# Assessing the Reliability of a Multi-Class Classifier

**Luca Frigau, Claudio Conversano, and Francesco Mola**

**Abstract** Multi-class learning requires a classifier to discriminate among a large set of $L$ classes in order to define a classification rule able to identify the correct class for new observations. The resulting classification rule could not always be robust, particularly when imbalanced classes are observed or the data size is not large.

In this paper a new approach is presented aimed at evaluating the reliability of a classification rule. It uses a standard classifier but it evaluates the reliability of the obtained classification rule by re-training the classifier on resampled versions of the original data. User-defined misclassification costs are assigned to the obtained confusion matrices and then used as inputs in a Beta regression model which provides a cost-sensitive weighted classification index. The latter is used jointly with another index measuring dissimilarity in distribution between observed classes and predicted ones. Both indices are defined in $[0, 1]$ so that their values can be graphically represented in a $[0, 1]^2$ space. The visual inspection of the points for each classifier allows us to evaluate its reliability on the basis of the relationship between the values of both indices obtained on the original data and on resampled versions of it.

## 1 Introduction

In a classification problem it is common practice testing a wide variety of learning algorithms by varying threshold values and by using different tuning parameters. In that way different classifiers are obtained which can be compared in order to evaluate their predictive ability, which is usually evaluated starting from the confusion matrix. This is a contingency table in which each column represents the observations in a predicted class, while each row represents those in an actual class. Notationally, given a classification problem on $L$ classes observed on $n$ cases, let

L. Frigau (✉) • F. Mola
Department of Business and Economics, University of Cagliari, Cagliari, Italy
e-mail: frigau@unica.it; mola@unica.it

C. Conversano
Department of Mathematics and Informatics, University of Cagliari, Cagliari, Italy
e-mail: conversa@unica.it

$Q$ be a confusion matrix resulting from a classifier $k$. In this framework rows of $Q$ refer to the true classes, and columns of $Q$ to the predicted ones. By checking rows, the elements $q_{\ell j}$ indicate how many cases have been classified in each predicted class $\hat{\ell}_j$ ($j = 1, \ldots, L$). By checking columns, the elements $q_{i\ell}$ indicate how many cases of each predicted class have been classified as $\ell_i$ ($i = 1, \ldots, L$). Starting from the confusion matrix $Q$ several measures and approaches have been proposed to evaluate classifier performance (accuracy, sensitivity, specificity, etc.). Likewise, the confusion entropy index (Wei et al. 2010), the global performance index (Freitas et al. 2007), the entropy of a confusion matrix (Van Son 1995), the transmitted information of the classifier (Abramson 1963), and the relative classifier information (Sindhwani et al. 2001) are all measures that have been defined in order to compare classifiers performance on the basis of the misclassification cells obtained from confusion matrices. Among all these measures, accuracy is the most known. It refers to the proportion of true results (both true positives and true negatives) among the total number of cases examined. This measure is very plain, overlooking a lot of information about the costs of different elements of misclassification (Hand and Till 2001).

The goal of this paper is to propose a new approach that enables us to compare performances of several classifiers in the framework of multi-class learning (i.e., when a new observation has to be classified into one, and only one, of $L$ non-overlapping classes). The output is a bivariate classifier performance index obtained from two different measures. The first one refers to a cost-sensitive weighted classification accuracy index. The second one refers to an index measuring the similarity in distribution between the $n$ observations which have been classified in one of the $L$ classes by a classifier and the original distribution of the $n$ cases among the $L$ classes. Both indices are defined in $[0, 1] \in \mathbb{R}$, so that a comparison of different classifier performance can be represented in a $[0, 1]^2$ space. Additionally, introducing a measure which is not one-dimensional allows us to study the reliability of each classifier by re-training the classifier on resampled versions of the original data and computing the convex hull of the area obtained in the 2 dimensions in which values of the bivariate classifier performance index are projected.

The rest of the paper is organized as follows. Section 2 presents the main features of the proposed bivariate classifier performance index and describes the three steps characterizing it, while Sect. 3 concentrates on reliability. Section 4 presents the results of the performance of the proposed approach on real data and Sect. 5 ends the paper with some concluding remarks.

## 2 The Bivariate Classifier Performance Index

The bivariate classifier performance index derives from a three steps procedure to be carried out for each candidate classifier. They can be briefly identified with: (1) the model-based measurement of classification accuracy; (2) the measurement of the similarity in distribution between observed classes and predicted ones; (3) the

visualization of the results of the previous steps in order to assess global classifier performance.

## 2.1 Model-Based Measurement of Classification Accuracy

In this section, we present a model-based and cost-sensitive index for measuring accuracy of a multi-class classifier. The basic idea is to use the cells of the observed confusion matrix, i.e., the confusion matrix obtained from training a classifier on the original data, within a regression model in order to derive the estimated cost-sensitive classification accuracy. The regression model is firstly estimated using data obtained from simulated confusion matrices which present the same marginal frequencies of the observed confusion matrix but they refer to situations in which a perfect or random classification is observed. Next, cells of the observed confusion matrix are used together with the estimated regression parameters to derive the value of the index. Let $\pi \in [0, 1]$ be a misclassification level, so that $1 - \pi$ is the classification accuracy level. If $K$ different classifiers are considered, $K$ values of $\pi$ can be observed and those values, defined in $[0, 1]$, can be modeled on the basis of other information related to each classifier. The model specified for $\pi$ allows us to assess classifier performance through a model-based classification accuracy index.

In a regression modeling framework characterized by a continuous response variable $Y$ defined in $[0, 1]$, data are usually transformed in order to map the domain of $Y$ in the real line and then a standard linear regression analysis is applied. This approach has some shortcomings (see Cribari-Neto and Zeileis 2010), such as heteroskedasticity and difficulties in the interpretation of estimated parameters, which are expressed in terms of the transformed variable instead of the original one. Ferrari and Cribari-Neto (2004) proposed a regression model for continuous variables that assumes values in $[0, 1]$, called *Beta Regression Model*. The assumption of this model is that the response variable is beta-distributed, $Y \sim \text{Beta}(a, b)$ with $a, b > 0$. The authors proposed a particular parameterization of the beta density in order to obtain a regression structure for the mean of the response along with a precision parameter. They showed that, through setting $\mu = a/(a + b)$ and $\phi = a + b$, it is possible to express expectation and variance of $Y$ as $E(Y) = \mu$ and $\text{VAR}(Y) = \mu(1 - \mu)/(1 + \phi)$, respectively. The parameter $\phi$ conveys a rate of precision because for larger $\phi$ $\text{VAR}(Y)$ decreases.

The Beta regression model introduced in Ferrari and Cribari-Neto (2004) is applied in the framework of the present study in order to estimate $\pi$ and, indirectly, $1 - \pi$. Specifically, the goal is to estimate a Beta regression model using a large number of simulated confusion matrices weighted by some proximity measures and misclassification costs, in order to obtain estimated regression parameters and associated $\pi$ values. Weighting is very important in this framework, because it conveys essential information to the model about the different importance attributed to possible different misclassification levels. Once the model is estimated, it is applied to the confusion matrix resulting from each classifier in order to estimate a *cost-sensitive (model-based) weighted classification index*. For a classifier $k$ ($k =$

$1, \ldots, K$) and assuming $\pi_k \sim \text{Beta}(\mu_k, \phi)$, the Beta regression model is defined as

$$g(\mu_k) = \sum_{i=1}^{L} \sum_{j=1}^{L} \beta_{ij} q_{ij}^k d(\ell_i, \ell_j) = \eta_k \tag{1}$$

where $d(\ell_i, \ell_j)$ is a cost-weighted proximity measure as defined in Eq. (2), $q_{ij}^k$ is the frequency of the cell of the $i$-th row and $j$-th column of the confusion matrix resulting from the classifier $k$, and $\beta_{ij}$ is the model coefficient that expresses the contribution of $q_{ij}^k$ to global misclassification of classifier $k$. Finally, $g(\cdot)$ is a link function. In Eq. (1) the probit distribution is chosen for specifying the link function $g(\cdot)$, so that the expectation of $\pi_k$ can be defined as $\mu_k = g^{-1}(\eta_k) = \Phi(\eta_k)$, where $\Phi(\cdot)$ is the cumulative distribution function of a standard normal distribution. As already mentioned, for estimating the $\beta_{ij}$ in Eq. (1) a large number $B$ of confusion matrices are simulated. A proportion $\alpha$ with $\pi = 0$ and non-zero elements in the diagonal only, and the other proportion $1-\alpha$ with random assigned elements in order to simulate random classifications, so that $\pi = 1$. A random classified confusion matrix is quite simple to obtain. All confusion matrices stemmed by classifiers have the same marginal row frequencies. In fact, since they come from the same dataset the number of true classes is fixed for all matrices. Hence, it is sufficient to simulate matrices with uniformly distributed rows by setting their marginal row frequencies equal to those of the confusion matrices resulting from the classifiers. Next step consists in excluding diagonal cells from simulated matrices, leaving just cells that convey misclassification information. Additionally, the cells of the simulated confusion matrices are weighted by some proximity measures, which are defined, for all entries $q_{ij}$ (with $i \neq j$) corresponding to off-diagonal elements of confusion matrix, as

$$d(\ell_i, \ell_j) = \begin{cases} \dfrac{\ell_L - \ell_1}{|\ell_i - \ell_j|} w_{ij} & \text{if x is numerical} \\[2ex] \dfrac{L-1}{|i-j|} w_{ij} & \text{if x is ordinal} \\[2ex] w_{ij} & \text{if x is nominal} \end{cases} \tag{2}$$

where $w_{ij}$ is a weight, fixed by the researcher, that specifies the importance in terms of misclassification cost attributed to the proximity level between $\ell_i$ and $\ell_j$. As such, weighting is motivated by the idea of adding information deriving from expert knowledge. Once the simulated matrices are weighted, the model could be fitted through them in order to derive the estimated value $\hat{\mu}_k$ of $\pi_k$ for the $k$-th classifier as

$$\hat{\mu}_k = \Phi \left( \sum_{i=1}^{L} \sum_{j=1}^{L} \hat{\beta}_{ij} q_{ij}^k d\left(\ell_i, \ell_j\right) \right) \tag{3}$$

$\hat{\mu}_k$ is the model-based classification accuracy index used in the rest of the paper.

## 2.2  Similarity in Distribution Index

One of the main problem in the framework of classifier performance measurement is the choice of the best classifier once that two (or more) classifiers present the same value of the classification accuracy $1 - \pi$ but the latter derives from different confusion matrices. To define a classifier performance measure that also considers information about the difference in distribution among classifier confusion matrices, a *normalized similarity in distribution index* is considered. It derives from a dissimilarity index introduced by Gini and used, among others, in Rachev (1985). In general, for a $L$-class classification problem $D$, the Gini index of dissimilarity in distribution, is defined as

$$D = \sqrt{\frac{1}{L^2 - 1} \sum_{h=1}^{L^2-1} |F_h^{v_1} - F_h^{v_2}|^2} \tag{4}$$

where $F_h^{v_1}$ and $F_h^{v_2}$ are the cumulative frequencies in $h$ of the vectors $v_1$ and $v_2$, whereas $\sqrt{L^2 - 1}$ is equal to the maximum value of this index, and it is used to normalize it. $D$ is defined in $[0, 1]$ and is susceptible to change in values as long as one or more observations are assigned to the class $j$ instead of the true class $i$ ($i \neq j$ and $i, j \in \{1, \ldots, L\}$).

  In the framework of the bivariate classifier performance index described so far, the dissimilarity in distribution index introduced in Eq. (4) is reformulated in terms of a similarity in distribution index. To this aim, let us consider two confusion matrices, $Q_{k_1}$ and $Q_{k_2}$, corresponding to classifiers $k_1$ and $k_2$, respectively. They refer to a situation in which the value of classification accuracy is the same for both classifiers, even if the two confusion matrices are clearly different. Measuring similarity between $Q_{k_1}$ and $Q_{k_2}$ requires the comparison of each element of the two matrices with those of a common reference matrix $Q_{max}$. The latter is the matrix which refers to the situation of maximum accuracy so that all predicted values correspond to observed ones. To make such a comparison, the matrices $Q_{max}$, $Q_{k_1}$ and $Q_{k_2}$ are transformed into vectors $v_{max}$, $v_{k_1}$, and $v_{k_2}$ by writing the matrix elements in row-major order. To compute the similarity in distribution for $Q_{k_1}$ and $Q_{k_2}$, it is necessary to compare the distribution of $v_{k_1}$ and $v_{k_2}$ with that of $v_{max}$. Considering the difference $1 - D$, where $D$ has been defined in Eq. (4), we define a similarity in distribution index for $Q_{k_1}$ and $Q_{k_2}$ whose values are in $[0, 1]$ as

$$S_{Q_{k_i}} = 1 - \sqrt{\frac{\sum_{h=1}^{L^2-1} |F_h^{v_{k_i}} - F_h^{v_{max}}|^2}{L^2 - 1}}, \qquad \forall i = 1, 2 \tag{5}$$

## 2.3 Visualization

Once both values of the *cost-sensitive (model-based) weighted classification index* introduced in Sect. 2.1 and the *normalized similarity in distribution index* introduced in Sect. 2.2 are available for each classifier, their values can be projected in a $[0, 1]^2$ space in order to evaluate their performance from the perspective of both classification accuracy and similarity in distribution. The possibility of analyzing classifier performance in a two-dimensional space is very useful since it facilitates the comparison among different classifiers and allows the user to understand which of the two considered items (weighted classification and similarity in distribution) mostly influences classifier performance. Of course, the two-dimensional representation is particularly helpful when the number of considered classifiers is very large.

## 3 Assessing Reliability

Besides measuring the performance of a classifier on the basis of classification accuracy and similarity in distribution, it is very important to define its reliability. The *cost-sensitive (model-based) weighted classification index* can be used to accomplish this goal also. In fact, the measurement of the performance of a classifier can be used as a tool in order to define a measure of its reliability. To this purpose, the basic idea is that applying the same classifier to slightly modified versions of the original data, we expect that its results are rather similar, so that the closer they are to each other the more reliable the classifier can be considered. Thus, the proximity of the results obtained from the same classifier by resampling and measured by the bivariate classifier performance index of Sect. 2 is considered as a measure of classifier performance reliability. Formally, if we have $p$ different measures of classification accuracy of a classifier $k$ (including $\hat{\mu}_k$ and $S_{Q_{k_i}}$) we can measure such a proximity as the convex hull of a set of points $\mathscr{P}$ in $p$ dimensions. The convex hull is computed by measuring the intersection of all convex sets containing $\mathscr{P}$. For $N$ points $p_1, \ldots, p_N$, the convex hull $\mathscr{C}$ is then given by:

$$\mathscr{C} = \left\{ \sum_{j=1}^{N} \lambda_j p_j : \lambda_j \geq 0 \quad \forall j \quad \text{and} \quad \sum_{j=1}^{N} \lambda_j = 1 \right\} \tag{6}$$

In the case of a bivariate index, like the one introduced in Sects. 2.1 and 2.2, this proximity is measured by the convex hull of a set of points defined in the Euclidean space obtained with respect to the two dimensions of the bivariate classifier performance index. In order to obtain this measure of reliability three steps are necessary:

1. Re-train the classifier $B$ times on resampled versions of the original data;

2. Use the resulting *B* confusion matrices as inputs for the two indices measuring cost-weighted accuracy and similarity in distribution;
3. Measure the classifier reliability as the area of the convex hull $\mathscr{C}$ of the set of points $\mathscr{P}$ defined by the values of two indices obtained over the *B* runs.

In our computations, the area of $\mathscr{C}$ is measured with the function `convhulln` implemented in the R package `geometry` (Habel et al. 2014).

## 4 Real Data Example: Classification of Botany Seeds

During the last decades, one of the most important target for botanists is to call a halt to the loss of plant diversity. To achieve that, two strategies are possible: *In situ* and *ex situ* plant conservation. *In situ* conservation consists in protecting threatened plant species in their natural habitat, whereas *ex situ* conservation consists in protecting them outside their natural habitat. Although the *in situ* conservation strategy is considered the best one for preserving plant diversity, its measures are more expensive than *ex situ* ones. For this reason, in the last two decades, the latter conservation approach has been used more often. Among all *ex situ* methods, the most effective is storage of plant seeds in seed banks. It allows us to save large amounts of genetic material in a small space and with minimum risk of genetic damage. Therefore, several seed banks and other structures have been established. Due to the increasing number of seeds gathered, more attention has been focused on classification of accessions in entry. Manual classification of seeds is still a common practice. It is labor-intensive, subjective, and suffers from inconsistencies and errors. It is also a time-consuming task even for highly specialized botanists, and the increasing number of seeds to classify is making the time spent for classification unbearable. For those reasons, application of statistical classifiers for seeds classification is ever more useful and common. Hence botanists require a tool that helps them to evaluate performance and reliability of classifiers, in order to be able to choose among them.

In this study a dataset containing seven variables and $n = 5712$ cases is considered. The response variable is plant family and has five classes (Cyperaceae, Dipsacaceae, Fabaceae, Iridaceae, Lamiaceae). The other six variables are used as predictors and consist in measurements of colorimetric characteristics of seeds. These are the mean of hue, the saturation, the luminance as well as the red channel, green channel, and blue channel intensity.

To measure classification accuracy and reliability the original data were randomly split into two subsets: a proportion of $0.5 \cdot n$ defines the training set and the remaining observations the test set. The experiment involves three different classifiers: CART-like recursive partitioning (CART), Random Forests (RF), and Support Vector Machines (SVM). The choice of these classifiers is based on the consideration that CART is notably known as unstable in terms of reliability of the classification outcome whereas the other two methods are presumably more

reliable and able to provide more accurate classification. The bivariate classification accuracy index and the classifier reliability measured and visualized through the convex hull are used to verify that the approach presented in Sects. 2 and 3 provides new insights for the analyzed dataset.

## 4.1 Results

When classifying botany seeds the goal was to measure the performance and reliability of three classifiers using the approach discussed above. It is worth to remember that the *cost-sensitive (model-based) weighted classification index* is made up of two measures: (1) the model-based measurement of classification accuracy and (2) the measurement of the similarity in distribution between observed classes and predicted ones.

To obtain the cost-sensitive weighted classification accuracy index as defined in Eq. (3) it is necessary to define a proximity measure between each pair of classes of the response variable. To this purpose, observations of the training set are standardized and the proximity is measured as the normalized Euclidean distance between the centroids related to pairs of response classes. Furthermore, for estimating the coefficients of the Beta regression model introduced in Eq. (1), $B = 1000$ confusion matrices were simulated, with a proportion $\alpha = 0.5$ of cases of perfect classification ($\pi = 0$) and the same proportion of cases of random classification ($\pi = 1$). The classifier (CART, SVM , or RF) was trained on the training set observations and predicted classes for the test set observations were used to obtain the confusion matrices, which are the input of the Beta regression model estimated according to the specification introduced in Eq. (1). As for the measurement of the similarity in distribution between observed classes and predicted ones, the Eq. (4) was applied to the three confusion matrices obtained by predicting the response classes of the test set observations for the classifiers CART, RF, and SVM , respectively.

Results are summarized in Table 1, where the two above-mentioned measures are compared with other measures which are frequently used to evaluate the accuracy of a classifier, namely: the proportion of data points in the main diagonal of the confusion matrix; the Rand index and the confusion entropy index (Wei et al. 2010). In order to assess reliability of the three classifiers we used the approach explained in Sect. 3. Firstly, we re-trained each classifier on 100 resampled versions of the training set. Next, we used the 100 confusion matrices obtained from each sample as inputs for the two considered accuracy indexes. Finally, we computed the convex hull $\mathscr{C}$ of the area defined by the values of two indexes obtained over the 100 runs as a measure of reliability.

**Table 1** Accuracy and reliability results for the Random Forest (RF), Support Vector Machine (SVM), and CART-like recursive partitioning classifiers

| Classifier | Diag | Rand | Cen | $(1 - \hat{\pi}_k)$ | $\hat{S}_Q$ | $\mathscr{C}$ |
|---|---|---|---|---|---|---|
| RF | 0.687 | 0.697 | 0.403 | 0.833 | 0.952 | 0.193 |
| | (0.667) | (0.686) | (0.423) | (0.778) | (0.950) | |
| SVM | 0.677 | 0.654 | 0.346 | 0.810 | 0.948 | 0.155 |
| | (0.674) | (0.653) | (0.350) | (0.802) | (0.947) | |
| CART | 0.623 | 0.618 | 0.408 | 0.602 | 0.943 | 0.409 |
| | (0.616) | (0.618) | (0.411) | (0.588) | (0.940) | |

*Notes*: diag is the proportion of data points in the main diagonal of the confusion matrix; rand is the Rand index; cen is the confusion entropy index; $(1 - \hat{\pi}_k)$ is the accuracy measure defined in Eq. (3); $\hat{S}_Q$ is the similarity in distribution as defined in Eq. (4); $\mathscr{C}$ is the reliability of a classifier as defined in Eq. (6). Each cell reports the value of the index obtained for test set observations and, in parentheses, the same value obtained as an average from 100 resampled versions of the original data
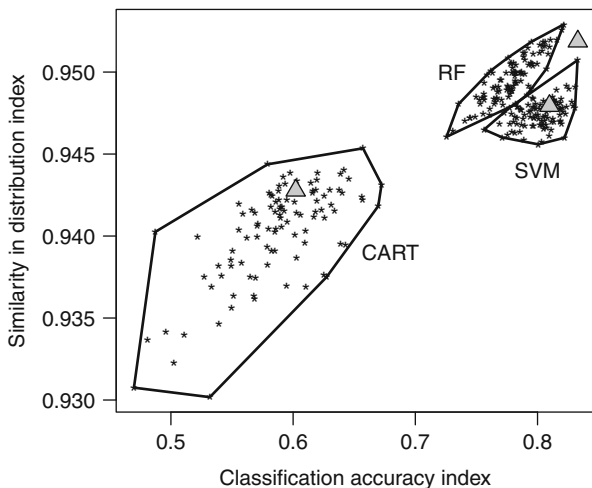


**Fig. 1** Accuracy and reliability of the Random Forests, Support Vector Machines, and CART-like recursive partitioning classifiers. The *triangles* correspond to the cost-sensitive (model-based) weighted classification index and the similarity in distribution index obtained from the original data, whereas the *stars* are values of the same indices obtained on resampled versions of the original data. Reliability is measured through the convex hull of the area defined by each set of points

As it is possible to note from both Table 1 and Fig. 1, Random Forest is the best classifier in this example with respect to accuracy. In fact, it has both the highest classification accuracy (0.833) and the highest similarity in distribution (0.952). In contrast, the most reliable classifier is SVM as it provides the smallest convex hull area ($\mathscr{C} = 0.155$). As expected, CART has to be considered as the worst one for both accuracy and reliability.

## 5    Concluding Remarks

Cost-sensitive classification is one of the mainstream research topics in data mining and machine learning that induces models from data with an unbalanced class distribution and impacts by quantifying and tackling the unbalance. In this paper a bivariate index based on a model-based accuracy measure and a similarity in distribution measure has been introduced. In addition, classifier performance reliability is also considered by computing the convex hull of the set of points in the two-dimensional space defined by the values of the above-mentioned bivariate index computed on resampled versions of the original data. Results obtained for a real data classification problem involving botanic seeds provide evidence about the effectiveness of the proposed approach, since they confirm the expectation that less accurate and less reliable classifiers (CART-like recursive partitioning) do not outperform more robust and accurate ones (SVM and Random Forest). Future research efforts will be directed to the identification and computation of other possible dimensions of accuracy and reliability (like those mentioned in Sect. 1). In addition, following the approach proposed in Müssell et al. (2012), the proposed measures will be framed within the context of Pareto dominance through the visualization of the relative Pareto fronts. Next, our method for measuring cost-sensitive classification accuracy and reliability will be tested on several datasets, with particular attention to multi-class learning problems characterized by an unbalanced distribution of the response classes and/or a reduced data size.

## References

Abramson, N. (1963). *Information theory and coding*. New York: McGraw-Hill.

Cribari-Neto, F., & Zeileis, A. (2010). Beta regression in R. *Journal of Statistical Software, 34*(2), 1–24.

Ferrari, S., & Cribari-Neto, F. (2004). Beta regression for modelling rates and proportions. *Journal of Applied Statistics, 31*(7), 799–815.

Freitas, C. O., De Carvalho, J. M., Oliveira, J. R., Aires, S. B., & Sabourin, R. (2007). Confusion matrix disagreement for multiple classifiers. In *Progress in pattern recognition, image analysis and applications* (pp. 387–396). Berlin: Springer.

Habel, K., Grasman, R., Stahel, A., & Sterrat, D. C. (2014). Geometry: Mesh generation and surface tesselation. R package version 0.3-5, http://CRAN.R-project.org/package=geometry

Hand, D. J., & Till, R. J. (2001). A simple generalisation of the area under the ROC curve for multiple class classification problems. *Machine Learning, 45*(2),171–186.

Müssell, C., Lausser, L., Maucher, M., & Kester, H. A. (2012). Multi-objective parameter selection for classifiers. *Journal of Statistical Software, 46*(5), 1–27.

Rachev, S. T. (1985). The Monge-Kantorovich mass transference problem and its stochastic applications. *Theory of Probability and Its Applications, 29*(4), 647–676.

Sindhwani, V., Bhattacharya, P., & Rakshit, S. (2001). Information theoretic feature crediting in multiclass support vector machines. In *Proceedings of the First SIAM International Conference on Data Mining* (pp. 5–7). Philadelphia, PA. SIAM.

Van Son, R. (1995). A method to quantify the error distribution in confusion matrices. In *Proceedings of Eurospeech 95*, Madrid, 22772280.

Wei, J.-M., Yuan, X.-J., Hu, Q.-H., & Wang, S.-Q. (2010). A novel measure for evaluating classifiers. *Expert Systems with Applications, 37(5),*3799–3809.