

## **Prediction and Analysis of Structural Brain Health Indicators Using Deep Learning Models with Functional Brain Images as Input**

Sakaki Shimojo and Hiroyuki Akama



v1

May 29, 2023

<https://doi.org/10.32388/RWZH4Y>

# Prediction and Analysis of Structural Brain Health Indicators Using Deep Learning Models with Functional Brain Images as Input

Sakaki Shimojo<sup>1</sup> and Hiroyuki Akama<sup>2, 3, \*</sup>

1. School of Life Science and Technology, Tokyo Institute of Technology, Tokyo, Japan

(Alumnus);

2. School of Life Science and Technology, Tokyo Institute of Technology, Tokyo, Japan;

3. Institute of Liberal Arts, Tokyo Institute of Technology, Tokyo, Japan

\* Corresponding author

## Abstract

There is growing emphasis on the importance of maintaining brain health to prevent disorders such as Alzheimer's disease, and one aspect of this challenge is measuring biomarkers of brain aging using magnetic resonance imaging (MRI). Previous studies have proposed the gray matter brain healthcare quotient (GM-BHQ) as a measure of brain aging and health, which is calculated using gray matter volume obtained from structural images of the brain. However, an index to evaluate brain health considering the functional aspect of the brain is needed, but has not yet been established. This is because resting-state functional connectivity MRI provides multivariate time-series data, which is difficult to reduce to a single feature or scalar like gray matter volume. Therefore, we used a large functional MRI (fMRI) dataset consisting of a wide age range and used the following three approaches: (1) We learned the relationship between resting-state fMRI data and the GM-BHQ, constructed a regression model between them to obtain the predictive value of a model based on functional information as a functional connectivity brain healthcare quotient (FC-BHQ), and tested its utility. (2) We verified the applicability of brain graph neural networks to regression tasks. (3) Finally, we identified brain regions that showed covariations in function and

structure with aging by analyzing the model parameters and interpreting the prediction results. The constructed model achieved moderate performance correlation ( $r > 0.6$ ) between the predictions and correct answers, and the clustering performed inside the model extracted brain regions and networks that reported significant changes with aging. Sparse modeling of the output clusters revealed brain regions strongly associated with the GM-BHQ, such as the amygdala, which is responsible for emotional processing, as well as the Rolandic operculum and superior temporal gyrus, which is characteristic of changes in connectivity with typical dedifferentiation associated with aging.

Keywords: resting-state functional connectivity MRI, brain healthcare quotient, brain graph neural network, group sparse lasso, age prediction

## 1. Introduction

Maintaining brain health is considered important for preventing disorders such as Alzheimer's disease. One aspect of this challenge is measuring the biomarkers of brain aging using medical and magnetic resonance imaging (MRI), also known as neuroimaging. In a previous study (Nemoto et al., 2017), the gray matter brain healthcare quotient (GM-BHQ) was proposed as a measure of brain aging and health. This index is calculated based on the volume of gray matter (GM) obtained from structural MRI images of the brain. The study concluded that the GM-BHQ should be useful as a health indicator because it is not only related to physical health indicators, such as age and BMI, but also to social indicators. However, there is a need for an index to evaluate brain health from the viewpoint of brain function, and no such method has yet been established. Therefore, our study used a large functional MRI (fMRI) dataset consisting of a wide range of age groups from young to old to investigate three aims: (1) to construct a regression model that learns the relationship between resting-state fMRI (rs-fMRI) data containing functional brain information and the GM-BHQ containing information on brain structure and health status; (2) to verify the applicability of graph neural networks (GNNs) to regression tasks using GNNs as regression models; and (3) to identify brain regions and networks that show co-variation of function and structure with aging by analyzing the parameters of the model.

In this study, we focused on the interpretability of the model and used a deep learning framework called BrainGNN (Li et al., 2021), which is a GNN that can extract sophisticated graph representations by performing a graph convolution that considers the positional specificity of brain regions. It also performs simultaneous clustering of brain regions to provide excellent interpretability. Using this approach, we visualized the clustering results and performed sparse modeling considering the cluster structure to identify the brain regions that were important for GM-BHQ prediction. The resulting prediction (named FC-BHQ) by the model using functional data as the input, showed a strong correlation with the true GM-BHQ. Sparse modeling also suggested the importance of functional networks between the amygdala-cortex and cerebellum related to emotional processing in predicting GM-BHQ, supporting the findings of previous studies investigating functional brain changes with aging. In addition to applying GNNs to regression tasks in neuroimaging, we believe our study is the first attempt to identify biomarkers associated with aging by investigating the relationship between brain function and structure using deep learning models.

The brain health quotient (BHQ) is "an index for brain health care, calculated by analyzing brain imaging data utilizing MRI," as proposed by Nemoto et al. (2017). This index consists of the GM-BHQ based on the GM volume (GMV) evaluated by voxel-based morphometry and the fractional anisotropy (FA) BHQ of white matter (WM) evaluated by diffusion tensor imaging. In GM, the desired state of health is considered to be a moderate spread of neuronal dendrites and moderate increase in synapses (Erickson et al., 2014). This state is reflected in the GMV (Ashburner & Friston, 2000) and is interpreted to lead to high synaptic plasticity, indicating learning flexibility (Holtmaat & Svoboda, 2009). In general, the GMV decreases with age. Moreover, GM-BHQ has been found to be negatively correlated with age, physical factors (BMI, blood pressure, and length of rest or relaxation time in daily activities), and social factors (such as socioeconomic status) (Nemoto et al., 2017).

Nemoto et al. (2017) focused on the structural aspects of the brain, such as GM and WM, but mentioned another possible coefficient of the BHQ that could incorporate functional information through the use of rs-fMRI. It has been suggested that rs-fMRI can provide a measure of the extent to which brain networks function. However,

because rs-fMRI is multivariate time-series data containing information on the entire brain, it is difficult to reduce it to a single feature, such as a scalar or GMV. Moreover, in such cases, the coefficient derived from rs-fMRI may lack universality and stability, which would deter local health status analysis. Therefore, in this study, we attempted to solve the above problem by defining FC-BHQ as the predicted value of GM-BHQ based on a deep learning model that uses a functional graph of the entire brain extracted from rs-fMRI data as the input. We hypothesized that this would make it possible to connect GM-BHQ information to functional information in the brain. In particular, we used GM data from the Nathan Kline Institute-Rockland Sample (NKI-RS) dataset and analyzed hidden layers of the deep learning model. Therefore, we identified brain regions and networks in which age-related covariations in structure and function most frequently appeared, and tested our hypothesis by comparing it with previous findings to investigate the usefulness of our FC-BHQ.

## 2. Materials and Methods

### 2.1. Overview

The procedure used in this study was as follows. First, we preprocessed the rs-fMRI time series and GM images for each subject obtained from the NKI-RS dataset; GM-BHQ was computed from the GM images, and the functional brain graph  $\mathcal{G}$ , as the input to the regression model, was constructed from the rs-fMRI time series. The BrainGNN used as the regression model in this study learned the mapping  $f: \mathcal{G} \mapsto y$  from the functional brain graph  $\mathcal{G}$  to the GM-BHQ  $y$  and output the predicted value  $\hat{y}$ . Henceforth, we refer to the output  $\hat{y}$  of BrainGNN as FC-BHQ using a BHQ prediction based on functional data.

### 2.2. Materials

#### 2.2.1. NKI-RS dataset

The dataset for this study was obtained from the enhanced NKI-RS, which is publicly available online. Several study codes were included in the publicly available data. We focused on structural MRI images and rs-fMRI (repetition time (TR) = 1400 ms,

multiband) data from subjects included in Baseline Visit BAS1 (one baseline visit). The dataset consisted of data from 1246 subjects (495 males, 750 females, and one unknown). The mean age was 39.2 years (standard deviation 21.7 years), and ranged from 6 to 85 years. MRI images were acquired using a 3.0 Tesla SIEMENS Trio Tim scanner with a 32-channel head coil. Structural MRI images were acquired using the magnetization-prepared rapid gradient-echo (MPRAGE) sequence with the following scan parameters: TR = 1900 ms, voxel size = 1 mm isotropic, time to echo (TE) = 2.52 ms, flip angle (FA) = 9°, thickness = 1.0 mm, slices = 176, matrix =  $256 \times 256$ , field of view (FOV) =  $256 \times 256$  mm. The rs-fMRI images were also acquired using an echo-planar imaging sequence, with the following scan parameters: TR = 1400 ms, TE = 30 ms, FA = 65°, FOV = 224 mm, matrix =  $112 \times 112$ , slices = 64, thickness = 2.0 mm, and volume = 404.

### 2.2.2. Preprocessing

The MRI images were preprocessed using the default (as of January 2020) configurable pipeline for analysis (<https://fcp-indi.github.io/>). Structural MRI preprocessing was performed as follows: (1) skull stripping using AFNI 3dSkullStrip; (2) registration to the Montreal Neurological Institute (MNI) 152 standard space with advanced normalization tools; and (3) GM, WM, and spinal fluid segmentation using FSL FAST. The GM images obtained in the above procedures were used for voxel-based morphometry, as described in Section 2.1.4.

The rs-fMRI images were preprocessed as follows: (1) slice timing correction, (2) functional-anatomical registration using the boundary-based registration method with AFNI 3dAutoMask, (3) registration to the MNI152 standard space, and (4) nuisance regression. The nuisance regression procedure was performed as follows: (4.1) temporal filtering, (4.2) cerebrospinal fluid regression, (4.3) global signal regression, (4.4) regression of motion parameters, (4.5) polynomial detrending, and (4.6) component-based noise reduction (aCompCor).

The region of interest (ROI) time series were extracted by averaging the blood-oxygen-level-dependent signals of voxels in 116 ROIs defined by automated

anatomical labeling (AAL) on the preprocessed fMRI images obtained from the above procedures.

## 2.3. Methods

### 2.3.1. Computational Environment

In this study, the environment was built using PyTorch and Pytorch Geometrics on a Docker with two GTX1080ti GPUs (11 GB VRAM) to implement the deep learning models based on the GitHub repository published by Li et al. (2021) ([https://github.com/xxlya/BrainGNN\\_Pytorch](https://github.com/xxlya/BrainGNN_Pytorch)). We implemented the sparse modeling code to evaluate the contribution of each node to the regression using scikit-learn (<https://scikit-learn.org/stable/>) and group Lasso (<https://github.com/yngvem/group-lasso/blob/master/docs/index.rst>). All the codes used in this research are available to the public through the following URL: <https://github.com/Skk5mj/masterthesis/>

### 2.3.2. Brain Graph Construction

In this section, we describe the construction of the brain graphs that were input into the GNN model. To generate a brain graph  $\mathcal{G}$ , it is necessary to calculate the feature  $h_i^{(0)}$  of node  $v_i$  corresponding to a certain ROI <sub>$i$</sub>  and the functional connection between  $v_i$  and  $v_j$ , that is, the weights of the edges. In this study, the Pearson correlation coefficients between nodes were used as node features and thresholded partial correlation coefficients were used for the edges. The computations for both indices are as follows:

$$\text{Pearson}(x, y) = \frac{\sum_{m=1}^M \left( x_m - \frac{1}{M} \sum_{m=1}^M x_m \right) \left( y_m - \frac{1}{M} \sum_{m=1}^M y_m \right)}{\sqrt{\sum_{m=1}^M \left( x_m - \frac{1}{M} \sum_{m=1}^M x_m \right)^2} \sqrt{\sum_{m=1}^M \left( y_m - \frac{1}{M} \sum_{m=1}^M y_m \right)^2}}, \quad (1)$$

$$\text{Partial}(x_1, x_2, \text{rest}) = \text{Pearson}(\boldsymbol{\varepsilon}_1, \boldsymbol{\varepsilon}_2). \quad (2)$$

Note that  $\boldsymbol{\varepsilon}_1$  and  $\boldsymbol{\varepsilon}_2$  represent the residuals from the linear regression of  $x_1$  and  $x_2$ , respectively, on the other factors. Positive values in the top 10% were selected to guarantee the absence of isolated nodes.

### 2.3.3. BHQ Computational Methods

#### GM-BHQ

The GM-BHQ calculation in this study was partially simplified from the procedure described by Nemoto et al. (2017). After calculating the total GMV for each subject from the segmented GM-masked images using the procedure described in Section 2.2.2, the GMV of all subjects were standardized. Finally, GM-BHQ was obtained by transforming the standardized GMV of each subject by defining the mean GM-BHQ of the subjects as 100 points and the standard deviation as 15 points, as in the calculation of intelligence quotient. This operation can be represented by the following equation:

$$\text{GM-BHQ} = 100 + 15 \times (\text{GMV}_{\text{individual}} - \text{mean}(\text{GMV}))/\text{std}(\text{GMV}). \quad (3)$$

#### FC-BHQ

In this study, we defined FC-BHQ as the value output by BrainGNN trained on a regression model  $f$  that maps the functional brain graph  $\mathcal{G}$  to GM-BHQ. Therefore, FC-BHQ is a function-based BHQ based on a model that learns the relationship between rs-fMRI and GMV. Because any BHQ must, by definition, have a mean of 100 and standard deviation of 15, standardization was applied to the correct answer  $y$  used to train GM-BHQ; therefore, the output value  $\hat{y}$  of the model was transformed inversely to the standardization to obtain the final predicted value FC-BHQ. This procedure is summarized by the following equation:

$$\text{FC-BHQ} = 100 + 15 \times (\hat{y} - \text{mean}(\text{GM-BHQ}))/\text{std}(\text{GM-BHQ}),$$

where  $\hat{y}$  is the value output from regression model BrainGNN that learned the mapping  $f(f: \mathcal{G} \mapsto y)$ .

### 2.3.4. BrainGNN

#### Overview

BrainGNN is a framework of the GNN proposed by Li et al. (2021), which was based on the concept of extracting graph structures from fMRI data to discover neuroscientific biomarkers. BrainGNN is characterized by the following three elements: (1) an ROI-aware graph convolutional (Ra-GConv) layer that utilizes fMRI topology and



functional information considering the location of ROIs in the brain graph, (2) an ROI selection pooling layer or R-pool layer, and (3) a loss function for the pooling results. In BrainGNN, modules (1) and (2) are designed as a single block for end-to-end learning. We followed the method described by Li et al. (2021) with a few limitations.

BrainGNN applies convolution operations to the node features of the input graph to update them. This procedure is performed in a single block. The convolution kernel (convolution weights) is the output from a two-layer multilayer perceptron (MLP) that uses node location information as an input. This allows for the use of a convolutional kernel that considers the location characteristics of the nodes. The nodes are then downsampled in the pooling layer. The new graph representation is used as the input to the next block. The graph representation obtained in this block is also summarized by maximum and average pooling to preserve information. Finally, the summarized information from each block is aggregated into a single vector, which is input into the MLP for the final regression. This last step differs from that of Li et al. (2021), in which the final output is a class.

#### Ra-GConv layer

In the operations of the Ra-GConv layer, the location information of the nodes is input to the 2-layer MLP that trains the convolutional kernel; the first layer uses this location information to calculate a membership score that indicates the degree to which each node belongs and assigns soft clusters based on that score. This allows for an interpretation of how the model understands the similarity between nodes. Node features are embedded based on the convolution kernel output by the MLP, and a new node representation is obtained.

In the two hidden layers of Ra-GConv, soft clustering of node  $v_i$  and a linear transformation of the membership score vector are performed, attributing node  $v_i$  to  $K$  distinct clusters. Thus, in a model that predicts structure from brain function, it is possible to analyze the weights of all ROIs that contribute to the prediction in relation to each other. This section details this mechanism. The graph convolution operation by the  $l^{\text{th}}$  Ra-GConv layer in the forward propagation is expressed as

$$\tilde{\mathbf{h}}_i^{(l+1)} = \text{relu}\left(W_i^{(l)} h_i^{(l)} + \sum_{j \in \mathcal{N}^{(l)}(i)} e_{ij} W_j^{(l)} h_j^{(l)}\right), \quad (4)$$

where  $h_i \in R^d$  is the feature value of node  $v_i$  before pooling,  $W_i^{(l)} \in R^{d^{(l+1)} \times d^{(l)}}$  is the graph kernel for node  $v_i$  in the  $l^{\text{th}}$  layer, and  $e_{ij}$  corresponds to the edge weight between adjacent nodes  $v_i$  and  $v_j$  ( $j \in \mathcal{N}(i)$ ). The convolution kernel  $W_i$  is learned by considering the specificity of each brain region (ROI) in the brain graph input in each block, and the ROI location information is represented by one-hot encoding vectors  $\mathbf{r}_i$ , which are input and embedded into the two-layer MLP. The output comprises a vector and is formatted such that it can be used in a convolution kernel to obtain  $W_i$ . The following is a mathematical formula illustrating this process:

$$\text{vec}(W_i) = f_{\text{MLP}}(\mathbf{r}_i) = \Theta_2 \text{relu}(\Theta_1 \mathbf{r}_i) + \mathbf{b}. \quad (5)$$

Equation (5) can be expressed as follows: First, let  $d^{(l)} = N^{(l)}$  be the number of nodes in the graph input to the  $l^{\text{th}}$  block and replace the weights of the first and second hidden layers of the MLP with the following respective notation:

$$\Theta_1^{(l)} = [\boldsymbol{\alpha}_1^{(l)}, \dots, \boldsymbol{\alpha}_{N^{(l)}}^{(l)}],$$

$$\Theta_2^{(l)} = [\boldsymbol{\beta}_1^{(l)}, \dots, \boldsymbol{\beta}_{K^{(l)}}^{(l)}].$$

Then, at the first hidden layer, the following conversions are performed.

$$\begin{aligned} \text{relu}\left(\Theta_1^{(l)} r_i^{(l)}\right) &= \text{relu}\left(\boldsymbol{\alpha}_i^{(l)}\right) = \text{relu}\left([\alpha_{i1}^{(l)}, \dots, \alpha_{iK^{(l)}}^{(l)}]^\top\right) \\ &= \left[(\alpha_{i1}^{(l)})^+, \dots, (\alpha_{iK^{(l)}}^{(l)})^+\right]^\top, \end{aligned}$$

where  $(\alpha_{iu}^{(l)})^+$  denotes the non-negative membership score of node  $v_i$  belonging to cluster  $u$  as follows:

$$(\alpha_{iu}^{(l)})^+ = \begin{cases} \alpha_{iu}^{(l)}, & (\text{if } \alpha_{iu}^{(l)} > 0) \\ 0, & (\text{if } \alpha_{iu}^{(l)} \leq 0). \end{cases}$$

From the above equation, soft clustering of node  $v_i$  is performed in the first hidden layer. The membership score vector is then linearly transformed in the second hidden

layer; thus, Equation (5) can be expressed as

$$\text{vec}\left(W_i^{(l)}\right) = \sum_{u=1}^{K^{(l)}} \left(\alpha_{iu}^{(l)}\right)^+ \boldsymbol{\beta}_u^{(l)} + \boldsymbol{b}^{(l)}. \quad (6)$$

In the above procedure, the convolution kernel is trained differently for each node  $v_i$ , which is simultaneously soft-clustered into  $K$  clusters.

#### R-pool Layer and Readout Layer

After the Ra-GConv layer, the R-pool Layer performs a pooling operation to maintain the important nodes based on the new graph representation output from the previous layer. Each node is mapped to a pooling score vector based on node features and assigned a score. In the readout layer, the node feature matrix of the newly obtained graph is flattened into vectors by pooling to preserve information. These summary vectors are further combined and used as inputs for the regression predictor. The details are as provided in Li et al. (2021) and omitted from this paper.

The model in this study was implemented by modifying parts of the BrainGNN code published on GitHub by Li et al. (2021). The number of blocks was set to  $L = 3$ , and the pooling ratio of the R-pool layer was set to 0.5, such that the number of nodes in each block was downsampled by half. The regularization parameter  $\lambda$  in the loss function (Equation 7) was set to 0.1.

#### 2.3.5. Loss Functions

To perform the regression task in this study, which was different from that in the Li et al. (2021) original model, we used the following loss function  $L_{\text{total}}$  to perform the regression task.

$$L_{\text{total}} = L_{\text{MSE}} + \lambda(L_{\text{unit}} + L_{\text{topk}}), \quad (7)$$

where  $\lambda$  is a regularization parameter that adjusts the importance of the loss function. The definitions of the terms on the right-hand side are as follows:

- Mean Squared Error (MSE) Loss

$$L_{\text{MSE}} = \sum_{m=1}^M \frac{1}{2M} (y_m - \hat{y}_m)^2, \quad (8)$$

where  $y_m$  and  $\hat{y}_m$  are the true and predicted values, respectively, of the true objective variables for subject  $m$ .

- Unit Loss Function

$$L_{\text{unit}} = \sum_{l=1}^L \left( \|\mathbf{w}^{(l)}\|_2 - 1 \right)^2 \quad (9)$$

This function provides a constraint on the learnable vector  $\mathbf{w}^{(l)} \in \mathbb{R}^{d^{(l)}}$  that projects the node features to obtain the pooling score vector  $\mathbf{s}^{(l)}$  in the R-pool layer.  $\mathbf{s}^{(l)}$  can be arbitrarily scaled using a real number  $a (\neq 0)$  as

$$\mathbf{s}^{(l)} = \tilde{H}^{(l+1)}(a\mathbf{w}^{(l)}) / \|a\mathbf{w}^{(l)}\|_2.$$

In other words,  $\mathbf{w}^{(l)}$  can assume arbitrary values and is no longer uniquely determined, which can render the learning process unstable. Therefore, the unit loss function  $L_{\text{unit}}$  adds the restriction that  $\mathbf{w}^{(l)}$  is a unit vector.

- Top-k Loss Function

$$L_{\text{topk}} = \sum_{l=1}^L \left\{ -\frac{1}{M} \sum_{m=1}^M \frac{1}{N^{(l)}} \left( \sum_{i=1}^{N^{(l)}k} \log(\hat{s}_{m,i}^{(l)}) \right. \right. \\ \left. \left. + \sum_{i=1}^{N^{(l)}(1-k)} \log(1 - \hat{s}_{m,i+N^{(l)}k}^{(l)}) \right) \right\} \quad (10)$$

The top-k loss function is intended to constrain the pooling score, such that the  $N^{(l)}k$  beneficial ROIs that are useful for prediction have scores that are distinctly farther apart than those of the unselected nodes. To achieve this, the top-k loss function ranks  $\hat{\mathbf{s}}_m^{(l)} = \text{sigmoid}(\tilde{\mathbf{s}}^{(l)})$  for the  $m$ -th instance in descending order and uses the binary cross-entropy function.

### 2.3.6. Learning and Cross-Validation

To explore the relationship between the brain function network and GM-BHQ,

the parameter to be adjusted was limited to  $K^{(l)}$ , which represents the number of clusters in the soft clustering performed in the Ra-GConv layer. The value of  $K^{(l)}$  was fixed in Li et al. (2021), all blocks ( $l = 1, 2, 3$ ) were fixed at eight, and no search for  $K^{(l)}$  was performed. Therefore, assuming that clustering yields clusters of known functional brain networks and that the performance of the model depends on the number of clusters, we conducted a search in the range  $K^{(l)} = \{7, 8, 9, 10\}$ . Hereafter, we refer to these models as Models 1 ( $K^{(l)}=7$ ), 2 ( $K^{(l)}=8$ ), 3 ( $K^{(l)}=9$ ), and 4 ( $K^{(l)}=10$ ) in order of  $K^{(l)}$  value. The models were evaluated using five-fold cross-validation (5-fold CV).

In the 5-fold CV, the data were first divided into five blocks: one for training and the others for testing the performance of the model after training. As the training progressed, the model was approximately fitted to the training data, and training was terminated when the prediction error  $L_{\text{total}}$  with respect to the test data showed no improvement for more than 10 epochs after the minimum value was recorded. At the epoch when the prediction error reached a minimum, the mean absolute error (MAE) between the model's predictions and true values was calculated, and this was used as the prediction accuracy in Fold 1. After Fold 1 was completed, the training and prediction of the test data were repeated in the same manner for Fold 2, and so on, to obtain the MAE of the test data five times. The values were averaged to determine the model with the best number of clusters for parameter  $K^{(l)}$ .

### 2.3.7. Group Sparse Lasso Applied to the Ra-GConv Layer

As mentioned earlier, when training the convolution kernel used in the Ra-GConv layer, a matrix of nonnegative membership scores ( $[(\alpha_{iu})^+]$ ), indicating the degree to which  $\text{ROI}_i$  belongs to cluster  $u$ , was obtained. Based on the membership scores in the first block of each model, we computed soft clustering in which each ROI simultaneously belonged to several clusters. We also performed hard clustering in which each ROI was assigned only to the cluster with the highest score.

Based on the clusters shown by BrainGNN, sparse modeling by group-sparse Lasso (Simon et al., 2013) was performed to interpret which functional indicators of brain regions were useful in predicting GM-BHQ as the objective variable. Here, not

only the information of each ROI but also the subject's personal information, such as gender and age, were included as explanatory variables. Subject information was added as a confounding factor to account for its contribution to the structural brain characteristics (GM-BHQ). Because the dimensions of the explanatory variables used for sparse modeling should be smaller than the number of samples (Cui & Gong, 2018), we used functional connectivity strength (FCS) as a functional measure. FCS was calculated for each ROI and corresponds to the centrality measure of the graph theory indicators. FCS( $i$ ) of ROI $_i$  was obtained by standardizing the Pearson correlation  $r_{ij}$  between all other ROI $_j$  ( $j \neq i$ ), using thresholding ( $> 0.2$ ) and summation (Li et al., 2021). The equation for the linear regression is given by

$$y_m = \boldsymbol{\beta}_0^T \mathbf{x}_{m,0} + \boldsymbol{\beta}_1^T \mathbf{x}_{m,1} + \boldsymbol{\beta}_2^T \mathbf{x}_{m,2} + \cdots + \boldsymbol{\beta}_K^T \mathbf{x}_{m,K},$$

where  $y_m$  is the objective variable of the subject  $m$ ,  $\boldsymbol{\beta}_0$  and  $\mathbf{x}_{m,0}$  are the partial regression coefficient and explanatory variable vectors, respectively, for the information (age and sex) of the subject  $m$ , and  $\boldsymbol{\beta}_{m,u}$  and  $\mathbf{x}_{m,u}$  ( $u = 1, \dots, K$ ) are the grouped partial regression coefficient and feature vectors, respectively. Standardization was applied to these features in advance.

Group-sparse Lasso is a regularization method that combines a group Lasso and Lasso. It minimizes the objective function that constrains the L1 and L2 norms of the partial regression coefficient vector while considering the clusters created by the elements of the explanatory variables. The objective function is expressed by the following equation.

$$S_\lambda(\boldsymbol{\beta}) = \frac{1}{2M} \|y - X\boldsymbol{\beta}\|_2^2 + \alpha\lambda \|\boldsymbol{\beta}\|_1 + (1 - \alpha)\lambda \sum_{u \in K} \sqrt{p_u} \|\boldsymbol{\beta}_u\|_2,$$

where the subscript  $u$  denotes the cluster number and  $K$  is the set of subscripts. The dimensions of the partial regression coefficient vector  $\boldsymbol{\beta}_u$  corresponding to each cluster are denoted by  $p_u$ . The second term on the right-hand side represents the constraint on the partial regression coefficients by the Lasso, and the third term refers to the constraint on the partial regression coefficients for each group by the group Lasso. These terms are controlled by the parameter  $\alpha$  ( $\in [0, 1]$ ). These regularization terms degenerate the estimates of some partial regression coefficients to zero, and explanatory variables with coefficients estimated to be zero can be interpreted as not contributing to

the objective variable. In particular, the third term promotes a group-wise reduction of the partial regression coefficients such that those for groups formed by the explanatory variables that have a low contribution are degenerated to zero. This allows for sparsity and may allow for a more concise interpretation of the explanatory and objective variables.

We performed 10-fold CV for each value of  $\lambda$ , moving it in the range  $[10^{-4}, 10^0]$  for  $\alpha = \alpha_{\text{fixed}}$  fixed between  $[0, 1]$ . The evaluation criterion was the mean between the folds of  $\text{MSE}(y_m, \hat{y}_m)$ , and the model with the smallest mean ( $\lambda = \lambda_{\text{best}}$ ) was selected. For  $\alpha$ , the search was performed in the range  $\alpha = \{0, 0.25, 0.5, 0.75, 1\}$ . When  $\alpha = 0$ , only the group-Lasso effect was applied to the objective function to be minimized. However, when  $\alpha = 1$ , only the effect of the Lasso was applied. In other words, the smaller the value of  $\alpha$ , the more importance was placed on the group structure of the explanatory variables; the larger the value of  $\alpha$ , the more importance was placed on the contribution of the individual explanatory variables. Using the above methods, we searched for the best parameter combination  $(\alpha_{\text{fixed}}, \lambda_{\text{best}})$  and examined the important features that contributed to the prediction of GM-BHQ by evaluating the average value of the partial regression coefficient for 10 models with 10-fold CV.

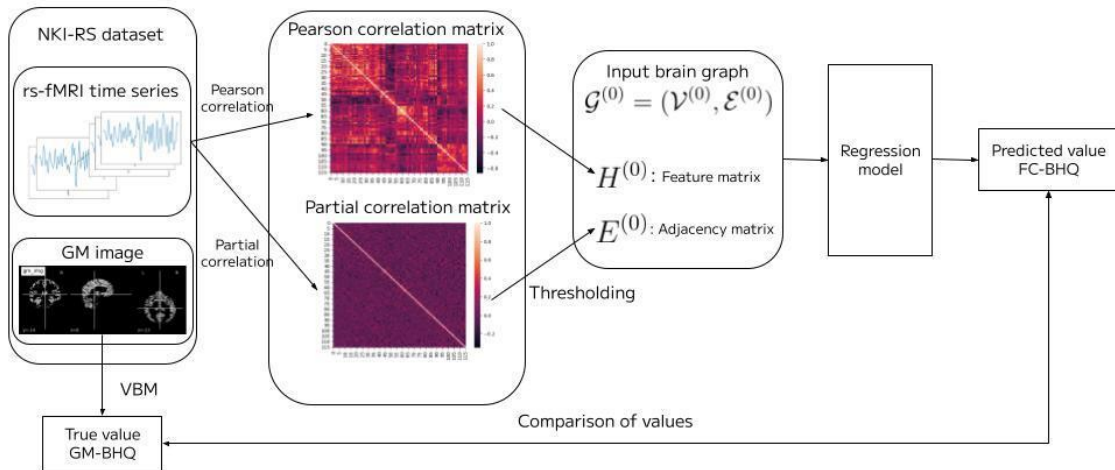


Figure 1. Overview of analyses. The rs-fMRI time series and GM images for each subject obtained from the NKI-RS dataset were preprocessed, GM-BHQ was computed from the GM

images, and the functional brain graph  $\mathcal{G}$  was constructed from the rs-fMRI time series. BrainGNN used as the regression model in this study learns the mapping  $f: \mathcal{G} \mapsto y$  from the functional brain graph  $\mathcal{G}$  to GM-BHQ  $y$  and outputs predicted value  $\hat{y}$ . The output  $\hat{y}$  of BrainGNN is referred to as FC-BHQ, the BHQ predicted from functional data.

### 3. Results

#### 3.1. Relationship Between GM-BHQ and Age

The following are the results of linear and polynomial regressions using the least squares method to explore the relationship between the GM-BHQ calculated from GM images and age. The order of the polynomial regression was determined based on the Akaike information criterion (AIC).

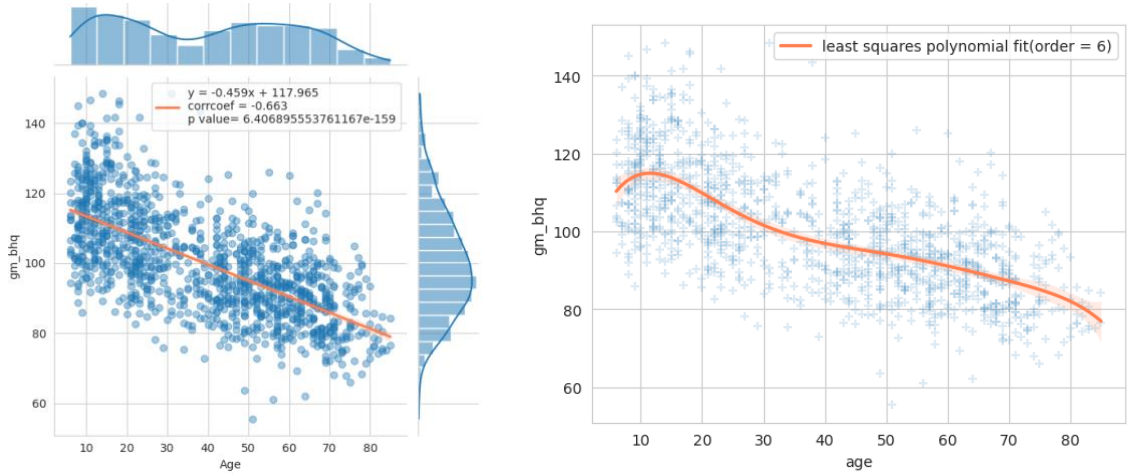


Figure 2. Relationship between GM-BHQ and age. The vertical and horizontal axes correspond to GM-BHQ and age, respectively. (Left) Regression lines (orange) were obtained by the least squares method. (Right) Results of curve fitting between GM-BHQ and age. The order of the polynomial regression curve (orange) was set to six based on the AIC.

GM-BHQ showed a strong negative correlation with age ( $r = -0.663$ ), reflecting age-related changes in GM volume. Polynomial regression revealed a curve that showed a slight increase until approaching the mid-teens and then a decreasing curve with an almost constant slope; after 60 years, the GM-BHQ values with respect to age decreased



faster.

### 3.2. GM-BHQ Prediction by BrainGNN

As stated earlier, we varied the number of clusters  $K$  in the Ra-GConv layer and compared the accuracies of the models.

Table 1: Prediction accuracy of each model by 5-fold CV

Models	Mean MAE
Model 1 ( $K = 7$ )	$9.149 \pm 0.197$
Model 2 ( $K = 8$ )	$9.359 \pm 0.190$
Model 3 ( $K = 9$ )	$9.302 \pm 0.124$
Model 4 ( $K = 10$ )	$9.186 \pm 0.191$

Among the four models, the MAE of Model 1 ( $K^{(l)}=7$ ) was the smallest, followed by that of Model 4 ( $K^{(l)}=10$ ). Therefore, only the results of Model 1 are discussed.

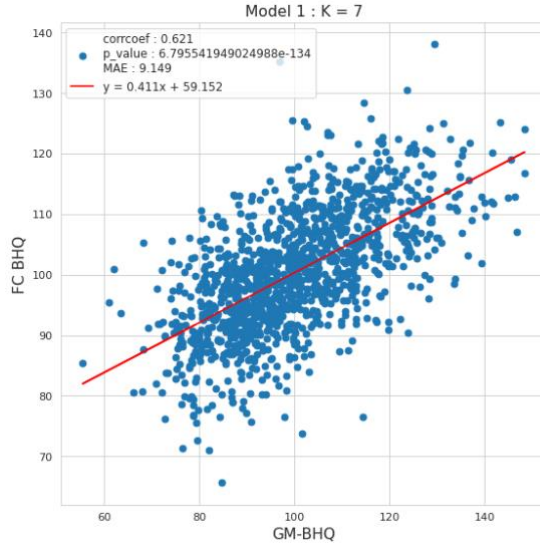


Figure 3. Prediction results for Model 1 with  $K=7$  clusters. The vertical axis corresponds to the predicted value of BrainGNN (FC-BHQ), and the horizontal axis corresponds to GM-BHQ. The regression lines shown in red were obtained by the least squares method. All p-values calculated by correlation analysis are uncorrected.

Again, it is necessary to emphasize that FC-BHQ refers to the GM-BHQ

predicted using the rs-fMRI information of the participants. The results predicted by Model 1 are shown in Figure 3. A correlation coefficient of 0.629 (p-value uncorrected) indicates that the model is highly significant; however, it underestimates GM-BHQ above 120 and overestimates it below 80, indicating that the prediction accuracy decreases significantly toward the base of the distribution. The predictions of the GM-BHQ using the BrainGNN model show that FC-BHQ reproduces the GM-BHQ features well.

Figure 4 shows the results of the polynomial regression using the least squares method to explore the relationship between the FC-BHQ and age. The polynomial regression curve (shown in orange) was determined based on the AIC and was almost identical to that obtained by curve fitting for age and GM-BHQ. The shape of the curve was similar to those of the other models for different  $K$  values, with a peak in the early teens, decrease until the 50s, plateau in the 50-70s, and a decrease thereafter.

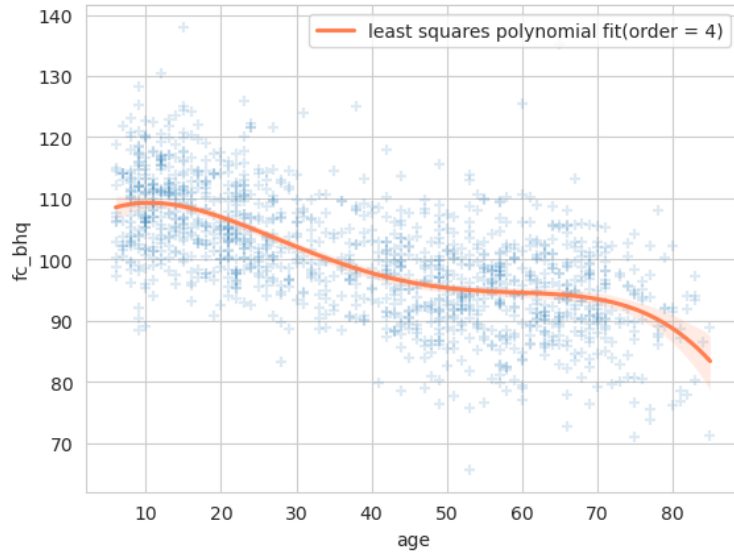


Figure 4. Results of curve fitting of FC-BHQ and age by Model 1 with  $K = 7$  clusters. The horizontal and vertical axes correspond to FC-BHQ and age, respectively. Curves from polynomial regression (orange) were determined for orders based on the AIC.

### 3.3. Results of Clustering and Sparse Modeling for ROI Evaluation

The results of the 10-fold CV for each value of  $\lambda$  in the range  $[10^{-4}, 10^0]$  for

$\alpha = \alpha_{\text{fixed}}$  fixed between  $[0, 1]$ , are now described. As shown in Figure 5, there was no significant difference in the mean MSE for different  $\alpha_{\text{fixed}}$ . This was also the case for  $\alpha_{\text{fixed}} = 1$ , that is, a simple Lasso 10-fold CV that did not consider the group structure of the explanatory variables.

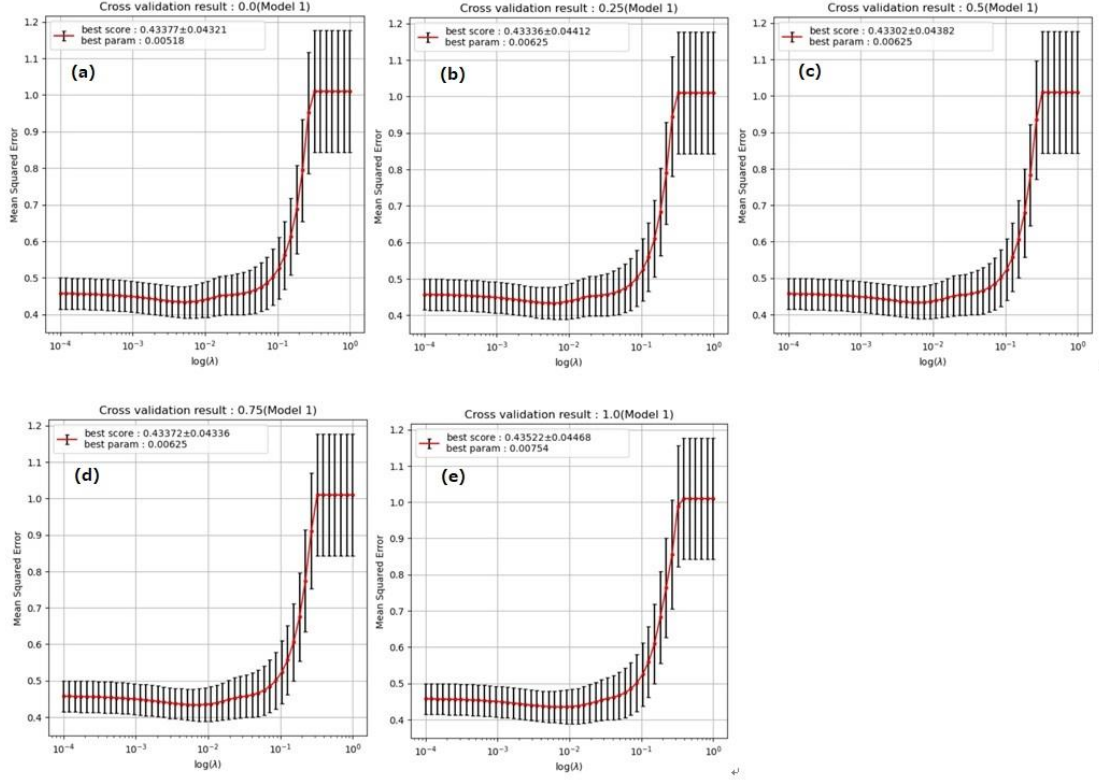


Figure 5. Results for group sparse Lasso (a-d) and Lasso 10-fold CV (e) considering the clustering results in Model 1: the horizontal axis shows the value of  $\lambda$  on a logarithmic scale, and the vertical axis is the average value of the mean squared errors (MAE) of the 10-fold CV performed for each value of  $\lambda$ . The mean is shown in red, and the black error bars represent the standard deviation.

We present the results of soft clustering in the Ra-GConv layer of the first block of Model 1. The distance between clusters was calculated using cosine similarity, and the resulting matrix was visualized on a two-dimensional plane using T-distributed stochastic neighbor embedding (t-SNE). As shown in Figure 6, the soft clustering results show that the clusters are clearly separated, which confirms the significance of the ROI-based convolutional kernel learning mechanism.

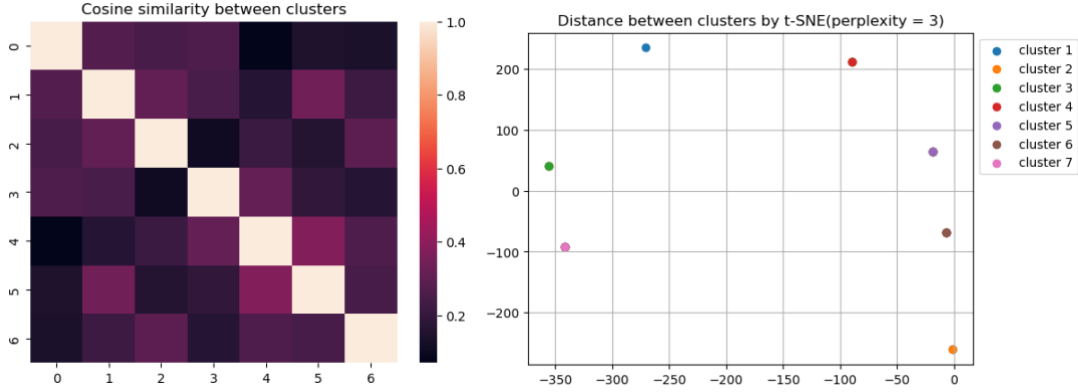


Figure 6. (Left) Heatmap showing distances between seven clusters using cosine similarity based on soft clustering of Ra-GConv layer values. (Right) Each cluster mapped on a two-dimensional plane using t-SNE.

The results of the hard clustering, in which each ROI was assigned to the cluster with the largest membership score, showed some variation depending on the computational conditions. However, the results were stable and similar for the top ROI contributing to the prediction. The names of the ROIs were based on AAL. Note that the 10-fold CV of the sparse group Lasso did not differ significantly among the models with respect to the behavior of the loss function; therefore, we only show the results for Model 1, which recorded the best performance for predicting GM-BHQ. The parameter  $\alpha_{\text{fixed}}$  was set to  $\alpha = \{0, 0.25, 0.5, 0.75, 1\}$ , as described above, where  $\alpha = 0$  means that only the group Lasso effect is applied in the objective function to be minimized. However, when  $\alpha = 1$ , only the effect of Lasso is applied and the group structure of the explanatory variables is not considered. When  $\alpha = 0$ , the sparse group Lasso coincides with the group Lasso and there is no partial regression coefficient that reduces to zero; this is also the case for  $\alpha = 0.25$ . Thereafter, there was an increase in the number of coefficients degenerated to 0 as the value of  $\alpha$  increased. There were six coefficients at  $\alpha = 0.5$ , 13 at  $\alpha = 0.75$ , and 32 at  $\alpha = 1$ , which is consistent with the Lasso.

Table 2 lists the absolute values of the partial regression coefficients under each  $\alpha$  value in increasing order. Figure 7 also shows the sum of the ranks under each condition in ascending order, with left amygdala (Amygdala\_L) and left Rolandic

operculum (Rolandic\_Oper\_L) in first and second places, respectively, under all conditions. In addition, the medial orbital parts of the right superior frontal gyrus (Frontal\_Med\_Orb\_L), right superior temporal gyrus (Temporal\_Sup\_R), and left cuneus (Cuneus\_L) recorded mean ranks in the top five, followed by the right hippocampus (Hippocampus\_R). Considering the hard clustering, it can be seen that the top-ranking AAL regions are concentrated in cluster numbers 1, 2, 3, and 6.

Table 2

The twenty areas recording the highest absolute partial regression coefficients under each of the  $\alpha$  values {0, 0.25, 0.5, 0.75, 1}.

$\alpha=0$ (cluster Lasso)			$\alpha=0.25$ (Sparse cluster Lasso)					
feature (ROI)	cluster_	$ \beta $	feature (ROI)	cluster_	$ \beta $			
Amygdala_L	3	0.053844	Amygdala_L	3	0.051917			
Rolandic_Oper_L	6	0.034503	Rolandic_Oper_L	6	0.031943			
Frontal_Med_Orb_R	3	0.03282	Frontal_Med_Orb_R	3	0.031875			
Cuneus_L	2	0.028488	Cuneus_L	2	0.026348			
Temporal_Sup_R	1	0.028117	Hippocampus_R	2	0.025015			
Hippocampus_R	2	0.026806	Temporal_Sup_R	1	0.024917			
Cerebelum_10_L	1	0.026343	Olfactory_L	3	0.024084			
Calcarine_L	6	0.025434	Cerebelum_10_L	1	0.023961			
Olfactory_L	3	0.025382	Vermis_7	2	0.02319			
Vermis_7	2	0.025124	Frontal_Sup_Orb_L	6	0.022916			
Frontal_Mid_Orb_L	2	0.024584	Frontal_Mid_Orb_L	2	0.022677			
Frontal_Sup_Orb_L	6	0.023416	Calcarine_L	6	0.022309			
Caudate_R	2	0.023351	Caudate_R	2	0.020365			
Occipital_Mid_R	1	0.020866	Occipital_Mid_R	1	0.017747			
Occipital_Inf_L	1	0.020318	Frontal_Inf_Oper_R	5	0.017433			
Frontal_Inf_Oper_R	5	0.019714	Temporal_Pole_Mid_R	5	0.017231			
Frontal_Inf_Orb_R	3	0.019403	Frontal_Inf_Orb_R	3	0.017103			
Frontal_Mid_R	2	0.019359	Occipital_Inf_L	1	0.016938			
Angular_L	6	0.018916	Angular_L	6	0.016888			
Temporal_Pole_Mid_R	5	0.018682	Heschl_R	6	0.015385			
$\alpha=0.5$ (Sparse cluster Lasso)			$\alpha=0.75$ (Sparse cluster Lasso)			$\alpha=1$ (Lasso)		
feature (ROI)	cluster_	$ \beta $	feature (ROI)	cluster_	$ \beta $	feature (ROI)	cluster_	$ \beta $
Amygdala_L	3	0.057708	Amygdala_L	3	0.063672	Amygdala_L	3	0.065531
Rolandic_Oper_L	6	0.038061	Rolandic_Oper_L	6	0.04747	Rolandic_Oper_L	6	0.056989
Frontal_Med_Orb_R	3	0.034758	Frontal_Med_Orb_R	3	0.037739	Temporal_Sup_R	1	0.041683
Temporal_Sup_R	1	0.030354	Temporal_Sup_R	1	0.037132	Hippocampus_R	2	0.039831
Cuneus_L	2	0.030264	Cuneus_L	2	0.035522	Cerebelum_10_L	1	0.038953
Hippocampus_R	2	0.028855	Hippocampus_R	2	0.034235	Frontal_Med_Orb_R	3	0.03881
Cerebelum_10_L	1	0.02861	Cerebelum_10_L	1	0.034154	Cuneus_L	2	0.037834
Olfactory_L	3	0.025281	Frontal_Mid_Orb_L	2	0.027958	Frontal_Mid_Orb_L	2	0.028788
Frontal_Mid_Orb_L	2	0.025094	Vermis_7	2	0.026775	Temporal_Pole_Mid_L	5	0.02772
Vermis_7	2	0.025	Calcarine_L	6	0.026171	Frontal_Inf_Oper_R	5	0.026839
Calcarine_L	6	0.024197	Olfactory_L	3	0.026078	Vermis_7	2	0.026604
Frontal_Sup_Orb_L	6	0.024088	Frontal_Inf_Oper_R	5	0.025333	Frontal_Sup_Medial_L	7	0.025314
Caudate_R	2	0.02197	Frontal_Sup_Orb_L	6	0.02506	Frontal_Sup_Orb_L	6	0.024523
Frontal_Inf_Oper_R	5	0.020786	Occipital_Inf_L	1	0.024427	Olfactory_L	3	0.022957
Occipital_Mid_R	1	0.020327	Temporal_Pole_Mid_R	5	0.023961	Calcarine_L	6	0.022486
Occipital_Inf_L	1	0.020257	Caudate_R	2	0.023675	Occipital_Inf_L	1	0.022342
Temporal_Pole_Mid_R	5	0.019889	Occipital_Mid_R	1	0.022493	Caudate_R	2	0.021725
Frontal_Sup_Medial_L	7	0.017205	Frontal_Sup_Medial_L	7	0.021499	Heschl_R	6	0.020727
Heschl_R	6	0.017171	Heschl_R	6	0.020022	Occipital_Sup_R	4	0.020512
Angular_L	6	0.016733	Lingual_L	4	0.019192	Lingual_L	4	0.019717

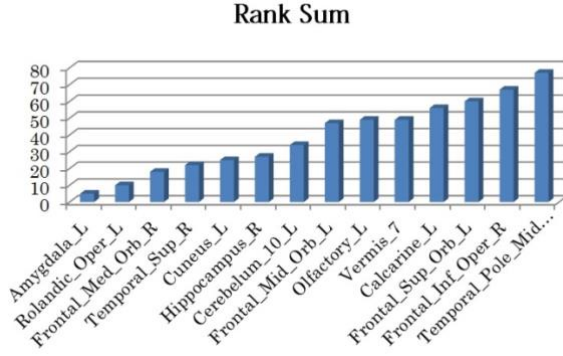


Figure 7. Bar chart of the rank sums in ascending order for each ROI recording the absolute values of the partial regression coefficients under each  $\alpha$  value.

## 4. Discussion

### 4.1. Interpretation of the Model

For the predicted (FC-BHQ) and correct (GM-BHQ) values obtained in this study, a moderate positive correlation ( $r > 0.6$ ) was achieved, although the slope of the regression line remained at approximately 0.4. This indicates that there is room for improvement in our model but suggests the possibility of applying GNNs to regression in the field of neuroimaging. Although the dataset used in this study consisted of a wide range of age groups and the GM-BHQ was unbiased, there was a tendency for the model to have significantly lower prediction accuracy at the base of the distribution; that is, the farther from the mean, the less accurate the model. It is possible that outliers prevented the model from learning useful graphical representations. Therefore, future studies should examine this possibility using outlier detection algorithms based on multivariate approaches, such as the k-nearest neighbor method (Su & Tsai, 2011).

Next, we investigate in detail the brain region clustering in the hidden layer of BrainGNN. As a result of soft clustering, the non-negative membership score matrix  $(\alpha_{iu}^{(1)})^+$  showed several areas in which the score at one cluster was significantly larger than those at the others; that is, where the belonging cluster was uniquely determined. However, there were also clear patterns of ambiguous cluster affiliation. It is worth noting here that Li et al. (2021) suggested the existence of ROIs that would not be confidently assigned to any cluster when the non-negative membership score matrix is

sparser than that obtained in this study. Although the prediction tasks and data were different in the two studies, it would be appropriate to say that they generally exhibited similar behavior.

Thus, the results of hard clustering were affected by subtle differences in membership scores by cluster. Simultaneously, each hard cluster was not necessarily composed of ROIs that were anatomically close. In addition, most clusters were not determined as a set of ROIs, as indicated by the default mode network (DMN) or other popular resting-state networks (RSNs), but rather as a complex set of ROIs simultaneously attributed to various RSNs. Therefore, this study does not focus on the interpretation of cluster attribution but interprets ROIs with large absolute partial regression coefficients by referring to the context in which they co-occur in prior aging science studies. In this section, we discuss Model 1 in particular, which was the most accurate in predicting GM-BHQ. We evaluated the absolute values of the partial regression coefficients in the sparse group Lasso ( $0 \leq \alpha < 1$ ) and Lasso ( $\alpha = 1$ ) and consistently found that the left amygdala (Amygdala\_L), left Rolandic operculum (Rolandic\_Oper\_L), and right superior temporal gyrus (Temporal\_Sup\_R) were the three brain regions that contributed most to the prediction of GM-BHQ (Table 2). Other brain regions with stable contributions were the left cerebellar lobule X (cerebellum\_10\_L), bilateral medial prefrontal and orbitofrontal cortices (Frontal\_Med\_Orb\_L and Frontal\_Med\_Orb\_R), and the right hippocampus (Hippocampus\_R).

Here, we first address the left amygdala, which consistently showed the largest contribution to prediction. The amygdala, located in the limbic system at the base of the forebrain, is part of the circuitry responsible for emotional processing, and many studies have shown that its function changes with age. For example, memories related to emotions tend to be retained better with age, whereas episodic memories deteriorate with age (Wright et al., 2006; Ritchey et al., 2011). This is called the “positive effect,” and various studies have suggested that age-related changes in functional connectivity (FC) with the amygdala might be a neurological biomarker for this effect (Sakaki et al., 2013; Addis et al., 2010). Aging and amygdala FC changes have not only been reported as progressing from younger to older adulthood, but also begin before adulthood, since FC between the amygdala and cortex is modulated from childhood through to the 20s

(McRae et al., 2012; Gee et al., 2013). This suggests that FC in the amygdala may be related to cognitive developmental processes. The amygdala is of interest in relation to the medial prefrontal cortex (Frontal\_Med\_Orb\_L), a brain region included in the DMN. The medial prefrontal cortex and amygdala have been suggested as important brain regions associated with the positive effects described above (Xiao et al., 2018). The FC between these brain regions belonging to the DMN and the amygdala have been reported to change with age, suggesting the involvement of the Frontal\_Med\_Orb\_L in the amygdala-centered aging-related network.

For other ROIs, the interpretation of the results can be facilitated by discussing them considering the functional networks to which they belong. Cognitive functions, including executive control, generally decline with age, even in the absence of a confirmed disease, and this modulation is characterized by a change in FC during the aging process (Ferreira & Busatto, 2013). DMN connectivity is significantly attenuated by senescence (Andrews-Hanna et al., 2007; Sambataro et al., 2010; Grady et al., 2010; Damoiseaux et al., 2008), and this phenomenon has been reported to be associated with reduced cognitive processing speed and performance in tasks related to working memory in elderly people. Such a disruption of the network configuration due to changes in connectivity is called *dedifferentiation*. In a review article on this topic, Koen and Rugg (2019) noted that although the mechanisms of reduced behavioral performance and various memory impairments may be explained by the loss of diversity in neural representation due to age-related dedifferentiation, this process does not necessarily imply a detrimental consequence of aging. When dedifferentiation is not accompanied by failure, it is specifically referred to as *degeneracy*, that is, an adaptive mechanism in which a dysfunctional network with some of its impaired nodes mobilizes a group of different nodes from another normally working network (Fornito et al., 2015).

After the amygdala, the next two areas of focus from the partial regression coefficients were the left Rolandic operculum (Rolandic\_Oper\_L) and right superior temporal gyrus (Temporal\_Sup\_R). Interestingly, in the context of age-induced changes in FC reported by Geerligs et al. (2014), Rolandic\_Oper\_L belongs to the cluster DAN-SMN (Damoiseaux et al., 2008), which consists of both the dorsal attention network (DAN) and somatomotor network (SMN). When seed regions were set in the



DMN, connectivity with the DAN-SMN containing Rolandic\_Oper\_L was enhanced in the elderly population. Furthermore, Geerligs et al. (2014) observed reduced connectivity between DAN-SMN seeds and Temporal\_Sup\_R in older participants. In addition, elderly people showed connectivity modulation within the DAN-SMN or with relationships to other networks. This change in connectivity can be considered as typical *dedifferentiation* associated with aging, and our model may have captured this phenomenon.

## 4.2. Limitations and Future Perspectives

In this study, we trained a deep learning model called BrainGNN to learn the relationship between rs-fMRI data and the brain health index GM-BHQ based on GMV and validated the output of the model by defining it as FC-BHQ. Thus, we created a model that expanded the functions of a GNN from discrimination to regression in the field of neuroimaging. Although the model had difficulty predicting the data at the base of the GM-BHQ distribution, it achieved a moderate positive correlation ( $r > 0.6$ ) between the predicted and correct values. Moreover, we confirmed that the parameter settings used by BrainGNN to learn the ROI-specific convolutional kernel did not cause a significant difference in performance. We also attempted sparse modeling to interpret the clustering generated at the hidden layers of BrainGNN. We found that the model emphasized changes in connectivity inside and outside of a functional network with age, and concluded that the FC-BHQ is a valid predictor of the GM-BHQ. Of particular interest was the extraction of network structures related to emotional processing in the cortex, particularly in the left amygdala. To further develop this method, a future challenge is to devise a mechanism to interpret the learning process that occurs inside the model and introduce it into BrainGNN so that it can be completed by deep learning modeling alone without resorting to additional analysis, especially sparse modeling.

The model constructed in this study achieved a moderate performance correlation ( $r > 0.6$ ) between the predicted and correct values, and the clustering performed inside the model extracted brain regions that have been reported to undergo significant changes with aging. However, there are some limitations and issues

regarding the analytical methods and data. First, it is unclear whether the model predicting GM-BHQ can provide meaningful information on the cognitive abilities of participants in the form of test score predictions at the individual level. Furthermore, this study only compared the prediction performance within BrainGNN, and the lack of similar research did not allow for adequate comparisons. These results alone do not fully demonstrate the usefulness of BrainGNN for regression tasks or the validity of the FC-BHQ. If prediction is considered at the individual level, the GNN framework is insufficient, and a hierarchical GNN, for example, needs to be constructed. This is a powerful framework (Bessadok et al., 2022) that uses a population graph constructed with brain graph representations built from brain image data as nodes and the pairwise similarity of phenotypic data (gender, age, genetic information, etc.) between subjects as edges. More concisely, a graph can be assumed in which each subject is a node and the similarities between subjects are constructed as edges. However, obtaining a wide range of brain images and phenotypic information simultaneously is difficult.

Regarding the algorithm used in this study, two issues must be addressed. First, we could not apply the group-level consistency loss function (GLC Loss), which has been implemented in previous studies, to the regression task in this study. GLC Loss constrains the nodes selected for pooling in the R-pool layer as close within a class. This loss function is easily applicable when the problem to be solved is classification, that is, the prediction of discrete values. However, the forecasting task in this study was the prediction of continuous values, so the data could not be directly categorized based on the predicted values; improving GLC Loss for the regression task could enhance the interpretability of the model and is expected to increase the prediction accuracy and make the results more robust, as it encourages common feature selection across similar data.

Another issue is the methodological limitation of the sparse modeling performed to analyze the clustering results of BrainGNN. The FCS used as features in sparse modeling is a further summary of the FC, which may miss some functional information. Modeling using different features is necessary to generalize the results. It should also be noted that the interpretation of results by the sparse group Lasso may be inherently different from the interpretation of nonlinear, complex, and highly abstract

data that BrainGNN performed in its model. The possible introduction of modules that make BrainGNN more interpretable alone would help to avoid these methodological limitations and problems. For example, attention maps (Huang et al., 2022), which visualize the brain function features on which the model has focused, could be used.

Finally, the limitations of the NKI-RS dataset used in this study should be mentioned. All fMRI data were collected at the same facility. In general, fMRI data are sensitive to imaging parameters; therefore, data collected at different imaging locations or with different imaging protocols may yield results different from those obtained in this study. Therefore, a new challenge is to introduce a mechanism that enables harmonization in the model. Another future challenge is to conduct experiments using an atlas other than AAL and to compare the results. The graph representation extracted by the GNN is also expected to change. Owing to the above limitations, future experiments using different datasets and atlases should be conducted for comparison and validation to generalize the results. It should also be noted that we were not able to confirm the association between GM-BHQ and health information in the NKI-RS dataset because we could not obtain information on the health status of the subjects in the dataset.

## 5. References

Andrews-Hanna, J. R., Snyder, A. Z., Vincent, J. L., Lustig, C., Head, D., Raichle, M. E., & Buckner, R. L. (2007). Disruption of large-scale brain systems in advanced aging. *Neuron*, 56(5), 924–935.

<https://doi.org/10.1016/j.neuron.2007.10.038>

Ashburner, J., & Friston, K. J. (2000). Voxel-based morphometry-the methods. *NeuroImage*, 11(6 Pt. 1), 805–821.

<https://doi.org/10.1006/nimg.2000.0582>

Bessadok, A., Mahjoub, M. A., & Rekik, I. (2022). Graph neural networks in network

neuroscience. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PP, Advance online publication.

<https://doi.org/10.1109/TPAMI.2022.3209686>

Cui, Z., & Gong, G. (2018). The effect of machine learning regression algorithms and sample size on individualized behavioral prediction with functional connectivity features. *NeuroImage*, 178, 622–637.

<https://doi.org/10.1016/j.neuroimage.2018.06.001>

Damoiseaux, J. S., Beckmann, C. F., Arigita, E. J., Barkhof, F., Scheltens, P., Stam, C. J., Smith, S. M., & Rombouts, S. A. (2008). Reduced resting-state brain activity in the "default network" in normal aging. *Cerebral Cortex* 18(8), 1856–1864.

<https://doi.org/10.1093/cercor/bhm207>

Erickson, K. I., Leckie, R. L., & Weinstein, A. M. (2014). Physical activity, fitness, and gray matter volume. *Neurobiology of Aging*, 35 Suppl. 2, S20–S28.

<https://doi.org/10.1016/j.neurobiolaging.2014.03.034>

Ferreira, L. K., & Busatto, G. F. (2013). Resting-state functional connectivity in normal brain aging. *Neuroscience and Biobehavioral Reviews*, 37(3), 384–400.

<https://doi.org/10.1016/j.neubiorev.2013.01.017>

Fornito, A., Zalesky, A., & Breakspear, M. (2015). The connectomics of brain disorders. *Nature reviews. Neuroscience*, 16(3), 159–172.

<https://doi.org/10.1038/nrn3901>

Gee, D. G., Humphreys, K. L., Flannery, J., Goff, B., Telzer, E. H., Shapiro, M., Hare, T. A., Bookheimer, S. Y., & Tottenham, N. (2013). A developmental shift from positive to negative connectivity in human amygdala-prefrontal circuitry. *The Journal of Neuroscience: the official journal of the Society for Neuroscience*, 33(10), 4584–4593.

<https://doi.org/10.1523/JNEUROSCI.3446-12.2013>

Geerligs, L., Maurits, N. M., Renken, R. J., & Lorist, M. M. (2014). Reduced specificity of functional connectivity in the aging brain during task performance. *Human Brain Mapping*, 35(1), 319–330.

<https://doi.org/10.1002/hbm.22175>

Grady, C. L., Protzner, A. B., Kovacevic, N., Strother, S. C., Afshin-Pour, B., Wojtowicz, M., Anderson, J. A., Churchill, N., & McIntosh, A. R. (2010). A multivariate analysis of age-related differences in default mode and task-positive networks across multiple cognitive domains. *Cerebral Cortex*, 20(6), 1432–1447.

<https://doi.org/10.1093/cercor/bhp207>

Holtmaat, A., Svoboda, K. (2009). Experience-dependent structural synaptic plasticity in the mammalian brain. *Nature Reviews Neuroscience* . 10, 647–658.

<https://doi.org/10.1038/nrn2699>

Huang, S. G., Xia, J., Xu, L., & Qiu, A. (2022). Spatio-temporal directed acyclic graph learning with attention mechanisms on brain functional time series and connectivity. *Medical Image Analysis*, 77, 102370.

<https://doi.org/10.1016/j.media.2022.102370>

Koen, J. D., & Rugg, M. D. (2019). Neural dedifferentiation in the aging brain. *Trends in*

Cognitive Sciences, 23(7), 547–559.

<https://doi.org/10.1016/j.tics.2019.04.012>

Li, X., Zhou, Y., Dvornek, N., Zhang, M., Gao, S., Zhuang, J., Scheinost, D., Staib, L. H., Ventola, P., & Duncan, J. S. (2021). BrainGNN: Interpretable brain graph neural network for fMRI analysis. *Medical Image Analysis*, 74, 102233.

<https://doi.org/10.1016/j.media.2021.102233>

McRae, K., Gross, J. J., Weber, J., Robertson, E. R., Sokol-Hessner, P., Ray, R. D., Gabrieli, J. D., & Ochsner, K. N. (2012). The development of emotion regulation: an fMRI study of cognitive reappraisal in children, adolescents and young adults. *Social Cognitive and Affective Neuroscience*, 7(1), 11–22.

<https://doi.org/10.1093/scan/nsr093>

Nemoto, K., Oka, H., Fukuda, H., & Yamakawa, Y. (2017). MRI-based brain healthcare quotients: A bridge between neural and behavioral analyses for keeping the brain healthy. *PloS One*, 12(10), e0187137.

<https://doi.org/10.1371/journal.pone.0187137>

Ritchey, M., Bessette-Symons, B., Hayes, S. M., & Cabeza, R. (2011). Emotion processing in the aging brain is modulated by semantic elaboration. *Neuropsychologia*, 49(4), 640–650.

<https://doi.org/10.1016/j.neuropsychologia.2010.09.009>

Sakaki, M., Nga, L., & Mather, M. (2013). Amygdala functional connectivity with medial prefrontal cortex at rest predicts the positivity effect in older adults' memory. *Journal of Cognitive Neuroscience*, 25(8), 1206–1224.

[https://doi.org/10.1162/jocn\\_a\\_00392](https://doi.org/10.1162/jocn_a_00392)

Sambataro, F., Murty, V. P., Callicott, J. H., Tan, H. Y., Das, S., Weinberger, D. R., & Mattay, V. S. (2010). Age-related alterations in default mode network: impact on working memory performance. *Neurobiology of Aging*, 31(5), 839–852.

<https://doi.org/10.1016/j.neurobiolaging.2008.05.022>

Simon, N., Friedman, J., Hastie, T., & Tibshirani, R. (2013). A sparse-group Lasso. *Journal of Computational and Graphical Statistics*, 22(2), 231–245.

<https://doi.org/10.1080/10618600.2012.681250>

Su, X., & Tsai, C. (2011). Outlier detection. *WIREs Data Mining and Knowledge Discovery*, 1(3), 261–268.

<https://doi.org/10.1002/widm.19>

Wright, C. I., Wedig, M. M., Williams, D., Rauch, S. L., & Albert, M. S. (2006). Novel fearful faces activate the amygdala in healthy young and elderly adults. *Neurobiology of Aging*, 27(2), 361–374.

<https://doi.org/10.1016/j.neurobiolaging.2005.01.014>

Xiao, T., Zhang, S., Lee, L. E., Chao, H. H., van Dyck, C., & Li, C. R. (2018). Exploring age-related changes in resting state functional connectivity of the amygdala: From young to middle adulthood. *Frontiers in Aging Neuroscience*, 10, 209.

<https://doi.org/10.3389/fnagi.2018.00209>

Note:

This paper is an abstract of the master's thesis of the first author, who is an alumnus of the Tokyo Institute of Technology, Japan.

COI:

There are no conflicts of interest.