

Ensemble deep learning: A review

M.A. Ganaie^a, Minghui Hu^b, A.K. Malik^a, M. Tanveer^{a,*}, P.N. Suganthan^{b,c,*}

^aDepartment of Mathematics, Indian Institute of Technology Indore, Simrol, Indore, 453552, India

^bSchool of Electrical & Electronic Engineering, Nanyang Technological University, Singapore

^cKINDI Center for Computing Research College of Engineering, Qatar University, Qatar

Abstract

Ensemble learning combines several individual models to obtain better generalization performance. Currently, deep learning architectures are showing better performance compared to the shallow or traditional models. Deep ensemble learning models combine the advantages of both the deep learning models as well as the ensemble learning such that the final model has better generalization performance. This paper reviews the state-of-art deep ensemble models and hence serves as an extensive summary for the researchers. The ensemble models are broadly categorised into bagging, boosting, stacking, negative correlation based deep ensemble models, explicit/implicit ensembles, homogeneous/heterogeneous ensemble, decision fusion strategies based deep ensemble models. Applications of deep ensemble models in different domains are also briefly discussed. Finally, we conclude this paper with some potential future research directions.

Keywords: Ensemble Learning, Deep Learning.

1. Introduction

Deep learning architectures have been successfully employed across a wide range of applications from image/video classification to the health care. The success of these models is attributed to the better feature representation via multi layer processing architectures. The deep learning models have been mainly used for classification, regression and clustering problems. Classification problem is defined as the categorization of the new observations based on the hypothesis

*Corresponding authors

Email addresses: phd1901141006@iiti.ac.in (M.A. Ganaie), minghui.hu@ntu.edu.sg (Minghui Hu), phd1801241003@iiti.ac.in (A.K. Malik), mtanveer@iiti.ac.in (M. Tanveer), epnsugan@ntu.edu.sg (P.N. Suganthan)

Preprint submitted to Elsevier

August 9, 2022

h learned from the set of training data. The hypothesis h represents a mapping of input data features to the appropriate target labels/classes. The main objective, while learning the hypothesis h , is that it should approximate the true unknown function as close as possible to reduce the generalization error. There exist several applications of these classification algorithms ranging from medical diagnosis to remote sensing. Mathematically,

$$O_c = h(x, \theta_c), O_c \in \mathbb{Z}, \quad (1)$$

where x is the input feature vector, O_c is the category of the sample x , θ_c is the set of learning parameters of the hypothesis h and \mathbb{Z} is the set of class labels.

Regression problems deal with the continuous decisions, instead of discrete categories. Mathematically,

$$O_r = h(x, \theta_r), O_r \in \mathbb{R}, \quad (2)$$

where x is the observation vector, O_r is the output, and θ_r is the set of learning parameters of the hypothesis h .

Broadly speaking, there are different approaches of classification like supervised, unsupervised classification, few-shot, one-shot and so on. Here, we only discuss supervised and unsupervised classification problems. In supervised learning, the building of hypothesis h is supervised based on the known output labels provided in the training data samples, while as in unsupervised learning hypothesis h is generated without any supervision as no known output values are available with the training data. This approach, also known as clustering, generates the hypothesis h based on the similarities and dissimilarities present in the training data.

Generally speaking, the goal of generating the hypothesis h in Machine learning area is that it should perform better when applied to unknown data. The performance of the model is measured with respect to the area in which the model is applied. Combining the predictions from several models has proven to be an elegant approach for increasing the performance of the models. Combination of several different predictions from different models to make the final prediction is known as ensemble learning or ensemble model. The ensemble learning involves multiple models combined in some fashion like averaging, voting such that the ensemble model is better than any of the individual models. To prove that average voting in an ensemble is better than individual model, Marquis de Condorcet proposed a theorem wherein he proved that if the probability of each voter being correct is above 0.5 and the voters are independent, then addition of more

voters increases the probability of majority vote being correct until it approaches 1 [1]. Although Marquis de Condorcet proposed this theorem in the field of political science and had no idea of the field of Machine learning, but it is the similar mechanism that leads to better performance of the ensemble models. Assumptions of Marquis de Condorcet theorem also holds true for ensembles [2]. The reasons for the success of ensemble learning include: statistical, computational and representation learning [3], bias-variance decomposition [4] and strength-correlation [5].

In this era of machine learning, deep learning automates the extraction of high-level features via hierarchical feature learning mechanism wherein the upper layer of features are generated on the previous set of layer/layers. Deep learning has been successfully applied across different fields since the ImageNet Large Scale Recognition Challenge (ILSVRC) competitions [6, 7] and has achieved state-of-art performance. It has obtained promising results in object detection, semantic segmentation, edge detection and number of other domains. However, given the computational cost, the training of deep ensemble models is an uphill task. Different views have been provided to understand how the deep learning models learn the features like learning through hierarchy of concepts via many levels of representation [8, 9, 10]. Given the advantages of deep learning models from deep architectures, there are several bottlenecks like vanishing/exploding gradients [11, 12] and degradation problem [13] which prevent to reach this goal. Recently, training deep network's has become feasible through the Highway networks [14] and Residual networks [13]. Both these networks enabled to train very deep networks. The ensemble learning has been recently known to be strong reason for enhancing the performance of deep learning models [15]. Thus, the objective of deep ensemble models is to obtain a model that has best of both the ensemble and deep models.

There exist multiple surveys in the literature which mainly focus on the review of ensemble learning like learning of ensemble models in classification problems [16, 17, 18, 19], regression problems [20, 21] and clustering [22]. Review of both the classification and regression models was given in [23]. Comprehensive review of the ensemble methods and the challenges were given in [24]. Though [24] provided some insight about the deep ensemble models but couldn't give the comprehensive review of the deep ensemble learning while as [25] reviewed the ensemble deep models in the context of bioinformatics. The past decade has successively evolved different deep learning strategies which have lead to the exploration and innovation of these models in multiple areas like health care, speech, image classification, forecasting and other

applications. Broadly speaking, ensemble learning approaches have followed classical methods, general methods and different fusion strategies for improving the performance of the models. Since deep learning models are computation and data extensive, hence, ensemble deep learning models need special attention while exploring the complementary information of multiple algorithms into a uniform framework. Ensemble deep learning models need to handle multiple questions like how to induce diversity among the baseline models, how to keep the training time as well the models complexity lower for the practical applications, how to fuse the predictions of the complementary algorithms. Multiple studies have handled these problems differently. In this review paper, we comprehensively review the different approaches used to handle the aforementioned problems. In this paper, we give a comprehensive review of deep ensemble models. **To the best of our knowledge, this is the first comprehensive review paper on deep ensemble models.**

The rest of this paper is organised as follows: Section-3 discusses the theoretical aspects of deep ensemble learning, Section-4 discusses the different approaches used in deep ensemble strategies, applications of deep ensemble methods are given in Section-5 and finally conclusions and future directions are given in Section-6.

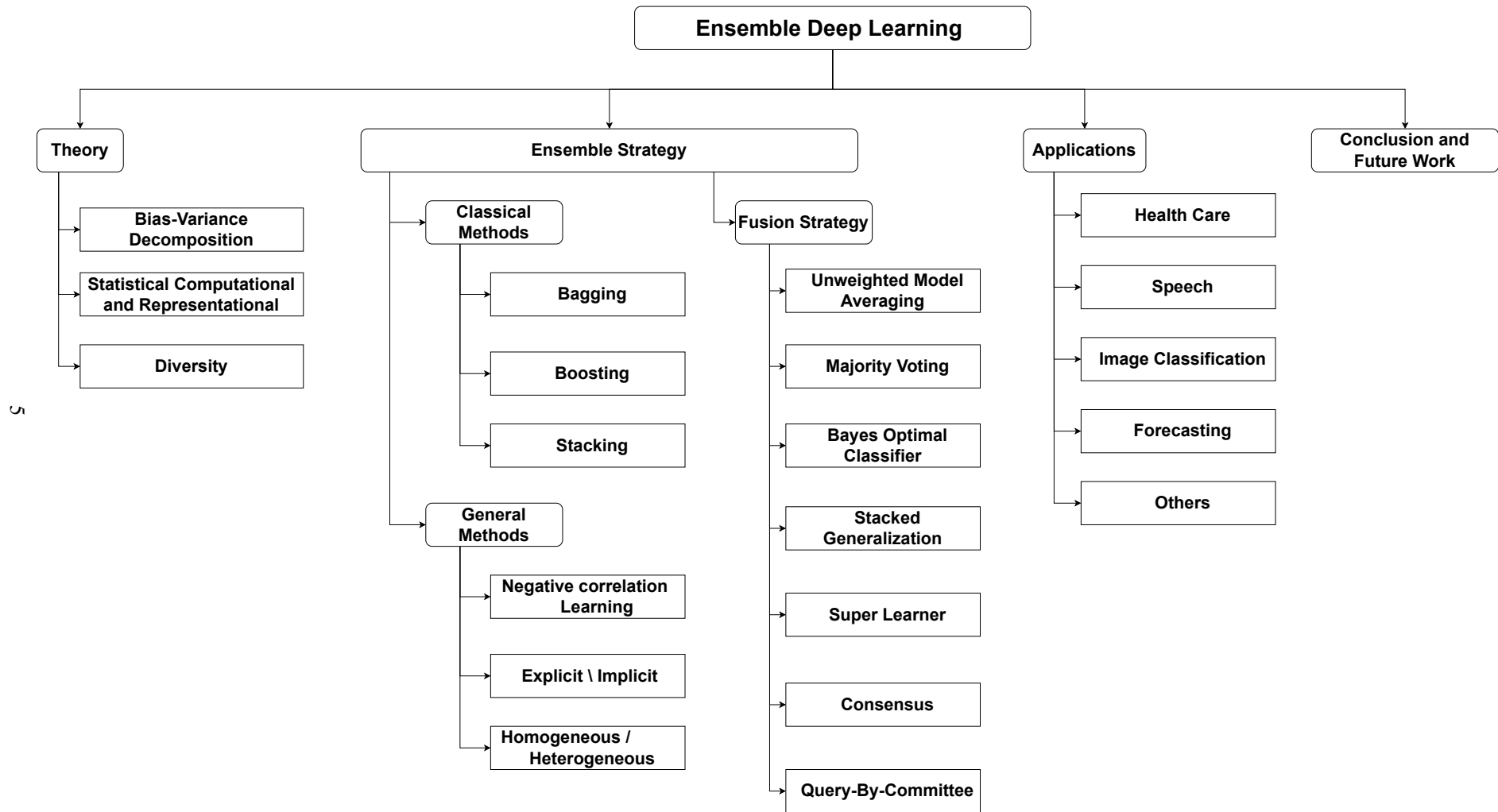


Figure 1: Layout of the paper.

2. Research Methodology

The studies in this review are searched from the Google Scholar and Scopus search engines. The papers are the result of ensemble learning, ensemble deep learning, deep ensemble learning, deep ensembles keywords. The articles were screened based on the title and abstract, followed by the screening of full-text version. The articles are elaborated based on the ensemble learning and deep learning approaches.

3. Theory

The various reasons which have been justified for the success of ensemble learning can be discussed under the following subheadings:

3.1. Bias-Variance Decomposition

Initially, the success of ensemble methods was theoretically investigated for regression problems. Krogh and Vedelsby [26], Brown et al. [27] proved via ambiguity decomposition that the proper ensemble classifier guarantees a smaller squared error as compared to the individual predictors of the classifier. Ambiguity decomposition was given for single dataset based ensemble methods, later on, multiple dataset bias-variance-covariance decomposition was introduced in [27, 28, 29, 30] and is given as:

$$\begin{aligned} E[o - t]^2 &= bias^2 + \frac{1}{M} var + \left(1 - \frac{1}{M}\right) covar, \\ bias &= \frac{1}{M} \sum_i (E[o_i] - t), \\ var &= \frac{1}{M} \sum_i E[o_i - E[o_i]]^2, \\ covar &= \frac{1}{M(M-1)} \sum_i \sum_{j \neq i} E[o_i - E[o_i]][o_j - E[o_j]], \end{aligned} \quad (3)$$

where t is target, o_i is the output of i^{th} model and M is the ensemble size. Here, $bias$ term measures the average difference between the base learner and the model output, var indicates their average variance, and $covar$ is the covariance term measuring the pairwise difference of the base learners.

Ensemble methods have been supported by several theories like bias-variance [4, 31], strength correlation [5], stochastic discrimination [32], and margin theory [33]. These theories provide the equivalent of bias-variance-covariance decomposition [34].

The above given equations of decomposition error can't be directly applied to the datasets with discrete class labels due to their categorical nature. However, alternate ways to decompose the error in classification problems are given in [4, 35, 36, 37, 38].

Multiple approaches like bagging, boosting have been proposed for generating the ensemble methods. Bagging reduces the variance among the base classifiers [39] while as boosting based ensembles lead to the bias and variance reduction [40, 41].

3.2. Statistical, Computational and Representational Aspects

Dietterich provided Statistical, Computational and Representational reasons [3] for success of ensemble models. The learning model is viewed as the search of the optimal hypothesis h among the several hypothesis in the search space. When the amount of data available for the training is smaller compared to the size of the hypothesis space, the statistical problem arises. Due to this statistical problem, the learning algorithm identifies the different hypothesis which gives same performance on the training samples. Ensembling of these hypothesis results in an algorithm which reduces the risk of being a wrong classifier. The second reason is computational wherein a learning algorithm sticks in a local optima due to some form of local search. Ensemble model overcomes this issue by performing some form of local search via different starting points which leads to better approximation of the true unknown function. Another reason is representational wherein none of the hypotheses among the set of hypothesis is able to represent the true unknown function. Hence, ensembling of these hypothesis via some weighting technique results into the hypothesis which expands the representable function space.

3.3. Diversity

One of the main reasons behind the success of ensemble methods is increasing the diversity among the base classifiers and the same thing was highlighted in [3]. Different approaches have been followed to generate diverse classifiers. Different methods like bootstrap aggregation (bagging) [39], Adaptive Boosting (AdaBoost) [42], random subspace [43], and random forest [5] approaches are followed for generating the multiple datasets from the original dataset to train the different predictors such that the outputs of predictors are diverse. Attempts have been made

to increase diversity in the output data wherein multiple outputs are created instead of multiple datasets for the supervision of the base learners. ‘Output smearing’ [44] is one of this kind which induces random noise to introduce diversity in the output space.

4. Ensemble Strategies:

The different ensemble strategies have evolved over a period of time which results in better generalization of the learning models. The ensemble strategies are broadly categorised as follows:

4.1. *Bagging*

Bagging [39], also known as bootstrap aggregating, is one of the standard techniques for generating the ensemble-based algorithms. Bagging is applied to enhance the performance of an ensemble classifier. The main idea in bagging is to generate a series of independent observations with the same size, and distribution as that of the original data. Given the series of observations, generate an ensemble predictor which is better than the single predictor generated on the original data. Bagging increases two steps in the original models: First, generating the bagging samples and passing each bag of samples to the base models and second, strategy for combining the predictions of the multiple predictors. Bagging samples may be generated with or without replacement. Combining the output of the base predictors may vary as mostly majority voting is used for classification problems while the averaging strategy is used in regression problems for generating the ensemble output. Figure 2 shows the diagram of the bagging technique. Here, D_i represents the bagged datasets, C_i represents the algorithms and F_{ens} calculates the final outcome.

Random Forest [5] is an improved version of the decision trees that uses the bagging strategy for improving the predictions of the base classifier which is a decision tree. The fundamental difference between these two methods is that at each tree split in Random Forest, only a subset of features is randomly selected and considered for splitting. The purpose of this method is to decorrelate the trees and prevent over-fitting. Breiman [5] showed heuristically that the variance of the bagged predictor is smaller than the original predictor and proposed that bagging is better in higher dimensional data. However, the analysis of the smoothing effect of bagging [45] revealed that bagging doesn’t depend on the data dimensionality.

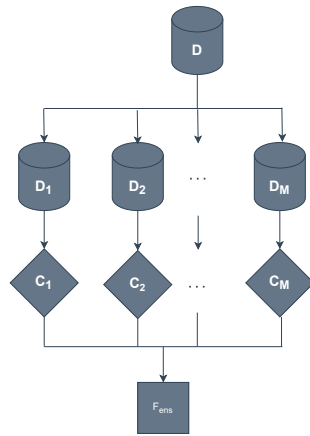


Figure 2: Bagging

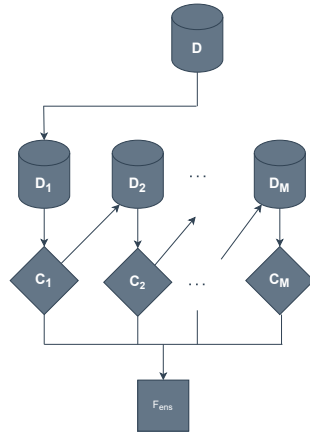


Figure 3: Boosting

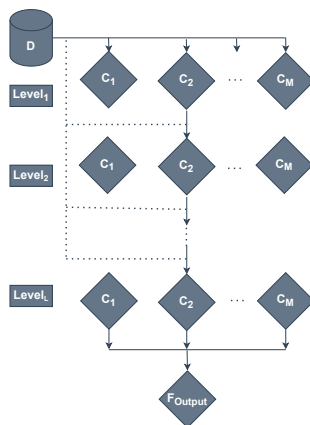


Figure 4: Stacking

Bühlmann and Yu [46] gave theoretical explanation of how bagging gives smooth hard decisions, small variance, and mean squared error. Since bagging is computationally expensive, hence subbagging and half subbagging [46] were introduced. Half subbagging, being computationally efficient, is as accurate as the bagging.

Several attempts tried to combine bagging with other machine learning algorithms. Kim et al. [47] used bagging method to generate multiple bags of the dataset and multiple support vector machines were trained independently with each bag as the input. The output of the models is combined via majority voting, least squares estimation weighting and double layer hierarchical approach. In the double layer hierarchical approach, another support vector machines (SVM) is used to combine the outcomes of the multiple SVM's efficiently. Tao et al. [48] used asymmetric bagging strategy to generate the ensemble model to handle the class imbalance problems. A case study of bagging, boosting and basic ensembles [49] revealed that at higher rejection rates of samples boosting is better as compared to bagging and basic ensembles. However, as the rejection rate increases the difference disappears among the boosting, bagging and basic ensembles. Bagging based multilayer perceptron [50] combined bagging to train multiple perceptrons with the corresponding bag and showed that bagging based ensemble models perform better as compared to individual multilayer perceptron. In [51], the analysis of the bagging approach and other regularisation techniques revealed that bagging regularized the neural networks and hence provide better generalization. In [52], bagged neural networks (BNNs) was proposed wherein each neural network was trained over different dataset sampled randomly with replacement from original dataset and was implemented for the short term load forecasting. Unlike Random forest [5] which uses majority voting for aggregating the ensemble of decision trees, bagging based survival trees [53] used Kaplan–Meier curve to predict the ensemble output for breast cancer and lymphoma patients. In [54], ensembles of stacked denoising autoencoders for classification showed that the bagging and switching technique in a general deep machine results in improved diversity.

Bagging has also been applied to solve the problem of imbalanced data. Roughly Balanced Bagging [55] tries to equalize each class's sampling probability in binary class problems wherein the negative class samples are sampled via negative binomial distribution, instead of keeping the sample size of each class the same number. Neighbourhood Balanced Bagging [56] incorporated the neighbourhood information for generating the bagging samples for the class imbalance

Years	Authors	Contribution
1996	Breiman [39]	Proposed the idea of Bagging
1998	Mao [49]	Case study of bagging, boosting and basic ensembles
2000	Buja and Stuetzle [45]	Theoretical analysis of bagging
2001	Breiman [5]	Bagging with random subspace Decision trees and ensembling outputs via majority voting
2001	Gençay and Qi [51]	Study of Bayesian regularization, early stopping and Bagging
2002	Kim et al. [47]	Bagging with SVM's and ensembling outputs via SVM's, majority voting and least squares estimation
2002	Bühlmann and Yu [46]	Theoretical justification of Bagging, proposed subbagging and half subbagging
2004	Hothorn et al. [53]	Bagging with decision trees and ensembling outputs via Kaplan–Meier curve
2005	Oza [57]	Theoretical and experimental analysis of online bagging and boosting
2006	Tao et al. [48]	Proposed asymmetric bagging with SVM's and ensembling outputs SVM's
2009	Hido et al. [55]	Roughly balanced bagging on decision trees and ensembling outputs via majority voting
2005, 2015	Ha et al. [50], Khwaja et al. [52]	Bagging with Neural networks and ensembling outputs via majority voting
2015	Błaszczczyński and Stefanowski [56]	Neighbourhood balanced bagging ensembling outputs via majority voting

Table 1: Bagging based ensemble models

problems. Błaszczczyński and Stefanowski [56] concluded that applying conventional diversification is more effective when applied at the last classification methods. Both roughly balanced Bagging and Neighbourhood Balanced Bagging have not been explored in deep learning architectures. Thus, these approaches can be exploited to handle the class imbalance problems via deep ensemble models.

The theoretical and experimental analysis of online bagging and boosting [57] showed that the online bagging algorithm can achieve similar accuracy as the batch bagging algorithm with only a little more training time. However, online bagging is an option when all training samples can't be loaded into the memory due to memory issues.

Although ensembling may lead to increase in the computational complexity, but bagging possesses the property that it can be paralleled and can lead to effective reduction in the training time subject to the availability of hardware for running the parallel models. Since deep learning models have high training time, hence optimization of multiple deep models on different training bags is not a feasible option.

4.2. *Boosting*

Boosting technique is used in ensemble models for converting a weak learning model into a learning model with better generalization. Figure 3 shows the diagram of the boosting technique. The techniques such as majority voting in case of classification problems or a linear combination of weak learners in the regression problems results in better prediction as compared to the

single weak learner. Boosting methods like AdaBoost [42] and Gradient Boosting [58] have been used across different domains. Adaboost uses a greedy technique for minimizing a convex surrogate function upper bounded by misclassification loss via augmentation, at each iteration, the current model with the appropriately weighted predictor. AdaBoost learns an effective ensemble classifier as it leverages the incorrectly classified sample at each stage of the learning. AdaBoost minimizes the exponential loss function while as the Gradient boosting generalized this framework to the arbitrary differential loss function.

Boosting, also known as forward stagewise additive modelling, was originally proposed to improve the performance of the classification trees. It has been recently incorporated in the deep learning models to further improve their performance.

Boosted deep belief network (DBN) [59] for facial expression recognition unified the boosting technique and multiple DBN's via objective function which results in a strong classifier. The model learns complex feature representation to build a strong classifier in an iterative manner. Deep boosting [60] is an ensemble model that uses the deep decision trees. It can also be used in combination with any other rich family classifier to improve the generalization performance. In each stage of the deep boosting, the decisions of which classifier to add and what weights should be chosen depends on the (data-dependent) complexity of the classifier to which it belongs. The interpretation of the deep boosting classifier is given via structural risk minimization principle at each stage of the learning. Multiclass Deep boosting [61] extended the Deep boosting [60] algorithm to theoretical, algorithmic, and empirical results to the multiclass problems. Due to the limitation of the training data in each mini batch, Boosting CNN may overfit the data. To avoid overfitting, incremental Boosting CNN (IBCNN) [62] accumulated the information of multiple batches of the training data samples. IBCNN uses decision stumps on the top of single neurons as the weak learners and learns weights via AdaBoost method in each mini batch. Unlike DBN [59] which uses image patch for learning the weak classifiers, IBCNN trains the weak classifiers from the fully connected layer i.e. the whole image is used for learning the weak classifiers. To make the IBCNN model more efficient, the weak learners loss functions are combined with the global loss function.

Boosted CNN [63] used boosting for training the deep CNN. Instead of averaging, least squares objective function was used to incorporate the boosting weights into CNN. Moghimi et al. [63] also showed that CNN can be replaced by network structure within their boosting

framework for improving the performance of the base classifier. Boosting increases the complexity of training the networks, hence the concept of dense connections was introduced in a deep boosting framework to overcome the problem of vanishing gradient problem for image denoising [64]. Deep boosting framework was extended to image restoration in [65] wherein the dilated dense fusion network was used to boost the performance.

The convolutional channel features [66] generated the high level features via CNN and then used boosted forest for final classification. Since CNN has high number of hyperparameters than the boosted forest, hence the model proved to be efficient than end-to-end training of CNN models both in terms of performance and time. Yang et al. [66] showed its application in edge detection, object proposal generation, pedestrian and face detection. A stagewise boosting deep CNN [67] trains several models of the CNNs within the offline paradigm boosting framework. To extend the concept of boosting in online scenario's wherein only a chunk of data is available at given time, Boosting Independent Embeddings Robustly (BIER) [68] was proposed to cope up the online scenario's. In BIER, a single CNN model is trained end-to-end with an online boosting technique. The training set in the BIER is reweighed via the negative gradient of the loss function to project the input spaces (images) into a collection of independent output spaces. To make BIER more robust, Hierarchical Boosted deep metric learning [69] incorporated the hierarchical label information into the embedding ensemble which improves the performance of the model on the large scale image retrieval application. Using deep boosting results in higher training time, to reduce the warm-up phase of training which trains the classifier from scratch deep incremental boosting [70] used transfer learning approach. This approach leveraged the initial warm-up phase of each incremental base model of the ensemble during the training of the network. To reduce the training time of boosting based ensembles, snapshot boosting [71] combined the merits of snapshot ensembling and boosting to improve the generalization without increasing the cost of training. Snapshot boosting trains each base network and combines the outputs via meta learner to combine the output of base learners more efficiently.

Literature shows that the boosting concept is the backbone behind well-known architectures like Deep Residual networks [13, 72], AdaNet [73]. The theoretical background for the success of the Deep Residual networks (DeepResNet) [13] was explained in the context of boosting theory [74]. The authors proposed multi-channel telescoping sum boosting learning framework, known as BoostResNet, wherein each channel is a scalar value updated during rounds of boosting

Years	Authors	Contribution
2014	Liu et al. [59]	Boosted deep belief network (DBN) as base classifiers for facial expression recognition.
2014	Cortes et al. [60]	Decision trees as base classifiers for binary class classification problems.
2014	Kuznetsov et al. [61]	Decision trees as base classifiers for multiclass classification problems.
2015	Yang et al. [66]	Ensemble of CNN and boosted forest for edge detection, object proposal generation, pedestrian and face detection.
2016	Moghimi et al. [63]	Boosted CNN
2016	Walach and Wolf [67]	CNN Boosting applied to bacteria cell images and crowd counting.
2017	Opitz et al. [68]	Boosted deep independent embedding model for online scenarios.
2017	Mosca and Magoulas [70]	Transfer learning based deep incremental boosting.
2017	Han et al. [62]	Boosting based CNN with incremental approach for facial action unit recognition.
2018	Chen et al. [64]	Deep boosting for image denoising with dense connections.
2019	Chen et al. [65]	Deep boosting for image restoration and image denoising.
2019	Waltner et al. [69]	Hierarchical boosted deep metric learning with hierarchical label embedding.
2020	Zhang et al. [71]	Snapshot boosting.

Table 2: Boosting based ensemble models

to minimize the multi-class error rate. The fundamental difference between the AdaNet and BoostResnet is that the former maps the feature vectors to classifier space and boosts weak classifiers while the latter used multi-channel representation boosting. Moreover, BoostResNet is more efficient than DeepResnet in terms of computational time.

The theory of boosting was extended to online boosting in [75] and provided theoretical convergence guarantees. Online boosting shows improved convergence guarantees for batch boosting algorithms.

The ensembles of bagging and boosting have been evaluated in [76]. The study evaluated the different algorithms based on the concept of bagging and boosting along with the availability of software tools. The study highlighted the practical issues and opportunities of their feasibility in ensemble modeling.

4.3. Stacking

Ensembling can be done either by combining outputs of multiple base models in some fashion or using some method to choose the “best” base model. Figure 4 shows the stacking technique. Stacking is one of the integration techniques wherein the meta-learning model is used to integrate the output of base models. If the final decision part is a linear model, the staking is often referred to as “model blending” or simply “blending”. The concept of stacking or stacked regression was initially given by [77]. In this technique, the dataset is randomly split into J equal parts. For the j^{th} -fold cross-validation one set is used for testing and the rest are used for training. With these

training testing pair subsets, we obtain the predictions of different learning models which are used as the meta-data to build the meta-model. Meta-model makes the final prediction, which is also called the winner-takes-all strategy.

Stacking is a bias reducing technique [78]. Following [77], Deep convex net (DCN) [79] was proposed which is a deep learning architecture composed of a variable number of modules stacked together to form the deep architecture. Each learning module in DCN is convex. DCN is a stack of several modules consisting of linear input units, hidden layer non-linear units, and the second linear layer with the number of units as that of target classification classes. The modules are connected layerwise as the output of the lower module is given as input to the adjacent higher module in addition to the original input data. The deep stacking network (DSN) enabling parallel training on very large scale datasets was proposed in [80], the network was named stacking based as it shared the concept of “stacked generalization” [77]. The kernelized version of DCN, known as kernel deep convex networks (K-DCN), was given in [81], here the number of hidden layer approach infinity via kernel trick. Deng et al. [81] showed that K-DCN performs better as compared to the DCN. However, due to kernel trick the memory requirements increase and hence may not be scalable to large scale datasets. Also, we need to optimize the hyperparameters like the number of levels in the stacked network, the kernel parameters to get the optimal performance of the network. To leverage the memory requirements, random Fourier feature-based kernel deep convex network [82] approximated the Gaussian kernel which reduces the training time and helps in the evaluation of K-DCN over large scale datasets. A framework for parameter estimation and model selection in kernel deep stacking networks [83] is based on the combination of model-based optimization and hill-climbing approaches. Welchowski and Schmid [83] used data-driven framework for parameter estimation, hyperparameter tuning and model selection in kernel deep stacking networks. Another improvement over DSN was Tensor Deep Stacking Network (T-DSN) [84], here in each block of the stacked network, large single hidden layer was split into two smaller ones and then mapped bilinearly to capture the higher-order interactions among the features. Comprehensive evaluation, the more detailed analysis of the learning algorithm and T-DSN implementation is given in [85]. Sparse coding is another popular method that is used in the deep learning area. The advantage of sparse representation is numerous, including robust to noise, effective for learning useful features, etc. Sparse Deep Stacking Network (S-DSN) is proposed for image classification and abnormal detection [86, 87]. Li et al. [86], Sun

et al. [87] stacked many sparse simplified neural network modules (SNNM) with mixed-norm regularization, in which weights are solved by using the convex optimization and the gradient descent algorithm. In order to make sparse SNNM learning the local dependencies between hidden units, Li et al. [88] split the hidden units or representations into different groups, which is termed as group sparse DSN (GS-DSN). The DSN idea is also utilized in the Deep Reinforcement Learning field. Zhang et al. [89] employed DSN method to integrate the observations from the formal network: Grasp network and Stacking network based on Q-learning algorithm to make an integrated robotic arm system do grasp and place actions. Wang et al. [90] stacked blocks multiple times to increase the performance of the neural architecture search task. Zhang et al. [91] presents a deep hierarchical multi-patch network for image deblurring via stacking approach.

Since there is no temporal representation of the data in DSNs, they are less effective to the problems where temporal dependencies exist in the input data. To embed the temporal information in DSNs, Recurrent Deep Stacking Networks (R-DSNs) [92] combined the advantages of DSNs and Recurrent neural networks (RNN). Unlike RNN which uses Back Propagation through time for training the network, R-DSNs use Echo State Network (ESN) to initialize the weights and then fine-tuning them via batch-mode gradient descent. A stacked extreme learning machine was proposed in [93]. Here, at each level of the network ELM with the reduced number of hidden nodes was used to solve the large scale problems. The number of hidden nodes was reduced via the principal component analysis (PCA) reduction technique. Keeping in view the efficiency of stacked models, the number of stacked models based on support vector machine have been proposed [94, 95, 96]. Traditional models like Random Forests have also been extended to deep architecture, known as deep forests [97], via stacking concept.

In addition to DSNs, there are some novel network architectures proposed based on the stacking method, Low et al. [98] contributed a stacking-based deep neural network (S-DNN) which is trained without a backpropagation algorithm. Kang et al. [99] presented a model by stacking conditionally restricted Boltzmann machine and deep neural network, which achieved significant superior performance with fewer parameters and fewer training samples.

4.4. Negative Correlation Based Deep Ensemble Methods

Negative correlation learning (NCL) [100] is an important technique for training the learning algorithms. The main concept behind the NCL is to encourage diversity among the individual

models of the ensemble to learn the diverse aspects of the training data. NCL minimizes the empirical risk function of the ensemble model via minimization of error functions of the individual networks. NCL [100] was evaluated for regression as well as classification tasks. The evaluation used different measures like simple averaging and winner-takes-all measures on classification tasks and simple average combination methods for regression problems. The authors figured out that winner-takes-all is better as compared to simple averaging in NCL ensemble models.

Shi et al. [101] proposed deep negative correlation learning architecture for crowd counting known as D-ConvNet i.e. decorrelated convolutional networks. Here, counting is done based on regression-based ensemble learning from a pool of convolutional feature mapped weak regressors. The main idea behind this is to introduce the NCL concept in deep architectures. Robust regression via deep NCL [102] is an extension of [101] in which theoretical insights about the Rademacher complexity are given and extended to more regression-based problems.

Buschjäger et al. [103] formulated a generalized bias-variance decomposition method to control the diversity and smoothly interpolates. They present the Generalized Negative Correlation Learning (GNCL) algorithm, which can encapsulate many existing works in literature and achieve superior performance.

The NCL can also be employed for incremental learning tasks. Muhlbaier and Polikar [104] employed a dynamically modified weighted majority voting strategy to combine the sub-classifiers. Tang et al. [105] proposed a negative correlation learning (NCL) based approach for ensemble incremental learning.

4.5. *Explicit / Implicit Ensembles*

Ensembling of deep neural networks doesn't seem to be an easy option as it may lead to increase in computational cost heavily due to the training of multiple neural networks. High performance hardware's with GPU acceleration may take weeks of weeks to train the deep networks. Implicit/Explicit ensembles obtain the contradictory goal wherein a single model is trained in such a manner that it behaves like ensemble of training multiple neural networks without incurring additional cost or to keep the additional cost as minimum as possible. Here, the training time of an ensemble is same as the training time of a single model. In implicit ensembles, the model parameters are shared and the single unthinned network at test times approximates the model averaging of the ensemble models. However, in explicit ensembles model parameters are not

shared and the ensemble output is taken as the combination of the predictions of the ensemble models via different approaches like majority voting, averaging and so on.

Dropout [106] creates an ensemble network by randomly dropping out hidden nodes from the network during the training of the network. During the time of testing, all nodes are active. Dropout provides regularization of the network to avoid overfitting and introduces sparsity in the output vectors. Overfitting is reduced as it trains exponential number of models with shared weights and provides an implicit ensemble of networks during testing. Dropping the units randomly avoids coadaptation of the units by making the presence of a particular unit unreliable. The network with dropout takes 2 – 3 times more time for training as compared to a standard neural network. Hence, a balance is to be set appropriately between the training time of the network and the overfitting. Generalization of DropOut was given in DropConnect [107]. Unlike DropOut which drops each output unit, DropConnect randomly drops each connection and hence, introduces sparsity in the weight parameters of the model. Similar to DropOut, DropConnect creates an implicit ensemble during test time by dropping out the connections (setting weights to zero) during training. Both DropOut and DropConnect suffer from high training time. To alleviate this problem, deep networks with Stochastic depth [108] aimed to reduce the network depth during training while keeping it unchanged during testing of the network. Stochastic depth is an improvement on ResNet [13] wherein residual blocks are randomly dropped during training and bypassing these transformation blocks connections via skip connections. Swapout [109] is a generalization of DropOut and Stochastic depth. Swapout involves dropping of individual units or to skip the blocks randomly. Embarking on a distinctive approach of reducing the test time, distilling the knowledge in a network [110] transferred the “knowledge” from ensembles to a single model. Gradual DropIn or regularised DropIn [111] of layers starts from a shallow network wherein the layers are added gradually. DropIN trains the exponential number of thinner networks, similar to DropOut, and also shallower networks.

All the aforementioned methods provided an ensemble of networks by sharing the weights. There have been attempts to explore explicit ensembles in which models do not share the weights. Snapshot ensembling [112] develops an explicit ensemble without sharing the weights. The authors exploited good and bad local minima and let the stochastic gradient descent (SGD) converge M -times to local minima along the optimization path and take the snapshots only when the model reaches the minimum. These snapshots are then ensembled by averaging at multiple local

Year	Authors	Contribution
2013	Wan et al. [107]	Introduced DropConnect (Random skipping of connections)
2014	Srivastava et al. [106]	Introduced Dropout (Random skipping of units)
2016	Huang et al. [108]	Deep networks with Stochastic depth (Random skipping of blocks)
2016	Singh et al. [109]	Introduced Swapout (Hybrid of Dropout and Stochastic depth approach)

Table 3: Implicit / Explicit ensembles

minima for object recognition. The training time of the ensemble is the same as that of the single model. The ensemble out is taken as the average of the output of the snapshot outputs at multiple local minimas. Random vector functional link network [113, 114] has also been explored for creating the explicit ensembles [115] where different random initialization of the hidden layer weights in a hierarchy diversifies the ensemble predictions.

Explicit/implicit produce ensembles out of a single network at the expense of base model diversity [25] as the lower level features across the models are likely to be the same. To alleviate this issue, branching based deep models [116] branch the network to induce more diversity. Motivated by different initializations of the neural networks leads to different local minima, Xue et al. [117] proposed deep ensemble model wherein ensemble of fully convolution neural network over multiloss module with coarse fine compensation module resulted in better segmentation of central serous chorioretinopathy lesion. Multiple neural networks with different initializations, multiple loss functions resulted in better diversity in an ensemble.

4.6. *Homogeneous & Heterogeneous ensembles*

Homogeneous ensemble (HOE) and heterogeneous ensemble (HEE) involve training a group of base learners either from the same family or different families, as shown in Fig. 5 and Fig. 6, respectively. Hence, each model of an ensemble must be as diverse as possible, and each base model must perform better than the random guess. The base learner can be a decision tree, neural network, or any other learning model.

In homogeneous ensembles, the same base learner is used multiple times to generate the family of base classifiers. However, the key issue is to train each base model such that the ensemble model is as diverse as possible, i.e. no two models are making the same error on a particular data sample. The two most common ways of inducing randomness in a homogeneous ensemble

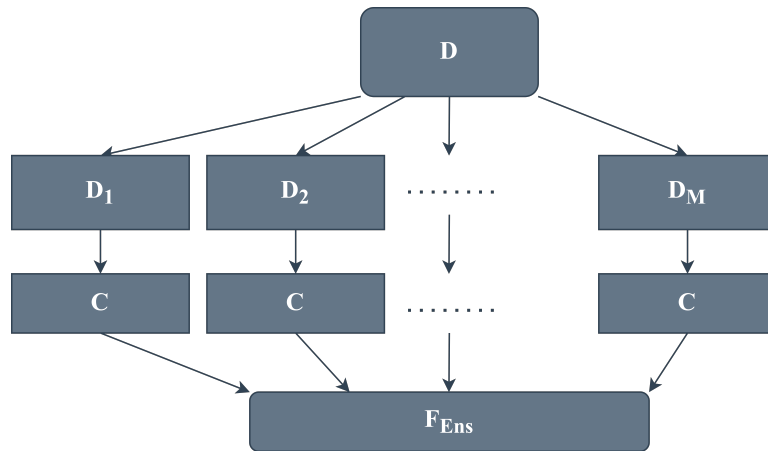


Figure 5: Homogeneous ensemble (HOE) has models based on the same algorithm, but each individual model are fed with distinct datasets.

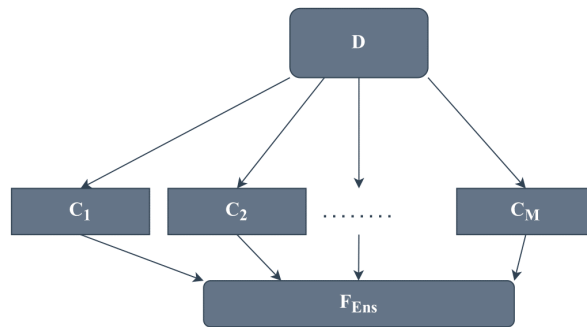


Figure 6: The components in heterogeneous ensemble (HEE) share the same dataset but consists of various algorithms.

are either sampling of the training set multiple times, thereby training each model on a different bootstrapped sample of the training data or sampling the feature space of the training data and train each model on a different feature subset of the training data. In some ensemble models like Random forest [5] used both these techniques for introducing diversity in the ensemble of decision trees. In neural networks, training models independently with different initialization of the models also induces diversity. However, deep learning models have high training costs

and hence, training of multiple deep learning models is not a feasible option. Some attempts, like horizontal vertical voting of deep ensembles [118] have been made to obtain ensembles of deep models without independent training. Temporal ensemble [119] trains multiple models with different input augmentation, different regularisation and different training epochs. Training of multiple deep neural networks for image classification [120] and for disease prediction [121] showed that better performance is achieved via an ensemble of multiple networks and averaging the outputs. Despite these models, training multiple deep learning models for ensemble is an uphill task as millions or billions of parameters need to be optimized. Hence, some studies have used deep learning in combination with traditional models to build heterogeneous ensemble models, enjoying the benefits of lower computation and higher diversity. Heterogeneous ensemble for default prediction [122] is an ensemble of the extreme gradient boosting, deep neural network and logistic regression. Heterogeneous ensemble for text classification [123] is an ensemble of multivariate Bernoulli naïve Bayes (MVNB), multinomial naïve Bayes (MNB), support vector machine (SVM), random forest (RF), and convolutional neural network (CNN) learning algorithms. Using different perspectives of data, model and decision fusion, heterogeneous deep network fusion [124] showed that complex heterogeneous fusion architectures are more diverse and hence, show better generalization performance. Furthermore, Seijo-Pardo et al. [125] employed both homogeneous and heterogeneous ensembles for feature selection. Zhao et al. [126] suggested that the heterogeneous bagging based ensemble strategy performs better than boosting based Learn++ algorithms and some other NCL methods. Other examples that employed homogeneous ensemble methods were used to deal with the presence of incremental tasks, such as concept drift [127], power load forecasting [128, 129], myoelectric prosthetic hands surface electromyogram characteristics [130], etc. Das et al. [131] proposed an ensemble incremental learning with pseudo-outer-product fuzzy neural network for traffic flow prediction, real-life stock price, and volatility predictions, etc.

4.7. Decision Fusion Strategies

Ensemble learning trains several base learners and aggregates the outputs of base learners using some rules. The rule used to combine the outputs determines the effective performance of an ensemble. Most of the ensemble models focus on the ensemble architecture followed by their naive averaging to predict the ensemble output. However, naive averaging of the models, followed in most of the ensemble models, is not data adaptive and leads to less optimal perfor-

mance [132] as it is sensitive to the performance of the biased learners. As there are billions of hyperparameters in deep learning architecture, the issue of overfitting may lead to the failure of some base learners. Hence, to overcome these issues, approaches like Bayes optimal classifier and super learner have been followed [132].

The different approaches followed in the literature for combining the outputs of the ensemble models are:

4.7.1. *Unweighted Model Averaging*

Unweighted averaging of the outputs of the base learners in an ensemble is the most followed approach for fusing the decisions in the literature. Here, the outcomes of the base learners are averaged to get the final prediction of the ensemble model. Deep learning architectures have high variance and low bias, thus, simple averaging of the ensemble models improve the generalization performance due to the reduction of the variance among the models.

The averaging of the base learners is performed either on the outputs of the base learners directly or on the predicted probabilities of the classes via softmax function:

$$P_i^j = \text{softmax}^j(O_i) = \frac{O_i^j}{\sum_{k=1}^K \exp(O_k^j)} \quad (4)$$

where P_i^j is the probability outcome of the i^{th} unit on the j^{th} base learner, O_i^j is the output of the i^{th} unit of the j^{th} base learner and K is the number of the classes.

Unweighted averaging is a reasonable choice when the performance of the base learners is comparable, as suggested in [13, 133, 134]. However, when the ensemble contains heterogeneous base learners naive unweighted averaging may result in suboptimal performance as it is affected by the performance of the weak learners and the overconfident learners [132]. The adaptive metalearner should be good enough to adaptively combine the strengths of the base learners as some learners may have lower overall performance but maybe good at the classification of certain subclasses and hence, leading to better overall performance.

4.7.2. *Majority Voting*

Similar to unweighted averaging, majority voting combines the outputs of the base learners. However, instead of taking the average of the probability outcomes, majority voting counts the votes of the base learners and predicts the final labels as the label with the majority of votes. In comparison to unweighted averaging, majority voting is less biased towards the outcome of

a particular base learner as the effect is mitigated by majority vote count. However, favouring of a particular event by most of the similar base learners or dependent base learners leads to the dominance of the event in the ensemble model. In majority voting, the analysis by Kuncheva et al. [135] showed that the pairwise dependence among the base learners plays an important role and for the classification of images, the prediction of shallow networks is more diverse as compared to the deeper networks [136]. Hence, Ju et al. [132] hypothesised that the performance of the majority voting based shallows ensemble models is better as compared to the majority based deep ensemble models.

Voting methods have also started to be integrated with semi-supervised deep learning. Li et al. [137] proposed an ensemble semi-supervised deep acoustic models for in automatic speech recognition. Wang et al. [138] explored an ensemble self-learning method to enhance semi-supervised performance and extract adverse drug events from social media in [139]. In the semi-supervised classification area, the author proposed a deep coupled ensemble learning method which is combined with complementary consistency regularization and gets the state of the art performance in [140]. Some results have also been achieved with semi-supervised ensemble learning on some datasets where the annotation is costly. Pio et al. [141] employed an ensemble method to improve the reliability of miRNA:miRNA predicted interactions.

Furthermore, the multi-label classification [142] problem is also a major point addressed by the voting method, a typical application is the RANdom k -labELsets (RAKEL) algorithm [143]. The author trained several single-label classifiers using small random subsets of actual labels. Then the final output is carried out by a voting scheme based on the predictions of these single classifiers. There are also many variants of RAKEL proposed in recent years [144, 145, 146]. Shi et al. [147] proposed a solution for multi-label ensemble learning problem, which construct several accurate and diverse multi-label based basic classifiers and employ two objective functions to evaluate the accuracy and diversity of multi-label base learners. Another work [148] proposed an ensemble multi-label classification framework based on variable pairwise constraint projection. Xia et al. [149] proposed a weighted stacked ensemble scheme that employs the sparsity regularization to facilitate classifier selection and ensemble construction. Besides, there are many applications of ensemble multi-label methods. Some publications employ multi-label ensemble classifiers to explore the protein, such as protein subcellular localization [150], protein function prediction [151], etc. The Muli-label classifier is also utilized in predicting the drug side

effects [152], predicting the gene prediction[153], etc. Moreover, there is another critical ensemble multi-label algorithm called ensemble classifier chains (ECC) [154]. This method involves binary classifiers linked along a chain. The first classifier is trained using only the input data, and then each subsequent classifier is trained on the input space and all previous classifiers in the chain. The final prediction is obtained by the integration of the predictions and selection above a manually set threshold. Chen et al. [155] propose an ensemble application of convolutional and recurrent neural networks to capture both the global and local textual semantics and to model high-order label correlations.

4.7.3. *Bayes Optimal Classifier*

In Bayesian method, hypothesis h_j of each base learner with the conditional distribution of target label t given x . Let h_j be the hypothesis generated on the training data D evaluated on test data (x, t) , mathematically, $h_j(t|x) = P[y|x, h_j, D]$. With Bayes rule, we have

$$P(t|x, D) \propto \sum_{h_j} P[t|h_j, x, D]P[D|h_j]P[h_j] \quad (5)$$

and the Bayesian Optimal classifier is given as:

$$\underset{t}{\operatorname{argmax}} \sum_{h_j} P[t|h_j, x, D]P[D|h_j]P[h_j], \quad (6)$$

where $P[D|h_j] = \prod_{(t,x) \in D} h_j(t|x)$ is the likelihood of the data under h_j . However, due to overfitting issues this might be not a good measure. Hence, training data is divided into two sets-one for training the model and the other for evaluating the model. Usually validation set is used to tune the hyperparameters of the model.

Choosing prior probabilities in Bayes optimal classifier is difficult and hence, usually set to uniform distribution for simplicity. With a large sample size, one hypothesis tends to give larger posterior probabilities than others and hence the weight vector is dominated by a single base learner and hence Bayes optimal classifier would behave as the discrete superlearner with a negative likelihood loss function.

4.7.4. *Stacked Generalization*

Stacked generalization [77] works by deducing the biases of the generalizer(s) with respect to a provided learning set. To obtain the good linear combination of the base learners in regression, cross-validation data and least squares under non-negativity constraints was used to get the

optimal weights of combination [156]. Consider the linear combination of the predictions of the base learners f_1, f_2, \dots, f_m given as:

$$f_{stacking}(x) = \sum_{j=1}^m w_j f_j(x) \quad (7)$$

where w is the optimal weight vector learned by the meta learner.

4.7.5. *Super Learner*

Inspired by the cross validation for choosing the optimal classifier, Van der Laan et al. [157] proposed super learner which is weighted combination of the predictions of the base learner. Unlike the stacking approach, it uses cross validation approach to select the optimal weights for combining the predictions of the base learners.

With smaller datasets, cross validation approach can be used to optimize the weights. However, with the increase in the size of the data and the number of base learners in the model, it may not be a feasible option. Instead of optimizing the V-fold cross validation, single split cross validation can also be used for optimizing the weights for optimal combination [158]. In deep learning models, usually, a validation set is used to evaluate the performance instead of using the cross validation.

Another application field for super learner is in Reinforcement Learning. With the development of Deep learning, some researchers have implemented deep reinforcement learning, which combines deep learning with a Q-learning algorithm [159]. Ensemble methods in deep Q learning have decent performance. Chen et al. [160] proposed an ensemble network architecture for deep reinforcement learning. The integrated network includes Temporal Ensemble and Target Values Ensemble. Develop a human-like chat robot is a challenging job, by incorporating deep reinforcement learning and ensemble method, Cuayáhuitl et al. [161] integrated 100 deep reinforcement learning agents, the agents are trained based on clustered dialogues. They also demonstrate the ensemble of DRL agents has better performance than the single variant or Seq2Seq model. Stock trading is another topic where ensemble deep reinforcement learning has achieved a promising result. Carta et al. [162] found the single supervised classifier is inadequate to deal with the complex and volatile stock market. They employed hundreds of neural networks to pre-process the data, then they combined several reward-based meta learners as a trading agency. Moreover, Yang et al. [163] trained an ensemble trading agency based on three different metrics: Proximal Policy Optimization (PPO), Advantage Actor-Critic (A2C), and Deep

Deterministic Policy Gradient (DDPG). The ensemble strategy combines the advantages of the three different algorithms. Besides, some researchers try to use ensemble strategy to solve the disease-prediction problem. The proposed model in [164] consists of several sub-models which are in response to different anatomical parts.

4.7.6. Consensus

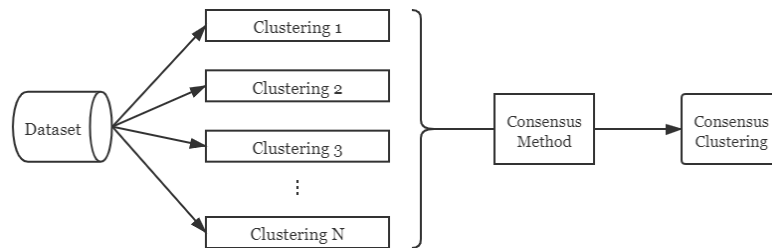


Figure 7: The process of consensus clustering. An ensemble of different clustering results can be combined by a consensus approach.

Unsupervised learning is another group of machine learning techniques. The fundamental difference between it and supervised learning is that unsupervised learning usually handles training samples without corresponding labels. Therefore, the primary usage of unsupervised learning is to do clustering. The reason why ensemble methods are employed is to combine some weak clusters into strong one. To create diverse clusters, several approaches can be applied: using different sampling data, using different subsets of the original features, and employing different clustering methods [165]. Sometimes, even some random noise can be added to these base models to increase randomness, which is good for ensemble methods according to [166]. After receiving all the outputs from each cluster, various consensus functions can be chosen to obtain the final output based on the user’s requirement [22]. The ensemble clustering is also known as consensus clustering Fig. 7.

Zhou and Tang [167] explored ensemble methods for unsupervised learning and developed four different approaches to combine the outputs of these clusters. In recent years, some new ensemble clustering methods have been proposed that illustrated the priority of ensemble learning [168, 169, 170]. Most of the clustering ensemble method is based on the co-association matrix solution, which can be regarded as a graph partition problem. Besides, there is some research

focus on integrating the deep structure and ensemble clustering method. Liu et al. [171, 172] firstly showed that ensemble unsupervised representation learning with deep structure can be applied in large scale data. Then the author combined the method with auto-encoder and extends it to the vision field. Shaham et al. [173] first demonstrated that some crowdsourcing algorithms can be replaced by a Restricted Boltzmann Machine with a single hidden neuron, then propose an RBM-based Deep Neural Net (DNN) used for unsupervised ensemble learning. The unsupervised ensemble method also makes some contribution to the field of Natural Language Processing. Alami et al. [174] demonstrated that the ensemble of unsupervised deep neural network models that use Sentence2Vec representation as the input has the best performance according to the experiments. Hassan et al. [175] proposed a module that includes four semantic similarity measures, which improves the performance on the semantic textual similarity (STS) task. The unsupervised ensemble method is also widely used for tasks that lack annotation, such as the medical image. Ahn et al. [176] proposed an unsupervised feature learning method integrated ensemble approach with a traditional convolutional neural network. Lahiri et al. [177] employed unsupervised hierarchical feature learning with ensemble sparsely autoencoder on retinal blood vessels segmentation task, meanwhile, Liu et al. [178] also propose an unsupervised ensemble architecture to automatically segment retinal vessel. Besides, there are also some ensemble deep methods working on localization predicting for long non-coding RNAs [179]. Hu and Suganthan [180] extended the ensemble random vector functional link to unsupervised tasks. The authors employ manifold regularization to re-represent the original features, and then use the Kuhn-Munkre algorithm with consensus clustering to ensemble the clustering results from multiple hidden layers.

4.7.7. Query-By-Committee

Active Learning is another popular topic in the deep learning area, which is also often used in conjunction with semi-supervised learning and ensemble learning. The key sight of this is to make the algorithm learning from less annotated data. Some conventional active learning algorithms, such as Query-By-Committee (as shown in Fig 8), have already adopted the idea of ensemble learning. Melville and Mooney [181, 182] explored an ensemble method that builds a diverse committee. Beluch et al. [183] discussed the power of ensembles for active learning is significantly better than Monte-Carlo Dropout and geometric approaches. Sharma and Rani [184] show some applications in drug-target interaction prediction. Ensemble active learning is

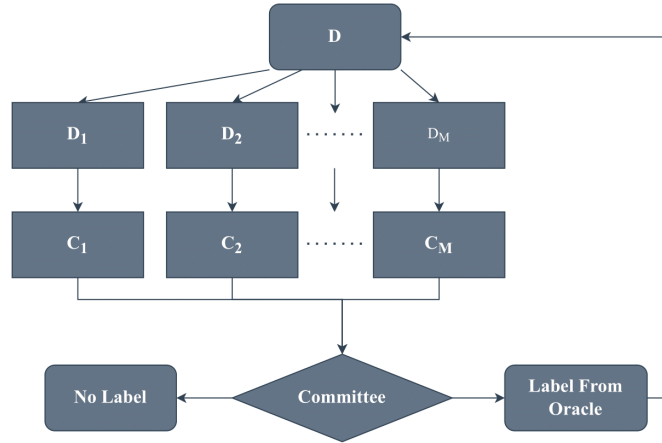


Figure 8: Query-by-committee in Active Learning. Sampling with replacement is used to partition the labeled training data set into training splits. The committee determines whether to label the data based on the output of several algorithms.

also available to conquer the concept drift and class imbalance problem [185].

5. Applications

In this section, we briefly present the applications of deep ensemble models across different domains in a tabular form. Ensemble deep models have been implemented in several domains and therefore, in broad sense, we have classified the application domains into five categories, i.e., health care, speech, image classification, forecasting and the rest models are listed in others category. Table 4 gives the information about the ensemble deep models that have been implemented in health care domain. Here, several papers are based on heterogeneous ensemble technique. It reveals that using different family’s models into a single frame perform better in health care domain. Recently, ensemble deep techniques have been successful and have shown good performance in health care domain. The models which have been implemented for speech task have been given in Table 5 and most of the ensemble approaches are based on stacking technique. Table 6 contains the ensemble deep models that have been implemented in speech areas. Models that have been implemented in forecasting and other domains have been given in Table 7 and Table 8, respectively. Fig. 9 shows the percentages of the application domains.

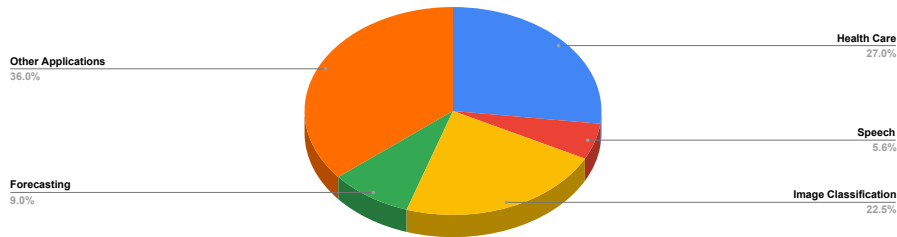


Figure 9: Ensemble-based approach in different areas. Data from Tables 4 to 8.

The statistics reveals that different ensemble deep techniques have been used in different areas. A larger number of models, i.e. 27% of the ensemble deep models, have been implemented in health care domain and 5.6% percent of the models for speech application and 22.5% of the models for image classification task. Moreover, 9% models have been used in forecasting and 36% in other applications areas, i.e., information retrieval, emotion recognition, text categorization and so on. Fig. 10 shows the ensemble strategies in percentage. In ensemble learning, there are several ways to integrate the outcomes of the models in an ensemble. In the literature, researchers have proposed different techniques of decision fusion according to different areas of application. Bagging, boosting and stacking are the classical ensemble techniques. Based on these three techniques, researchers have developed several other techniques also. Boosting (18.2%), stacking (12.5%) and bagging (4.5%) techniques have been implemented in ensemble deep framework. Heterogeneous and implicit ensemble are also popular for making an efficient ensemble model and their contribution are as follows 11.4% and 10.2%, respectively. The rest ensemble techniques are, i.e. unsupervised (3.4%), NCL (3.4%), reinforcement (1.1%), active learning (1.1%), explicit ensemble (1.1%) and homogeneous ensemble (3.4%).

6. Conclusions and future works

In this paper, we reviewed the recent developments of ensemble deep learning models. The theoretical background of ensemble learning has been elaborated to understand the success of ensemble learning. The various approaches ranging from traditional ones like bagging, boosting to the recent novel approaches like implicit/explicit ensembles, heterogeneous ensembles, have

Year	Author	Title	Approach	Area
2014	Zheng et al. [186]	HIBAG—HLA genotype imputation with attribute bagging	Bagging	Genotype Imputation
2014	Cortes et al. [60]	Deep Boosting	Boosting	Classification
2015	Zhang et al. [152]	Predicting drug side effects by multi-label learning and ensemble learning	Decision Fusion	Predict the drug side effects
2016	Guo et al. [150]	Human protein subcellular localization with integrated source and multi-label ensemble classifier.	Decision Fusion	Protein subcellular localization prediction
2016	Lahiri et al. [177]	Deep neural ensemble for retinal vessel segmentation in fundus images towards achieving label-free angiography	Decision Fusion	Medical image segmentation
2017	Cabria and Gondra [187]	MRI segmentation fusion for brain tumor detection	Heterogeneous ensemble	MRI segmentation
2018	Grassmann et al. [121]	A deep learning algorithm for prediction of age-related eye disease study severity scale for age-related macular degeneration from color fundus photography	Homogeneous ensemble	Disease prediction
2018	Cao et al. [179]	The Inclocator: a subcellular localization predictor for long non-coding RNAs based on a stacked ensemble classifier	Decision Fusion	Subcellular localization predictor
2018	Sharma and Rani [184]	Be-dti: Ensemble framework for drug target interaction prediction using dimensionality reduction and active learning	Active learning	Drug target interaction prediction
2019	Ahn et al. [176]	Unsupervised feature learning with k-means and an ensemble of deep convolutional neural networks for medical image classification	Decision Fusion	Medical image classification
2019	Liu et al. [178]	Unsupervised ensemble strategy for retinal vessel segmentation.	Unsupervised	Medical image classification
2020	Shalhaf and Vafaeezadeh [188]	Automated detection of COVID-19 using ensemble of transfer learning with deep convolutional neural network based on CT scans	Heterogeneous Ensemble	Detection of COVID-19
2020	Ali et al. [189]	A smart healthcare monitoring system for heart disease prediction based on ensemble deep learning and feature fusion	Boosting	Heart disease prediction
2021	Zhou et al. [190]	The ensemble deep learning model for novel COVID-19 on CT images	Heterogeneous Ensemble	Detection of COVID-19
2021	Li et al. [191]	Intelligent Fault Diagnosis by Fusing Domain Adversarial Training and Maximum Mean Discrepancy via Ensemble Learning	Heterogeneous Ensemble	Fault diagnosis
2021	Das et al. [192]	Automatic COVID-19 detection from X-ray images using ensemble learning with convolutional neural network	Heterogeneous Ensemble	Detection of COVID-19
2022	Sukegawa et al. [193]	Identification of osteoporosis using ensemble deep learning model with panoramic radiographs and clinical covariates	Decision Fusion	Identification of osteoporosis
2022	Gao et al. [194]	Vessel segmentation for X-ray coronary angiography using ensemble methods with deep learning and filter-based features	Boosting	Vessel segmentation
2022	Rath et al. [195]	Improved heart disease detection from ECG signal using deep learning based ensemble model	Heterogeneous Ensemble	Heart disease detection
2022	Tanveer et al. [196]	Classification of Alzheimer's Disease Using Ensemble of Deep Neural Networks Trained Through Transfer Learning	Heterogeneous Ensemble	Classification of Alzheimer's Disease
2022	Rai and Chatterjee [197]	Hybrid CNN-LSTM deep learning model and ensemble technique for automatic detection of myocardial infarction using big ECG data	Heterogeneous Ensemble	Detection of myocardial infarction
2022	Ganaie and Tanveer [198]	Ensemble deep random vector functional link network using privileged information for Alzheimer's disease diagnosis	Implicit ensemble	Diagnosis of Alzheimer's disease

Table 4: Applications in health care

Year	Author	Title	Approach	Area
2012	Tur et al. [199]	Towards deeper understanding: Deep convex networks for semantic utterance classification	Stacking	Semantic Utterance Classification
2012	Deng et al. [200]	Use of kernel deep convex networks and end-to-end learning for spoken language understanding	Stacking	Spoken Language Understanding
2014	Deng and Platt [201]	Ensemble deep learning for speech recognition	Stacking	Speech Recognition
2014	Palangi et al. [92]	Recurrent Deep-Stacking Networks for sequence classification	Stacking	Sequence classification
2017	Li et al. [137]	Semi-supervised ensemble DNN acoustic model training	Decision Fusion	Speech Recognition

Table 5: Applications in speech

Year	Author	Title	Approach	Area
2012	Ciregan et al. [120]	Multi-column deep neural networks for image classification	Homogeneous ensemble	Classification
2013	Wan et al. [107]	Regularization of Neural Networks using DropConnect	Implicit ensemble	Image recognition
2014	Srivastava et al. [106]	Dropout: a simple way to prevent neural networks from overfitting	Implicit Ensemble	Computer vision, speech recognition document classification and computational biology
2014	Liu et al. [59]	Facial expression recognition via a boosted deep belief network	Boosting	Facial expression recognition
2015	Li et al. [86]	Sparse deep stacking network for image classification	Stacking	Image Classification
2015	Yang et al. [66]	Convolutional channel features	Boosting	Pedestrian detection, face detection, edge detection and object proposal generation
2016	Moghimi et al. [63]	Boosted Convolutional Neural Networks	Boosting	Classification
2016	Huang et al. [108]	Deep networks with stochastic depth	Implicit ensemble	Classification
2016	He et al. [13]	Deep residual learning for image recognition	Implicit ensemble	classification, and object detection
2016	Singh et al. [109]	Swapout: Learning an ensemble of deep architectures	Implicit ensemble	classification
2016	Smith et al. [111]	Gradual dropin of layers to train very deep neural networks	Implicit ensemble	Classification
2016	Laine and Aila [119]	Temporal ensembling for semi-supervised learning	Homogeneous ensemble	Classification
2016	Tang et al. [164]	Inquire and diagnose: Neural symptom checking ensemble using deep reinforcement learning	Decision Fusion	Inquire symptoms and diagnose diseases
2017	Huang et al. [112]	Snapshot ensembles: train 1, get M for free	Explicit ensemble	Classification
2017	Mosca and Magoulas [70]	Deep incremental boosting	Boosting	Classification
2018	Beluch et al. [183]	The power of ensembles for active learning in image classification	Decision Fusion	Image classification
2019	Amin-Naji et al. [202]	Ensemble of CNN for multi-focus image fusion	Decision Fusion	Image Classification
2019	Li et al. [140]	Semi-supervised deep coupled ensemble learning with classification landmark exploration.	Decision Fusion	Image classification
2020	Wang et al. [90]	Particle swarm optimisation for evolving deep neural networks for image classification by evolving and stacking transferable blocks.	Stacking	Image Classification

Table 6: Applications in image classification

Year	Author	Title	Approach	Area
2014	Qiu et al. [203]	Ensemble deep learning for regression and time series forecasting	Decision Fusion	Regression and Time Series Forecasting
2016	Grmanová et al. [129]	Incremental ensemble learning for electricity load forecasting	Decision Fusion	Electricity load forecasting
2017	Qiu et al. [204]	Empirical Mode Decomposition based ensemble deep learning for load demand time series forecasting	Decision Fusion	Load demand forecasting
2017	Liu et al. [205]	A Flood Forecasting Model Based on Deep Learning Algorithm via Integrating Stacked Autoencoders with BP Neural Network	Stacking	Flood Forecasting
2018	Qiu et al. [128]	Ensemble incremental learning random vector functional link network for short-term electric load forecasting.	Decision Fusion	Electric load forecasting
2020	Carta et al. [162]	A multi-layer and multi-ensemble stock trader using deep learning and deep reinforcement learning	Implicit ensemble	Stock trader
2020	Yang et al. [163]	Deep reinforcement learning for automated stock trading: An ensemble strategy	Decision Fusion	Stock trading agency
2021	Bhusal et al. [206]	Deep ensemble learning-based approach to real-time power system state estimation	Stacking	Electric Power
2022	Singla et al. [207]	An ensemble method to forecast 24-h ahead solar irradiance using wavelet decomposition and BiLSTM deep learning network	Decision Fusion	Forecasting

Table 7: Applications in forecasting

Year	Author	Title	Approach	Area
2013	Deng et al. [208]	Deep stacking networks for information retrieval	Stacking	Information Retrieval
2014	Kuznetsov et al. [61]	Multi-class deep boosting	Boosting	Classification
2014	Wang et al. [209]	Sentiment classification The contribution of ensemble learning	Bagging, Boosting	Sentiment classification
2015	Zarepoor and Shamsolmoali [210]	Application of Credit Card Fraud Detection: Based on Bagging Ensemble Classifier	Bagging	Credit Card Fraud Detection
2016	Yin et al. [211]	Recognition of emotions using multimodal physiological signals and an ensemble deep learning model	Decision Fusion	Emotions Recognition
2016	Liu et al. [172]	A deep learning approach to unsupervised ensemble learning	Decision Fusion	Clustering
2016	Walach and Wolf [67]	Learning to count with CNN boosting	Boosting	Object counting in images
2017	Han et al. [62]	Incremental boosting convolutional neural network for facial action unit recognition	Boosting	Facial action unit recognition
2017	Chen et al. [155]	Ensemble application of convolutional and recurrent neural networks for multi-label text categorization	Decision Fusion	Text Categorization.
2017	Opitz et al. [68]	Bier-boosting independent embeddings robustly	Boosting	Image retrieval
2018	Shi et al. [101]	Crowd Counting with Deep Negative Correlation Learning	Negative correlation learning	Crowd Counting
2018	Kazemi et al. [212]	Novel genetic-based negative correlation learning for estimating soil temperature	Negative correlation learning	Soil Temperature Estimation
2018	Randhawa et al. [213]	Credit Card Fraud Detection Using AdaBoost and Majority Voting	Boosting	Credit Card Fraud Detection
2018	Sun et al. [87]	Sparse Deep Stacking Network for Fault Diagnosis of Motor	Stacking	Fault Diagnosis
2018	Chen et al. [64]	Deep boosting for image denoising	Boosting	Image denoising
2018	Li et al. [122]	Heterogeneous ensemble for default prediction of peer-to-peer lending in China	Heterogeneous ensemble	Default prediction
2018	Kilimci and Akyokus [123]	Deep learning-and word embedding-based heterogeneous classifier ensembles for text classification	Heterogeneous ensemble	Classification
2018	Liu et al. [139]	SSEL-ADE: a semi-supervised ensemble learning framework for extracting adverse drug events from social media	Decision Fusion	Extracting adverse drug events
2019	Martín et al. [214]	Android malware detection through hybrid features fusion and ensemble classifiers: The AndroPyTool framework and the OmniDroid dataset	Decision Fusion	Android malware detection
2019	Cuayáhuitl et al. [161]	Ensemble-based deep reinforcement learning for chatbots	Reinforcement	Chat robot
2019	Chen et al. [65]	Real-world image denoising with deep boosting	Boosting	Image denoising
2019	Walmer et al. [69]	HiBsteR: Hierarchical Boosted Deep Metric Learning for Image Retrieval	Boosting	Image Retrieval
2019	Wang et al. [215]	Adaboost-based security level classification of mobile intelligent terminals	Adaboost	Security Level Classification
2019	Chen et al. [216]	Novel Hybrid Integration Approach of Bagging-Based Fisher's Linear Discriminant Function for Groundwater Potential Analysis	Bagging	Groundwater Potential Analysis
2019	Shi et al. [115]	Random vector functional link neural network based ensemble deep learning	Implicit ensemble	Classification
2019	Zhang et al. [91]	Deep stacked hierarchical multi-patch network for image deblurring	Stacking	Deblurring Image
2019	Alami et al. [174]	Enhancing unsupervised neural networks based text summarization with word embedding and ensemble learning	Decision Fusion	Text summarization
2019	Hassan et al. [175]	Uests: An unsupervised ensemble semantic textual similarity method	Decision Fusion	Semantic textual similarity
2020	Zhang et al. [89]	Grasp for stacking via deep reinforcement learning	Stacking	Robotic arm control
2020	Zhang et al. [71]	Snapshot boosting: a fast ensemble framework for deep neural networks	Boosting	Computer vision (CV) and the natural language processing (NLP) tasks
2021	Tsogbaatar et al. [217]	Del-IoT: A deep ensemble learning approach to uncover anomalies in IoT	Decision Fusion	IoT
2022	Wen et al. [218]	A new ensemble convolutional neural network with diversity regularization for fault diagnosis	Snapshot ensemble learning	Fault diagnosis
2022	Hu and Suganthan [180]	Representation Learning Using Deep Random Vector Functional Link Networks for Clustering	Decision Fusion	Clustering

Table 8: Other applications

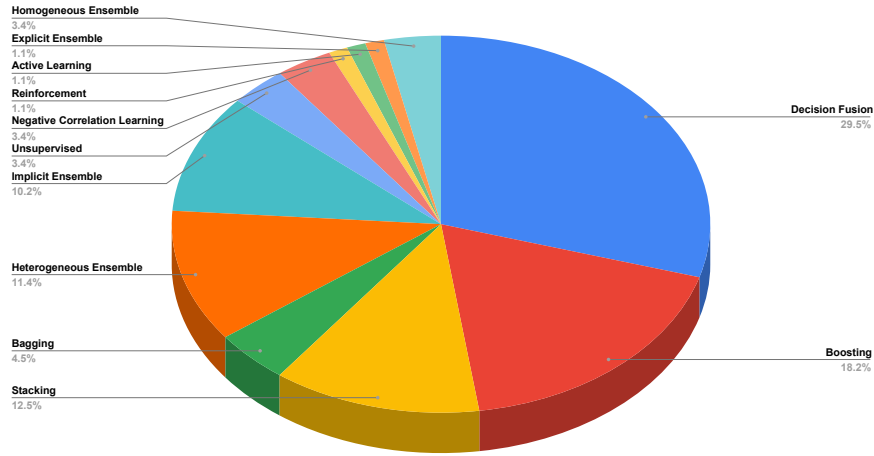


Figure 10: Analysis of the applications of various ensemble methods. Data from Tables 4 to 8.

led to better performance of deep ensemble models. We also reviewed the applications of the deep ensemble models in different domains.

Although deep ensemble models have been applied across different domains, there are several open problems which can be explored in the future to fill the gap. Big data [219] is still a challenging problem, one can explore the benefits of deep ensemble models for learning the patterns using the techniques like implicit deep ensemble to maximize the performance in both time and generalization aspects.

Deep learning models are difficult to train than shallow models as large number of weights corresponding to different layers need to be tuned. Creating deep ensemble models may further complicate the problem. Hence, randomized models can be explored to overcome the training cost. Bagging based deep ensemble may incur heavy training time for optimizing the ensemble models. Hence, one can investigate the alternate ways of inducing diversity in the base models with lesser training cost. Randomized learning modules like random vector functional link network [113] are best suited for creating the ensemble models as randomized models lead to a significant variance reduction. Also, the hidden layers are randomly initialized, hence, can be used to create deep ensembles without incurring any additional cost of training [115]. Random-

ized modules can be further explored using different techniques like implicit / explicit ensembles [115], stacking based ensembles [220]. However, there are still open directions which can be worked upon like negative correlation learning, heterogeneous ensembles and so on.

Implicit/explicit ensembles are faster compared to training of multiple deep models. However, creating diversity within a single model is a big challenge. One can explore the methods to induce more diversity among the learners within these ensembles like branching based deep models [116]. Investigate the extension of explicit/implicit ensembles to traditional models.

Following the stacking based approach, Deep convex net (DCN) [79], traditional methods like random forest [5, 97], support vector machines [94, 95, 96] have been extended to deep learning architectures which resulted in improved performance. One can investigate these traditional models for creating the deep ensemble models.

Another big challenge of ensemble deep learning lies in model selection for building the ensemble architecture, homogeneous and heterogeneous ensembles represent two different ways for choosing the model. However, to answer how many different algorithms, and how many base learners in the ensemble architecture, are still problem-dependent. Finding a criterion for model selection in ensemble deep learning should be an important target for researchers in the next few years. Since most of the models focus on developing the architectures with little attention towards how to combine the base learners prediction is still unanswered. Hence, one can investigate the effect of different fusion strategies on the prediction of an ensemble output.

For unsupervised ensemble learning or consensus clustering, the ensemble approaches include but are not limited to: Hyper-graph partitioning, Voting approach, Mutual information, etc. Consensus clustering is a powerful tool and it can improve performance in most cases. However, there are many concerns remain to be tackled, it is exquisitely sensitive, which might assert as an apparent structure without obvious demarcation or declared cluster stable without cluster resistance. Besides, current method cannot handle some complex but possible scenarios, such as the boundary samples are assigned to the single cluster, clusters do not intersect and the methods are not able to represent outliers. These are the possible research directions for future work.

The problem of semi-supervised ensemble domains has not been extensively studied yet, and most of the literature shows that semi-supervised ensemble methods are mainly used in cases where there is insufficient labeling data. Also, combining the semi-supervision with some other machine learning methods, such as active learning, is a direction for future research.

Reinforcement learning is another popular topic recently. The idea of integrating model-based reinforcement learning with ensemble learning has been used with promising results in many applications, but there is little integration of planning & learning-based reinforcement learning with ensemble learning methods.

Acknowledgment

The funding for this work is provided by the National Supercomputing Mission under DST and Miety, Govt. of India under Grant No. DST/NSM/ R&D_HPC_Appl/2021/03.29, as well as the Department of Science and Technology under Interdisciplinary Cyber Physical Systems (ICPS) Scheme grant no. DST/ICPS/CPS-Individual/2018/276. Mr. Ashwani Kumar Malik acknowledges the financial support (File no - 09/1022 (0075)/2019-EMR-I) given as scholarship by Council of Scientific and Industrial Research (CSIR), New Delhi, India. We are grateful to IIT Indore for the facilities and support being provided.

References

- [1] M. d. Condorcet, *Essay on the application of analysis to the probability of majority decisions*, Paris: Imprimerie Royale (1785).
- [2] L. K. Hansen, P. Salamon, Neural network ensembles, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 12 (1990) 993–1001.
- [3] T. G. Dietterich, Ensemble methods in machine learning, in: *International Workshop on Multiple Classifier Systems*, Springer, 2000, pp. 1–15.
- [4] R. Kohavi, D. H. Wolpert, Bias plus variance decomposition for zero-one loss functions, in: *ICML*, volume 96, 1996, pp. 275–83.
- [5] L. Breiman, Random forests, *Machine Learning* 45 (2001) 5–32.
- [6] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, Imagenet large scale visual recognition challenge, *International Journal of Computer Vision* 115 (2015) 211–252.
- [7] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional neural networks, in: *Advances in Neural Information Processing Systems*, 2012, pp. 1097–1105.
- [8] L. Deng, D. Yu, Deep learning: methods and applications, *Foundations and Trends® in Signal Processing* 7 (2014) 197–387.
- [9] I. Goodfellow, Y. Bengio, A. Courville, *Deep learning*. book in preparation for MIT press, URL: <http://www.deeplearningbook.org> (2016).
- [10] Y. LeCun, Y. Bengio, G. Hinton, Deep learning, *Nature* 521 (2015) 436–444.
- [11] S. Hochreiter, *Untersuchungen zu dynamischen neuronalen netzen*, Diploma, Technische Universität München 91 (1991).

- [12] X. Glorot, Y. Bengio, Understanding the difficulty of training deep feedforward neural networks, in: Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics, 2010, pp. 249–256.
- [13] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 770–778.
- [14] R. K. Srivastava, K. Greff, J. Schmidhuber, Training very deep networks, in: Advances in Neural Information Processing Systems, 2015, pp. 2377–2385.
- [15] A. Veit, M. J. Wilber, S. Belongie, Residual networks behave like ensembles of relatively shallow networks, in: Advances in Neural Information Processing Systems, 2016, pp. 550–558.
- [16] Y. Zhao, J. Gao, X. Yang, A survey of neural network ensembles, in: 2005 International Conference on Neural Networks and Brain, volume 1, IEEE, 2005, pp. 438–442.
- [17] L. Rokach, Ensemble-based classifiers, *Artificial Intelligence Review* 33 (2010) 1–39.
- [18] D. Gopika, B. Azhagusundari, An Analysis on Ensemble Methods In Classification Tasks, *International Journal of Advanced Research in Computer and Communication Engineering* 3 (2014) 7423–7427.
- [19] P. Yang, Y. Hwa Yang, B. B Zhou, A. Y Zomaya, A review of ensemble methods in bioinformatics, *Current Bioinformatics* 5 (2010) 296–308.
- [20] J. Mendes-Moreira, C. Soares, A. M. Jorge, J. F. D. Sousa, Ensemble approaches for regression: A survey, *ACM Computing Surveys (csur)* 45 (2012) 10.
- [21] Y. Ren, P. Suganthan, N. Srikanth, Ensemble methods for wind and solar power forecasting—a state-of-the-art review, *Renewable and Sustainable Energy Reviews* 50 (2015) 82–91.
- [22] S. Vega-Pons, J. Ruiz-Shulcloper, A survey of clustering ensemble algorithms, *International Journal of Pattern Recognition and Artificial Intelligence* 25 (2011) 337–372.
- [23] Y. Ren, L. Zhang, P. N. Suganthan, Ensemble classification and regression-recent developments, applications and future directions, *IEEE Computational Intelligence Magazine* 11 (2016) 41–53.
- [24] O. Sagi, L. Rokach, Ensemble learning: A survey, *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 8 (2018) e1249.
- [25] Y. Cao, T. A. Geddes, J. Y. H. Yang, P. Yang, Ensemble deep learning in bioinformatics, *Nature Machine Intelligence* 2 (2020) 500–508.
- [26] A. Krogh, J. Vedelsby, Neural network ensembles, cross validation, and active learning, in: Advances in Neural Information Processing Systems, 1995, pp. 231–238.
- [27] G. Brown, J. Wyatt, R. Harris, X. Yao, Diversity creation methods: a survey and categorisation, *Information Fusion* 6 (2005) 5–20.
- [28] S. Geman, E. Bienenstock, R. Doursat, Neural networks and the bias/variance dilemma, *Neural Computation* 4 (1992) 1–58.
- [29] G. Brown, J. L. Wyatt, P. Tiño, Managing diversity in regression ensembles, *Journal of Machine Learning Research* 6 (2005) 1621–1650.
- [30] D. Pedro, A unified bias-variance decomposition and its applications, in: 17th International Conference on Machine Learning, 2000, pp. 231–238.
- [31] D. H. Wolpert, On bias plus variance, *Neural Computation* 9 (1997) 1211–1243.
- [32] E. Kleinberg, Stochastic discrimination, *Annals of Mathematics and Artificial Intelligence* 1 (1990) 207–239.

- [33] R. E. Schapire, Y. Freund, P. Bartlett, W. S. Lee, Boosting the margin: A new explanation for the effectiveness of voting methods, *The annals of statistics* 26 (1998) 1651–1686.
- [34] V. Pisetta, *New Insights into Decision Trees Ensembles*, Ph.D. thesis, Lyon 2, 2012.
- [35] E. B. Kong, T. G. Dietterich, Error-correcting output coding corrects bias and variance, in: *Machine Learning Proceedings 1995*, Elsevier, 1995, pp. 313–321.
- [36] J. H. Friedman, On bias, variance, 0/1—loss, and the curse-of-dimensionality, *Data Mining and Knowledge Discovery* 1 (1997) 55–77.
- [37] L. Breiman, Arcing classifier (with discussion and a rejoinder by the author), *The Annals of Statistics* 26 (1998) 801–849.
- [38] G. M. James, Variance and bias for general loss functions, *Machine Learning* 51 (2003) 115–135.
- [39] L. Breiman, Bagging predictors, *Machine Learning* 24 (1996) 123–140.
- [40] L. Breiman, Bias, variance, and arcing classifiers (1996).
- [41] C.-X. Zhang, J.-S. Zhang, RotBoost: a technique for combining rotation forest and AdaBoost, *Pattern recognition letters* 29 (2008) 1524–1536.
- [42] Y. Freund, R. E. Schapire, Experiments with a new boosting algorithm, in: *icml*, volume 96, Citeseer, 1996, pp. 148–156.
- [43] I. Barandiaran, The random subspace method for constructing decision forests, *IEEE Trans. Pattern Anal. Mach. Intell* 20 (1998).
- [44] L. Breiman, Randomizing outputs to increase prediction accuracy, *Machine Learning* 40 (2000) 229–242.
- [45] A. Buja, W. Stuetzle, Smoothing effects of bagging, Preprint. AT&T Labs-Research (2000).
- [46] P. Bühlmann, B. Yu, Analyzing bagging, *The Annals of Statistics* 30 (2002) 927–961.
- [47] H.-C. Kim, S. Pang, H.-M. Je, D. Kim, S.-Y. Bang, Support vector machine ensemble with bagging, in: *International Workshop on Support Vector Machines*, Springer, 2002, pp. 397–408.
- [48] D. Tao, X. Tang, X. Li, X. Wu, Asymmetric bagging and random subspace for support vector machines-based relevance feedback in image retrieval, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 28 (2006) 1088–1099.
- [49] J. Mao, A case study on bagging, boosting and basic ensembles of neural networks for ocr, in: *1998 IEEE International Joint Conference on Neural Networks Proceedings. IEEE World Congress on Computational Intelligence (Cat. No. 98CH36227)*, volume 3, IEEE, 1998, pp. 1828–1833.
- [50] K. Ha, S. Cho, D. MacLachlan, Response models based on bagging neural networks, *Journal of Interactive Marketing* 19 (2005) 17–30.
- [51] R. Gençay, M. Qi, Pricing and hedging derivative securities with neural networks: Bayesian regularization, early stopping, and bagging, *IEEE Transactions on Neural Networks* 12 (2001) 726–734.
- [52] A. Khwaja, M. Naeem, A. Anpalagan, A. Venetsanopoulos, B. Venkatesh, Improved short-term load forecasting using bagged neural networks, *Electric Power Systems Research* 125 (2015) 109–115.
- [53] T. Hothorn, B. Lausen, A. Benner, M. Radespiel-Tröger, Bagging survival trees, *Statistics in Medicine* 23 (2004) 77–91.
- [54] R. F. Alvear-Sandoval, A. R. Figueiras-Vidal, On building ensembles of stacked denoising auto-encoding classifiers and their further improvement, *Information Fusion* 39 (2018) 41–52.

- [55] S. Hido, H. Kashima, Y. Takahashi, Roughly balanced bagging for imbalanced data, *Statistical Analysis and Data Mining: The ASA Data Science Journal* 2 (2009) 412–426.
- [56] J. Błaszczyński, J. Stefanowski, Neighbourhood sampling in bagging for imbalanced data, *Neurocomputing* 150 (2015) 529–542.
- [57] N. C. Oza, Online bagging and boosting, in: *2005 IEEE International Conference on Systems, Man and Cybernetics*, volume 3, Ieee, 2005, pp. 2340–2345.
- [58] J. H. Friedman, Greedy function approximation: a gradient boosting machine, *Annals of Statistics* (2001) 1189–1232.
- [59] P. Liu, S. Han, Z. Meng, Y. Tong, Facial expression recognition via a boosted deep belief network, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 1805–1812.
- [60] C. Cortes, M. Mohri, U. Syed, Deep boosting, in: *31st International Conference on Machine Learning, ICML 2014*, 2014.
- [61] V. Kuznetsov, M. Mohri, U. Syed, Multi-class deep boosting, *Advances in Neural Information Processing Systems* 3 (2014) 2501–2509.
- [62] S. Han, Z. Meng, A. S. Khan, Y. Tong, Incremental boosting convolutional neural network for facial action unit recognition, *Advances in Neural Information Processing Systems* 29 (2016).
- [63] M. Moghimi, S. J. Belongie, M. J. Saberian, J. Yang, N. Vasconcelos, L.-J. Li, Boosted convolutional neural networks., in: *BMVC*, 2016, pp. 24–1.
- [64] C. Chen, Z. Xiong, X. Tian, F. Wu, Deep boosting for image denoising, in: *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 3–18.
- [65] C. Chen, Z. Xiong, X. Tian, Z.-J. Zha, F. Wu, Real-world image denoising with deep boosting, *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2019).
- [66] B. Yang, J. Yan, Z. Lei, S. Z. Li, Convolutional channel features, in: *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 82–90.
- [67] E. Walach, L. Wolf, Learning to count with cnn boosting, in: *European Conference on Computer Vision*, Springer, 2016, pp. 660–676. doi:10.1007/978-3-319-46475-6_41.
- [68] M. Opitz, G. Waltner, H. Possegger, H. Bischof, Bier-boosting independent embeddings robustly, in: *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 5189–5198.
- [69] G. Waltner, M. Opitz, H. Possegger, H. Bischof, Hibster: Hierarchical boosted deep metric learning for image retrieval, in: *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, IEEE, 2019, pp. 599–608.
- [70] A. Mosca, G. D. Magoulas, Deep incremental boosting, *arXiv preprint arXiv:1708.03704* (2017).
- [71] W. Zhang, J. Jiang, Y. Shao, B. Cui, Snapshot boosting: a fast ensemble framework for deep neural networks, *Science China Information Sciences* 63 (2020) 112102.
- [72] C. Siu, Residual networks behave like boosting algorithms, in: *2019 IEEE International Conference on Data Science and Advanced Analytics (DSAA)*, IEEE, 2019, pp. 31–40.
- [73] C. Cortes, X. Gonzalvo, V. Kuznetsov, M. Mohri, S. Yang, Adanet: Adaptive structural learning of artificial neural networks (2017) 874–883.
- [74] F. Huang, J. Ash, J. Langford, R. Schapire, Learning deep resnet blocks sequentially using boosting theory (2018)

2058–2067.

- [75] A. Beygelzimer, E. Hazan, S. Kale, H. Luo, Online gradient boosting, *Advances in neural information processing systems* 28 (2015).
- [76] S. González, S. García, J. Del Ser, L. Rokach, F. Herrera, A practical tutorial on bagging and boosting based ensembles for machine learning: Algorithms, software tools, performance study, practical perspectives and opportunities, *Information Fusion* 64 (2020) 205–237.
- [77] D. H. Wolpert, Stacked generalization, *Neural Networks* 5 (1992) 241–259.
- [78] M. Leblanc, R. Tibshirani, Combining Estimates in Regression and Classification, *Journal of the American Statistical Association* 91 (1996) 1641–1650. doi:10.1080/01621459.1996.10476733.
- [79] L. Deng, D. Yu, Deep convex net: A scalable architecture for speech pattern classification, *Proceedings of the Annual Conference of the International Speech Communication Association, INTERSPEECH (2011)* 2285–2288.
- [80] L. L. Deng, D. Yu, J. Platt, Scalable stacking and learning for building deep architectures, *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings (2012)* 2133–2136. doi:10.1109/ICASSP.2012.6288333.
- [81] L. Deng, G. Tur, X. He, D. Hakkani-Tur, Use of kernel deep convex networks and end-to-end learning for spoken language understanding, in: *2012 IEEE Workshop on Spoken Language Technology, SLT 2012 - Proceedings, IEEE, 2012*, pp. 210–215. doi:10.1109/SLT.2012.6424224.
- [82] P.-S. Huang, L. Deng, M. Hasegawa-Johnson, X. He, Random features for Kernel Deep Convex Network, in: *2013 IEEE International Conference on Acoustics, Speech and Signal Processing, 2, IEEE, 2013*, pp. 3143–3147. doi:10.1109/ICASSP.2013.6638237.
- [83] T. Welchowski, M. Schmid, A framework for parameter estimation and model selection in kernel deep stacking networks, *Artificial Intelligence in Medicine* 70 (2016) 31–40. doi:10.1016/j.artmed.2016.04.002.
- [84] B. Hutchinson, L. Deng, D. Yu, A deep architecture with bilinear modeling of hidden representations: Applications to phonetic recognition, in: *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, 2012*, pp. 4805–4808. doi:10.1109/ICASSP.2012.6288994.
- [85] B. Hutchinson, L. L. Deng, D. Yu, Tensor deep stacking networks, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 35 (2013) 1944–1957. doi:10.1109/TPAMI.2012.268.
- [86] J. Li, H. Chang, J. Yang, Sparse deep stacking network for image classification, in: *Twenty-Ninth AAAI Conference on Artificial Intelligence, 2015*.
- [87] C. Sun, M. Ma, Z. Zhao, X. Chen, Sparse deep stacking network for fault diagnosis of motor, *IEEE Transactions on Industrial Informatics* 14 (2018) 3261–3270.
- [88] J. Li, H. Chang, J. Yang, W. Luo, Y. Fu, Visual representation and classification by learning group sparse deep stacking network, *IEEE Transactions on Image Processing* 27 (2017) 464–476.
- [89] J. Zhang, W. Zhang, R. Song, L. Ma, Y. Li, Grasp for stacking via deep reinforcement learning (2020) 2543–2549.
- [90] B. Wang, B. Xue, M. Zhang, Particle swarm optimisation for evolving deep neural networks for image classification by evolving and stacking transferable blocks, in: *2020 IEEE Congress on Evolutionary Computation (CEC), IEEE, 2020*, pp. 1–8.
- [91] H. Zhang, Y. Dai, H. Li, P. Koniusz, Deep stacked hierarchical multi-patch network for image deblurring, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019*, pp. 5978–5986.

- [92] H. Palangi, L. Deng, R. K. Ward, Recurrent Deep-Stacking Networks for sequence classification, 2014 IEEE China Summit and International Conference on Signal and Information Processing, IEEE ChinaSIP 2014 - Proceedings (2014) 510–514. doi:10.1109/ChinaSIP.2014.6889295.
- [93] H. Zhou, G. B. Huang, Z. Lin, H. Wang, Y. C. Soh, Stacked extreme learning machines, IEEE Transactions on Cybernetics 45 (2015) 2013–2025. doi:10.1109/TCYB.2014.2363492.
- [94] G. Wang, G. Zhang, K. S. Choi, J. Lu, Deep Additive Least Squares Support Vector Machines for Classification with Model Transfer, IEEE Transactions on Systems, Man, and Cybernetics: Systems 49 (2019) 1527–1540. doi:10.1109/TSMC.2017.2759090.
- [95] J. Wang, K. Feng, J. Wu, SVM-Based Deep Stacking Networks, Proceedings of the AAAI Conference on Artificial Intelligence 33 (2019) 5273–5280. doi:10.1609/aaai.v33i01.33015273.
- [96] X. Li, Y. Yang, H. Pan, J. Cheng, J. Cheng, A novel deep stacking least squares support vector machine for rolling bearing fault diagnosis, Computers in Industry 110 (2019) 36–47. doi:10.1016/j.compind.2019.05.005.
- [97] Z.-H. Zhou, J. Feng, Deep forest, arXiv preprint arXiv:1702.08835 (2017).
- [98] C.-Y. Low, J. Park, A. B.-J. Teoh, Stacking-based deep neural network: Deep analytic network for pattern classification, IEEE Transactions on Cybernetics 50 (2019) 5021–5034.
- [99] T. Kang, P. Chen, J. Quackenbush, W. Ding, A novel deep learning model by stacking conditional restricted boltzmann machine and deep neural network, in: Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, 2020, pp. 1316–1324.
- [100] Y. Liu, X. Yao, Ensemble learning via negative correlation, Neural Networks 12 (1999) 1399–1404. doi:10.1016/S0893-6080(99)00073-8.
- [101] Z. Shi, L. Zhang, Y. Liu, X. Cao, Y. Ye, M.-M. Cheng, G. Zheng, Crowd counting with deep negative correlation learning, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2018, pp. 5382–5390.
- [102] L. Zhang, Z. Shi, M.-M. Cheng, Y. Liu, J.-W. Bian, J. T. Zhou, G. Zheng, Z. Zeng, Nonlinear Regression via Deep Negative Correlation Learning, IEEE Transactions on Pattern Analysis and Machine Intelligence PP (2019) 1–1. doi:10.1109/tpami.2019.2943860.
- [103] S. Buschjäger, L. Pfahler, K. Morik, Generalized negative correlation learning for deep ensembling, arXiv preprint arXiv:2011.02952 (2020).
- [104] M. D. Muhlbaier, R. Polikar, An ensemble approach for incremental learning in nonstationary environments, in: M. Haindl, J. Kittler, F. Roli (Eds.), Multiple Classifier Systems, Springer Berlin Heidelberg, Berlin, Heidelberg, 2007, pp. 490–500.
- [105] K. Tang, M. Lin, F. L. Minku, X. Yao, Selective negative correlation learning approach to incremental learning, Neurocomputing 72 (2009) 2796 – 2805. doi:https://doi.org/10.1016/j.neucom.2008.09.022, hybrid Learning Machines (HAIS 2007) / Recent Developments in Natural Computation (ICNC 2007).
- [106] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, R. Salakhutdinov, B. Mele, G. Altarelli, Dropout: a simple way to prevent neural networks from overfitting, The journal of machine learning research 15 (2014) 1929–1958. doi:10.1016/0370-2693(93)90272-J.
- [107] L. Wan, M. Zeiler, S. Zhang, Y. L. Cun, R. Fergus, Regularization of Neural Networks using DropConnect, in: S. Dasgupta, D. McAllester (Eds.), Proceedings of the 30th International Conference on Machine Learning,

- volume 28 of *Proceedings of Machine Learning Research*, PMLR, Atlanta, Georgia, USA, 2013, pp. 1058–1066. doi:10.1109/TPAMI.2017.2703082.
- [108] G. Huang, Y. Sun, Z. Liu, D. Sedra, K. Q. Weinberger, Deep networks with stochastic depth, in: European Conference on Computer Vision, Springer, 2016, pp. 646–661. doi:10.1007/978-3-319-46493-0_39.
- [109] S. Singh, D. Hoiem, D. Forsyth, Swapout: Learning an ensemble of deep architectures (2016) 28–36.
- [110] G. Hinton, O. Vinyals, J. Dean, Distilling the knowledge in a neural network, arXiv preprint arXiv:1503.02531 (2015).
- [111] L. N. Smith, E. M. Hand, T. Doster, Gradual dropin of layers to train very deep neural networks, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 4763–4771.
- [112] G. Huang, Y. Li, G. Pleiss, Z. Liu, J. E. Hopcroft, K. Q. Weinberger, Snapshot ensembles: Train 1, get M for free, arXiv preprint arXiv:1704.00109 (2017).
- [113] Y.-H. Pao, G.-H. Park, D. J. Sobajic, Learning and generalization characteristics of the random vector functional-link net, *Neurocomputing* 6 (1994) 163–180.
- [114] A. K. Malik, R. Gao, M. A. Ganaie, M. Tanveer, P. N. Suganthan, Random vector functional link network: recent developments, applications, and future directions, arXiv preprint arXiv:2203.11316 (2022).
- [115] Q. Shi, R. Katuwal, P. N. Suganthan, M. Tanveer, Random vector functional link neural network based ensemble deep learning, *Pattern Recognition* 117 (2021) 107978.
- [116] B. Han, J. Sim, H. Adam, Branchout: Regularization for online ensemble tracking with convolutional neural networks, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 3356–3365.
- [117] J. Xue, Z. Wang, D. Kong, Y. Wang, X. Liu, W. Fan, S. Yuan, S. Niu, D. Li, Deep ensemble neural-like p systems for segmentation of central serous chorioretinopathy lesion, *Information Fusion* 65 (2021) 84–94.
- [118] J. Xie, B. Xu, Z. Chuang, Horizontal and vertical ensemble with deep representation for classification, arXiv preprint arXiv:1306.2759 (2013).
- [119] S. Laine, T. Aila, Temporal ensembling for semi-supervised learning, arXiv preprint arXiv:1610.02242 (2016).
- [120] D. Ciregan, U. Meier, J. Schmidhuber, Multi-column deep neural networks for image classification, in: 2012 IEEE Conference on Computer Vision and Pattern Recognition, IEEE, 2012, pp. 3642–3649.
- [121] F. Grassmann, J. Mengelkamp, C. Brandl, S. Harsch, M. E. Zimmermann, B. Linkohr, A. Peters, I. M. Heid, C. Palm, B. H. Weber, A deep learning algorithm for prediction of age-related eye disease study severity scale for age-related macular degeneration from color fundus photography, *Ophthalmology* 125 (2018) 1410–1420.
- [122] W. Li, S. Ding, Y. Chen, S. Yang, Heterogeneous ensemble for default prediction of peer-to-peer lending in china, *IEEE Access* 6 (2018) 54396–54406.
- [123] Z. H. Kilimci, S. Akyokus, Deep learning-and word embedding-based heterogeneous classifier ensembles for text classification, *Complexity* 2018 (2018).
- [124] S. Tabik, R. F. Alvear-Sandoval, M. M. Ruiz, J.-L. Sancho-Gómez, A. R. Figueiras-Vidal, F. Herrera, Mnist-net10: A heterogeneous deep networks fusion based on the degree of certainty to reach 0.1% error rate. ensembles overview and proposal, *Information Fusion* 62 (2020) 73 – 80. doi:10.1016/j.inffus.2020.04.002.
- [125] B. Seijo-Pardo, I. Porto-Díaz, V. Bolón-Canedo, A. Alonso-Betanzos, Ensemble feature selection: homogeneous and heterogeneous approaches, *Knowledge-Based Systems* 118 (2017) 124–139.

- [126] Q. L. Zhao, Y. H. Jiang, M. Xu, Incremental learning by heterogeneous bagging ensemble, in: L. Cao, J. Zhong, Y. Feng (Eds.), *Advanced Data Mining and Applications*, Springer Berlin Heidelberg, Berlin, Heidelberg, 2010, pp. 1–12.
- [127] L. L. Minku, A. P. White, X. Yao, The impact of diversity on online ensemble learning in the presence of concept drift, *IEEE Transactions on Knowledge and Data Engineering* 22 (2009) 730–742.
- [128] X. Qiu, P. N. Suganthan, G. A. Amaratunga, Ensemble incremental learning random vector functional link network for short-term electric load forecasting, *Knowledge-Based Systems* 145 (2018) 182 – 196. doi:<https://doi.org/10.1016/j.knsys.2018.01.015>.
- [129] G. Grmanová, P. Laurinec, V. Rozinajová, A. B. Ezzeddine, M. Lucká, P. Lacko, P. Vrablecová, P. Návrát, Incremental ensemble learning for electricity load forecasting, volume 13, 2016, pp. 97–117.
- [130] F. Duan, L. Dai, Recognizing the gradual changes in sEMG characteristics based on incremental learning of wavelet neural network ensemble, *IEEE Transactions on Industrial Electronics* 64 (2017) 4276–4286. doi:10.1109/TIE.2016.2593693.
- [131] R. T. Das, K. K. Ang, C. Quek, Ierspop: A novel incremental rough set-based pseudo outer-product with ensemble learning, *Applied Soft Computing* 46 (2016) 170 – 186. doi:<https://doi.org/10.1016/j.asoc.2016.04.015>.
- [132] C. Ju, A. Bibaut, M. van der Laan, The relative performance of ensemble methods with deep convolutional neural networks for image classification, *Journal of Applied Statistics* 45 (2018) 2800–2818.
- [133] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, *arXiv preprint arXiv:1409.1556* (2014).
- [134] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, A. Rabinovich, Going deeper with convolutions, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.
- [135] L. I. Kuncheva, C. J. Whitaker, C. A. Shipp, R. P. Duin, Limits on the majority vote accuracy in classifier fusion, *Pattern Analysis & Applications* 6 (2003) 22–31.
- [136] A. Choromanska, M. Henaff, M. Mathieu, G. B. Arous, Y. LeCun, The loss surfaces of multilayer networks, in: *Artificial Intelligence and Statistics*, 2015, pp. 192–204.
- [137] S. Li, X. Lu, S. Sakai, M. Mimura, T. Kawahara, Semi-supervised ensemble DNN acoustic model training, in: *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, IEEE, 2017, pp. 5270–5274.
- [138] X. Wang, D. Kihara, J. Luo, G.-J. Qi, ENAET: self-trained ensemble autoencoding transformations for semi-supervised learning, *arXiv preprint arXiv:1911.09265* (2019).
- [139] J. Liu, S. Zhao, G. Wang, Ssel-ade: a semi-supervised ensemble learning framework for extracting adverse drug events from social media, *Artificial Intelligence in Medicine* 84 (2018) 34–49.
- [140] J. Li, S. Wu, C. Liu, Z. Yu, H.-S. Wong, Semi-supervised deep coupled ensemble learning with classification landmark exploration, *IEEE Transactions on Image Processing* 29 (2019) 538–550.
- [141] G. Pio, D. Malerba, D. D’Elia, M. Ceci, Integrating microrna target predictions for the discovery of gene regulatory networks: a semi-supervised ensemble learning approach, *BMC Bioinformatics* 15 (2014) S4.
- [142] G. Tsoumakas, I. Katakis, Multi-label classification: An overview, *International Journal of Data Warehousing*

- and Mining (IJDM) 3 (2007) 1–13.
- [143] G. Tsoumakas, I. Vlahavas, Random k-labelsets: An ensemble method for multilabel classification, in: *European Conference on Machine Learning*, Springer, 2007, pp. 406–417.
- [144] J. M. Moyano, E. L. Gibaja, K. J. Cios, S. Ventura, An evolutionary approach to build ensembles of multi-label classifiers, *Information Fusion* 50 (2019) 168–180.
- [145] K. Kimura, M. Kudo, L. Sun, S. Koujaku, Fast random k-labelsets for large-scale multi-label classification, in: *2016 23rd International Conference on Pattern Recognition (ICPR)*, IEEE, 2016, pp. 438–443.
- [146] R. Wang, S. Kwong, X. Wang, Y. Jia, Active k-labelsets ensemble for multi-label classification, *Pattern Recognition* 109 (2021) 107583.
- [147] C. Shi, X. Kong, P. S. Yu, B. Wang, Multi-label ensemble learning, in: D. Gunopulos, T. Hofmann, D. Malerba, M. Vazirgiannis (Eds.), *Machine Learning and Knowledge Discovery in Databases*, Springer Berlin Heidelberg, Berlin, Heidelberg, 2011, pp. 223–239.
- [148] P. Li, H. Li, M. Wu, Multi-label ensemble based on variable pairwise constraint projection, *Information Sciences* 222 (2013) 269 – 281. doi:<https://doi.org/10.1016/j.ins.2012.07.066>, including Special Section on New Trends in Ambient Intelligence and Bio-inspired Systems.
- [149] Y. Xia, K. Chen, Y. Yang, Multi-label classification with weighted classifier selection and stacked ensemble, *Information Sciences* 557 (2021) 421–442.
- [150] X. Guo, F. Liu, Y. Ju, Z. Wang, C. Wang, Human protein subcellular localization with integrated source and multi-label ensemble classifier, *Scientific Reports* 6 (2016) 28087.
- [151] G. Yu, C. Domeniconi, H. Rangwala, G. Zhang, Z. Yu, Transductive multi-label ensemble classification for protein function prediction, in: *Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2012, pp. 1077–1085.
- [152] W. Zhang, F. Liu, L. Luo, J. Zhang, Predicting drug side effects by multi-label learning and ensemble learning, *BMC bioinformatics* 16 (2015) 365.
- [153] L. Schietgat, C. Vens, J. Struyf, H. Blockeel, D. Kocev, S. Džeroski, Predicting gene function using hierarchical multi-label decision tree ensembles, *BMC Bioinformatics* 11 (2010) 2.
- [154] J. Read, B. Pfahringer, G. Holmes, E. Frank, Classifier chains for multi-label classification, *Machine Learning* 85 (2011) 333.
- [155] G. Chen, D. Ye, Z. Xing, J. Chen, E. Cambria, Ensemble application of convolutional and recurrent neural networks for multi-label text categorization, in: *2017 international joint conference on neural networks (IJCNN)*, IEEE, 2017, pp. 2377–2383.
- [156] L. Breiman, Stacked regressions, *Machine Learning* 24 (1996) 49–64.
- [157] M. J. Van der Laan, E. C. Polley, A. E. Hubbard, Super learner, *Statistical Applications in Genetics and Molecular Biology* 6 (2007).
- [158] C. Ju, M. Combs, S. D. Lendle, J. M. Franklin, R. Wyss, S. Schneeweiss, M. J. van der Laan, Propensity score prediction for electronic healthcare databases using super learner and high-dimensional propensity score methods, *Journal of Applied Statistics* 46 (2019) 2216–2236.
- [159] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, M. Riedmiller, Playing atari with deep reinforcement learning, *arXiv preprint arXiv:1312.5602* (2013).

- [160] X.-l. Chen, L. Cao, C.-x. Li, Z.-x. Xu, J. Lai, Ensemble network architecture for deep reinforcement learning, *Mathematical Problems in Engineering* 2018 (2018).
- [161] H. Cuayáhuitl, D. Lee, S. Ryu, Y. Cho, S. Choi, S. Indurthi, S. Yu, H. Choi, I. Hwang, J. Kim, Ensemble-based deep reinforcement learning for chatbots, *Neurocomputing* 366 (2019) 118–130.
- [162] S. Carta, A. Corriga, A. Ferreira, A. S. Podda, D. R. Recupero, A multi-layer and multi-ensemble stock trader using deep learning and deep reinforcement learning, *Applied Intelligence* (2020) 1–17.
- [163] H. Yang, X.-Y. Liu, S. Zhong, A. Walid, Deep reinforcement learning for automated stock trading: An ensemble strategy (2020) 1–8.
- [164] K.-F. Tang, H.-C. Kao, C.-N. Chou, E. Y. Chang, Inquire and diagnose: Neural symptom checking ensemble using deep reinforcement learning, in: *Proceedings of NIPS Workshop on Deep Reinforcement Learning*, 2016.
- [165] Y. Şenbabaoğlu, G. Michailidis, J. Z. Li, Critical limitations of consensus clustering in class discovery, *Scientific Reports* 4 (2014) 1–13.
- [166] S. Bian, W. Wang, On diversity and accuracy of homogeneous and heterogeneous ensembles, *International Journal of Hybrid Intelligent Systems* 4 (2007) 103–128.
- [167] Z.-H. Zhou, W. Tang, Clusterer ensemble, *Knowledge-Based Systems* 19 (2006) 77–83.
- [168] D. Huang, C.-D. Wang, J.-H. Lai, Locally weighted ensemble clustering, *IEEE Transactions on Cybernetics* 48 (2017) 1460–1473.
- [169] L. Zheng, T. Li, C. Ding, Hierarchical ensemble clustering, in: *2010 IEEE International Conference on Data Mining*, IEEE, 2010, pp. 1199–1204.
- [170] D. Huang, J. Lai, C.-D. Wang, Ensemble clustering using factor graph, *Pattern Recognition* 50 (2016) 131–142.
- [171] H. Liu, T. Liu, J. Wu, D. Tao, Y. Fu, Spectral ensemble clustering, in: *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2015, pp. 715–724.
- [172] H. Liu, M. Shao, S. Li, Y. Fu, Infinite ensemble for image clustering, in: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2016, pp. 1745–1754.
- [173] U. Shaham, X. Cheng, O. Dror, A. Jaffe, B. Nadler, J. Chang, Y. Kluger, A deep learning approach to unsupervised ensemble learning, in: *International Conference on Machine Learning*, 2016, pp. 30–39.
- [174] N. Alami, M. Meknassi, N. En-nahnahi, Enhancing unsupervised neural networks based text summarization with word embedding and ensemble learning, *Expert systems with applications* 123 (2019) 195–211.
- [175] B. Hassan, S. E. Abdelrahman, R. Bahgat, I. Farag, Uests: An unsupervised ensemble semantic textual similarity method, *IEEE Access* 7 (2019) 85462–85482.
- [176] E. Ahn, A. Kumar, D. Feng, M. Fulham, J. Kim, Unsupervised feature learning with k-means and an ensemble of deep convolutional neural networks for medical image classification, *arXiv preprint arXiv:1906.03359* (2019).
- [177] A. Lahiri, A. G. Roy, D. Sheet, P. K. Biswas, Deep neural ensemble for retinal vessel segmentation in fundus images towards achieving label-free angiography, in: *2016 38th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, IEEE, 2016, pp. 1340–1343.
- [178] B. Liu, L. Gu, F. Lu, Unsupervised ensemble strategy for retinal vessel segmentation, in: *International Conference on Medical Image Computing and Computer-Assisted Intervention*, Springer, 2019, pp. 111–119.
- [179] Z. Cao, X. Pan, Y. Yang, Y. Huang, H.-B. Shen, The Inclocator: a subcellular localization predictor for long non-coding RNAs based on a stacked ensemble classifier, *Bioinformatics* 34 (2018) 2185–2194.

- [180] M. Hu, P. Suganthan, Representation learning using deep random vector functional link networks for clustering, *Pattern Recognition* (2022) 108744.
- [181] P. Melville, R. J. Mooney, Constructing diverse classifier ensembles using artificial training examples, in: *IJCAI*, volume 3, 2003, pp. 505–510.
- [182] P. Melville, R. J. Mooney, Diverse ensembles for active learning, in: *Proceedings of the twenty-first international conference on Machine learning*, 2004, p. 74.
- [183] W. H. Beluch, T. Genewein, A. Nürnberger, J. M. Köhler, The power of ensembles for active learning in image classification, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 9368–9377.
- [184] A. Sharma, R. Rani, Be-dti': Ensemble framework for drug target interaction prediction using dimensionality reduction and active learning, *Computer Methods and Programs in Biomedicine* 165 (2018) 151–162.
- [185] H. Zhang, W. Liu, J. Shan, Q. Liu, Online active learning paired ensemble for concept drift and class imbalance, *IEEE Access* 6 (2018) 73815–73828. doi:10.1109/ACCESS.2018.2882872.
- [186] X. Zheng, J. Shen, C. Cox, J. C. Wakefield, M. G. Ehm, M. R. Nelson, B. S. Weir, HIBAG—HLA genotype imputation with attribute bagging, *The Pharmacogenomics Journal* 14 (2014) 192–200.
- [187] I. Cabria, I. Gondra, Mri segmentation fusion for brain tumor detection, *Information Fusion* 36 (2017) 1–9.
- [188] A. Shalhaf, M. Vafaeezadeh, Automated detection of COVID-19 using ensemble of transfer learning with deep convolutional neural network based on CT scans, *International journal of computer assisted radiology and surgery* 16 (2021) 115–123.
- [189] F. Ali, S. El-Sappagh, S. R. Islam, D. Kwak, A. Ali, M. Imran, K.-S. Kwak, A smart healthcare monitoring system for heart disease prediction based on ensemble deep learning and feature fusion, *Information Fusion* 63 (2020) 208–222.
- [190] T. Zhou, H. Lu, Z. Yang, S. Qiu, B. Huo, Y. Dong, The ensemble deep learning model for novel COVID-19 on CT images, *Applied Soft Computing* 98 (2021) 106885.
- [191] Y. Li, Y. Song, L. Jia, S. Gao, Q. Li, M. Qiu, Intelligent fault diagnosis by fusing domain adversarial training and maximum mean discrepancy via ensemble learning, *IEEE Transactions on Industrial Informatics* 17 (2020) 2833–2841.
- [192] A. K. Das, S. Ghosh, S. Thunder, R. Dutta, S. Agarwal, A. Chakrabarti, Automatic covid-19 detection from x-ray images using ensemble learning with convolutional neural network, *Pattern Analysis and Applications* 24 (2021) 1111–1124.
- [193] S. Sukegawa, A. Fujimura, A. Taguchi, N. Yamamoto, A. Kitamura, R. Goto, K. Nakano, K. Takabatake, H. Kawai, H. Nagatsuka, Y. Furuki, Identification of osteoporosis using ensemble deep learning model with panoramic radiographs and clinical covariates, *Scientific reports* 12 (2022) 1–10.
- [194] Z. Gao, L. Wang, R. Sorousmehr, A. Wood, J. Gryak, B. Nallamothu, K. Najarian, Vessel segmentation for x-ray coronary angiography using ensemble methods with deep learning and filter-based features, *BMC Medical Imaging* 22 (2022) 1–17.
- [195] A. Rath, D. Mishra, G. Panda, S. C. Satapathy, K. Xia, Improved heart disease detection from ECG signal using deep learning based ensemble model, *Sustainable Computing: Informatics and Systems* 35 (2022) 100732.
- [196] M. Tanveer, A. Rashid, M. Ganaie, M. Reza, I. Razzak, K.-L. Hua, Classification of Alzheimer's disease using

- ensemble of deep neural networks trained through transfer learning, *IEEE Journal of Biomedical and Health Informatics* 26 (2022) 1453 – 1463. doi:10.1109/JBHI.2021.3083274.
- [197] H. M. Rai, K. Chatterjee, Hybrid cnn-lstm deep learning model and ensemble technique for automatic detection of myocardial infarction using big ecg data, *Applied Intelligence* 52 (2022) 5366–5384.
- [198] M. Ganaie, M. Tanveer, Ensemble deep random vector functional link network using privileged information for Alzheimer’s disease diagnosis, *IEEE/ACM Transactions on Computational Biology and Bioinformatics* (2022). doi:10.1109/TCBB.2022.3170351.
- [199] G. Tur, L. Deng, D. Hakkani-Tür, X. He, Towards deeper understanding: Deep convex networks for semantic utterance classification, in: 2012 IEEE international conference on acoustics, speech and signal processing (ICASSP), IEEE, 2012, pp. 5045–5048.
- [200] L. Deng, G. Tur, X. He, D. Hakkani-Tur, Use of kernel deep convex networks and end-to-end learning for spoken language understanding, in: 2012 IEEE Spoken Language Technology Workshop (SLT), IEEE, 2012, pp. 210–215.
- [201] L. Deng, J. C. Platt, Ensemble deep learning for speech recognition, in: Fifteenth Annual Conference of the International Speech Communication Association, 2014.
- [202] M. Amin-Naji, A. Aghagolzadeh, M. Ezoji, Ensemble of cnn for multi-focus image fusion, *Information Fusion* 51 (2019) 201–214.
- [203] X. Qiu, L. Zhang, Y. Ren, P. N. Suganthan, G. Amaratunga, Ensemble deep learning for regression and time series forecasting, in: 2014 IEEE Symposium on Computational Intelligence in Ensemble Learning (CIEL), IEEE, 2014, pp. 1–6.
- [204] X. Qiu, Y. Ren, P. N. Suganthan, G. A. Amaratunga, Empirical mode decomposition based ensemble deep learning for load demand time series forecasting, *Applied Soft Computing* 54 (2017) 246–255.
- [205] F. Liu, F. Xu, S. Yang, A flood forecasting model based on deep learning algorithm via integrating stacked autoencoders with BP neural network, in: 2017 IEEE Third International Conference on Multimedia Big Data (BigMM), IEEE, 2017, pp. 58–61.
- [206] N. Bhusal, R. M. Shukla, M. Gautam, M. Benidris, S. Sengupta, Deep ensemble learning-based approach to real-time power system state estimation, *International Journal of Electrical Power & Energy Systems* 129 (2021) 106806.
- [207] P. Singla, M. Duhan, S. Saroha, An ensemble method to forecast 24-h ahead solar irradiance using wavelet decomposition and BiLSTM deep learning network, *Earth Science Informatics* (2021) 1–16.
- [208] L. Deng, X. He, J. Gao, Deep stacking networks for information retrieval, in: 2013 IEEE International Conference on Acoustics, Speech and Signal Processing, IEEE, 2013, pp. 3153–3157.
- [209] G. Wang, J. Sun, J. Ma, K. Xu, J. Gu, Sentiment classification: The contribution of ensemble learning, *Decision Support Systems* 57 (2014) 77–93.
- [210] M. Zareapoor, P. Shamsolmoali, Application of credit card fraud detection: Based on bagging ensemble classifier, *Procedia Computer Science* 48 (2015) 679–685.
- [211] Z. Yin, M. Zhao, Y. Wang, J. Yang, J. Zhang, Recognition of emotions using multimodal physiological signals and an ensemble deep learning model, *Computer Methods and Programs in Biomedicine* 140 (2017) 93–110.
- [212] S. Kazemi, B. Minaei Bidgoli, S. Shamshirband, S. M. Karimi, M. A. Ghorbani, K.-w. Chau, R. Kazem Pour,

- Novel genetic-based negative correlation learning for estimating soil temperature, *Engineering Applications of Computational Fluid Mechanics* 12 (2018) 506–516.
- [213] K. Randhawa, C. K. Loo, M. Seera, C. P. Lim, A. K. Nandi, Credit card fraud detection using adaboost and majority voting, *IEEE Access* 6 (2018) 14277–14284.
- [214] A. Martín, R. Lara-Cabrera, D. Camacho, Android malware detection through hybrid features fusion and ensemble classifiers: The andropytool framework and the omnidroid dataset, *Information Fusion* 52 (2019) 128–142.
- [215] F. Wang, D. Jiang, H. Wen, H. Song, Adaboost-based security level classification of mobile intelligent terminals, *The Journal of Supercomputing* 75 (2019) 7460–7478.
- [216] W. Chen, B. Pradhan, S. Li, H. Shahabi, H. M. Rizeei, E. Hou, S. Wang, Novel hybrid integration approach of bagging-based fisher's linear discriminant function for groundwater potential analysis, *Natural Resources Research* 28 (2019) 1239–1258.
- [217] E. Tsogbaatar, M. H. Bhuyan, Y. Taenaka, D. Fall, K. Gonchigsumlaa, E. Elmroth, Y. Kadobayashi, Del-iot: A deep ensemble learning approach to uncover anomalies in iot, *Internet of Things* 14 (2021) 100391.
- [218] L. Wen, X. Xie, X. Li, L. Gao, A new ensemble convolutional neural network with diversity regularization for fault diagnosis, *Journal of Manufacturing Systems* 62 (2020) 964 – 971. doi:10.1016/j.jmsy.2020.12.002.
- [219] Z.-H. Zhou, N. V. Chawla, Y. Jin, G. J. Williams, Big data opportunities and challenges: Discussions from data analytics perspectives [discussion forum], *IEEE Computational Intelligence Magazine* 9 (2014) 62–74.
- [220] R. Katuwal, P. N. Suganthan, Stacked autoencoder based deep random vector functional link neural network for classification, *Applied Soft Computing* 85 (2019) 105854.