: 391–400.
- 736 [40] Pimentel RS, Niewiadomska-Bugaj M, Wang JC. Association of zero-inflated continuous  
737 variables. *Statistics & Probability Letters* 2015; 96: 61–67.
- 738 [41] Finos L, Pesarin F. On zero-inflated permutation testing and some related problems.  
739 *Statistical Papers* 2020; 61: 2157–2174.
- 740 [42] Alfons A, Croux C, Filzmoser P. Robust maximum association between data sets: The  
741 R package ccaPP. *Austrian Journal of Statistics* 2016; 45: 71–79.
- 742 [43] Warner B, Misra M. Understanding neural networks as statistical tools. *The American*  
743 *Statistician* 1996; 50: 284–293.
- 744 [44] Friedman JH. Multivariate adaptive regression splines. *The Annals of Statistics* 1991;  
745 19: 1–67.
- 746 [45] Breiman L. Random forests. *Machine Learning* 2001; 45: 5–32.
- 747 [46] Friedman JH. Stochastic gradient boosting. *Computational Statistics & Data Analysis*  
748 2002; 38: 367–378.
- 749 [47] Bühlmann P, Yu B. Boosting with the L<sub>2</sub> loss: regression and classification. *Journal of*  
750 *the American Statistical Association* 2003; 98: 324–339.
- 751 [48] Georgescu V, Soubeyrand S, Kretzschmar A, Laine AL. Exploring spatial and multitype  
752 assemblages of species abundances. *Biometrical Journal* 2009; 51: 979–995.
- 753 [49] Cauchemez S, Bhattarai A, Marchbanks TL, et al. Role of social networks in shaping  
754 disease transmission during a community outbreak of 2009 H1N1 pandemic influenza.  
755 *Proceedings of the National Academy of Sciences* 2011; 108: 2825–2830.
- 756 [50] Jørgensen GHM, Borsheim L, Mejdell CM, Søndergaard E, Bøe KE. Grouping horses  
757 according to gender—effects on aggression, spacing and injuries. *Applied Animal Be-*  
758 *haviour Science* 2009; 120: 94–99.



- 759 [51] Jombart T, Cori A, Didelot X, Cauchemez S, Fraser C, Ferguson N. Bayesian recon-  
760 struction of disease outbreaks by combining epidemiologic and genomic data. *PLoS*  
761 *Computational Biology* 2014; 10: e1003457.
- 762 [52] Lau MS, Marion G, Streftaris G, Gibson G. A systematic Bayesian integration of epi-  
763 demiological and genetic data. *PLoS Computational Biology* 2015; 11: e1004633.
- 764 [53] Lau MS, Gibson GJ, Adrakey H, et al. A mechanistic spatio-temporal framework for  
765 modelling individual-to-individual transmission—With an application to the 2014-2015  
766 West Africa Ebola outbreak. *PLoS Computational Biology* 2017; 13: e1005798.
- 767 [54] Mollentze N, Nel LH, Townsend S, et al. A Bayesian approach for inferring the dy-  
768 namics of partially observed endemic infectious diseases from space-time-genetic data.  
769 *Proceedings of the Royal Society B* 2014; 281: 20133251.
- 770 [55] Morelli MJ, Thébaud G, Chadœuf J, King DP, Haydon DT, Soubeyrand S. A Bayesian  
771 inference framework to reconstruct transmission trees using epidemiological and genetic  
772 data. *PLoS Computation Biology* 2012; 8: e1002768.
- 773 [56] Soubeyrand S. Construction of semi-Markov genetic-space-time SEIR models and infer-  
774 ence. *Journal de la Société Française de Statistique* 2016; 157(1): 129–152.
- 775 [57] Didelot X, Fraser C, Gardy J, Colijn C. Genomic infectious disease epidemiology in  
776 partially sampled and ongoing outbreaks. *Molecular Biology and Evolution* 2017; 34(4):  
777 997–1007.
- 778 [58] Hall M, Woolhouse M, Rambaut A. Epidemic reconstruction in a phylogenetics frame-  
779 work: transmission trees as partitions of the node set. *PLoS Computational Biology*  
780 2015; 11: e1004613.
- 781 [59] Leitner T, Romero-Severson E. Phylogenetic patterns recover known HIV epidemiolog-  
782 ical relationships and reveal common transmission of multiple variants. *Nature Micro-*  
783 *biology* 2018; 3: 983.
- 784 [60] Pybus OG, Suchard MA, Lemey P, et al. Unifying the spatial epidemiology and molec-  
785 ular evolution of emerging epidemics. *Proceedings of the National Academy of Sciences*  
786 2012; 109: 15066–15071.
- 787 [61] Rakotomalala M, Vrancken B, Pinel-Galzi A, et al. Comparing patterns and scales of  
788 plant virus phylogeography: Rice yellow mottle virus in Madagascar and in continental  
789 Africa. *Virus Evolution* 2019; 5: vez023.

- 790 [62] Valdazo-González B, Kim JT, Soubeyrand S, et al. The impact of within-herd genetic  
791 variation upon inferred transmission trees for foot-and-mouth disease virus. *Infection,*  
792 *Genetics and Evolution* 2015; 32: 440–448.
- 793 [63] Picard C, Dallot S, Brunker K, et al. Exploiting genetic information to trace plant virus  
794 dispersal in landscapes. *Annual Review of Phytopathology* 2017; 55: 139–160.
- 795 [64] Dellicour S, Rose R, Pybus OG. Explaining the geographic spread of emerging epi-  
796 demics: a framework for comparing viral phylogenies and environmental landscape data.  
797 *BMC Bioinformatics* 2016; 17: 1–12.