# A. Supplementary Material

## Contents

| A.1. List of Symbols   | 12 |
|--|----|
| A.1.1 Mathematical Symbols   | 12 |
| A.1.2 Symbols Related to Our Mathematical Framework  | 12 |
| A.1.3 Symbols used for Operations on Performances  | 13 |
| A.1.4 Symbols used in the Performance Ordering and Performance-Based Ranking Theory                                | 13 |
| A.1.5 Symbols used for the Particular Case of Two-Class Crisp Classifications                                      | 13 |
| A.2 How to Use our Framework: a Little Catalog of Problems   | 14 |
| A.2.1 Multi-Class Classification (with a Note on Micro- and Macro-Averaging)                                       | 14 |
| A.2.2 Regression (with a Note on the Mean Squared Error and the Mean Absolute Error)                               | 15 |
| A.2.3 Information Retrieval  | 15 |
| A.2.4 Detection (with a Note about the Intersection-over-Union and the F-Score)                                    | 16 |
| A.2.5 Clustering (with a Note about Fowlkes-Mallows Index)   | 17 |
| A.2.6 Ranking (with a Note about Kendall's $\tau$ )  | 17 |
| A.3 Supplementary Material about Sec. 3.2  | 19 |
| A.3.1 Reminders of Classical Definitions.  | 19 |
| A.3.2 The 4 Cases in the Comparison of Two Performances with a Preorder $\leq \dots \dots \dots \dots \dots \dots$ | 19 |
| A.3.3 Implications of the Transitivity of $\lesssim$ .   | 19 |
| A.3.4 Properties of $\sim, >, <, \nleq, \lesssim$ , and $\gtrsim$ .  | 20 |
| A.4. Supplementary Material about Sec. 3.6   | 22 |
| A.4.1 The Visual Inspection of Formulas, to Determine the Importance given by Scores, can be Misleading!           | 22 |
| A.4.2 So, How can we Determine the Importance given by Scores?   | 22 |
| A.5 Supplementary Material about Sec. 4.2  | 23 |
| A.5.1 Proof of Theorem 1   | 23 |
| A.5.2 Proof of Theorem 2   | 23 |
| A.5.3 Proof of Theorem 3   | 24 |
| A.6 Supplementary Material about Sec. 4.3  | 26 |
| A.6.1 All Ranking Scores can be Used to Rank Performances (for $\Phi = \text{conv}$ )                              | 26 |
| A.6.2 On the Properties of Ranking Scores.   | 28 |
| A.7. Supplementary Material about Sec. 5.2   | 31 |
| A.7.1 Link between Classical Formulation and Ours.   | 31 |
| A.7.2 Custom Optimization Algorithm to Estimate Kendall's $\tau$   | 31 |
| A.7.3 Scores Perfectly Correlated with a Ranking Score, for all Performances                                       | 31 |
| A.7.4 Scores Perfectly Correlated with a Ranking Score, for the Performances Corresponding to Given                |    |
| Class Priors $\pi_{-} \neq 0$ and $\pi_{+} \neq 0$   | 31 |

## A.1. List of Symbols

## A.1.1 Mathematical Symbols

- $\mathbf{1}_U$ : the 0-1 indicator function of subset U
- $\mathbb{R}$ : the real numbers
- $\mathcal{R}$ : a relation
- conv: the set of convex combinations
- $\forall$ : the *inclusive disjunction* (*i.e.*, logical or)
- $\wedge$ : the *conjunction* (*i.e.*, logical and)
- o: the composition of functions, *i.e.*  $(q \circ f)(x) = q(f(x))$
- E: the mathematical expectation

## A.1.2 Symbols Related to Our Mathematical Framework

We organize these symbols according to the 6 pillars depicted in Fig. 1, which correspond to the 6 subsections of Sec. 3.

## Symbols related to the 1<sup>st</sup> pillar (Sec. 3.1)

- $\Omega$ : the sample space (universe)
- $\omega$ : a sample (*i.e.*, an element of  $\Omega$ )
- $\Sigma$ : the event space (a  $\sigma$ -algebra on  $\Omega$ , *e.g.*  $2^{\Omega}$ )
- E: an event (*i.e.*, an element of  $\Sigma$ )
- $(\Omega, \Sigma)$ : the measurable space
- $\mathbb{P}_{(\Omega,\Sigma)}$ : all performances on  $(\Omega,\Sigma)$
- $\Pi$ : a set of performances ( $\Pi \subseteq \mathbb{P}_{(\Omega,\Sigma)}$ )
- *P*: a performance (*i.e.*, an element of  $\mathbb{P}_{(\Omega, \Sigma)}$ )

## Symbols related to the 2<sup>nd</sup> pillar (Sec. 3.2)

- ≤: binary relation *worse or equivalent* on P<sub>(Ω,Σ)</sub>
  ≥: binary relation *better or equivalent* on P<sub>(Ω,Σ)</sub>
- ~: binary relation *equivalent* on  $\mathbb{P}_{(\Omega,\Sigma)}$
- >: binary relation *better* on  $\mathbb{P}_{(\Omega,\Sigma)}$
- <: binary relation *worse* on  $\mathbb{P}_{(\Omega,\Sigma)}$
- $\not\equiv$ : binary relation *incomparable* on  $\mathbb{P}_{(\Omega,\Sigma)}$

## Symbols related to the 3<sup>rd</sup> pillar (Sec. 3.3)

- S: the random variable Satisfaction
- $s_{min,\Omega}$ : the minimum satisfaction value
- $s_{max,\Omega}$ : the maximum satisfaction value

## Symbols related to the 4<sup>th</sup> pillar (Sec. 3.4)

- $\mathbb{E}$ : the set of entities to rank
- $\epsilon$ : an entity, *i.e.* an element of  $\mathbb{E}$
- eval: the performance evaluation function
- $\Phi$ : some performances that are for sure achievable

## Symbols related to the 5<sup>th</sup> pillar (Sec. 3.5)

- X: a score
- dom(X): the domain of the score X
- $X_V^E$ : the *expected value score* parameterized by the random variable V•  $X_{E_1|E_2}^P$ : the *probabilistic score* parameterized by the events  $E_1$  and  $E_2$

## Symbols related to the 6<sup>th</sup> pillar (Sec. 3.6)

• *I*: the random variable *Importance* 

## A.1.3 Symbols used for Operations on Performances

• filter<sub>*I*</sub>: the *filtering* operation

#### A.1.4 Symbols used in the Performance Ordering and Performance-Based Ranking Theory

- $\operatorname{rank}_{\mathbb{E}}$ : the *ranking* function, w.r.t. the set of entities  $\mathbb{E}$
- $\leq_X$ : the ordering induced by the score X (cf. Theorem 1)
- $R_I$ : the *ranking score* parameterized by the importance I
- $\tau$ : the rank correlation coefficient of Kendall [13]

#### A.1.5 Symbols used for the Particular Case of Two-Class Crisp Classifications

#### Particularization of the mathematical framework

- *tn*: the sample *true negative*
- *fp*: the sample *false positive*, *a.k.a.* type I error
- *fn*: the sample *false negative*, *a.k.a.* type II error
- *tp*: the sample *true positive*

#### Extensions to the mathematical framework

- ROC: the *Receiver Operating Characteristic* space, *i.e.*  $FPR \times TPR$
- PR: the *Precision-Recall* space, *i.e.*  $TPR \times PPV$
- *Y*: the random variable for the ground truth
- $\hat{Y}$ : the random variable for the prediction
- $\mathbb{C}$ : the set of classes
- c: a class (*i.e.*, an element of  $\mathbb{C}$ )
- *c*\_: the negative class
- $c_+$ : the positive class

#### Scores

- *A*: the *accuracy*
- *TNR*: the *true negative rate*
- *FPR*: the *false positive rate*
- *TPR*: the *true positive rate*
- *PPV*: the *positive predictive value*
- $F_{\beta}$ : the F-scores
- $\pi_+$ : the prior of the positive class
- $\pi_{-}$ : the prior of the negative class

#### A.2. How to Use our Framework: a Little Catalog of Problems

Throughout the paper, we have exemplified our theory with the problem of two-class classification. This section aims at showing the universality of our theory. It presents a little catalog of other problems, together with discussions on how to use our framework for them. These discussions are introductions. As shown by Sec. 5 and two recent works [12, 21], an in-depth analysis and particularization of our theory to the various problems (*e.g.*, to highlight their distinctive features, to review the current ranking practices from the literature and the consistency of popular scores, and to establish practical tools tailored to different user needs) may require substantial work that is out of the scope of this supplementary material.

In the following, we adopt a systematic approach: for each problem, we start by specifying a thought random experiment for the evaluation (this is an arbitrary choice since the random experiment is not unique for each problem; do not hesitate to use a different random experiment, the important is to specify it explicitly!), then we discuss the possible choices for the sample space  $\Omega$  (and what the set of all performances is), for the modeling of tasks with the satisfaction S, for the modeling of the knowledge we have about the evaluation with the function  $\Phi$ , and for the modeling of the application-specific preferences with the importance I.

#### A.2.1 Multi-Class Classification (with a Note on Micro- and Macro-Averaging)

Let us consider the following thought experiment to evaluate classifiers predicting classes in a finite and non-empty set C.

**Random Experiment 1.** (1) Draw a sample  $s \in \mathbb{S}$  at random from a given source S. (2) Apply the oracle  $\mathcal{O}$  on s to obtain the ground-truth class  $y(s) \in \mathbb{C}$ . (3) Apply a descriptor  $\mathcal{D}$  on s to obtain the features (a.k.a. attributes)  $x(s) \in \mathbb{X}$ . (4) Feed the classifier  $\mathcal{C}$  with x(s) to obtain the predicted class  $\hat{y}(x(s)) \in \mathbb{C}$ . (5) Set the outcome of the experiment to the pair  $(y, \hat{y}) \in \mathbb{C}^2$ .

- **Choice for**  $\Omega$  and  $\mathbb{P}_{(\Omega,\Sigma)}$ . Our theory applies with  $\Omega = \mathbb{C}^2$ . By definition, the function eval gives, for any classifier  $\mathcal{C} : \mathbb{X} \to \mathbb{C}$  (the evaluated entity, either deterministic or not), the distribution of outcomes resulting from this random experiment:  $P_{\mathcal{C}} = \text{eval}(\mathcal{C})$ . Note that it is implicit that the performances are specific for some given source  $\mathcal{S}$  (*e.g.*, evaluation dataset), oracle  $\mathcal{O}$ , and descriptor  $\mathcal{D}$ . By convenience, one can manipulate the ground-truth and predicted classes with, respectively, the random variables Y and  $\hat{Y}$  defined in such a way that  $\omega = (Y(\omega), \hat{Y}(\omega)) \forall \omega \in \Omega$ .
- Choice for S. Several classification tasks can be distinguished, as the following two examples show. (1) One can consider that all erroneous classifications are unsatisfactory and that correct classifications are satisfactory. For this task, the satisfaction is then binary and given by  $S = \mathbf{1}_{Y=\hat{Y}}$ . The expected value of the satisfaction, which is a particular case of ranking scores, is then equal to the multi-class accuracy. (2) One can also consider the similarity sim :  $\mathbb{C}^2 \to \mathbb{R}$  between classes, and choose the satisfaction accordingly:  $S(\omega) = \sin(Y(\omega); \hat{Y}(\omega))$ . A wide variety of tasks can be considered by tuning sim. In general, the expected value of the satisfaction is different from the multi-class accuracy.
- **Choice for**  $\Phi$ . As the classifier is used once and only once during the execution of the evaluation, we know that if the performances  $P_1$  and  $P_2$  are achievable by some classifiers  $C_1$  and  $C_2$ , then any performance  $\overline{P} = \lambda_1 P_1 + \lambda_2 P_2$  (with  $\lambda_1 \ge 0, \lambda_2 \ge 0$ , and  $\lambda_1 + \lambda_2 = 1$ ) is achievable by a classifier  $\overline{C}$  obtained by a non-deterministic combination of  $C_1$  and  $C_2$  that chooses them with respective probabilities  $\lambda_1$  and  $\lambda_2$ . Thus,  $\Phi = \text{conv}$  makes sense, and all ranking scores can be used to rank classifiers. However, it would be possible to go further, by considering other functions  $\Phi$  that would include the knowledge that we can predict the performance achievable by composing the classifier with any of the  $|\mathbb{C}|^{|\mathbb{C}|}$  functions  $f : \mathbb{C} \to \mathbb{C}$ . This would lead to other performance orderings suitable for ranking classifiers.
- **Choice for** *I*. When  $\Phi = \text{conv}$ , we have demonstrated that all rankings induced by the ranking scores  $R_I$  satisfy all our three axioms. This leaves a great flexibility for the users to fine-tune the ranking w.r.t. their application-specific preferences through the random variable *I*. As we have seen, the only constraints are that  $I \neq 0$  and  $I(\omega) \ge 0 \forall \omega \in \Omega$ .
- Note on micro- and macro-averaging. Micro- and macro-averaging are commonly used techniques to build scores for multi-class classification from scores for two-class classification [22]. We warn that they have pitfalls. In general, micro- and macro-averaging scores suitable for ranking two-class classifiers do not lead to scores suitable for ranking multi-class classifiers. The accuracy put aside, the performance orderings induced from the micro-averaged versions of the scores put in green in Tab. 1 are incompatible with  $S = \mathbf{1}_{Y=\hat{Y}}$ : our 2<sup>nd</sup> axiom is not satisfied. Moreover, the accuracy put again aside, the performance orderings induced from the macro-averaged versions of the scores put in green in Tab. 1 are incompatible with  $\Phi = \text{conv}$ : our 3<sup>rd</sup> axiom is not satisfied. A solution consists in using directly ranking scores defined for multi-class classification.

#### A.2.2 Regression (with a Note on the Mean Squared Error and the Mean Absolute Error)

Let us consider the following thought experiment to evaluate regressors.

**Random Experiment 2.** (1) Draw a sample  $s \in \mathbb{S}$  at random from a given source S. (2) Apply the oracle  $\mathcal{O}$  on s to obtain the ground-truth value  $y(s) \in \mathbb{R}$ . (3) Apply a descriptor  $\mathcal{D}$  on s to obtain the features (a.k.a. attributes)  $x(s) \in \mathbb{X}$ . (4) Feed the regressor  $\mathcal{R}$  with x(s) to obtain the predicted value  $\hat{y}(x(s)) \in \mathbb{R}$ . (5) Set the outcome of the experiment to the pair  $(y, \hat{y}) \in \mathbb{R}^2$ .

- **Choice for**  $\Omega$  and  $\mathbb{P}_{(\Omega,\Sigma)}$ . Our theory applies with  $\Omega = \mathbb{R}^2$ . By definition, the function eval gives, for any regressor  $\mathcal{R} : \mathbb{X} \to \mathbb{R}$  (the evaluated entity, either deterministic or not), the distribution of outcomes resulting from this random experiment. It is the performance  $P_{\mathcal{R}} = \operatorname{eval}(\mathcal{R})$  of  $\mathcal{R}$ . By convenience, one can manipulate the ground-truth and predicted values with, respectively, the random variables Y and  $\hat{Y}$  defined in such a way that  $\omega = (Y(\omega), \hat{Y}(\omega)) \forall \omega \in \Omega$ .
- Choice for S. Clearly, it is not a good idea to choose  $S = \mathbf{1}_{Y=\hat{Y}}$ , similarly as one can do in classification. In practice, a regressor as no chance to predict the same value as the oracle, so this unfortunate choice for S would lead to performances P such that P(S = 0) = 1. In other words, all performances observed about real regressors would belong to the set of the worst performances (see Corollary 2), and their ranking would be of little interest (see Corollary 1). A better option consists in specifying a tolerance  $\epsilon > 0$  and choosing  $S = \mathbf{1}_{|Y \hat{Y}| \leq \epsilon}$ . An even more flexible option, which takes advantage of the fact that satisfaction values do not necessarily have to be positive, is to choose  $S = f(|Y \hat{Y}|)$  with any arbitrarily chosen monotonically decreasing function f. The plethora of choices that can be made for S makes it clear that there is an infinity of tasks related to the regression problem.
- **Choice for**  $\Phi$ . As the regressor is used once and only once during the execution of the evaluation, we know that if the performances  $P_1$  and  $P_2$  are achievable by some regressors  $\mathcal{R}_1$  and  $\mathcal{R}_2$ , then any performance  $\overline{P} = \lambda_1 P_1 + \lambda_2 P_2$  (with  $\lambda_1 \ge 0, \lambda_2 \ge 0$ , and  $\lambda_1 + \lambda_2 = 1$ ) is achievable by a regressor  $\overline{\mathcal{R}}$  obtained by a non-deterministic combination of  $\mathcal{R}_1$  and  $\mathcal{R}_2$  that chooses them with respective probabilities  $\lambda_1$  and  $\lambda_2$ . Thus,  $\Phi = \text{conv}$  makes sense, and all ranking scores can be used to rank regressors. However, it would be possible to go further, by considering other functions  $\Phi$  that would include the knowledge that we can predict the performance achievable by adding noise, or applying a transformation on the output of the regressor. This would lead to other performance orderings suitable for ranking regressors.
- **Choice for** *I*. When  $\Phi = \text{conv}$ , we have demonstrated that all rankings induced by the ranking scores  $R_I$  satisfy all our three axioms. This leaves a great flexibility for the users to fine-tune the ranking w.r.t. their application-specific preferences through the random variable *I*. As we have seen, the only constraints are that  $I \neq 0$  and  $I(\omega) \ge 0 \forall \omega \in \Omega$ .
- Note on the mean squared error and the mean absolute error. If we choose  $S = -|Y \hat{Y}|^2$ , the ranking score  $R_I$  corresponding to uniform importance values, *i.e.*, the expected satisfaction  $X_S^E$ , yields a ranking that minimizes the *mean squared error* (MSE). If we choose  $S = -|Y \hat{Y}|$ , the ranking score  $R_I$  corresponding to uniform importance values, *i.e.*, the expected satisfaction  $X_S^E$ , yields a ranking that minimizes the *mean squared error* (MSE). If we choose  $S = -|Y \hat{Y}|$ , the ranking score  $R_I$  corresponding to uniform importance values, *i.e.*, the expected satisfaction  $X_S^E$ , yields a ranking that minimizes the *mean absolute error* (MAE).

#### A.2.3 Information Retrieval

Let us consider the following thought experiment to evaluate information retrieval systems. We denote by  $\mathbb{Q}$  the set of all possible queries.

**Random Experiment 3.** (1) Draw a query  $q \in \mathbb{Q}$  at random from a given source S. (2) Apply the oracle  $\mathcal{O}$  on q to obtain the ground-truth set of results  $\mathbb{Y}$ . (3) Apply the evaluated information retrieval system S on q to obtain the predicted set of results  $\hat{\mathbb{Y}}$ . (4) If  $\mathbb{Y} = \emptyset$  and  $\hat{\mathbb{Y}} = \emptyset$ , restart the experiment, otherwise draw a result r at random in  $\mathbb{Y} \cup \hat{\mathbb{Y}}$ . (5) Choose the outcome as follows: fp if  $r \notin \mathbb{Y}$  and  $r \in \hat{\mathbb{Y}}$ , fn if  $r \in \mathbb{Y}$  and  $r \notin \hat{\mathbb{Y}}$ , and tp if  $r \in \mathbb{Y}$  and  $r \in \hat{\mathbb{Y}}$ .

- Choice for  $\Omega$  and  $\mathbb{P}_{(\Omega,\Sigma)}$ . Our theory applies with  $\Omega = \{fp, fn, tp\}$ . By definition, the function eval gives, for any retrieval system S defined on  $\mathbb{Q}$  (the evaluated entity, either deterministic or not), the distribution of outcomes resulting from this random experiment:  $P_S = \text{eval}(S)$ .
- **Choice for** S. Intuitively, everyone certainly agrees that S(fp) < S(tp) and S(fn) < S(tp). But we expect different opinions regarding whether the outcome (sample) fp gives less, equal, or more satisfaction than fn.

- Choice for  $\Phi$ . This random experiment is very interesting as, during its execution, the evaluated entity (the information retrieval system S) can be used multiple times. In such a case, we have to discuss whether  $\Phi = \text{conv}$  is adequate. Let us consider two systems  $S_1$ ,  $S_2$  and their respective performances  $P_1 = \text{eval}(S_1)$ ,  $P_2 = \text{eval}(S_2)$ . It is possible to show that the performance of a retrieval system  $\overline{S}$  obtained by a non-deterministic combination of  $S_1$  and  $S_2$ , that chooses them with respective probabilities  $\lambda_1$  and  $\lambda_2$ , is some interpolated performance  $\overline{P} = \mu_1 P_1 + \mu_2 P_2$  (with  $\mu_1 \ge 0$ ,  $\mu_2 \ge 0$ , and  $\mu_1 + \mu_2 = 1$ ). Unless being in very particular cases,  $\mu \ne \lambda$ . In other words, we know that the performances that are convex combinations of achievable performances are also achievable (for any  $\lambda$ , there exists  $\mu$ ), but we do not know in general how to achieve them (for most  $\mu$  it is not possible to determine  $\lambda$ ). This contrasts with the other kinds of problems discussed in this catalog. In fact, the question we raise here is not specific to the information retrieval problem: it is peculiar to the random experiment that we have chosen for it. By slightly modifying the thought experiment, the question vanishes: instead of restarting the experiment when  $\mathbb{Y} \cup \hat{\mathbb{Y}} = \emptyset$ , one could yield a fourth outcome (and add it to the sample space  $\Omega$ ). By doing so, the evaluated entity S is used only once and  $\Phi = \text{conv}$  makes sense for sure.
- **Choice for** *I*. If  $\Phi = \text{conv}$  is considered as adequate, then we have demonstrated that all rankings induced by the ranking scores  $R_I$  satisfy all our three axioms. This leaves a great flexibility for the users to fine-tune the ranking w.r.t. their application-specific preferences through the random variable *I*. As we have seen, the only constraints are that  $I \neq 0$  and  $I(\omega) \geq 0 \forall \omega \in \Omega$ .

#### A.2.4 Detection (with a Note about the Intersection-over-Union and the F-Score)

Different types of detections are present in the literature. An example of spatial detection aims at predicting the axis-aligned bounding boxes around all the objects that match some given properties (*i.e.*, a semantic class) in input images. Examples of temporal detections include the detection of events in video streams and in audio recordings. By definition, such detection problems are called *action spotting* when the temporal window is small, and *activity detection* otherwise. Let us consider the following, generic, thought experiment to evaluate detectors.

**Random Experiment 4.** (1) Draw an input at random from a given source S (e.g., dataset). (2) Apply the oracle O on it to obtain a set  $\hat{\mathbb{Y}}$  of ground-truth detections. (3) Also apply the detector D on it to obtain a set  $\hat{\mathbb{Y}}$  of predicted detections. (4) If  $\hat{\mathbb{Y}} = \emptyset$  and  $\hat{\mathbb{Y}} = \emptyset$ , then end the experiment with the outcome  $\oplus$ . Otherwise: (5) Apply a matching criterion between  $\mathbb{Y}$  and  $\hat{\mathbb{Y}}$  such that to any detection in  $\mathbb{Y}$  should be associated at most a detection in  $\hat{\mathbb{Y}}$  and vice versa. (6) Draw a detection d at random in  $\mathbb{Y} \cup \hat{\mathbb{Y}}$ . (7) Give as outcome fp, fn or tp depending on whether d is a prediction, a ground truth, or both (i.e., a match).

- Choice for  $\Omega$  and  $\mathbb{P}_{(\Omega,\Sigma)}$ . Our theory applies with  $\Omega = \{ \otimes, fp, fn, tp \}$ . By definition, the function eval gives, for any detector (the evaluated entity), the distribution of outcomes resulting from this random experiment. It is the performance  $P_{\mathcal{D}} = \operatorname{eval}(\mathcal{D})$  of  $\mathcal{D}$ .
- **Choice for** S. Intuitively, everyone certainly agrees that  $\otimes$  and tp give entire satisfaction. Moreover, we expect agreement on S(fp) < S(tp) and S(fn) < S(tp). But we expect different opinions regarding whether fp gives less, equal, or more satisfaction than fn.
- **Choice for**  $\Phi$ . As the detector is used once and only once during the execution of the evaluation, we know that if the performances  $P_1$  and  $P_2$  are achievable by some detectors  $D_1$  and  $D_2$ , then any performance  $\overline{P} = \lambda_1 P_1 + \lambda_2 P_2$  (with  $\lambda_1 \ge 0$ ,  $\lambda_2 \ge 0$ , and  $\lambda_1 + \lambda_2 = 1$ ) is achievable by a detector  $\overline{D}$  obtained by a non-deterministic combination of  $D_1$  and  $D_2$  that chooses them with respective probabilities  $\lambda_1$  and  $\lambda_2$ . Thus,  $\Phi = \text{conv}$  makes sense, and all ranking scores can be used to rank detectors.
- **Choice for** *I*. If  $\Phi = \text{conv}$  is considered as adequate, then we have demonstrated that all rankings induced by the ranking scores  $R_I$  satisfy all our three axioms. This leaves a great flexibility for the users to fine-tune the ranking w.r.t. their application-specific preferences through the random variable *I*. As we have seen, the only constraints are that  $I \neq 0$  and  $I(\omega) \ge 0 \forall \omega \in \Omega$ .
- Note about the Intersection-over-Union and the F-score. Traditionally, in the literature, as soon as one has symbols fp, fn, and tp, regardless of their very fine meaning, one defines quantities  $IoU = \frac{P(tp)}{P(fp)+P(fn)+P(tp)}$  and  $F_1 = \frac{2P(tp)}{P(fp)+P(fn)+2P(tp)}$ , and name them *Intersection-over-Union* (or *Jaccard*) and *F-one*, respectively. The exact meaning

of these quantities is not well standardized. In particular, the random experiment supporting the evaluation, if it exists, is rarely specified explicitly. For this reason, we cannot give the guarantee that these quantities are suitable to rank detectors in all works in which they have been used. However, with the random experiment given here-above, with  $\Phi = \text{conv}$ , and with  $S = \mathbf{1}_{\{\oplus, tp\}}$ , we can guarantee that IoU and  $F_1$  are suitable to rank detectors because they are equal to the ranking scores with, respectively,  $I = \mathbf{1}_{\{fp, fn, tp\}}$  and  $I = \mathbf{1}_{\{fp, tp\}} + \mathbf{1}_{\{fn, tp\}}$ . Thus, the performance orderings induced by them fulfill our three axioms.

#### A.2.5 Clustering (with a Note about Fowlkes-Mallows Index)

Let us consider the following thought experiment to evaluate clustering methods. These methods aim to place in different clusters (groups) dissimilar objects and in the same cluster (group) objects that are similar to each other. We denote by  $\mathbb{E}$  the set of elements that these methods have to deal with. For the sake of simplicity, we do not consider hierarchical clustering.

**Random Experiment 5.** (1) Apply both the clustering method C and the oracle O on  $\mathbb{E}$  to obtain, respectively, the predicted and ground-truth clusterings. (2) Randomly draw two distinct elements,  $\epsilon_1$  and  $\epsilon_2$ , from  $\mathbb{E}$ . (3) Consider that the pair ( $\epsilon_1$ ;  $\epsilon_2$ ) is a negative or a positive, in a given clustering, when  $\epsilon_1$  and  $\epsilon_2$  are in different clusters or in the same cluster, respectively. (4) Choose the outcome as follows: tn when ( $\epsilon_1$ ;  $\epsilon_2$ ) is negative in both the predicted and ground-truth clusterings, fp when ( $\epsilon_1$ ;  $\epsilon_2$ ) is negative in the ground-truth clustering and positive in the predicted clustering, fn when ( $\epsilon_1$ ;  $\epsilon_2$ ) is positive in the ground-truth clustering and negative in the predicted clustering, and tp when ( $\epsilon_1$ ;  $\epsilon_2$ ) is positive in both the predicted and ground-truth clusterings.

- Choice for  $\Omega$  and  $\mathbb{P}_{(\Omega,\Sigma)}$ . Our theory applies with  $\Omega = \{tn, fp, fn, tp\}$ . By definition, the function eval gives, for any clustering method (the evaluated entity), the distribution of outcomes resulting from this random experiment. It is the performance  $P_{\mathcal{C}} = \operatorname{eval}(\mathcal{C})$  of  $\mathcal{C}$ .
- **Choice for** S. When S is chosen such that S(fp) = S(fn) = 0 and S(tn) = S(tp) = 1, we are in the same setting as the one we studied for the two-class classification in Sec. 5. This is indeed not because we use the same symbols for the elements of  $\Omega$  —this is just a convention—, but because, in both settings, we have  $|\Omega| = 4$ , |S = 0| = 2, and |S = 1| = 2. This implies that the performance orderings that satisfy our three axioms for ranking two-class classifiers can also be used for ranking clustering methods, and vice versa.
- **Choice for**  $\Phi$ . As the clustering method is used once and only once during the execution of the evaluation, we know that if the performances  $P_1$  and  $P_2$  are achievable by some clustering methods  $C_1$  and  $C_2$ , then any performance  $\overline{P} = \lambda_1 P_1 + \lambda_2 P_2$  (with  $\lambda_1 \ge 0$ ,  $\lambda_2 \ge 0$ , and  $\lambda_1 + \lambda_2 = 1$ ) is achievable by a clustering method  $\overline{C}$  obtained by a non-deterministic combination of  $C_1$  and  $C_2$  that chooses them with respective probabilities  $\lambda_1$  and  $\lambda_2$ . Thus,  $\Phi = \text{conv}$  makes sense, and all ranking scores can be used to rank clustering methods.
- **Choice for** *I*. When  $\Phi = \text{conv}$ , we have demonstrated that all rankings induced by the ranking scores  $R_I$  satisfy all our three axioms. This leaves a great flexibility for the users to fine-tune the ranking w.r.t. their application-specific preferences through the random variable *I*. As we have seen, the only constraints are that  $I \neq 0$  and  $I(\omega) \ge 0 \forall \omega \in \Omega$ .
- Note about Fowlkes-Mallows index. The score FMI known as Fowlkes-Mallows index [9] and cosine coefficient [1], which is commonly used for clustering methods and defined as the geometric mean of the positive predictive value PPV (a.k.a. precision) and true positive rate TPR (a.k.a. sensitivity and recall), does not satisfy our 3<sup>rd</sup> axiom: the performance ordering induced by FMI is incompatible with  $\Phi = \text{conv}$ . More precisely, the clustering method C obtained by randomly choosing between some methods  $C_1$  or  $C_2$  can be such that  $FMI(\text{eval}(C)) < \min(FMI(\text{eval}(C_1)), FMI(\text{eval}(C_2)))$ , while it makes no sense to say that C can be worse than  $C_1$  or  $C_2$ . From this perspective, it is advisable to use any ranking score instead of FMI, for example, those in green in Tab. 1.

#### A.2.6 Ranking (with a Note about Kendall's $\tau$ )

Let us consider the following thought experiment to evaluate ranking methods. We denote by  $\mathbb{E}$  the set of elements that these methods have to rank. For the sake of simplicity, we prefer to deal only with the case with no tie hereafter.

**Random Experiment 6.** (1) Apply both the ranking method  $\mathcal{R}$  and the oracle  $\mathcal{O}$  on  $\mathbb{E}$  to obtain, respectively, the predicted and ground-truth sequences of elements. (2) Randomly draw two distinct elements,  $\epsilon_1$  and  $\epsilon_2$ , from  $\mathbb{E}$ . (3) Four cases can

occur depending on whether  $\epsilon_1$  is before or after  $\epsilon_2$  in the predicted sequence and whether  $\epsilon_1$  is before or after  $\epsilon_2$  in the ground-truth sequence. Nevertheless, two outcomes are enough: choose  $\odot$  if  $\epsilon_1$  and  $\epsilon_2$  appear in the same order in both sequences,  $\odot$  otherwise.

**Choice for**  $\Omega$  and  $\mathbb{P}_{(\Omega,\Sigma)}$ . Our theory applies with  $\Omega = \{ \odot, \odot \}$ . By definition, the function eval gives, for any ranking method (the evaluated entity), the distribution of outcomes resulting from this random experiment. It is the performance  $P_{\mathcal{R}} = \operatorname{eval}(\mathcal{R})$  of  $\mathcal{R}$ . The probability of drawing a *discordant pair* is given by  $P(\{\odot\})$ , and the probability of drawing a *concordant pair* is given by  $P(\{\odot\})$ .

**Choice for** S. Clearly, S(o) < S(o) is wanted.

- **Choice for**  $\Phi$ . As the ranking method is used once and only once during the execution of the evaluation, we know that if the performances  $P_1$  and  $P_2$  are achievable by some ranking methods  $\mathcal{R}_1$  and  $\mathcal{R}_2$ , then any performance  $\overline{P} = \lambda_1 P_1 + \lambda_2 P_2$  (with  $\lambda_1 \ge 0$ ,  $\lambda_2 \ge 0$ , and  $\lambda_1 + \lambda_2 = 1$ ) is achievable by a ranking method  $\overline{\mathcal{R}}$  obtained by a non-deterministic combination of  $\mathcal{R}_1$  and  $\mathcal{R}_2$  that chooses them with respective probabilities  $\lambda_1$  and  $\lambda_2$ . Thus,  $\Phi = \text{conv}$  makes sense, and all ranking scores can be used to rank ranking methods.
- **Choice for** *I*. Because  $|\Omega| = 2$  and  $S(\odot) \neq S(\odot)$ , we are in a particular case in which all ranking scores rank the ranking methods in the same way. From this point of view, fine-tuning *I* is useless.
- Note about Kendall's  $\tau$ . When  $S(\odot) = -1$  and  $S(\odot) = 1$ , the expected value of the satisfaction is given by  $X_S^E(P) = 1 2P(\{\odot\}) = P(\{\odot\}) P(\{\odot\}) = \tau(P)$ . In other words, with the task corresponding to this choice for the satisfaction, Kendall's correlation coefficient  $\tau$  [13] is the ranking score corresponding to uniform importance values.

#### A.3. Supplementary Material about Sec. 3.2

This section is devoted to reminders about the order theory.

#### A.3.1 Reminders of Classical Definitions.

Let  $\mathcal{R}$  be a homogeneous binary relation on  $\mathbb{P}_{(\Omega,\Sigma)}$ . It is said:

- reflexive iif  $PRP \forall P$ ;
- *irreflexive iif*  $\nexists P : P\mathcal{R}P$ ;
- transitive iif  $P_1 \mathcal{R} P_2 \wedge P_2 \mathcal{R} P_3 \Rightarrow P_1 \mathcal{R} P_3 \forall P_1, P_2, P_3;$
- symmetric iif  $P_1 \mathcal{R} P_2 \Leftrightarrow P_2 \mathcal{R} P_1 \forall P_1, P_2;$
- asymmetric iif  $\nexists(P_1, P_2) : P_1 \mathcal{R} P_2 \wedge P_2 \mathcal{R} P_1;$
- and antisymmetric iif  $P_1 \mathcal{R} P_2 \wedge P_2 \mathcal{R} P_1 \Rightarrow P_1 = P_2$ .

Two homogeneous binary relations  $\mathcal{R}_a$  and  $\mathcal{R}_b$  on  $\mathbb{P}_{(\Omega,\Sigma)}$  are said *converse iif*  $P_1\mathcal{R}_aP_2 \Leftrightarrow P_2\mathcal{R}_bP_1 \forall P_1, P_2$ . A relation  $\mathcal{R}$  is:

- an *equivalence iif* it is reflexive, transitive, and symmetric;
- a preorder iif it is reflexive and transitive;
- and an *order iif* it is reflexive, transitive, and antisymmetric.

An order  $\mathcal{R}$  is said *total iif*  $\nexists(P_1, P_2) : P_1 \not \mathcal{R}P_2 \land P_2 \not \mathcal{R}P_1$ . It is said *partial* otherwise.

#### A.3.2 The 4 Cases in the Comparison of Two Performances with a Preorder $\leq$ .

Let us now consider a preorder  $\leq$  and derive the homogeneous binary relations  $\sim, >, <, \notin$  as follows:

$$P_1 \sim P_2 \Leftrightarrow P_1 \lesssim P_2 \land P_2 \lesssim P_1 \tag{1}$$

$$P_1 > P_2 \Leftrightarrow P_1 \nleq P_2 \land P_2 \lesssim P_1 \tag{2}$$

$$P_1 < P_2 \Leftrightarrow P_1 \lesssim P_2 \land P_2 \nleq P_1 \tag{3}$$

$$P_1 \nleq P_2 \Leftrightarrow P_1 \nleq P_2 \land P_2 \nleq P_1 . \tag{4}$$

Indeed, we have:

$$P_1 \lesssim P_2 \Leftrightarrow P_1 < P_2 \lor P_1 \sim P_2 \,. \tag{5}$$

Similarly, one can derive other binary relations taking unions of  $\sim$ , >, <, or  $\nleq$ . For example,

$$P_1 \gtrsim P_2 \Leftrightarrow P_1 > P_2 \lor P_1 \sim P_2 \,. \tag{6}$$

#### A.3.3 Implications of the Transitivity of $\leq$ .

We can easily check, for each  $\mathcal{R}_{ab} \in \{\leq, \not\leq\}$ , each  $\mathcal{R}_{ba} \in \{\leq, \not\leq\}$ , each  $\mathcal{R}_{bc} \in \{\leq, \not\leq\}$ , each  $\mathcal{R}_{cb} \in \{\leq, \not\leq\}$ , each  $\mathcal{R}_{ca} \in \{\leq, \not\leq\}$ , and each  $\mathcal{R}_{ac} \in \{\leq, \not\leq\}$ , if there exists  $(P_a, P_b, P_c)$  such that  $P_a \mathcal{R}_{ab} P_b$ ,  $P_b \mathcal{R}_{ba} P_a$ ,  $P_b \mathcal{R}_{bc} P_c$ ,  $P_c \mathcal{R}_{cb} P_b$ ,  $P_c \mathcal{R}_{ca} P_a$ , and  $P_a \mathcal{R}_{ac} P_c$ . Because of the assumed transitivity of  $\leq$ , there are only 29 possible cases out of the 2<sup>6</sup>:

| 1. $P_a \nleq P_b, P_b \nleq P_c, P_a \nleq P_c$   | 8. $P_a \nleq P_b, P_b < P_c, P_a < P_c$   |
|--|--|
| 2. $P_a \nleq P_b, P_b \nleq P_c, P_a > P_c$       | 9. $P_a \nleq P_b, P_b \sim P_c, P_a \nleq P_c$  |
| 3. $P_a \nleq P_b, P_b \nleq P_c, P_a < P_c$       | 10. $P_a > P_b, P_b \nleq P_c, P_a \nleq P_c$  |
| 4. $P_a \nleq P_b, P_b \nleq P_c, P_a \sim P_c$    | 11. $P_a > P_b, P_b \nleq P_c, P_a > P_c$  |
| 5. $P_a \notin P_b, P_b > P_c, P_a \notin P_c$     | 12. $P_a > P_b, P_b > P_c, P_a > P_c$  |
| 6. $P_a \notin P_b, P_b > P_c, P_a > P_c$          | 13. $P_a > P_b, P_b < P_c, P_a \notin P_c$<br>14. $P_a > P_b, P_b < P_c, P_a \notin P_c$ |
| 7. $P_a \not\leq P_b, P_b < P_c, P_a \not\leq P_c$ | 15. $P_a > P_b, P_b < P_c, P_a < P_c$  |

| 16. $P_a > P_b, P_b < P_c, P_a \sim P_c$      | 23. $P_a < P_b, P_b > P_c, P_a \sim P_c$   |
|---|--|
| 17. $P_a > P_b, P_b \sim P_c, P_a > P_c$      | 24. $P_a < P_b, P_b < P_c, P_a < P_c$  |
| 18. $P_a < P_b, P_b \nleq P_c, P_a \nleq P_c$ | 25. $P_a < P_b, P_b \sim P_c, P_a < P_c$   |
| 19. $P_a < P_b, P_b \notin P_c, P_a < P_c$    | 26. $P_a \sim P_b, P_b \notin P_c, P_a \notin P_c$                                       |
| 20. $P_a < P_b, P_b > P_c, P_a \nleq P_c$     | 21. $P_a \sim P_b, P_b > P_c, P_a > P_c$<br>28. $P \sim P_b, P_b < P_c > P_c$            |
| 21. $P_a < P_b, P_b > P_c, P_a > P_c$         | $20. I_a \sim I_b, I_b \sim I_c, I_a \sim I_c$ $29 P_c \sim P_c P_c \sim P_c P \sim P_c$ |
| 22. $P_a < P_b, P_b > P_c, P_a < P_c$         | $2$ , $1_a$ , $1_b$ , $1_b$ , $1_c$ , $1_a$ , $1_c$                                      |

From this list, we can derive some rules for manipulating the binary relations  $\sim$ , >, <, and  $\nleq$ . First, we can see that  $\sim$ , >, and < are transitive:

$$P_1 \sim P_2 \wedge P_2 \sim P_3 \Rightarrow P_1 \sim P_3 \tag{7}$$

$$P_1 > P_2 \land P_2 > P_3 \Rightarrow P_1 > P_3 \tag{8}$$

$$P_1 < P_2 \land P_2 < P_3 \Rightarrow P_1 < P_3.$$

$$\tag{9}$$

Second, we can also see how  $\sim$  can be combined with the other 3 relations:

$$(P_1 \sim P_2 \land P_2 > P_3) \lor (P_1 > P_2 \land P_2 \sim P_3) \Rightarrow P_1 > P_3$$

$$(10)$$

$$P_1 \sim P_2 \wedge P_2 < P_3) \lor (P_1 < P_2 \wedge P_2 \sim P_3) \Rightarrow P_1 < P_3 \tag{11}$$

$$(P_1 \sim P_2 \land P_2 \nleq P_3) \lor (P_1 \gneqq P_2 \land P_2 \sim P_3) \Rightarrow P_1 \gneqq P_3.$$

$$(12)$$

And, third, we can see how  $\nleq$  can be combined with > and <:

$$(P_1 \nleq P_2 \land P_2 > P_3) \lor (P_1 > P_2 \land P_2 \nleq P_3) \Rightarrow P_1 > P_3 \lor P_1 \gneqq P_3$$

$$(13)$$

$$(P_1 \nleq P_2 \land P_2 < P_3) \lor (P_1 < P_2 \land P_2 \nleq P_3) \Rightarrow P_1 < P_3 \lor P_1 \nleq P_3 .$$

$$(14)$$

A.3.4 Properties of  $\sim$ , >, <,  $\nleq$ ,  $\lesssim$ , and  $\gtrsim$ .

**Lemma 1.** When  $\leq$  is a preorder,  $\sim$  is reflexive.

| Proof. | This results from the reflexivity of | $\lesssim$ and from Eq. (1 | ): $P \sim P \Leftrightarrow P \leq$ | $\leq P \land P \lesssim P \Leftrightarrow true.$ |  |
|--------|--------------------------------------|----------------------------|--------------------------------------|---|--|
|--------|--------------------------------------|----------------------------|--------------------------------------|---|--|

**Lemma 2.** When  $\leq$  is a preorder,  $\sim$  in transitive.

*Proof.* This results from the transitivity of  $\leq (cf. \text{ Eq. } (7))$ .

**Lemma 3.** When  $\leq$  is a preorder,  $\sim$  is symmetric.

*Proof.* This results from the fact that the conjunction is symmetric and from Eq. (1):  $P_1 \sim P_2 \Leftrightarrow P_1 \lesssim P_2 \land P_2 \lesssim P_1 \Leftrightarrow P_2 \lesssim P_1 \land P_1 \lesssim P_2 \Leftrightarrow P_2 \sim P_1$ .

**Lemma 4.** When  $\leq$  is a preorder, > and < are converse.

*Proof.* This results from the fact that the conjunction is symmetric and from Eqs. (2) and (3):  $P_1 > P_2 \Leftrightarrow P_1 \not\leq P_2 \land P_2 \leq P_1 \Leftrightarrow P_2 \leq P_1 \Leftrightarrow P_2 < P_1$ .

**Lemma 5.** When  $\leq$  is a preorder, > and < are irreflexive.

*Proof.* For >, this results from the reflexivity of  $\leq$  and from Eq. (2):  $P > P \Leftrightarrow P \nleq P \land P \leq P \Leftrightarrow false \land true = false$ . For <, the proof is similar.

**Lemma 6.** When  $\leq$  is a preorder, > and < are asymmetric.

 $\begin{array}{l} \textit{Proof. For >, this is because } P_1 > P_2 \land P_2 > P_1 \Leftrightarrow (P_1 \not\lesssim P_2 \land P_2 \lesssim P_1) \land (P_2 \not\lesssim P_1 \land P_1 \lesssim P_2) \Leftrightarrow (P_1 \not\lesssim P_2 \land P_1 \lesssim P_2) \land (P_2 \not\lesssim P_1 \land P_2 \lesssim P_1) \Leftrightarrow false \land false \Leftrightarrow false. \ \text{For <, the proof is similar.} \end{array}$ 

**Lemma 7.** When  $\leq$  is a preorder, > and < are transitive.

*Proof.* This results from the transitivity of  $\leq (cf. \text{ Eqs. } (8) \text{ and } (9)).$ 

**Lemma 8.** When  $\leq$  is a preorder,  $\nleq$  is irreflexive.

*Proof.* This results from the reflexivity of  $\leq$  and from Eq. (4):  $P \not\leq P \Leftrightarrow P \not\leq P \land P \not\leq P \Leftrightarrow false \land false \Leftrightarrow false$ .  $\Box$ 

**Lemma 9.** When  $\leq$  is a preorder,  $\nleq$  is symmetric.

*Proof.* This results from the fact that the conjunction is symmetric and from Eq. (4):  $P_1 \nleq P_2 \Leftrightarrow P_1 \nleq P_2 \land P_2 \nleq P_1 \Leftrightarrow P_2 \nleq P_1 \land P_2 \not \lesssim P_1$ .

**Lemma 10.** When  $\leq$  is a preorder,  $\leq$  and  $\geq$  are converse.

*Proof.* From Eqs. (5) and (6), as > and < are converse, we have  $P_1 \gtrsim P_2 \Leftrightarrow P_1 > P_2 \lor P_1 \sim P_2 \Leftrightarrow P_2 < P_1 \lor P_2 \sim P_1 \Leftrightarrow P_2 \lesssim P_1$ .

**Lemma 11.** When  $\leq$  is a preorder,  $\leq$  and  $\geq$  are reflexive.

*Proof.* For  $\leq$ , it is by definition of preorders. For  $\geq$ , from Eqs. (1), (2) and (6), we have  $P \geq P \Leftrightarrow P > P \lor P \sim P \Leftrightarrow (P \leq P \land P \leq P) \lor (P \leq P \land P \leq P) \Leftrightarrow (false \land true) \lor (true \land true) \Leftrightarrow true.$ 

**Lemma 12.** When  $\leq$  is a preorder,  $\leq$  and  $\geq$  are transitive.

*Proof.* For  $\leq$ , it is by definition of preorders. For  $\geq$ , as  $\leq$  and  $\geq$  are converse,  $P_1 \geq P_2 \land P_2 \geq P_3 \Leftrightarrow P_3 \leq P_2 \land P_2 \leq P_1 \Rightarrow P_3 \leq P_1 \Rightarrow P_1 \geq P_3$ .

#### A.4. Supplementary Material about Sec. 3.6

#### A.4.1 The Visual Inspection of Formulas, to Determine the Importance given by Scores, can be Misleading!

Let us consider the example of two-class classification, with  $P(\{tn\})$ ,  $P(\{fn\})$ ,  $P(\{fn\})$ , and  $P(\{tp\})$  denoting, respectively, the probability (or proportion) of true negatives, false positives, false negatives, and true positives. Here are two classical scores, the accuracy and the true positive rate:

$$A = P(\{tn\}) + P(\{tp\}) \qquad TPR = \frac{P(\{tp\})}{P(\{fn\}) + P(\{tp\})}$$

The formula for the accuracy gives the illusion that the same importance is given to  $\{tn\}$  and  $\{tp\}$  and that no importance at all is given to  $\{fn\}$  and  $\{fn\}$ . For the true positive rate, the formula might give the impression that the same importance is given to  $\{fn\}$  and  $\{tp\}$  and that no importance at all is given to  $\{tn\}$  and  $\{tp\}$ . In fact, the visual inspection of formulas like these is not reliable at all to judge the importance given by a score to the various events. To see it, consider rewriting the previous equations as

$$\begin{split} A &= (1-\alpha)P(\{tn\}) - \alpha P(\{fp\}) - \alpha P(\{fn\}) + (1-\alpha)P(\{tp\}) + \alpha & \forall \alpha \\ TPR &= \frac{-\alpha P(\{tn\}) - \alpha P(\{fp\}) - \alpha P(\{fn\}) + (1-\alpha)P(\{tp\}) + \alpha}{-\beta P(\{tn\}) - \beta P(\{fp\}) + (1-\beta)P(\{fn\}) + (1-\beta)P(\{tp\}) + \beta} & \forall \alpha, \beta \end{split}$$

A visual inspection of such formulas would lead, indeed, to other illusions about the events that have no importance. This observation, however, should not stop us from thinking in terms of importance. In fact, we do it in this paper, but we do it in a mathematical framework that allows to do it rigorously.

#### A.4.2 So, How can we Determine the Importance given by Scores?

One cannot determine the application-specific preferences implicitly considered by any score X. However, one can determine those that are consistent, or have the best consistency, with X. We detail these notions hereafter.

For a given set of performances. Consider a score X and a ranking score  $R_I$ . If, for a given set  $\Pi \subseteq \mathbb{P}_{(\Omega,\Sigma)}$ , the scores X and  $R_I$  are linked by a strict monotonic relationship on  $\Pi \cap \operatorname{dom}(X) \cap \operatorname{dom}(R_I)$  and if  $\Pi \cap \operatorname{dom}(X) = \Pi \cap \operatorname{dom}(R_I)$ , then the performance orderings  $\leq_X$  and  $\leq_{R_I}$  (induced by X and  $R_I$  in the way specified in Theorem 1) are identical on  $\Pi$ . We say that the score X is consistent with the application-specific preferences I, on this set. A score can be consistent with different importance values (*e.g.*, as consequence of Properties 3 and 4).

For a given distribution of performances. Let  $\mathring{\mathbb{P}} = \{P \in \mathbb{P}_{(\Omega,\Sigma)} : P(\{\omega\}) > 0 \forall \omega \in \Omega)\}$ . All ranking scores are defined on this set. Consider a score X and the set  $\Pi = \text{dom}(X) \cap \mathring{\mathbb{P}}$ . We say that the score X has the best consistency with the application-specific preferences I when I maximizes the rank correlation between X and  $R_I$ , on  $\Pi$ , for the given distribution of performances. See Appendix A.7.2 for computational details.

#### A.5. Supplementary Material about Sec. 4.2

#### A.5.1 Proof of Theorem 1.

For convenience, we provide a reminder of Theorem 1 and Axiom 1 below.

**Theorem 1** (Sufficient condition for Axiom 1). A binary relation  $\leq_X$  on  $\mathbb{P}_{(\Omega,\Sigma)}$  induced by a score X as  $P_1 \leq_X P_2$  iff either  $P_1 = P_2$  or  $P_1 \in \text{dom}(X)$  and  $P_2 \in \text{dom}(X)$  and  $X(P_1) \leq X(P_2)$ , is a preorder satisfying Axiom 1.

**Axiom 1.** The ranking function  $\operatorname{rank}_{\mathbb{E}} : \mathbb{E} \to [1, |\mathbb{E}|] : \epsilon \mapsto \operatorname{rank}_{\mathbb{E}}(\epsilon)$  satisfies  $|\{\epsilon' \in \mathbb{E} : \operatorname{eval}(\epsilon) < \operatorname{eval}(\epsilon')\}| + 1 \leq \operatorname{rank}_{\mathbb{E}}(\epsilon) \leq |\{\epsilon' \in \mathbb{E} : \operatorname{eval}(\epsilon) \lesssim \operatorname{eval}(\epsilon')\}|$ , where  $\lesssim$  is a preorder on  $\mathbb{P}_{(\Omega, \Sigma)}$ .

*Proof.* To establish that  $\lesssim_X$  is a preorder, we have to show that it is (1) reflexive and (2) transitive.

- (1) The reflexivity of  $\leq_X$  is trivial to establish, since  $P_1 = P_2 \Rightarrow P_1 \leq_X P_2$ .
- (2) The transitivity of  $\leq_X$  can be shown as follows.  $P_1 \leq_X P_2 \land P_2 \leq_X P_3$  implies that:
  - either  $P_1 \in \text{dom}(X)$ ,  $P_2 \in \text{dom}(X)$ ,  $P_3 \in \text{dom}(X)$ , and  $X(P_1) \leq X(P_2) \wedge X(P_2) \leq X(P_3) \Rightarrow X(P_1) \leq X(P_3) \Rightarrow P_1 \lesssim_X P_3$ ;
  - or  $P_1 \notin \operatorname{dom}(X)$ ,  $P_2 \notin \operatorname{dom}(X)$ ,  $P_3 \notin \operatorname{dom}(X)$ , and  $P_1 = P_2 \wedge P_2 = P_3 \Rightarrow P_1 = P_3 \Rightarrow P_1 \lesssim_X P_3$ .

We conclude that, in all cases,  $P \leq_X P$  and  $P_1 \leq_X P_2 \land P_2 \leq_X P_3 \Rightarrow P_1 \leq_X P_3$ . The orderings  $\leq_X$  induced by scores X are thus preorders.

Summary. If the homogeneous binary relations  $\sim, >, <$ , and  $\notin$  on  $\mathbb{P}_{(\Omega,\Sigma)}$  are derived from the ordering  $\lesssim_X$  as explained above, and if the  $\lesssim_X$  is derived from the score X, then the comparison between performances  $P_1$  and  $P_2$  can be summarized as follows.

|  | $P_1 \in \operatorname{dom}(X)$                | $P_1 \not\in \operatorname{dom}(X)$          |
|--|--|--|
|  | $X(P_1) < X(P_2) \Leftrightarrow P_1 < P_2$    |  |
| $P_2 \in \operatorname{dom}(X)$          | $X(P_1) = X(P_2) \Leftrightarrow P_1 \sim P_2$ | $P_1 \nleq P_2$                              |
|  | $X(P_1) > X(P_2) \Leftrightarrow P_1 > P_2$    |  |
| $\mathbf{D}$ d $\mathbf{d}$ $\mathbf{v}$ | $p \leq p$                                     | $P_1 = P_2 \Leftrightarrow P_1 \sim P_2$     |
| $P_2 \not\in \operatorname{dom}(X)$      | $P_1 \not\equiv P_2$                           | $P_1 \neq P_2 \Leftrightarrow P_1 \nleq P_2$ |

#### A.5.2 Proof of Theorem 2

For convenience, we provide a reminder of Theorem 2 and Axiom 2 below.

**Theorem 2** (Sufficient condition for Axiom 2). If a score X satisfies  $\min_{\omega \in E} S(\omega) \leq X(P) \leq \max_{\omega \in E} S(\omega)$  for all events  $E \in \Sigma$  and all performances  $P \in \text{dom}(X)$  such that P(E) = 1, then the ordering  $\lesssim_X$  satisfies Axiom 2.

**Axiom 2.** For  $P_1, P_2 \in \mathbb{P}_{(\Omega,\Sigma)}$  such that  $P_1(S \leq s) = 1$  and  $P_2(S \geq s) = 1$  for some s, then  $P_1 \leq P_2$  or  $P_1 \notin P_2$ .

*Proof.* Axiom 2 is satisfied when  $P_1 \notin \text{dom}(X)$  or  $P_2 \notin \text{dom}(X)$ .

- Either  $P_1 = P_2 \Leftrightarrow P_1 \sim P_2 \Rightarrow P_1 \lesssim P_2$ ,
- or  $P_1 \neq P_2 \Leftrightarrow P_1 \nleq P_2$ .

Axiom 2 is also satisfied when  $P_1 \in \text{dom}(X)$  and  $P_2 \in \text{dom}(X)$ .

- On the one hand, the axiom stipulates that the event  $E_1 = \{\omega \in \Omega : S(\omega) \le s\}$  and the performance  $P_1$  are such that  $P_1(E_1) = 1$ . Trivially, we have  $\max_{\omega \in E_1} S(\omega) \le s$ . On the other hand, the theorem stipulates that, as  $P_1(E_1) = 1$ ,  $X(P_1) \le \max_{\omega \in E_1} S(\omega)$ . Putting all together, we have  $X(P_1) \le s$ .
- On the one hand, the axiom stipulates that the event  $E_2 = \{\omega \in \Omega : S(\omega) \ge s\}$  and the performance  $P_2$  are such that  $P_2(E_2) = 1$ . Trivially, we have  $s \le \min_{\omega \in E_2} S(\omega)$ . On the other hand, the theorem stipulates that, as  $P_2(E_2) = 1$ ,  $\min_{\omega \in E_2} S(\omega) \le X(P_2)$ . Putting all together, we have  $s \le X(P_2)$ .
- As we have established that  $X(P_1) \leq s$  and  $s \leq X(P_2)$ , we have  $X(P_1) \leq X(P_2) \Leftrightarrow P_1 \leq P_2$ .

#### A.5.3 Proof of Theorem 3

For convenience, we provide a reminder of Theorem 3 and Axiom 3 below.

**Theorem 3** (Sufficient condition for Axiom 3). If a score X is such that  $\Pi \subseteq \operatorname{dom}(X) \Rightarrow \Phi(\Pi) \subseteq \operatorname{dom}(X)$  and  $\min_{P \in \Pi} X(P) \leq X(\overline{P}) \leq \max_{P \in \Pi} X(P)$  for all  $\Pi \subseteq \operatorname{dom}(X)$  and all  $\overline{P} \in \Phi(\Pi)$ , then the ordering  $\lesssim_X$  satisfies Axiom 3.

**Axiom 3.** Let P be a performance, and  $\Pi$  be a set of performances on  $\mathbb{P}_{(\Omega,\Sigma)}$  such that  $P' \lesssim P \lor P \lesssim P' \forall P' \in \Pi$ .

- $P' \lesssim P \forall P' \in \Pi \Rightarrow \overline{P} \lesssim P \forall \overline{P} \in \Phi(\Pi);$   $P' \lesssim P \forall P' \in \Pi \Rightarrow \overline{P} \lesssim P \forall \overline{P} \in \Phi(\Pi);$   $P \lesssim P' \forall P' \in \Pi \Rightarrow P \lesssim \overline{P} \forall \overline{P} \in \Phi(\Pi);$  and  $P \lesssim P' \forall P' \in \Pi \Rightarrow P \lesssim \overline{P} \forall \overline{P} \in \Phi(\Pi);$

*Proof.* We take  $\leq = \leq_X$  and  $\Pi \neq \emptyset$ .

Remainder of the conditions. The first condition of Theorem 3 is:

$$\Pi \subseteq \operatorname{dom}(X) \Rightarrow \Phi(\Pi) \subseteq \operatorname{dom}(X).$$
(15)

The second condition of Theorem 3 is:

$$\min_{P \in \Pi} X(P) \le X(\overline{P}) \le \max_{P \in \Pi} X(P) \qquad \forall \Pi \subseteq \operatorname{dom}(X) \qquad \forall \overline{P} \in \Phi(\Pi) \,.$$
(16)

The condition of Axiom 3 is that P is comparable to all performances in the set  $\Pi$ :

$$P' \lesssim P \lor P \lesssim P' \qquad \forall P' \in \Pi.$$
 (17)

On the domain of X. By Theorem 1, this last condition implies that

$$P \in \operatorname{dom}(X), \tag{18}$$

and

$$\Pi \subseteq \operatorname{dom}(X) \,. \tag{19}$$

Taking Eq. (15) and Eq. (19) together, we have

$$\Phi(\Pi) \subseteq \operatorname{dom}(X) \qquad \Leftrightarrow \qquad \overline{P} \in \operatorname{dom}(X) \qquad \forall \overline{P} \in \Phi(\Pi) \,. \tag{20}$$

**Proof that**  $P' \leq P \forall P' \in \Pi \Rightarrow \overline{P} \leq P \forall \overline{P} \in \Phi(\Pi)$ . On the one hand, we have, by Theorem 1,

$$P' \lesssim P \,\forall P' \in \Pi \Leftrightarrow X(P') \leq X(P) \qquad \forall P' \in \Pi$$
$$\Leftrightarrow \max_{P' \in \Pi} X(P') \leq X(P) \,.$$

On the other hand, Eq. (16) implies that

$$X(\overline{P}) \le \max_{P' \in \Pi} X(P') \qquad \forall \overline{P} \in \Phi(\Pi)$$

Considering the last two equations together, we obtain

$$\begin{split} X(\overline{P}) &\leq \max_{P' \in \Pi} X(P') \leq X(P) \qquad \forall \, \overline{P} \in \Phi(\Pi) \\ \Rightarrow &X(\overline{P}) \leq X(P) \qquad \forall \, \overline{P} \in \Phi(\Pi) \,. \end{split}$$

By Theorem 1, we have thus  $\overline{P} \leq P$ .

**Proof that**  $P' \not\leq P \forall P' \in \Pi \Rightarrow \overline{P} \not\leq P \forall \overline{P} \in \Phi(\Pi)$ . On the one hand, we have, by Eq. (17) and Theorem 1,

$$P' \not\lesssim P \forall P' \in \Pi \Rightarrow P < P' \forall P' \in \Pi$$
  
$$\Leftrightarrow X(P) < X(P') \qquad \forall P' \in \Pi$$
  
$$\Leftrightarrow X(P) < \min_{P' \in \Pi} X(P')$$

On the other hand, Eq. (16) implies that

$$\min_{P'\in\Pi} X(P') \le X(\overline{P}) \qquad \forall \, \overline{P} \in \Phi(\Pi)$$

Considering the last two equations together, we obtain

$$\begin{split} X(P) &< \min_{P' \in \Pi} X(P') \leq X(\overline{P}) \qquad \forall \, \overline{P} \in \Phi(\Pi) \\ \Rightarrow &X(P) < X(\overline{P}) \qquad \forall \, \overline{P} \in \Phi(\Pi) \,. \end{split}$$

By Theorem 1, we have thus  $P < \overline{P}$ , and by Eq. (17),  $\overline{P} \leq P$ . **Proof that**  $P \leq P' \forall P' \in \Pi \Rightarrow P \leq \overline{P} \forall \overline{P} \in \Phi(\Pi)$ . On the one hand, we have, by Theorem 1,

$$P \lesssim P' \,\forall P' \in \Pi \Leftrightarrow X(P) \le X(P') \qquad \forall P' \in \Pi$$
$$\Leftrightarrow X(P) \le \min_{P' \in \Pi} X(P') \,.$$

On the other hand, Eq. (16) implies that

$$\min_{P'\in\Pi} X(P') \le X(\overline{P}) \qquad \forall \, \overline{P} \in \Phi(\Pi)$$

Considering the last two equations together, we obtain

$$\begin{split} X(P) &\leq \min_{P' \in \Pi} X(P') \leq X(\overline{P}) \qquad \forall \, \overline{P} \in \Phi(\Pi) \\ \Rightarrow &X(P) \leq X(\overline{P}) \qquad \forall \, \overline{P} \in \Phi(\Pi) \,. \end{split}$$

By Theorem 1, we have thus  $P \lesssim \overline{P}$ .

**Proof that**  $P \not\leq P' \forall P' \in \Pi \Rightarrow P \not\leq \overline{P} \forall \overline{P} \in \Phi(\Pi)$ . On the one hand, we have, by Eq. (17) and Theorem 1,

$$\begin{split} P \not\lesssim P' \,\forall \, P' \in \Pi \Rightarrow & P' < P \,\forall \, P' \in \Pi \\ \Leftrightarrow & X(P') < X(P) \qquad \forall \, P' \in \Pi \\ \Leftrightarrow & \max_{P' \in \Pi} X(P') < X(P) \end{split}$$

On the other hand, Eq. (16) implies that

$$X(\overline{P}) \le \max_{P' \in \Pi} X(P') \qquad \forall \overline{P} \in \Phi(\Pi)$$

Considering the last two equations together, we obtain

$$X(\overline{P}) \le \max_{P' \in \Pi} X(P') < X(P) \qquad \forall \overline{P} \in \Phi(\Pi)$$
  
$$\Rightarrow X(\overline{P}) < X(P) \qquad \forall \overline{P} \in \Phi(\Pi) .$$

By Theorem 1, we have thus  $\overline{P} < P$ , and by Eq. (17),  $P \not\leq \overline{P}$ .

#### A.6. Supplementary Material about Sec. 4.3

#### A.6.1 All Ranking Scores can be Used to Rank Performances (for $\Phi = \text{conv}$ )

To show that all ranking scores can be used to rank performances, for  $\Phi = \text{conv}$ , we show that these scores satisfy the conditions of Theorems 1, 2, and 3.

All ranking scores satisfy the conditions of Theorem 1, and thus Axiom 1. For convenience, we provide a reminder of Theorem 1 and Axiom 1 below.

**Theorem 1** (Sufficient condition for Axiom 1). A binary relation  $\leq_X$  on  $\mathbb{P}_{(\Omega,\Sigma)}$  induced by a score X as  $P_1 \leq_X P_2$  iff either  $P_1 = P_2$  or  $P_1 \in \text{dom}(X)$  and  $P_2 \in \text{dom}(X)$  and  $X(P_1) \leq X(P_2)$ , is a preorder satisfying Axiom 1.

**Axiom 1.** The ranking function  $\operatorname{rank}_{\mathbb{E}} : \mathbb{E} \to [1, |\mathbb{E}|] : \epsilon \mapsto \operatorname{rank}_{\mathbb{E}}(\epsilon)$  satisfies  $|\{\epsilon' \in \mathbb{E} : \operatorname{eval}(\epsilon) < \operatorname{eval}(\epsilon')\}| + 1 \leq \operatorname{rank}_{\mathbb{E}}(\epsilon) \leq |\{\epsilon' \in \mathbb{E} : \operatorname{eval}(\epsilon) \lesssim \operatorname{eval}(\epsilon')\}|$ , where  $\lesssim$  is a preorder on  $\mathbb{P}_{(\Omega, \Sigma)}$ .

**Theorem 4.** All ranking scores satisfy the conditions of Theorem 1.

*Proof.* For all ranking scores  $R_I$ , it is possible to induce an ordering  $\leq_{R_I}$  satisfying the requirements of Theorem 1.

All ranking scores satisfy the conditions of Theorem 2, and thus Axiom 2. For convenience, we provide a reminder of Theorem 2 and Axiom 2 below.

**Theorem 2** (Sufficient condition for Axiom 2). If a score X satisfies  $\min_{\omega \in E} S(\omega) \leq X(P) \leq \max_{\omega \in E} S(\omega)$  for all events  $E \in \Sigma$  and all performances  $P \in \text{dom}(X)$  such that P(E) = 1, then the ordering  $\lesssim_X$  satisfies Axiom 2.

**Axiom 2.** For  $P_1, P_2 \in \mathbb{P}_{(\Omega,\Sigma)}$  such that  $P_1(S \leq s) = 1$  and  $P_2(S \geq s) = 1$  for some s, then  $P_1 \leq P_2$  or  $P_1 \notin P_2$ .

**Theorem 5.** All ranking scores satisfy the conditions of Theorem 2.

*Proof.* We take  $X = R_I$ . When P(E) = 1, we have

$$R_{I}(P) = \frac{\sum_{\omega \in \Omega} I(\omega)S(\omega)P(\{\omega\})}{\sum_{\omega \in \Omega} I(\omega)P(\{\omega\})} = \frac{\sum_{\omega \in E} I(\omega)S(\omega)P(\{\omega\})}{\sum_{\omega \in E} I(\omega)P(\{\omega\})}$$

with  $\sum_{\omega \in \Omega} I(\omega)P(\{\omega\}) > 0$  when  $P \in \operatorname{dom}(R_I)$ . • Let  $M = \max_{\omega \in E} S(\omega)$ . We have

$$\begin{split} S(\omega) &\leq M \quad \forall \omega \in E \\ \Leftrightarrow S(\omega) - M \leq 0 \quad \forall \omega \in E \\ \Rightarrow &\sum_{\omega \in E} I(\omega) \left[ S(\omega) - M \right] P(\{\omega\}) \leq 0 \quad \text{as } I(\omega) \geq 0 \text{ and } P(\{\omega\}) \geq 0 \\ \Leftrightarrow &\sum_{\omega \in E} I(\omega) S(\omega) P(\{\omega\}) \leq \sum_{\omega \in E} I(\omega) M P(\{\omega\}) \\ \Leftrightarrow &\sum_{\omega \in E} I(\omega) S(\omega) P(\{\omega\}) \leq M \underbrace{\sum_{\omega \in E} I(\omega) P(\{\omega\})}_{>0} \\ \Leftrightarrow &\underbrace{\sum_{\omega \in E} I(\omega) S(\omega) P(\{\omega\})}_{\sum_{\omega \in E} I(\omega) P(\{\omega\})} \leq M \\ \Leftrightarrow &R_I(P) \leq M \end{split}$$

• Let  $m = \min_{\omega \in E} S(\omega)$ . We have

$$\begin{split} S(\omega) &\geq m \quad \forall \omega \in E \\ \Leftrightarrow S(\omega) - m \geq 0 \quad \forall \omega \in E \\ \Rightarrow &\sum_{\omega \in E} I(\omega) \left[ S(\omega) - m \right] P(\{\omega\}) \geq 0 \quad \text{as } I(\omega) \geq 0 \text{ and } P(\{\omega\}) \geq 0 \\ \Leftrightarrow &\sum_{\omega \in E} I(\omega) S(\omega) P(\{\omega\}) \geq \sum_{\omega \in E} I(\omega) m P(\{\omega\}) \\ \Leftrightarrow &\sum_{\omega \in E} I(\omega) S(\omega) P(\{\omega\}) \geq m \underbrace{\sum_{\omega \in E} I(\omega) P(\{\omega\})}_{>0} \\ \Leftrightarrow &\underbrace{\sum_{\omega \in E} I(\omega) S(\omega) P(\{\omega\})}_{\sum_{\omega \in E} I(\omega) P(\{\omega\})} \geq m \\ \Leftrightarrow &R_I(P) \geq m \end{split}$$

• Putting all together, when P(E) = 1, we have  $m \leq R_I(P) \leq M$ , and so,

$$\min_{\omega \in E} S(\omega) \le R_I(P) \le \max_{\omega \in E} S(\omega) \,. \tag{21}$$

All ranking scores satisfy the conditions of Theorem 3, and thus Axiom 3 (for  $\Phi = \text{conv}$ ). For convenience, we provide a reminder of Theorem 3 and Axiom 3 below.

**Theorem 3** (Sufficient condition for Axiom 3). If a score X is such that  $\Pi \subseteq \operatorname{dom}(X) \Rightarrow \Phi(\Pi) \subseteq \operatorname{dom}(X)$  and  $\min_{P \in \Pi} X(P) \leq X(\overline{P}) \leq \max_{P \in \Pi} X(P)$  for all  $\Pi \subseteq \operatorname{dom}(X)$  and all  $\overline{P} \in \Phi(\Pi)$ , then the ordering  $\lesssim_X$  satisfies Axiom 3.

**Axiom 3.** Let P be a performance, and  $\Pi$  be a set of performances on  $\mathbb{P}_{(\Omega,\Sigma)}$  such that  $P' \leq P \lor P \leq P' \forall P' \in \Pi$ .

- $\begin{array}{l} \bullet \ P' \lesssim P \ \forall P' \in \Pi \Rightarrow \overline{P} \lesssim P \ \forall \overline{P} \in \Phi(\Pi); \\ \bullet \ P' \lesssim P \ \forall P' \in \Pi \Rightarrow \overline{P} \ \lesssim P \ \forall \overline{P} \in \Phi(\Pi); \\ \bullet \ P \ \lesssim P' \ \forall P' \in \Pi \Rightarrow P \ \lesssim \overline{P} \ \forall \overline{P} \in \Phi(\Pi); \\ \bullet \ and \ P \ \lesssim P' \ \forall P' \in \Pi \Rightarrow P \ \lesssim \overline{P} \ \forall \overline{P} \in \Phi(\Pi). \end{array}$

**Theorem 6.** All ranking scores satisfy the conditions of Theorem 3 (for  $\Phi = \text{conv}$ ).

*Proof.* The proof is in two parts.

• First, let us show that  $\Pi \subseteq \operatorname{dom}(R_I) \Rightarrow \operatorname{conv}(\Pi) \subseteq \operatorname{dom}(R_I)$ . For any  $\overline{P} \in \operatorname{conv}(\Pi)$  there exists a weighting function  $\lambda_{\Pi,\overline{P}}:\Pi\to\mathbb{R}_{\geq 0}:P\mapsto\lambda_{\Pi,\overline{P}}(P)\text{ such that }\sum_{P\in\Pi}\lambda_{\Pi,\overline{P}}(P)=1\text{ and }\sum_{P\in\Pi}\lambda_{\Pi,\overline{P}}(P)P=\overline{P}.\text{ For all }\overline{P}\in\operatorname{conv}(\Pi)\text{, we}$ have:

$$\Pi \subseteq \operatorname{dom}(R_I) \Leftrightarrow \sum_{\omega \in \Omega} I(\omega) P(\{\omega\}) \neq 0 \qquad \forall P \in \Pi$$
$$\Rightarrow \sum_{P \in \Pi} \lambda_{\Pi, \overline{P}}(P) \sum_{\omega \in \Omega} I(\omega) P(\{\omega\}) \neq 0$$
$$\Leftrightarrow \sum_{\omega \in \Omega} I(\omega) \sum_{P \in \Pi} \lambda_{\Pi, \overline{P}} P(\{\omega\}) \neq 0$$
$$\Leftrightarrow \sum_{\omega \in \Omega} I(\omega) \overline{P}(\{\omega\}) \neq 0$$
$$\Leftrightarrow \overline{P} \in \operatorname{dom}(R_I) .$$

• Second, let us show that, for all  $\overline{P} \in \operatorname{conv}(\Pi)$ ,  $\min_{P \in \Pi} R_I(P) \leq R_I(\overline{P}) \leq \max_{P \in \Pi} R_I(P)$ . Let us pose  $l = \min_{P \in \Pi} R_I(P)$  and  $u = \max_{P \in \Pi} R_I(P)$ . We have:

$$\begin{split} l &\leq R_{I}(P) \leq u \quad \forall P \in \Pi \\ \Leftrightarrow l \leq \frac{\sum_{\omega \in \Omega} I(\omega)S(\omega)P(\{\omega\})}{\sum_{\omega \in \Omega} I(\omega)P(\{\omega\})} \leq u \quad \forall P \in \Pi \\ \Leftrightarrow l \sum_{\omega \in \Omega} I(\omega)P(\{\omega\}) \leq \sum_{\omega \in \Omega} I(\omega)S(\omega)P(\{\omega\}) \leq u \sum_{\omega \in \Omega} I(\omega)P(\{\omega\}) \quad \forall P \in \Pi \\ \Rightarrow l \sum_{\omega \in \Omega} I(\omega)\overline{P}(\{\omega\}) \leq \sum_{\omega \in \Omega} I(\omega)S(\omega)\overline{P}(\{\omega\}) \leq u \sum_{\omega \in \Omega} I(\omega)\overline{P}(\{\omega\}) \\ \Leftrightarrow l \leq \frac{\sum_{\omega \in \Omega} I(\omega)S(\omega)\overline{P}(\{\omega\})}{\sum_{\omega \in \Omega} I(\omega)\overline{P}(\{\omega\})} \leq u \\ \Leftrightarrow l \leq R_{I}(\overline{P}) \leq u \end{split}$$

#### A.6.2 On the Properties of Ranking Scores.

*Proof of Property 1.* Let us demonstrate that we have  $R_I(P) = X_S^E(P')$  with  $P' = \text{filter}_I(P)$ . For all  $\omega \in \Omega$ , we have:

$$P'(\{\omega\}) = \frac{P(\{\omega\})I(\omega)}{\sum_{\omega'\in\Omega} P(\{\omega'\})I(\omega')}$$

Thus,

$$\begin{split} X_{S}^{E}(P') &= \sum_{\omega \in \Omega} P'(\{\omega\})S(\omega) \\ &= \sum_{\omega \in \Omega} \frac{P(\{\omega\})I(\omega)}{\sum_{\omega' \in \Omega} P(\{\omega'\})I(\omega')}S(\omega) \\ &= \frac{\sum_{\omega \in \Omega} P(\{\omega\})I(\omega)S(\omega)}{\sum_{\omega' \in \Omega} P(\{\omega'\})I(\omega')} \\ &= R_{I}(P) \end{split}$$

| 1 |  | _ | 1 |
|---|--|---|---|
|   |  |   | L |
|   |  |   | L |

*Proof of Property 2.* Let  $S' = \alpha S + \beta$  with  $\alpha, \beta \in \mathbb{R}$ .

$$\frac{\sum_{\omega \in \Omega} P(\{\omega\}) S'(\omega) I(\omega)}{\sum_{\omega \in \Omega} P(\{\omega\}) I(\omega)} = \alpha \frac{\sum_{\omega \in \Omega} P(\{\omega\}) S(\omega) I(\omega)}{\sum_{\omega \in \Omega} P(\{\omega\}) I(\omega)} + \beta$$

Proof of Property 3. 
$$R_{kI} = \frac{\sum_{\omega \in \Omega} kI(\omega)S(\omega)P(\{\omega\})}{\sum_{\omega \in \Omega} kI(\omega)P(\{\omega\})} = \frac{\sum_{\omega \in \Omega} I(\omega)S(\omega)P(\{\omega\})}{\sum_{\omega \in \Omega} I(\omega)P(\{\omega\})} = R_I$$

*Proof of Property 4.* Let us consider a binary satisfaction, that is  $S(\omega) \in \{0, 1\} \forall \omega \in \Omega$ . Let us define the events  $E_0 = \{\omega \in \Omega : S(\omega) = 0\}$  and  $E_1 = \{\omega \in \Omega : S(\omega) = 1\}$ . If  $I' = (\mathbf{1}_{S=0}\alpha_0 + \mathbf{1}_{S=1}\alpha_1)I$  with  $\alpha_0 > 0$  and  $\alpha_1 > 0$ , then

$$R_{I}(P) = \frac{\sum_{\omega \in \Omega} I(\omega)S(\omega)P(\{\omega\})}{\sum_{\omega \in \Omega} I(\omega)P(\{\omega\})}$$
$$= \frac{\sum_{\omega \in E_{1}} I(\omega)P(\{\omega\})}{\sum_{\omega \in E_{0}} I(\omega)P(\{\omega\}) + \sum_{\omega \in E_{1}} I(\omega)P(\{\omega\})}$$

and

$$R_{I'}(P) = \frac{\sum_{\omega \in \Omega} I'(\omega)S(\omega)P(\{\omega\})}{\sum_{\omega \in \Omega} I'(\omega)P(\{\omega\})}$$
$$= \frac{\sum_{\omega \in E_1} I'(\omega)P(\{\omega\})}{\sum_{\omega \in E_0} I'(\omega)P(\{\omega\}) + \sum_{\omega \in E_1} I'(\omega)P(\{\omega\})}$$
$$= \frac{\sum_{\omega \in E_1} \alpha_1 I(\omega)P(\{\omega\})}{\sum_{\omega \in E_0} \alpha_0 I(\omega)P(\{\omega\}) + \sum_{\omega \in E_1} \alpha_1 I(\omega)P(\{\omega\})}$$
$$= \frac{\alpha_1 \sum_{\omega \in E_1} I(\omega)P(\{\omega\})}{\alpha_0 \sum_{\omega \in E_0} I(\omega)P(\{\omega\}) + \alpha_1 \sum_{\omega \in E_1} I(\omega)P(\{\omega\})}$$

Thus,  $R_{I'} = \frac{\alpha_1 R_I}{\alpha_0 (1-R_I) \alpha_1 R_I}$  and  $\frac{\partial R_{I'}}{\partial R_I} = \frac{\alpha_0 \alpha_1}{(\alpha_0 (1-R_I) \alpha_1 R_I)^2} > 0$ . This leads immediately to the conclusion that  $\leq_{R_{I'}} = \leq_{R_I}$ .

*Proof of Properties 5 and 6.* Let us consider a binary satisfaction and the events  $E_0 = \{\omega \in \Omega : S(\omega) = 0\}$  and  $E_1 = \{\omega \in \Omega : S(\omega) = 1\}$ . Let  $I_1$  and  $I_2$  be two random variables and  $I = \lambda_1 I_1 + \lambda_2 I_2$  with  $\lambda_1, \lambda_2 \in \mathbb{R}$  such that  $\lambda_1 + \lambda_2 = 1$ . • When the random variables  $I_1$  and  $I_2$  are such that  $I_1(\omega) = I_2(\omega) = I(\omega) \forall \omega \in E_1$ , if we take  $f : x \mapsto x^{-1}$ ,

$$\lambda_{1}f\left(R_{I_{1}}(P)\right) + \lambda_{2}f\left(R_{I_{1}}(P)\right)$$

$$=\lambda_{1}\left(1 + \frac{\sum_{\omega \in E_{0}} I_{1}(\omega)P(\{\omega\})}{\sum_{\omega \in E_{1}} I_{1}(\omega)P(\{\omega\})}\right) + \lambda_{2}\left(1 + \frac{\sum_{\omega \in E_{0}} I_{2}(\omega)P(\{\omega\})}{\sum_{\omega \in E_{1}} I_{2}(\omega)P(\{\omega\})}\right)$$

$$=\lambda_{1}\left(1 + \frac{\sum_{\omega \in E_{0}} I_{1}(\omega)P(\{\omega\})}{\sum_{\omega \in E_{1}} I(\omega)P(\{\omega\})}\right) + \lambda_{2}\left(1 + \frac{\sum_{\omega \in E_{0}} I_{2}(\omega)P(\{\omega\})}{\sum_{\omega \in E_{1}} I(\omega)P(\{\omega\})}\right)$$

$$=(\lambda_{1} + \lambda_{2}) + \frac{\lambda_{1}\sum_{\omega \in E_{0}} I_{1}(\omega)P(\{\omega\}) + \lambda_{2}\sum_{\omega \in E_{0}} I_{2}(\omega)P(\{\omega\})}{\sum_{\omega \in E_{1}} I(\omega)P(\{\omega\})}$$

$$=1 + \frac{\sum_{\omega \in E_{0}} (\lambda_{1}I_{1} + \lambda_{2}I_{2})(\omega)P(\{\omega\})}{\sum_{\omega \in E_{1}} I(\omega)P(\{\omega\})}$$

$$=1 + \frac{\sum_{\omega \in E_{0}} I(\omega)P(\{\omega\})}{\sum_{\omega \in E_{1}} I(\omega)P(\{\omega\})}$$

$$=f\left(R_{I}(P)\right)$$

• When the random variables  $I_1$  and  $I_2$  are such that  $I_1(\omega) = I_2(\omega) = I(\omega) \ \forall \omega \in E_0$ , if we take  $f: x \mapsto (1-x)^{-1}$ ,

$$\begin{split} \lambda_1 f\left(R_{I_1}(P)\right) &+ \lambda_2 f\left(R_{I_1}(P)\right) \\ = \lambda_1 \left(1 + \frac{\sum_{\omega \in E_1} I_1(\omega) P(\{\omega\})}{\sum_{\omega \in E_0} I_1(\omega) P(\{\omega\})}\right) + \lambda_2 \left(1 + \frac{\sum_{\omega \in E_1} I_2(\omega) P(\{\omega\})}{\sum_{\omega \in E_0} I_2(\omega) P(\{\omega\})}\right) \\ = \lambda_1 \left(1 + \frac{\sum_{\omega \in E_1} I_1(\omega) P(\{\omega\})}{\sum_{\omega \in E_0} I(\omega) P(\{\omega\})}\right) + \lambda_2 \left(1 + \frac{\sum_{\omega \in E_1} I_2(\omega) P(\{\omega\})}{\sum_{\omega \in E_0} I(\omega) P(\{\omega\})}\right) \\ = (\lambda_1 + \lambda_2) + \frac{\lambda_1 \sum_{\omega \in E_1} I_1(\omega) P(\{\omega\}) + \lambda_2 \sum_{\omega \in E_1} I_2(\omega) P(\{\omega\})}{\sum_{\omega \in E_0} I(\omega) P(\{\omega\})} \\ = 1 + \frac{\sum_{\omega \in E_1} (\lambda_1 I_1 + \lambda_2 I_2)(\omega) P(\{\omega\})}{\sum_{\omega \in E_0} I(\omega) P(\{\omega\})} \\ = 1 + \frac{\sum_{\omega \in E_1} I(\omega) P(\{\omega\})}{\sum_{\omega \in E_0} I(\omega) P(\{\omega\})} \\ = f\left(R_I(P)\right) \end{split}$$

*Proof of Property 7.* Let  $\mathcal{R} \in \{<, \leq, =, \geq, >\}$  and  $v = R_I(P)$ . For all  $P' \in \mathbb{P}_{(\Omega, \Sigma)}$ , we have:

$$R_I(P')\mathcal{R}R_I(P) \tag{22}$$

$$\Leftrightarrow R_I(P')\mathcal{R}v \tag{23}$$

$$\Leftrightarrow \frac{\sum_{\omega \in \Omega} I(\omega) S(\omega) P'(\{\omega\})}{\sum_{\omega \in \Omega} I(\omega) P'(\{\omega\})} \mathcal{R}v$$
(24)

$$\Leftrightarrow \left[\sum_{\omega \in \Omega} I(\omega) S(\omega) P'(\{\omega\})\right] \mathcal{R}\left[v \sum_{\omega \in \Omega} I(\omega) P'(\{\omega\})\right]$$
(25)

$$\Leftrightarrow \sum_{\omega \in \Omega} I(\omega) \left[ S(\omega) - v \right] P'(\{\omega\}) \mathcal{R}0$$
(26)

This is either a linear equality or a linear inequality constraint. Thus,

$$\phi_{\mathcal{R}}(P) = \left\{ P' \in \mathbb{P}_{(\Omega,\Sigma)} : R_I(P') \mathcal{R}R_I(P) \right\}$$
(27)

$$= \left\{ P' \in \mathbb{P}_{(\Omega,\Sigma)} : \sum_{\omega \in \Omega} I(\omega) \left[ S(\omega) - v \right] P'(\{\omega\}) \mathcal{R}0 \right\}$$
(28)

is a convex subset of  $\mathbb{P}_{(\Omega,\Sigma)}$ .

Figure A.7.1. Passages between two formulations (left: classical, right: ours) for the performance analysis of two-class classification problems.

#### A.7. Supplementary Material about Sec. 5.2

#### A.7.1 Link between Classical Formulation and Ours.

Fig. A.7.1 shows the connections between the classical formulation of the two-class classification task and our formulation, as explained in Sec. 5.

#### A.7.2 Custom Optimization Algorithm to Estimate Kendall's $\tau$ .

For any score X, our algorithm aims at determining the minimum and maximum values that a rank correlation between Xand our ranking scores  $R_I$  can take over all possible importances I. Note that this algorithm is not specific to Kendall's  $\tau$  [13] and could also be used with any other rank correlation, for example Spearman's  $\rho$  [23].

- **Variables.** Leveraging Properties 3 and 4, we know that the rank-correlation between X and a ranking score  $R_{I_1}$  is equal to the rank-correlation between X and another ranking score  $R_{I_2}$  if  $\frac{I_1(tp)}{I_1(tn)+I_1(tp)} = \frac{I_2(tp)}{I_2(tn)+I_2(tp)}$  and  $\frac{I_1(fn)}{I_1(fp)+I_1(fn)} = \frac{I_2(fn)}{I_2(fp)+I_2(fn)}$ . For this reason, we consider only two variables:  $a = \frac{I(tp)}{I(tn)+I(tp)} \in [0,1]$  and  $b = \frac{I(fn)}{I(fp)+I(fn)} \in [0,1]$ .
- **Objective function.** We optimize the function  $\tau(a, b)$  that gives the rank correlation between X and  $R_{I^*}$  with  $I^*(tp) = 1-a$ ,  $I^{*}(tp) = 1 - b$ ,  $I^{*}(tp) = b$ , and  $I^{*}(tp) = a$ . In practice, this is an estimation based on a finite set of performances on which X and  $R_{I^*}$  are applied. Note that  $\tau(a, b)$  is not a continuous function when estimated on a finite set of performances. The chosen optimization technique circumvents the difficulties related to that.
- **Optimization technique.** We implemented a custom coarse-to-fine grid-based direct search [5]: we compute  $\tau(a, b)$  on a coarse grid over the unit square, locate the maximum on the grid, center a smaller square and a finer grid around that point, and iterate until the square is small enough.

#### A.7.3 Scores Perfectly Correlated with a Ranking Score, for all Performances

- The accuracy:  $A = R_I$  with  $I(tn) = \frac{1}{2}$ ,  $I(fp) = \frac{1}{2}$ ,  $I(fn) = \frac{1}{2}$ , and  $I(tp) = \frac{1}{2}$ .
- The F-score with  $\beta = 0.5$ :  $F_{0.5} = R_I$  with I(tn) = 0,  $I(fp) = \frac{4}{5}$ ,  $I(fn) = \frac{1}{5}$ , and I(tp) = 1.
- The F-score with  $\beta = 1.0$ :  $F_1 = R_I$  with I(tn) = 0,  $I(fp) = \frac{1}{2}$ ,  $I(fn) = \frac{1}{2}$ , and I(tp) = 1.
- The F-score with  $\beta = 2.0$ :  $F_2 = R_I$  with I(tn) = 0,  $I(fp) = \frac{1}{5}$ ,  $I(fn) = \frac{4}{5}$ , and I(tp) = 1.
- The negative predictive value:  $NPV = R_I$  with I(tn) = 1, I(fp) = 0, I(fn) = 1, and I(tp) = 0.
- The positive predictive value:  $PPV = R_I$  with I(tn) = 0, I(fp) = 1, I(fn) = 0, and I(tp) = 1.
- The true negative rate:  $TNR = R_I$  with I(tn) = 1, I(fp) = 1, I(fn) = 0, and I(tp) = 0.
- The true positive rate:  $TPR = R_I$  with I(tn) = 0, I(fp) = 0, I(fn) = 1, and I(tp) = 1.
- A.7.4 Scores Perfectly Correlated with a Ranking Score, for the Performances Corresponding to Given Class Priors  $\pi_{-} \neq 0$  and  $\pi_{+} \neq 0$
- The balanced accuracy:  $BA = R_I$  with  $I(tn) = \pi_+$ ,  $I(fp) = \pi_+$ ,  $I(fn) = \pi_-$ , and  $I(tp) = \pi_-$ . Cohen's kappa:  $\kappa = \frac{R_I 2\pi_- \pi_+}{\pi_-^2 + \pi_+^2}$  with  $I(tn) = \frac{\pi_+^2}{\pi_-^2 + \pi_+^2}$ ,  $I(fp) = \frac{1}{2}$ ,  $I(fn) = \frac{1}{2}$ , and  $I(tp) = \frac{\pi_-^2}{\pi_-^2 + \pi_+^2}$ . Thus,  $\frac{\partial \kappa}{R_I} > 0$ .

- The informedness (a.k.a. Youden's J):  $J_Y = 2R_I 1$  with  $I(tn) = \pi_+$ ,  $I(fp) = \pi_+$ ,  $I(fn) = \pi_-$ , and  $I(tp) = \pi_-$ . Thus,  $\frac{\partial J_Y}{R_I} > 0$ .

- The negative likelihood ratio:  $NLR = \frac{1-R_I}{R_I}$  with I(tn) = 1, I(fp) = 0, I(fn) = 1, and I(tp) = 0. Thus,  $\frac{\partial NLR}{R_I} < 0$ . The positive likelihood ratio:  $PLR = \frac{R_I}{1-R_I}$  with I(tn) = 0, I(fp) = 1, I(fn) = 0, and I(tp) = 1. Thus,  $\frac{\partial PLR}{R_I} > 0$ . The probability of the elementary event *true negative*:  $PTN = \pi_-R_I$  with I(tn) = 1, I(fp) = 1, I(fn) = 0, and I(tp) = 1. Thus,  $\frac{\partial PTR}{R_I} > 0$ . The probability of the elementary event *true positive*:  $PTP = \pi_+R_I$  with I(tn) = 0, I(fp) = 0, I(fn) = 1, and I(tp) = 1. Thus,  $\frac{\partial PTP}{R_I} > 0$ .