

Analyze Informant-based Questionnaire for The Early Diagnosis of Senile Dementia Using Deep Learning

Fubao Zhu, Ph.D.¹, Xiaonan Li, B.S.¹, Daniel McGonigle, B.S.², Haipeng Tang, M.S.², Zhuo He, B.S.³,
Chaoyang Zhang, Ph.D.², Guang-Uei Hung, M.D.⁴, Pai-Yi Chiu, M.D.^{5*},
Weihua Zhou, Ph.D.^{3*}

¹School of Computer and Communication Engineering, Zhengzhou University of Light Industry, Zhengzhou, Henan

²School of Computing Sciences and Computer Engineering, University of Southern Mississippi, Long Beach, MS

³College of Computing, Michigan Technological University, Houghton, MI

³Department of Neurology, Show Chwan Memorial Hospital, Changhua, Taiwan

⁴Department of Nuclear Medicine, Show Chwan Memorial Hospital, Changhua, Taiwan

*Corresponding authors:

Pai-Yi Chiu, MD, PhD

Department of Neurology, Show Chwan Memorial Hospital, Changhua, Taiwan

E-mail: paiyibox@gmail.com

Weihua Zhou, PhD

College of Computing, Michigan Technological University, Houghton, MI 49931

E-Mail: whzhou@mtu.edu

Abstract

Objective: This paper proposes a multiclass deep learning method for the classification of dementia using an informant-based questionnaire.

Methods: A deep neural network classification model based on Keras framework is proposed in this paper. To evaluate the advantages of our proposed method, we compared the performance of our model with industry-standard machine learning approaches. We enrolled 6,701 individuals, which were randomly divided into training data sets (6030 participants) and test data sets (671 participants). We evaluated each diagnostic model in the test set using accuracy, precision, recall, and F1-Score.

Results: Compared with the seven conventional machine learning algorithms, the DNN showed higher stability and achieved the best accuracy with 0.88, which also showed good results for identifying normal (F1-score=0.88), mild cognitive impairment (MCI) (F1-score=0.87), very mild dementia (VMD) (F1-score=0.77) and Severe (F1-score=0.94).

Conclusion: The deep neural network (DNN) classification model can effectively help doctors accurately screen patients who have normal cognitive function, mild cognitive impairment (MCI), very mild dementia (VMD), mild dementia (Mild), moderate dementia (Moderate), and severe dementia (Severe).

Keywords: dementia, information gain, deep neural network, machine learning

1. INTRODUCTION

Alzheimer's disease is a neurodegenerative disease characterized by cognitive and intellectual impairment,

often later accompanied by decreased bodily functions and life expectancy. It is also the most common form of the syndrome known as dementia, which affects 47.5 million people worldwide according to the World Health Organization[1]. In vivo diagnostic methods currently face

limitations in accurately identifying neurodegenerative disorders, with confirmatory diagnosis often occurring with the post-mortem autopsy or biopsy[2]. Clinical/cognitive function metrics such as the Mini-Mental State Examination scale and the Alzheimer’s disease assessment scale assess the cognitive status of patients, indicating potential cases of Alzheimer’s disease. Cognitive and/or memory disorders are currently the primary clinical markers used to identify subjects at risk of developing dementia.

Applying artificial intelligence technology to the early diagnosis of Alzheimer’s disease has become an important research topic. Examples of include neuroimaging[3]-[5] with structural magnetic resonance imaging (tissue density, cortical surface and hippocampal measurements), functional MRI (slight coherence and functional connectivity of different brain regions), diffusion tensor imaging (along the pattern of white matter fibers), positron emission tomography (PET) and other criteria for assessing the status of subjects. This study focuses on helping doctors make preliminary judgments from simple cognitive screening surveys. In this paper, we proposed a deep neural network (DNN) classification model to assist the preliminary diagnosis of normal, mild cognitive impairment (MCI), very mild dementia (VMD), mild dementia (Mild), moderate dementia (Moderate), and severe dementia (Severe) using informant-based questionnaire.

In recent years, deep neural network methods for developing classifiers have shown amazing performance in many recognition tasks, which has garnered clinical interest toward cognitive diagnostic applications. This paper focuses on analyzing the predictive power of cognitive tests in diagnosing patients with Alzheimer’s disease. A popular research approach to studying this illness at present focused on neuroimaging, such as the MRI, SPECT or PET model, which is the gold standard for assessing the status of subjects. Yong et al.[6] summarized the latest advances in AD brain network research. The evidence collected from the study indicates that AD patients are associated with a comprehensive anomaly in the distributed neuron network. These findings may provide a new way for clinical diagnosis and monitoring of AD progress, and help us

uncover imaging-based biomarkers for diagnosis and monitoring of disease. Jain et al.[7] proposed a transfer learning approach for accurately classifying brain sMRI slices amongst 3 different classes: AD, CN and MCI. They adopt a mathematical model P F S E C TL based on transfer learning, adopt the CNN architecture, and use the image network data set as the training object VGG-16 as the feature extractor of the classification task. For the validation set, the accuracy of the three-way classification using the method is 95.73%. Orimaye et al.[8] proposed a combination of deep neural network and deep language models (D2NNLM) for classifying the disease. The experimental results show that the model can accurately predict MCI and AD type dementia on a very sparse clinical datasets. Donghuan et al.[9] proposed a novel deep-learning-based framework to discriminate individuals with AD utilizing a multimodal and multiscale deep neural network. They have obtained 82.4% accuracy in identifying the individuals with mild cognitive impairment (MCI). Dolph et al.[10] studied multiclass deep learning classification of Alzheimer’s disease (AD) using novel texture and other associated features extracted from structural MRI. They have got classification accuracies of 51.4% and 56.8%. Our future work will be carried out from neuroimaging to further improve our diagnostic model.

II. MATERIALS AND METHODS

There are two major steps in the proposed framework: (1) feature selection: using feature selection algorithms to optimize or even reduce the number of neuropsychological tests; (2) classification: training a deep neural network to classify the participants into normal cognitive function, MCI, VMD, Mild, Moderate, and Severe categories.

A. Patient sample collection

In this work, the study used data collected from the three centers of the Show Chwan Healthcare System. The data selected from the register-based database of the Show Chwan Health System were analyzed anonymously with the informed consent from all participants, and the study was designed retrospectively in accordance with relevant

guidelines and regulations. The project was reviewed by the Medical Research Ethics Committee of Show Chwan Memorial Hospital, and the study was approved by the Data Inspectorate.

The data for the study consisted of samples of clinical and neuropsychological assessment obtained from 6701 patients. For detailed neuropsychological tests, we assessed the history of cognitive status and objective assessments including the Clinical Dementia Ratings (CDR), Mini Mental Status Examination (MMSE), Cognitive Abilities Screening Instrument (CASI) and Montreal Cognitive Assessment (MoCA) performed to evaluate memory, executive function, orientation, visual-spatial ability, and language function[11]. Along with the current scales such as CDR, MMSE, CASI, MoCA, we used a newly designed Informant-based questionnaire named HAICDDS which is applied in dementia registration in a health system with 9 regional hospitals in Taiwan. Clinical application of the HAICDDS had been published in journals [11]-[13] or conferences[14]-[15]. The CDR determined the severity of dementia. Experienced neurologists evaluated the participants based on their clinical symptoms and reviews of medical/medication history, neuropsychological test results, and then classified the participants into six diagnostic groups: normal (535 participants), MCI (1687 participants), VMD (678 participants), Mild (1812 participants), Moderate (1309 participants), and Severe (680 participants). The six diagnostic groups were defined using the CDR staging. Among CDR 0.5, participants without significantly impaired activities of daily living were divided as CDR 0.5 MCI and those with significantly impaired activities of daily living were divided as CDR 0.5 very mild dementia (VMD). Therefore, the 6 groups were CDR 0, CDR 0.5 MCI, CDR 0.5 VMD, CDR 1, CDR 2, and CDR 3. The operational diagnosis of a significant interfere with ADL is the IADL total score <7.

To provide an estimate of generalization error, we used the `cross_validation.train_test_split` function provided by `scikit-learn`[16] to randomly split the data into a training data sets (6030 participants), and a test data sets (671 participants). We finally obtained 10 groups of different training set (6030 participants) and test set (671

participants). The reason is that 10 groups is the appropriate choice to obtain the best error estimate through a large number of tests[17]-[18]00.

B. Feature selection

Neuropsychologists selected 50 items from neuropsychological tests to form an optimal questionnaire for screening patients with varying severities of dementia. While the proposed algorithm still works with the entirety of those features, utilizing feature selection lowers the computational requirements0 and improves interpretability as we are able to see more clearly which items are directly correlated with a patient's mental condition[11]. In order to retain the features with higher prediction performance, provide faster and more cost-effective predictors, reduce the curse of dimensionality problem and the possibility of overfitting during the training phase, we used information gain feature selection algorithms to rank the importance score of all 50 features and then the low ranking features were filtered out. We discard features with a lower score one by one, and input the remaining features into the DNN model to observe the change of the classification accuracy in order to identify the feature set with a smaller number of features but only a minor drop-off of classification accuracy.

Information gain is an information theory method widely used in data mining00. The information gain measures how much information the feature can provide the classification model. If a feature has a larger information gain value for a class, the feature contains more classification information for that class. We used the information gain algorithm provided by Weka0, which is an open source machine learning and data mining software based on the Java environment.

C. Overview of methods

We proposed a deep neural network (DNN) classification model based on the Keras framework. In order to study the performance of the DNN model for discriminating normal, MCI, VMD, Mild, Moderate, and Severe, we compared its results with others well-known classification models (MLP, GCForest, random forest, AdaBoost, LogitBoost, Naïve Bayes and SVM). First, we tested the performance of the

model respectively using 50 features selected by the Neuropsychologists and the top 44 features selected by information gain score. Then, for testing the stability of the classification model, we tested 10 groups of different training sets and test sets separately. Finally, we evaluated the model's classification of each of the 10 holdout sets by measuring accuracy and F1-Score. TP is the number of positive samples predicted by the classifier in the number of true positive samples, FP is the number of positive samples predicted by the classifier in the number of true negative samples, TN is the number of negative samples predicted by the classifier in the number of true negative samples, FN is the number of true positive samples predicted by the classifier as negative samples. The accuracy is the evaluation of the correct rate of the classifier as a whole. It is defined as $\text{Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{FN} + \text{FP} + \text{TN})$, generally speaking, the higher the accuracy, the better the classifier. F1-Score is a kind of statistic, which is also called F-measure. F1-Score is the weighted harmonic average of Precision and Recall. It is defined as $\text{F1-Score} = 2\text{TP} / (2\text{TP} + \text{FP} + \text{FN})$ and carries a range of 0 to 1, with higher scores indicating a more robust classification model. It is a commonly used evaluation criterion in the field of IR (Information Retrieval). It is often used to evaluate the quality of classification models. Precision, also called Positive Predictive Value in clinical settings, refers to how many of the samples that the model is positive are true positive samples, which is defined as $\text{Precision} = \text{TP} / (\text{TP} + \text{FP})$. Recall, also called Sensitivity in clinical settings, refers to how many positive samples are classified as positive by the model, which is defined as $\text{Recall} = \text{TP} / (\text{TP} + \text{FN})$.

1) DNN

DNN is a multi-hidden layer feedforward neural network, which has a total of $L+1$ layers, the 0th layer is the input layer, the 1st to $L-1$ layers are hidden layers, the L th layer is the output layer. The nodes of adjacent layers are connected by links and the weights of all links form a feedforward weight matrix. As shown in Eqs.(1)-(2),

suppose there are n_l neurons in the l th layer, and the input

vector of these neurons is $a^{(l)}$, the output vector is $z^{(l)}$. At the same time, we distinguish the final output of DNN on the output of the hidden layer by $u = y^{(L)}$. Given the characteristic x of a training sample, there is $a^{(0)} = z^{(0)} = x$.

$$a^{(l)} = W^{(l)} * a^{(l-1)} + b^{(l)}, l=1,2,\dots,L \quad (1)$$

$$z^{(l)} = h(a^{(l)}) \dots\dots\dots (2)$$

where $W^{(l)}$ is the weight matrix of the $l-1$ layer to the l th layer, b is the offset vector of the l th layer, $h()$ is the activation function of the l th layer.

As a feedforward neural network, given an input vector, DNN can get an output vector immediately, that is to say, the output of DNN only depends on the current input, so DNN is suitable for pattern classification problem. This paper adjusts the network parameters of the whole network through supervised training. After repeated extensive training, we got relatively optimal parameters. The optimal parameters of the constructed DNN model use three layers, the first layer uses the relu activation function, the second layer uses the tanh activation function, the third layer uses the softmax activation function. The epoch is set to 40, the dropout rate is set to 0.2, the batch size is set to 32, the learning rate is set to 0.004, and the number of neurons is set to 20 in each layer. The specific structure of the DNN model is shown in the figure below. where l is the number of layers, x is the input feature, b is the offset vector of the

l th layer, $W^{(l)}$ is the weight matrix of the $l-1$ layer to the

l th layer, $a^{(l)}$ is the input vector of the l th layer, $h()$ is the

activation function of the l th layer, $z^{(l)}$ is the input vector of the next layer, $f()$ is the output activation function and y is the output vector.

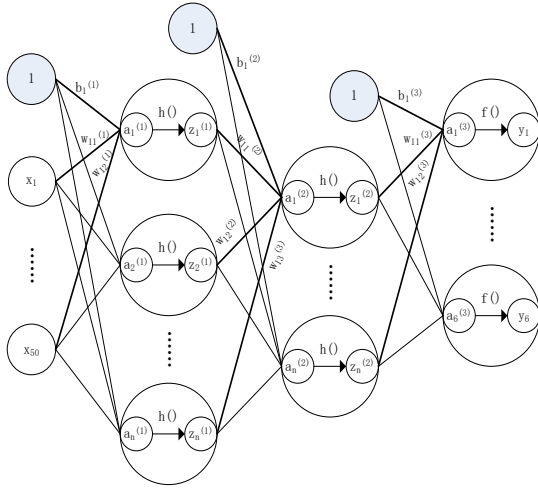


Fig 1. DNN model structure

2) Other investigated classification models

We investigated other commonly used classification models (MLP, GCForest, random forest, AdaBoost, LogitBoost, Naïve Bayes and SVM) in Python toolbox[16], which is a set of freeware academic software packages. The following briefly describes the basic principles of the classification models and the further details can refer to cited literatures.

An MLP can be seen as a directed graph, consisting of multiple node layers, each layer connected to the next layer. In addition to input nodes, each node is a neuron (or processing unit) with a non-linear activation function. Unlike the DNN model, the output activation function is not used here. GCForest is a model of a deep forest, which is mainly divided into two parts, multi-grained scanning, and cascade forest structure. GCForest performs well in small sample data. Random forest is an algorithm that integrates multiple trees by the idea of ensemble learning. Its basic unit is a decision tree, which is a subclass of ensemble learning. It depends on the voting choice of a decision tree to determine the final classification results. In the basic Adaboost algorithm, each weak classifier has the right to weight, and the weighted sum of the weak classifier prediction results forms the final prediction result. In training, training samples have also weight, which dynamically adjusts during the training process. The samples that are misclassified by the previous weak classifier will increase the weight, so the algorithm will focus on the difficult samples. The Logitboost algorithm is a discriminant classification algorithm based on machine

learning. LogitBoost belongs to the AdaBoost system. The LogitBoost structure is similar in general, but its loss function uses the maximum logarithmic likelihood function. The basic method of Naïve Bayes is to calculate the probability that the current feature samples belong to a certain classification based on the statistical data and the conditional probability formula, and select the maximum probability classification. Support Vector Machine (SVM) has achieved the best performance in many classification problems. The kernel function subtly transforms the linear indivisible problem into a linear separable problem, and has very good generalization performance.

III. RESULTS

A. Feature analyses

We find that the classification accuracy decreases with the decrease of the number of features. When the number of features decreases to 44 features, the classification accuracy dropped down. So we discarded the corresponding features by setting the threshold of the information gain score to 0.16.

Figure 2 shows the trend of classification accuracy by our DNN model as the features with lower scores are discarded one by one. The classification accuracy is the average of ten experiments. With a decreasing number of features, the classification accuracy decreases. After reducing to 44 features, the subsequent classification accuracy has declined by a certain extent.

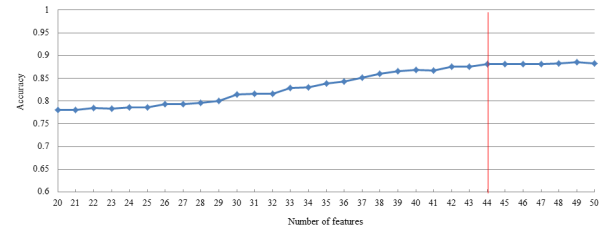


Fig 2. Classification accuracy with decreasing number of features.

Figure 3 shows the features ranked in descending significance with respect to the information gain scores. The cut-off is shown reducing the number of features to the 88% with setting the threshold of the information gain score to 0.16, thus reducing the feature number from 50 to 44 features. Among the top 44 selected features, the feature

‘H01’ has the highest ranking score of 0.902, and the feature ‘L01’ has the lowest ranking score of 0.1665.

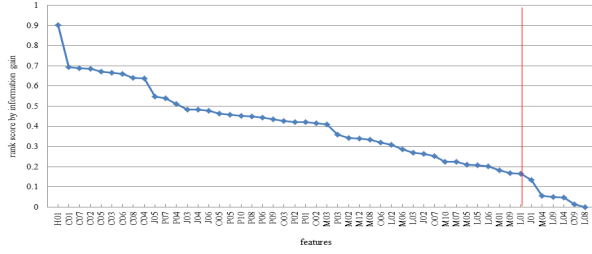


Fig 3. Features ranked according to their information gain scores.

B. Performance of classification models

Figure 4 shows the classification performance for each of the 10 rounds individually. (a) shows the accuracy analysis results using 50 features selected by the Neuropsychologists. (b) shows the accuracy analysis results using the top 44 features selected by information gain score. The accuracy performance of the DNN reaches a plateau in the 10 rounds, which are better than other algorithms. The performance of the model reduces when lower features are used as input into the classifier.

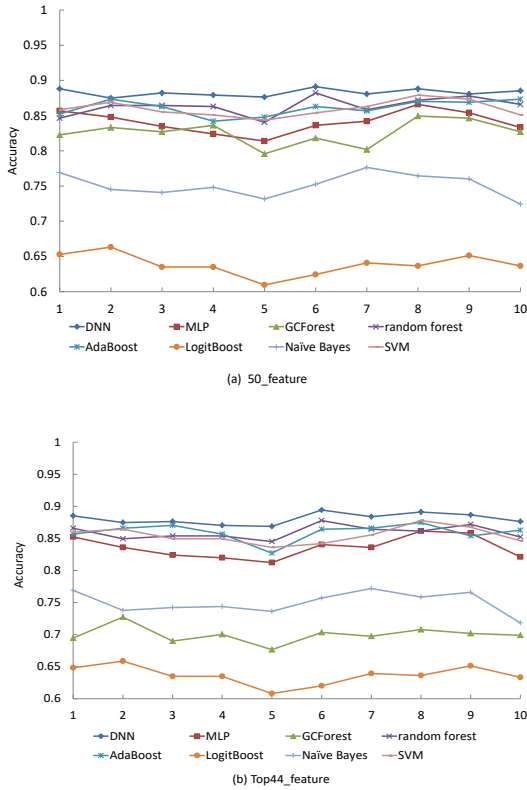


Fig 4. Performance of accuracy for each of the 10 rounds.

Figure 5 shows the average results in 10 rounds of the comparison of our DNN accuracy performance results and

other well-known classifiers (MLP, GCForest, random forest, AdaBoost, LogitBoost, Naïve Bayes and SVM) for the same dataset. When using 50 features, the best accuracy was obtained by the DNN classifier (accuracy=0.8748), followed by the MLP classifier (accuracy=0.851). When using the top 44 features, the DNN classifier performs the best (accuracy=0.8808), followed by the AdaBoost(accuracy=0.8599).

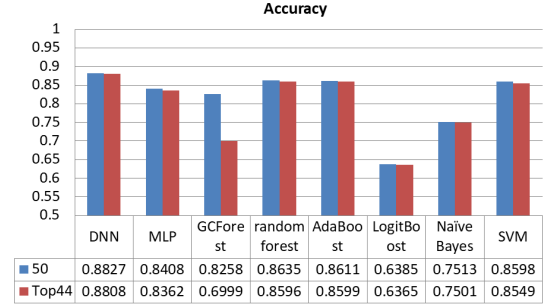


Fig 5. Comparison of the accuracy obtained by DNN and other classifiers.

C. Multi-class classification

Figure 6 compares the F1-score performance of the 10 rounds in the classification of normal, MCI, VMD, Mild, Moderate and Severe using DNN, MLP, GCForest, random forest, AdaBoost, LogitBoost, Naïve Bayes and SVM. As shown in Figure 6, when using all the 50 features, the DNN algorithm effectively improved the overall performance in classifying normal (F1-score=0.89), MCI (F1-score=0.89), VMD (F1-score=0.74), Mild (F1-score=0.85), Moderate (F1-score=0.88) and Severe (F1-score=0.92). When using the top 44 features selected by information gain score. The DNN algorithm performed best result in screening the normal (F1-score=0.88), MCI (F1-score=0.87), VMD (F1-score=0.77) and Severe (F1-score=0.94), and poorest in Mild (F1-score=0.83) and Moderate (F1-score=0.83) categories.

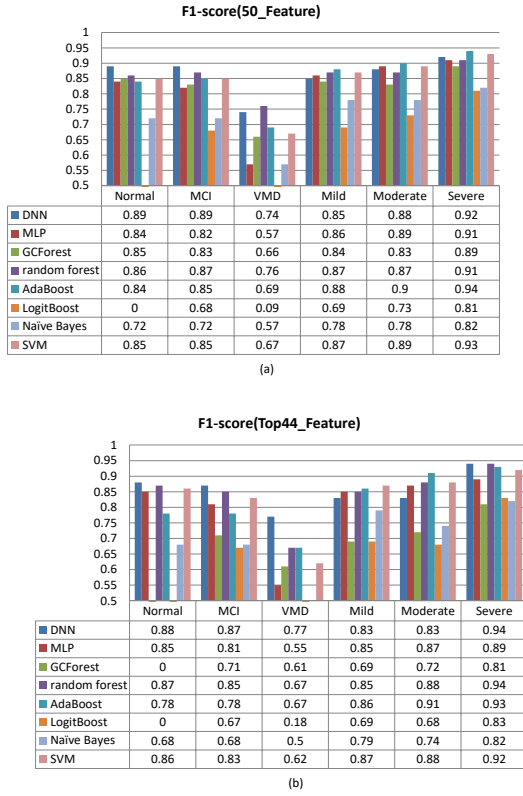


Fig 6 Performance of F1-score in the classification of normal, MCI, VMD, Mild, Moderate and Severe using classifiers.

IV. DISCUSSION

In this study, we proposed a deep neural network classification model based on the Keras framework. In order to evaluate the advantages of our proposed method, we compared two indicators, accuracy and F1-score. In addition, we compared our method with other well-known machine learning methods. The results showed our DNN method had a stable classification performance, higher classification accuracy and performed well in dealing with class imbalance problems. It has great potential for clinical application. We will discuss these in detail below.

From the perspective of model classification stability and accuracy, when using 50 features by the Neuropsychologists, the DNN model shows higher stability and the classification accuracy is the highest compared with the other seven algorithms (MLP, GCForest, random forest, AdaBoost, LogitBoost, Naive Bayes and SVM), basically stable at around 0.88. When it comes to the classification accuracy of each category, our results show that the DNN

model improved the overall performance of the classification accuracy of each category.

We further studied the classification performance of the DNN model after reducing features by information gain feature selection, which is to simplify the procedure of diagnosis and enhance the practicality in clinic. The overall classification performance of the model has decreased after the reduction. In order to ensure the classification accuracy of the model, we set the threshold of information gain fraction to 0.16, thus discarding some redundant features. Compared with the six classification models, the DNN performed the best accuracy with 0.88, which also showed good results for identifying normal (F1-score=0.88), MCI (F1-score=0.87), VMD (F1-score=0.77) and Severe (F1-score=0.94).

V. CONCLUSIONS

We proposed a new approach to diagnosing normal, MCI, VMD, Mild, Moderate, and Severe using a deep learning approach, more specifically, a deep neural network classification model based on the Keras framework. By using the real-world dataset, i.e., the register-based database in the Show Chwan Health System, we tested and validated our method. Overall, the results of this project show that the proposed DNN model provides a tool with accurate and stable performance for clinicians to diagnose the early stages of dementia.

ACKNOWLEDGMENTS

This project was supported by a grant from Henan Science and Technology Department (Project Number: 182102210157, PI: Fubao Zhu). This research was also supported in part by the American Heart Association under Award Number 17AIREA33700016.

REFERENCES

- [1] The Lancet. Dementia burden coming into focus. Lancet. 2017;390(10113):2606. doi:10.1016/S0140-6736(17)33304-4.
- [2] Petrella JR, Coleman RE, Doraiswamy PM. Neuroimaging and Early Diagnosis of Alzheimer Disease: A Look to the Future. Radiology. 2003;226(2):315-336. doi:10.1148/radiol.2262011600.

- [3] Rathore S, Habes M, Ifitikhar MA, Shacklett A, Davatzikos C. A review on neuroimaging-based classification studies and associated feature extraction methods for Alzheimer's disease and its prodromal stages. *Neuroimage*. 2017;155(March):530-548. doi:10.1016/j.neuroimage.2017.03.057.
- [4] Duchesne S, Caroli A, Geroldi C, Collins DL, Frisoni GB. Relating one-year cognitive change in mild cognitive impairment to baseline MRI features. *Neuroimage*. 2009;47(4):1363-1370. doi:10.1016/j.neuroimage.2009.04.023.
- [5] Stonnington CM, Chu C, Klöppel S, Jack CR, Ashburner J, Frackowiak RSJ. Predicting clinical scores from magnetic resonance scans in Alzheimer's disease. *Neuroimage*. 2010;51(4):1405-1413. doi:10.1016/j.neuroimage.2010.03.051.
- [6] He, Y., Chen, Z., Gong, G. & Evans, A. Neuronal networks in Alzheimer's disease. *Neuroscientist* 15, 333-350, doi:10.1177/1073858409334423 (2009).
- [7] Jain, R., Jain, N., Aggarwal, A. & Hemanth, D. J. Convolutional neural network based Alzheimer's disease classification from magnetic resonance brain images. *Cogn Syst Res* 57, 147-159, doi:10.1016/j.cogsys.2018.12.015 (2019).
- [8] Orimaye, S. O., Wong, J. S. & Wong, C. P. Deep language space neural network for classifying mild cognitive impairment and Alzheimer-type dementia. *PLoS One* 13, e0205636, doi:10.1371/journal.pone.0205636 (2018).
- [9] Lu D, Popuri K, Ding GW, Balachandar R. Multimodal and Multiscale Deep Neural Networks for the Early Diagnosis of Alzheimer's Disease using structural MR and FDG-PET images. 2018;(October 2017):1-13. doi:10.1038/s41598-018-22871-z.
- [10] Dolph C V, Alam M, Shboul Z, Samad MD, Iftekharuddin KM. Deep Learning of Texture and Structural Classification. 2017 Int Jt Conf Neural Networks. 2017;(1310353):2259-2266. doi:10.1109/IJCNN.2017.7966129.
- [11] Chiu P, Tang H, Wei C, Zhang C, Hung G, Zhou W. NMD-12: A new machine-learning derived screening instrument to detect mild cognitive impairment and dementia. *PloS one* 14 (3), e0213430 2019:1-11.
- [12] Chiu PY, Wei CY, Hung GU. Preliminary Study of the History-based Artificial Intelligent Clinical Dementia Diagnostic System. *Show Chwan Medical Journal* 2019. DOI: 10.3966/156104972019061801003
- [13] Lin CM, Hung GU, Wei CY, Tzeng RC, Chiu PY*. An Informant-Based Simple Questionnaire for Language Assessment in Neurodegenerative Disorders. *Dement Geriatr Cogn Disord*. 2018 Oct 18;46(3-4):207-216. doi: 10.1159/000493540.
- [14] Chiu PY, Hung GU, Wei CY. History-based Questionnaire for the Diagnosis of Severity and Subtypes of Dementia: Design and Verify. *The Alzheimer's Association International Conference 2019 (AAIC 2019)*.
- [15] Chiu PY, Tsai CF. Freezing of Speech Single Questionnaire as a Screening Tool for Dementia with Lewy Bodies. *The Alzheimer's Association International Conference 2019 (AAIC 2019)*.
- [16] Pedregosa F, Varoquaux G, Gramfort A, et al. *Scikit-learn: Machine Learning in Python*. 2012;(January). doi:10.1007/s13398-014-0173-7.2.
- [17] Beheshti I, Demirel H, Neuroimaging D. Probability distribution function-based classification of structural MRI for the detection of Alzheimer's disease. *Comput Biol Med*. 2015;64:208-216. doi:10.1016/j.combiomed.2015.07.006
- [18] Ergin S, Kilic O. A new feature extraction framework based on wavelets for breast cancer diagnosis. *Comput Biol Med*. 2014;51:171-182. doi:10.1016/j.combiomed.2014.05.008