# Combination of linear classifiers using score function – analysis of possible combination strategies

Pawel Trajdos and Robert Burduk

Department of Systems and Computer Networks, Wroclaw University of Science and Technology,
Wybrzeze Wyspianskiego 27, 50-370 Wroclaw, Poland
`pawel.trajdos@pwr.edu.pl`

**Abstract.** In this work, we addressed the issue of combining linear classifiers using their score functions. The value of the scoring function depends on the distance from the decision boundary. Two score functions have been tested and four different combination strategies were investigated. During the experimental study, the proposed approach was applied to the heterogeneous ensemble and it was compared to two reference methods – majority voting and model averaging respectively. The comparison was made in terms of seven different quality criteria. The result shows that combination strategies based on simple average, and trimmed average are the best combination strategies of the geometrical combination.

**Keywords:** binary classifiers, linear classifiers, geometrical space, potential function

## 1 Introduction

The combination of multiple base classifiers has been an important issue in machine learning for about twenty years [8], [35]. The ensembles of classifiers (EoC) or multiple classifiers systems (MCSs) [5], [21], [11], [26], [34] are popular in supervised classification algorithms where single classifiers are often unstable (small changes in input data may result in creation of very different decision boundaries) or are often more accurate than any of the base classifiers.

The task of constructing MCSs can be generally divided into three steps: generation, selection and integration [2]. In the first step a set of base classifiers is trained using manipulation of the training patterns, manipulation of the training parameters or manipulation of the feature space.

The second phase of building EoCs is related to the choice of a set or one classifier from the whole available pool of base classifiers. It is popular to use the diversity measure to select one classifier or a subset of all base classifiers. In the literature, there are many approaches to the selection phase of building EoCs [17], [3], [28], [27].

The integration process is the last stage of constructing EoCs and it is widely discussed in the pattern recognition literature [24], [32]. Generally, supervised learning methods produce a classifier whose output is represented as a score function. This function is mapping to a function that is interpreted as a posteriori probability, rank level function or directly as a class label. Depending on the type of mapping, many methods for integrating base classifiers can be distinguished [19], [25], [31].

In this paper we propose the concept of the classifier integration process which uses score functions without their further transformation. In this paper we examined two forms of the score function that is called the potential function and four different combination strategies were investigated.

The remainder of this paper is organized as follows. Section 2 presents the proposed method of EoC integration using two types of the potential function. The experimental evaluation is presented in Section 3. The discussion and conclusions from the experiments are presented in Section 4.

## 2    Proposed Method

In this section, the proposed approach is explained. Additionally, this section introduces the notation used in this paper.

### 2.1    Linear Binary Classifiers

In this paper, it is assumed that the input space $\mathbb{X}$ is a $d-$dimensional Euclidean space $\mathbb{X} = \mathbb{R}^d$. Each object from the input space $x \in \mathbb{X}$ belongs to one of two available classes, so the output space is: $\mathbb{M} = \{-1; 1\}$. It is assumed that there exists an unknown mapping $f : \mathbb{X} \mapsto \mathbb{M}$ that assigns each input space coordinates into a proper class. A classifier $\psi : \mathbb{X} \mapsto \mathbb{M}$ is a function that is designed to provide an approximation of the unknown mapping $f$. A linear classifier makes its decision according to the following rule:

$$\psi(x) = \text{sign}\left(\omega(x)\right), \tag{1}$$

where $\omega(x) = \langle n; x \rangle + b$ is the so called *discriminant function* of the classifier $\psi$ [19], $n$ is a unit normal vector of the decision hyperplane ($\|n\| = 1$), $b$ is the distance from the hyperplane to the origin and $\langle \cdot; \cdot \rangle$ is a dot product defined as follows:

$$\langle a; b \rangle = \sum_{i=1}^{d} a_i b_i, \ \forall a, b \in \mathbb{X}. \tag{2}$$

In this paper, we use a norm of the vector $x$ defined using the dot product:

$$\|x\| = \sqrt{\langle x; x \rangle}. \tag{3}$$

When the normal vector of the plane is a unit vector, the absolute value of the discriminant function equals to the distance from the decision hyperplane to

point $x$. The sign of the discriminant function depends on the site of the plane where the instance $x$ lies.

Now, let us define an ensemble classifier:

$$\Psi = \left\{ \psi^{(1)}, \psi^{(2)}, \cdots, \psi^{(N)} \right\} \tag{4}$$

that is a set of $N$ classifiers that work together in order to produce a more robust result [19]. In this paper, it is assumed that only linear, binary classifiers are employed. There are multiple strategies to combine the classifiers constituting the ensemble. The simplest strategy to combine the outcomes of multiple classifiers is to apply the majority voting scheme [19]:

$$\omega_{\text{MV}}(x) = \sum_{i=1}^{N} \text{sign}(\omega^{(i)}(x)), \tag{5}$$

where $\omega^{(i)}(x)$ is the value of the discriminant function provided by the classifier $\psi^{(i)}$ for point $x$. However, this simple yet effective strategy completely ignores the distance of the instance $x$ from the decision planes.

Another strategy is model averaging [29]. The output of the averaged model may be calculated by simply averaging the values of the discriminant functions:

$$\omega_{\text{MA}}(x) = \frac{1}{N} \sum_{i=1}^{N} \omega^{(i)}(x) \tag{6}$$

After combining the base classifiers, the final prediction of the ensemble is obtained according to the rule (1).
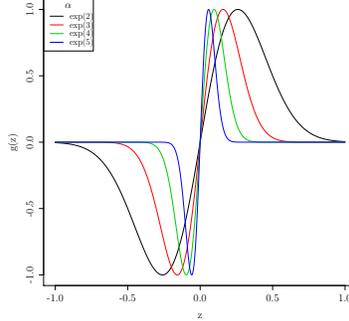
## 2.2 The Proposed Method

In this paper, an approach similar to the softmax [19] normalization is proposed. Contrary to the softmax normalization, our goal is not to provide a probabilistic interpretation of the linear classifier but to provide a fusion technique that works in the geometrical space. The idea is to span a potential field around the decision plane. The potential field may be constructed by applying a transformation on the value of the discriminant function. The transformation must meet the following properties:

$$\text{sign}(g(\omega^{(i)}(x))) = \text{sign}(\omega^{(i)}(x)) \forall x \in \mathbb{X}, \tag{7}$$

$$g(\omega^{(i)}(x)) \in [-1; 1] \, \forall z \in \mathbb{R}, \tag{8}$$

$$g(0) = 0. \tag{9}$$

Property (7) assures that the crisp decision based on the transformed value is the same as the decision based on the unmodified discriminant function. Property (8) bounds $g$ in interval $[-1; 1]$. However, contrary to the softmax normalization the transformation does not have to be a sigmoid function. Property (9) assures that

**Fig. 1.** Potential function $g$.

the potential is 0 at the surface of the decision plane. In this paper, the following transformation function is used:

$$g(z) = z \exp(-\gamma z^2 + 0.5)\sqrt{2\gamma}, \tag{10}$$

where $\gamma$ is a coefficient that determines the position and steepness of the peak. The translation constant 0.5 and the scaling factor $\sqrt{2\gamma}$ guarantee that the maximum and minimum values are 1 and $-1$ respectively. The function is visualised in the figure 1.

All models in the ensemble share the same shape coefficient $\gamma$. The shape coefficient is tuned in order to achieve the best quality of the entire ensemble.

After transforming the values of discriminant functions for the entire ensemble, there is a need to combine the outcomes to produce the final decision. In this paper, we analyze four different combination rules. The first one is a simple average of the transformed values of discriminant functions:

$$\omega_{\text{TA}}(x) = \frac{1}{N} \sum_{i=1}^{N} g(\omega^{(i)}(x)). \tag{11}$$

The other one is to apply the trimmed mean approach:

$$\omega_{\text{TME}}(x) = \frac{1}{N-2} \sum_{i=1}^{N} \left[ g(\omega^{(i)}(x)) - \max_{i \in \{1,2,\cdots,N\}} \omega^{(i)}(x) - \min_{i \in \{1,2,\cdots,N\}} \omega^{(i)}(x) \right]. \tag{12}$$

Before the remaining combination rules are defined, let us introduce subsets of negative and positive values of the transformed ensemble outcomes:

$$\mathcal{G}_-(x) = \left\{ g(\omega^{(i)}(x)) \mid g(\omega^{(i)}(x)) < 0 \right\}, \tag{13}$$

$$\mathcal{G}_+(x) = \left\{ g(\omega^{(i)}(x)) \mid g(\omega^{(i)}(x)) \geq 0 \right\}. \tag{14}$$

Then, the remaining rules are as follows:

$$\omega_{\text{MAX}}(x) = \max(\mathcal{G}_+(x)) + \min(\mathcal{G}_-(x)), \tag{15}$$

$$\omega_{\text{MIN}}(x) = \min(\mathcal{G}_+(x)) + \max(\mathcal{G}_-(x)), \tag{16}$$

$$\omega_{\text{GME}}(x) = \left(\prod_{z \in \mathcal{G}_+(x))} |z|\right)^{|\mathcal{G}_+(x))|^{-1}} - \left(\prod_{z \in \mathcal{G}_-(x))} |z|\right)^{|\mathcal{G}_-(x))|^{-1}}, \tag{17}$$

where $|\mathcal{G}_-(x))|$ and $|z|$ are cardinality of set $\mathcal{G}_-(x))$ and the absolute value of $z$ respectively.

The proposed algorithm is able to deal only with the binary classification problems. However, any multi-class problem can be decomposed into multiple binary problems. In the experimental stage the One-vs-One strategy was used [16]. This strategy builds a separate binary classifier for each pair of classes. In our method, a single pair-specific is replaced by the above-described ensemble classifier.

## 3 Experimental Setup

In the conducted experimental study, the proposed approach was used to combine classifiers in the heterogeneous ensemble of classifier. The following base classifiers were employed:

- $\psi_{\text{FLDA}}$ – Fisher LDA[22]
- $\psi_{\text{MLP}}$ – single layer MLP classifier[12]
- $\psi_{\text{NC}}$ – nearest centroid (Nearest Prototype)[20,18]
- $\psi_{\text{SVM}}$ – SVM classifier with linear kernel (no kernel) [4],
- $\psi_{\text{LR}}$ – logistic regression classifier [7].

The classifiers implemented in WEKA framework [13] were used. The classifier parameters were set to their defaults. The multi-class problems were dealt with using One-vs-One decomposition [16]. The experimental code was implemented using WEKA framework [13].The source code of the algorithms is available on-line [1]. The heterogeneous ensemble employs one copy of each of the above-mentioned base classifiers. Each classifier is learned using the entire dataset.

During the experimental evaluation the following combination methods were compared:

1. $\Psi_{\text{MV}}$ – the ensemble combined using the majority voting approach,
2. $\Psi_{\text{MA}}$ – the ensemble combined using the model averaging approach,
3. $\Psi_{\text{TA}}$ – the ensemble combined using the rule described in (11).
4. $\Psi_{\text{MAX}}$ – the ensemble combined using the rule described in (15).
5. $\Psi_{\text{MIN}}$ – the ensemble combined using the rule described in (16).
6. $\Psi_{\text{TME}}$ – the ensemble combined using the rule described in (17).

---

[1] https://github.com/ptrajdos/piecewiseLinearClassifiers/tree/master

7. $\Psi_{\mathrm{GME}}$ – the ensemble combined using the rule described in (17).

The coefficient $\gamma$ for transformation and $g$ was tuned using the grid search approach. The following set of parameter values were investigated:

$$\{\gamma = \exp(i)|i \in \{2, \cdots, 10\}\}.$$

The parameter is chosen in such a way that it provides the maximum value of the macro-averaged $F_1$ criterion.

To evaluate the proposed methods the following classification-quality criteria are used [30]: Zero-one loss (Accuracy); Macro-averaged FDR, FNR, $F_1$;Micro-averaged FDR, FNR, $F_1$.

Following the recommendations of [6] and [10], the statistical significance of the obtained results was assessed using the two-step procedure. The first step is to perform the Friedman test [9] for each quality criterion separately. Since the multiple criteria were employed, the familywise errors (FWER) should be controlled [36]. To do so, the Bergman-Hommel [1] procedure of controlling FWER of the conducted Friedman tests was employed. When the Friedman test shows that there is a significant difference within the group of classifiers, the pairwise tests, which use the Wilcoxon signed-rank test [33], [6] were employed. To control FWER of the Wilcoxon-testing procedure, the Bergman-Hommel approach was employed [15]. For all the tests the significance level was set to $\alpha = 0.05$.

Table 1 displays the collection of the 64 benchmark sets that were used during the experimental evaluation of the proposed algorithms. The table is divided into two columns. Each column is organized as follows. The first column contains the names of the datasets. The remaining ones contain the set-specific characteristics of the benchmark sets: the number of instances in the dataset ($|S|$); dimensionality of the input space ($d$); the number of classes ($C$);average imbalance ratio (IR).

The datasets come from the Keel [2] repository or are generated by us. The datasets are available online [3].

During the dataset-preprocessing stage, a few transformations on datasets were applied. That is, features are selected using the correlation-based approach [14]. Then, the PCA method was applied [23] and the percentage of variance was set to 0.95. The attributes were also scaled to fit the interval $[0;1]$. Additionally, in order to ensure the dot product to be in the interval $[-1;1]$, vectors in each dataset were scaled using the factor $\frac{1}{d^2}$. This normalization makes it easier to find proper $\gamma$.

## 4   Results and Discussion

To compare multiple algorithms on multiple benchmark sets the average ranks approach [6] is used. In the approach, the winning algorithm achieves rank equal

---

**Table 1.** The characteristics of the benchmark sets

| Name | $|S|$ | $d$ | $C$ | IR | Name | $|S|$ | $d$ | $C$ | IR | Name | $|S|$ | $d$ | $C$ | IR |
|------|------|-----|-----|------|------|------|-----|-----|------|------|------|-----|-----|------|
| appendicitis | 106 | 7 | 2 | 2.52 | housevotes | 435 | 16 | 2 | 1.29 | shuttle | 57999 | 9 | 7 | 1326.03 |
| australian | 690 | 14 | 2 | 1.12 | ionosphere | 351 | 34 | 2 | 1.39 | sonar | 208 | 60 | 2 | 1.07 |
| balance | 625 | 4 | 3 | 2.63 | iris | 150 | 4 | 3 | 1.00 | spambase | 4597 | 57 | 2 | 1.27 |
| banana2D | 2000 | 2 | 2 | 1.00 | led7digit | 500 | 7 | 10 | 1.16 | spectfheart | 267 | 44 | 2 | 2.43 |
| bands | 539 | 19 | 2 | 1.19 | lin1 | 1000 | 2 | 2 | 1.01 | spirals1 | 2000 | 2 | 2 | 1.00 |
| Breast Tissue | 105 | 9 | 6 | 1.29 | lin2 | 1000 | 2 | 2 | 1.83 | spirals2 | 2000 | 2 | 2 | 1.00 |
| check2D | 800 | 2 | 2 | 1.00 | lin3 | 1000 | 2 | 2 | 2.26 | spirals3 | 2000 | 2 | 2 | 1.00 |
| cleveland | 303 | 13 | 5 | 5.17 | magic | 19020 | 10 | 2 | 1.42 | texture | 5500 | 40 | 11 | 1.00 |
| coil2000 | 9822 | 85 | 2 | 8.38 | mfdig fac | 2000 | 216 | 10 | 1.00 | thyroid | 7200 | 21 | 3 | 19.76 |
| dermatology | 366 | 34 | 6 | 2.41 | movement libras | 360 | 90 | 15 | 1.00 | titanic | 2201 | 3 | 2 | 1.55 |
| diabetes | 768 | 8 | 2 | 1.43 | newthyroid | 215 | 5 | 3 | 3.43 | twonorm | 7400 | 20 | 2 | 1.00 |
| Faults | 1940 | 27 | 7 | 4.83 | optdigits | 5620 | 62 | 10 | 1.02 | ULC | 675 | 146 | 9 | 2.17 |
| gauss2DV | 800 | 2 | 2 | 1.00 | page-blocks | 5472 | 10 | 5 | 58.12 | vehicle | 846 | 18 | 4 | 1.03 |
| gauss2D | 4000 | 2 | 2 | 1.00 | penbased | 10992 | 16 | 10 | 1.04 | Vertebral Column | 310 | 6 | 3 | 1.67 |
| gaussSand2 | 600 | 2 | 2 | 1.50 | phoneme | 5404 | 5 | 2 | 1.70 | wdbc | 569 | 30 | 2 | 1.34 |
| gaussSand | 600 | 2 | 2 | 1.50 | pima | 767 | 8 | 2 | 1.44 | wine | 178 | 13 | 3 | 1.23 |
| glass | 214 | 9 | 6 | 3.91 | ring2D | 4000 | 2 | 2 | 1.00 | winequality-red | 1599 | 11 | 6 | 20.71 |
| haberman | 306 | 3 | 2 | 1.89 | ring | 7400 | 20 | 2 | 1.01 | winequality-white | 4898 | 11 | 7 | 82.94 |
| halfRings1 | 400 | 2 | 2 | 1.00 | saheart | 462 | 9 | 2 | 1.44 | wisconsin | 699 | 9 | 2 | 1.45 |
| halfRings2 | 600 | 2 | 2 | 1.50 | satimage | 6435 | 36 | 6 | 1.66 | yeast | 1484 | 8 | 10 | 17.08 |
| hepatitis | 155 | 19 | 2 | 2.42 | Seeds | 210 | 7 | 3 | 1.00 | | | | | |
| HillVall | 1212 | 100 | 2 | 1.01 | segment | 2310 | 19 | 7 | 1.00 | | | | | |

'1', the second achieves rank equal '2', and so on. In the case of ties, the ranks of algorithms that achieve the same results, are averaged. To provide a visualisation of the average ranks, the radar plots are employed. In the plots, the data is visualised in such a way that the lowest ranks are closer to the centre of the graph. The radar plots related to the experimental results are shown in figure 2.

Due to the page limit, the full results are published online [4]

The numerical results are given in Table 2. The table is structured as follows. The first row contains names of the investigated algorithms. Then, the table is divided into seven sections – one section is related to a single evaluation criterion. The first row of each section is the name of the quality criterion investigated in the section. The second row shows the p-value of the Friedman test. The third one shows the average ranks achieved by algorithms. The following rows show p-values resulting from pairwise Wilcoxon test. The p-value which is equal to 0.000 informs that the p-values are lower than $10^{-3}$ and p-value is equal to 1.000 informs that the value is higher than 0.999.

The analysis of the radar plot suggests that two groups of classification criteria can be distinguished. The first group contains micro-averaged criteria and the zero-one criterion, the second one contains macro-averaged criteria. Evaluation of the classifiers carried out with the use of criteria belonging to a specific group reveals different relationships between classifiers. These differences are a consequence of the properties of the quality criteria used. This means that the zero-one criterion and micro-averaged criteria give us information related to the classification quality for the majority classes. On the other hand, the macro-averaged criteria put more emphasis on classification quality for minority classes [30].

For the zero-one criterion and micro-averaged criteria, three main groups of classifiers can be seen. The first group contains $\Psi_{\text{MIN}}$ and $\Psi_{\text{GME}}$ classifiers that perform significantly worse than the other analysed classifiers. What is more, classifier $\Psi_{\text{MIN}}$ is significantly worse than $\Psi_{\text{GME}}$ for all quality criteria belonging to the investigated group. The second group contains only one classifier – $\Psi_{\text{MV}}$.

---

[4] https://github.com/ptrajdos/MLResults/blob/master/Boundaries/bounds_hetero_15.01.2019E4_m_R.zip

According to average ranks, this classifier is the best performing one for the investigated set of quality criteria. According to the statistical analysis, this classifier outperforms the remaining classifiers except for $\Psi_{\mathrm{TA}}$ and $\Psi_{\mathrm{TME}}$. The third group consisted of classifiers $\Psi_{\mathrm{MA}}$, $\Psi_{\mathrm{TA}}$, $\Psi_{\mathrm{MAX}}$, and $\Psi_{\mathrm{TME}}$. There are no significant differences between the classifiers within this group.

For macro-averaged measures, the situation changes significantly. First of all, it may be noticed that average ranks of reference methods ($\Psi_{\mathrm{MV}}$ and $\Psi_{\mathrm{MA}}$) increase, whereas the average ranks of the proposed methods decrease. That is, the model-averaging classifier $\Psi_{\mathrm{MA}}$ becomes the worst one except for $\Psi_{\mathrm{MIN}}$ according to macro-averaged $F_1$ and FNR criteria. The majority voting classifier $\Psi_{\mathrm{MV}}$ also deteriorates significantly. Now it is comparable to $\Psi_{\mathrm{MAX}}$, $\Psi_{\mathrm{MIN}}$ and $\Psi_{\mathrm{GME}}$ classifiers. What is more, $\Psi_{\mathrm{MV}}$ classifier is outperformed by $\Psi_{\mathrm{TA}}$ and $\Psi_{\mathrm{TME}}$ classifiers in terms of macro-averaged FNR and $F_1$ criteria. The reason for the above-mentioned deterioration of the reference methods is the fact that they are not tuned to perform better on minority classes, whereas the investigated methods were tuned to do so.

Now let us investigate the differences inside the group of the proposed combination criteria. First of all, classifiers $\Psi_{\mathrm{TA}}$ and $\Psi_{\mathrm{TME}}$ offer the best classification quality under macro-averaged $F_1$ measure. It means that these classifiers offer the best trade-off between macro-averaged precision and recall. Under macro-averaged FDR ($1 - \text{precision}$) measure, these algorithms outperform only $\Psi_{\mathrm{MIN}}$ and $\Psi_{\mathrm{GME}}$ classifiers. For macro-averaged FNR ($1 - \text{recall}$) the investigated classifiers outperform all but $\Psi_{\mathrm{MIN}}$ classifiers. On the other hand, under the macro-averaged measures, there are no significant differences between $\Psi_{\mathrm{TA}}$ and $\Psi_{\mathrm{TME}}$.

**Table 2.** Statistical evaluation. Wilcoxon test for the heterogeneous ensemble – p-values for paired comparisons of the investigated methods.

**Zero-One**

| | $\Psi_{\mathrm{MV}}$ | $\Psi_{\mathrm{MA}}$ | $\Psi_{\mathrm{TA}}$ | $\Psi_{\mathrm{MAX}}$ | $\Psi_{\mathrm{MIN}}$ | $\Psi_{\mathrm{TME}}$ | $\Psi_{\mathrm{GME}}$ |
|---|---|---|---|---|---|---|---|
| Frd | 5.729e-14 | | | | | | |
| Rnk | 2.98 | 3.78 | 3.36 | 3.73 | 5.72 | 3.56 | 4.87 |
| $\Psi_{\mathrm{MV}}$ | | .007 | .091 | .002 | .000 | .161 | .000 |
| $\Psi_{\mathrm{MA}}$ | | | .968 | .968 | .000 | .968 | .007 |
| $\Psi_{\mathrm{TA}}$ | | | | .080 | .968 | .968 | .000 |
| $\Psi_{\mathrm{MAX}}$ | | | | | .000 | .846 | .000 |
| $\Psi_{\mathrm{MIN}}$ | | | | | | .000 | .000 |
| $\Psi_{\mathrm{TME}}$ | | | | | | | .000 |

**MaFDR**

| | $\Psi_{\mathrm{MV}}$ | $\Psi_{\mathrm{MA}}$ | $\Psi_{\mathrm{TA}}$ | $\Psi_{\mathrm{MAX}}$ | $\Psi_{\mathrm{MIN}}$ | $\Psi_{\mathrm{TME}}$ | $\Psi_{\mathrm{GME}}$ |
|---|---|---|---|---|---|---|---|
| Frd | 2.873e-04 | | | | | | |
| Rnk | 3.76 | 4.45 | 3.41 | 3.93 | 4.75 | 3.31 | 4.39 |
| $\Psi_{\mathrm{MV}}$ | | .016 | .295 | .969 | .155 | .279 | .673 |
| $\Psi_{\mathrm{MA}}$ | | | .001 | .025 | .878 | .002 | .295 |
| $\Psi_{\mathrm{TA}}$ | | | | .056 | .013 | .878 | .025 |
| $\Psi_{\mathrm{MAX}}$ | | | | | .028 | .056 | .155 |
| $\Psi_{\mathrm{MIN}}$ | | | | | | .004 | .295 |
| $\Psi_{\mathrm{TME}}$ | | | | | | | .003 |

**MaFNR**

| | $\Psi_{\mathrm{MV}}$ | $\Psi_{\mathrm{MA}}$ | $\Psi_{\mathrm{TA}}$ | $\Psi_{\mathrm{MAX}}$ | $\Psi_{\mathrm{MIN}}$ | $\Psi_{\mathrm{TME}}$ | $\Psi_{\mathrm{GME}}$ |
|---|---|---|---|---|---|---|---|
| Frd | 1.791e-08 | | | | | | |
| Rnk | 4.27 | 5.32 | 3.32 | 3.58 | 4.00 | 3.09 | 4.42 |
| $\Psi_{\mathrm{MV}}$ | | .000 | .003 | .505 | 1.00 | .000 | 1.00 |
| $\Psi_{\mathrm{MA}}$ | | | .000 | .000 | .018 | .000 | .002 |
| $\Psi_{\mathrm{TA}}$ | | | | .049 | .139 | 1.00 | .008 |
| $\Psi_{\mathrm{MAX}}$ | | | | | .601 | .049 | .016 |
| $\Psi_{\mathrm{MIN}}$ | | | | | | .139 | 1.00 |
| $\Psi_{\mathrm{TME}}$ | | | | | | | .001 |

**MaF1**

| | $\Psi_{\mathrm{MV}}$ | $\Psi_{\mathrm{MA}}$ | $\Psi_{\mathrm{TA}}$ | $\Psi_{\mathrm{MAX}}$ | $\Psi_{\mathrm{MIN}}$ | $\Psi_{\mathrm{TME}}$ | $\Psi_{\mathrm{GME}}$ |
|---|---|---|---|---|---|---|---|
| Frd | 2.641e-09 | | | | | | |
| Rnk | 3.96 | 5.10 | 3.23 | 3.59 | 4.81 | 2.96 | 4.35 |
| $\Psi_{\mathrm{MV}}$ | | .000 | .017 | .548 | .117 | .000 | .340 |
| $\Psi_{\mathrm{MA}}$ | | | .000 | .000 | .315 | .000 | .017 |
| $\Psi_{\mathrm{TA}}$ | | | | .017 | .002 | .454 | .001 |
| $\Psi_{\mathrm{MAX}}$ | | | | | .007 | .014 | .011 |
| $\Psi_{\mathrm{MIN}}$ | | | | | | .000 | .185 |
| $\Psi_{\mathrm{TME}}$ | | | | | | | .000 |

**MiFDR**

| | $\Psi_{\mathrm{MV}}$ | $\Psi_{\mathrm{MA}}$ | $\Psi_{\mathrm{TA}}$ | $\Psi_{\mathrm{MAX}}$ | $\Psi_{\mathrm{MIN}}$ | $\Psi_{\mathrm{TME}}$ | $\Psi_{\mathrm{GME}}$ |
|---|---|---|---|---|---|---|---|
| Frd | 5.729e-14 | | | | | | |
| Rnk | 2.98 | 3.78 | 3.36 | 3.73 | 5.72 | 3.56 | 4.87 |
| $\Psi_{\mathrm{MV}}$ | | .007 | .091 | .002 | .000 | .161 | .000 |
| $\Psi_{\mathrm{MA}}$ | | | .968 | .968 | .000 | .968 | .007 |
| $\Psi_{\mathrm{TA}}$ | | | | .080 | .000 | .968 | .000 |
| $\Psi_{\mathrm{MAX}}$ | | | | | .000 | .846 | .000 |
| $\Psi_{\mathrm{MIN}}$ | | | | | | .000 | .000 |
| $\Psi_{\mathrm{TME}}$ | | | | | | | .000 |

**MiFNR**

| | $\Psi_{\mathrm{MV}}$ | $\Psi_{\mathrm{MA}}$ | $\Psi_{\mathrm{TA}}$ | $\Psi_{\mathrm{MAX}}$ | $\Psi_{\mathrm{MIN}}$ | $\Psi_{\mathrm{TME}}$ | $\Psi_{\mathrm{GME}}$ |
|---|---|---|---|---|---|---|---|
| Frd | 5.729e-14 | | | | | | |
| Rnk | 2.98 | 3.78 | 3.36 | 3.73 | 5.72 | 3.56 | 4.87 |
| $\Psi_{\mathrm{MV}}$ | | .007 | .091 | .002 | .000 | .161 | .000 |
| $\Psi_{\mathrm{MA}}$ | | | .968 | .968 | .000 | .968 | .007 |
| $\Psi_{\mathrm{TA}}$ | | | | .080 | .000 | .968 | .000 |
| $\Psi_{\mathrm{MAX}}$ | | | | | .000 | .846 | .000 |
| $\Psi_{\mathrm{MIN}}$ | | | | | | .000 | .000 |
| $\Psi_{\mathrm{TME}}$ | | | | | | | .000 |

**MiF1**

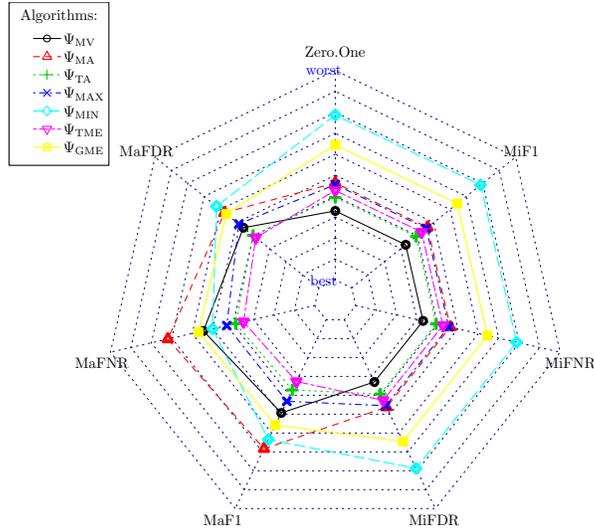| | $\Psi_{\mathrm{MV}}$ | $\Psi_{\mathrm{MA}}$ | $\Psi_{\mathrm{TA}}$ | $\Psi_{\mathrm{MAX}}$ | $\Psi_{\mathrm{MIN}}$ | $\Psi_{\mathrm{TME}}$ | $\Psi_{\mathrm{GME}}$ |
|---|---|---|---|---|---|---|---|
| Frd | 5.729e-14 | | | | | | |
| Rnk | 2.98 | 3.78 | 3.36 | 3.73 | 5.72 | 3.56 | 4.87 |
| $\Psi_{\mathrm{MV}}$ | | .007 | .091 | .002 | .000 | .161 | .000 |
| $\Psi_{\mathrm{MA}}$ | | | .968 | .968 | .000 | .968 | .007 |
| $\Psi_{\mathrm{TA}}$ | | | | .080 | .000 | .968 | .000 |
| $\Psi_{\mathrm{MAX}}$ | | | | | .000 | .846 | .000 |
| $\Psi_{\mathrm{MIN}}$ | | | | | | .000 | .000 |
| $\Psi_{\mathrm{TME}}$ | | | | | | | .000 |

**Fig. 2.** Average ranks of for the heterogeneous ensemble.

## 5   Conclusions

In this paper, a geometric combination scheme was proposed. Four different methods of producing the final output of EoC were investigated. The goal of this paper is to determine the best combination strategy for the given potential-function-induced geometrical space. The experimental comparison shows that $\Psi_{TA}$ and $\Psi_{TME}$ algorithms are the best choice. This is because under macro-averaged measures they are outperforming the other proposed strategies and reference methods. What is more, under the micro-averaged criteria they are comparable to the majority voting procedure. According to the outcome of the statistical evaluation, these algorithms perform equally well. However, under macro-averaged measures, $\Psi_{TME}$ achieves a slightly lower average rank. This suggests that $\Psi_{TME}$ may be slightly better since the truncated mean combination rule removes extreme values of the potential function so it may be less influenced by outliers.

The obtained results are very interesting, so we are willing to continue our research in the field of combining classifiers in the geometrical space. An interesting direction to explore may be the application of the potential function whose shape is not given arbitrary but is created considering data distribution.

# References

1. Bergmann, B., Hommel, G.: Improvements of general multiple test procedures for redundant systems of hypotheses. In: Multiple Hypothesenprüfung / Multiple Hypotheses Testing, pp. 100–115. Springer Berlin Heidelberg (1988). https://doi.org/10.1007/978-3-642-52307-6_8, https://doi.org/10.1007/978-3-642-52307-6_8

2. Britto, A.S., Sabourin, R., Oliveira, L.E.: Dynamic selection of classifiers——a comprehensive review. Pattern Recognition **47**(11), 3665–3680 (2014)

3. Burduk, R., Walkowiak, K.: Static classifier selection with interval weights of base classifiers. In: Asian Conference on Intelligent Information and Database Systems. pp. 494–502. Springer (2015)

4. Cortes, C., Vapnik, V.: Support-vector networks. Mach Learn **20**(3), 273–297 (Sep 1995). https://doi.org/10.1007/bf00994018

5. Cyganek, B.: One-class support vector ensembles for image segmentation and classification. Journal of Mathematical Imaging and Vision **42**(2-3), 103–117 (2012)

6. Demšar, J.: Statistical comparisons of classifiers over multiple data sets. The Journal of Machine Learning Research **7**, 1–30 (2006)

7. Devroye, L., Györfi, L., Lugosi, G.: A Probabilistic Theory of Pattern Recognition. Springer New York (1996). https://doi.org/10.1007/978-1-4612-0711-5, http://dx.doi.org/10.1007/978-1-4612-0711-5

8. Drucker, H., Cortes, C., Jackel, L.D., LeCun, Y., Vapnik, V.: Boosting and other ensemble methods. Neural Computation **6**(6), 1289–1301 (1994)

9. Friedman, M.: A comparison of alternative tests of significance for the problem of $m$ rankings. Ann. Math. Statist. **11**(1), 86–92 (Mar 1940). https://doi.org/10.1214/aoms/1177731944

10. Garcia, S., Herrera, F.: An extension on"statistical comparisons of classifiers over multiple data sets"for all pairwise comparisons. Journal of Machine Learning Research **9**, 2677–2694 (Dec 2008)

11. Giacinto, G., Roli, F.: An approach to the automatic design of multiple classifier systems. Pattern Recognition Letters **22**, 25–33 (2001)

12. Gurney, K.: An introduction to neural networks. Taylor & Francis, London (1997). https://doi.org/10.4324/9780203451519

13. Hall, M., Frank, E., Holmes, G., Pfahringer, B., Reutemann, P., Witten, I.H.: The WEKA data mining software. SIGKDD Explor. Newsl. **11**(1), 10 (Nov 2009). https://doi.org/10.1145/1656274.1656278

14. Hall, M.A.: Correlation-based feature selection for machine learning. Ph.D. thesis, The University of Waikato (1999)

15. Holm, S.: A Simple Sequentially Rejective Multiple Test Procedure. Scandinavian Journal of Statistics **6**(2), 65–70 (1979). https://doi.org/10.2307/4615733

16. Hüllermeier, E., Fürnkranz, J.: On predictive accuracy and risk minimization in pairwise label ranking. Journal of Computer and System Sciences **76**(1), 49–62 (feb 2010). https://doi.org/10.1016/j.jcss.2009.05.005

17. Ko, A.H., Sabourin, R., Britto Jr, A.S.: From dynamic classifier selection to dynamic ensemble selection. Pattern recognition **41**(5), 1718–1731 (2008)

18. Kuncheva, L., Bezdek, J.: Nearest prototype classification: clustering, genetic algorithms, or random search? IEEE Transactions on Systems, Man and Cybernetics, Part C (Applications and Reviews) **28**(1), 160–164 (1998). https://doi.org/10.1109/5326.661099, https://doi.org/10.1109/5326.661099

19. Kuncheva, L.I.: Combining Pattern Classifiers: Methods and Algorithms. Wiley-Interscience, 1 edn. (Jul 2004)
20. Manning, C.D., Raghavan, P., Schutze, H.: Introduction to Information Retrieval. Cambridge University Press, New York (2008). https://doi.org/10.1017/cbo9780511809071
21. Markiewicz, A., Forczmański, P.: Detection and classification of interesting parts in scanned documents by means of adaboost classification and low-level features verification. In: International Conference on Computer Analysis of Images and Patterns. pp. 529–540. Springer (2015)
22. McLachlan, G.J.: Discriminant Analysis and Statistical Pattern Recognition. Wiley Series in Probability and Mathematical Statistics: Applied Probability and Statistics, John Wiley & Sons, Inc. (Mar 1992). https://doi.org/10.1002/0471725293, a Wiley-Interscience Publication
23. Pearson, K.: LIII. on lines and planes of closest fit to systems of points in space. The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science $2$(11), 559–572 (Nov 1901). https://doi.org/10.1080/14786440109462720
24. Ponti Jr, M.P.: Combining classifiers: from the creation of ensembles to the decision fusion. In: Graphics, Patterns and Images Tutorials (SIBGRAPI-T), 2011 24th SIBGRAPI Conference on. pp. 1–10. IEEE (2011)
25. Przybyła-Kasperek, M., Wakulicz-Deja, A.: Comparison of fusion methods from the abstract level and the rank level in a dispersed decision-making system. International Journal of General Systems $46$(4), 386–413 (2017)
26. Przybyła-Kasperek, M.: Three conflict methods in multiple classifiers that use dispersed knowledge. International Journal of Information Technology & Decision Making pp. 1–45 (2018)
27. Reif, M., Shafait, F., Goldstein, M., Breuel, T., Dengel, A.: Automatic classifier selection for non-experts. Pattern Analysis and Applications $17$(1), 83–96 (2014)
28. Rejer, I., Burduk, R.: Classifier selection for motor imagery brain computer interface. In: IFIP International Conference on Computer Information Systems and Industrial Management. pp. 122–130. Springer (2017)
29. Skurichina, M., Duin, R.P.: Bagging for linear classifiers. Pattern Recognit. $31$(7), 909–930 (Jul 1998). https://doi.org/10.1016/s0031-3203(97)00110-6, `https://doi.org/10.1016/s0031-3203(97)00110-6`
30. Sokolova, M., Lapalme, G.: A systematic analysis of performance measures for classification tasks. Information Processing & Management $45$(4) (Jul 2009). https://doi.org/10.1016/j.ipm.2009.03.002, `http://dx.doi.org/10.1016/j.ipm.2009.03.002`
31. Trawiński, B., Lasota, T., Kempa, O., Telec, Z., Kutrzyński, M.: Comparison of ensemble learning models with expert algorithms designed for a property valuation system. In: International Conference on Computational Collective Intelligence. pp. 317–327. Springer (2017)
32. Tulyakov, S., Jaeger, S., Govindaraju, V., Doermann, D.: Review of classifier combination methods. In: Machine Learning in Document Analysis and Recognition, pp. 361–386. Springer (2008)
33. Wilcoxon, F.: Individual comparisons by ranking methods. Biometrics Bulletin $1$(6), 80 (Dec 1945). https://doi.org/10.2307/3001968
34. Woźniak, M., Graña, M., Corchado, E.: A survey of multiple classifier systems as hybrid systems. Information Fusion $16$, 3–17 (2014)
35. Xu, L., Krzyzak, A., Suen, C.Y.: Methods of combining multiple classifiers and their applications to handwriting recognition. IEEE transactions on systems, man, and cybernetics $22$(3), 418–435 (1992)

36. Yekutieli, D., Benjamini, Y.: The control of the false discovery rate in multiple testing under dependency. Ann. Statist. **29**(4), 1165–1188 (Aug 2001). https://doi.org/10.1214/aos/1013699998, `https://doi.org/10.1214/aos/1013699998`