

A nearest-neighbor-based ensemble classifier and its large-sample optimality

Mojirsheibani, Majid; Pouliot, William

DOI:

[10.1080/00949655.2021.1882458](https://doi.org/10.1080/00949655.2021.1882458)

License:

Creative Commons: Attribution-NonCommercial-NoDerivs (CC BY-NC-ND)

Document Version

Peer reviewed version

Citation for published version (Harvard):

Mojirsheibani, M & Pouliot, W 2021, 'A nearest-neighbor-based ensemble classifier and its large-sample optimality', *Journal of Statistical Computation and Simulation*, vol. 91, no. 10, pp. 1-17.
<https://doi.org/10.1080/00949655.2021.1882458>

[Link to publication on Research at Birmingham portal](#)

Publisher Rights Statement:

This is an Accepted Manuscript version of the following article, accepted for publication in *Journal of Statistical Computation and Simulation*. Majid Mojirsheibani & William Pouliot (2021) A nearest-neighbor-based ensemble classifier and its large-sample optimality, *Journal of Statistical Computation and Simulation*, DOI: 10.1080/00949655.2021.1882458. It is deposited under the terms of the Creative Commons Attribution-NonCommercial-NoDerivatives License (<http://creativecommons.org/licenses/by-nc-nd/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is properly cited, and is not altered, transformed, or built upon in any way.

General rights

Unless a licence is specified above, all rights (including copyright and moral rights) in this document are retained by the authors and/or the copyright holders. The express permission of the copyright holder must be obtained for any use of this material other than for purposes permitted by law.

- Users may freely distribute the URL that is used to identify this publication.
- Users may download and/or print one copy of the publication from the University of Birmingham research portal for the purpose of private study or non-commercial research.
- User may use extracts from the document in line with the concept of 'fair dealing' under the Copyright, Designs and Patents Act 1988 (?)
- Users may not further distribute the material nor use it for the purposes of commercial gain.

Where a licence is displayed above, please note the terms and conditions of the licence govern your use of this document.

When citing, please reference the published version.

Take down policy

While the University of Birmingham exercises care and attention in making items available there are rare occasions when an item has been uploaded in error or has been deemed to be commercially or otherwise sensitive.

If you believe that this is the case for this document, please contact UBIRA@lists.bham.ac.uk providing details and we will remove access to the work immediately and investigate.

A nearest-neighbor-based ensemble classifier and its large-sample optimality

Majid Mojirsheibani¹ and William Pouliot²

Department of Mathematics, California State University Northridge, CA, 91330, USA¹

Department of Economics, University of Birmingham, B15 2TT, UK²

Abstract

A nonparametric approach is proposed to combine several individual classifiers in order to construct an asymptotically more accurate classification rule in the sense that its misclassification error rate is, asymptotically, at least as low as that of the best individual classifier. The proposed method uses a nearest neighbor type approach to estimate the conditional expectation of the class associated with a new observation (conditional on the vector of individual predictions). Both mechanics and the theoretical validity of the proposed approach are discussed. As an interesting byproduct of our results, it is shown that the proposed method can also be applied to any single classifier in which case the resulting new classifier will be at least as good as the original one. Several numerical examples, involving both real and simulated data, are also given. These numerical studies further confirm the superiority of the proposed classifier.

Keywords: Nonparametric, asymptotics, classification.

1 Introduction

Consider the following standard two-group classification problem. Let (χ, Y) be a random pair, where χ takes values in a metric space (\mathcal{F}, ρ) and $Y \in \{0, 1\}$, called the class label, must be predicted based on χ . Here, \mathcal{F} is not necessarily confined to \mathbb{R}^d . In classification one wants to find a classifier, $g : \mathcal{F} \rightarrow \{0, 1\}$, whose misclassification error, $\mathbb{P}\{g(\chi) \neq Y\}$, is as small as possible. The classifier with the lowest misclassification error, called the Bayes classifier, is given by $g_B(\chi) = 1$ if $\mathbb{P}\{Y = 1 | \chi = \chi\} > 1/2$, and

¹Corresponding author. Email: majid.mojirsheibani@csun.edu

This work was supported by the NSF under Grant DMS-1916161 of Majid Mojirsheibani.

²Email: w.pouliot@bham.ac.uk

$g_{\mathbb{B}}(\chi) = 0$ otherwise; see, for example, Cérou and Guyader [9], Abraham et al. [1], and Devroye, et al. [12]. Although our setup is expressed for the popular two-class problem, all the results in this paper can be extended in a straightforward manner to the multi-class classification problem; see Remark 2.

Of course, in practice the distribution of (χ, Y) is virtually always unknown and, typically, one only has access to a training sample of n independent and identically distributed (iid) observations $\mathbb{T}_n = \{(\chi_1, Y_1), \dots, (\chi_n, Y_n)\}$ from the distribution of (χ, Y) . Much of the theory of classification deals with the construction of sample-based classification rules \hat{g}_n based on \mathbb{T}_n whose error rates are somehow as small as possible. Of course, the choice of \hat{g}_n is at the practitioner's discretion and there may be several different options; therefore, let $\hat{g}_{n,1}, \dots, \hat{g}_{n,J}$ be $J \geq 1$ different classification rules for predicting Y . Here, $\hat{g}_{n,1}$ may be a linear classifier, $\hat{g}_{n,2}$ a kernel classifier, while $\hat{g}_{n,3}$ may be Breiman's [8] random forest classifier, etc. The aim is then to combine these individual classifiers in such a way that the resulting classifier is in some sense at least as good as the best individual classifier.

There is a vast literature on combined or ensemble methods and there are many different approaches available; this is particularly true for the important problems of classification and regression function estimation. One may divide the existing methods into roughly two types: (a) those approaches that involve a large number of similar or homogeneous base models. Relevant examples here include Breiman's [7, 8] random forest, Lin and Jeon [17], Biau et al. [3], and Rahman et al. [19]. (b) Those approaches that combine a number of models or estimators that are constructed based on different theories or estimation methods. Results under (b) include Fischer and Mougeot [13], Biau et al. [4], Cholaquidis et al. [10], Balakrishnan and Mojirsheibani [2], Mojirsheibani [18], and LeBlanc and Tibshirani [16]. The methods employed in the cited papers under (b) are mainly nonlinear in nature, which is also the framework of the current paper. There is also a large body of literature on linear and convex aggregation methods; in fact, Chapter

3 of the monograph by Giraud [14] presents a detailed account of such methods along with many relevant references.

In passing, we also note that there are other taxonomies for characterizing combined classifiers. In fact, as explained in Rokach [20, 21], combined classification methods can be put into two main categories: *weighting* methods and *non-weighting* or *meta-learning*. Popular weighting methods include the majority voting employed by Breiman [7, 8] in the context of tree classification, and by Xu et al. [24] for the problem of handwriting recognition. Weighted-averaging of estimated class conditional probabilities that are produced by each classifier, has also been studied by several researchers; results along these lines include the work of Xu et al. [24], Breiman [6], and LeBlanc and Tibshirani [16]. There are also other weighting methods that can be found in Rokach [20, 21]. Meta-learning methods typically work by using the predicted values of the individual classifiers on the data. Relevant results along these lines include the stacked generalization of Wolpert [23], Breiman's [6] stacked method, and the nonlinear methods of Mojirsheibani [18], Balakrishnan and Mojirsheibani [2], Biau et al. [4], and Cholaquidis et al. [10]. For more on meta-learning methods, one may refer to Rokach [20].

In the next section, we consider the problem of combining several individual classifiers in such a way that the resulting ensemble is, asymptotically, at least as good as the best individual one. The paper is organized as follows. Section 2 presents the main results, where both the mechanics and the theoretical validity of the proposed approach are discussed. Numerical studies involving both simulated as well as real data are carried out in section 3; these studies further confirm the good finite-sample performance of the proposed approach. All proofs are deferred to section 4.

2 Main results

In order to motivate our proposed method, consider the following hypothetical oversimplified setup. Let g_1, \dots, g_J be J classifiers for predicting Y based on $\boldsymbol{\chi}$ (no data yet). Here, each g_j is a map of the form $g_j : \mathcal{F} \rightarrow \{0, 1\}$. Define the combined classifier $G^* : \{0, 1\}^J \rightarrow \{0, 1\}$, for predicting the same Y , by

$$G^*(g_1(\boldsymbol{\chi}), \dots, g_J(\boldsymbol{\chi})) = \begin{cases} 1 & \text{if } \mathbb{E}[(2Y - 1) | g_1(\boldsymbol{\chi}), \dots, g_J(\boldsymbol{\chi})] > 0 \\ 0 & \text{otherwise.} \end{cases} \quad (1)$$

Then we have the following elementary result

Proposition 1 *The combined classifier G^* in (1) is optimal, i.e.,*

$$\mathbb{P}\{G^*(g_1(\boldsymbol{\chi}), \dots, g_J(\boldsymbol{\chi})) \neq Y\} = \inf_{G: \{0,1\}^J \rightarrow \{0,1\}} \mathbb{P}\{G(g_1(\boldsymbol{\chi}), \dots, g_J(\boldsymbol{\chi})) \neq Y\}.$$

We also observe that in view of Proposition 1, and without further ado, one has

$$\mathbb{P}\{G^*(g_1(\boldsymbol{\chi}), \dots, g_J(\boldsymbol{\chi})) \neq Y\} \leq \min_{1 \leq j \leq J} \mathbb{P}\{g_j(\boldsymbol{\chi}) \neq Y\}; \quad (2)$$

in other words, G^* is at least as good as the best classifier among g_1, \dots, g_J . When $J = 1$, we may simply write g instead of g_1 in which case (2) reduces to

$$\mathbb{P}\{G^*(g(\boldsymbol{\chi})) \neq Y\} \leq \mathbb{P}\{g(\boldsymbol{\chi}) \neq Y\},$$

where the equality holds when g is the Bayes classifier, i.e., $g(\boldsymbol{\chi}) = 1$ if $\mathbb{E}[(2Y - 1) | \boldsymbol{\chi}] > 0$, otherwise $g(\boldsymbol{\chi}) = 0$.

Next, suppose that we have J individual classifiers, $\hat{g}_{n,1}, \dots, \hat{g}_{n,J}$, constructed based on the data \mathbb{T}_n for predicting Y . As explained in the introduction, these could be J very different classifiers; for example, $\hat{g}_{n,1}$ may be a linear classifier, $\hat{g}_{n,2}$ a kernel classifier, $\hat{g}_{n,3}$ may be a random forest classifier or the support vector machine, etc. Then (1) in conjunction with Proposition 1 suggests considering a combined classifier, G_n^* , of the following form

$$G_n^*(\hat{g}_{n,1}(\boldsymbol{\chi}), \dots, \hat{g}_{n,J}(\boldsymbol{\chi})) = \begin{cases} 1 & \text{if } \mathbb{E}[(2Y - 1) | \hat{g}_{n,1}(\boldsymbol{\chi}), \dots, \hat{g}_{n,J}(\boldsymbol{\chi})] > 0 \\ 0 & \text{otherwise.} \end{cases} \quad (3)$$

Although (3) is not available in practice, the following counterpart of Proposition 1 shows that (3) is in fact theoretically optimal in the important sense that its overall error rate is the smallest:

Proposition 2 *Let G_n^* be the combined classifier in (3). Then*

$$\mathbb{P}\left\{G_n^*(\hat{g}_{n,1}(\boldsymbol{x}), \dots, \hat{g}_{n,J}(\boldsymbol{x})) \neq Y\right\} = \inf_{G: \{0,1\}^J \rightarrow \{0,1\}} \mathbb{P}\left\{G(\hat{g}_{n,1}(\boldsymbol{x}), \dots, \hat{g}_{n,J}(\boldsymbol{x})) \neq Y\right\},$$

and in particular $\mathbb{P}\left\{G_n^*(\hat{g}_{n,1}(\boldsymbol{x}), \dots, \hat{g}_{n,J}(\boldsymbol{x})) \neq Y\right\} \leq \min_{1 \leq j \leq J} \mathbb{P}\left\{\hat{g}_{n,j}(\boldsymbol{x}) \neq Y\right\}$.

Obviously, the theoretically optimal classifier G_n^* is not useful in practice because the conditional expectation on the right hand side of (3) is virtually always unknown. Therefore, in what follows, the aim is to construct estimates of (3) whose error rates can be arbitrarily close to that of G_n^* , as the sample size n grows larger and larger. In what follows, we propose a rather simple-to-implement nearest neighbor (NN) type method that works as follows. Randomly split the data \mathbb{T}_n into a training sample \mathbb{T}_m of size m and a testing sequence \mathbb{T}_ℓ of size $\ell = n - m$, where $\mathbb{T}_m \cup \mathbb{T}_\ell = \mathbb{T}_n$ and $\mathbb{T}_m \cap \mathbb{T}_\ell = \emptyset$. Let $\hat{g}_{m,1}, \dots, \hat{g}_{m,J}$ be the J individual classifiers constructed based on \mathbb{T}_m only, and consider the k -NN type combined classifier $G_{n,k}$, $1 \leq k \leq \ell$, given by

$$G_{n,k}(\hat{g}_{m,1}(\boldsymbol{x}), \dots, \hat{g}_{m,J}(\boldsymbol{x})) = \begin{cases} 1 & \text{if } \sum_{i: (\boldsymbol{x}_i, Y_i) \in \mathbb{T}_\ell} (2Y_i - 1) \cdot \mathcal{I}_m(k, \boldsymbol{x}, \boldsymbol{x}_i) > 0 \\ 0 & \text{otherwise,} \end{cases} \quad (4)$$

where $\mathcal{I}_m(k, \boldsymbol{x}, \boldsymbol{x}_i) = \mathbb{I}\{(\hat{g}_{m,1}(\boldsymbol{x}_i), \dots, \hat{g}_{m,J}(\boldsymbol{x}_i)) \text{ is among the } k \text{ nearest neighbors of } (\hat{g}_{m,1}(\boldsymbol{x}), \dots, \hat{g}_{m,J}(\boldsymbol{x}))\}$; (5)

here, the distance between two J -dimensional vectors in $\{0, 1\}^J$ is measured with respect to the Hamming distance, i.e., the number of discrepancies between the corresponding components of the two vectors. In the case of ties, the nearest neighbor to be selected is determined by random chance; thus, for example, if $(\hat{g}_{m,1}(\boldsymbol{x}_i), \dots, \hat{g}_{m,J}(\boldsymbol{x}_i))$ is the third nearest neighbor of $(\hat{g}_{m,1}(\boldsymbol{x}), \dots, \hat{g}_{m,J}(\boldsymbol{x}))$ for $i = 2, 7$, and 10 , then we randomly choose one of these three candidates (and its corresponding Y_i) to be used as the third nearest

neighbor of $(\hat{g}_{m,1}(\boldsymbol{x}), \dots, \hat{g}_{m,J}(\boldsymbol{x}))$. Then, the proposed combined classifier, denoted by $G_{n,\hat{k}}$, is given by (4) with k replaced by \hat{k} that minimizes the re-substitution error, i.e.,

$$\hat{k} = \operatorname{argmin}_{1 \leq k \leq \ell} \frac{1}{\ell} \sum_{i: (\boldsymbol{x}_i, Y_i) \in \mathbb{D}_\ell} \mathbb{I}\{G_{n,k}(\hat{g}_{m,1}(\boldsymbol{x}_i), \dots, \hat{g}_{m,J}(\boldsymbol{x}_i)) \neq Y_i\} \quad (6)$$

To study the asymptotic optimality of $G_{n,\hat{k}}$ we first state two assumptions. Define the quantity

$$\mathcal{S}_{m,\ell}(\boldsymbol{x}) = \sum_{i: \boldsymbol{x}_i \in \mathbb{T}_\ell} \mathbb{I}\{(\hat{g}_{m,1}(\boldsymbol{x}_i), \dots, \hat{g}_{m,J}(\boldsymbol{x}_i)) = (\hat{g}_{m,1}(\boldsymbol{x}), \dots, \hat{g}_{m,J}(\boldsymbol{x}))\}, \quad (7)$$

and consider the following assumptions:

Assumption A.

The quantity $\mathcal{S}_{m,\ell}(\boldsymbol{x})$ diverges with n : $\mathcal{S}_{m,\ell}(\boldsymbol{x}) \rightarrow \infty$, in probability, as n (and thus ℓ) $\rightarrow \infty$.

Assumption B.

For the classifier G_n^* in (3), one has $\mathbb{P}\{G_n^*(\hat{g}_{n,1}(\boldsymbol{x}), \dots, \hat{g}_{n,J}(\boldsymbol{x})) \neq Y\} \rightarrow c$, for some constant c in $[0, 1]$, as $n \rightarrow \infty$.

Assumption A above is not unrealistic at all; to appreciate this, observe that conditional on \mathbb{T}_m and \boldsymbol{x} , the quantity $\mathcal{S}_{m,\ell}(\boldsymbol{x})$ merely represents the total number of “successes” in ℓ independent Bernoulli trials, which is, intuitively, expected to diverge as $\ell \rightarrow \infty$. In fact, alternative versions of this assumption have already been used in the literature (e.g., Devroye et al. [12]; p. 94)). The following result summarizes the asymptotic optimality of the proposed combined classifier $G_{n,\hat{k}}$.

Theorem 1 *Let $G_{n,\hat{k}}$ be the nearest neighbor combined classifier defined in (4), where \hat{k} is the minimizer of the empirical error in (6). If Assumptions A and B hold then*

$$\mathbb{P}\{G_{n,\hat{k}}(\hat{g}_{m,1}(\boldsymbol{x}), \dots, \hat{g}_{m,J}(\boldsymbol{x})) \neq Y\} - \inf_{G: \{0,1\}^J \rightarrow \{0,1\}} \mathbb{P}\{G(\hat{g}_{n,1}(\boldsymbol{x}), \dots, \hat{g}_{n,J}(\boldsymbol{x})) \neq Y\} \rightarrow 0,$$

as $n \rightarrow \infty$. In particular, $G_{n,\hat{k}}$ is asymptotically at least as good as the best individual classifier, i.e.,

$$\limsup_{n \rightarrow \infty} \max_{1 \leq j \leq J} \left[\mathbb{P}\{G_{n,\hat{k}}(\hat{g}_{m,1}(\boldsymbol{x}), \dots, \hat{g}_{m,J}(\boldsymbol{x})) \neq Y\} - \mathbb{P}\{\hat{g}_{n,j}(\boldsymbol{x}) \neq Y\} \right] \leq 0,$$

where $\hat{g}_{n,j}$ is the j -th individual classifier constructed based on the full data \mathbb{T}_n .

Remark 1 Let $G_{n,k}$ and $\mathcal{S}_{m,\ell}(\boldsymbol{\chi})$ be as in (4) and (7), respectively. Now, if we choose k to be equal to $\mathcal{S}_{m,\ell}(\boldsymbol{\chi})$ in $G_{n,k}$, provided that $\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) > 0$, then our proposed combined classifier reduces to that of Balakrishnan and Mojirsheibani [2]. To appreciate this, observe that in this case (4) becomes

$$\begin{aligned}
& G_{n, \mathcal{S}_{m,\ell}(\boldsymbol{\chi})}(\widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi})) \\
&= \begin{cases} 1 & \text{if } \mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) > 0\}} \cdot \sum_{i: (\boldsymbol{\chi}_i, Y_i) \in \mathbb{T}_\ell} (2Y_i - 1) \cdot \mathcal{I}_m(\mathcal{S}_{m,\ell}(\boldsymbol{\chi}), \boldsymbol{\chi}, \boldsymbol{\chi}_i) > 0 \\ 0 & \text{otherwise,} \end{cases} \\
&\quad \left(\text{where } \mathcal{I}_m(\mathcal{S}_{m,\ell}(\boldsymbol{\chi}), \boldsymbol{\chi}, \boldsymbol{\chi}_i) \text{ is as in (5) with } k \text{ replaced by } \mathcal{S}_{m,\ell}(\boldsymbol{\chi}) \right) \\
&= \begin{cases} 1 & \text{if } \mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) > 0\}} \sum_{i: (\boldsymbol{\chi}_i, Y_i) \in \mathbb{T}_\ell} (2Y_i - 1) \cdot \mathbb{I}_{\{(\widehat{g}_{m,1}(\boldsymbol{\chi}_i), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}_i)) = (\widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}))\}} > \frac{0}{0} \\ 0 & \text{otherwise,} \end{cases} \tag{8}
\end{aligned}$$

where we have used the fact that, in view of the definition of $\mathcal{S}_{m,\ell}(\boldsymbol{\chi})$,

$$\mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) > 0\}} \times \mathbb{I}_{\{(\widehat{g}_{m,1}(\boldsymbol{\chi}_i), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}_i)) \text{ is among the } \mathcal{S}_{m,\ell}(\boldsymbol{\chi}) \text{ nearest neighbors of } (\widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}))\}}$$

is equal to 1 if and only if $\mathbb{I}_{\{(\widehat{g}_{m,1}(\boldsymbol{\chi}_i), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}_i)) = (\widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}))\}} = 1$. It is straightforward to see that (8) is equivalent to

$$\begin{cases} 1 & \text{if } \frac{\mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) > 0\}}}{\mathcal{S}_{m,\ell}(\boldsymbol{\chi})} \left[\sum_{i: (\boldsymbol{\chi}_i, Y_i) \in \mathbb{T}_\ell} Y_i \cdot \mathbb{I}_{\{(\widehat{g}_{m,1}(\boldsymbol{\chi}_i), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}_i)) = (\widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}))\}} \right] > \frac{1}{2} \\ 0 & \text{otherwise,} \end{cases} \tag{9}$$

which is the combined classifier of Balakrishnan and Mojirsheibani [2], where, by convention, $0/0 := 0$. The classifier in (9) is essentially a weighted average of all $Y_i \in \mathbb{T}_\ell$, where the weights are indicator functions $\mathbb{I}_{\{(\widehat{g}_{m,1}(\boldsymbol{\chi}_i), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}_i)) = (\widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}))\}}$, each of which will be 1 if and only if $\widehat{g}_{m,j}(\boldsymbol{\chi}_i) = \widehat{g}_{m,j}(\boldsymbol{\chi})$ for all $j \in \{1, 2, \dots, J\}$. Unfortunately, from a practical point of view, if there are a few weak/poor classifiers among $\widehat{g}_{m,1}, \dots, \widehat{g}_{m,J}$, then one could end up with $\widehat{g}_{m,j}(\boldsymbol{\chi}_i) \neq \widehat{g}_{m,j}(\boldsymbol{\chi})$ for a large number of $\boldsymbol{\chi}_i$'s in \mathbb{T}_ℓ and this can hold true even when $\boldsymbol{\chi}_i$ and $\boldsymbol{\chi}$ belong to the same class. This means that many of the weights (i.e., the indicator functions) in (9) will be zero, which in turn reduces the predictive performance of (9). Our proposed method in this paper circumvents this vulnerability of the combined classifier of Balakrishnan and Mojirsheibani [2] in (9) by allowing a few weak classifiers to “misbehave” or be incorrect in their predictions. Thus, unlike (9), the new classifier $G_{n,\widehat{k}}$ is not seriously affected by the impact of a few poor classifiers.

Remark 2 The results of this section can be extended to the M -group ($M \geq 2$) classification problem in a straightforward manner. More specifically, let $(\boldsymbol{\chi}, Y)$ be a random pair where $\boldsymbol{\chi}$ is as before, but $Y \in \{1, \dots, M\}$. A classifier of the form $g^*(\boldsymbol{\chi}) := \operatorname{argmax}_{1 \leq k \leq M} \mathbb{P}\{Y = k | \boldsymbol{\chi} = \boldsymbol{\chi}\}$ is optimal in the sense that $\mathbb{P}\{g^*(\boldsymbol{\chi}) \neq$

$Y\} = \min_{g: \mathcal{F} \rightarrow \{1, \dots, M\}} \mathbb{P}\{g(\boldsymbol{\chi}) \neq Y\}$; see, for example, Devroye and Györfi [11], ch. 10. In this case, with $\mathcal{I}_m(k, \boldsymbol{\chi}, \boldsymbol{\chi}_i)$ as in (5), we have the following counterpart of (4):

$$G_{n,k}(\hat{g}_{m,1}(\boldsymbol{\chi}), \dots, \hat{g}_{m,J}(\boldsymbol{\chi})) = \operatorname{argmax}_{1 \leq j \leq M} \sum_{i: (\boldsymbol{\chi}_i, Y_i) \in \mathbb{T}_\ell} \mathbb{I}_{\{Y_i=j\}} \cdot \mathcal{I}_m(k, \boldsymbol{\chi}, \boldsymbol{\chi}_i)$$

and the proposed combined classifier is given by $G_{n,\hat{k}}$, where \hat{k} is as in (6). It can be shown that, under assumption A and the version of B corresponding to the M -group problem, the conclusion of Theorem 1 continues to hold in the general M -group problem in the sense that

$$\begin{aligned} & \mathbb{P}\left\{G_{n,\hat{k}}(\hat{g}_{m,1}(\boldsymbol{\chi}), \dots, \hat{g}_{m,J}(\boldsymbol{\chi})) \neq Y\right\} \\ & - \inf_{G: \{1, \dots, M\}^J \rightarrow \{1, \dots, M\}} \mathbb{P}\left\{G(\hat{g}_{m,1}(\boldsymbol{\chi}), \dots, \hat{g}_{m,J}(\boldsymbol{\chi})) \neq Y\right\} \longrightarrow 0, \text{ as } n \rightarrow \infty. \end{aligned}$$

Remark 3 In the case where $J = 1$, Theorem 1 essentially implies that, given any initial classifier \hat{g}_n , the error of the new classifier $G_{n,\hat{k}}$ defined via (4) and (6) can always be asymptotically less than or equal to that of \hat{g}_n . To the best of our knowledge, this is a new result in the literature.

Remark 4 Since, for each $k = 1, \dots, \ell$, the nearest neighbor type combined classifier $G_{n,k}$ in (4) is constructed based on one sample split, its performance can be affected by the particular split used. Thus, unless n is very large, a ‘‘bad’’ split can in practice result in a poor choice of \hat{k} in (6) which will lead to a poor corresponding classifier $G_{n,\hat{k}}$. This practical issue suggests using several random splits and taking their average. More precisely, with N sample splits, each split will produce an estimate \hat{k}_b of k , where $b = 1, \dots, N$, and the corresponding predicted class Y (corresponding to $\boldsymbol{\chi}$) is given by $\hat{Y}_b = G_{n,\hat{k}_b}(\hat{g}_{m,1}(\boldsymbol{\chi}), \dots, \hat{g}_{m,J}(\boldsymbol{\chi}))$. Since each \hat{Y}_b is either 0 or 1, the overall predicted value of Y is taken to be 1 if $N^{-1} \sum_{b=1}^N \hat{Y}_b > 1/2$, and 0 otherwise. As for the choice of N , our experience shows that one can expect good results with N as small as 15 or 20.

3 Numerical examples

In what follows, we study the prediction of the class variable, Y ($= 0$ or 1 , corresponding to the random covariate $\boldsymbol{\chi}$, using the ensemble methods proposed in this paper.

Example A (Simulated data).

In this example we consider the prediction of $Y = 0$ or 1 , based on $\boldsymbol{\chi} \in \mathbb{R}^{10}$, where $\mathbb{P}\{Y = 1\} = \mathbb{P}\{Y = 0\} = 0.5$. Here, we have taken $\boldsymbol{\chi} \sim N_{10}(\mathbf{0}, 8\boldsymbol{\Sigma})$ whenever

$Y = 1$ (i.e., class 1), where $\Sigma = (\sigma_{ij})_{i,j=1,\dots,10}$, with $\sigma_{ij} = 2^{|i-j|}$. On the other hand, if $Y = 0$ then χ has a 10-dim standard Cauchy distribution with independent components, i.e., the components of χ are iid random variables with the pdf $f(x) = (\pi(1+x^2))^{-1}$, $-\infty < x < \infty$. As for the choice of the individual classifiers, we have considered the following six classifiers: (i) a 1-Nearest Neighbor (1-NN) classifier, (ii) a 7-NN classifier, (iii) the Support Vector Machine (SVM) of Boser et al. [5], (iv) Breiman's [8] Random Forest, (v) a Gaussian kernel classifier with a bandwidth of $n^{-0.2}$, and (vi) the Linear Discriminant Analysis (LDA). Next we considered three different combined classifiers: The proposed combined classifier $G_{n,\hat{k}}$ defined in (4), where \hat{k} is the minimizer of the empirical error in (6), the combined classifier of Balakrishnan and Mojirsheibani [2] given in (9), and the majority-vote classifier. To trace the performance of various methods, several sample sizes were used: $n = 50, 100, 200, \dots, 800, 900$, with $\frac{n}{2}$ observations from each population. As for the sample splits, we took $m = 0.65n$ and $\ell = n - m$, but any other fraction in the 0.55 to 0.85 range seems to work just as well. Next, for each sample size, we constructed the six individual classifiers based on \mathbb{T}_m , which were then used to construct the above combined classifiers based on \mathbb{T}_ℓ . For each value of n , a total of 25 sample splits were used; this is in view of Remark 4. To assess the performance of various classifiers, we also generated 1000 additional observations, with 500 from each of the two populations; these were used as test samples for each classifier. Finally, the whole process above was repeated a total of 300 times, yielding 300 estimates of the misclassification errors of each classifiers discussed above. The average errors (over 300 Monte Carlo runs) are summarized in Table 1 along with their standard errors in parentheses.

As Table 1 shows, the SVM and random forest classifiers perform very well for smaller sample sizes, but as n reaches 100, the proposed combined classifier $G_{n,\hat{k}}$ (as well as B-M) start performing better than all other classifiers. The boldface values represent the smallest error rates for each n . We also note that as n gets larger and larger, most classifiers start performing better and better (except for the LDA which is a wrong classifier in the presence of Cauchy populations); however, the proposed combined classifier is consistently superior

Table 1: Misclassification errors of the three combined classifiers and the six individual classifiers for the simulated data of Example A. Here $\widehat{G}_{n,\widehat{k}}$ is the proposed classifier, B-M is the combined classifier of Balakrishnan and Mojirsheibani [2], and Vote is the combined classifier based on majority voting. The boldfaced values represent the lowest errors for each n .

n	$\widehat{G}_{n,\widehat{k}}$	B-M	Vote	SVM	Forest	1-NN	7-NN	Kernel	LDA
50	.3342 (.0029)	.3404 (.0023)	.3732 (.0015)	.3020 (.0016)	.3180 (.0018)	.3616 (.0015)	.4698 (.0013)	.3339 (.0014)	.4864 (.0013)
100	.2802 (.0017)	.2828 (.0015)	.3366 (.0016)	.2847 (.0010)	.2875 (.0017)	.3329 (.0016)	.4374 (.0017)	.3126 (.0013)	.4906 (.0012)
200	.1709 (.0014)	.1821 (.0013)	.2646 (.0010)	.2506 (.0009)	.2181 (.0011)	.2696 (.0011)	.3341 (.0011)	.2635 (.0010)	.4849 (.0013)
300	.1491 (.0013)	.1623 (.0011)	.2463 (.0010)	.2428 (.0009)	.2078 (.0011)	.2534 (.0009)	.2897 (.0012)	.2510 (.0009)	.4873 (.0013)
400	.1452 (.0010)	.1514 (.0009)	.2248 (.0013)	.2385 (.0008)	.1909 (.0009)	.2399 (.0009)	.2570 (.0011)	.2387 (.0009)	.4843 (.0015)
500	.1394 (.0009)	.1405 (.0008)	.2105 (.0009)	.2262 (.0012)	.1739 (.0011)	.2275 (.0010)	.2354 (.0010)	.2265 (.0010)	.4945 (.0013)
600	.1360 (.0008)	.1360 (.0007)	.2001 (.0009)	.2212 (.0007)	.1676 (.0008)	.2221 (.0007)	.2195 (.0010)	.2213 (.0007)	.4845 (.0014)
700	.1310 (.0010)	.1343 (.0008)	.1961 (.0007)	.2199 (.0010)	.1656 (.0009)	.2230 (.0009)	.2097 (.0009)	.2221 (.0009)	.4891 (.0012)
800	.1249 (.0008)	.1278 (.0007)	.1861 (.0009)	.2119 (.0008)	.1583 (.0008)	.2137 (.0009)	.2066 (.0009)	.2133 (.0009)	.4851 (.0013)
900	.1190 (.0009)	.1243 (.0008)	.1856 (.0006)	.2076 (.0010)	.1505 (.0010)	.2122 (.0008)	.1969 (.0008)	.2118 (.0008)	.4920 (.0013)

to all the other ones. Furthermore, in some cases this superiority is quite notable; see, for example, the row corresponding to $n=300$, where the error of $\widehat{G}_{n,\widehat{k}}$ is only 0.1491 as compared to the best individual classifier, the random forest, with an error of 0.2078. Such rather large discrepancies can also be noticed for many other values of n in Table 1.

Another feature of the proposed classifier $\widehat{G}_{n,\widehat{k}}$ is that, unlike linear combined classifiers, it can even be used to improve the predictive performance of any single classifier (although the improvement may be quite incremental in some cases). To appreciate this, observe that according to Theorem 1, the proposed combined classifier $\widehat{G}_{n,\widehat{k}}$ can asymptotically outperform each constituent classifier $\widehat{g}_{n,1}, \dots, \widehat{g}_{n,J}$, where $J \geq 1$. Now, taking $J=1$, this theorem states that, given a single classifier \widehat{g}_n , the proposed approach can produce an improved version of \widehat{g}_n . To put this to the test, we applied $\widehat{G}_{n,\widehat{k}}$ to some of the six

classifiers in Table 1; in fact, we applied it to each of the two classifiers that have already performed quite well and are difficult to outperform, i.e., the SVM and random forest classifiers. The results appear in Table 2. This table shows that as n increases, the

Table 2: Effects of applying the combined classifier $\widehat{G}_{n,\widehat{k}}$ to $J=1$ classifier only. Here \widehat{SVM} and \widehat{Forest} represent the classifiers obtained by applying $\widehat{G}_{n,\widehat{k}}$ to SVM and random forest, respectively. For each n , the boldfaced values represent the smaller of the two error rates when comparing SVM and \widehat{SVM} in columns 3 and 4, and when comparing random forest and \widehat{Forest} in columns 4 and 5.

n	SVM	\widehat{SVM}	Forest	\widehat{Forest}
50	0.3020 (.0016)	0.3076 (.0017)	0.3180 (.0018)	0.3193 (.0024)
100	0.2847 (.0010)	0.2871 (.0011)	0.2875 (.0017)	0.2844 (.0016)
200	0.2506 (.0009)	0.2411 (.0008)	0.2181 (.0011)	0.2032 (.0013)
300	0.2428 (.0009)	0.2380 (.0009)	0.2078 (.0011)	0.1993 (.0008)
400	0.2385 (.0009)	0.2302 (.0007)	0.1909 (.0009)	0.1815 (.0009)
500	0.2262 (.0012)	0.2265 (.0010)	0.1739 (.0011)	0.1732 (.0010)
600	0.2212 (.0007)	0.2207 (.0008)	0.1676 (.0008)	0.1768 (.0008)
700	0.2199 (.0010)	0.2188 (.0008)	0.1656 (.0009)	0.1643 (.0009)
800	0.2119 (.0008)	0.2031 (.0007)	0.1583 (.0008)	0.1492 (.0008)
900	0.2076 (.0010)	0.2070 (.0007)	0.1505 (.0010)	0.1480 (.0009)

proposed classifier can still improve upon the performance of each of these two classifiers, individually, despite the fact that both SVM and random forest are well known to be superb classifiers. It is true that the improvement is rather minimal, but the main message here is that $\widehat{G}_{n,\widehat{k}}$ is more than just a combined classifier, it can also improve the predictive power of a single classifier.

Example B (*Wisconsin Breast Cancer Data*).

This real data set has 683 fully observed instances, 444 of which have been labeled *benign*,

which is class 1, and the rest are *malignant*, i.e., class 0. There are also 9 numerical covariates associated with each instance. A full description of this data set is available from the UCI Machine Learning Repository of data sets: <https://archive.ics.uci.edu/ml/datasets.php>. Also, see Wolberg and Mangasarian (1990).

To carry out the analysis, 500 of the 683 instances were randomly selected to be used as the training data, whereas the remaining 183 were set aside as the test sequence to be used to estimate the error rates of different classifiers. To study the performance of various classifiers as a function of the sample size n , we considered 5 different sample sizes $n = 100, \dots, 500$ (since n can only go up to 500 here) and, for each value of n , six individual classifiers were constructed which were then used to construct the proposed combined classifier. Here, as in Example A, the sample splits were taken to be $m = 0.65n$ and $\ell = n - m$. Finally, the error rates of all classifiers were estimated using the test sequence of 183 instances that were set aside. This whole process was repeated 100 times. Table 3 reports the average error rates of various classifiers over 100 runs; the standard errors appear in parentheses. The boldfaced values represent the smallest error rates for each n . The combined classifier B-M of Table 1 is not included here for the simple reason that it is always inferior to $\widehat{G}_{n,\widehat{k}}$. As Table 3 shows, the proposed combined classifier can outperform the individual classifiers. This can be noticed by comparing the error of the best individual classifier (the random forest in this case) with that of $\widehat{G}_{n,\widehat{k}}$ that appears in the first column.

Example C (*German Credit Data*).

Here we consider a real data set consisting of 1000 individuals, 700 of whom have been labeled as having “good credit”, i.e., class 1, whereas the remaining 300 have “bad credit”, which is class 0. There are 24 numerical covariates associated with each person. A full description of this data set is available from the UCI repository of machine learning data sets at <https://archive.ics.uci.edu/ml/index.php>.

Table 3: Misclassification errors of various classifiers for the *Wisconsin Breast Cancer* data of Example B. Here $\widehat{G}_{n,\widehat{k}}$ is our proposed classifier and Vote is the combined classifier based on majority voting. The boldfaced values represent the lowest errors for each n .

n	$\widehat{G}_{n,\widehat{k}}$	Vote	SVM	Forest	1-NN	15-NN	Kernel	LDA
100	.0362 (.0013)	.0387 (.0012)	.0541 (.0017)	.0387 (.0014)	.0469 (.0016)	.0449 (.0012)	.0459 (.0016)	.0453 (.0016)
200	.0312 (.0012)	.0347 (.0013)	.0511 (.0018)	.0337 (.0014)	.0427 (.0015)	.0384 (.0014)	.0422 (.0016)	.0392 (.0013)
300	.0279 (.0011)	.0304 (.0011)	.0439 (.0015)	.0303 (.0012)	.0396 (.0012)	.0346 (.00012)	.0391 (.0011)	.0375 (.0012)
400	.0277 (.0010)	.0316 (.0012)	.0445 (.0015)	.0298 (.0010)	.0408 (.0014)	.0338 (.0012)	.0409 (.0013)	.0385 (.0012)
500	.0272 (.0010)	.0327 (.0011)	.0422 (.0014)	.0295 (.0010)	.0416 (.0012)	.0336 (.0012)	.0416 (.0012)	.0395 (.0012)

To carry out the analysis, first we randomly selected, and set aside, 300 of the 1000 observations to be used as the test sequence to estimate the error rates of various classifiers. As for the training sample size, seven values were considered: $n = 100, 200, \dots, 700$, (since n cannot go beyond $700=1000-300$). This grid of values of n allows us to somewhat monitor the performance of different classifiers as n increases. Then, given a sample of size n , each of the six individual classifiers of Example B were constructed and used to construct the proposed combined classifiers. Here, once again, the sample splits were taken to be $m = 0.65n$ and $\ell = n - m$. Finally, the error rates of various classifiers were estimated using the test sample of 300 observations. The entire process above was repeated 100 times and the average misclassification error rates were calculated. The results are summarized in Table 4. As this table shows, the combined classifier $\widehat{G}_{n,\widehat{k}}$ has the ability to perform well and, in fact, slightly outperform the best individual classifier (which is random forest) as the sample size increases to about 600.

4 Proofs

PROOF OF THEOREM 1

Let G_n^* and $G_{n,k}$ be as in (3) and (4), respectively. Similarly, let G_m^* be as in (3), but

Table 4: Misclassification errors of different classifiers for the *German Credit* data of Example C. Here $\widehat{G}_{n,\widehat{k}}$ is our proposed classifier and Vote is the combined classifier based on majority voting. The boldfaced values represent the lowest errors for each n .

n	$\widehat{G}_{n,\widehat{k}}$	Vote	SVM	Forest	1-NN	15-NN	Kernel	LDA
100	.2914 (.0028)	.3012 (.0024)	.3053 (.0022)	.2755 (.0024)	.3746 (.0032)	.3061 (.0023)	.3718 (.0032)	.3150 (.0033)
200	.2691 (.0026)	.2922 (.0026)	.2975 (.0022)	.2613 (.0023)	.3590 (.0025)	.3030 (.0023)	.3572 (.0025)	.2974 (.0030)
300	.2608 (.0024)	.2838 (.0022)	.2954 (.0024)	.2523 (.0020)	.3534 (.0026)	.2984 (.00023)	.3521 (.0027)	.2893 (.0022)
400	.2472 (.0024)	.2794 (.0029)	.2909 (.0023)	.2458 (.0027)	.3494 (.0027)	.2973 (.0024)	.3467 (.0026)	.2860 (.0023)
500	.2438 (.0021)	.2738 (.0024)	.2886 (.0024)	.2414 (.0021)	.3417 (.0023)	.2955 (.0023)	.3399 (.0024)	.2775 (.0022)
600	.2398 (.0023)	.2749 (.0022)	.2914 (.0023)	.2411 (.0022)	.3424 (.0023)	.2994 (.0021)	.3395 (.0024)	.2774 (.0024)
700	.2372 (.0021)	.2756 (.0022)	.2869 (.0026)	.2376 (.0021)	.3407 (.0023)	.2990 (.0022)	.3385 (.0022)	.2800 (.0024)

with n replaced with m everywhere in (3). Also, let $G_{n,\widehat{k}}$ be the combined classifier given by (4) and (6), and define the quantities

$$L(G_n^*) = \mathbb{P}\left\{G_n^*(\widehat{g}_{n,1}(\boldsymbol{x}), \dots, \widehat{g}_{n,J}(\boldsymbol{x})) \neq Y\right\} \quad (10)$$

$$L(G_m^*) = \mathbb{P}\left\{G_m^*(\widehat{g}_{m,1}(\boldsymbol{x}), \dots, \widehat{g}_{m,J}(\boldsymbol{x})) \neq Y\right\} \quad (11)$$

$$\widehat{L}_\ell(G_{n,k}) = \frac{1}{\ell} \sum_{i: (\boldsymbol{x}_i, Y_i) \in \mathbb{T}_\ell} \mathbb{I}\{G_{n,k}(\widehat{g}_{m,1}(\boldsymbol{x}_i), \dots, \widehat{g}_{m,J}(\boldsymbol{x}_i)) \neq Y_i\} \quad (12)$$

$$L(G_{n,k}|\mathbb{T}_n) = \mathbb{P}\left\{G_{n,k}(\widehat{g}_{m,1}(\boldsymbol{x}), \dots, \widehat{g}_{m,J}(\boldsymbol{x})) \neq Y \mid \mathbb{T}_n\right\} \quad (13)$$

$$L(G_{n,\widehat{k}}) = \mathbb{P}\left\{G_{n,\widehat{k}}(\widehat{g}_{m,1}(\boldsymbol{x}), \dots, \widehat{g}_{m,J}(\boldsymbol{x})) \neq Y\right\} \quad (14)$$

$$L(G_{n,\widehat{k}}|\mathbb{T}_n) = \mathbb{P}\left\{G_{n,\widehat{k}}(\widehat{g}_{m,1}(\boldsymbol{x}), \dots, \widehat{g}_{m,J}(\boldsymbol{x})) \neq Y \mid \mathbb{T}_n\right\}, \quad (15)$$

where \widehat{k} is as in (6) and $1 \leq k \leq \ell$. Therefore, in view of Proposition 2, and the notation in (10) and (14), one must show $L(G_{n,\widehat{k}}) - L(G_n^*) \rightarrow 0$, as $n \rightarrow \infty$. Now, start by writing

$$\begin{aligned} L(G_{n,\widehat{k}}) - L(G_n^*) &= \mathbb{E}\left[L(G_{n,\widehat{k}}|\mathbb{T}_n) - \min_{1 \leq k \leq \ell} L(G_{n,k}|\mathbb{T}_n)\right] \\ &\quad + \left[\mathbb{E} \min_{1 \leq k \leq \ell} L(G_{n,k}|\mathbb{T}_n) - L(G_m^*)\right] + [L(G_n^*) - L(G_m^*)] \end{aligned}$$

$$:= R_{n,1} + R_{n,2} + R_{n,3}. \quad (16)$$

Then, fix \mathbb{T}_n and define the classes of sets $\mathcal{A}_{n,1} = \{A_{n,1,1}, \dots, A_{n,1,\ell}\}$ and $\mathcal{A}_{n,0} = \{A_{n,0,1}, \dots, A_{n,0,\ell}\}$ where, for $1 \leq k \leq \ell$,

$$A_{n,1,k} = \left\{ (\hat{g}_{m,1}(\chi), \dots, \hat{g}_{m,J}(\chi)), \text{ as } \chi \text{ varies over } \mathcal{F} \mid G_{n,k}(\hat{g}_{m,1}(\chi), \dots, \hat{g}_{m,J}(\chi)) = 1 \right\} \times \{0\}$$

$$A_{n,0,k} = \left\{ (\hat{g}_{m,1}(\chi), \dots, \hat{g}_{m,J}(\chi)), \text{ as } \chi \text{ varies over } \mathcal{F} \mid G_{n,k}(\hat{g}_{m,1}(\chi), \dots, \hat{g}_{m,J}(\chi)) = 0 \right\} \times \{1\}$$

Furthermore, for $1 \leq k \leq \ell$, let

$$\lambda(A_{n,b,k} | \mathbb{T}_n) = \mathbb{P} \left\{ (\hat{g}_{m,1}(\mathbf{X}), \dots, \hat{g}_{m,J}(\mathbf{X}), Y) \in A_{n,b,k} \mid \mathbb{T}_n \right\}, \quad b = 0, 1$$

$$\hat{\lambda}_\ell(A_{n,b,k}) = \ell^{-1} \sum_{i: (\mathbf{X}_i, Y_i) \in \mathbb{T}_\ell} \mathbb{I} \left\{ (\hat{g}_{m,1}(\mathbf{X}_i), \dots, \hat{g}_{m,J}(\mathbf{X}_i), Y_i) \in A_{n,b,k} \right\}, \quad b = 0, 1,$$

and observe that

$$\begin{aligned} & L(G_{n,\hat{k}} | \mathbb{T}_n) - \min_{1 \leq k \leq \ell} L(G_{n,k} | \mathbb{T}_n) \\ &= L(G_{n,\hat{k}} | \mathbb{T}_n) - \hat{L}_\ell(G_{n,\hat{k}}) + \hat{L}_\ell(G_{n,\hat{k}}) - \min_{1 \leq k \leq \ell} L(G_{n,k} | \mathbb{T}_n) \\ &\leq 2 \max_{1 \leq k \leq \ell} \left| \hat{L}_\ell(G_{n,k}) - L(G_{n,k} | \mathbb{T}_n) \right| \\ &\leq 2 \sum_{b=0,1} \max_{1 \leq k \leq \ell} \left| \hat{\lambda}_\ell(A_{n,b,k}) - \lambda(A_{n,b,k} | \mathbb{T}_n) \right| \\ &\leq 4 \sup_{B \in \mathcal{B}} \left| \hat{\lambda}_\ell(B) - \lambda(B | \mathbb{T}_m) \right| \\ &\quad (\text{where } \mathcal{B} \text{ is the collection of the Borel sets of } \mathbb{R}^{J+1}) \\ &\leq 4 \sum_{\mathbf{z} \in \{0,1\}^{J+1}} \left| \hat{\lambda}_\ell(\{\mathbf{z}\}) - \lambda(\{\mathbf{z}\} | \mathbb{T}_m) \right|, \end{aligned} \quad (17)$$

where, for $\mathbf{z} \in \mathbb{R}^{J+1}$, we have $\hat{\lambda}_\ell(\{\mathbf{z}\}) = \ell^{-1} \sum_{i: (\mathbf{X}_i, Y_i) \in \mathbb{T}_\ell} \mathbb{I} \left\{ (\hat{g}_{m,1}(\mathbf{X}_i), \dots, \hat{g}_{m,J}(\mathbf{X}_i), Y_i) = \mathbf{z} \right\}$ and $\lambda(\{\mathbf{z}\} | \mathbb{T}_m) = \mathbb{P} \left\{ (\hat{g}_{m,1}(\mathbf{X}), \dots, \hat{g}_{m,J}(\mathbf{X}), Y) = \mathbf{z} \mid \mathbb{T}_m \right\}$. Therefore, in view of (17), for every $\epsilon > 0$, one has

$$\begin{aligned} & \mathbb{P} \left\{ L(G_{n,\hat{k}} | \mathbb{T}_n) - \min_{1 \leq k \leq \ell} L(G_{n,k} | \mathbb{T}_n) > \epsilon \right\} \\ & \leq \mathbb{P} \left\{ \sum_{\mathbf{z} \in \{0,1\}^{J+1}} \left| \hat{\lambda}_\ell(\{\mathbf{z}\}) - \lambda(\{\mathbf{z}\} | \mathbb{T}_m) \right| > \epsilon/4 \right\} \end{aligned}$$

$$\begin{aligned}
&\leq \sum_{\mathbf{z} \in \{0,1\}^{J+1}} \mathbb{E} \left[\mathbb{P} \left\{ \left| \widehat{\lambda}_\ell(\{\mathbf{z}\}) - \lambda(\{\mathbf{z}\} | \mathbb{T}_m) \right| > 2^{-(J+3)} \epsilon \mid \mathbb{T}_m \right\} \right] \\
&\leq 2^{J+1} \left[2 e^{-C \ell \epsilon^2} \right], \quad \text{with } C = 2^{-(2J+5)}, \tag{18}
\end{aligned}$$

where (18) follows from an application of Hoeffding's (1963) inequality in view of the fact that, conditional on \mathbb{T}_m , the term $\widehat{\lambda}_\ell(\{\mathbf{z}\})$ is the average of ℓ independent indicator functions, $\mathbb{I}_{\{(\widehat{g}_{m,1}(\mathbf{x}_i), \dots, \widehat{g}_{m,J}(\mathbf{x}_i), Y_i) = \mathbf{z}\}}$, corresponding to the ℓ pairs $(\mathbf{x}_i, Y_i) \in \mathbb{T}_\ell$. Now, (18) together with the Borel-Cantelli lemma yields $L(G_{n,\widehat{k}} | \mathbb{T}_n) - \min_{1 \leq k \leq \ell} L(G_{n,k} | \mathbb{T}_n) \xrightarrow{a.s.} 0$, as $n \rightarrow \infty$. Thus, by Lebesgue's dominated convergence theorem,

$$R_{n,1} := \mathbb{E} \left[L(G_{n,\widehat{k}} | \mathbb{T}_n) - \min_{1 \leq k \leq \ell} L(G_{n,k} | \mathbb{T}_n) \right] \rightarrow 0, \quad \text{as } n \rightarrow \infty.$$

To deal with the term $R_{n,2}$ in (16), start by defining the quantity

$$\nu_n(\boldsymbol{\chi}) = \mathcal{S}_{m,\ell}(\boldsymbol{\chi}) + \ell \cdot \mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi})=0\}}, \tag{19}$$

where $\mathcal{S}_{m,\ell}(\boldsymbol{\chi})$ is as in (7). Here, we note that $\nu_n(\boldsymbol{\chi})$ takes values in $\{1, \dots, \ell\}$. Also put

$$\widehat{P}_{m,\ell}(\boldsymbol{\chi}) = \frac{\sum_{i: \mathbf{x}_i \in \mathbb{T}_\ell} Y_i \mathbb{I}_{\{(\widehat{g}_{m,1}(\mathbf{x}_i), \dots, \widehat{g}_{m,J}(\mathbf{x}_i)) = (\widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}))\}}}{\mathcal{S}_{m,\ell}(\boldsymbol{\chi})} \tag{20}$$

and observe that upon replacing k by $\nu_n(\boldsymbol{\chi})$ in (4), we can write

$$\begin{aligned}
&G_{n,\nu_n(\boldsymbol{\chi})}(\widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi})) \\
&= \begin{cases} 1 & \text{if } \frac{1}{\nu_n(\boldsymbol{\chi})} \sum_{i: (\mathbf{x}_i, Y_i) \in \mathbb{T}_\ell} (2Y_i - 1) \cdot \mathcal{I}_m(\nu_n(\boldsymbol{\chi}), \boldsymbol{\chi}, \mathbf{x}_i) > 0 \\ 0 & \text{otherwise,} \end{cases} \\
&\quad \text{(where the term } \mathcal{I}_m(\nu_n(\boldsymbol{\chi}), \boldsymbol{\chi}, \mathbf{x}_i) \text{ is defined via (5))} \\
&= \begin{cases} 1 & \text{if } \frac{\mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) > 0\}}}{\mathcal{S}_{m,\ell}(\boldsymbol{\chi})} \sum_{i: (\mathbf{x}_i, Y_i) \in \mathbb{T}_\ell} (2Y_i - 1) \cdot \mathcal{I}_m(\mathcal{S}_{m,\ell}(\boldsymbol{\chi}), \boldsymbol{\chi}, \mathbf{x}_i) \\ & \quad + \frac{\mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) = 0\}}}{\ell} \sum_{i: (\mathbf{x}_i, Y_i) \in \mathbb{T}_\ell} (2Y_i - 1) \cdot \mathcal{I}_m(\ell, \boldsymbol{\chi}, \mathbf{x}_i) > 0 \\ 0 & \text{otherwise,} \end{cases} \tag{21}
\end{aligned}$$

(where (21) follows from the definition of $\nu_n(\boldsymbol{\chi})$ in (19))

$$= \begin{cases} 1 & \text{if } \frac{\mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) > 0\}}}{\mathcal{S}_{m,\ell}(\boldsymbol{\chi})} \sum_{i: (\mathbf{x}_i, Y_i) \in \mathbb{T}_\ell} (2Y_i - 1) \mathbb{I}_{\{(\widehat{g}_{m,1}(\mathbf{x}_i), \dots, \widehat{g}_{m,J}(\mathbf{x}_i)) = (\widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}))\}} \\ & \quad + (2\bar{Y}_\ell - 1) \mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) = 0\}} > 0 \\ 0 & \text{otherwise,} \end{cases} \tag{22}$$

$$= \begin{cases} 1 & \text{if } \widehat{\mathcal{P}}_{m,\ell}(\boldsymbol{\chi}) \cdot \mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) > 0\}} + \bar{Y}_\ell \cdot \mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) = 0\}} > 1/2 \\ 0 & \text{otherwise,} \end{cases} \quad (23)$$

where $\widehat{\mathcal{P}}_{m,\ell}(\boldsymbol{\chi})$ is as in (20) and $\bar{Y}_\ell = \ell^{-1} \sum_{i: Y_i \in \mathbb{T}_\ell} Y_i$. Here, (22) follows from (21) because of the following simple facts:

(i) For each $\boldsymbol{\chi}_i \in \mathbb{T}_\ell$, the product of the two indicator functions, $\mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) > 0\}}$ and

$$\mathbb{I}_{\{(\widehat{g}_{m,1}(\boldsymbol{\chi}_i), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}_i)) \text{ is among the } \mathcal{S}_{m,\ell}(\boldsymbol{\chi}) \text{ neighbors of } (\widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}))\}},$$

(where the latter indicator function is just the term $\mathcal{I}_m(\mathcal{S}_{m,\ell}(\boldsymbol{\chi}), \boldsymbol{\chi}, \boldsymbol{\chi}_i)$ in (21)), will be equal to 1 if and only if $\mathbb{I}_{\{(\widehat{g}_{m,1}(\boldsymbol{\chi}_i), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}_i)) = (\widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}))\}} = 1$.

(ii) For each $\boldsymbol{\chi}_i \in \mathbb{T}_\ell$, the term $\mathcal{I}_m(\ell, \boldsymbol{\chi}, \boldsymbol{\chi}_i)$ in (21), which is just the indicator function

$$\mathbb{I}_{\{(\widehat{g}_{m,1}(\boldsymbol{\chi}_i), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}_i)) \text{ is among the } \ell \text{ nearest neighbors of } (\widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}))\}},$$

is always equal to 1; this is because here the number of nearest neighbors, ℓ , is the same as the entire sample size, ℓ , (the size of \mathbb{T}_ℓ).

To complete the proof, we also need the following lemma which puts bounds on the term $R_{n,2}$ in (16) based on the expected value of the expression that appears on the right side of (23). More specifically,

Lemma 1 *Let $R_{n,2}$ be as in (16) and put $\bar{Y}_\ell = \ell^{-1} \sum_{i: Y_i \in \mathbb{T}_\ell} Y_i$. Then*

$$0 \leq R_{n,2} \leq 2 \mathbb{E} \left| \widehat{\mathcal{P}}_{m,\ell}(\boldsymbol{\chi}) \cdot \mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) > 0\}} + \bar{Y}_\ell \cdot \mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) = 0\}} - \mathcal{P}_m(\boldsymbol{\chi}) \right|,$$

where $\mathcal{P}_m(\boldsymbol{\chi}) = \mathbb{P} \{Y = 1 \mid \widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi})\}$.

Now, to show $R_{n,2} \rightarrow 0$, let $\mathcal{P}_m(\boldsymbol{\chi})$ be as in Lemma 1 and observe that by Lemma 1 and the fact that $0 \leq \bar{Y}_\ell \leq 1$, one has

$$\begin{aligned} R_{n,2} &\leq 2 \mathbb{E} \left| \widehat{\mathcal{P}}_{m,\ell}(\boldsymbol{\chi}) \cdot \mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) > 0\}} - \mathcal{P}_m(\boldsymbol{\chi}) \right| + 2 \mathbb{E} \left(\mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) = 0\}} \right) \\ &\leq 2 \sqrt{\mathbb{E} \left[\mathbb{E} \left(\left| \widehat{\mathcal{P}}_{m,\ell}(\boldsymbol{\chi}) \cdot \mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) > 0\}} - \mathcal{P}_m(\boldsymbol{\chi}) \right|^2 \middle| \mathbb{T}_m, \boldsymbol{\chi}, \{\boldsymbol{\chi}_i\}_i \in \mathbb{T}_\ell \right) \right]} \quad (24) \end{aligned}$$

$$+ 2\mathbb{P}\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) = 0\}.$$

However, conditional on \mathbb{T}_m , $\boldsymbol{\chi}$, and $\{\boldsymbol{\chi}_i\}_{i \in \mathbb{T}_\ell}$, the random variable $\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) \cdot \widehat{\mathcal{P}}_{m,\ell}(\boldsymbol{\chi})$, which is equal to the numerator of the right side of (20), has the binomial distribution $\text{Bin}(\mathcal{S}_{m,\ell}(\boldsymbol{\chi}), \mathcal{P}_m(\boldsymbol{\chi}))$ whenever $\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) > 0$. Therefore, the expression under the square-root sign in (24) can be bounded as follows

$$\begin{aligned} & \mathbb{E} \left[\mathbb{E} \left(\left| \widehat{\mathcal{P}}_{m,\ell}(\boldsymbol{\chi}) \cdot \mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) > 0\}} - \mathcal{P}_m(\boldsymbol{\chi}) \right|^2 \middle| \mathbb{T}_m, \boldsymbol{\chi}, \{\boldsymbol{\chi}_i\}_{i \in \mathbb{T}_\ell} \right) \right] \\ & \leq \mathbb{E} \left[\mathbb{E} \left(\left| \frac{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) \cdot \widehat{\mathcal{P}}_{m,\ell}(\boldsymbol{\chi}) \cdot \mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) > 0\}}}{\mathcal{S}_{m,\ell}(\boldsymbol{\chi})} - \mathcal{P}_m(\boldsymbol{\chi}) \right|^2 \right. \right. \\ & \quad \left. \left. \times \mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) > 0\}} \middle| \mathbb{T}_m, \boldsymbol{\chi}, \{\boldsymbol{\chi}_i\}_{i \in \mathbb{T}_\ell} \right) + \mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) = 0\}} \right] \\ & = \mathbb{E} \left[\mathbb{E} \left(\left| \frac{\text{Bin}(\mathcal{S}_{m,\ell}(\boldsymbol{\chi}), \mathcal{P}_m(\boldsymbol{\chi}))}{\mathcal{S}_{m,\ell}(\boldsymbol{\chi})} - \mathcal{P}_m(\boldsymbol{\chi}) \right|^2 \right. \right. \\ & \quad \left. \left. \times \mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) > 0\}} \middle| \mathbb{T}_m, \boldsymbol{\chi}, \{\boldsymbol{\chi}_i\}_{i \in \mathbb{T}_\ell} \right) + \mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) = 0\}} \right] \\ & = \mathbb{E} \left[\frac{\mathcal{P}_m(\boldsymbol{\chi})(1 - \mathcal{P}_m(\boldsymbol{\chi}))}{\mathcal{S}_{m,\ell}(\boldsymbol{\chi})} \times \mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) > 0\}} \right] + \mathbb{P}\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) = 0\}, \end{aligned}$$

where the last line follows from the usual binomial variance formula. This last expression, in conjunction with (24), and the fact that $\mathcal{P}_m(\boldsymbol{\chi})(1 - \mathcal{P}_m(\boldsymbol{\chi})) \leq 1/4$ immediately yields

$$R_{n,2} \leq 2 \sqrt{\mathbb{E} \left[(4\mathcal{S}_{m,\ell}(\boldsymbol{\chi}))^{-1} \mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) > 0\}} \right]} + \mathbb{P}\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) = 0\} + 2\mathbb{P}\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) = 0\}. \quad (25)$$

But, upon replacing $\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) > 0\}$ by $\{1 \leq \mathcal{S}_{m,\ell}(\boldsymbol{\chi}) \leq k\} \cup \{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) > k\}$, for any arbitrary integer $k \geq 1$, one finds $\mathbb{E} \left[\frac{1}{\mathcal{S}_{m,\ell}(\boldsymbol{\chi})} \mathbb{I}_{\{\mathcal{S}_{m,\ell}(\boldsymbol{\chi}) > 0\}} \right] \leq \mathbb{P}\{1 \leq \mathcal{S}_{m,\ell}(\boldsymbol{\chi}) \leq k\} + k^{-1}$ holds for all $k \geq 1$. Therefore, by first choosing k large enough and then applying Assumption A, the bound in (25) can be made as small as desired. The proof of Theorem 1 now follows since, in view of Assumption B, we have $R_{n,3} \rightarrow 0$, as n (and thus m) $\rightarrow \infty$, where $R_{n,3}$ is as in (16). This completes the proof of Theorem 1. \square

PROOF OF LEMMA 1

Let $L(G_m^*)$ be as defined in (11). Also, let $G_{n,\nu_n(\boldsymbol{\chi})}$ be as given by the right side of (23) and put $L(G_{n,\nu_n(\boldsymbol{\chi})}|\mathbb{T}_n) = \mathbb{P}\{G_{n,\nu_n(\boldsymbol{\chi})}(\widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi})) \neq Y | \mathbb{T}_n\}$. Now, observe that $\mathbb{E}[\min_{1 \leq k \leq \ell} L(G_{n,k}|\mathbb{T}_n)] \leq \mathbb{E}[L(G_{n,\nu_n(\boldsymbol{\chi})}|\mathbb{T}_n)] = L(G_{n,\nu_n(\boldsymbol{\chi})})$, from which one obtains

$$R_{n,2} \leq L(G_{n,\nu_n(\boldsymbol{\chi})}) - L(G_m^*). \quad (26)$$

Furthermore, $R_{n,2} \geq 0$ which follows from the fact that

$$\begin{aligned} \mathbb{E}\left[\min_{1 \leq k \leq \ell} L(G_{n,k}|\mathbb{T}_n)\right] &= \mathbb{E}\left[L\left(\operatorname{argmin}_{G_{n,k} \in \{G_{n,1}, \dots, G_{n,\ell}\}} L(G_{n,k}|\mathbb{T}_n) \middle| \mathbb{T}_n\right)\right] \\ &= L\left(\operatorname{argmin}_{G_{n,k} \in \{G_{n,1}, \dots, G_{n,\ell}\}} L(G_{n,k}|\mathbb{T}_n)\right) \\ &\geq L(G_m^*), \end{aligned}$$

where the last line follows from Proposition 2 with n replaced by m . Next, observe that

$$\begin{aligned} R_{n,2} &\leq L(G_{n,\nu_n(\boldsymbol{\chi})}) - L(G_m^*), \quad (\text{by (26)}) \\ &= \mathbb{P}\{G_m^*(\widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi})) = Y\} - \mathbb{P}\{G_{n,\nu_n(\boldsymbol{\chi})}(\widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi})) = Y\} \\ &= \sum_{k=0,1} \mathbb{E}\left(\mathbb{I}_{\{G_m^*(\widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}))=k\}} \cap [Y=k]\} - \mathbb{I}_{\{G_{n,\nu_n(\boldsymbol{\chi})}(\widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}))=k\}} \cap [Y=k]\}\right) \\ &= \sum_{k=0,1} \mathbb{E}\left[\mathbb{E}\left(\mathbb{I}_{\{G_m^*(\widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}))=k\}} \cdot \mathbb{I}_{\{Y=k\}}\right.\right. \\ &\quad \left.\left. - \mathbb{I}_{\{G_{n,\nu_n(\boldsymbol{\chi})}(\widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}))=k\}} \cdot \mathbb{I}_{\{Y=k\}} \middle| \widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \right.\right. \\ &\quad \left.\left. \widehat{g}_{m,J}(\boldsymbol{\chi}), \left(\widehat{g}_{m,1}(\boldsymbol{\chi}_i), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}_i), Y_i\right)_{i: \boldsymbol{\chi}_i \in \mathbb{T}_\ell}\right)\right] \\ &= \sum_{k=0,1} \mathbb{E}\left[\left(\mathbb{I}_{\{G_m^*(\widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}))=k\}} - \mathbb{I}_{\{G_{n,\nu_n(\boldsymbol{\chi})}(\widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}))=k\}}\right)\right. \\ &\quad \left. \times \mathbb{E}\left(\mathbb{I}_{\{Y=k\}} \middle| \widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi})\right)\right] \\ &\quad (\text{because } Y \text{ is independent of } \widehat{g}_{m,1}(\boldsymbol{\chi}_i), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}_i), Y_i, \quad i: \boldsymbol{\chi}_i \in \mathbb{T}_\ell) \\ &= \mathbb{E}\left[\left(\mathbb{I}_{\{G_m^*(\widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}))=1\}} - \mathbb{I}_{\{G_{n,\nu_n(\boldsymbol{\chi})}(\widehat{g}_{m,1}(\boldsymbol{\chi}), \dots, \widehat{g}_{m,J}(\boldsymbol{\chi}))=1\}}\right) \cdot \mathcal{P}_m(\boldsymbol{\chi})\right] \end{aligned}$$

$$\begin{aligned}
& - \mathbb{E} \left[\left(\mathbb{I}_{\{G_m^*(\hat{g}_{m,1}(\boldsymbol{x}), \dots, \hat{g}_{m,J}(\boldsymbol{x}))=1\}} - \mathbb{I}_{\{G_{n,\nu_n}(\hat{g}_{m,1}(\boldsymbol{x}), \dots, \hat{g}_{m,J}(\boldsymbol{x}))=1\}} \right) \cdot (1 - \mathcal{P}_m(\boldsymbol{x})) \right] \\
& \quad (\text{where, as in Lemma 1, } \mathcal{P}_m(\boldsymbol{x}) = \mathbb{P}\{Y = 1 \mid \hat{g}_{m,1}(\boldsymbol{x}), \dots, \hat{g}_{m,J}(\boldsymbol{x})\}.) \\
& = \mathbb{E} \left[\left(\mathbb{I}_{\{G_m^*(\hat{g}_{m,1}(\boldsymbol{x}), \dots, \hat{g}_{m,J}(\boldsymbol{x}))=1\}} - \mathbb{I}_{\{G_{n,\nu_n}(\hat{g}_{m,1}(\boldsymbol{x}), \dots, \hat{g}_{m,J}(\boldsymbol{x}))=1\}} \right) \cdot (2\mathcal{P}_m(\boldsymbol{x}) - 1) \right] \\
& = 2 \mathbb{E} \left[\mathbb{I}_{\{G_m^*(\hat{g}_{m,1}(\boldsymbol{x}), \dots, \hat{g}_{m,J}(\boldsymbol{x})) \neq G_{n,\nu_n}(\hat{g}_{m,1}(\boldsymbol{x}), \dots, \hat{g}_{m,J}(\boldsymbol{x}))\}} \cdot \left| \mathcal{P}_m(\boldsymbol{x}) - 0.5 \right| \right] \\
& \leq 2 \mathbb{E} \left[\hat{\mathcal{P}}_{m,\ell}(\boldsymbol{x}) \cdot \mathbb{I}_{\{S_{m,\ell}(\boldsymbol{x}) > 0\}} + \bar{Y}_\ell \cdot \mathbb{I}_{\{S_{m,\ell}(\boldsymbol{x}) = 0\}} - \mathcal{P}_m(\boldsymbol{x}) \right],
\end{aligned}$$

where the last line follows from the definition of $G_{n,\nu_n}(\hat{g}_{m,1}(\boldsymbol{x}), \dots, \hat{g}_{m,J}(\boldsymbol{x}))$ in the far right side of (23), the definition of $G_m^*(\hat{g}_{m,1}(\boldsymbol{x}), \dots, \hat{g}_{m,J}(\boldsymbol{x}))$, and the fact that the inequality $|\mathcal{P}_m(\boldsymbol{x}) - 0.5| \leq \left| \hat{\mathcal{P}}_{m,\ell}(\boldsymbol{x}) \cdot \mathbb{I}_{\{S_{m,\ell}(\boldsymbol{x}) > 0\}} + \bar{Y}_\ell \mathbb{I}_{\{S_{m,\ell}(\boldsymbol{x}) = 0\}} - \mathcal{P}_m(\boldsymbol{x}) \right|$ holds on the set $\left\{ G_m^*(\hat{g}_{m,1}(\boldsymbol{x}), \dots, \hat{g}_{m,J}(\boldsymbol{x})) \neq G_{n,\nu_n}(\hat{g}_{m,1}(\boldsymbol{x}), \dots, \hat{g}_{m,J}(\boldsymbol{x})) \right\}$. \square

PROOF OF PROPOSITION 1

The proof Proposition 1 is similar to (and, in fact, much simpler than) that of Proposition 2, and will not be given.

PROOF OF PROPOSITION 2

The proof of this Proposition is similar to (and easier than) the proof of Lemma 1 and goes as follows. Let $G : \{0, 1\}^J \rightarrow \{0, 1\}$ be any combined classifier with error $L(G) = \mathbb{P}\{G(\hat{g}_{n,1}(\boldsymbol{x}), \dots, \hat{g}_{n,J}(\boldsymbol{x})) \neq Y\}$. Then, with G_n^* as in (3), one has

$$\begin{aligned}
& L(G) - L(G_n^*) \\
& = \mathbb{P}\{G_n^*(\hat{g}_{n,1}(\boldsymbol{x}), \dots, \hat{g}_{n,J}(\boldsymbol{x})) = Y\} - \mathbb{P}\{G(\hat{g}_{n,1}(\boldsymbol{x}), \dots, \hat{g}_{n,J}(\boldsymbol{x})) = Y\} \\
& = \sum_{k=0,1} \mathbb{E} \left[\mathbb{E} \left[\mathbb{I}_{\{G_n^*(\hat{g}_{n,1}(\boldsymbol{x}), \dots, \hat{g}_{n,J}(\boldsymbol{x}))=k\}} \cdot \mathbb{I}_{\{Y=k\}} \right. \right. \\
& \quad \left. \left. - \mathbb{I}_{\{G(\hat{g}_{n,1}(\boldsymbol{x}), \dots, \hat{g}_{n,J}(\boldsymbol{x}))=k\}} \cdot \mathbb{I}_{\{Y=k\}} \right| \hat{g}_{n,1}(\boldsymbol{x}), \dots, \hat{g}_{n,J}(\boldsymbol{x}) \right] \\
& = \mathbb{E} \left[\left(\mathbb{I}_{\{G_n^*(\hat{g}_{n,1}(\boldsymbol{x}), \dots, \hat{g}_{n,J}(\boldsymbol{x}))=1\}} - \mathbb{I}_{\{G(\hat{g}_{n,1}(\boldsymbol{x}), \dots, \hat{g}_{n,J}(\boldsymbol{x}))=1\}} \right) \cdot (2\mathcal{P}_n(\boldsymbol{x}) - 1) \right]
\end{aligned}$$

$$\begin{aligned}
&= 2 \mathbb{E} \left[\mathbb{I} \left\{ G(\hat{g}_{n,1}(\boldsymbol{x}), \dots, \hat{g}_{n,J}(\boldsymbol{x})) \neq G_n^*(\hat{g}_{n,1}(\boldsymbol{x}), \dots, \hat{g}_{n,J}(\boldsymbol{x})) \right\} \times \left| \mathcal{P}_n(\boldsymbol{x}) - 0.5 \right| \right] \\
&\geq 0, \quad \text{where } \mathcal{P}_n(\boldsymbol{x}) = \mathbb{P}\{Y = 1 \mid \hat{g}_{n,1}(\boldsymbol{x}), \dots, \hat{g}_{n,J}(\boldsymbol{x})\}.
\end{aligned}$$

□

Acknowledgements

This work was supported by the NSF under Grant DMS-1916161 of Majid Mojirsheibani.

Disclosure statement

The authors report that there are no potential conflicts of interest.

References

- [1] Abraham, C., Biau, G., and Cadre, B. On the kernel rule for functional classification. *Annals of the Institute of Statistical Mathematics*, 2006;58:619-633.
- [2] Balakrishnan, N. and Mojirsheibani, M. A simple method for combining estimates to improve the overall error rates in classification. *Computational Statistics*, 2015;30:1033-1049.
- [3] Biau, G., Devroye, L. and Lugosi, G. Consistency of random forests and other averaging classifiers. *Journal of Machine Learning Research*, 2008;9:2015-2033.
- [4] Biau, G., Fischer, A., Guedj, B. and Malley, J. D. COBRA: A combined regression strategy. *Journal of Multivariate Analysis*, 2016;146:18-28 (with supplementary material).
- [5] Boser, B., Guyon, I., and Vapnik, V. (1992) A training algorithm for optimal margin classifiers. In *Proceedings of the Fifth Annual Workshop on Computational Learning Theory*, Vol. 5, pp. 144 - 152; 1992.

- [6] Breiman, L. Stacked regression. *Machine Learning*, 1995;24:49-64.
- [7] Breiman, L. Bagging predictors. *Machine Learning*, 1996;24:123-140.
- [8] Breiman, L. Random Forests. *Machine Learning*, 2001;45:5-32.
- [9] Cérou, F. and Guyader, A. Nearest neighbor classification in infinite dimensions. *ESAIM: Probability and Statistics*, 2006;10:340-355.
- [10] Cholaquidis, A., Fraiman, R., Kalemkerian, J., and Llop, P. A nonlinear aggregation type classifier. *Journal of Multivariate Analysis*, 146, 269-281.
- [11] Devroye, L. and Györfi, L. *Nonparametric density estimation: The L_1 view*. NY: John Wiley & Sons; 1985.
- [12] Devroye, L., Györfi, L., and Lugosi, G. *A Probabilistic Theory of Pattern Recognition*, NY: Springer; 1996.
- [13] Fischer A. and Mougeot, M. Aggregation using input-output trade-off. *Journal of Statistical Planning and Inference*, 2019;200:1-19.
- [14] Giraud, C. *Introduction to High-Dimensional Statistics*, Chapman & Hall/CRC; 2014.
- [15] Hoeffding, W. Probability Inequalities for Sums of Bounded Random Variables. *Journal of the American Statistical Association*, 1963;58:13-30.
- [16] LeBlanc, M. and Tibshirani, R. Combining estimates in regression and classification. *Journal of the American Statistical Association*, 1996;91:1641-1650.
- [17] Lin, Y. and Jeon, Y. Random forests and adaptive nearest neighbors. *Journal of the American Statistical Association*, 2006;101:578-590.
- [18] Mojirsheibani, M. Combining classifiers based on discretization. *Journal of the American Statistical Association*, 1999;94:600-609.

- [19] Rahman, R., Dhruva, S.R., Ghosh S., and Pal, R. Functional random forest with applications in dose-response predictions. *Scientific Reports*, 2019;9(Article number: 1628).
- [20] Rokach, L. Taxonomy for characterizing ensemble methods in classification tasks: a review and annotated bibliography. *Computational Statistics & Data Analysis*, 2009;53:4046-4072.
- [21] Rokach, L. Ensemble-based classifiers. *Artificial Intelligence Review*, 2010;33:1-39.
- [22] Wolberg, W.H. and Mangasarian, O.L. Multisurface method of pattern separation for medical diagnosis applied to breast cytology. In *Proceedings of the National Academy of Sciences*, 1990;87:9193-9196.
- [23] Wolpert, D. Stacked Generalization. *Neural Networks*, 1992;5:241-259.
- [24] Xu, L., Kryzak, A., and Suen, C.Y. Methods of combining multiple classifiers and their applications to handwriting recognition. *IEEE Transactions on Systems, Man and Cybernetics*, 1992;22:418-435.