

TWIN CLASSIFICATION IN RESTING-STATE BRAIN CONNECTIVITY

Andrey Gritsenko[†] Martin Lindquist[§] Moo K. Chung[‡]

[†] Northeastern University, Department of Electrical and Computer Engineering, Boston, USA

[§] Johns Hopkins University, Department of Biostatistics, Baltimore, USA

[‡] University of Wisconsin, Department of Biostatistics and Medical Informatics, Madison, USA

ABSTRACT

Twin study is one of the major parts of human brain research that reveals the importance of environmental and genetic influences on different aspects of brain behavior and disorders. Accurate characterization of identical and fraternal twins allows us to infer on the genetic influence in a population. In this paper, we propose a novel pair-wise classification pipeline to identify the zygosity of twin pairs using the resting state functional magnetic resonance images (rs-fMRI). The new feature representation is utilized to efficiently construct brain network for each subject. Specifically, we project the fMRI signal to a set of cosine series basis and use the projection coefficients as the compact and discriminative feature representation of noisy fMRI. The pair-wise relation is encoded by a set of twin-wise correlations between functional brain networks across brain regions. We further employ hill climbing variable selection to identify the most genetically affected brain regions. The proposed framework has been applied to 208 twin pairs in Human Connectome Project (HCP) and we achieved 92.23(±4.43)% classification accuracy.

Index Terms— Resting-state fMRI, brain connectivity, twin study, zygosity, neural networks, hill climbing variable selection

1. INTRODUCTION

The development and functioning of human brain are influenced by both genetic and environmental factors. Though, the extent by which these factors shape the brain structure and function is still unknown and has been a research interest for decades. Twins provide a valuable source of information, as their unique relationship allows researchers to pull apart and examine genetic and environmental influences *in-vivo*. The power of twin study arises from the fact that there are only two types of twins, identical (or *monozygotic*, MZ) and fraternal (or *dizygotic*, DZ), that share different amount of genetic information: MZ twins share 100% of genes, and DZ twins on average share only 50% of genes [1]. Comparing the similarity of MZ and DZ twins, we can draw inference on the genetic influence in a population. In previous studies [2, 3], it has been shown that task-related brain activity is strongly affected by genetics, and the purpose of this paper is to demonstrate that genes also shape activity of human brain in resting state. Additionally, we deduce what brain regions are the most genetically affected. To achieve this goal, we compare the similarity of MZ and DZ twins using rs-fMRI as it allows to investigate the baseline functional connectivity of the brain and has a broader use for clinical applications, in contrast to task-fMRI [4].

This study was supported by NIH Grant EB022856 and NCATS Grant UL1TR002373. Correspondence should be sent to Andrey Gritsenko (email: agritsenko@ece.neu.edu). Andrey Gritsenko performed the work while at University of Wisconsin.

In this paper, we propose a unified classification pipeline that automatically determines the zygosity of a paired twin fMRI (Figure 1). The proposed pipeline utilizes 2-layer neural network to achieve very high accuracy (> 90%) on zygosity classification task. Shallow structure of the network also allows to preserve interpretability of the results, compared to a more sophisticated deep neural networks [5]. The interpretability of results is achieved with a novel feature representation that converts twin fMRI into a single algebraic representation defined on the automated anatomical labeling (AAL) human brain atlas with 116 predefined regions of interest (ROI) [6]. First, we propose a simultaneous temporal smoothing and dimension reduction approach to construct a new compact feature representation for the fMRI signal at each voxel using the cosine series representation (CSR). To properly correlate twin fMRI we define functional brain network for each subject. The connectivity between voxels is then represented using the lower-dimensional data representation involving CSR. Next, we gain vector representation of twin data computing the correlation between networks. Using this vector representation in a two-layer artificial neural network (ANN), we achieve high classification accuracy. Finally, we use hill climbing, a variable selection procedure, to determine the significance of information provided by AAL regions and discover the most informative regions. Because these regions contribute to the discrimination of MZ and DZ pairs, we imply that the most significant ROI are the most genetically affected.

The main contributions of the paper are: (1) we introduce new feature representation to characterize voxel-wise functional connectivity in the whole brain that substantially improves the classification performance; (2) we present a solution to identify the zygosity of twin pairs using only rs-fMRI data, achieving 92.23% classification accuracy; (3) we propose a principled way to examine the genetic influence in each brain region. The proposed framework has proved its robustness in one of the largest publicly available twin datasets – HCP that contains resting state functional magnetic resonance imaging of 131 MZ and 77 DZ twin pairs.

2. TWIN RS-FMRI IN HUMAN CONNECTOME PROJECT (HCP)

We have used HCP [7] dataset with rs-fMRI scans of 208 same-sex twin pairs. The scans undergone spatial and temporal preprocessing including: spatial distortions and motion correction, frequency filtering, artifact removal, etc. as part of HCP preprocessing pipeline [8]. The information regarding the zygosity of twin pairs is given: there are 131 MZ twin pairs (age 29.3 ± 3.3 , 56M/75F) and 77 DZ twin pairs (age 29.1 ± 3.5 , 30M/47F). fMRI were collected with 3T scanner at 2.0 mm voxel resolution ($91 \times 109 \times 91$ spatial dimensionality), 1200 frames at a rate of one frame every 720 ms.

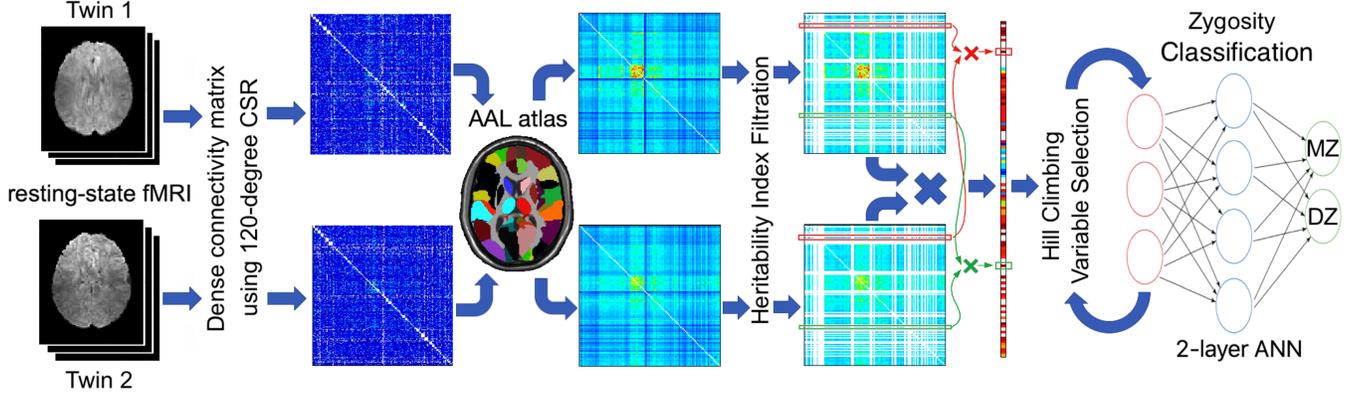


Fig. 1. Outline of the proposed framework. We construct brain networks using temporally smoothed data with 120-degree CSR. We use AAL atlas to break down voxel-level whole brain connectivity matrices into functional networks with 116 nodes. Region-level correlations between twin brain networks result in vector representations used in ANN.

The concept of *heritability* comes from genetics where it was used for decades to estimate the influence of genes on the phenotype (or, variance) of different traits. The population-level variance of any trait is due to genetic effect A and environmental effect C . While both MZ and DZ twins share the common environment, MZ twins share 100% of genes and DZ twins share on average only 50%. At each network node, the twin-wise correlation within MZ and DZ twin pairs is given by $\rho_{MZ} = A + C$ and $\rho_{DZ} = A/2 + C$. The influence of genetic effects is then quantified as *heritability index* (HI), which is estimated as

$$\mathcal{H} = 2(\rho_{MZ} - \rho_{DZ}), \quad (1)$$

according to Falconer's formula [9].

3. METHODS

3.1. Cosine Series Representation

Highly noisy fMRI data is often subjected to spatial smoothing using Gaussian kernel or Fourier transform [10]. Here, we propose to smooth fMRI temporally without spatial smoothing using CSR to increase the localization power [11, 12].

Given functional time series $\zeta(t) = \mu(t) + \epsilon(t)$ of the variable $t \in [0, 1]$, where $\epsilon(t)$ is a zero-mean noise, we want to estimate the unknown signal $\mu(t)$. We assume $\zeta, \mu \in \mathcal{L}^2[0, 1]$, the space of square integrable functions in $[0, 1]$. We can always transform data such that they are defined in $[0, 1]$. Consider the eigenfunction problem $\Delta\psi + \lambda\psi = 0$ with the periodic $\psi(t) = \psi(t + 2m)$, $m \in \mathbb{Z}$ and evenness $\psi(-t) = \psi(t)$ constraints. We can show that eigenfunctions ψ_j form an orthonormal basis in $\mathcal{L}^2[0, 1]$ with respect to the inner product $\langle \psi_i, \psi_j \rangle = \int_0^1 \psi_i(t)\psi_j(t)dt = \delta_{ij}$, where δ_{ij} is the Kronecker's delta. The unique solution satisfying above-mentioned eigenfunction problem is the basis $\psi_0(t) = 1$, $\psi_l(t) = \sqrt{2} \cos(l\pi t)$, $l = 1, 2, \dots$

We estimate $\mu(t)$ in a subspace \mathcal{M}_k spanned by up to the k -th degree basis functions. The least squares estimation (LSE) of μ in \mathcal{M}_k is given by $\hat{\mu}(t) = \sum_{l=0}^k \hat{c}_l \psi_l(t)$, where $\hat{c}_l = \arg \min_{c_l \in \mathbb{R}} \left| \sum_{l=0}^k c_l \psi_l(t) - \zeta(t) \right|^2 = \langle \zeta, \psi_l \rangle$.

In fMRI, we deal with a collection of functional time series $\zeta_i(t)$, $i = 1, \dots, n$ recorded from n voxels at time points t_j , $j = 1, \dots, p$. CSR of ζ_i can be represented as matrix multiplication $Y = \Psi \cdot C$, where $Y_{p \times n} = (\zeta_i(t_j))$, $\Psi_{p \times k} = (\psi_l(t_j))$, and $C_{k \times n} = (c_{li})$ is the matrix of CSR coefficients to be estimated.

The LSE of C is found as $\hat{C} = \Psi^\dagger Y = (\Psi^T \Psi)^{-1} \Psi^T Y$, where Ψ^\dagger is the Moore-Penrose inverse of matrix Ψ . This matrix formulation speeds up the computation dramatically for large-scale computation since the inverse Ψ^\dagger is computed only once and the size of Ψ is much smaller than the number of voxels in the whole brain fMRI.

3.2. Computing Correlations using CSR

Consider CSR of fMRI obtained at two voxels $\zeta(t) = \sum_{i=0}^k \zeta_i \psi_i(t)$ and $\eta(t) = \sum_{j=0}^k \eta_j \psi_j(t)$. fMRI signals are often translated to the mean to compensate for varying baselines across subjects. Thus, we can assume $\mathbb{E}\zeta(t) = \int_0^1 \zeta(t) dt = \mathbb{E}\eta(t) = \int_0^1 \eta(t) dt = 0$. Let $\sigma^2 \zeta(t) = \int_0^1 \zeta^2(t) dt$, $\sigma^2 \eta(t) = \int_0^1 \eta^2(t) dt$ be the variances of $\zeta(t)$ and $\eta(t)$. The correlation $\hat{\rho}(\zeta, \eta)$ between ζ and η represents functional connectivity between fMRI at two corresponding voxels, and is given by CSR coefficients

$$\hat{\rho}(\zeta, \eta) = \frac{\int_0^1 \zeta(t)\eta(t)dt}{\sigma\zeta(t) \cdot \sigma\eta(t)} = \frac{\sum_{l=0}^k \zeta_l \eta_l}{\left[\sum_{i=0}^k \zeta_i^2 \sum_{j=0}^k \eta_j^2 \right]^{\frac{1}{2}}} = \zeta^T \eta, \quad (2)$$

where ζ and η are scaled CSR coefficients. The degree $k = 120$ is selected empirically.

3.3. Infer Zygosity via Brain Network Correlation

Given paired twin images, we use ANN to identify if they are MZ or DZ twins. At first, we obtain functional brain connectivity for each subject: eq. (2) is used to compute voxel-wise correlations in subject's brain across all voxels. We employ AAL parcellation template to construct a functional brain network $A_{N \times N} = (a_{ij})$ for each subject, where $N = 116$ is the number of ROI in AAL atlas, and a_{ij} represents functional connectivity between i -th and j -th brain regions computed as the average of all voxel-wise correlations between corresponding parcels. Since we cannot average correlations by taking the arithmetic mean without biasing, we transform correlations using the Fisher z -transform: $z = F(\rho) = \frac{1}{2} \ln \frac{1+\rho}{1-\rho}$. Obtained correlations are averaged and back projected via the inverse Fisher transform [13]. Finally, we obtain a vector representation $\vartheta_{N \times 1} = (\text{corr}(\mathbf{a}_i^1, \mathbf{a}_i^2))$ of twin data by correlating functional networks of twins: each element represents correlation between i -th ROI of twins brain networks with respect to their connectivity with other regions, where $\mathbf{a}_i^1, \mathbf{a}_i^2$ are i -th rows of connectivity matrices A^1, A^2 of the first and the second twin subjects, respectively.

A two-layer feed-forward artificial neural network with 50 hidden and 1 output neuron is used to classify zygosity of twins. Because the input of the network is vectors of region-wise correlations between twins brain networks, using sigmoid function as activation of hidden neurons is sufficient to capture the whole range of input values. We employ classification accuracy, false-positive rate (FPR) and false-negative rate (FNR) to evaluate the performance of the network (MZ is considered as ‘positive’ and DZ as ‘negative’ class). We use the holdout method to split dataset randomly into training, validation and test subsets in 70:15:15 proportion [14]. The averaged results of 1000 randomly initialized ANN models are presented.

Results. Using 116 feature vector representation of twin data, obtained as region-wise correlation between twin networks after temporal smoothing with 120-degree CSR at the voxel level, we achieved 76.94(± 8.01)% classification accuracy, with false-positive rate of 38.23(± 15.37)% and 13.76(± 7.75)% false-negative rate (Table 1, Base Model).

3.4. Boosting with Heritability and Variable Selection

The core of the proposed framework is the conversion from a pair of twin fMRI to a single algebraic expression representing the similarity between twins across different brain regions. So far, we assumed that all regions contribute equally in the classification. However, it is more reasonable to assume that connectivity between some regions are likely to be more heritable so the corresponding regions contribute more to the classification accuracy.

Heritability. Our goal is to identify a collection of the most genetically affected regions as the ones with the greatest contribution to the classification accuracy based on functional activity in twin pairs brain networks. We determine the genetic significance by incorporating the node level computation of the heritability index \mathcal{H} . Correlating twins functional networks across all nodes representing AAL parcels yields insufficiently high classification accuracy. We expect that not all regions in AAL parcellation contribute equally to classification accuracy. To boost overall accuracy, we filter out ROI with $\mathcal{H} < 0.05$ and use correlation of brain connectivity computed only across significant regions ($\mathcal{H} \geq 0.05$). Here, \mathcal{H} is computed according to eq. (1) as the doubled difference between average correlation vectors of MZ and DZ populations, respectively, and the threshold $\mathcal{H} = 0.05$ is selected as the 95% lower confidence bound.

Results. We achieved 82.68(± 8.68)% classification accuracy, with 28.27(± 17.03)% FPR and 10.75(± 7.48)% FNR classifying zygosity in variable space constructed of only significant regions, (Table 1, HI Filtration). We performed the HI filtration *only* on the training subset, and used the identified significant regions in validation and test subsets.

Hill Climbing. Because we leave 30% of population out when computing HI, it would be too optimistic to weight ROI contribution solely based on its \mathcal{H} value. We use the hill climbing variable selection to identify the most genetically affected regions by considering each AAL parcel as a variable [15]. We start hill climbing with an empty variable space (*selected* variables) and a pool of candidate variables. At each step, we test a candidate variable by adding it to the variable space of the model and estimating the classification accuracy. In other words, we use truncated vectors of region-wise correlations as input to NN, keeping only a part of it corresponding to already selected and currently testing ROI. When all candidates are tested, a variable with the highest accuracy is removed from the pool of candidates and added to the selected variables. The process iteratively continues until there are no candidate variables.

Results. Figure 2 shows the result of the hill climbing at each

Table 1. Performance of the proposed classification pipeline at different settings

Method	Accuracy (%)	FPR (%)	FNR (%)
Model without CSR	61.23(± 9.56)	55.72(± 17.23)	22.13(± 9.85)
Base model	76.94(± 8.01)	38.23(± 15.37)	13.76(± 7.75)
HI Filtration	82.68(± 8.68)	28.27(± 17.03)	10.75(± 7.48)
HI Filtration + HC	92.23(± 4.43)	16.86(± 7.81)	5.29(± 2.94)

iteration. White color represents the presence of a region in the variable space at a given iteration. Variable space, consisted of 15 variables (bold fonts), represents the most genetically affected AAL parcels (Figure 3). It also presents classification accuracy when a corresponding AAL region is added to the variable space. We achieve 92.23(± 4.43)% accuracy, with 16.86(± 7.81)% FPR and 5.29(± 2.94)% FNR (Table 1, HI Filtration+Hill Climbing).

4. DISCUSSION

In this paper, we present a framework to solve the problem of classifying the zygosity of a pair of twin resting state fMRI. This is a more complex problem than the usual classification problem of labeling each image into distinct classes. Here, we are interested in learning if the relationship between pairs of images is associated with the zygosity of twins. This was done within the proposed pipeline through sequential preprocessing and classification of the data. Using HI filtration followed by hill climbing variable selection, we found that 15 AAL parcels, namely Angular_L, Hippocampus_L, Frontal_Sup_R, Heschl_L, Cerebelum_Crus1_L, Vermis_9, Cerebelum_9_R, Vermis_1_2, Frontal_Mid_Orb_R, Frontal_Inf_Orb_R, Amygdala_R, Thalamus_L, Olfactory_L, Cerebelum_10_R and Frontal_Sup_Orb_R, are the most genetically affected (Figure 3). These regions contributed to 92.23(± 4.43)% classification accuracy in determining the zygosity of twins.

There were numerous advantages for using CSR. To determine the effect of temporal smoothing, we classified twin fMRI without incorporating CSR in the pipeline and obtained 61.23(± 9.56)% accuracy (Table 1, first row). Denoising fMRI with CSR increased the classification accuracy to 76.94(± 8.01)%, the boost of 15%. Furthermore, CSR with $k = 120$ drastically reduced the computational time by 90%, since the correlation of any two fMRI was simplified as the product of CSR coefficient vectors.

For variable selection, we applied the hill climbing algorithm – a greedy search algorithm that considers AAL regions as independent variables and tests one variable at a time. Due to its greedy nature, hill climbing algorithm cannot guarantee convergence to the global optimum, so we performed it for each of 1000 ANN models to achieve reliable local optimum approximation. To get a more accurate inference on the optimal variable space, one may consider taking advantage of the hypernetwork structure of human brain [16] and performing a group variable selection, i.e., combine variables in groups and test each group as a single instance [15]. It has to be highlighted that HI filtration can be considered as a preliminary step for variable selection. Though it is possible to omit this step in the pipeline, its purpose is to utilize domain knowledge, thus helping to exclude brain regions with knowingly low heritability index from further consideration.

Our proposed twin classification pipeline is very general and applicable to other type of paired image settings such as longitudinal studies with two scans or multimodal studies with two different imaging modalities if appropriate handling of functional brain networks in this setting is provided.

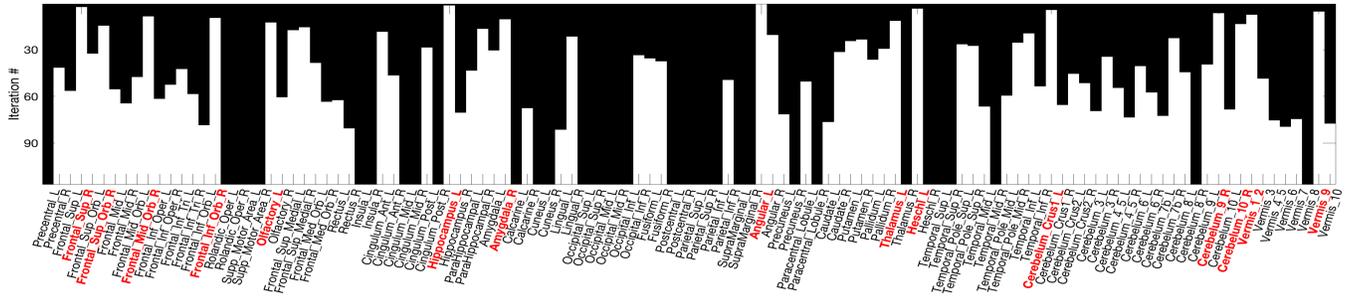


Fig. 2. Results of heritability index filtration with hill climbing. Height of white bars represent iteration, at which region is introduced to the variable space (the higher, the earlier). Parcels in bold font (also red in color version) constitute the optimal variable space of 15 variables with the highest classification accuracy, $92.23(\pm 4.43)$.

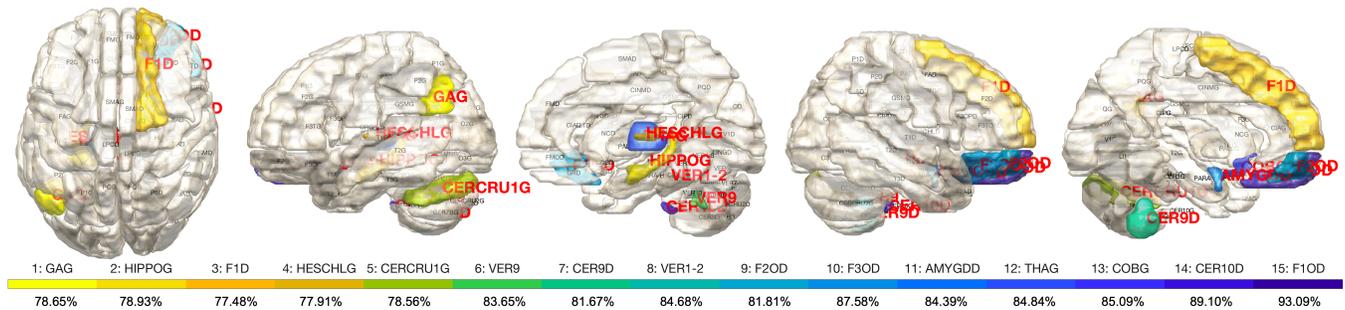


Fig. 3. AAL regions of the optimal variable space discovered by hill climbing. Colors represent iterations when ROI are added to the variable space. Third and fifth projections illustrate right and left hemispheres from the center-of-brain, respectively.

Acknowledgements. We would like to thank Guorong Wu of University of North Carolina for valuable discussions and feedback, and Gregory Kirk of University of Wisconsin for logistical support.

5. REFERENCES

- [1] M.C. Neale and L.R. Cardon, *Methodology for Genetic Studies of Twins and Families*, vol. 67, Springer Science & Business Media, 2013.
- [2] G.I. de Zubicaray et al., “Meeting the Challenges of Neuroimaging Genetics,” *Brain Imaging and Behavior*, vol. 2, no. 4, pp. 258, 2008.
- [3] A.G. Jansen et al., “What Twin Studies Tell Us About the Heritability of Brain Development, Morphology, and Function: a Review,” *Neuropsychology Review*, vol. 25, no. 1, pp. 27–46, 2015.
- [4] R.M.S. Panchuelo, M.C. Stephenson, S.T. Francis, and P.G. Morris, “Neural Brain Activation Imaging,” in *Biomedical Imaging*, Peter Morris, Ed., pp. 112 – 162. Woodhead Publishing, 2014.
- [5] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-Based Learning Applied to Document Recognition,” *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [6] N. Tzourio-Mazoyer et al., “Automated Anatomical Labeling of Activations in SPM Using a Macroscopic Anatomical Parcellation of the MNI MRI Single-Subject Brain,” *NeuroImage*, vol. 15, no. 1, pp. 273 – 289, 2002.
- [7] D.C. Van Essen et al., “The Human Connectome Project: A Data Acquisition Perspective,” *NeuroImage*, vol. 62, no. 4, pp. 2222–2231, 2012.
- [8] S.M. Smith et al., “Resting-state fMRI in the Human Connectome Project,” *NeuroImage*, vol. 80, pp. 144–168, 2013.
- [9] D.S. Falconer, T.F.C. Mackay, and R. Frankham, *Introduction to Quantitative Genetics (4th edn)*, Pearson Education, 1996.
- [10] R.A. Poldrack, J.A. Mumford, and T.E. Nichols, *Handbook of Functional MRI Data Analysis*, Cambridge University Press, 2011.
- [11] M.K. Chung et al., “Cosine Series Representation of 3D Curves and Its Application to White Matter Fiber Bundles in Diffusion Tensor Imaging,” *Stat Interface*, vol. 3, no. 1, pp. 69–80, 2010.
- [12] S. Huang, A. Gritsenko, M.A. Lindquist, and M.K. Chung, “Circular Pearson Correlation Using Cosine Series Expansion,” in *2019 IEEE 16th International Symposium on Biomedical Imaging (ISBI 2019)*. IEEE, 2019, pp. 1774–1777.
- [13] J. Vrbik, “Population Moments of Sampling Distributions,” *Computational Statistics*, vol. 20, no. 4, pp. 611–621, 2005.
- [14] M.T. Hagan, H.B. Demuth, M.H. Beale, and O. De Jess, *Neural Network Design*, Martin Hagan, USA, 2nd edition, 2014.
- [15] S.J. Russell and P. Norvig, *Artificial Intelligence: A Modern Approach*, Pearson Education, 2003.
- [16] B. Jie et al., “Brain Connectivity Hyper-Network for MCI Classification,” in *Medical Image Computing and Computer-Assisted Intervention*, 2014, pp. 724–732.