# Ensemble Performance Through the Lens of Linear Independence of Classifier Votes in Data Streams

Enes Bektas ⬤ and Fazli Can ⬤

*Abstract*—Ensemble learning improves classification performance by combining multiple base classifiers. While increasing the number of classifiers generally enhances accuracy, excessively large ensembles can lead to computational inefficiency and diminishing returns. This paper investigates the relationship between ensemble size and performance through the lens of linear independence among classifier votes in data streams. We propose that ensembles composed of linearly independent classifiers maximize representational capacity, particularly under a geometric model. We then generalize the importance of linear independence to the weighted majority voting problem. By modeling the probability of achieving linear independence among classifier outputs, we derive a theoretical framework that explains the trade-off between ensemble size and accuracy. Our analysis leads to a theoretical estimate of the ensemble size required to achieve a user-specified probability of linear independence. We validate our theory through experiments on both real-world and synthetic datasets using two ensemble methods, OzaBagging and GOOWE. Our results confirm that this theoretical estimate effectively identifies the point of performance saturation for robust ensembles like OzaBagging. Conversely, for complex weighting schemes like GOOWE, our framework reveals that high theoretical diversity can trigger algorithmic instability. Our implementation is publicly available to support reproducibility and future research[1].

*Index Terms*—Data stream classification, ensemble size, ensemble cardinality, law of diminishing returns, weighed majority voting

## I. INTRODUCTION

**E**NSEMBLE learning is a widely used technique in machine learning and statistics for improving the performance and robustness of predictive models. An ensemble combines the predictions of multiple models, often referred to as "base classifiers" or "learners," to form a more accurate and resilient predictive model [1]–[6]. In recent years, with the exponential increase in data volume and complexity, ensemble models have gained popularity over single classifiers. The foundational principle behind their success is the ability to mitigate the bias-variance trade-off more effectively than individual models [2]. This approach has been successfully applied in a variety of fields, including statistics, machine learning, pattern recognition, and knowledge discovery in databases [7]. Numerous studies have consistently demonstrated that ensembles generally outperform single classifiers in classification tasks [2], [5], [6], [8], [9].

[1]Click here to visit GitHub repository.

The efficacy of an ensemble is intrinsically linked to its construction methodology, which comprises two principal stages: the generation of diverse base classifiers and the subsequent combination of their predictions. Foundational generation methods documented in the literature include Bagging [10], Boosting [11], [12], and Random Forests [13]. The methods for combining predictions are broadly classified into weighting schemes and meta-combination techniques [4]. Weighting methods—such as Majority Voting, Performance Weighting, and Bayesian Combination—assign a specific weight to each base classifier, reflecting its contribution to the final prediction. In contrast, meta-combination techniques—such as Stacking and Arbiter Trees—employ a meta-learning approach, where a secondary model learns to optimally combine the base classifiers' outputs [2], [4]. Among these, majority voting is a simple yet remarkably effective technique, particularly when the base classifiers are diverse and independent. The theoretical underpinnings that explain the effectiveness of majority voting are explored in detail in [14]–[17]. This paper focuses specifically on weighting methods, with an emphasis on weighted majority voting.

While the fundamental concept of ensemble learning encourages the use of numerous classifiers, deploying an excessively large ensemble introduces significant computational overhead in terms of memory and processing time. Furthermore, the performance gains often diminish beyond a certain number of classifiers. Research into the relationship between ensemble size and performance has led to two main perspectives. The first posits that a smaller, carefully pruned subset of classifiers can perform comparably to the full ensemble, a concept often framed as "many could be similar to all" [18]–[21]. The second perspective, characterized as "many could be better than all," argues that a subset of classifiers can even outperform the original, complete ensemble [22], [23].

Conversely, other studies suggest that ensemble performance monotonically increases with ensemble size, although the marginal improvements progressively decrease [24]. Our work aligns with this latter view, providing a theoretical foundation to explain the observed diminishing returns. To our knowledge, the theoretical principles governing this trade-off between ensemble size and predictive performance remain an open area of investigation. This paper introduces a theoretical framework to clarify this relationship, positing that the concept of linear independence among classifier votes is central to understanding and optimizing ensemble performance.

The contributions of this study are the following. We

- Theoretically demonstrate the importance of linear independence among classifier votes, providing a new algebraic perspective on ensemble diversity and its impact on

performance.

- Develop a probabilistic framework to model the trade-off between ensemble size and linear independence, culminating in two metrics to estimate the point of full diversity: a precise formula (`INC`) and a computationally efficient, closed-form approximation (`SINC`).
- Validate our theoretical estimates through comprehensive experiments, showing they successfully identify the point of full performance saturation for robust ensembles (OzaBagging) and, in contrast, reveal the instability and performance degradation of complex weighting schemes (GOOWE) at high diversity.

The remainder of this paper is organized as follows. Section II reviews prior work related to ensemble size, and introduces the geometric framework for our analysis. In Section III, we develop our core theory on ensemble construction. Section IV discusses the practical implications of this theory. Section V presents experimental results to validate our claims. Finally, Section VI provides concluding remarks.

## II. RELATED WORK

This section reviews prior research relevant to our study. We first provide a broad overview of existing strategies for determining ensemble size, including ensemble selection, pruning, and other theoretical approaches. We then narrow our focus to the geometric framework for ensemble classification, discussing the conflicting recommendations that motivate our work.

### A. Determining Ensemble Size

The optimal size of an ensemble has been a subject of extensive investigation, leading to various strategies for its determination. A significant body of research has focused on *ensemble selection*, which involves constructing an optimal subset of classifiers from a larger pool. For instance, Ulaş et al. proposed an incremental construction method using accuracy, statistically significant improvement, and diversity as selection criteria [25]. Similarly, Xiao et al. introduced a dynamic selection approach that considers both accuracy and diversity, particularly for noisy data [26], while Yang employed Q-statistics as a diversity measure alongside accuracy for classifier selection [27].

Another prominent research direction is *ensemble pruning*, which aims to reduce the size of a pre-existing ensemble to improve computational efficiency without degrading predictive performance. The underlying principle is to first generate a large ensemble and subsequently eliminate redundant or underperforming members [20], [28]–[34]. However, a limitation of many pruning algorithms is their singular focus on performance, often neglecting the computational cost. Bhardwaj et al. highlighted this gap, arguing that the cost-effectiveness of an ensemble must be evaluated as a function of both its size and its accuracy [35].

Specific studies have also addressed the optimal number of trees in Random Forests. Latinne et al. applied the McNemar test to determine *a priori* the minimum ensemble size required to achieve a performance level comparable to that of much larger forests [36]. Oshiro et al. empirically demonstrated that while performance generally increases with the number of trees, the improvements become marginal beyond a certain threshold [37]. Probst and Boulesteix provided a theoretical perspective, showing that the expected error rate of a random forest is not necessarily a monotonic function of the number of trees [38].

Other researchers have explored theoretical and analytical approaches to determine ensemble size. Bax analyzed the majority voting mechanism and concluded that any odd number of classifiers could be optimal, proposing a validation-based selection method [39]. Hernàndez-Lobato et al. developed a method to estimate the number of classifiers needed for a parallel ensemble's vote to approximate the vote of an infinite-sized ensemble with a user-defined confidence level [24]. Fumera et al. conducted a theoretical analysis of bagging, modeling it as a linear combination of classifiers to derive the expected error as a function of ensemble size [40], [41]. In the context of data streams, Jackowski introduced a diversity measure based on classifier reactions to concept drift, which can also be used to control the ensemble size [42].

TABLE I
DEFINITION OF FREQUENTLY USED SYMBOLS

| Symbol | Definition |
|---|---|
| $\mathcal{D} = \{I_1, I_2, ...\}$ | The data stream. |
| $E = \{C_1, \ldots, C_n\}$ | The ensemble with $n$ component classifiers. |
| $I_k$ | An instance of data stream. |
| $m$ | Number of class labels in the data stream. |
| $n$ | Number of classifiers in the ensemble. |
| $o_k$ | The ideal vector in the form $(0, \ldots, 0, 1, 0, \ldots, 0)$. |
| $p_l$ | Probability of a classifier's vote being linearly dependent on an existing $l$-dimensional space. |
| $S_i$ | Vote vector of $i^{th}$ classifier in the ensemble. |
| $S_{ij}$ | $j^{th}$ component of $i^{th}$ classifier's vote. |
| $V_k$ | Vote of the ensemble for instance $I_k$ |
| $W_i$ | Weight of $i^{th}$ classifier in the ensemble. |
| PLI | Probability of Linear Independence. |
| INC | Ideal Number of Classifiers (calculated from Theorem 2). |
| SINC | Simplified Ideal Number of Classifiers. |

### B. The Geometric Framework and the Ensemble Size Debate

Our work builds upon a geometric framework for data fusion, originally introduced by Wu and Crestani for information retrieval systems [43]. This framework was subsequently adapted for ensemble classification by Bonab and Can, who treated each classifier's vote as a vector in a high-dimensional space and used weighted majority voting [44], [45]. Wu and Ding later extended this model to dataset-level classification [46]. We formally define the key components of this framework, with our notation summarized in Table I, before discussing its relevance to our study.

In this geometric model, for a data stream $\mathcal{D}$ with $m$ class labels and an ensemble $E$ with $n$ classifiers, each classifier's vote on an instance $I_k$ is represented as a vector $S_i$ in an $m$-dimensional space. Each vote vector is normalized such that the sum of its components is unity ($\sum_{j=1}^{m} S_{ij} = 1$). An

TABLE II
COMPARISON OF THEORETICAL FRAMEWORKS FOR ENSEMBLE SIZING

| Approach | Core Concept / Model | Key Goal / Finding |
|---|---|---|
| Hernández-Lobato et al. [24] | Statistical approximation | Estimates size $n$ needed to approximate the vote of an infinite-sized ensemble. |
| Fumera et al. [40], [41] | Models bagging as a linear combination | Derives the expected error as a function of ensemble size $n$. |
| Bax [39] | Majority vote analysis | Proposes that any odd number can be optimal; suggests a validation-based selection. |
| Wu & Crestani [43] | Geometric framework | Proves that adding more classifiers generally improves performance (monotonicity); implies larger is better. |
| Bonab & Can [44], [45] | Geometric framework | Minimizes Euclidean distance between ensemble vote and ideal vector; suggests a heuristic of $n = m$. |
| **This Paper** | **Geometric framework + Linear Algebra** | **Models the probability (PLI) of achieving $m$ linearly independent votes; provides a theoretical basis for performance saturation.** |

ideal vector, $o_k$, is defined as a vector that indicates the true class label of instance $I_k$. The ensemble's final vote, $V$, is a weighted linear combination of the individual classifier votes. The optimal weights, $W_i$, are calculated to minimize a loss function defined as the Euclidean distance between the ideal vector $o_k$ and the ensemble's vote $V$ [43], [45].

Within this framework, conflicting recommendations regarding the optimal ensemble size have emerged. Bonab and Can suggested that the number of classifiers should ideally equal the number of class labels ($n = m$) for optimal weight assignment [44], [45]. In contrast, Wu and Crestani proved that adding more classifiers generally improves performance, or at worst, leaves it unchanged [43]. This contradiction highlights a critical gap in understanding the principles that govern ensemble size. To resolve this ambiguity, we propose a probabilistic approach grounded in the concept of linear independence.

Table II provides a summary of these prior theoretical frameworks and contrasts them with the approach proposed in this paper.

## III. A THEORETICAL FRAMEWORK BASED ON LINEAR INDEPENDENCE

Within the geometric framework, the process of calculating optimal weights reveals that the linear independence of classifier votes is a critical property. The ensemble's final prediction is a linear combination of its constituent classifiers' votes. Consequently, the degree of linear dependence among these vote vectors directly constrains the vector space spanned by the ensemble, which in turn determines its expressive power and predictive capacity.

Before detailing our theorems, it is important to clarify the relationship between our general theory and the data stream context. The following framework, based on the linear independence of vote vectors, is a general algebraic model. It is applicable to any weighted majority voting ensemble, whether the data is processed in a batch or as a stream.

However, we center our analysis on data streams for two critical reasons. First, the problem of ensemble size and computational overhead is acute in stream mining, where

models must operate under strict resource constraints. Second, our instance-based analysis—which examines the properties of vote vectors for a given instance $I_k$—maps naturally to the online, instance-by-instance learning model required by data streams. Therefore, we develop this general theory to solve a problem to the data stream domain and validate it using standard stream classification methods and protocols. A summary of the theoretical contributions developed in this section is presented in Table III.

TABLE III
SUMMARY OF THEORETICAL CONTRIBUTIONS IN SECTION III

| Theorem | Key Contribution |
|---|---|
| Theorem 1 | Establishes that $m$ linearly independent votes are sufficient for an ensemble to perfectly represent the true class label for a single instance. |
| Theorem 2 | Provides a formula to calculate the probability of achieving $m$ linearly independent votes within an ensemble of size $n$, given the dependence probabilities $p_l$. |
| Theorem 3 | Proves that this probability converges to 1 as the ensemble size $n$ increases (assuming $p_l < 1$), theoretically justifying the diminishing returns in ensemble performance. |

The following theorem formalizes the importance of achieving a set of linearly independent votes.

### A. Perfect Classification Under Linear Independence

**Theorem 1:** For any instance $I_k$ from a data stream $\mathcal{D}$ with $m$ classes, if an ensemble $E$ contains at least $m$ classifiers that produce linearly independent vote vectors, then there exists a set of weights $W = \{W_1, ..., W_n\}$ such that $\sum_{i=1}^{n} W_i = 1$ and the resulting ensemble vote $V$ is identical to the ideal vector $o$.

**Proof:** Let us first consider the case where the ensemble size $n$ is equal to the number of classes $m$, and all $m$ classifiers provide linearly independent vote vectors $\{S_1, ..., S_m\}$ for instance $I_k$. Because these vectors form a basis for the $m$-

dimensional space, there exists a unique solution for the weights $W$ in the linear system:

$$
\begin{bmatrix} S_{11} & S_{21} & \cdots & S_{n1} \\ S_{12} & S_{22} & \cdots & S_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ S_{1m} & S_{2m} & \cdots & S_{nm} \end{bmatrix} \begin{bmatrix} W_1 \\ W_2 \\ \vdots \\ W_n \end{bmatrix} = \begin{bmatrix} o_1 \\ o_2 \\ \vdots \\ o_m \end{bmatrix} = o \quad (1)
$$

This system of equations is equivalent to $\sum_{i=1}^{n} W_i S_i = o$. Summing the components of the resulting vector equation from $j = 1$ to $m$ gives:

$$
\sum_{j=1}^{m} \left( \sum_{i=1}^{n} W_i S_{ij} \right) = \sum_{j=1}^{m} o_j = 1
$$

By rearranging the order of summation, we get:

$$
\sum_{i=1}^{n} W_i \left( \sum_{j=1}^{m} S_{ij} \right) = 1
$$

Given that each vote vector is normalized such that $\sum_{j=1}^{m} S_{ij} = 1$ for all classifiers $i$, the equation simplifies to:

$$
\sum_{i=1}^{n} W_i = 1
$$

For the case where $n > m$, if a subset of $m$ classifiers provides linearly independent votes, a solution can be found by assigning the ideal weights to this subset as described above, while setting the weights of the remaining $n$ - $m$ classifiers to zero. This completes the proof. $\square$

*Discussion:* Theorem 1 establishes that with $m$ linearly independent votes, an ensemble can perfectly classify any given instance $I_k$. However, the set of linearly independent classifiers may change from one instance to another. This dynamic nature necessitates a shift from a deterministic to a probabilistic perspective. We, therefore, seek to determine the probability of obtaining at least $m$ linearly independent votes within an ensemble of size $n$.

## B. Probability of Achieving Linear Independence

To model the ensemble construction process, we make a simplifying assumption that the probability of a classifier's vote being linearly dependent on the existing vote space depends solely on the current dimension of that space, and is uniform across all classifiers.

*Definition 1:* Let $p_l$ denote the probability that a new classifier's vote vector lies within the subspace spanned by $l$ existing linearly independent vote vectors. In other words, $p_l$ is the probability that adding a new classifier fails to increase the dimension from $l$ to $l+1$. We assume $p_l$ is constant for any classifier added to an $l$-dimensional span. These probabilities can be estimated empirically from a given dataset.

*Definition 2:* Let $\chi_k$ be the set of $(m-1)$-tuples of non-negative integers that sum to $k$:

$$
\chi_k = \{(x_1, \ldots, x_{m-1}) \mid \sum_{j=1}^{m-1} x_j = k; \; x_j \in \mathbb{N}_0\}
$$

*Theorem 2:* Given the probabilities $p_1, \ldots, p_{m-1}$, the probability of obtaining at least $m$ linearly independent votes from an ensemble of $n$ classifiers ($n \geq m$) is given by:

$$
P(n, m) = \left( \prod_{i=1}^{m-1} (1 - p_i) \right) \left( \sum_{k=0}^{n-m} \sum_{(x_1, \ldots, x_{m-1}) \in \chi_k} \prod_{j=1}^{m-1} p_j^{x_j} \right)
\tag{2}
$$

*Proof:* To understand the structure of this formula, consider the process of constructing a set of votes incrementally. We can model the growth of independence as a branching process, as illustrated in Figure 1. The process of achieving $m$ linearly independent votes requires successfully increasing the dimension of the spanned subspace from 1 to $m$. The probability of a new vote increasing the dimension from $i$ to $i+1$ is $(1 - p_i)$. Therefore, the probability of achieving $m$ linearly independent votes in the first $m$ attempts (i.e., with the first $m$ classifiers) is the product of these success probabilities, $\prod_{i=1}^{m-1} (1 - p_i)$. This term represents the base probability of the most direct path to success.
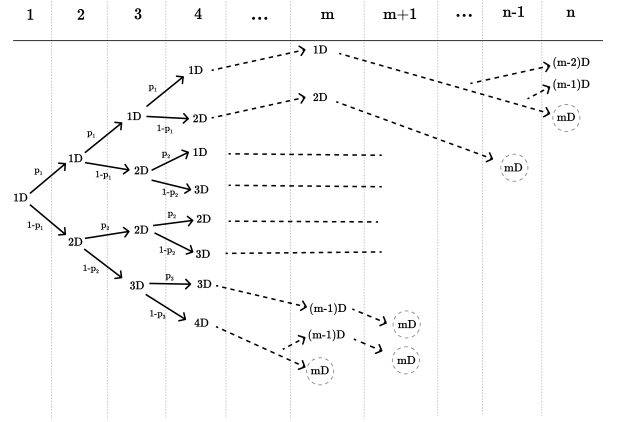


Fig. 1. A conceptual probability tree illustrating the paths to achieving an $m$-dimensional vote space with $n$ classifiers. Each branch represents the addition of a classifier, which either increases the dimension (D) or is linearly dependent.

If this is not achieved within the first $m$ classifiers, additional classifiers are needed. The term $\sum_{k=0}^{n-m} (\ldots)$ accounts for all scenarios where success is achieved using up to $n$ classifiers. The outer sum over $k$ considers the total number of "wasted" classifiers—those that did not increase the dimension of the spanned space. The inner summation over $\chi_k$ considers all possible ways these $k$ dependent votes can be distributed. Specifically, $x_j$ represents the number of times a new vote was found to be dependent on the existing $j$-dimensional subspace. The term $\prod_{j=1}^{m-1} p_j^{x_j}$ is the probability of one such specific sequence of dependencies occurring. Summing over all possible combinations of these dependent votes gives the complete probability of achieving success within an ensemble of size $n$. $\square$

*Discussion:* Theorem 2 provides a quantitative tool to evaluate the marginal benefit of adding a new classifier to the ensemble. It directly models the trade-off between ensemble size ($n$) and the likelihood of achieving a sufficiently diverse set of votes. To understand the long-term behavior of this

trade-off, we analyze the limit of this probability as the ensemble size grows infinitely large.

## C. Asymptotic Behavior of Ensemble Diversity

**Theorem 3:** Provided that $p_l \neq 1$ for all $l \in \{1, ..., m-1\}$, the probability of obtaining $m$ linearly independent votes converges to 1 as the ensemble size $n$ approaches infinity.

$$\lim_{n \to \infty} P(n, m) = 1 \tag{3}$$

**Proof:** The convergence depends on the values of $p_l$.

*Case 1: At least one $p_l = 1$.* If there exists an $l$ such that $p_l = 1$, it is impossible to surpass an $l$-dimensional space. The term $\prod_{i=1}^{m-1}(1 - p_i)$ becomes 0, and thus the limit is 0.

*Case 2: All $p_l = 0$.* If all $p_l = 0$, every classifier is guaranteed to be linearly independent of the previous ones. The term $\prod_{i=1}^{m-1}(1-p_i)$ becomes 1. Regarding the summation, any term where $k > 0$ implies at least one dependency occurred (some $x_j \geq 1$), causing the term to vanish since $p_j = 0$. Therefore, the entire summation collapses to the single term for $k = 0$ (where all $x_j = 0$), which equals 1. This results in a total probability of 1 for any $n \geq m$.

*Case 3: $0 \leq p_l < 1$ for all $l$, with at least one $p_l > 0$.* We prove this by induction on $m$. Let $S_m = \sum_{k=0}^{\infty} \sum_{(x_1,...,x_{m-1}) \in \chi_k} \prod_{j=1}^{m-1} p_j^{x_j}$. We must show that $\left(\prod_{i=1}^{m-1}(1 - p_i)\right) S_m = 1$.

*Base Case (m=2):* The expression becomes $(1 - p_1) \sum_{k=0}^{\infty} p_1^k$. Since $0 \leq p_1 < 1$, this is a geometric series which converges to $\frac{1}{1-p_1}$. Thus, $(1 - p_1)\frac{1}{1-p_1} = 1$. The theorem holds.

*Inductive Step:* Assume the theorem holds for $m - 1 = d$, i.e., $S_d = \frac{1}{\prod_{j=1}^{d-1}(1-p_j)}$. We must show it holds for $m = d + 1$. Since the outer sum extends over all possible values of $k$ (from 0 to $\infty$), the constraint $\sum x_j = k$ covers the entire space of non-negative integers $\mathbb{N}_0^d$. This allows us to decouple $x_d$ from the other variables and factor out the terms involving $p_d$:

$$S_{d+1} = \sum_{k=0}^{\infty} \sum_{(x_1,...,x_d) \in \chi_k} \prod_{j=1}^{d} p_j^{x_j}$$

$$= \sum_{x_d=0}^{\infty} p_d^{x_d} \left( \sum_{k'=0}^{\infty} \sum_{(x_1,...,x_{d-1}) \in \chi_{k'}} \prod_{j=1}^{d-1} p_j^{x_j} \right)$$

$$= \left( \sum_{x_d=0}^{\infty} p_d^{x_d} \right) (S_d)$$

$$= \left( \frac{1}{1 - p_d} \right) S_d$$

By the inductive hypothesis, $S_d = \frac{1}{\prod_{j=1}^{d-1}(1-p_j)}$. Substituting this in, we get:

$$S_{d+1} = \frac{1}{1 - p_d} \cdot \frac{1}{\prod_{j=1}^{d-1}(1-p_j)} = \frac{1}{\prod_{j=1}^{d}(1-p_j)}$$

Therefore, $\left(\prod_{i=1}^{d}(1-p_i)\right) S_{d+1} = 1$. The induction holds.
□

**Discussion:** Theorems 2 and 3 provide theoretical justification for the empirical observation that larger ensembles tend to perform better, but with diminishing returns. Theorem 3 proves that, under the reasonable assumption that no feature space is inherently impassable ($p_l \neq 1$), achieving a full-rank set of votes is a probabilistic certainty given a large enough ensemble. This convergence allows us to reframe the problem of finding the optimal ensemble size: instead of seeking a universal number, we can determine the minimum ensemble size $n$ required to achieve a desired target probability (e.g., 99.9%) of obtaining $m$ linearly independent votes. This provides a principled and practical method for guiding ensemble construction.

## IV. PRACTICAL IMPLICATIONS

The preceding theoretical development has direct practical applications for ensemble design. A foundational principle in ensemble learning is that diversity among base classifiers is essential for robust performance [47]. Our framework contributes to this understanding by providing a formal, algebraic interpretation and a mathematical measure for diversity: The linear independence of classifier vote vectors. We claim that a higher probability of achieving linear independence directly corresponds to increased ensemble diversity and, consequently, to a more powerful ensemble, up to the point of performance saturation.

Theorem 2 offers a quantitative metric to assess this form of diversity. This allows for a quantitative approach to determining the required ensemble size, $n$, moving beyond the fixed heuristic of $n = m$ previously derived from this geometric framework [44], [45]. Our methodology involves specifying a desired confidence threshold, $T$, which represents the minimum acceptable probability of obtaining $m$ linearly independent votes. By applying the formula from Theorem 2, one can then calculate the minimum number of classifiers, $n$, required to ensure this probability meets or exceeds the threshold $T$. This procedure enables the construction of ensembles that are provably diverse with a user-specified degree of confidence, providing a solution to the trade-off between ensemble size and performance.

The remainder of this section details these practical implications. First, Subsection A establishes the generalizability of our theorems beyond the specific geometric framework, demonstrating their relevance to all weighted majority voting schemes. Subsection B then provides a worked example illustrating how to determine the minimum ensemble size for a given probability threshold. Following this, Subsection C analyzes the computational complexity of the formula presented in Theorem 2. To address the potential for high computational costs, Subsection D introduces a simplified, closed-form approximation of the formula. Finally, Subsection E presents a practical algorithm for empirically estimating the crucial $p_l$ parameters from a dataset, making the entire framework applicable to real-world problems.

## A. Generalizing Beyond the Geometric Model

While our theorems were derived within a specific geometric framework, their central conclusion—the critical importance of linear independence—extends to the broader class of all weighted majority voting ensembles.

In any weighted majority voting scheme, the ensemble's final decision is based on an aggregated vote vector, $V$, which is a linear combination of the vote vectors, $\{S_1, S_2, ..., S_n\}$, produced by the $n$ base classifiers. Each classifier's output, whether a probability distribution across classes or a one-hot encoding of its predicted class, can be naturally represented as a vector within an $m$-dimensional class space. The set of all possible aggregated vote vectors, $V$, that the ensemble can produce is therefore confined to the span of the base classifiers' vote vectors: $V \in \mathrm{span}(\{S_1, S_2, ..., S_n\})$.

The classification rule in such a scheme is to select the class corresponding to the index of the maximum component in the vector $V$. The fundamental goal is to select weights such that this maximum component aligns with the true class label of a given instance. However, the ability of the ensemble to achieve this is fundamentally constrained by the dimensionality of its spanned space.

If the set of base classifiers fails to produce at least $m$ linearly independent vote vectors, their span will be a subspace of dimension less than $m$. Consequently, there will exist regions of the $m$-dimensional class space that are unreachable by the ensemble. For an instance whose true class corresponds to an outcome in one of these unreachable regions, it is impossible to find a set of weights that will result in a vector $V$ where the true class's component is maximal. Regardless of the weighting strategy, the ensemble is representationally incapable of making the correct prediction, leading to an irreducible error. Common constraints on the weights, such as requiring them to sum to one ($\sum W_i = 1$), do not alter this conclusion, as it is just a scaling factor that does not change the index of the maximum value.

This demonstrates that achieving $m$ linearly independent votes is a necessary condition for an ensemble to have the capacity to correctly classify all possible outcomes. As established by Theorem 1, the presence of $m$ linearly independent votes guarantees this representational capacity is met. Therefore, the principles of linear independence and the probabilistic framework developed in this paper are not confined to the initial geometric model but are of universal importance to the weighted majority voting problem.

### B. Example: Computing Minimum Ensemble Size for Target Confidence

To demonstrate the practical utility of our theoretical framework, we present a concrete example of determining the minimum ensemble size required to meet a specific confidence level.

Consider a binary classification problem where the number of classes, $m$, is 2. Let us assume that through empirical analysis of a given dataset, the probability of a new classifier's vote being linearly dependent on a single existing vote vector is found to be $p_1 = 0.5$. Our objective is to calculate the minimum ensemble size, $n$, required to achieve a probability of at least 99% ($T = 0.99$) of obtaining two linearly independent votes.

For the case where $m = 2$, the general probability formula from Theorem 2 simplifies to the sum of a geometric series:

$$P(n, 2) = (1 - p_1) \sum_{k=0}^{n-2} p_1^k$$

We apply this formula iteratively, increasing the ensemble size $n$ and calculating the resulting probability until it meets or exceeds our target threshold of 0.99. With $p_1 = 0.5$:

- For $n = 2$: $P(2, 2) = (1 - 0.5) \cdot (0.5^0) = 0.5$
- For $n = 3$: $P(3, 3) = (1 - 0.5) \cdot (0.5^0 + 0.5^1) = 0.75$
- For $n = 4$: $P(4, 4) = (1 - 0.5) \cdot (0.5^0 + 0.5^1 + 0.5^2) = 0.875$

Continuing this process, we find that for an ensemble of size $n = 7$, the probability is approximately $0.9844$, which is still below our target. Incrementing the size one final time:

- For $n = 8$: $P(8, 8) = (1 - 0.5) \cdot \sum_{k=0}^{6} 0.5^k \approx 0.9922$

Since $P(8, 2) > 0.99$, we conclude that a minimum of 8 classifiers are required to satisfy the condition of having two linearly independent votes with at least 99% confidence. This example illustrates how our framework provides a methodology for ensemble sizing.

### C. Computational Complexity Analysis

This section analyzes the computational complexity of iteratively calculating the probability derived in Theorem 2. The primary computational burden lies within the nested summation:

$$\sum_{k=0}^{n-m} \left( \sum_{(x_1, ..., x_{m-1}) \in \chi_k} \prod_{j=1}^{m-1} p_j^{x_j} \right) \quad (4)$$

To determine the probability for an ensemble of size $n$, we must evaluate the inner sum for each value of $k$ from $0$ to $n - m$. The complexity of this operation is determined by the number of terms in the inner sum, which corresponds to the number of ways the integer $k$ can be expressed as the sum of $m - 1$ non-negative integers. This is a classic combinatorial problem known as "stars and bars." The number of terms for a given $k$ is given by the multiset coefficient:

$$|\chi_k| = \binom{k + (m-1) - 1}{(m-1) - 1} = \binom{k + m - 2}{m - 2} \quad (5)$$

The total number of product terms that must be computed to find the probability for an ensemble of size $n$ is the sum of these binomial coefficients over the full range of $k$. By applying the hockey-stick identity, this summation simplifies to a single binomial coefficient:

$$\text{Total Terms} = \sum_{k=0}^{n-m} \binom{k + m - 2}{m - 2} = \binom{n - 1}{m - 1} \quad (6)$$

The complexity is therefore determined by the asymptotic behavior of this resulting term. The binomial coefficient $\binom{n-1}{m-1}$ can be expanded as:

$$\binom{n - 1}{m - 1} = \frac{(n - 1)(n - 2) \cdots (n - m + 1)}{(m - 1)!}$$

For a fixed number of classes $m$, the denominator $(m-1)!$ is a constant, while the numerator is a polynomial in $n$ of degree $m-1$. This leads to an overall time complexity of $O(n^{m-1})$ for calculating the probability up to an ensemble size of $n$.

As this analysis indicates, the direct application of the formula becomes computationally intensive for problems with a large number of classes $(m)$, necessitating the development of a more efficient computational approach.

### D. Closed-Form Approximation for Uniform Dependence Probability

The $O(n^{m-1})$ complexity of the formula derived from Theorem 2 renders it computationally expensive for datasets with a large number of classes $(m)$. To overcome this barrier, we introduce a simplifying assumption: the probability of linear dependency is uniform across all dimensions. That is, we assume $p_l = p$ for all $l \in \{1, ..., m-1\}$, where $p$ represents a single, average probability that any new classifier's vote will be linearly dependent on the existing vote space.

This assumption allows us to transform the original nested summation into a closed-form expression. The derivation proceeds as follows. First, we substitute the uniform probability $p$ into the general formula:

$$P(n, m) = (1-p)^{m-1} \sum_{k=0}^{n-m} \left( \sum_{(x_1,...,x_{m-1}) \in \chi_k} \prod_{j=1}^{m-1} p^{x_j} \right)$$

$$= (1-p)^{m-1} \sum_{k=0}^{n-m} p^k \left( \sum_{(x_1,...,x_{m-1}) \in \chi_k} 1 \right)$$

The inner sum is now simply the number of terms, $|\chi_k|$, which is given by the binomial coefficient $\binom{k+m-2}{m-2}$. This yields:

$$P(n, m) = (1-p)^{m-1} \sum_{k=0}^{n-m} \binom{k+m-2}{m-2} p^k \qquad (7)$$

The key insight for further simplification is to recognize the term $\binom{k+m-2}{m-2} p^k$ as being related to the derivative of a geometric series. Specifically, the polynomial part can be generated by differentiation:

$$\binom{k+m-2}{m-2} p^k = \frac{(k+m-2)!}{k!(m-2)!} p^k$$

$$= \frac{1}{(m-2)!} \frac{d^{m-2}}{dp^{m-2}} (p^{k+m-2})$$

By substituting this back into the summation and interchanging the sum and derivative operators, we get:

$$P(n, m) = \frac{(1-p)^{m-1}}{(m-2)!} \sum_{k=0}^{n-m} \frac{d^{m-2}}{dp^{m-2}} (p^{k+m-2})$$

$$= \frac{(1-p)^{m-1}}{(m-2)!} \frac{d^{m-2}}{dp^{m-2}} \left( \sum_{k=0}^{n-m} p^{k+m-2} \right)$$

The summation inside the derivative is now a standard geometric series. Summing this series gives the final closed-form expression:

$$P(n, m) = \frac{(1-p)^{m-1}}{(m-2)!} \frac{d^{m-2}}{dp^{m-2}} \left( p^{m-2} \frac{1-p^{n-m+1}}{1-p} \right) \qquad (8)$$

This simplified formula allows the probability to be calculated without explicit summation over $n$. The complexity is now dominated by the $(m-2)$-th derivative of a rational function, which, for a fixed $m$, is independent of the ensemble size $n$.

This reduces the complexity from $O(n^{m-1})$ to $O(1)$ with respect to $n$, making the calculation efficient even for very large ensembles.

### E. Empirical Estimation of Linear Dependence Probabilities ($p_l$)

---
**Algorithm 1** Algorithm for Empirically Estimating $p_l$ Values
---
1: Let $m \leftarrow$ number of class labels
2: $counts\_dependent \leftarrow$ zero array of size $m-1$
3: $counts\_total \leftarrow$ zero array of size $m-1$
4: **for** each instance $I$ in the dataset **do**
5:      $voteMatrix \leftarrow n \times m$ matrix of classifier votes for $I$
6:      $curMatrix \leftarrow$ empty matrix
7:      $prevRank \leftarrow 0$
8:      **for** $i = 1 \rightarrow n$ **do**
9:          Append $i$-th row of $voteMatrix$ to $curMatrix$
10:          $curRank \leftarrow \text{rank}(curMatrix)$
11:          **if** $curRank < m$ **then**
12:             $counts\_total[prevRank] \leftarrow counts\_total[prevRank] + 1$
13:             **if** $curRank == prevRank$ **then**
14:                 $counts\_dependent[prevRank] \leftarrow counts\_dependent[prevRank] + 1$
15:             **end if**
16:          **else**
17:             **break**     ▷ Matrix has reached full rank
18:          **end if**
19:          $prevRank \leftarrow curRank$
20:      **end for**
21: **end for**
22: $p\_values \leftarrow$ zero array of size $m-1$
23: **for** $l = 0 \rightarrow m-2$ **do**
24:      **if** $counts\_total[l] > 0$ **then**
25:          $p\_values[l] \leftarrow counts\_dependent[l]/counts\_total[l]$
26:      **else**
27:          $p\_values[l] \leftarrow 1.0$     ▷ If dimension $l+1$ was never reached
28:      **end if**
29: **end for**
30: **return** $p\_values$ ▷ $p\_values[l]$ is the estimate for $p_{l+1}$
---

To apply the theoretical framework to real-world data, the linear dependence probabilities, $p_l$, must be estimated. We propose an empirical method, detailed in Algorithm 1, to calculate these values from a given dataset and ensemble. The fundamental principle is to simulate the incremental construction of a basis of vote vectors for each data instance and to aggregate the observed frequencies of linear dependence at each dimensional step.

The algorithm initializes two arrays, `counts_dependent` and `counts_total` (Lines 2-3), to store counts for each potential dimension from 0 to

$m - 2$. It then iterates through every instance $I$ in the dataset (Line 4). For each instance, it gathers the $n \times m$ matrix of vote vectors (Line 5) and initializes an empty matrix `curMatrix` and a rank tracker `prevRank` (Lines 6-7). The core logic resides in the inner loop (Lines 8-20), which progressively builds `curMatrix` by adding one vote vector at a time (Line 9) and computing its rank, `curRank` (Line 10).

Crucially, each time a vector is added when the previous rank was `prevRank`, the algorithm increments the total count for attempts to expand that dimension (`counts_total[prevRank]`, Line 12). If adding the vector does not increase the rank (`curRank == prevRank`, Line 13), it signifies linear dependence, and the corresponding `counts_dependent[prevRank]` is incremented (Line 14). The loop breaks early if the matrix reaches full rank $m$ (Line 17). The `prevRank` is updated at the end of each iteration (Line 19).

After processing all instances, the algorithm calculates the final probabilities (Lines 22-30). For each dimension $l$ from 0 to $m - 2$, the probability $p_{l+1}$ (stored in `p_values[l]`) is computed as the ratio of dependent counts to total counts (Line 25), provided that attempts were actually made to expand the $(l + 1)$-dimensional space (Line 24). If a dimension $(l + 1)$ was never reached (i.e., `counts_total[l]` is 0), $p_{l+1}$ is conservatively set to 1.0 (Line 27). The resulting array `p_values` contains the empirical estimates for $p_1, \ldots, p_{m-1}$. This procedure provides the necessary empirical parameters for applying our theoretical model.

## V. EXPERIMENTS

In this section, we present a series of experiments designed to empirically validate our theoretical framework. The primary objectives are threefold: 1) to investigate the relationship between the Probability of Linear Independence (PLI) and ensemble accuracy, 2) to demonstrate the practical utility of our proposed method for determining ensemble size, and 3) to explore the behavior of our framework across different ensemble methods and dataset characteristics.

To this end, Subsection A details our experimental setup, including the ensemble methods, datasets, and evaluation protocols used. Subsection B presents the empirical results, and Subsection C provides a detailed discussion of these findings, validating our theoretical model as a practical heuristic for ensemble sizing.

### A. Experimental Setup

*1) Ensemble Methods and Base Learners:* Two distinct ensemble methods were employed: OzaBagging [48] and GOOWE [49]. OzaBagging is a variant of bootstrap aggregating that combines classifier votes via majority voting. It is important to note that majority voting is a special case of weighted majority voting where all weights are equal, thus making OzaBagging a suitable method for validating our theoretical claims regarding the generalizability of linear independence across weighted voting schemes. GOOWE was selected as it directly implements the geometric weighting

framework that inspired our theoretical development. To support the theoretical assumption of consistent probabilities of dependence ($p_l$) across classifiers, all ensembles are constructed using Hoeffding Trees as the base classifier [50].

*2) Datasets and Evaluation:* Performance was assessed on twelve datasets, detailed in Table IV. These datasets were selected to provide a diverse and representative evaluation, covering a variety of data types and patterns. This includes six real-world benchmarks chosen from different contexts (e.g., Airlines, Poker, Electricity) to ensure variety in the number of attributes, instances, and class label distributions. This real-world set is complemented by six synthetic datasets, generated using the scikit-multiflow library's random RBF generator [51]. This synthetic data allows us to systematically control for and isolate the effect of the number of classes ($m$), a key parameter in our theoretical framework.

TABLE IV
DATASET SPECIFICATIONS. THE UPPER HALF IS REAL-LIFE, THE LOWER HALF IS SYNTHETIC DATASETS

| Dataset | #Instance | #Attr. | #Class Labels (m) |
|---|---|---|---|
| Airlines | 539,383 | 7 | 2 |
| Click Prediction | 399,482 | 11 | 2 |
| Electricity | 45,312 | 6 | 2 |
| Covtype | 581,012 | 54 | 7 |
| Poker | 829,201 | 10 | 10 |
| Rialto | 82,250 | 27 | 10 |
| RBF2 | 1,000,000 | 20 | 2 |
| RBF4 | 1,000,000 | 20 | 4 |
| RBF8 | 1,000,000 | 20 | 8 |
| RBF16 | 1,000,000 | 20 | 16 |
| RBF32 | 1,000,000 | 20 | 32 |
| RBF64 | 1,000,000 | 20 | 64 |

All experiments were evaluated using the prequential (interleaved test-then-train) methodology, which is standard for data stream classification. For each dataset and ensemble method, the experiments were repeated 10 times with different random seeds, and we report the average accuracy and standard deviation.

*3) PLI and Ideal Ensemble Size Calculation:* The probabilities of linear dependence, $p_l$, were empirically estimated for each dataset and ensemble size using Algorithm 1. The Probability of Linear Independence (PLI) plotted for an ensemble of size $n$ was calculated using the specific $p_l$ values measured at that size. We also calculated two theoretical estimates for the "Ideal Number of Classifiers" required to achieve a PLI of 0.9999. We selected this threshold ($T = 0.9999$) to identify the point of practical convergence, ensuring that the ensemble has effectively maximized its algebraic representational capacity.

- INC (Ideal Number of Classifiers): Calculated using Theorem 2, with $p_l$ values that were averaged across all tested ensemble sizes (from 2 to 128) to provide a single, robust estimate for the dataset.
- SINC (Simplified INC): Calculated using our simplified formula, with a single probability $p$ derived from the average of all $p_l$ values.

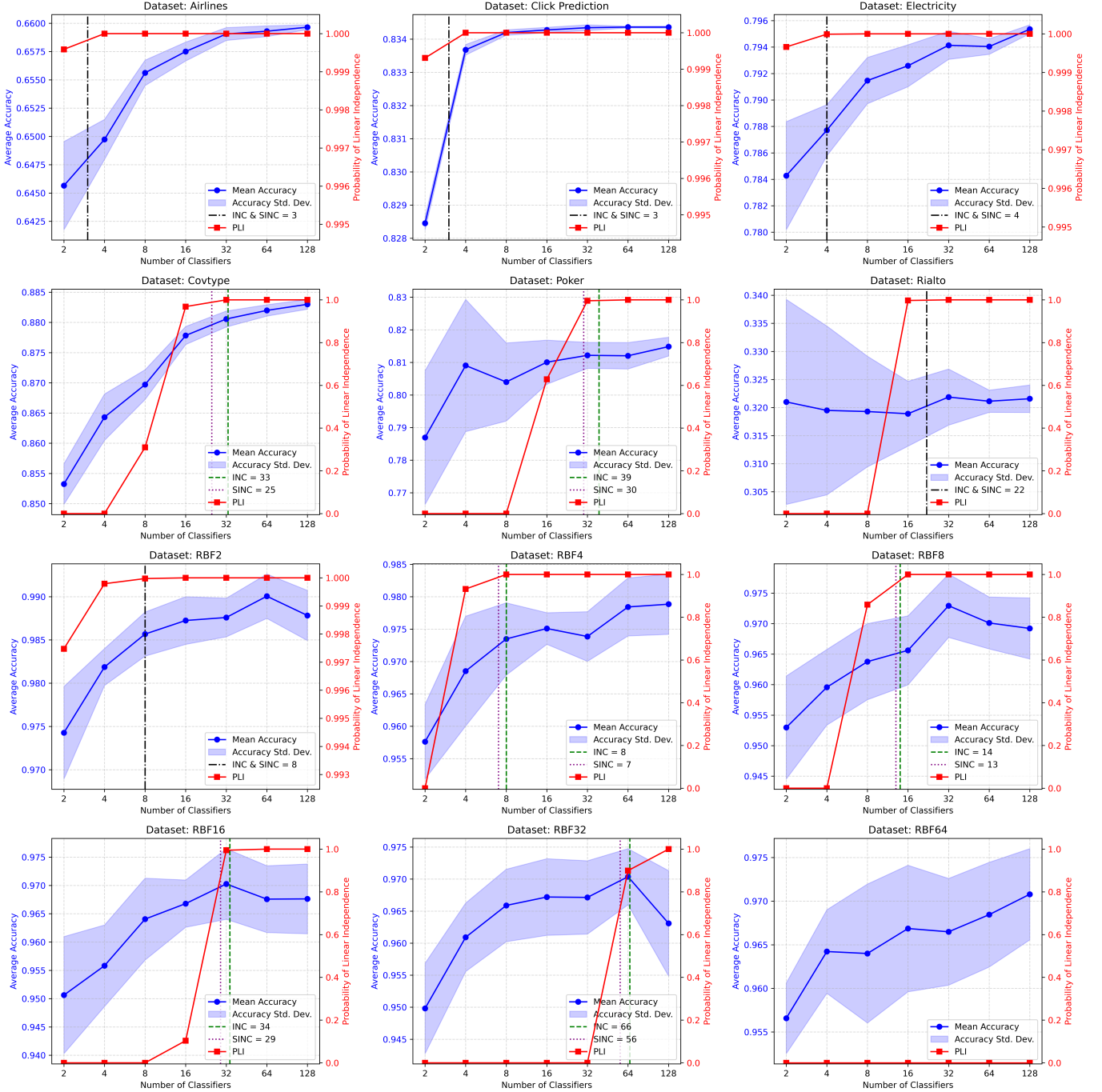All experimental code is publicly available for reproducibil-

Fig. 2. Performance of OzaBagging. Each subplot shows the average accuracy (blue line, left y-axis) with standard deviation represented by the shaded area, and the Probability of Linear Independence (PLI) (red, right y-axis) as a function of ensemble size. **Note that PLI is 0 for all** $n < m$ **(where** $m$ **is the number of classes), as** $m$ **linearly independent vectors cannot be obtained from an ensemble smaller than** $m$**.** The vertical dashed lines indicate the theoretically derived Ideal Number of Classifiers (INC, green; SINC, purple) for a PLI threshold of 0.9999. Note: If INC and SINC coincide, only a single black line is shown. No INC/SINC lines are shown for RBF64 as the PLI remains 0 for all tested ensemble sizes.
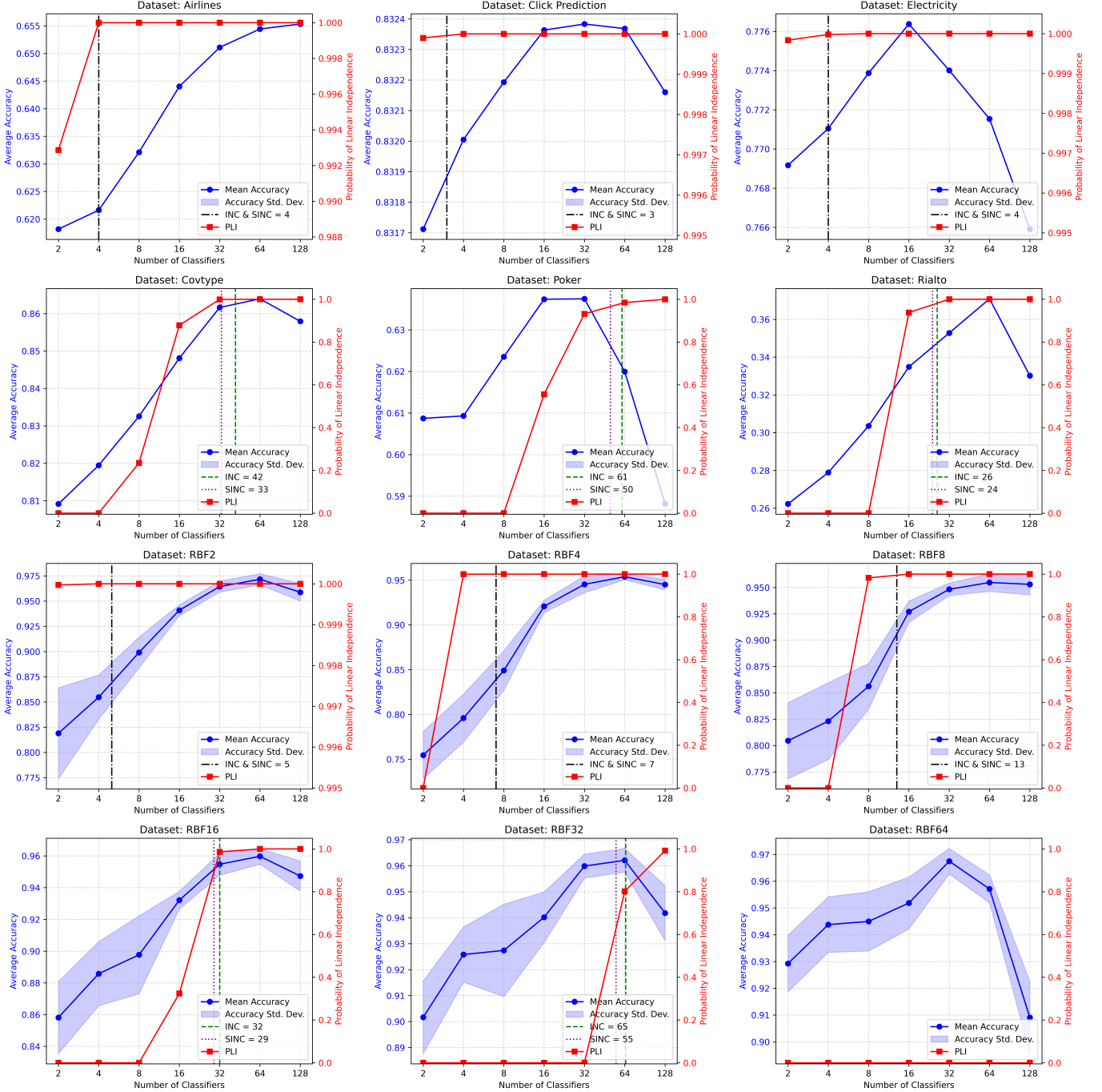
Fig. 3. Performance of GOOWE. Each subplot shows the average accuracy (blue line, left y-axis) with standard deviation represented by the shaded area, and the Probability of Linear Independence (PLI) (red, right y-axis) as a function of ensemble size. **Note that PLI is 0 for all** $n < m$ **(where** $m$ **is the number of classes), as** $m$ **linearly independent vectors cannot be obtained from an ensemble smaller than** $m$**.** The vertical dashed lines indicate the theoretically derived Ideal Number of Classifiers (INC, green; SINC, purple) for a PLI threshold of 0.9999. Note: If INC and SINC coincide, only a single black line is shown. No INC/SINC lines are shown for RBF64 as the PLI remains 0 for all tested ensemble sizes. Note the instances of performance degradation at larger ensemble sizes, which contrasts with the behavior of OzaBagging.

TABLE V
VALIDATION OF THEORETICAL ENSEMBLE SIZE AS A PRACTICAL HEURISTIC

| Dataset | Method | m | SINC | INC | $n_{INC}$ | Accuracy @$n_{INC}$ (% of Max) | Correlation |
|---|---|---|---|---|---|---|---|
| Airlines | OzaBagging | 2 | 3 | 3 | 4 | 98.4969 | 0.775 |
| Airlines | GOOWE | 2 | 4 | 4 | 4 | 94.8549 | 0.605 |
| Click Prediction | OzaBagging | 2 | 3 | 3 | 4 | 99.9181 | 0.994 |
| Click Prediction | GOOWE | 2 | 3 | 3 | 4 | 99.9545 | 0.824 |
| Electricity | OzaBagging | 2 | 4 | 4 | 4 | 99.0396 | 0.797 |
| Electricity | GOOWE | 2 | 4 | 4 | 4 | 99.3130 | 0.337 |
| Covtype | OzaBagging | 7 | 25 | 33 | 32 | 99.7248 | 0.940 |
| Covtype | GOOWE | 7 | 33 | 42 | 32 | 99.7320 | 0.971 |
| Poker | OzaBagging | 10 | 30 | 39 | 32 | 99.6704 | 0.698 |
| Poker | GOOWE | 10 | 50 | 61 | 64 | 97.2533 | 0.034 |
| Rialto | OzaBagging | 10 | 22 | 22 | 16 | 99.0780 | 0.418 |
| Rialto | GOOWE | 10 | 24 | 26 | 32 | 95.1252 | 0.897 |
| RBF2 | OzaBagging | 2 | 8 | 8 | 8 | 99.5581 | 0.913 |
| RBF2 | GOOWE | 2 | 5 | 5 | 4 | 87.9696 | 0.714 |
| RBF4 | OzaBagging | 4 | 7 | 8 | 8 | 99.4517 | 0.905 |
| RBF4 | GOOWE | 4 | 7 | 7 | 8 | 89.0144 | 0.687 |
| RBF8 | OzaBagging | 8 | 13 | 14 | 16 | 99.2480 | 0.885 |
| RBF8 | GOOWE | 8 | 13 | 13 | 16 | 97.1040 | 0.857 |
| RBF16 | OzaBagging | 16 | 29 | 34 | 32 | 100.0000 | 0.709 |
| RBF16 | GOOWE | 16 | 29 | 32 | 32 | 99.4795 | 0.912 |
| RBF32 | OzaBagging | 32 | 56 | 66 | 64 | 100.0000 | 0.307 |
| RBF32 | GOOWE | 32 | 55 | 65 | 64 | 100.0000 | 0.448 |
| RBF64 | OzaBagging | 64 | – | – | – | – | – |
| RBF64 | GOOWE | 64 | – | – | – | – | – |

*Note:* INC (Ideal Number of Classifiers) and SINC (Simplified Ideal Number of classifiers) are the theoretical ensemble sizes calculated to achieve a PLI threshold of 0.9999. $n_{INC}$ is the closest tested ensemble size (from $n = 2$ to $n = 128$) to the INC value. The Accuracy @$n_{INC}$ (% of Max) column shows the average accuracy achieved at $n_{INC}$, expressed as a percentage of the maximum average accuracy observed for that dataset. Correlation is the Pearson correlation between PLI and Average Accuracy across all tested ensemble sizes.

ity[2].

### B. Results

The experimental results for OzaBagging and GOOWE are presented visually in Figure 2 and Figure 3, respectively. We begin with OzaBagging to establish a baseline under equal-weighted majority voting, followed by GOOWE to examine performance under a complex, geometry-driven weighting scheme.

Visually, OzaBagging (Figure 2) displays a monotonic increase in accuracy as the ensemble size ($n$) grows, typically stabilizing at a high value. In contrast, GOOWE (Figure 3) exhibits more complex behavior; while accuracy initially rises, it frequently plateaus earlier or degrades at larger ensemble sizes (e.g., Poker, RBF32). In both figures, the Probability of Linear Independence (PLI), shown in red, consistently rises from 0 toward 1 as $n$ increases, except for the high-dimensional RBF64 dataset where it remains zero.

A quantitative summary is provided in Table V. To link these empirical results to our theory, we utilize the $n_{INC}$ metric—the tested ensemble size ($n \in \{2, 4, ..., 128\}$) closest to the theoretical INC value. The table reports the performance at this point as Accuracy @$n_{INC}$ (% of Max), expressing the result as a percentage of the maximum accuracy achieved for

[2]Click here to visit GitHub repository.

that dataset. Finally, we report the Pearson correlation between the Probability of Linear Independence (PLI) and the average accuracy across all tested ensemble sizes ($n = 2$ to 128), providing a statistical measure of the relationship between diversity and performance.

### C. Discussion

In this section, we interpret these empirical findings to evaluate the validity and practical utility of our theoretical framework. We first examine the fundamental correlation between algebraic diversity (PLI) and accuracy. We then assess the effectiveness of INC as a sizing heuristic for both robust (OzaBagging) and complex (GOOWE) ensemble methods, discuss the implications of high-dimensional spaces, and finally refine the interpretation of our theoretical estimators.

*1) Confirmation of the PLI-Accuracy Relationship:* As a foundational check, we first confirm the relationship between PLI and accuracy. Across nearly all datasets and for both ensemble methods, a clear positive correlation is observed in the figures. As $n$ increases, the PLI (red curve) rises, empirically validating Theorem 3. As expected, the PLI is 0 for all $n < m$, as it is mathematically impossible to obtain $m$ linearly independent vectors from fewer than $m$ classifiers. This increase in PLI is mirrored by a rise in classification accuracy (blue curve).

This visual trend is quantitatively supported by the generally high Pearson correlation coefficients listed in Table V. Excluding the RBF64 dataset where PLI is universally 0, we observe a strong positive correlation ($> 0.6$) in 17 out of the 22 experimental cases. This confirms the strong statistical link between achieving linear independence and improving accuracy, motivating our central hypothesis: that the saturation of this algebraic diversity should correspond to the saturation of performance.

The few exceptions where correlation is low (e.g., Poker with GOOWE, 0.034) are notably instructive. In these cases, the low correlation is not due to a lack of diversity, but rather the behavior of the weighting mechanism. As seen in Figure 3, for Poker, the accuracy actually degrades at larger ensemble sizes even as the PLI continues to rise toward 1.0. This observation aligns with the established principle in ensemble theory that diversity is a necessary but not sufficient condition for performance. While high PLI ensures the potential to represent complex decision boundaries, it does not guarantee accuracy if the combination mechanism (in this case, GOOWE's weighting) fails to effectively aggregate that diversity.

*2) OzaBagging/Validating INC as a Practical Performance Heuristic:* Table V provides a direct validation of INC as a practical heuristic. For OzaBagging, the results are compelling. In nearly all cases, the ensemble size closest to our theoretical INC ($n_{INC}$) achieves over 99% of the maximum possible accuracy (and 100% in some cases).

- For RBF16 and RBF32, targeting the INC (34 and 66, respectively) leads to $n_{INC}$ values (32 and 64) that achieve 100.00% of the maximum accuracy.
- For Click Prediction ($m = 2$), our $n_{INC} = 4$ achieves 99.91% of the maximum.
- For Covtype ($m = 7$), our $n_{INC} = 32$ achieves 99.72% of the maximum.

This provides powerful evidence that INC is a highly effective and practical heuristic for robust, majority-voting ensembles. It reliably identifies the point of full performance saturation, confirming that achieving a high degree of theoretical linear independence is a direct proxy for achieving peak empirical performance.

*3) GOOWE/The Instability of Complex Weighting:* In contrast, the results for GOOWE reveal that achieving representational capacity is not always sufficient for optimal performance when using complex weighting schemes. While GOOWE shows high performance on some datasets (e.g., Covtype at 99.73%, RBF32 at 100%), it often fails to saturate at the theoretical diversity point identified by INC.

- On RBF2 and RBF4, targeting INC (5 and 7) yields only 87.97% and 89.01% of the maximum accuracy.
- Similarly, on Airlines and Rialto, the performance at $n_{INC}$ is roughly 95% of the maximum, with peak accuracy occurring at larger ensemble sizes (e.g., $n = 64$).

This highlights a critical distinction. While INC guarantees that the ensemble *can* algebraically represent the solution (necessary condition), GOOWE's geometric optimization may require additional factors to realize this potential. As Figure 3 shows, GOOWE's accuracy often peaks and then *degrades* at larger ensemble sizes (e.g., Poker, RBF32). In some cases (e.g., Airlines), the method benefits from redundancy, requiring ensemble sizes larger than INC to stabilize the weights. In other cases (e.g., Poker, RBF32), the method succumbs to instability, where performance degrades as the ensemble grows. Unlike OzaBagging, GOOWE's performance is dependent on complex dynamics that decouple the point of peak accuracy from the point of algebraic saturation.

*4) The Challenge of High-Dimensionality (RBF64):* The RBF64 dataset ($m = 64$) provides an edge-case. For both methods, the PLI (red curve) effectively remains 0 for all tested sizes, as spanning a 64-dimensional space is exceptionally difficult. Consequently, INC and SINC are incalculable.

However, despite this complete *lack* of full linear independence, the ensembles still achieve high practical performance (as seen in Figure 2 and 3, with OzaBagging achieving over 96% accuracy). This offers a crucial nuance to our theory: while Theorem 1 and Section III establish that $m$ independent votes are necessary for universal representational capacity (the ability to classify *any* possible instance), the RBF64 results show that this full capacity is not always required for high practical accuracy.

*5) Re-evaluating INC and SINC as Estimators:* This analysis leads to a more nuanced interpretation of our theoretical metrics, INC and SINC. They should be understood as a theoretical benchmark for full representational capacity. The key takeaway is twofold:

1) For robust methods like OzaBagging, our framework is validated: INC serves as a reliable, conservative target that guarantees full performance saturation (over 99% of max accuracy).
2) For complex methods like GOOWE, INC identifies the point of algebraic sufficiency, but not necessarily algorithmic optimality. Our results show that complex weighting schemes may require sizes larger than INC or smaller than INC to achieve their peak.

This analysis also clarifies the role of SINC. As shown in Table V, SINC consistently provides a smaller, more optimistic estimate for the point of saturation compared to INC (i.e., SINC $\leq$ INC). This is a valuable practical feature, as SINC is derived from a simplified, closed-form approximation that is computationally far more efficient than the full INC model.

We also note that due to the discrete nature of our tested ensemble sizes (e.g., 16, 32, 64), targeting the SINC value would have resulted in the same $n_{INC}$ values in our experiments. Therefore, SINC can be effectively used as a computationally cheaper lower-bound heuristic for INC.

Our framework is thus successful not just in predicting performance, but in establishing a theoretical landmark that helps explain and quantify the trade-offs between diversity, complexity, and performance saturation.

## VI. CONCLUSION AND FUTURE WORK

This paper introduced a theoretical framework that explains the trade-off between ensemble size and performance by modeling the linear independence of classifier vote vectors. We established that achieving a set of $m$ linearly independent votes is a necessary condition for an ensemble's representational capacity and developed a probabilistic model to determine the ensemble size required to meet this condition with a specified confidence level.

Our empirical results confirmed a strong correlation between the Probability of Linear Independence (PLI) and accuracy, validating that our method can effectively identify the point of diminishing returns for majority voting ensembles and reveal limitations in more complex weighting schemes. Ultimately, this work provides a principled methodology for ensemble sizing that moves beyond simple heuristics.

A key simplification in our model is the assumption that the dependency probabilities, $p_l$, are uniform across all classifiers. This "homogeneity" assumption implies that classifiers are, on average, indistinguishable from one another. This first-order approximation allowed us to establish a clean and effective theoretical framework, which, as our experiments show, successfully models the behavior of robust ensembles like OzaBagging.

For the future work, this framework could be extended to a "heterogeneous" model. Investigating individual classifier-specific dependency probabilities ($p_{i,l}$) could provide a more detailed and accurate theoretical framework, potentially explaining the performance nuances of individual classifiers or more complex, non-linear weighting schemes.

### REFERENCES

[1] H. M. Gomes, J. P. Barddal, F. Enembreck, and A. Bifet, "A survey on ensemble learning for data stream classification," *ACM Comput. Surv.*, vol. 50, no. 2, mar 2017. [Online]. Available: https://doi.org/10.1145/3054925

[2] X. Dong, Z. Yu, W. Cao, Y. Shi, and Q. Ma, "A survey on ensemble learning," *Frontiers of Computer Science*, vol. 14, 08 2019. [Online]. Available: https://doi.org/10.1007/s11704-019-8208-z

[3] B. Krawczyk, L. L. Minku, J. Gama, J. Stefanowski, and M. Woźniak, "Ensemble learning for data stream analysis: A survey," *Information Fusion*, vol. 37, pp. 132–156, 2017. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S1566253516302329

[4] L. Rokach, "Ensemble-based classifiers," *Artif. Intell. Rev.*, vol. 33, pp. 1–39, 02 2010. [Online]. Available: https://doi.org/10.1007/s10462-009-9124-7

[5] K. Ali and M. Pazzani, "Error reduction through learning multiple descriptions," *Machine Learning*, vol. 24, 11 1997. [Online]. Available: https://doi.org/10.1007/BF00058611

[6] T. G. Dietterich, "Ensemble methods in machine learning," in *International workshop on multiple classifier systems*. Springer, 2000, pp. 1–15. [Online]. Available: https://doi.org/10.1007/3-540-45014-9_1

[7] G. Tsoumakas, I. Partalas, and I. Vlahavas, "A taxonomy and short review of ensemble selection," *ECAI 2008, Workshop on Supervised and Unsupervised Ensemble Methods and Their Applications*, 01 2008.

[8] R. Polikar, "Ensemble based systems in decision making," *IEEE Circuits and Systems Magazine*, vol. 6, no. 3, pp. 21–45, 2006. [Online]. Available: https://doi.org/10.1109/MCAS.2006.1688199

[9] D. Opitz and R. Maclin, "Popular ensemble methods: an empirical study," *J. Artif. Int. Res.*, vol. 11, no. 1, p. 169–198, Jul. 1999. [Online]. Available: https://doi.org/10.1613/jair.614

[10] L. Breiman, "Bagging predictors," *Mach. Learn.*, vol. 24, no. 2, p. 123–140, Aug. 1996. [Online]. Available: https://doi.org/10.1023/A:1018054314350

[11] R. Schapire, "The strength of weak learnability," *Machine Learning*, vol. 5, no. 2, pp. 197–227, Jun. 1990. [Online]. Available: https://doi.org/10.1023/A:1022648800760

[12] Y. Freund and R. E. Schapire, "A decision-theoretic generalization of on-line learning and an application to boosting," *Journal of Computer and System Sciences*, vol. 55, no. 1, pp. 119–139, 1997. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S002200009791504X

[13] L. Breiman, "Random forests," *Machine Learning*, vol. 45, pp. 5–32, 10 2001. [Online]. Available: https://doi.org/10.1023/A:1010950718922

[14] L. Lam and S. Suen, "Application of majority voting to pattern recognition: An analysis of its behavior and performance," *IEEE Transactions on Systems, Man, and Cybernetics - Part A: Systems and Humans*, vol. 27, no. 5, pp. 553–568, 1997. [Online]. Available: https://doi.org/10.1109/3468.618255

[15] L. Kuncheva, C. Whitaker, C. Shipp, and R. Duin, "Limits on the majority vote accuracy in classier fusion," *Formal Pattern Analysis & Applications*, vol. 6, pp. 22–31, 04 2003. [Online]. Available: https://doi.org/10.1007/s10044-002-0173-7

[16] A. Narasimhamurthy, "Theoretical bounds of majority voting performance for a binary classification problem," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 27, no. 12, pp. 1988–1995, 2005. [Online]. Available: https://doi.org/10.1109/TPAMI.2005.249

[17] D. Ruta and B. Gabrys, "A theoretical analysis of the limits of majority voting errors for multiple classifier systems," *Pattern Anal. Appl.*, vol. 5, pp. 333–350, 10 2002. [Online]. Available: https://doi.org/10.1007/s100440200030

[18] H. Liu, A. Mandvikar, and J. Mody, "An empirical study of building compact ensembles," in *Advances in Web-Age Information Management: 5th International Conference, WAIM 2004, Dalian, China, July 15-17, 2004 5*. Springer, 2004, pp. 622–627. [Online]. Available: https://doi.org/10.1007/978-3-540-27772-9_63

[19] A. Lazarevic and Z. Obradovic, "Effective pruning of neural network classifier ensembles," in *IJCNN'01. International Joint Conference on Neural Networks. Proceedings (Cat. No.01CH37222)*, vol. 2, 2001, pp. 796–801 vol.2. [Online]. Available: https://doi.org/10.1109/IJCNN.2001.93946

[20] D. D. Margineantu and T. G. Dietterich, "Pruning adaptive boosting," in *ICML*, vol. 97. Citeseer, 1997, pp. 211–218. [Online]. Available: https://dl.acm.org/doi/10.5555/645526.757762

[21] W. Fan, F. Chu, H. Wang, and P. S. Yu, "Pruning and dynamic scheduling of cost-sensitive ensembles," in *AAAI/IAAI*, 2002, pp. 146–151.

[22] Z.-H. Zhou, J. Wu, and W. Tang, "Ensembling neural networks: Many could be better than all," *Artificial Intelligence*, vol. 137, no. 1, pp. 239–263, 2002. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S000437020200190X

[23] Z.-H. Zhou and W. Tang, "Selective ensemble of decision trees," in *Rough Sets, Fuzzy Sets, Data Mining, and Granular Computing: 9th International Conference, RSFDGrC 2003, Chongqing, China, May 26–29, 2003 Proceedings 9*. Springer, 2003, pp. 476–483. [Online]. Available: https://doi.org/10.1007/3-540-39205-X_81

[24] D. Hernández-Lobato, G. Martínez-Muñoz, and A. Suárez, "How large should ensembles of classifiers be?" *Pattern Recognition*, vol. 46, no. 5, pp. 1323–1336, 2013. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0031320312004554

[25] A. Ulaş, M. Semerci, O. T. Yıldız, and E. Alpaydın, "Incremental construction of classifier and discriminant ensembles," *Information Sciences*, vol. 179, no. 9, pp. 1298–1318, 2009. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0020025509000061

[26] J. Xiao, C. He, X. Jiang, and D. Liu, "A dynamic classifier ensemble selection approach for noise data," *Information Sciences*, vol. 180, no. 18, pp. 3402–3421, 2010. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0020025510002240

[27] L. Yang, "Classifiers selection for ensemble learning based on accuracy and diversity," *Procedia Engineering*, vol. 15, pp. 4266–4270, 12 2011. [Online]. Available: https://doi.org/10.1016/j.proeng.2011.08.800

[28] T. Windeatt and C. Zor, "Ensemble pruning using spectral coefficients," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 24, no. 4, pp. 673–678, 2013. [Online]. Available: https://doi.org/10.1109/TNNLS.2013.2239659

[29] H. Ykhlef and D. Bouchaffra, "An efficient ensemble pruning approach based on simple coalitional games," *Information Fusion*, vol. 34, pp. 28–42, 2017. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S1566253516300562

[30] X. Zhu, Z. Ni, L. Ni, F. Jin, M. Cheng, and J. Li, "Improved discrete artificial fish swarm algorithm combined with margin distance minimization for ensemble pruning," *Computers & Industrial Engineering*, vol. 128, 12 2018. [Online]. Available: https://doi.org/10.1016/j.cie.2018.12.021

[31] M. N. Adnan and M. Z. Islam, "Optimizing the number of trees in a decision forest to discover a subforest with high ensemble accuracy using a genetic algorithm," *Knowledge-Based Systems*, vol. 110, pp. 86–97, 2016. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0950705116302301

[32] L. Rokach, "Collective-agreement-based pruning of ensembles," *Computational Statistics & Data Analysis*, vol. 53, pp. 1015–1026, 02 2009. [Online]. Available: https://doi.org/10.1016/j.csda.2008.12.001

[33] S. Abadifard, S. Bakhshi, S. Gheibuni, and F. Can, "Dyned: Dynamic ensemble diversification in data stream classification," in *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*, ser. CIKM '23. New York, NY, USA: Association for Computing Machinery, 2023, p. 3707–3711. [Online]. Available: https://doi.org/10.1145/3583780.3615266

[34] S. Elbasi, A. Büyükçakir, H. R. Bonab, and F. Can, "On-the-fly ensemble pruning in evolving data streams," *CoRR*, vol. abs/2109.07611, 2021. [Online]. Available: https://arxiv.org/abs/2109.07611

[35] M. Bhardwaj, V. Bhatnagar, and K. Sharma, "Cost-effectiveness of classification ensembles," *Pattern Recognition*, vol. 57, pp. 84–96, 2016. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0031320316300024

[36] P. Latinne, O. Debeir, and C. Decaestecker, "Limiting the number of trees in random forests," vol. 2096, 07 2001, pp. 178–187. [Online]. Available: https://doi.org/10.1007/3-540-48219-9_18

[37] T. Oshiro, P. Perez, and J. Baranauskas, "How many trees in a random forest?" vol. 7376, 07 2012. [Online]. Available: https://doi.org/10.1007/978-3-642-31537-4_13

[38] P. Probst and A.-L. Boulesteix, "To tune or not to tune the number of trees in random forest?" *Journal of Machine Learning Research*, vol. 18, 05 2017. [Online]. Available: https://doi.org/10.48550/arXiv.1705.05654

[39] E. Bax, "Selecting a number of voters for a voting ensemble," *CoRR*, vol. abs/2104.11833, 2021. [Online]. Available: https://arxiv.org/abs/2104.11833

[40] G. Fumera and F. Roli, "A theoretical and experimental analysis of linear combiners for multiple classifier systems," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 27, pp. 942–56, 07 2005. [Online]. Available: https://doi.org/10.1109/TPAMI.2005.109

[41] G. Fumera, R. Fabio, and S. Alessandra, "A theoretical analysis of bagging as a linear combination of classifiers," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 30, pp. 1293–9, 08 2008. [Online]. Available: https://doi.org/10.1109/TPAMI.2008.30

[42] K. Jackowski, "New diversity measure for data stream classification ensembles," *Eng. Appl. Artif. Intell.*, vol. 74, no. C, p. 23–34, Sep. 2018. [Online]. Available: https://doi.org/10.1016/j.engappai.2018.05.006

[43] S. Wu and F. Crestani, "A geometric framework for data fusion in information retrieval," *Information Systems*, vol. 50, pp. 20–35, 2015. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0306437915000113

[44] H. R. Bonab and F. Can, "A theoretical framework on the ideal number of classifiers for online ensembles in data streams," in *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*, ser. CIKM '16. New York, NY, USA: Association for Computing Machinery, 2016, p. 2053–2056. [Online]. Available: https://doi.org/10.1145/2983323.2983907

[45] H. Bonab and F. Can, "Less is more: A comprehensive framework for the number of components of ensemble classifiers," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 30, no. 9, pp. 2735–2745, 2019. [Online]. Available: https://doi.org/10.1109/TNNLS.2018.2886341

[46] S. Wu and W. Ding, "A dataset-level geometric framework for ensemble classifiers," *CoRR*, vol. abs/2106.08658, 2021. [Online]. Available: https://arxiv.org/abs/2106.08658

[47] R. E. Banfield, L. O. Hall, K. W. Bowyer, and W. Kegelmeyer, "Ensemble diversity measures and their application to thinning," *Information Fusion*, vol. 6, no. 1, pp. 49–62, 2005, diversity in Multiple Classifier Systems. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S1566253504000387

[48] N. Oza, "Online bagging and boosting," in *2005 IEEE International Conference on Systems, Man and Cybernetics*, vol. 3, 2005, pp. 2340–2345 Vol. 3. [Online]. Available: https://doi.org/10.1109/ICSMC.2005.1571498

[49] H. R. Bonab and F. Can, "Goowe: Geometrically optimum and online-weighted ensemble classifier for evolving data streams," *ACM Trans. Knowl. Discov. Data*, vol. 12, no. 2, Jan 2018. [Online]. Available: https://doi.org/10.1145/3139240

[50] P. Domingos and G. Hulten, "Mining high-speed data streams," in *Proceedings of the Sixth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ser. KDD '00. New York, NY, USA: Association for Computing Machinery, 2000, p. 71–80. [Online]. Available: https://doi.org/10.1145/347090.347107

[51] J. Montiel, J. Read, A. Bifet, and T. Abdessalem, "Scikit-multiflow: A multi-output streaming framework," *Journal of Machine Learning Research*, vol. 19, no. 72, pp. 1–5, 2018. [Online]. Available: http://jmlr.org/papers/v19/18-251.html