

Automatic Digital Analysis System to Grade Diabetic Retinopathy by Integrated Stacking Model Concept



Thangavel Manivel^{1*}, Umathurai Saravanakumar²

¹ Department of Information Technology, Muthayammal Engineering College, Rasipuram 637408, India ² Department of Electronics and Communications Engineering, Muthayammal Engineering College, Rasipuram 637408, India

Corresponding Author Email: manivel.t.it@mec.edu.in

Copyright: ©2024 The authors. This article is published by IIETA and is licensed under the CC BY 4.0 license (http://creativecommons.org/licenses/by/4.0/).

https://doi.org/10.18280/ts.410643

ABSTRACT

Received: 14 November 2023 Revised: 1 June 2024 Accepted: 25 September 2024 Available online: 31 December 2024

Keywords:

diabetic retinopathy, image analysis, deep learning, Integrated Stacking Model (ISM), fundus image, soft computing

Ophthalmic diagnosis is based primarily on visual information (photographs of the retina), and the recent availability of digital fundus images allows such quantization of parameters or classification to be derived from computerized image processing. Thus, an Automated Digital Analysis System (ADAS) is designed for grading Diabetic Retinopathy (DR) in this work. The Integrated Stacking Model (ISM) concept is employed to design a single large multi-headed deep learning model. Each sub-model has three layers with a predefined number of 3×3 convolution filters (128, 256 and 512) and a pooling layer to abstract deep features. The outputs of sub-models are fed to the meta-learner for grading the DR. The performance of the ADAS for grading DR is evaluated using MESSIDOR-1 and Kaggle datasets. For the images in MESSIDOR-1, the proposed system is considered a four-class problem, and it is a five-class problem for Kaggle dataset images. Results show that the proposed ISM-DR classification system provides promising results with an average classification accuracy of 99.2% (four-class) and 99.1% (five-class) using MESSIDOR-1 and Kaggle dataset images respectively when stacking five sub-models. A clear trend can be seen from the experimental results towards better performance when more sub-models are stacked.

1. INTRODUCTION

The eye is subject to several pathologies which affect it in different ways. Accurate classification of disease is an important measure of its current state, and contributes to an assessment of temporal progression, both of which can determine when surgical interventions should take place. Particular attention is given to DR, due to its high prevalence amongst diabetics. The classification of DR is important in assessment of the patient but is difficult due to the timeconsuming processes involved. In clinical practice therefore, qualitative analysis is the norm, although several research programmes have developed to grade the DR from 0 (low) to 3 (high).

A wide variety of algorithms have been created using knearest neighbor, random forest, support vector machine, artificial neural networks, decision trees, and naive Bayes. These approaches rely on the textural characteristics and their performances depend on the retrieved features and the classifiers. The salient attributes may be compromised because of the inadequate classifier, and vice versa. Furthermore, the establishment of a precise classification system requires the use of two separate modules. Recent soft computing approaches for the diagnosis of various diseases do not require segmentation and directly work on the image characteristic to classify the abnormality [1-3]. Though these algorithms can be effective in some specific applications, they have a limited perspective on the data and are prone to overfitting. They are also unable to capture complex patterns particularly in medical images and may lack robustness. To overcome these difficulties, a single large multi-headed deep learning model is designed based on ISM concept. It is a powerful technique capable of providing accurate prediction.

The ISM for grading DR is a machine learning technique that combines various submodels to improve predictive accuracy. This involves training multiple individual models using different algorithms or neural network architectures on the same dataset. The predictions from these models are then combined using a higher-level model (the meta-learner) to generate more accurate predictions. The ISM aims to utilize the different perspectives of individual models to produce more robust and accurate results. By allowing the meta-learner to learn the best way to combine the submodels' predictions, the ISM reduces the weaknesses of each model.

2. LITERATURE SURVEY

Convolutional Neural Networks (CNNs) are becoming more popular in the field of medical image processing [4]. The most recent techniques of state-of-the-art detection and classification of DR color fundus images utilizing deep learning algorithms have been studied and analyzed. Challenges with automated disease diagnosis in medical images with CNNs, and pre-processing methods, training, performance metrics for unbalanced classes are discussed in study [5]. The objective of grading and extracting diabetic retinopathy severity stages using an optimized deep learning framework is provided in study [6]. It involves several processes, such as backdrop segmentation, feature set extraction, feature optimization via the use of cuckoo search, and CNN for severity grade classification.

DR lesions are detected automatically and classified in study [7]. The first step involves collecting data on diabetic retinopathy from a variety of sources into a single pool for use in the subsequent stages. Lesions may be identified with the use of Faster Region CNN (FRCNN) and then categorized via the use of the transfer learning and attention. A Deep Graph Convolutional Network (DGCN) is used in study [8] to exploit inherent correlations from separate retinal image characteristics learnt by CNN. The optimization and the use of the DGCN in an automated task of DR grading is constrained by the presence of three defined loss functions: a graph-center function, a transformation-invariant function, and a pseudocontrastive function.

A domain adaptation method called Multi-Model Domain Adaptation (MMDA) is discussed in study [9] for DR classification. A weight mechanism is introduced into the MMDA by evaluating the relevance of each source domain, and a weighted pseudo-labeling approach is connected to the source feature extractors for training the target DR classification model. A comparison of the performances of deep-learning methods is discussed in study [10] for automated DR detection. Multiple methods are used to categorize retinal images in study [11] such as, training CNN models to learn the characteristics that correspond to each grade and identifying and segmenting lesions by making use of information regarding their location. An algorithm which can perform automatic detection and classification of DR images is discussed in study [12] based on deep learning. The image must initially go through a preprocessing step and then a histogram-based segmentation followed by a technique called synergic deep learning to classify DR fundus images. A variety of pre-trained CNN models is described for DR classification in study [13]. It also provides two CNN architectures for binary and multiclass classification. A hybrid model for diagnosing DR is discussed in study [14] using fundus images. The combination of morphological image processing with Inception v3 allows for the classification of DR.

A method for the automated diagnosis of DR using fuzzy image processing method is discussed in study [15]. A combination of fuzzy technique and the circular Hough transform is utilized. A machine learning strategy based on the human extraction of features and the automated extraction of features using a CNN is described in study [16]. Recent developments in deep learning, such as capsule networks, and their significant advantages over conventional machine learning methods in a variety of applications inspired the researchers to utilize them for the detection of diabetic retinopathy.

A system which utilizes the multiple instance learning paradigm to classify DR from fundus images is discussed in study [17]. It extracts local information from rectangular image patches independently, and then combines it in an efficient manner through an attention mechanism. The resultant final image representation is suitable for DR classification. In addition to attention mechanism, heatmaps are constructed to show the abnormal regions. Table 1 shows a comprehensive analysis of the proposed ISM-DR classification system with existing systems.

Architecture	Model Type	Segmentation	Classification	Performance Measures	Databases Used
InceptionResNet [5]	Single	No	Binary	Accuracy, Precision, Recall	EyePACS, MESSIDOR-1 & 2
Optimized CNN [6]	Single	No	Binary	Accuracy, Sensitivity, Specificity, Precision	MESSIDOR-1 & IDRiD
Faster RCNN + CNN [7]	Single	Yes	Binary	Accuracy, Sensitivity, Specificity, Area under the curve	MESSIDOR-1, Kaggle dataset
DGCN [8]	Single	No	Multiclass	Accuracy, Sensitivity, Specificity,	EyePACS, MESSIDOR -2
MMDA [9]	Multimodal	No	Binary	Accuracy, Sensitivity	DDR, IDRiD, MESSIDOR, and MESSIDOR -2
CNN [10]	Single	No	Multiclass	Accuracy, Sensitivity, Specificity, Positive and Negative Predictive value	MESSIDOR -2
Synergic CNN [12]	Single	Yes	Multiclass	Accuracy, Sensitivity, Specificity	MESSIDOR -2
CNN [13]	Single	No	Multiclass	Accuracy, Sensitivity, Precision, F1 Score	EyePACS, MESSIDOR-1 & 2
Inception V3 [14]	Hybrid	Yes	Multiclass	Accuracy, Sensitivity, Specificity	MESSIDOR-1, Kaggle dataset
Capsule Network [16]	Single	No	Multiclass	Accuracy, Precision, Recall, F1 Score	MESSIDOR -1
Proposed ISM-DR	Hybrid (Integrated Stacking Model)	No	Multiclass	Accuracy, Precision, Recall	MESSIDOR-1, Kaggle dataset

Table 1. Comprehensive analysis of the Proposed ISM-DR classification system with existing systems

3. MATERIALS AND METHODS

When employing neural networks as submodels, it can be advantageous to employ a neural network as a meta-learner. In this context, the sub-networks can be integrated into a larger multi-headed neural network, which combines predictions from each individual sub-model. This configuration transforms the stacking ensemble into a unified extensive model. This strategy offers the advantage of directly providing the outputs of the submodels to the meta-learner [18]. Additionally, if required, it becomes feasible to concurrently adjust the weights of the submodels alongside the meta-learner model. Figure 1 shows the proposed ISM-DR classification system with five sub-models. Brief descriptions of each layer in the proposed ISM-DR classification system are discussed in the following sub-sections.



Figure 1. Proposed ISM-DR classification system

3.1 Convolution and max pooling layers

The convolutional layer is the primary constituent of CNN. It basically does a dot product, or matrix multiplication, between two matrices [19]. The receptive field, or specific area of the input image, is the first matrix. The second matrix has a set of learnable parameters, called kernel or convolution filter. The method used to perform this matrix multiplication (dot product) involves swiping the kernel across various image regions, namely the height and breadth of the image. Although the kernel's breadth and height are smaller than those of the image, its spatial size is still equal to the image's depth. For example, if an image contains three (RGB) channels, the depth will span into R, G, and B channels, but the height and width of the input image.

To create an image representation of every receptive region, the kernel traverses all image components, going across the height and width of the image. The kernel's response at each spatial location in the image produces a two-dimensional feature map. The kernel's stepping over the image is determined by a hyperparameter called stride (slide size). Extraction of high-level characteristics from input images is the aim of the convolution process. A CNN design should ideally consist of several convolutional layers. Low-level features like edges are mostly produced by the first convolutional layer. The network then acquires high-level characteristics by adding further convolutional layers. An example of a 2D convolution operation that happens when an input image is subjected to a convolution operation is shown in Figure 2.

The pooling layer is the subsequent common layer which comes after the convolutional layer. The spatial size of feature map generated by the convolution layer is mostly reduced by the pooling layer. Through dimensionality reduction, this step helps to lower the amount of computing power required for data processing. Additionally, it helps to capture dominant features that are position- and rotation-invariant, which makes it possible for the model to be effectively trained. Figure 3 shows the example of max pooling layer.



Figure 2. 2D convolution operations



Figure 3. Max-pooling operations

A deep learning architecture's pooling and convolutional layers can process the input image and capture its features. In the next step, a typical neural network block is filled with the flattened output from the subsequent convolutional and pooling layers to classify an input image.

3.2 Dense layer

The dense layer is used for the DR classification based on the characteristics retrieved by the previous layers. Its main function is to utilize the high-level characteristics produced by the convolutional and pooling layers to make final predictions about the input data. The function of dense layer is like a neural network. The proposed system uses one hidden layer feedforward multilayer perceptron (MLPs) [20] and the backpropagation technique is used to train each sub-model. The MLP is trained in each sub-model by varying both the number of hidden neurons and the learning rate within a range. Each of these types of models is trained ten times using an initialization of random starting weight. To generate probability distributions across several classes for classification tasks, the dense layer often uses activation functions like SoftMax. The neuron in the dense layer's output that is most activated correlates to the class that is expected.

3.3 Synchronizing the sub-models

The synchronization sub-models in Figure 1 can be facilitated through the application of the Keras functional interface, a framework suitable for constructing models. Once the models are imported and organized into a list, it becomes possible to establish an extended stacking ensemble model. In this configuration, every loaded model serves as an independent input head to the overarching model. It's crucial to ensure that all layers within each of the loaded models are explicitly set as non-trainable. This precaution prevents any updates to their weights during the training of the new larger model. For seamless integration within Keras, it's essential for each layer to possess a unique identifier. Therefore, the nomenclature of each layer in every loaded model will need to be adjusted to indicate their association with a specific ensemble member.

Once the submodels are ready, each submodels input layer will be treated as an independent input source for the proposed ISM-DR model. Consequently, for each of the input models, the input data needs to be duplicated. The predictions generated by each individual model can then be brought together. For simplicity, a basic concatenation method is employed. This involves combining the four class probabilities from each of the five models, resulting in a unified 20-element vector.

Following this, a hidden layer that interprets this amalgamated input for the meta-learner is established. Subsequently, an output layer is defined to provide its own probabilistic predictions. Then the proposed stacked model encapsulates this process, generating a neural network model for stacked generalization. This function takes a collection of trained submodels as input and returns the corresponding stacked model. Upon defining the model, the next step is to train it. It's possible to directly fit the model using the withheld test dataset. As the submodels remain non-trainable, their weights remain unaltered during the training process. Only the weights of the newly introduced hidden and output layers are updated. The system is designed to train the stacking neural network with 300 epochs, ensuring comprehensive learning and convergence. Once the proposed model is trained, the newly assembled stacked model is used to generate predictions for classifying the fundus image for DR.

To provide better performance, the proposed ISM-DR model's weights are updated during training based on crossentropy loss. It is defined by:

$$Loss = -\sum_{n=1}^{m} t_n \log(s_n)$$
(1)

where, *m* is the total number of classes, t_n is the true label and s_n the sigmoid function output for the class *n*. The activation function used in the hidden layer is Rectified Linear Unit (ReLU) activation function. It is defined by:

$$U = \max(0, I) \tag{2}$$

where, I, U are the input and output value respectively. It is observed from Eq. (2) that this function only passes the positive values through the layers and neglects the negative values. Softmax function is defined by:

$$O(y_x) = \frac{e^{o_x}}{\sum_k e^{o_k}}$$
(3)

where, k is the number of previous layer's outputs. The x^{th} layer output is o_x . Table 2 shows the parameter settings to train each sub-model in the ISM-DR classification system.

 Table 2. Parameters setting to train the proposed ISM-DR classification system

Network Parameters	Settings
Epochs	300
Activation function (Hidden layer)	Rectified Linear
Learning rate	0.01
Loss function	Cross-entropy
Optimization function	Gradient Descent
Momentum	0.9
Prediction function (Output layer)	Softmax

4. RESULTS AND DISCUSSION



Figure 4. Number of images in fundus image datasets (a) MESSIDOR-1 (b) Kaggle dataset

(b)

The standard set of fundus images available in Methods to Evaluate Segmentation and Indexing Techniques in the field of Retinal Ophthalmology (MESSIDOR-1) (1200 images) [21] and Kaggle dataset (35126 images) [22] are utilized for performance evaluation of the proposed ISM for DR classification. MESSIDOR-1 images are in TIFF format with three different resolutions: 2240×1488, 1440×960 and 2304×1536pixels.Kaggle dataset images are in JPEG format different resolutions: 4752×3168. with 2592×1944. 4928×3264, and 3888×2592. Figure 4 shows the number of images in both datasets such as MESSIDOR-1 and Kaggle dataset for each DR grade. Figure 5 shows sample images from both datasets.

It can be seen from Figure 4 that class imbalance arises in situations in which the samples of one class are higher than another. This results in biased models that perform badly when tested in the minority class. To avoid class imbalance problems, many techniques such as re-sampling and under sampling can be employed. The number of instances in the minority classes is increased by generating more samples by data augmentation with the help of rotation and flipping. Also, the number of instances in the majority class is reduced by randomly removing samples. After re-sampling and under sampling, 2000 images from each class are selected from both datasets for the performance evaluation.

The experiments on MESSIDOR and Kaggle datasets are evaluated in terms of sensitivity or recall (R_c), precision or positive predictive value (P_r), and accuracy (A_c). Table 3 shows how the confusion matrix is formed to obtain the abovementioned performance measures and their formula. Accordingly, a true positive (TP) result is achieved when both the ground truth (actual) and the system output are 1. In all other cases, a false negative (FN) result is produced. If both the actual and system output are equal to zero, then a true negative (TN) is achieved; in any other circumstance, a false positive (FP) is produced.

To construct training, validation and testing sets, datasets were divided in the ratio of 60:20:20 respectively. This division utilizes a split method which is determined at random. Training is carried out on each run using the training set (1200 images/class), validation (400 images/class) and testing of the proposed system is carried out with the help of the validation set (400 images/class). The performance measures discussed in Table 1 and their corresponding 95% confidence interval (95% Cl) are computed for performance evaluation. The performance of the ISM-DR classification system is analyzed using different numbers of sub-models and the obtained confusion matrices are shown in Figure 6.



Figure 5. Sample images (a) MESSIDOR-1 (b) Kaggle dataset

Performance Measures	$A_c = \frac{TP + TN}{TP + FN + TN + FP} \qquad \qquad R_c$	$=\frac{TP}{TP+FN}$	$P_r = \frac{TP}{TP + FP}$
	Confusion Matrix		
System	Ground	Truth	
Output	Class-M		Class-N
Close M	#prediction of Class-M samples as Class-M samples	s. #prediction of	f Class-N samples as Class-M
Class-IVI	(True Positive - TP)	sample	es (False Positive - <i>FP</i>)
Close N	#prediction of Class-M samples as Class-N samples	s #prediction of	f Class-N samples as Class-N
Class-IV	(False Negative-FN)	sample	es (True Negative-TN)

Table 3. Performance metrics of the proposed ISM-DR classification system

20.070	21.070	21.070	21.070	21.170
79.5%	78.8%	78.2%	79.0%	78.9%
25	31	22	316	80.2%
1.6%	1.9%	1.4%	19.8%	19.8%
26	29	313	27	79.2%
1.6%	1.8%	19.6%	1.7%	20.8%
31	315	29	28	78.2%
1.9%	19.7%	1.8%	1.8%	21.8%
318	25	36	29	77.9%
19.9%	1.6%	2.2%	1.8%	22.1%

(a) 2 sub-models



Γ					
1	377	8	11	12	92.4%
	23.6%	0.5%	0.7%	0.8%	7.6%
2	9	375	9	6	94.0%
	0.6%	23.4%	0.6%	0.4%	6.0%
	8	6	372	8	94.4%
	0.5%	0.4%	23.2%	0.5%	5.6%
	6	11	8	374	93.7%
	0.4%	0.7%	0.5%	23.4%	6.3%
	94.2%	93.8%	93.0%	93.5%	93.6%
	5.8%	6.2%	7.0%	6.5%	6.4%
L	~	r	3	>	

(b) 3 sub-models

(c) 4 sub-models

Γ					
1	399	0	2	2	99.0%
	24.9%	0.0%	0.1%	0.1%	1.0%
2	0	397	3	1	99.0%
	0.0%	24.8%	0.2%	0.1%	1.0%
3	1	1	395	1	99.2%
	0.1%	0.1%	24.7%	0.1%	<mark>0.8%</mark>
4	0	2	0	396	99.5%
	0.0%	0.1%	0.0%	24.8%	0.5%
	99.8%	99.2%	98.8%	99.0%	99.2%
	0.2%	0.7%	1.2%	1.0%	0.8%
	~	r		>	

(d) 5 sub-models

Figure 6. Obtained confusion matrices of the proposed ISM-DR classification system for MESSIDOR-1 dataset

It is evident from Figure 6 that the average accuracy of the system improves as the number of sub-models increases. The first sub-model exhibits an average accuracy of 78.9% (95% CI: 75.5%-82.3%), while the subsequent ones show significant enhancements in performance. The third sub-model achieves an average accuracy of 89.8% (95% CI: 87.6% -92%), further improving to 93.6% (95% CI: 92.1%-95.1%) for the fourth sub-model. The fifth sub-model marks the highest accuracy level, reaching an impressive 99.2% (95% CI: 98.6%-99.8%). This demonstrates that the incorporation of additional submodels enhances the overall accuracy and effectiveness of the ISM-DR system, making it a more robust and reliable tool for DR classification. Table 4 shows the proposed ISM-DR classification system's performances for each class on MESSIDOR-1 images. Further, the performances of the ISM-DR classification system are analyzed using Kaggle dataset images and the obtained confusion matrices are shown in Figure 7.

The performance comparison of the sub-models in Table 4 displays unique patterns across the Pr and Rc metrics by comparing their respective performance. Specifically, there is a steady trend of improvement in both Pr and Rc scores as the number of sub-models grows from two to five. When the classification is applied on normal images, the model obtains a Pr and Rc of 0.779 and 0.795 respectively with two submodels. However, when there are five sub-models applied, these metrics considerably increase to 0.990 and 0.998 respectively. When it comes to classification tasks, it is evident that there is a definite pattern of performance with a greater number of sub-models across all levels of DR. This highlights the benefit of using a five-ensemble model for DR classification of MESSIDOR-1 dataset images. When it comes to constructing and improving ensemble-based classification systems, these findings highlight how important it is to take into consideration the number of sub-models.

Table 4. Precision and recall measures of the proposed ISM-DR classification system for each class for MESSIDOR-1 dataset

]	Performan	ce Analysi	s		
Class	2 sub-	models	3 sub-	models	4 sub-	models	5 sub-	models
	P_r	R_c	P_r	R_c	P_r	R_c	P_r	R_c
No DR	0.779	0.795	0.926	0.905	0.924	0.942	0.990	0.998
Mild DR	0.782	0.788	0.902	0.898	0.940	0.938	0.990	0.992
Moderate DR	0.792	0.782	0.897	0.892	0.944	0.930	0.992	0.998
PDR	0.802	0.790	0.869	0.898	0.937	0.935	0.995	0.990

Table 5. Precision and recall measures of the proposed ISM-DR classification system for each class for Kaggle dataset

				Performan	ce Analysis			
Class	2 sub-	models	3 sub-	models	4 sub-	models	5 sub-	models
	P_r	R_c	P_r	R_c	P_r	R_c	P_r	R_c
No DR	0.737	0.742	0.888	0.892	0.925	0.928	0.993	0.995
Mild DR	0.758	0.735	0.890	0.888	0.915	0.918	0.990	0.988
Moderate DR	0.737	0.742	0.891	0.895	0.920	0.922	0.988	0.992
Severe DR	0.738	0.738	0.903	0.888	0.918	0.918	0.992	0.990
PDR	0.729	0.740	0.874	0.882	0.919	0.912	0.990	0.988

			Confusi	on Matrix		
1	297	22	31	25	28	73.7%
	14.8%	1.1%	1.6%	1.2%	1.4%	26.3%
2	25	294	18	28	23	75.8%
	1.2%	14.7%	0.9%	1.4%	1.1%	24.2%
3	27	33	297	21	25	73.7%
	1.4%	1.7%	14.8%	1.1%	1.2%	26.3%
4	22	27	28	295	28	73.8%
	1.1%	1.4%	1.4%	14.8%	1.4%	26.2%
5	29	24	26	31	296	72.9%
	1.5%	1.2%	1.3%	1.6%	14.8%	27.1%
	74.2%	73.5%	74.2%	73.8%	74.0%	74.0%
	25.7%	26.5%	25.7%	26.2%	26.0%	26.0%
1	~	r	'n	Þ	6	

(a) 2 sub-models

		1200	100		200	2.2 2.2
	357	10	11	13	11	88.8%
	17.8%	0.5%	0.5%	0.7%	0.5%	11.2%
	9	355	7	11	17	89.0%
	0.4%	17.8%	0.4%	0.5%	0.9%	11.0%
	16	10	358	10	8	89.1%
	0.8%	0.5%	17.9%	0.5%	0.4%	10.9%
	6	11	10	355	11	90.3%
	0.3%	0.5%	0.5%	17.8%	0.5%	9.7%
	12	14	14	11	353	87.4%
	0.6%	0.7%	0.7%	0.5%	17.6%	12.6%
	89.2%	88.8%	89.5%	88.8%	88.2%	88.9%
	10.8%	11.3%	10.5%	11.3%	11.8%	11.1%
L	~	2	<u>~</u>	Þ	6	
	1000		Tarnet	Class		

(b) 3 sub-models

074	-		-		00 50
18.6%	0.4%	9 0.4%	5 0.2%	9 0.4%	92.5%
			_	<u></u>	
10 0.5%	367 18.4%	10 0.5%	0.4%	0.4%	8.5%
7	8	369	7	10	92.0%
0.4%	0.4%	18.4%	0.4%	0.5%	8.0%
6	11	7	367	9	91.8%
0.3%	0.5%	0.4%	18.4%	0.4%	8.3%
6	7	5	14	365	91.9%
0.3%	0.4%	0.2%	0.7%	18.2%	8.1%
92.8%	91.8%	92.2%	91.8%	91.2%	92.0%
7.3%	8.3%	7.8%	8.3%	8.8%	8.1%
~	r	3	Þ	5	

	308	0	1	1	1	00 3%
1	19.9%	0.0%	0.1%	0.1%	0.1%	0.7%
2	1	395	0	1	2	99.0%
	0.1%	19.8%	0.0%	0.1%	0.1%	1.0%
3	0	2	397	2	1	98.8%
	0.0%	0.1%	19.9%	0.1%	0.1%	1.2%
4	0	2	0	396	1	99.2%
	0.0%	0.1%	0.0%	19.8%	0.1%	0.8%
5	1	1	2	0	395	99.0%
	0.1%	0.1%	0.1%	0.0%	19.8%	1.0%
	99.5%	98.8%	99.2%	99.0%	98.8%	99.1%
	0.5%	1.2%	0.7%	1.0%	1.2%	0.9%
<u></u>	~	r	3	⊳	6	

(c) 4 sub-models

(d) 5 sub-models

Figure 7. Obtained confusion matrices of the proposed ISM-DR classification system for Kaggle dataset

In can be seen from Figure 7 that the number of sub-models integrated into the stacking model increases, there is a clear trend of improved performance, as measured by average accuracy. The first sub-model, with two sub-models, achieves an accuracy of 74% (95% CI: 68.7%-79.3%), third sub-model reaches an accuracy of 88.9% (95% CI: 85.1%-92.7%), while the fourth one further improves to 92% (95% CI: 89.6%-94.4%). The fifth sub-model records the highest accuracy at 99.1% (95% CI: 98.6%-99.6%). This pattern suggests that the stacking model approach effectively leverages the strengths of multiple sub-models, resulting in a highly accurate and robust system. Table 5 shows the proposed ISM-DR classification system's performances for each class on Kaggle images.

It can be seen from Table 5 that the ISM-DR classification system is increased when the number of sub-models increases. There is a consistent improvement in both *Pr* and *Rc* metrics for all classes of DR images. It is also demonstrated that the system can reach higher levels of *Pr* and *Rc* by using five sub-models in the proposed ISM-DR systems which is similar to the performances on MESSIDOR-1 images. Figure 8 and Figure 9 show the comparative analysis of the ISM-DR classification system with conventional deep learning architectures, such as VGG [23], ResNet [24], AlexNet [25], GoogleNet [26] using MESSIDOR-1 and Kaggle dataset respectively.



Figure 8. Comparative analysis of the proposed ISM-DR classification system for MESSIDOR-1 dataset



Figure 9. Comparative analysis of the proposed ISM-DR classification system for Kaggle dataset

Among the architectures considered for performance comparison in Figures 8 and 9, the ISM-DR system demonstrates superior performance to others. The ISM-DR system with an impressive accuracy of 99.2% (MESSIDOR-1) and 99.1% (Kaggle) followed by GoogleNet with 97.9% (MESSIDOR-1) and 96.5% (Kaggle). In terms of precision, both GoogleNet and the ISM-DR system maintain high precision scores of 97.5% and 99.2% for MESSIDOR-1 and 96% and 99.1% (Kaggle), respectively. Similarly, regarding recall, these models excel with GoogleNet achieving a recall of 98.5% and the ISM-DR system reaching an outstanding recall of 99.5%. VGG, ResNet, and AlexNet also exhibit

slightly lower compared to GoogleNet and the ISM-DR system, with accuracy ranging from 87.2% to 93.7% and precision and recall scores ranging from 85.7% to 95% and 90% to 95%, respectively. Overall, while all models demonstrate significant capabilities, the ISM-DR system has exceptional accuracy, precision, and recall, making them preferred choices for tasks requiring high-performance image classification.

5. CONCLUSIONS

An efficient deep learning architecture for DR classification is proposed in this paper using ISM concept. Stacking is an ensemble learning approach to create a meta-model that combines the multiple base models' predictions, aiming to improve overall predictive accuracy. The outcomes of the proposed ISM-DR classification system are more reliable and accurate by making use of the many viewpoints offered by each sub-model. The inherent flaws in each model are mitigated by the ISM, which gives the meta-learner the ability to discover the optimal approach to integrate the predictions of the submodels for DR classification. Experimental results on MESSIDOR-I and Kaggle datasets show it superior performance than other architectures and state-of-art techniques in the literature. It is clearly demonstrated in the results section that the ISM concept is a powerful strategy for boosting predictive accuracy and making the system more reliable for DR classification.

REFERENCES

- [1] Ramitha, M.A., Mohanasundaram, N. (2021). Classification of pneumonia by modified deeply supervised ResNet and SEnet using chest X-Ray images. International Journal of Advances in Signal and Image Sciences, 7(1): 30-37. https://doi.org/10.29284/ijasis.7.1.2021.30-37
- [2] Senthil Kumar, M., Azath, H., Velmurugan, A.K., Padmanaban, K., Subbiah, M. (2023). Prediction of Alzheimer's disease using hybrid machine learning technique. AIP Conference Proceedings, 2523(1): 1-12. https://doi.org/10.1063/5.0110283
- [3] Olufunso, O.S., Evwiekpaefe, A.E., Irhebhude, M.E. (2022). Ethnicity classification using a dynamic horizontal voting ensemble approach based on fingerprint. International Journal of Advances in Signal and Image Sciences, 8(2): 36-47. https://doi.org/10.29284/ijasis.8.2.2022.36-47
- [4] Alyoubi, W.L., Shalash, W.M., Abulkhair, M.F. (2020). Diabetic retinopathy detection through deep learning techniques: A review. Informatics in Medicine Unlocked, 20: 1-11. https://doi.org/10.1016/j.imu.2020.100377
- [5] Saxena, G., Verma, D. K., Paraye, A., Rajan, A., Rawat, A. (2020). Improved and robust deep learning agent for preliminary detection of diabetic retinopathy using public datasets. Intelligence-Based Medicine, 3-4: 1-11. https://doi.org/10.1016/j.ibmed.2020.100022
- [6] Zhang, Q.M., Luo, J., Cengiz, K. (2021). An optimized deep learning-based technique for grading and extraction of diabetic retinopathy severities. Informatica, 45(5): 1-8. https://doi.org/10.31449/inf.v4515.3561
- [7] Erciyas, A., Barışçı, N. (2021). An effective method for

detecting and classifying diabetic retinopathy lesions based on deep learning. Computational and Mathematical Methods in Medicine, 2021: 1-13. https://doi.org/10.1155/2021/9928899

- [8] Zhang, G., Sun, B., Chen, Z., Gao, Y., Zhang, Z., Li, K., Yang, W. (2022). Diabetic retinopathy grading by deep graph correlation network on retinal images without manual annotations. Frontiers in Medicine, 9: 1-9. https://doi.org/10.3389/fmed.2022.872214
- [9] Zhang, G., Sun, B., Zhang, Z., Pan, J., Yang, W., Liu, Y. (2022). Multi-model domain adaptation for diabetic retinopathy classification. Frontiers in Physiology, 13: 1-10. https://doi.org/10.3389/fphys.2022.918929
- [10] Abràmoff, M.D., Lou, Y., Erginay, A., Clarida, W., Amelon, R., Folk, J.C., Niemeijer, M. (2016). Improved automated detection of diabetic retinopathy on a publicly available dataset through integration of deep learning. Investigative Ophthalmology & Visual Science, 57(13): 5200-5206. https://doi.org/10.1167/iovs.16-19964
- [11] Tajudin, N.M.A., Kipli, K., Mahmood, M.H., Lim, L.T., Awang Mat, D.A., Sapawi, R., Hoque, M.E. (2022). Deep learning in the grading of diabetic retinopathy: A review. IET Computer Vision, 16(8): 667-682. https://doi.org/10.1049/cvi2.12116
- [12] Shankar, K., Sait, A.R.W., Gupta, D., Lakshmanaprabu, S.K., Khanna, A., Pandey, H.M. (2020). Automated detection and classification of fundus diabetic retinopathy images using synergic deep learning model. Pattern Recognition Letters, 133: 210-216. https://doi.org/10.1016/j.patrec.2020.02.026
- [13] Berbar, M.A. (2022). Diabetic Retinopathy Detection and Grading using Deep learning. Menoufia Journal of Electronic Engineering Research, 31(2): 11-20. https://doi.org/10.21608/mjeer.2022.138003.1057
- [14] Mahmood, M.A.I., Aktar, N., Kader, M.F. (2023). A hybrid approach for diagnosing diabetic retinopathy from fundus image exploiting deep features. Heliyon, 9(9): 1-14. https://doi.org/10.1016/heliyon.2023.e119625
- [15] Rahim, S.S., Palade, V., Shuttleworth, J., Jayne, C. (2016). Automatic screening and classification of diabetic retinopathy and maculopathy using fuzzy image processing. Brain Informatics, 3: 249-267. https://doi.org/10.1007/s40708-016-0045-3
- [16] Kalyani, G., Janakiramaiah, B., Karuna, A., Prasad,

L.V.N. (2023). Diabetic retinopathy detection and classification using capsule networks. Complex and Intelligent Systems, 9(3): 2651-2664. https://doi.org/10.1007/s40747-021-00318-9

- [17] Papadopoulos, A., Topouzis, F., Delopoulos, A. (2021). An interpretable multiple-instance approach for the detection of referable diabetic retinopathy in fundus images, Scientific Reports, 11(1): 1-15. https://doi.org/10.1038/s41598-021-93632-8
- [18] Vogiatzis, A., Orfanoudakis, S., Chalkiadakis, G., Moirogiorgou, K., Zervakis, M. (2022). Novel metalearning techniques for the multiclass image classification problem. Sensors, 23(1): 1-23. https://doi.org/10.3390/s23010009
- [19] Manikandan, S.P., Karthikeyan, V., Nalinashini, G. (2022). DermICNet: Efficient dermoscopic image classification network for automated skin cancer diagnosis. Revue d'Intelligence Artificielle, 36(5): 1-7. https://doi.org/10.18280/ria.360519
- [20] Lakshmi, V.V., Jasmine, J.L. (2021). A hybrid artificial intelligence model for skin cancer diagnosis. Computer Systems Science and Engineering, 37(2): 233-245. https://doi.org/10.32604/csse.2021.015700
- [21] Messidor-ADCIS. http://www.adcis.net/en/thirdparty/messidor/, accessed on 20 Oct. 2022.
- [22] Kaggle DR Dataset. https://www.kaggle.com/c/diabeticretinopathy-detection/data, accessed on 20 Oct. 2022.
- [23] Simonyan, K., Zisserman, A. (2015). Very deep convolutional networks for large-scale image recognition. Internation Conference on Learning Representations, pp. 1-14. https://doi.org/10.45550/arXiv.1409.1556
- [24] He, K., Zhang, X., Ren, S., Sun, J. (2016). Deep residual learning for image recognition. In IEEE Conference on Computer Vision and Pattern Recognition, pp. 770-778. https://doi.org/10.48550/arXiv.1512.03385
- [25] Krizhevsky, A., Sutskever, I., Hinton, G.E. (2012). Imagenet classification with deep convolutional neural networks. Proceedings of Advances in Neural Information Processing Systems, pp. 1-9. https://doi.org/10.1145/3065386
- [26] Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S. (2014). Going deeper with convolutions. In IEEE Conference on Computer Vision and Pattern Recognition, pp. 1-9. https://doi.org/10.48550/arXiv.1409.4842