

Visual-Based Analysis of Classification Measures with Applications to Imbalanced Data

Dariusz Brzezinski*, Jerzy Stefanowski, Robert Susmaga, Izabela Szczęch

*Institute of Computing Science, Poznan University of Technology,
ul. Piotrowo 2, 60-965 Poznan, Poland*

Abstract

With a plethora of available classification performance measures, choosing the right metric for the right task requires careful thought. To make this decision in an informed manner, one should study and compare general properties of candidate measures. However, analysing measures with respect to complete ranges of their domain values is a difficult and challenging task. In this study, we attempt to support such analyses with a specialized visualization technique, which operates in a barycentric coordinate system using a 3D tetrahedron. Additionally, we adapt this technique to the context of imbalanced data and put forward a set of properties which should be taken into account when selecting a classification performance measure. As a result, we compare 22 popular measures and show important differences in their behaviour. Moreover, for parametric measures such as the F_β and $IBA_\alpha(G\text{-mean})$, we analytically derive parameter thresholds that change measure properties. Finally, we provide an online visualization tool that can aid the analysis of complete domain ranges of performance measures.

Keywords: classification, performance measures, visualization, barycentric system, class imbalance

1. Introduction

Classification is one of the most important machine learning tasks, commonly applied to many real-world problems. One of the crucial ingredients of this supervised learning task is the selection of a performance measure that allows the user to discern good classifiers from bad ones. An appropriate measure should support choosing the best classifier among several candidates and help tune its parameters. As a result, the selected performance measure is responsible for the optimization of the learning process [8].

Although researchers often consider performance measures that promote predicting correctly the highest number of instances, many applications require other ways of handling errors referring to particular subsets of examples. This is especially true for imbalanced data [18, 23], where classifiers are biased towards the majority classes yet the under-represented minority class is usually of more value to human experts.

Since typical performance measures, such as classification accuracy, are not appropriate for imbalanced data [10, 27], several more relevant measures have been considered. The most popular ones include *precision*, *recall* (*sensitivity*), *specificity*, and their aggregates, e.g. *G-mean* or *F₁-score*. These and other measures for imbalanced data are typically defined on the basis of confusion matrices summarizing the predictions of a binary classifier. Looking into the related studies, one can notice that the number of such measures is relatively high and each represents different aspects of classification performance, often leading to quite different interpretations [20]. This shows that there is no one measure that would be the best choice in

*Tel.: +48 61 665 30 57

Email addresses: `dariusz.brzezinski@cs.put.poznan.pl` (Dariusz Brzezinski), `jerzy.stefanowski@cs.put.poznan.pl` (Jerzy Stefanowski), `robert.susmaga@cs.put.poznan.pl` (Robert Susmaga), `izabela.szczuch@cs.put.poznan.pl` (Izabela Szczęch)

all situations. However, which measure is used in a given problem seems to be, to a large extent, dictated simply by the measure’s popularity rather than a thorough discussion of its properties.

Although there are a few systematic studies on different properties of classifier performance measures [19, 16, 11, 30], we still postulate the need for thorough analysis of the measures’ behaviour. In particular, methods for: interpreting and comparing measures with respect to whole domain ranges, analysing their nature for different class and prediction distributions, and detecting the presence of unusual values are much needed. Theoretical investigations of these aspects are often very laborious and time consuming, especially when multi-dimensional aspects, provided by the confusion matrices, need to be taken into account. Due to these difficulties, such an analysis could be alternatively carried out with visual techniques to aid the understanding and interpretability of various measure properties.

In this paper, we put forward a new visualization technique for analysing entire domains of classification performance measures, which depicts all possible configurations of predictions in a confusion matrix, regardless of the used classifier. For this purpose, we adapt an approach originally created for rule interestingness measures to the context of classification [31]. Contrary to existing performance measure visualizations, such as ROC space [11], the proposed approach presents measures in a space which is defined directly on elements of the confusion matrix, is easily interpretable in 3D, and remains defined for all elements of the domain. Moreover, based on the devised visualization, we propose ten properties which should be taken into account while selecting evaluation measures, particularly for class imbalanced data. Consequently, we compare 22 popular classifier performance measures (both non-parametric and parametric) and highlight important differences in their behaviour. Finally, we demonstrate that the proposed approach can lead to concrete results by deriving property thresholds for the parametrized F_β and $IBA_\alpha(G\text{-mean})$ measures.

The main contributions of our paper are as follows:

- In Section 3, we adapt a technique for visualizing classification performance measures using the barycentric coordinate system and discuss its characteristics. Additionally, we present an online tool that implements the proposed technique and allows for the analysis of several predefined and custom user-defined 4D measures.
- In Section 4, we put forward ten properties, providing knowledge on the behaviour of the classifier performance measures for class biased problems. The introduced properties involve analysing maxima, minima, elements of symmetry, monotonicity, and undefined values.
- In Section 5.1, using the proposed visualization technique we analyse and compare 22 classification measures with respect to the proposed properties.
- In Section 5.2, we analyse how the proposed properties change for parametric measures. More precisely, we study the effect of internal parametrization on the F_β measure and external parametrization for $IBA_\alpha(G\text{-mean})$. Apart from visual inspection, we analytically derive threshold parameter values for the selected measures.
- In Section 6, we discuss the most important issues in analysing classification performance measures and draw lines of further investigations.

2. Related Works

2.1. Classifier performance measures

Classifiers can be assessed in many aspects, such as their predictive ability, training time, memory usage, model complexity, interpretability, or other criteria [20]. In this paper, we consider predictive performance only and focus on measures that evaluate crisp binary classifier predictions; measures specific to only rankers or probabilistic classifiers are out of the scope of this study. Furthermore, we concentrate mainly on measures which take into account the binary class imbalance problem.

As discussed in [18], when dealing with imbalanced data measures should focus on the more interesting minority class. Such measures are defined as functions of the confusion matrix for two-class problems, with

the minority class typically referred to as *positive* (P), while the remaining majority class as *negative* (N) [20, 17] (multiple non-positive classes, if present, are usually aggregated into one).

Table 1: Confusion matrix for two-class classification

Actual \ Predicted	Positive	Negative	total
	Positive	TP	FN
Negative	FP	TN	N
total	\hat{P}	\hat{N}	n

Table 1 illustrates a two-class confusion matrix, which may be regarded as a special case of a contingency table that can be multi-dimensional in general. The TP (*True Positive*) and TN (*True Negative*) entries denote the number of examples classified correctly by the classifier as positive and negative, while the FN (*False Negative*) and FP (*False Positive*) indicate the number of misclassified positive and negative examples, respectively. Based on these values, the most common performance measures are defined as:

$$accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (1) \qquad precision = \frac{TP}{TP + FP} \quad (2)$$

$$specificity = \frac{TN}{FP + TN} \quad (3) \qquad sensitivity (recall) = \frac{TP}{TP + FN} \quad (4)$$

Many other classification performance measures were proposed based on values from the confusion matrix; for their reviews see [18, 20, 19, 16, 2]. In this study, we analyse the properties of 22 measures, listed and defined in the supplementary material.¹ Below, we highlight four measures, chosen for diversity of their characteristics, which we will analyse and compare in more detail:

$$F_\beta = \frac{(1 + \beta) \cdot precision \cdot recall}{\beta \cdot precision + recall}, \text{ where } \beta \geq 0, \quad (5)$$

$$G\text{-mean} = \sqrt{sensitivity \cdot specificity}, \quad (6)$$

$$MCC = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{\hat{P} \cdot P \cdot N \cdot \hat{N}}}, \quad (7)$$

$$OP = accuracy - \frac{|specificity - sensitivity|}{specificity + sensitivity}. \quad (8)$$

The F_β combines *precision* and *recall* as a weighted harmonic mean, with the β parameter as their relative weight. Commonly $\beta = 1$ and then the measure is referred to as F_1 -score. G -mean [21] is the geometric mean of *sensitivity* and *specificity*, which takes into account the relative balance of recognition of both positive and negative classes. The *Matthews Correlation Coefficient* (MCC) expresses a correlation between the actual and predicted classification and returns a value between -1 (total disagreement) and $+1$ (perfect classification). We highlight MCC in our study as it was considered by some authors as one of the recommended measures for imbalanced data [2, 3]. *Optimized precision* (OP) combines *sensitivity* and *specificity* in a more complex way, also producing values in the $[-1, +1]$ range [28].

Apart from these “closed-formula” measures, we shall also analyze in more detail a representative of what may be thought of as “open-formula” measures, in this case $IBA_\alpha(M)$. This particular measure-wrapper is aimed at applying more weight to minority class predictions in a given measure M , according to a user-defined parameter α [13, 14].

These and other measures were compared in such surveys as e.g. [18, 16, 17, 2], however usually with respect to discussing the main differences in their definitions. Additionally, the F_1 -score was thoroughly analysed by Powers [26] who claimed that some of its properties, such as focusing only on the minority

¹<https://dabrze.shinyapps.io/Tetrahedron/>

class and assuming that actual and predicted distributions are identical, may be critical flaws. Another theoretical study showed that aggregating *sensitivity* and *specificity* presented more suitable behaviour than measures aggregating *precision* and *recall* [19]. Nevertheless, theoretical analyses of measures with respect to complete ranges of domain values are very laborious and have been done only for a few classifier performance measures.

2.2. Visualization of measures

In this paper, we focus on visualizing measures defined on a binary confusion matrix. We note that this should not be confused with visualizations of classifier performance, e.g. using ROC graphs [9], precision-recall curves [7], lift charts [25], or other attempts to graphically present experimental comparisons of classifiers [33, 1, 5]. Our intention is to study general properties of measures rather than visualize the predictive performance of a classifier on a given dataset.

The 3D visualizations of 2×2 sum-constrained matrices, applicable in particular to confusion matrices, have already been considered in different papers. Below, we recapitulate shortly three approaches, which bear some relation to the (regular) tetrahedron visualization used throughout this paper [22, 6, 11].

Le Bras et al. [22] introduce a system of 3D spaces (referred to as the *Formal Framework*), in which the contents of sum-constrained 2×2 matrices can be represented. Because of the three actual degrees of freedom of a sum-constrained 2×2 matrix, domains consisting of three variables are required and sufficient to express the matrix entries. However, the choice of a particular domain, with three particular variables, may vary depending on the application at hand.

While the representations with three variables might be used to produce 3D visualizations of measures, the paper of Le Bras et al. [22] does not exploit this fact in too much a detail, as its focus lies elsewhere. The authors introduce three very particular, application-driven, 3D domains referred to as: *confidence*, *examples* and *counterexamples*. In its central part, the paper recalls 38 measures related to association rules and defines them consistently in terms of the matrix entries, as well as in terms of the three proposed domains. This allows for conducting dedicated analyses of the measures (e.g. expressing the Piatetsky-Shapiro recommendations [24] in the *examples* domain), with the main objective of identifying measures most relevant to association rule pruning. The introduced and in detail scrutinized properties include: all-monotonicity, generalized universal existential upward closure, and opti-monotonicity [22].

As far as the tetrahedron-based visualization is concerned, the *examples* and *counterexamples* 3D spaces introduced in [22] assume the shapes of tetrahedra. However, contrary to the approach presented in our paper, the domains are designed for analysing rule interestingness measures. Moreover, the tetrahedra of the *examples* and *counterexamples* domains are irregular, since these domains are assumed to have two orthogonal variables each, implying shapes with two orthogonal edges incident with one vertex, a feature unattainable in the regular tetrahedron.

Celotto [6] has introduced 2D visualization spaces that are very natural to the considered measures, i.e., Bayesian confirmation measures. The primary space, suitably referred to as the *confirmation space*, consists of: $P(H|E)$ (x -axis) and $P(H)$ (y -axis). As noted within the paper, the 2D representation of 2×2 sum-constrained matrices is incomplete, and thus aptly called a *fingerprint* of the measure. The incompleteness results from the fact that the fingerprint changes as some third parameter which defines the third dimension, in this case chosen to be $P(E)$, is varied.

The confirmation space is initially set side by side with its analogue, denoted as *dual confirmation space*, and another 2D space, i.e. the ROC space, which consists of false positive rate $fpr = FP/N$ (x -axis) and true positive rate $tpr = TP/P$ (y -axis). However, because confirmation measures remain the main focus of the study of Celotto [6], presented analyses are basically confined to the confirmation space and its dual, which are used to analyse 19 measures. The measure analyses, principally concerned with identifying measures most relevant to classification rule pruning, include visualizations of some ordinal equivalence aspects and a multitude of symmetry aspects. The latter also include visually-assisted design and synthesis of measures possessing desired symmetries.

The 2D confirmation spaces introduced by Celotto correspond to rectangular cross-sections of the 3D tetrahedron presented in this paper. However, contrary to the presented approach, in [6] these originally

non-independent variables are presented as orthogonal and of unified ranges, which thus requires some amount of orthonormalization.

Flach [11] mentions several possible definitions of variables suitable for 3D visualizations of 2×2 confusion matrices, but focuses primarily on *3D ROC space*, a generalization of traditional 2D ROC space [20]. The 3D ROC space consists of the false positive rate $fpr = FP/N$ (x -axis) and true positive rate $tpr = TP/P$ (y -axis), which basically constitute traditional ROC space, together with the frequency of positives $pos = P/n$ (z -axis). This choice had been dictated by the general topic of the paper, which was the analysis of classifier performance measures and their behaviour in ROC spaces. Notice that the three variables are selected so that the resulting XY -plane hosts the ROC space, while the third co-ordinate varies with the actual class distribution. In result, the 3D ROC space is thus a collection of stacked-up ROC spaces, with the z -coordinate corresponding to the proportion of the positive class. Owing to the variable mutual orthogonality and similar ranges ($[0, 1]$ for x and y and $(0, 1)$ for z) the total domain shape is thus a $[0, 1] \times [0, 1] \times (0, 1)$ pseudo-cube, i.e. a cube with both the lowermost layer, corresponding to $FN + TP = 0$, and the uppermost layer, corresponding to $FP + TN = 0$, removed.

In the cited study [11], Flach combines the proportion of classes with misclassification costs, generally referred to as skew, and focuses on analysing 8 selected measures in terms of sensitivity to skew. The considered key notions involve: skew-equivalence and weak/strong skew-insensitivity of the measures. We also note that a similar techniques have been used to analyse rule quality measures. The most well known are coverage spaces, introduced by Fürnkranz and Flach [12], which plot the number of positive training examples and negative ones covered by the rule in the given data. Coverage spaces can be considered similar to ROC spaces in analysing isometrics of evaluation measures.

The stacked 2D spaces considered by Flach [11] basically correspond to the rectangle-shaped cross-sections of the tetrahedron presented in this paper. However, ROC spaces are presented in square form, which requires some amount of orthogonal rescaling compared to the approach presented in this paper. Furthermore, contrary to the visualization technique introduced in this paper, 3D ROC space remains undefined for confusion matrices with $FN + TP = 0$ or $FP + TN = 0$.

3. The barycentric visualization technique

As presented in Table 1, a confusion matrix for binary classification consists of four entries: TP , FP , FN , TN . However, for a dataset of n examples these four entries are constrained, as $n = TP + FP + FN + TN$. Therefore, for a given constant n , any three values in the confusion matrix uniquely define the fourth value. This property allows to visualize any classification performance measure based on the two-class confusion matrix using a 4D barycentric coordinate system [35].

In the *barycentric coordinate system* point locations are specified relatively to vertices of a simplex (a triangle, tetrahedron, etc.). A 4D barycentric coordinate system is a tetrahedron, where each dimension is represented as one of the four vertices. Choosing vectors that represent TP , FP , FN , TN as vertices of a regular tetrahedron in a 3D space, one arrives at a barycentric coordinate system as in Fig. 1.

In this system, every confusion matrix $\begin{bmatrix} TP & FN \\ FP & TN \end{bmatrix}$ is represented as a point of the tetrahedron. Let us illustrate this fact with a few examples. Figure 1 shows a skeleton of a tetrahedron with 4 exemplary points:

- one located in vertex TP, which represents $\begin{bmatrix} n & 0 \\ 0 & 0 \end{bmatrix}$,
- one located in the middle of edge TP–FP, which represents $\begin{bmatrix} n/2 & 0 \\ n/2 & 0 \end{bmatrix}$,
- one located in the middle of face $\triangle TP$ – FP – FN , which represents $\begin{bmatrix} n/3 & n/3 \\ n/3 & 0 \end{bmatrix}$,
- one located in the middle of the tetrahedron, which represents $\begin{bmatrix} n/4 & n/4 \\ n/4 & n/4 \end{bmatrix}$.

One way of understanding this representation is to imagine a point in the tetrahedron as the centre of mass of the examples in a confusion matrix. If all n examples are true positives, then the entire mass of the predictions is at TP and the point coincides with vertex TP. If all examples are false negatives, the point lies on vertex FN, etc. Generally, whenever $a > b$ ($a, b \in \{TP, FN, FP, TN\}$) then the point is closer to the vertex corresponding to a rather than b .

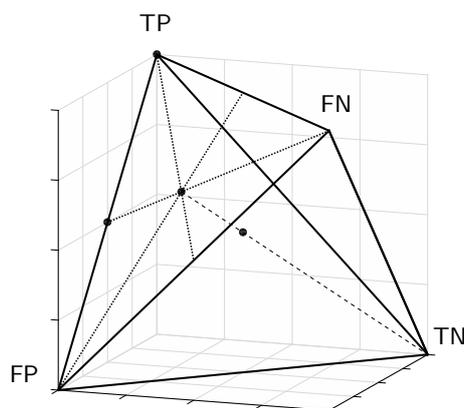


Figure 1: A skeleton visualization of the tetrahedron with four exemplary points

Using the barycentric coordinate system makes it possible to depict the originally 4D data (two-class confusion matrices) as points in 3D. Moreover, as in [31, 32], an additional variable based on the depicted four values may be rendered as colour. Although any colour map can be used, in the following paragraphs we utilize the map shown in Fig. 2: dark blue — minimum values, to dark brown — maximum values. Areas of the same colour signify then the same values of the variable. The shape of such areas is determined by the nature of the visualized variable and usually occurs as lines in 2D (isolines) and surfaces in 3D (isosurfaces). Undefined values of the measures will be rendered in magenta, i.e., a colour not occurring in the map.

Here, we adapt this procedure to colour-code the values of classification performance measures, which remain the principal focus of this paper. In this respect, the presented approach is different from [31] and [32], in which Bayesian confirmation measures were mainly addressed. In particular, this paper introduces and discusses those aspects of the tetrahedron-based visualization that are especially useful for the analysis of classification performance measures.



Figure 2: The color map

The described visualization technique has been implemented as an interactive web application, available at: <https://dabrze.shinyapps.io/Tetrahedron/>. The application can visualize 86 predefined 4D measures, including the 22 classification performance measures described further. The user can also visualize custom measures by providing their formulae. For the remainder of the paper, the reader is encouraged to use this tool to interactively analyse the described properties of various classification measures.

Since classification *accuracy* is one of the simplest and most often used performance measures, let us use it for an exemplary visualization in Fig. 3. Its values range from 0 to 1, and there are no undefined ones. One can notice that confusion matrices with a high number of *FP* and *FN* result in low *accuracy*, whereas high *TP* and *TN* yield high *accuracy*. The visualization in Fig. 3a is only partially comprehensive, as it only shows the externals of the tetrahedron which correspond to very specific confusion matrices. However, both external as well as internal areas can be shown, e.g. by padding tetrahedron points (Fig. 3b), using “under the skin” views (Fig. 3c) or performing cross-sections (Fig. 4).

The indicated cross-sections are of particular interest in the context of analysing measures for class imbalance problems. Notice that traversing the tetrahedron alongside the vertical axis (up-down in Fig. 4a) corresponds to changing the proportions between sums $TP + FN = P$ and $FP + TN = N$, which specify the cardinalities of the actual classes. If $P = N$, then a situation of balanced classes is reproduced; otherwise the classes are imbalanced.

How a measure behaves for a particular class proportion may be visualized by producing a cross-section

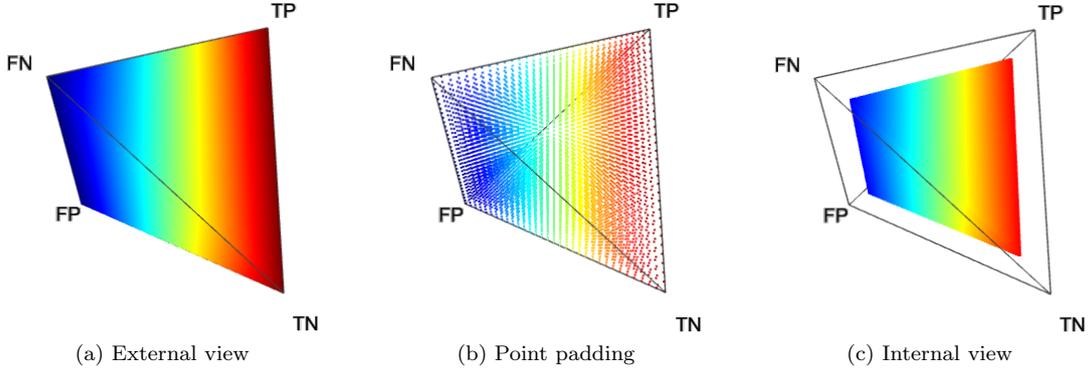


Figure 3: Visualizations of classification *accuracy*

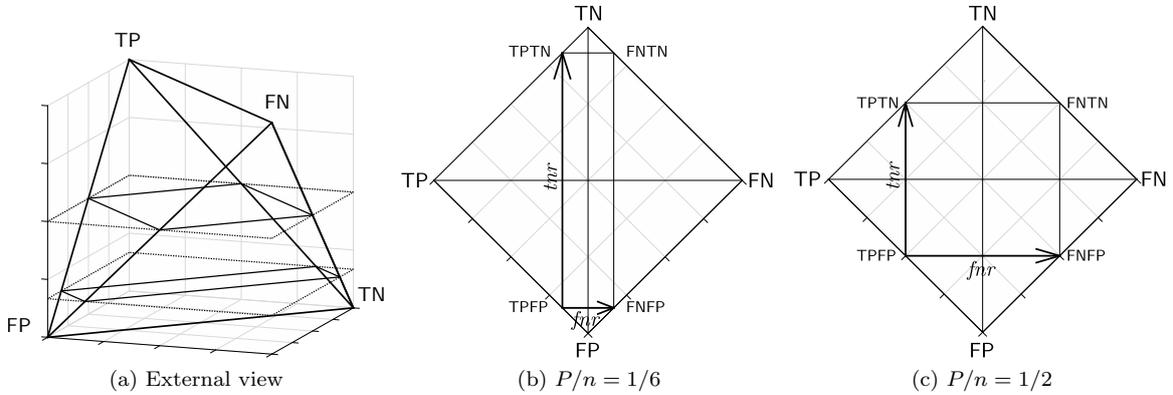


Figure 4: Skeleton visualizations of the tetrahedron and top-down view depictions of rectangular cross-sections for two selected values of positive class rate (P/n)

of the tetrahedron with a horizontal plane that cuts its vertical height. Figures 4b and 4c show the two cross-sections visible in Fig. 4a), one at $P/n = 1/6$ (positive class as the minority class) and one at $P/n = 1/2$ (class balance), as seen from above the tetrahedron. Cutting the shape with a horizontal plane at $P/n = 1/6$ produces the lower, rectangular cross-section (4b), while at $P/n = 1/2$ — the upper, square one (4c). In their corresponding figures, the cross-sections are oriented so that their sides incident with face $\triangle TP-FP-FN$ of the tetrahedron are positioned at the bottom, while those incident with face $\triangle TP-FN-TN$ — at the top. It is additionally worth noting that at every section the proportion of the rectangle's side lengths follows that of P (the horizontal side) and N (the vertical side), i.e. the class cardinalities.

Accordingly to the notation of the vertices of the tetrahedron, the sides and vertices of a cross-section rectangle are labelled as follows:

- sides: \overline{TP} (left), \overline{TN} (upper), \overline{FN} (right), \overline{FP} (lower),
- vertices: TPTN (upper-left), FNTN (upper-right), FNFP (lower-right), TFPF (lower-left).

The two axes, fnr and tnr , of the 2D space in which all cross-sections are represented (including those for $P/n = 1/6$ and $P/n = 1/2$), correspond to the false negative rate, $fnr = \frac{FN}{FN+TP} = 1 - recall$, and the true negative rate, $tnr = \frac{TN}{TN+FP} = specificity$. The orientation of the axes results from the fact that traversing the rectangle left-to-right corresponds to increasing fnr from 0 to 1, whereas traversing the rectangle down-up corresponds to increasing tnr from 0 to 1. The resulting 2D space of the presented cross-section is thus

an analogue of 2D ROC space, where, somewhat reversely, the false positive rate, $fpr = \frac{FP}{FP+TN} = 1 - specificity$, and the true positive rate, $tpr = \frac{TP}{FP+TP} = recall$, are used as x and y axes, respectively.

The presented rectangular cross-sections and 2D ROC space constitute the same, though seen from different angles, cross-sections of the tetrahedron. However, contrary to 3D ROC space [11], the presented technique does not involve any non-linear transformations of the elements of the confusion matrix and remains defined for all elements of the domain. Furthermore, because the proposed barycentric coordinates directly correspond to elements of the confusion matrix, the visualization is easily interpretable also in 3D, which helps analysing the whole range of possible domain values.

In the following sections, we demonstrate the usage of the visualization technique in some analyses of the considered classifier performance measures for imbalanced data. The technique, including the cross-sections, was particularly used to visualize several postulated properties of the measures.

4. Properties of Measures for Imbalanced Data

With a visualization technique at hand, it is much easier to define and interpret potentially desirable measure properties. In this section we put forward and discuss ten properties designed to highlight characteristic features of classifier performance measures designed for imbalanced data. The proposed properties can aid researchers in the selection of measures suitable for a given context and raise much needed discussion on the applicability of measures in certain domains.

Recall that the interpretation of the rectangular cross-section discussed in Section 3 is as follows.

- side $\overline{TP} / \overline{FN}$: full/null recognition of the positive class (Fig. 5),
- side $\overline{TN} / \overline{FP}$: full/null recognition of the negative class (Fig. 6).

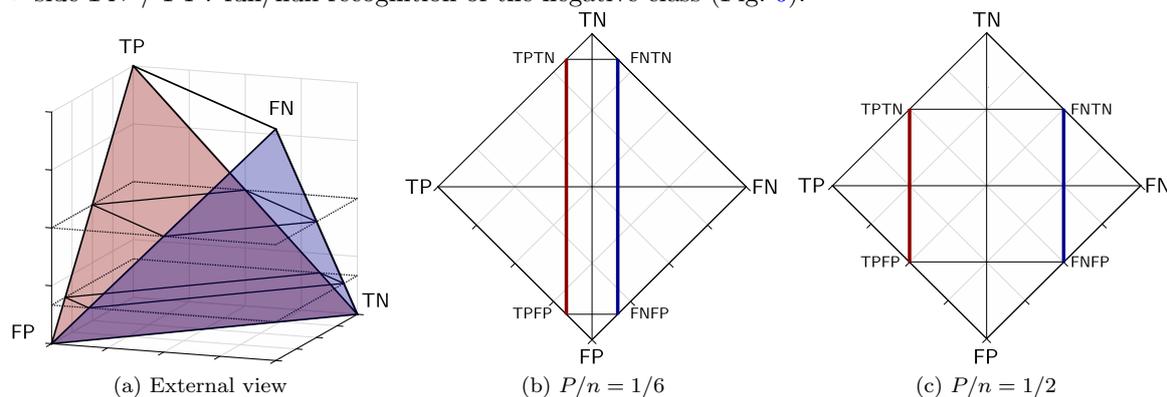


Figure 5: Illustration of full/null recognition of the positive class

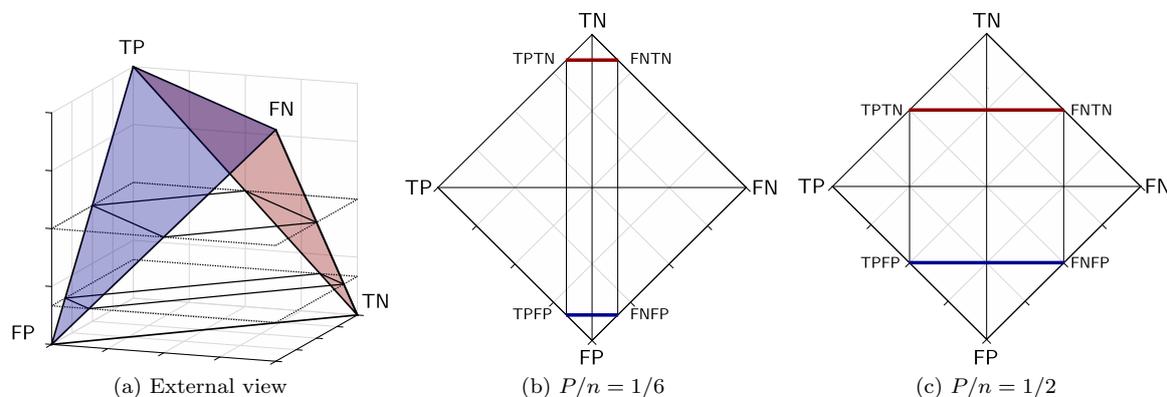


Figure 6: Illustration of full/null recognition of the negative class

In this context, we postulate to analyse classifier performance measures with respect to ten properties:

$TPTN_{max}$: vertex TPTN maximal value,

\overline{FN}_{min} : side \overline{FN} minimal value,

\overline{FP}_{min} : side \overline{FP} minimal value,

TP_{\nearrow} : horizontal lines weakly monotonic value growth (from \overline{FN} to \overline{TP}),

TN_{\nearrow} : vertical lines weakly monotonic value growth (from \overline{FP} to \overline{TN}),

$\overline{TN}_{\neq max}$: side \overline{TN} less than maximal value except for vertex TPTN,

$\overline{TP}_{\neq max}$: side \overline{TP} less than maximal value except for vertex TPTN,

ACE : for any two corresponding points on sides \overline{TP} and \overline{TN} (e.g. middle points) the value on side \overline{TP} is greater or equal to that on \overline{TN} ,

ACH : values invariant under exchange of TP with TN and FN with FP ,

$UnDefs$: the existence (and the location) of undefined values.

If present, undefined measure values are excluded from the above considerations, except for the last property, which is directly concerned with those values. Similarly, all but the last two properties are analysed only for ‘non-degenerated’ rectangular cross-sections, i.e. cross-sections corresponding to $P > 0$ and $N > 0$. On the other hand, the ‘degenerated’ cross-section, i.e. cross-sections that result in rectangles of either zero breadth or zero width, are taken into account only in the ACH and $UnDefs$ properties. The presented properties may be regarded as a basic ‘check-list’, providing knowledge on the behaviour of classifier performance measures for imbalanced data.

Notice that when all feasible rectangular cross-sections of the considered type are taken into account, the properties naturally extend from 2D in the rectangles to 3D in the tetrahedron. For example, points TPTN of all rectangles form edge \overline{TP} – \overline{TN} of the tetrahedron, sides \overline{TP} of all rectangles form face $\triangle FP$ – \overline{TP} – \overline{TN} of the tetrahedron, etc. This multidimensional nature of the measures renders the analytical process of their property verification harder, emphasizing the usefulness of the introduced visual-based 3D analyses.

Recall that the analysed measures are functions of $TP \geq 0$, $FN \geq 0$, $FP \geq 0$ and $TN \geq 0$, $TP + FN + FP + TN = n$, which constitute the elements of the confusion matrix $\begin{bmatrix} TP & FN \\ FP & TN \end{bmatrix}$ (see Table 1). In this context, $f(\begin{bmatrix} TP & FN \\ FP & TN \end{bmatrix})$ denotes the value of any of the considered classification performance measures.

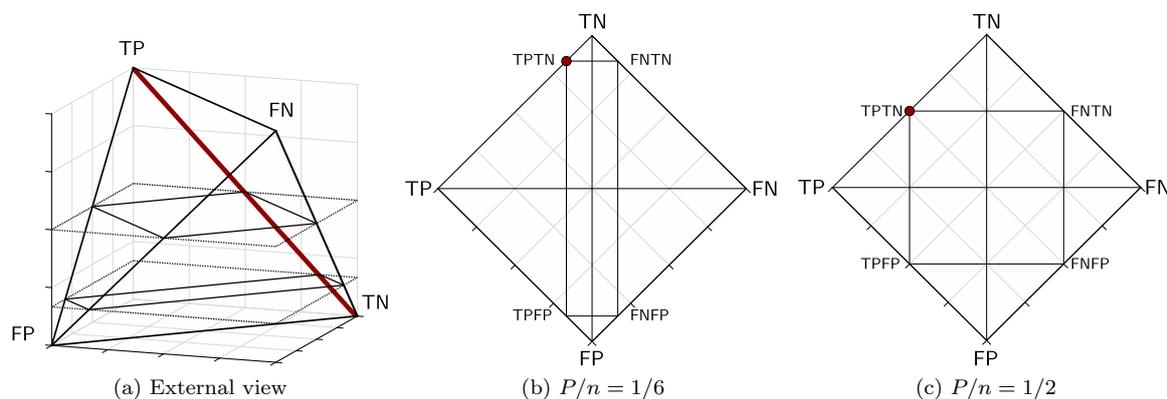


Figure 7: Illustration of property $TPTN_{max}$

Property $TPTN_{max}$ ensures that perfect predictions of both classes always render the best measure value (see Fig. 7). Notice that vertex TPTN, being the common part of both side \overline{TP} and side \overline{TN} , is actually

the only point of full recognition of both the positive and the negative class. Because TPTN corresponds to $\begin{bmatrix} P & 0 \\ 0 & N \end{bmatrix}$, this implies $f(\begin{bmatrix} P & 0 \\ 0 & N \end{bmatrix}) = \max$.

Properties \overline{FN}_{min} and \overline{FP}_{min} state that not recognizing one of the classes should correspond to the worst possible measure value (see Fig. 8). Recall that side \overline{FN} and side \overline{FP} correspond to null recognition of the positive and the negative class, respectively. In binary classification, a null recognition of any of the two classes (which concerns the minority class in most cases) is certainly insufficient. Thus, it is naturally required that measures should obtain minimal values on sides \overline{FN} and \overline{FP} . This boils down to $f(\begin{bmatrix} 0 & P \\ FP & TN \end{bmatrix}) = \min$ and $f(\begin{bmatrix} TP & FN \\ N & 0 \end{bmatrix}) = \min$.

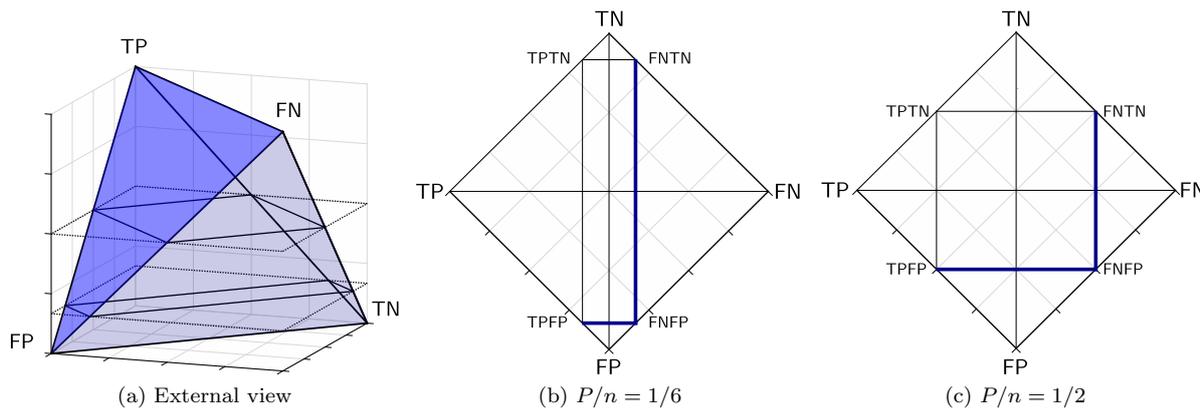


Figure 8: Illustration of properties \overline{FN}_{min} and \overline{FP}_{min}

Properties TP_{\nearrow} and TN_{\nearrow} require that growing TP and TN values should coincide with a weakly monotonic growth of the measure's value (see Fig. 9). As far as TP_{\nearrow} is concerned, observe that the greater TP is in the confusion matrix, the closer we move from side \overline{FN} to side \overline{TP} in the rectangular cross-section, which translates directly to increased recognition of the positive class. Naturally, it would be counter-intuitive if such increased recognition resulted in decreasing values of the measure. Thus, its weakly monotonic growth is expected. As opposed to requirements \overline{FN}_{min} and \overline{FP}_{min} , which concern merely the borders of the cross-section, TP_{\nearrow} concerns the entirety of the cross-section. In particular, also side \overline{FP} , where the value is required to be minimal (according to property \overline{FP}_{min}), satisfies the weak monotonicity. Property TP_{\nearrow} boils down to the following condition: if $TP_1 \geq TP_2$, then $f(\begin{bmatrix} TP_1 & FN_1 \\ FP & TN \end{bmatrix}) \geq f(\begin{bmatrix} TP_2 & FN_2 \\ FP & TN \end{bmatrix})$. Analogously property TN_{\nearrow} , which boils down to the condition: if $TN_1 \geq TN_2$, then $f(\begin{bmatrix} TP & FN \\ FP_1 & TN_1 \end{bmatrix}) \geq f(\begin{bmatrix} TP & FN \\ FP_2 & TN_2 \end{bmatrix})$.

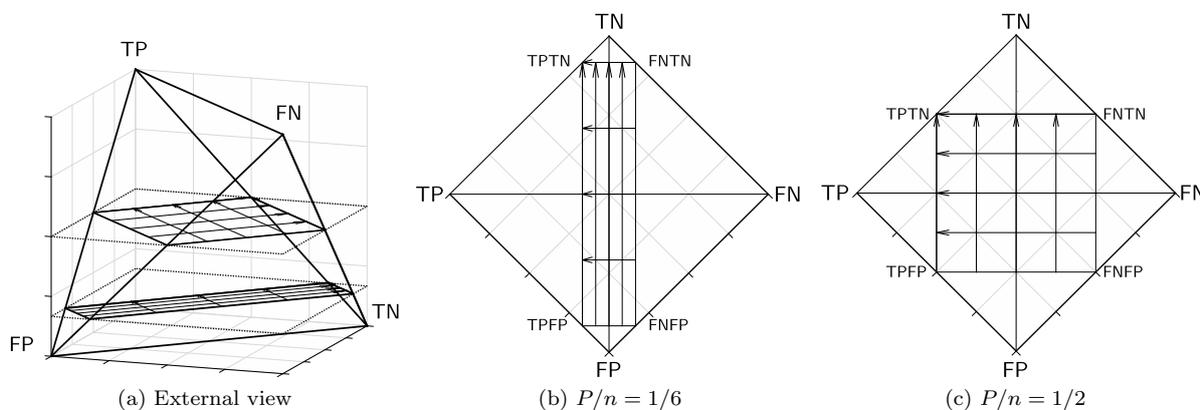


Figure 9: Illustration of properties TP_{\nearrow} and TN_{\nearrow}

Properties $\overline{TN}_{\neq max}$ and $\overline{TP}_{\neq max}$ tackle the problem of maximal values of the measure. Observe that in a two-class problem, the full recognition of just one class (only positive or only negative), which can be achieved trivially, should not render the highest value of the measure. Only the full recognition of both classes should be rewarded with the maximum, as stated by $TPTN_{max}$. Thus, properties $\overline{TN}_{\neq max}$ and $\overline{TP}_{\neq max}$ require that the measure's values on sides \overline{TN} and \overline{TP} should be less than maximal, except for the very vertex $TPTN$. If a classification measure fulfils this property, a simple majority or minority stub will never be mistaken with the best possible classifier. This boils down to: if $FN + FP > 0$, then $f(\begin{bmatrix} TP & FN \\ FP & TN \end{bmatrix}) < max$.

Property *ACE* reveals the class bias resulting from *asymmetric class evaluation*, typical for class imbalance problems. It is introduced to guarantee that full recognition of only the negative class is never rewarded with a higher value than the full recognition of only the positive one (assuming the respective other class is recognized to the same degree). In particular, since the recognition of the positive class is of high importance, the middle point of side \overline{TP} (i.e. when the whole positive and half of the negative class is recognized) should not be assessed with a lower value than the middle point of side \overline{TN} (i.e. when the whole negative and half of the positive class is recognized). Similarly for all other pairs of corresponding points on sides \overline{TP} and \overline{TN} (three of which are depicted in Fig. 10). In terms of the entries of the confusion matrix, property *ACE* boils down to: $f(\begin{bmatrix} P & 0 \\ \gamma N & (1-\gamma)N \end{bmatrix}) \geq f(\begin{bmatrix} (1-\gamma)P & \gamma P \\ 0 & N \end{bmatrix})$, where $\gamma \in [0, 1]$ (in Fig. 10 γ takes on values 1/4, 2/4 and 3/4). Notice that the weak nature of the property is implied by the fact that it does not specify by how much the full recognition of the positive class should be favoured over the full recognition of the negative class. On the other hand the unsatisfied *ACE* reveals instantly, however, that the measure favours (in the above sense) the negative over the positive.

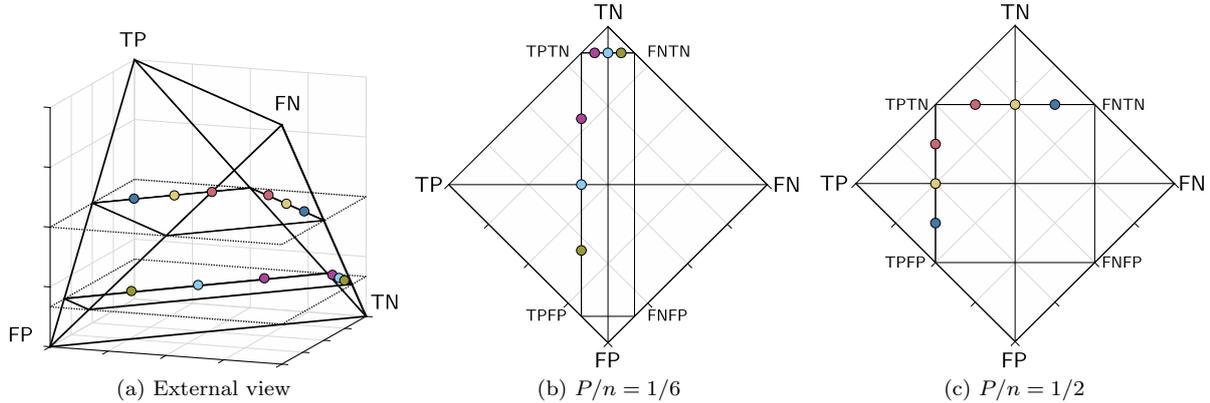


Figure 10: Illustration of property *ACE*

Much the same, property *ACH* deals with the issue of *asymmetric class handling*. It tests if the classes can be exchanged without influencing the measure's behaviour. This could be especially relevant in highly dynamic situations, e.g. in data streams plagued by concept drift [4], in which the percentage of the positive class may increase to make it actually (albeit temporarily) the majority class [34]. Expressed with the confusion matrix, *ACH* boils down to: $f(\begin{bmatrix} TP & FN \\ FP & TN \end{bmatrix}) = f(\begin{bmatrix} TN & FP \\ FN & TP \end{bmatrix})$.

Finally, property *UnDefs* pinpoints the existence and the location of undefined values ($+\infty$ whenever a positive value is divided by 0, $-\infty$ whenever a negative value is divided by 0, and NaN whenever 0 is divided by 0). As such, it highlights potential numerical pitfalls that can arise when calculating the measure values. While occurring fairly seldom with real life data, such undefined values are needed to fully characterize and thoroughly compare the considered measures.

In the following section, we use the proposed ten properties to compare various classification measures.

5. Visual-based Analysis of Selected Measures

Having presented the visualization technique in Section 3 and having defined the properties to be researched in Section 4, now we use the proposed tools to analyse 22 classification measures. The selected set of measures includes the most popular ones defined using elements from a two-class confusion matrix, and comprises non-parametric as well as parametric indices. Table 2 presents the analysis results for the selected measures, whereas their definitions are available in the supplementary data².

Table 2: Properties of selected classification measures; *: contains NaN (undefined value); †: NaN side, ^s: strong monotonicity

Measure	\overline{TP}_{max}	\overline{FN}_{min}	\overline{FP}_{min}	TP_{\nearrow}	TN_{\nearrow}	$\overline{TN}_{\neq max}$	$\overline{TP}_{\neq max}$	ACE	ACH	$UnDefs$
<i>Accuracy</i>	✓	×	×	✓ ^s	✓ ^s	✓	✓	×	✓	none
<i>Area Under Lift</i>	×	×	×	✓ ^s	✓ ^s	✓	✓	✓	×	TN-FP; TP-FN
<i>Balanced accuracy</i>	✓	×	×	✓ ^s	✓ ^s	✓	✓	✓	✓	TN-FP; TP-FN
<i>F₁-score</i>	✓	× [†]	×	✓ ^{s*}	✓ ^{s*}	✓ [*]	✓	×	×	△FP-FN-TN
<i>False negative rate</i>	×	×	×	×	✓	×	✓	×	×	TN-FP
<i>False positive rate</i>	×	×	×	✓	×	✓	×	✓	×	TP-FN
<i>F_β, β ∈ [0, ∞)</i>	✓	× [†]	×	✓ ^{s*}	✓ ^{s*}	✓ [*]	✓	×	×	△FP-FN-TN
<i>G-mean</i>	✓	✓	✓	✓	✓	✓	✓	✓	✓	TN-FP; TP-FN
<i>IBA_α(Accuracy), α ∈ (0, ∞)</i>	×	×	×	✓ ^s	×	✓	×	×	×	TN-FP; TP-FN
<i>IBA_α(F₁-score), α ∈ (0, ∞)</i>	×	× [†]	×	×	×	✓ [*]	×	×	×	△FP-FN-TN; TP-FN
<i>IBA_α(G-mean), α ∈ (0, ∞)</i>	×	✓	✓	✓	×	✓	×	✓	×	TN-FP; TP-FN
<i>IBA_α(F_β), α, β ∈ (0, ∞)</i>	×	× [†]	×	×	×	✓ [*]	×	×	×	△FP-FN-TN; TP-FN
<i>Jaccard coefficient</i>	✓	✓	×	✓ ^s	✓	✓	✓	×	×	TN
<i>Kappa</i>	✓	×	×	✓ ^s	✓ ^s	✓	✓	×	×	TN; TP
<i>Log odds-ratio</i>	✓	✓ [*]	✓ [*]	✓ [*]	✓ [*]	×	×	× [†]	✓	TN-FN; TN-FP; TP-FN; TP-FP
<i>MCC</i>	✓	×	×	✓ ^{s*}	✓ ^{s*}	✓ [*]	✓ [*]	×	✓	TN-FN; TN-FP; TP-FN; TP-FP
<i>Neg. predictive value</i>	✓	×	✓ [*]	✓ [*]	✓ [*]	✓	×	✓	×	TP-FP
<i>OP</i>	✓	×	×	×	×	✓	✓	×	✓	TN-FP; FP-FN; TP-FN
<i>Pointwise AUC-ROC</i>	✓	✓	✓	✓	✓	✓	✓	✓	✓	TN-FP; TP-FN
<i>Precision</i>	✓	✓ [*]	×	✓ [*]	✓ [*]	×	✓	×	✓	TN-FN
<i>Recall</i>	✓	✓	×	✓ ^s	✓	✓	×	✓	×	TN-FP
<i>Specificity</i>	✓	×	✓	✓	✓ ^s	×	✓	×	×	TP-FN

Looking at the entries of Table 2, one can notice that the proposed properties clearly differentiate the analysed measures. Having realized the differences in the measures' behaviour, one can more accurately choose the measures for the application at hand.

Let us start with having a closer look at one exemplary property listed in Table 2, i.e. the existence and location of undefined values ($UnDefs$). The undefined measure values, usually resulting from division by zero, and commonly neglected, may well occur with imbalanced data, e.g. during unstratified cross-validation procedures when one of the two classes happens to be unrepresented in the learning or the testing set. The problem becomes aggravated for multi-class problems when the measure is macro-averaged for all classes, since the resulting average becomes undefined if at least one of the averaged values is undefined.

²<https://dabrze.shinyapps.io/Tetrahedron/>

An interesting observation is that, except for *accuracy*, all of the considered measures contain undefined values. In particular, the *Kappa* statistic is undefined when there exist only positive or only negative examples in the dataset and none of them is misclassified, which translates to two different locations in the tetrahedron, namely vertex TP and vertex TN. Even worse, *balanced accuracy* is undefined when there are only positive or only negative examples in the dataset, which translates directly to whole edges TP–FN and TN–FP in the tetrahedron. Worst of all, *F₁-score* (as well as its generalizations) exhibits undefined values in the whole face $\triangle FP\text{--}FN\text{--}TN$, which occurs when all positive examples are misclassified (even when both classes are represented).

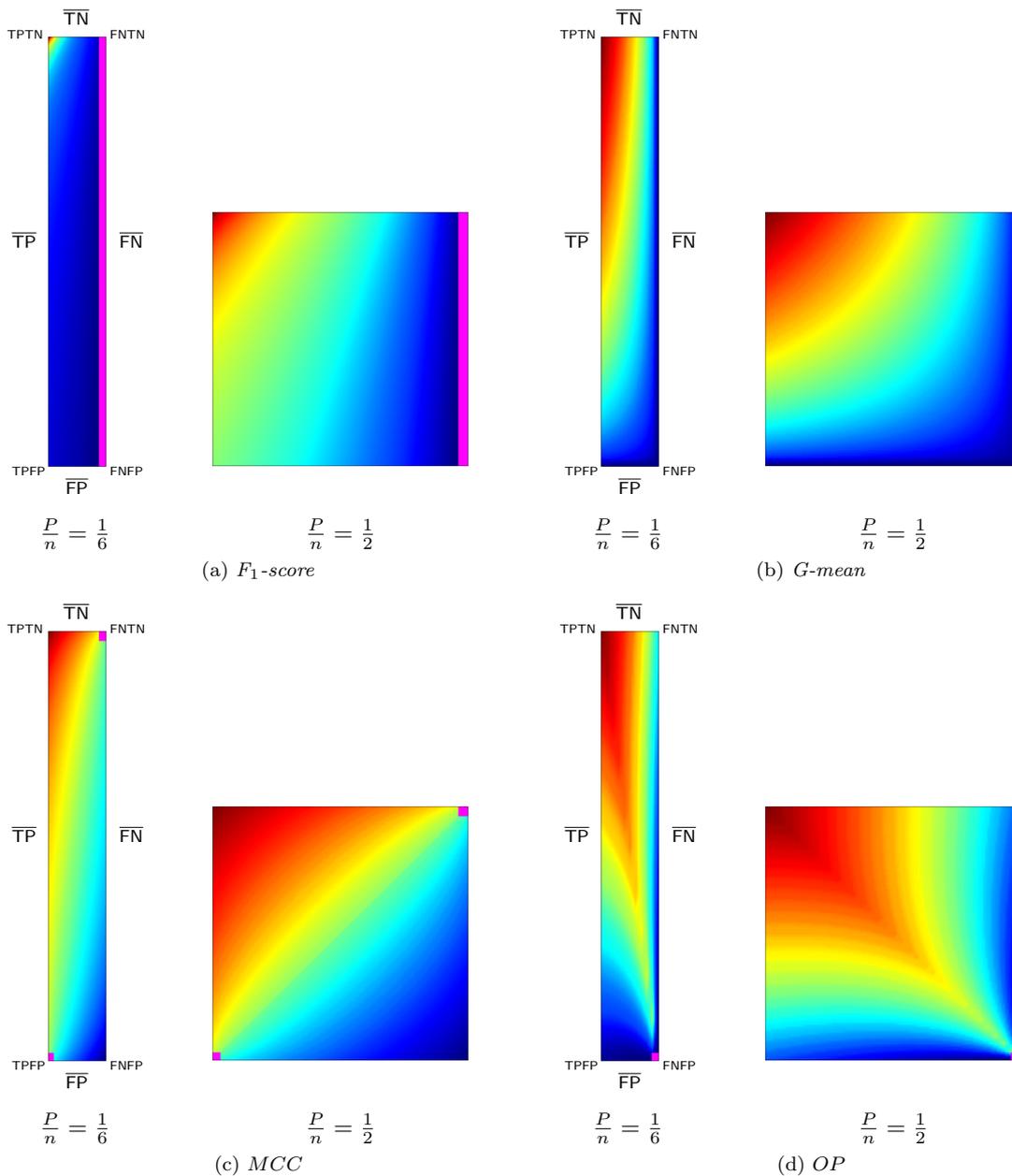


Figure 11: Cross-sections of selected measures for $P/n = 1/6$ and $P/n = 1/2$

5.1. Non-parametric Measures

Let us now conduct a more detailed visual analysis of the four highlighted measures from Section 2: F_1 -score, G -mean, Mathews Correlation Coefficient (MCC), and $Optimized\ precision$ (OP), putting particular emphasis on their behaviour with respect to imbalanced data. Due to the page limit, in this paper we present only cross-sections produced for $P/n = 1/6$ and $P/n = 1/2$. However, other cross-sections of the tetrahedron, including cross-sections produced for higher levels of class imbalance, can be viewed in the online visualization tool.

The analysis of F_1 -score, visualized in Fig. 11a, show that the growth (although monotonic) of the measure along side \overline{TP} is very slow and does not fulfil the ACE property when the data are imbalanced. To illustrate this, consider Fig. 11a (left), which corresponds to class imbalance, and a point located in the middle of \overline{TP} . The value there is much lower than the corresponding point on side \overline{TN} . Taking into account the fact that the middle point of \overline{TP} corresponds to full recognition of the positive class and 50% recognition of the negative class, this shows that with class imbalance high values corresponding to full recognition of the positive class are harder to obtain. Expressed in terms of values in the confusion matrix: $f(\begin{bmatrix} P & 0 \\ \gamma N & (1-\gamma)N \end{bmatrix}) < f(\begin{bmatrix} (1-\gamma)P & \gamma P \\ 0 & N \end{bmatrix})$, where $\gamma = 1/2$. Notice that while F_1 -score fulfils ACE for $P/n = 1/2$ (11a (right)), it does not for the above mentioned $P/n = 1/6$ (11a (left)), which means that the property is not satisfied in general (i.e. throughout the tetrahedron). Evidently, the property cannot be verified using only one selected cross-section. As may be observed using the online tool (in particular, by animating P/n from 1/2 down to 0), this flawed feature of the measure aggravates for increasing class imbalance (i.e. when P/n drops). This may be quite surprising as the F_1 -score is often brought out in the literature as especially suited for the positive class. Generalizations of F_1 -score will be discussed in subsection 5.2 devoted to parametric measures.

The visual-based analysis of G -mean (Fig. 11b) reveals that the measure satisfies the devised properties. In particular, it satisfies some important properties not fulfilled by F_1 -score, MCC and OP . First, as opposed to the other three measures, G -mean features minimal values on whole sides \overline{FN} and \overline{FP} . Additionally, it enjoys the ACE property, which makes the measure especially useful in the contexts of imbalanced data: for any two corresponding points on sides \overline{TP} and \overline{TN} , the value on side \overline{TP} happens to be equal (and thus not smaller) to that on \overline{TN} . This means that for any γ : $f(\begin{bmatrix} P & 0 \\ \gamma N & (1-\gamma)N \end{bmatrix}) = f(\begin{bmatrix} (1-\gamma)P & \gamma P \\ 0 & TN \end{bmatrix})$.

As to the behaviour of MCC (Fig. 11c), one can observe that its values on sides \overline{FP} and \overline{FN} are not minimal, which violates properties \overline{FN}_{min} and \overline{FP}_{min} . Even worse, comparing cross-sections for $P/n = 1/2$ and $P/n = 1/6$, one can observe that small values are harder to obtain with the increase of class imbalance. Furthermore, similarly to the F_1 -score, MCC does not satisfy property ACE . Even though for balanced classes (Fig. 11c(right)) the corresponding points in \overline{TP} and \overline{TN} feature equal values, this deteriorates with growing disproportion between classes (Fig. 11c(left)). In other words, for imbalanced data it is easier to obtain undue high values by recognizing the negative class.

Finally, let us consider measure OP (Fig. 11d). The visual-based analysis reveals that OP is the only of the selected measures that does not satisfy properties TP_{\nearrow} and TN_{\nearrow} . Observe that traversing the cross-sections horizontally right-to-left or vertically bottom-up (thus increasing the recognition of one of the classes while keeping the recognition of the second one constant) the values of the measure first increase and then decrease. In fact, the visual analysis discloses that the measure is designed to increase its values monotonically only when the recognition of both classes increases. Undeniably, the increase of the recognition of both classes at the same time is highly desirable and should imply increasing measure values, however, the observed behaviour of OP in (acceptable) cases when the classifier increases the recognition of one class, while keeping the recognition of the other constant is rather surprising and counter-intuitive.

5.2. Parametric Measures

Recalling that the classifier performance measures are functions of the four entries of the confusion matrix, it may be observed that as far as their analytical forms are concerned, the various measures may be divided into unparametrized (e.g. G -mean measure) and parametrized (e.g. F_β measure). This parametrization process has been designed to lend the measures some amount of universality, as is the case with F_β , where the β parameter is supposed to control the class bias. As such control is much desired, external parametrization

procedures have also been developed to modify the measures' behaviour, e.g. by adapting them to problems with imbalanced data. One such procedure, called *Index of Balanced Accuracy* (IBA_α) [13, 14, 15], produces a parametrized measure, in which the α controls the amount by which the original measure is actually modified.

The above approaches allow us to focus on two following parametrization types:

- internal parametrization, (e.g. F_β),
- external parametrization, (e.g. $IBA_\alpha(G\text{-mean})$),

though also a kind of a simultaneous parametrization, e.g. $IBA_\alpha(F_\beta)$, is feasible.

Observe that measure parametrization actually increases the number of available degrees of freedom, making the inherently complex analyses of such measures even more challenging. The principal question is: how are the particular parameter values to be established? And further, what are their applicability ranges?

Procedures adapted to answer these questions vary from simple trial-and-error approaches to more intricate ones, in which parameter values are possibly gleaned from accessible data. In all cases visualization seems indispensable, providing valuable insights as to the measures' behaviour throughout their multidimensional, parametrized domains.

Let us now conduct a more detailed visual analysis of measures, representing both types of parametrization: F_β (internal) and $IBA_\alpha(G\text{-mean})$ (external), which illustrates the impact of the parametrization upon the measures' behaviour with respect to imbalanced data. Consistently, we present only cross-sections produced for $P/n = 1/6$ and $P/n = 1/2$, while other cross-sections as well as the entire tetrahedrons can be viewed in our online visualization tool.

5.2.1. Internal parametrization: F_β

While F_1 -score is a regular harmonic mean of *precision* and *recall*, F_β originated as a weighed version of this mean. In F_β λ and $1 - \lambda$ act as non-negative ($0 \leq \lambda \leq 1$) weights of *precision* and *recall*, respectively. This means that λ may be chosen to produce any convex combination of $\frac{1}{\text{precision}}$ and $\frac{1}{\text{recall}}$ to be actually used in the mean. Let p denote *precision* and r denote *recall*, the weighted harmonic mean of p and r is: $(\frac{\lambda \frac{1}{p} + (1-\lambda) \frac{1}{r}}{\lambda + (1-\lambda)})^{-1} = \frac{\lambda + (1-\lambda)}{\lambda \frac{1}{p} + (1-\lambda) \frac{1}{r}} = \frac{1}{\lambda \frac{1}{p} + (1-\lambda) \frac{1}{r}} = \frac{1}{\lambda \frac{r}{pr} + (1-\lambda) \frac{p}{pr}} = \frac{pr}{\lambda r + (1-\lambda)p} = \frac{\frac{1}{\lambda} pr}{r + \frac{1-\lambda}{\lambda} p}$ (from now on: $\lambda > 0$).

After setting³ $\beta = \frac{1-\lambda}{\lambda}$, one gets $\frac{1}{\lambda} = \beta + 1$, which finally produces: $F_\beta = \frac{(\beta+1)pr}{\beta p+r}$. Notice that in this scheme:

- $\lambda \rightarrow 0$ corresponds to $\beta \rightarrow \infty$ (emphasis on *precision*),
- $\lambda = 0.5$ corresponds to $\beta = 1.0$ (equal emphasis),
- $\lambda \rightarrow 1$ corresponds to $\beta \rightarrow 0$ (emphasis on *recall*).

Of course, for $\beta = 1.0$, measure F_β becomes $\frac{(1+1)pr}{1 \cdot p+r} = 2 \frac{pr}{p+r} = F_1$ -score, which is thus the regular (unweighed) harmonic mean of *precision* and *recall*.

The harmonic mean, used in this context happens to be the most conservative of the three popular Pythagorean means: arithmetic (A), geometric (G) and harmonic (H), as they satisfy $A \geq G \geq H$, but it is also easy to visualize the two others in this role. To what extent and in which regions of the domain these three different means of *precision* and *recall* actually diverge from one another may be observed e.g. in Figs 12 and 13, where both *precision* and *recall* as well as their three means (arithmetic: $A(p, r)$, geometric: $G(p, r)$ and harmonic: $H(p, r)$) are shown. This visualization illustrates well the concave isolines of $A(p, r)$ and $G(p, r)$, which means that they obtain excessively high values for increasingly divergent recognition of classes, making $H(p, r)$ the best choice out of three in this respect.

³Some authors set $\beta = \sqrt{\frac{1-\lambda}{\lambda}}$ instead, resulting in $\frac{1-\lambda}{\lambda} = \beta^2$, which allows for some further interpretation of such β [29]; not to be pursued in this paper.

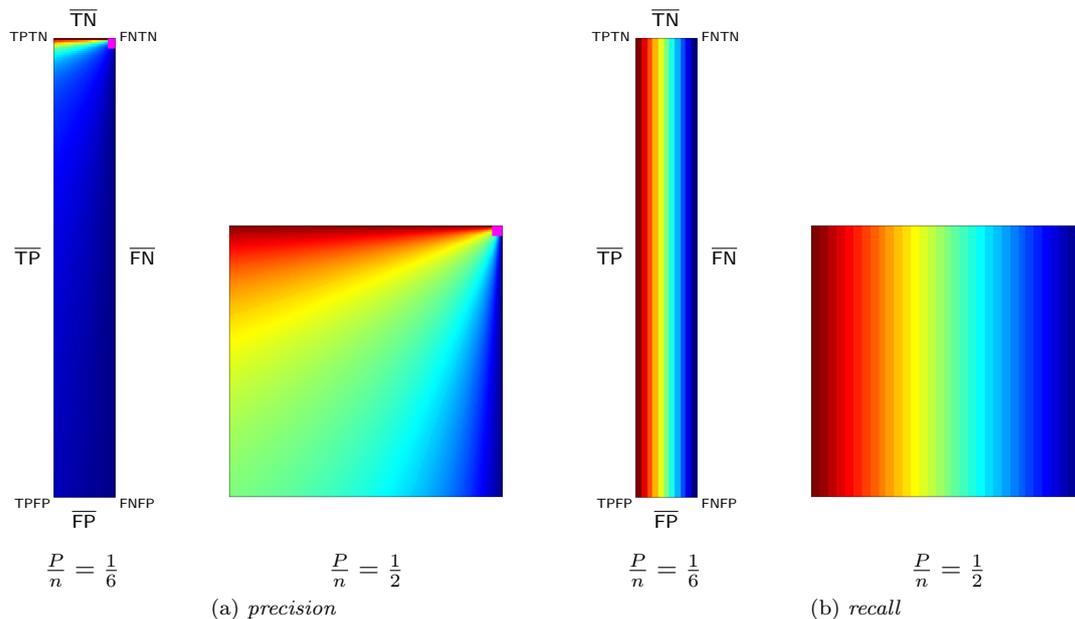


Figure 12: Cross-sections of *precision* and *recall*

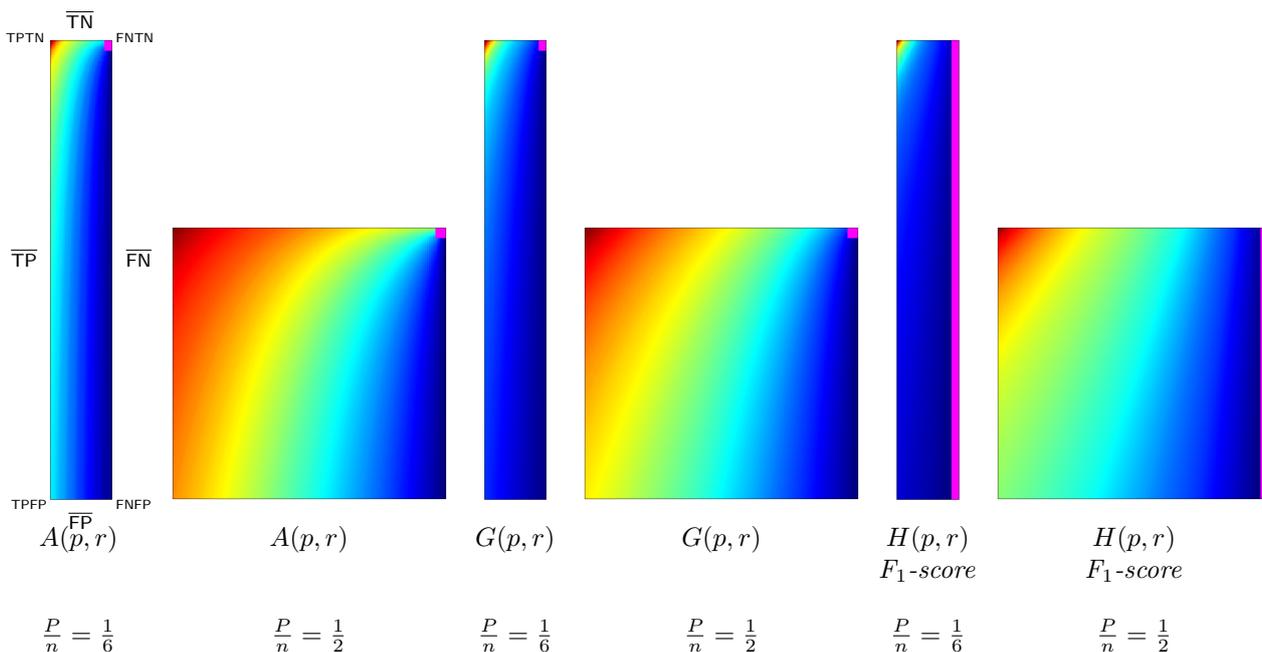


Figure 13: Cross-sections of different means of *precision* and *recall*

Deciding on the mean, however, is not enough, as the remaining problem regards changes in the measure's behaviour across the differing P/n . Unfortunately, for the harmonic mean as well as for the other two means, the measure's values gradually shift away from the positive class as P/n decreases, making all the three (regular) means of *precision* and *recall* (and thus the F_1 -score in particular) less and less suited for imbalanced data. This is where the weighed means, in particular F_β (the weighed harmonic mean of *precision* and *recall*) may actually turn out to be more useful.

The arising question regards the appropriate value of β . Clearly, the desired bias towards the positive class requires $\beta > 1$, which corresponds to applying more weight to *recall*. The visual solution to this problem is provided in Fig. 14, which shows cross-section visualizations of F_β for three values of $\beta \in \{1, 3, 5\}$. The range of these values has been inspired by the accessible data, in this case the class ratios considered in previous sections: $P/n = 1/6$ and $P/n = 1/2$. These values may be assumed to directly express the $[0, 1]$ -based weights of *precision* and *recall*, i.e. $\lambda = 1/6$ and $\lambda = 1/2$, which translate to $\beta = 5$ and $\beta = 1$, respectively.

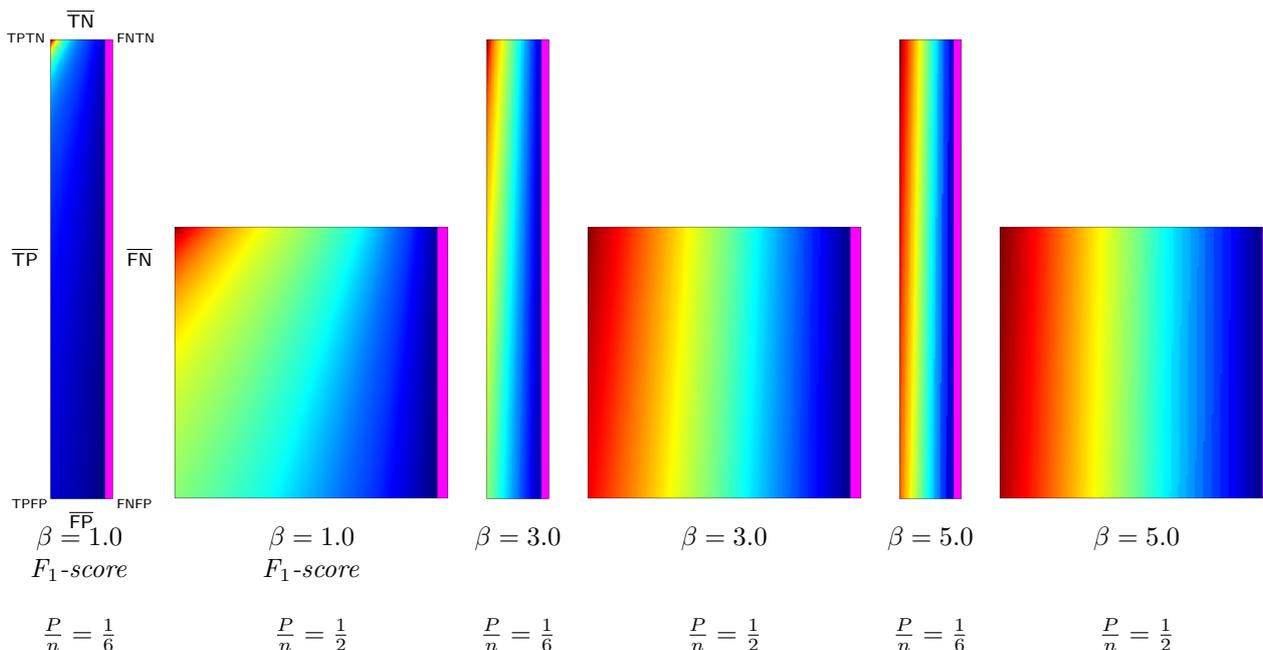


Figure 14: Cross-sections of F_β

Despite the fact that all ten of the earlier discussed properties of F_1 -score and F_β (for $\beta \in [0, \infty)$) are identical, as exemplified in Table 2, the increasing difference between F_1 -score and F_β resulting from the changing values of β is clearly visible in Fig. 14, revealing the truly multidimensional complexity of the measures' domains.

A slightly closer explanation may only be due for *ACE* property, as F_β 's particular visualization for $\beta = 5$ in Fig. 14 suggests that the measure satisfies the requirements of *ACE* (its values on the \overline{TP} side are not lower than their counterparts on the \overline{TN} side for both $P/n = 1/6$ and $P/n = 1/2$), whereas Table 2 states that *ACE* is not met by F_β . This is because the ten proposed properties are of general character, i.e. they concern the whole tetrahedron, which means that they must be satisfied in cross-sections corresponding to all feasible class proportions. In case of $F_{\beta=5}$, for some class ratios that are lower than those considered in the presented visualizations, e.g. for $P/n = 1/10$ (easily reproducible in the online visualization tool), the *ACE* conditions are actually not satisfied, thus justifying the contents of Table 2.

Nevertheless, for cases when the class ratio is known or predictable, the visualizations are of utmost practical value. In the discussed situation, the visual-based analysis may suggest non-trivial values of β for which F_β certainly satisfies selected properties, in this case the conditions of *ACE* for a particular P/n . A thorough analysis of cross-sections clearly suggested the existence of a particular dependency between β and the class proportion, which influences the *ACE* property. This observation inspired us to derive analytically the borderline value of β that ensures that *ACE* is met by F_β .

Proposition 1. F_β satisfies *ACE* property for $\beta \geq N/P$ (for proof see the Appendix).

Practically this means that the user must bear in mind the class proportions and may use it to make F_β satisfy *ACE* property, if needed.

5.2.2. External parametrization: $IBA_\alpha(G\text{-mean})$

Applying any external parametrization, e.g. the IBA_α scheme [13, 14, 15], to different measures evokes the usual problems, first of all related to establishing the values of required parameters. Visualization provides a very practical solution to these issues, as shall be demonstrated in this section.

Given a classifier performance measure M , a parameter $\alpha \geq 0$ and a tentative measure $Dom = \text{sensitivity} - \text{specificity}$, the formula:

$$IBA_\alpha(M) = (1 + \alpha Dom)M$$

defines the parametrization of M , in which this measure is multiplicatively combined with $(1 + \alpha Dom)$. Of course, $IBA_\alpha(M) = M$ for $\alpha = 0$. Simultaneously, when $Dom \in [-1, +1]$ and $\alpha \leq 1$ then $1 + \alpha Dom \geq 0$, which, together with $M \geq 0$, implies $IBA_\alpha(M) \geq 0$.

The scheme has been conceived to increase the measure orientation towards the positive class, which makes it a good choice in the imbalanced contexts. Notice, however, that neither Dom is a classic classifier performance measure (as its domain includes negative values), nor is $IBA_\alpha(M)$ a simple convex combination of Dom and M . This renders strictly analytical (without any visualization tool) analysis of $IBA_\alpha(M)$ very hard, especially for larger values of α . In result, while the general goal of reorienting the measure towards the positive class is certainly achieved by IBA_α , it is not instantly clear how this reorientation is practically manifested. In particular, one might be interested in identifying whether measure M subjected to $IBA_\alpha(M)$ satisfies any of the postulated properties, or not (and, if it does, which ones and for what ranges of α).

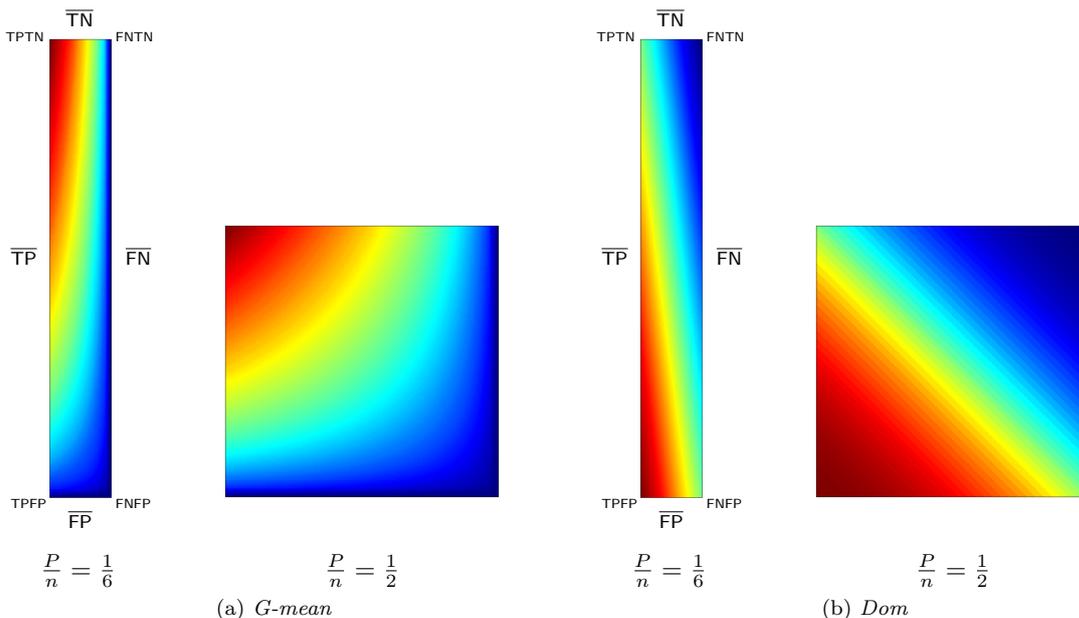


Figure 15: Cross-sections of $G\text{-mean}$ and Dom

Below, we visualize and analyse $G\text{-mean}$ externally parametrized according to IBA_α for $\alpha \in \{0, 0.5, 1\}$. The combination $IBA_\alpha(G\text{-mean})$ was particularly recommended and analytically studied for the aforementioned α values by García et al. [13]. Tracing the influence of Dom on $G\text{-mean}$ within the IBA_α approach may well be started with the visualization of the components of the parametrization procedure, see Fig. 15, as only having realized the behaviour of $G\text{-mean}$ and Dom , can one infer how the changing α impacts the parametrized measure. Clearly, for $\alpha \rightarrow 0$, $IBA_\alpha(G\text{-mean}) \rightarrow G\text{-mean}$, so only $\alpha > 0$ exerts any influence on the result. Notice that Dom features a rather unexpected growth towards vertex TFPF, implying the specific behaviour of $IBA_\alpha(G\text{-mean})$, see Fig. 16. Because the combination is multiplicative, the values

of G -mean are being ‘amplified’ by the corresponding values of $(1 + \alpha Dom)$, in particular: increased for $(1 + \alpha Dom) > 1$, and decreased for $(1 + \alpha Dom) < 1$.

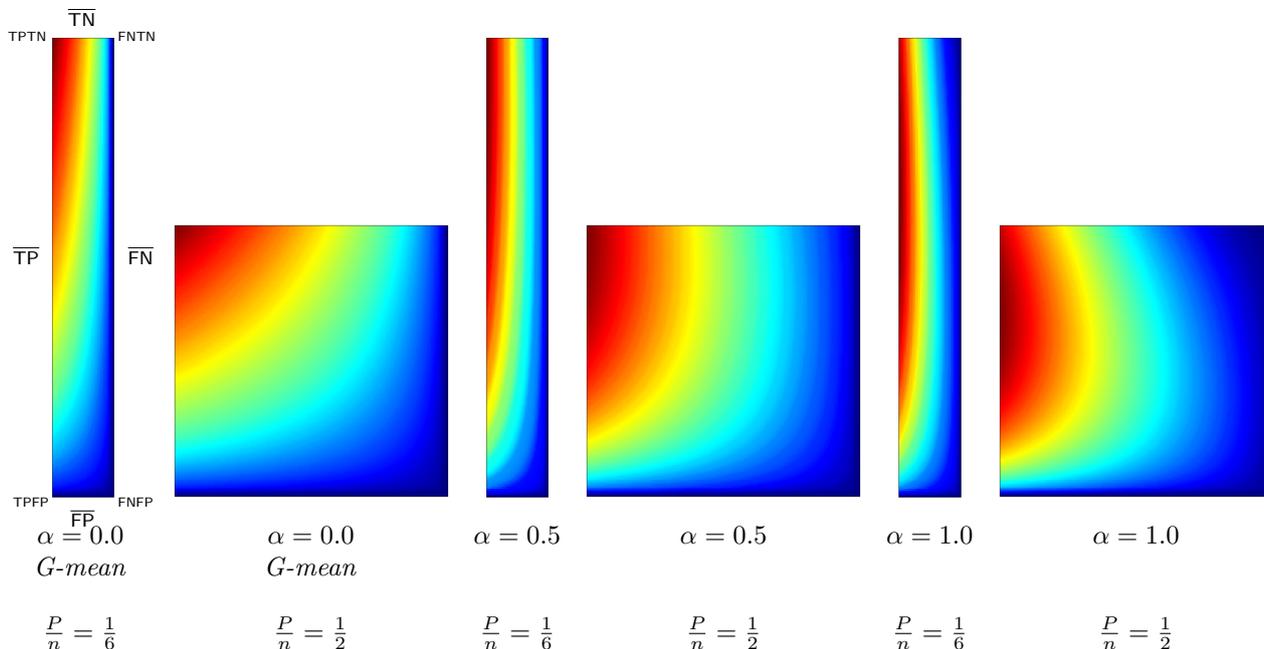


Figure 16: Cross-sections of $IBA_\alpha(G$ -mean)

As stated in Table 2 G -mean satisfies all of the proposed properties. The important question is how the application of external parametrization to the measure influences its properties, e.g. the ACE property (thoroughly discussed for internal parametrization). Unsurprisingly, $IBA_\alpha(G$ -mean) may be proven to satisfy the ACE property for all assumed values of α . Notably, this comes with a cost, as this parametrization of G -mean is not equally stable with respect to all other properties.

Proposition 2. $IBA_\alpha(G$ -mean) satisfies ACE property for $\alpha \geq 0$ (for proof see the Appendix).

Practically, this means that the $IBA_\alpha(G$ -mean) does not depart from the original G -mean in terms of ACE . However, a thorough visual-based analysis of the impact of the α parameter on satisfying TN_{\nearrow} by $IBA_\alpha(G$ -mean) suggested a border-line value of α . Inspired thereby, we derived the exact value analytically.

Proposition 3. $IBA_\alpha(G$ -mean) satisfies TN_{\nearrow} property for $\alpha \leq 1/3$ (for proof see the Appendix).

Table 3 gathers the results concerning the ten devised properties for particular intervals of the α parameter implied by its border-line value.

Table 3: Properties of G -mean and its parametrizations; *: contains NaN (undefined value); †: NaN side, *: strong monotonicity

Measure	$TPTN_{max}$	\overline{FN}_{min}	\overline{FP}_{min}	TP_{\nearrow}	TN_{\nearrow}	$\overline{TN}_{\neq max}$	$\overline{TP}_{\neq max}$	ACE	ACH	$UnDefs$
$IBA_0(G$ -mean) = G -mean	✓	✓	✓	✓	✓	✓	✓	✓	✓	TN-FP; TP-FN
$IBA_\alpha(G$ -mean), $\alpha \in (0, 1/3]$	✓	✓	✓	✓	✓	✓	✓	✓	×	TN-FP; TP-FN
$IBA_\alpha(G$ -mean), $\alpha \in (1/3, \infty)$	×	✓	✓	✓	×	✓	×	✓	×	TN-FP; TP-FN

In particular, the entries for G -mean and for $IBA_{\alpha \in (0, 1/3]}(G$ -mean) state that such external parametrization eliminates the symmetry of handling both classes. It is the result of incorporating the class-asymmetric Dom component in the IBA_α parametrization procedure. The parametrized G -mean becomes slightly (as

α does not exceed $1/3$) more oriented towards the positive class (see also Fig. 16), and thus does not satisfy the *ACH* property any more. Nevertheless, other properties remain satisfied as long as $\alpha \leq 1/3$. The behaviour of $IBA_\alpha(G\text{-mean})$ changes drastically, however, when α exceeds $1/3$ (see Table 3 and Fig. 16). On one hand, for $\alpha > 0$ one gets the much desired focus on the positive class, reflected by the *ACE* property (for any two corresponding points on sides \overline{TP} and \overline{TN} , the value on side \overline{TP} is strictly greater than that on \overline{TN}), however, for $\alpha > 1/3$ this comes with the inevitable cost of losing not only the above-mentioned *ACH* property, but also the TN_{\nearrow} property (manifested by non-monotonic growth of the measure from \overline{FP} to \overline{TN}). Additionally, for $\alpha > 1/3$ the maximal value of $IBA_\alpha(G\text{-mean})$ drifts away from vertex $TPTN$ (i.e. the full recognition of both classes is no longer rewarded with the maximal measure value) violating the $TPTN_{max}$ and $\overline{TP}_{\neq max}$ properties. In this context, the usability of $IBA_\alpha(G\text{-mean})$ for $\alpha > 1/3$ becomes questionable.

6. Conclusions

In this paper, we proposed a new visualization technique for analysing classification performance measures and contributed an interactive tool implementing it in the form of a web application. The technique uses a barycentric coordinate system by projecting values from the confusion matrix into a three-dimensional figure — a tetrahedron. Unlike simpler visualizations this technique:

- provides general interpretations in terms of the four values of the two-class confusion matrix,
- involves exclusively linear, and thus easily interpretable, $4D \rightarrow 3D$ transformations,
- allows for analysing full ranges of measure values with respect to all possible combinations of confusion matrix entries,
- naturally illustrates the $TP + FN + FP + TN = n$ constraint, manifested in the shape of the space (i.e. tetrahedron),
- remains defined for all possible combinations of the matrix entries,
- admits multiple cross-sections with natural interpretations in terms of simple measures, e.g. horizontal cross-sections, which correspond to the proportion of actual classes (i.e. the positive ($\frac{P}{n}$) and the negative ($\frac{N}{n}$) class) and are thus especially well suited for analysis of imbalanced data.

Using this visualization technique, we analysed 22 classifier performance measures in terms of ten purposefully defined properties, which can help assess the measures in the context of class imbalanced data. The analysis included non-parametric as well as parametric measures, which led to discovering property changes upon certain parametrizations for the latter. In particular, we have derived threshold values for selected properties of F_β and $IBA_\alpha(G\text{-mean})$. The detection of these non-trivial thresholds would be difficult without the proposed visualization technique.

The analysis of the selected measures illustrates how the proposed visualization can depict individual characteristics and potential caveats of each measure. It is worth stressing that it was not our intention to promote any single measure as the best, since the measure choice always finally depends on the user and the application at hand. Nevertheless, our visualization tool and the results gathered in Table 2 should support making this choice.

As future work, we plan to consider also other properties, such as gradients of measure as functions of the four arguments. Moreover, it would be interesting to analyse the effects of applying cost matrices to the visualized measures. Similarly, the effects of micro- and macro-averaging of binary measures in multi-class scenarios are worth studying. Finally, we hope that the visualization technique may be helpful in defining new classifier performance measures.

Acknowledgement

This research was partly supported by the FNP START scholarship (first author) and Institute of Computing Science Statutory Funds.

References

- [1] R. Alaíz-Rodríguez, N. Japkowicz, P. E. Tischer, A Visualization-Based Exploratory Technique for Classifier Comparison with Respect to Multiple Metrics and Multiple Domains, in: Proc. 19th European Conf. Mach. Learn., Part II, 660–665, 2008.
- [2] P. Baldi, S. Brunak, Y. Chauvin, C. Andersen, H. Nielsen, Assessing the Accuracy of Prediction Algorithms for Classification: An Overview, *Bioinformatics* 16 (2000) 412–424.
- [3] M. Bekkar, H. Djemaa, A. Taklit, Evaluation Measures for Models Assessment Over Imbalanced Data Sets, *Journal of Inform. Eng. and Appl.* 3 (10) (2013) 27–38.
- [4] D. Brzezinski, J. Stefanowski, Prequential AUC: Properties of the Area Under the ROC Curve for Data Streams with Concept Drift, *Knowledge and Information Systems* 52 (2) (2017) 531–562.
- [5] R. Caruana, A. Niculescu-Mizil, Data Mining in Metric Space: An Empirical Analysis of Supervised Learning Performance Criteria, in: Proc. 10th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 69–78, 2004.
- [6] E. Celotto, Visualizing the behavior and some symmetry properties of Bayesian confirmation measures, *Data Min. Knowl. Discov.* 31 (3) (2017) 739–773.
- [7] J. Davis, M. Goadrich, The relationship between Precision-Recall and ROC curves, in: Proc. 23rd Int. Conf. Mach. Learn., 233–240, 2006.
- [8] P. M. Domingos, A few useful things to know about machine learning, *Commun. ACM* 55 (10) (2012) 78–87.
- [9] T. Fawcett, An introduction to ROC analysis, *Pattern Recognition Letters* 27 (8) (2006) 861–874.
- [10] C. Ferri, J. Hernández-Orallo, R. Modroiu, An Experimental Comparison of Performance Measures for Classification, *Pattern Recognit. Lett.* 30 (1) (2009) 27–38.
- [11] P. A. Flach, The Geometry of ROC Space: Understanding Machine Learning Metrics through ROC Isometrics, in: Proc. 20th Int. Conf. Mach. Learn., 194–201, 2003.
- [12] J. Fürnkranz, P. A. Flach, An Analysis of Rule Evaluation Metrics, in: Proc. 20th Int. Conf. Mach. Learn., 202–209, 2003.
- [13] V. García, R. A. Mollineda, J. S. Sánchez, Index of Balanced Accuracy: A Performance Measure for Skewed Class Distributions, in: Proc. 4th Iberian Conf. Pattern Recognition Image Analysis, 441–448, 2009.
- [14] V. García, R. A. Mollineda, J. S. Sánchez, Theoretical Analysis of a Performance Measure for Imbalanced Data, in: Proc. 20th Int. Conf. Pattern Recognition, 617–620, 2010.
- [15] V. García, R. A. Mollineda, J. S. Sánchez, A bias correction function for classification performance assessment in two-class imbalanced problems, *Knowledge-Based Systems* 59 (2014) 66–74.
- [16] Q. Gu, L. Zhu, C. Z., Evaluation Measures of the Classification Performance of Imbalanced Data Set, in: Proc. ISICA, Springer, 461–471, 2009.
- [17] H. He, E. A. Garcia, Learning from Imbalanced Data, *IEEE Trans. Knowl. Data Eng.* 21 (9) (2009) 1263–1284.
- [18] H. He, Y. Ma (Eds.), *Imbalanced Learning: Foundations, Algorithms and Applications*, IEEE - Wiley, 2013.
- [19] B. Hu, W. Dong, A Study on Cost Behaviors of Binary Classification Measures in Class-imbalanced Problems, *CoRR* abs/1403.7100 .
- [20] N. Japkowicz, M. Shah, *Evaluating Learning Algorithms: A Classification Perspective*, Cambridge University Press, ISBN 9780521196000, 2011.
- [21] M. Kubat, R. Holte, S. Matwin, Machine Learning for the Detection of Oil Spills in Radar Images, *Machine Learning Journal* 30 (1998) 195–215.
- [22] Y. Le Bras, P. Lenca, S. Lallich, *Data Mining: Foundations and Intelligent Paradigms*, chap. Formal Framework for the Study of Algorithmic Properties of Objective Interestingness Measures, Springer, 77–98, 2012.
- [23] V. López, A. Fernández, S. García, V. Palade, F. Herrera, An insight into classification with imbalanced data: Empirical results and current trends on using data intrinsic characteristics, *Inf. Sci.* 250 (2013) 113–141.
- [24] G. Piatetsky-Shapiro, Discovery, Analysis, and Presentation of Strong Rules, in: *Knowledge Discovery in Databases*, AAAI/MIT Press, 229–248, 1991.
- [25] G. Piatetsky-Shapiro, B. M. Masand, Estimating Campaign Benefits and Modeling Lift, in: Proc. 5th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining, 185–193, 1999.
- [26] D. M. Powers, What the F-measure Doesn't Measure: Features, Flaws, Fallacies And Fixes, *CoRR* abs/1503.06410 .
- [27] F. J. Provost, T. Fawcett, R. Kohavi, The Case Against Accuracy Estimation for Comparing Induction Algorithms, in: Proc. 15th Int. Conf. Mach. Learn., 445–453, 1998.
- [28] R. Ranawana, V. Palade, Optimized Precision - A New Measure for Classifier Performance Evaluation, in: Proc. IEEE Cong. on Evol. Computation, 16–21, 2006.
- [29] Y. Sasaki, The truth of the F-measure, <http://www.cs.odu.edu/mukka/cs795sum10dm/Lecturenotes/Day3/F-measure-YS-26Oct07.pdf> .
- [30] M. Sokolova, G. Lapalme, A Systematic Analysis of Performance Measures for Classification Tasks, *Inf. Process. Manage.* 45 (4) (2009) 427–437.
- [31] R. Susmaga, I. Szczęch, Can Interestingness Measures Be Usefully Visualized?, *Int. J. Applied Math. Comp. Science* 25 (2) (2015) 323–336.
- [32] R. Susmaga, I. Szczęch, Visualization Support for the Analysis of Properties of Interestingness Measures, *Bulletin of the Polish Academy of Sciences Technical Sciences* 63 (1) (2015) 315–327.
- [33] S. Vanderlooy, I. G. Sprinkhuizen-Kuyper, E. N. Smirnov, H. J. van den Herik, The ROC isometrics approach to construct reliable classifiers, *Intell. Data Anal.* 13 (1) (2009) 3–37.
- [34] S. Wang, L. L. Minku, X. Yao, Resampling-Based Ensemble Methods for Online Class Imbalance Learning, *IEEE Trans. Knowl. Data Eng.* 27 (5) (2015) 1356–1368.

- [35] J. Warren, On the Uniqueness of Barycentric Coordinates, in: Contemporary Mathematics. Proceedings of AGGM '02, 93–99, 2003.

Appendix: Proofs of Propositions

Proof of Proposition 1

For $P > 0$ (the positive class) and $N \geq 0$ (the negative class), the (positive) class ratio is expressed as N/P . Given that, F_β satisfies the ACE property if $F_\beta\left(\begin{smallmatrix} P \\ \gamma N \end{smallmatrix} \begin{smallmatrix} 0 \\ (1-\gamma)N \end{smallmatrix}\right) \geq F_\beta\left(\begin{smallmatrix} (1-\gamma)^P & \gamma^P \\ 0 & N \end{smallmatrix}\right)$ for every $\gamma \in [0, 1]$, provided both sides of the inequality are defined.

Because F_β , a function of *precision* (everywhere below in this subsection: p) and *recall* (everywhere below in this subsection: r), is defined as $F_\beta = \frac{(1+\beta)pr}{\beta p+r}$ with $\beta \geq 0$, and

- on the left-hand side:

$$\begin{aligned} - p &= \frac{P}{P+\gamma N}, \\ - r &= \frac{P}{P+0}, \text{ so under the assumed } P > 0, r = 1, \end{aligned}$$

- on the right-hand side:

$$\begin{aligned} - p &= \frac{(1-\gamma)P}{(1-\gamma)P+0}, \text{ so under the assumed } P > 0, = 1 \text{ for } \gamma \neq 1, \\ - r &= \frac{(1-\gamma)P}{(1-\gamma)P+\gamma P} = \frac{(1-\gamma)P}{(1-\gamma+\gamma)P} = \frac{(1-\gamma)P}{P}, \text{ so under the assumed } P > 0, r = 1 - \gamma, \end{aligned}$$

the inequality $F_\beta\left(\begin{smallmatrix} P \\ \gamma N \end{smallmatrix} \begin{smallmatrix} 0 \\ (1-\gamma)N \end{smallmatrix}\right) \geq F_\beta\left(\begin{smallmatrix} (1-\gamma)^P & \gamma^P \\ 0 & N \end{smallmatrix}\right)$ is expressed as:

$$\begin{aligned} \frac{(1+\beta)\frac{P}{P+\gamma N} \cdot 1}{\beta\frac{P}{P+\gamma N}+1} &\geq \frac{(1+\beta)(1-\gamma)}{\beta \cdot 1+(1-\gamma)} \\ \frac{(1+\beta)\frac{P}{P+\gamma N}}{\beta\frac{P}{P+\gamma N}+1} &\geq \frac{(1+\beta)(1-\gamma)}{\beta+(1-\gamma)} \end{aligned}$$

The assumed $P > 0$, $N \geq 0$, $\gamma \geq 0$, $\beta \geq 0$ ensure $\beta\frac{P}{P+\gamma N} + 1 > 0$, so:

$$(1 + \beta)\frac{P}{P+\gamma N} \geq \frac{(1+\beta)(1-\gamma)}{\beta+(1-\gamma)}\left(\beta\frac{P}{P+\gamma N} + 1\right)$$

The assumed $\gamma \in [0, 1]$ and $\beta \geq 0$ ensure $\beta+(1-\gamma) \geq 0$, so assuming additionally $\gamma < 1$ ensures $\beta+(1-\gamma) > 0$, so:

$$(\beta + (1 - \gamma))(1 + \beta)\frac{P}{P+\gamma N} \geq (1 + \beta)(1 - \gamma)\left(\beta\frac{P}{P+\gamma N} + 1\right)$$

The assumed $\beta \geq 0$ ensures $1 + \beta > 0$, so:

$$\begin{aligned} (\beta + (1 - \gamma))\frac{P}{P+\gamma N} &\geq (1 - \gamma)\left(\beta\frac{P}{P+\gamma N} + 1\right) \\ \beta\frac{P}{P+\gamma N} + (1 - \gamma)\frac{P}{P+\gamma N} &\geq (1 - \gamma)\left(\beta\frac{P}{P+\gamma N} + 1\right) \\ \beta\frac{P}{P+\gamma N} + (1 - \gamma)\frac{P}{P+\gamma N} &\geq (1 - \gamma)\beta\frac{P}{P+\gamma N} + (1 - \gamma) \\ \beta\frac{P}{P+\gamma N} - (1 - \gamma)\beta\frac{P}{P+\gamma N} &\geq (1 - \gamma)\left(1 - \frac{P}{P+\gamma N}\right) \\ \beta\frac{P}{P+\gamma N}(1 - (1 - \gamma)) &\geq (1 - \gamma)\left(1 - \frac{P}{P+\gamma N}\right) \\ \beta\frac{P}{P+\gamma N}\gamma &\geq (1 - \gamma)\left(1 - \frac{P}{P+\gamma N}\right) \\ \beta\frac{P}{P+\gamma N}\gamma &\geq (1 - \gamma)\left(\frac{P+\gamma N}{P+\gamma N} - \frac{P}{P+\gamma N}\right) \\ \beta\frac{P}{P+\gamma N}\gamma &\geq (1 - \gamma)\frac{P+\gamma N-P}{P+\gamma N} \end{aligned}$$

The assumed $P > 0$, $N \geq 0$, $\gamma \geq 0$ ensure $P + \gamma N > 0$, so:

$$\beta P \gamma \geq (1 - \gamma)\gamma N$$

The assumed $P > 0$ allows for:

$$\beta \gamma \geq (1 - \gamma)\frac{\gamma N}{P}$$

Assuming additionally $\gamma > 0$ allows for:

$$\beta \geq (1 - \gamma)\frac{N}{P}$$

Intermediate conclusion: given $P > 0$ (the positive class) and $N \geq 0$ (the negative class) the inequality: $F_\beta\left(\begin{smallmatrix} P \\ \gamma N \end{smallmatrix} \begin{smallmatrix} 0 \\ (1-\gamma)N \end{smallmatrix}\right) \geq F_\beta\left(\begin{smallmatrix} (1-\gamma)^P & \gamma^P \\ 0 & N \end{smallmatrix}\right)$ is fully defined and holds for $\gamma \in (0, 1)$ if β is taken to satisfy $\beta \geq (1 - \gamma)\frac{N}{P}$.

The two remaining border cases (resulting from additionally assuming $\gamma < 1$ and $\gamma > 0$) are:

- $\gamma = 1$:

$$F_{\beta}(\begin{bmatrix} P & 0 \\ \gamma N & (1-\gamma)N \end{bmatrix}) \geq F_{\beta}(\begin{bmatrix} (1-\gamma)^P & \gamma^P \\ 0 & N \end{bmatrix})$$

$$F_{\beta}(\begin{bmatrix} P & 0 \\ N & 0 \end{bmatrix}) \geq F_{\beta}(\begin{bmatrix} 0 & P \\ 0 & N \end{bmatrix})$$

which cannot be established, as $p = \frac{0}{0+0}$, and thus F_{β} , is undefined on the right-hand side.

- $\gamma = 0$:

$$F_{\beta}(\begin{bmatrix} P & 0 \\ 0 & N \end{bmatrix}) \geq F_{\beta}(\begin{bmatrix} P & 0 \\ 0 & N \end{bmatrix})$$

which holds trivially, as the argument on both sides is the same

(so, on both sides, either F_{β} is undefined or it is defined and equal).

Final conclusion: given $P > 0$ (the positive class) and $N \geq 0$ (the negative class) the inequality:

$$F_{\beta}(\begin{bmatrix} P & 0 \\ \gamma N & (1-\gamma)N \end{bmatrix}) \geq F_{\beta}(\begin{bmatrix} (1-\gamma)^P & \gamma^P \\ 0 & N \end{bmatrix})$$

is fully defined and holds for every $\gamma \in [0, 1]$ if β is taken to satisfy $\beta \geq (1-\gamma)\frac{N}{P}$.

This result may be further simplified, because $(1-\gamma)\frac{N}{P}$ changes linearly with γ and for $\gamma = 1$: $(1-\gamma)\frac{N}{P} = 0$ (which means that β is required to satisfy condition $\beta \geq 0$), while for $\gamma = 0$: $(1-\gamma)\frac{N}{P} = \frac{N}{P}$ (which means that β is required to satisfy condition $\beta \geq \frac{N}{P}$). Notice that the assumed $P > 0$ and $N \geq 0$ ensure $(1-\gamma)\frac{N}{P} \geq 0$, which subsumes the assumed $\beta \geq 0$. Setting β to satisfy $\beta \geq \frac{N}{P}$ ensures satisfying both conditions.

Summarizing all the considered cases, $\beta \geq \frac{N}{P}$ ensures $F_{\beta}(\begin{bmatrix} P & 0 \\ \gamma N & (1-\gamma)N \end{bmatrix}) \geq F_{\beta}(\begin{bmatrix} (1-\gamma)^P & \gamma^P \\ 0 & N \end{bmatrix})$ for every $\gamma \in [0, 1]$ for which both sides of the inequality are defined, which proves that $F_{\beta \geq \frac{N}{P}}$ satisfies ACE.

Proof of Proposition 2

Let $P > 0$ (the positive class) and $N > 0$ (the negative class). $IBA_{\alpha}(G\text{-mean})$ satisfies the ACE property if $IBA_{\alpha}(G\text{-mean})(\begin{bmatrix} P & 0 \\ \gamma N & (1-\gamma)N \end{bmatrix}) \geq IBA_{\alpha}(G\text{-mean})(\begin{bmatrix} (1-\gamma)^P & \gamma^P \\ 0 & N \end{bmatrix})$ for every $\gamma \in [0, 1]$, provided both sides of the inequality are defined.

$IBA_{\alpha}(G\text{-mean})$ is a function of *recall* (everywhere below in this subsection: r) and *specificity* (everywhere below in this subsection: s), and

- on the left-hand side:

$$- r = \frac{P}{P+0}, \text{ so under the assumed } P > 0, r = 1,$$

$$- s = \frac{(1-\gamma)N}{\gamma N + (1-\gamma)N} = \frac{(1-\gamma)N}{(\gamma+1-\gamma)N} = \frac{(1-\gamma)N}{N} = \frac{(1-\gamma)N}{N}, \text{ so under the assumed } N > 0, s = 1 - \gamma,$$

- on the right-hand side:

$$- r = \frac{(1-\gamma)^P}{(1-\gamma)^P + \gamma^P} = \frac{(1-\gamma)^P}{(1-\gamma+\gamma)^P} = \frac{(1-\gamma)^P}{P}, \text{ so under the assumed } P > 0, r = 1 - \gamma,$$

$$- s = \frac{N}{N+0}, \text{ so under the assumed } N > 0, s = 1.$$

Given $r \in [0, 1]$, $s \in [0, 1]$ and $\alpha \geq 0$, $IBA_{\alpha}(G\text{-mean})$ is defined in terms of r and s as: $IBA_{\alpha}(G\text{-mean}) = (1 + \alpha(r - s))r^{\frac{1}{2}}s^{\frac{1}{2}}$.

In result, $IBA_{\alpha}(G\text{-mean})$ satisfies the ACE property if $IBA_{\alpha}(G\text{-mean})(\begin{bmatrix} P & 0 \\ \gamma N & (1-\gamma)N \end{bmatrix}) \geq IBA_{\alpha}(G\text{-mean})(\begin{bmatrix} (1-\gamma)^P & \gamma^P \\ 0 & N \end{bmatrix})$, which is also expressed as:

$$(1 + \alpha(1 - (1 - \gamma)))1^{\frac{1}{2}}(1 - \gamma)^{\frac{1}{2}} \geq (1 + \alpha((1 - \gamma) - 1))(1 - \gamma)^{\frac{1}{2}}1^{\frac{1}{2}}$$

$$(1 + \alpha\gamma)(1 - \gamma)^{\frac{1}{2}} \geq (1 - \alpha\gamma)(1 - \gamma)^{\frac{1}{2}}$$

Assuming additionally $\gamma < 1$, which implies $(1 - \gamma)^{\frac{1}{2}} > 0$, and dividing by $(1 - \gamma)^{\frac{1}{2}}$

$$1 + \alpha\gamma \geq 1 - \alpha\gamma$$

$$\alpha\gamma \geq -\alpha\gamma$$

$$2\alpha\gamma \geq 0$$

Assuming additionally $\gamma > 0$, which implies $2\gamma > 0$, and dividing by 2γ

$\alpha \geq 0$

Intermediate conclusion: given $P > 0$ (the positive class) and $N \geq 0$ (the negative class) the inequality: $IBA_\alpha(G\text{-mean})([\begin{smallmatrix} P & 0 \\ \gamma N & (1-\gamma)N \end{smallmatrix}]) \geq IBA_\alpha(G\text{-mean})([\begin{smallmatrix} (1-\gamma)^P & \gamma^P \\ 0 & N \end{smallmatrix}])$ is fully defined and holds for every $\gamma \in (0, 1)$ if α is taken to satisfy $\alpha \geq 0$.

The two remaining border cases (resulting from additionally assuming $\gamma < 1$ and $\gamma > 0$) are:

- $\gamma = 1$:

$$\begin{aligned} IBA_\alpha(G\text{-mean})([\begin{smallmatrix} P & 0 \\ \gamma N & (1-\gamma)N \end{smallmatrix}]) &\geq IBA_\alpha(G\text{-mean})([\begin{smallmatrix} (1-\gamma)^P & \gamma^P \\ 0 & N \end{smallmatrix}]) \\ IBA_\alpha(G\text{-mean})([\begin{smallmatrix} P & 0 \\ N & 0 \end{smallmatrix}]) &\geq IBA_\alpha(G\text{-mean})([\begin{smallmatrix} 0 & P \\ 0 & N \end{smallmatrix}]) \\ (1 + \alpha(1 - (1 - 1)))1^{\frac{1}{2}}(1 - 1)^{\frac{1}{2}} &\geq (1 + \alpha((1 - 1) - 1))(1 - 1)^{\frac{1}{2}}1^{\frac{1}{2}} \\ 0 &\geq 0 \text{ which holds trivially for every } \alpha, \end{aligned}$$

- $\gamma = 0$:

$$\begin{aligned} IBA_\alpha(G\text{-mean})([\begin{smallmatrix} P & 0 \\ 0 & N \end{smallmatrix}]) &\geq IBA_\alpha(G\text{-mean})([\begin{smallmatrix} P & 0 \\ 0 & N \end{smallmatrix}]) \\ \text{which holds trivially, as the argument on both sides is the same} \\ \text{(so, on both sides, either } IBA_\alpha(G\text{-mean}) \text{ is undefined or it is defined and equal).} \end{aligned}$$

Final conclusion: given $P > 0$ (the positive class) and $N > 0$ (the negative class) the inequality: $IBA_\alpha(G\text{-mean})([\begin{smallmatrix} P & 0 \\ \gamma N & (1-\gamma)N \end{smallmatrix}]) \geq IBA_\alpha(G\text{-mean})([\begin{smallmatrix} (1-\gamma)^P & \gamma^P \\ 0 & N \end{smallmatrix}])$ is fully defined and holds for every $\gamma \in [0, 1]$ if α is taken to satisfy $\alpha \geq 0$.

Summarizing all the considered cases, $\alpha \geq 0$ ensures

$$IBA_\alpha(G\text{-mean})([\begin{smallmatrix} P & 0 \\ \gamma N & (1-\gamma)N \end{smallmatrix}]) \geq IBA_\alpha(G\text{-mean})([\begin{smallmatrix} (1-\gamma)^P & \gamma^P \\ 0 & N \end{smallmatrix}])$$

for every $\gamma \in [0, 1]$ for which both sides of the inequality are defined, which proves that $IBA_{\alpha \geq 0}(G\text{-mean})$ satisfies ACE.

Proof of Proposition 3

Let $M(\alpha, r, s)$, where M is a function of *recall* (everywhere below in this subsection: r), *specificity* (everywhere below in this subsection: s) and α , denote $IBA_\alpha(G\text{-mean})$.

Given $r \in [0, 1]$, $s \in [0, 1]$ and $\alpha \geq 0$:

$$M(\alpha, r, s) = (1 + \alpha(r - s))(rs)^{\frac{1}{2}} = (1 + \alpha r)r^{\frac{1}{2}}s^{\frac{1}{2}} - \alpha r^{\frac{1}{2}}s^{\frac{3}{2}} = r^{\frac{1}{2}}((1 + \alpha r)s^{\frac{1}{2}} - \alpha s^{\frac{3}{2}})$$

$M(\alpha, r, s)$ satisfies the TP_{\nearrow} property if it features a weakly monotonic value growth along vertical lines in its cross-sections for $P/n \in (0, 1)$, which is equivalent to it being a weakly increasing function of $s \in [0, 1]$.

Calculating $\frac{\partial M}{\partial s} = r^{\frac{1}{2}}(\frac{1}{2}(1 + \alpha r)s^{-\frac{1}{2}} - \frac{3}{2}\alpha s^{\frac{1}{2}})$ allows for:

$$\frac{\partial M}{\partial s} \geq 0$$

$$r^{\frac{1}{2}}(\frac{1}{2}(1 + \alpha r)s^{-\frac{1}{2}} - \frac{3}{2}\alpha s^{\frac{1}{2}}) \geq 0$$

Assuming additionally $r > 0$, which implies $r^{\frac{1}{2}} > 0$, and dividing by $r^{\frac{1}{2}}$

$$\frac{1}{2}(1 + \alpha r)s^{-\frac{1}{2}} - \frac{3}{2}\alpha s^{\frac{1}{2}} \geq 0$$

Assuming additionally $s > 0$, which implies $2s^{\frac{1}{2}} > 0$, and multiplying by $2s^{\frac{1}{2}}$

$$2^{\frac{1}{2}}(1 + \alpha r)s^{-\frac{1}{2}}s^{\frac{1}{2}} - 2^{\frac{3}{2}}\alpha s^{\frac{1}{2}}s^{\frac{1}{2}} \geq 0$$

$$1 + \alpha r - 3\alpha s \geq 0$$

$$1 + \alpha(r - 3s) \geq 0$$

Let $F(\alpha, r, s) = 1 + \alpha r - 3\alpha s$. $F(\alpha, r, s)$ is defined and continuous for $r \in [0, 1]$, $s \in [0, 1]$ and $\alpha \geq 0$, and treats r and s independently (as indicated by $\frac{\partial F}{\partial r} = \alpha$ and $\frac{\partial F}{\partial s} = -3\alpha$, which are independent of r and s). Thus,

$$\begin{aligned} F(\alpha, r, s) &\geq 0 \\ 1 + \alpha(r - 3s) &\geq 0 \\ \alpha(r - 3s) &\geq -1 \end{aligned}$$

Consider $r - 3s$:

- case $r - 3s = 0$ produces $0 \geq -1$ (holds trivially),
- case $r - 3s > 0$ produces $\alpha \geq \frac{-1}{r-3s}$, with $\frac{-1}{r-3s} < 0$,
- case $r - 3s < 0$ produces $\alpha \leq \frac{-1}{r-3s}$, with $\frac{-1}{r-3s} > 0$, further resolved into:
 - for $r \rightarrow 0$ and $s \rightarrow 0$: $\frac{-1}{r-3s} \rightarrow \infty$, in which sub-case $\alpha \leq \infty$
 - for $r \rightarrow 0$ and $s = 1$: $\frac{-1}{r-3s} \rightarrow 1/3$, in which sub-case $\alpha \leq 1/3$
 - for $r = 1$ and $s = 1$: $\frac{-1}{r-3s} = 1/2$, in which sub-case $\alpha \leq 1/2$

The resulting conditions on α are: $\alpha \geq \frac{-1}{r-3s}$ with $\frac{-1}{r-3s} < 0$, $\alpha \leq \infty$, $\alpha \leq 1/3$ and $\alpha \leq 1/2$, while the assumed condition is $\alpha \geq 0$ (with some of them subsuming some others). Setting α to satisfy $\alpha \in [0, 1/3]$ ensures satisfying all those conditions.

Intermediate conclusion: given $r \in (0, 1]$: $F(\alpha, r, s)$ is non-negative function of $s \in (0, 1]$ and $\frac{\partial M}{\partial s}$ is non-negative function of $s \in (0, 1]$ and $M(\alpha, r, s)$ is a weakly increasing function of $s \in (0, 1]$ if α is taken to satisfy $\alpha \in [0, 1/3]$.

The two remaining border cases (resulting from additionally assuming $r > 0$ and $s > 0$) are:

- $r = 0$: $M(\alpha, 0, s) = 0$ for $\alpha \in [0, 1/3]$ and $s \in [0, 1]$, so $M(\alpha, r, s)$ is a weakly increasing function of $s \in (0, 1]$ (thus also for $r = 0$),
- $s = 0$: $M(\alpha, r, 0) = 0$ for $\alpha \in [0, 1/3]$ and $r \in [0, 1]$ (including the above considered $r = 0$), while simultaneously $M(\alpha, r, s) \geq 0$ for $\alpha \in [0, 1/3]$ and $r \in [0, 1]$ (including the above considered $r = 0$) and $s \in [0, 1]$ ⁴, so $M(\alpha, r, s)$ is a weakly increasing function of $s \in [0, 1]$ (thus also for $s = 0$).

Final conclusion: given $r \in [0, 1]$, $M(\alpha, r, s)$ is a weakly increasing function of $s \in [0, 1]$ if α is taken to satisfy $\alpha \in [0, 1/3]$.

Summarizing all the considered cases, $\alpha \in [0, 1/3]$ ensures the weakly increasing character of $M(\alpha, r, s) = IBA_\alpha(G\text{-mean})$ as a function of $s \in [0, 1]$ for any $r \in [0, 1]$, being equivalent to featuring a weakly monotonic value growth along vertical lines in its cross-sections for $P/n \in (0, 1)$, which proves that $IBA_\alpha(G\text{-mean})$ satisfies TP_\nearrow .

⁴Proving $M \geq 0$ is analogous to proving $\frac{\partial M}{\partial s} \geq 0$