# NORMATIVE MODELING VIA CONDITIONAL VARIATIONAL AUTOENCODER AND ADVERSARIAL LEARNING TO IDENTIFY BRAIN DYSFUNCTION IN ALZHEIMER'S DISEASE

*Xuetong Wang[1], Kanhao Zhao[2], Rong Zhou[1], Alex Leow[3], Ricardo Osorio[4], Yu Zhang[2,5], Lifang He[1]*

[1] Department of Computer Science and Engineering, Lehigh University, PA, USA
[2] Department of Bioengineering, Lehigh University, PA, USA
[3] Department of Psychiatry, University of Illinois at Chicago, IL, USA
[4] Department of Psychiatry, NYU Grossman School of Medicine, NY, USA
[5] Department of Electrical and Computer Engineering, Lehigh University, PA, USA

## ABSTRACT

Normative modeling is an emerging and promising approach to effectively study disorder heterogeneity in individual participants. In this study, we propose a novel normative modeling method by combining conditional variational autoencoder with adversarial learning (ACVAE) to identify brain dysfunction in Alzheimer's Disease (AD). Specifically, we first train a conditional VAE on the healthy control (HC) group to create a normative model conditioned on covariates like age, gender and intracranial volume. Then we incorporate an adversarial training process to construct a discriminative feature space that can better generalize to unseen data. Finally, we compute deviations from the normal criterion at the patient level to determine which brain regions were associated with AD. Our experiments on OASIS-3 database show that the deviation maps generated by our model exhibit higher sensitivity to AD compared to other deep normative models, and are able to better identify differences between the AD and HC groups.

***Index Terms—*** Normative modeling, conditional variational autoencoder, adversarial learning, Alzheimer's disease
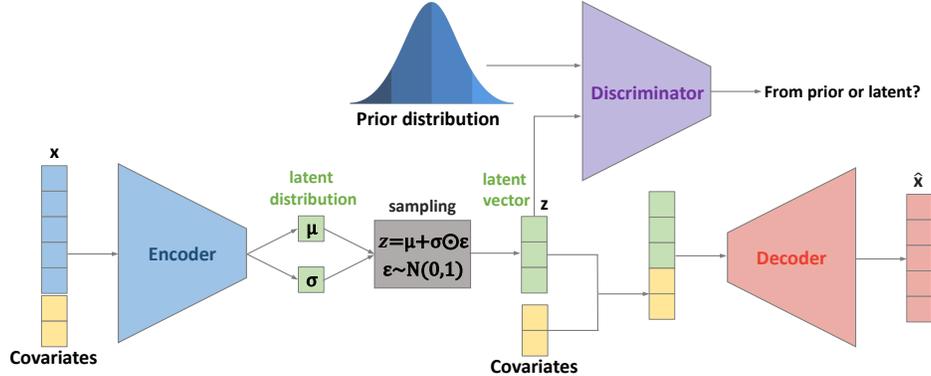
## 1. INTRODUCTION

Brain diseases such as Alzheimer's disease (AD) are usually highly heterogeneous, which exacerbates the difficulty of clinical treatment. The traditional case-control approaches assume a consistent pattern of abnormalities among individuals belonging to the same cohort, but ignore the heterogeneity of the disorder [1, 2]. In contrast to case-control studies, normative modeling can quantify how individual patients deviate from the expected normative range by learning in the healthy control (HC) group and thus explicitly modeling disease heterogeneity, which provides information about potential abnormalities in each particular individual. The normalized model is a two-step process in which the model is first trained on the HC cohort, and then the trained model is applied to a target cohort to quantify deviations [3].

Recently, deep learning techniques are very popular for normative modeling, especially autoencoder (AE) based methods [4, 5]. AE consists of two components: encoder and decoder. The encoder compresses the data from a high-dimensional space to a low-dimensional space also called latent code, and then the decoder converts the data from the latent code to a high-dimensional space like the input data. However, since AE mainly emphasizes the image reproduction function, this also leads to a drawback that is the lack of randomness, which results in a model that is not sufficiently generalized. When we randomly change the latent code, the output may not be related to the original data at all. In other words, the latent code for AE is non-regularized.

To solve the above issue, variational autoencoder (VAE) [6] is adopted in the normative modeling. VAE is a generative model that differs from the AE in that the encoder of VAE outputs the parameters of a pre-defined distribution rather than just a latent vector. It then forces this distribution to be a normal distribution, which ensures that the latent space is regularized. Moreover, to eliminate the influence of some confounding variables (e.g., age, gender, and intracranial volume) on the model in the learning process, conditional variational autoencoder (CVAE) has also been applied to the normative modeling [6]. Unlike the vanilla VAE, the encoder of CVAE can generate latent distribution parameters defined in advance based on the input data and confounding variables. Similarly, the decoder can reconstruct the original input based on the confounding variables and the vectors sampled in the latent space. On the other hand, there are some recent studies applying adversarial learning to autoencoders (AAE) in normative modeling [7, 8], but the learning does not sufficiently take into account the randomness in the latent space which may lead to unreliable result in the analysis.

In this paper, we propose a novel normative model named as ACVAE by integrating CVAE and adversarial learning that improves the generality of the model by adding randomness to

**Fig. 1**. The framework of the adversarial conditional variational autoencoder.

the latent space and effectively reduces the impact on model learning by passing confounding variables in the encoder and decoder. First, we used regional brain volume data from the HC group to train our model. Next, we tested the trained model in the target population, generating deviation maps for AD and HC groups. As a result, each patient's deviation from the norm was estimated. By comparison with other models, our model exhibited better separation in terms of deviation between HC and AD subjects. By observing different deviation plots for each patient, heterogeneity can be better understood, which can provide a more reliable reference basis for clinical diagnosis and treatment.

## 2. MATERIALS AND METHODS

### 2.1. fMRI Acquisition and Preprocessing

The fMRI data used in this study was obtained from the OASIS-3 database [9]. It included a total of 1098 subjects with 605 cognitively normal adults and 493 individuals at various stages of cognitive decline ranging in age from 42 to 95 years. For each session, the fMRI data were scanned during resting state for 6 min (164 volumes) using 16-channel head coil of scanners with parameters: TR=2.2 s, TE=27 ms, FOV=240×240 mm, and FA=90°. The acquired rs-fMRI data were preprocessed using the reproducible fM-RIPrep pipeline [10]. The T1-weighted image was corrected for intensity nonuniformity and then stripped skull as T1w-reference. Spatial normalization was done through nonlinear registration, with the T1w reference [11], followed by FSL-based segmentation. The BOLD reference was then transformed to the T1w reference with a boundary-based registration method, configured with nine degrees of freedom to account for distortion remaining in the BOLD reference [12]. BOLD signals were slice-time corrected and resampled onto the participant's original space with head-motion parameters [13], susceptibility distortion correction, and then resampled into standard space, generating a preprocessed BOLD run in MNI standard space. ICA-AROMA [14] was then performed

**Table 1**. Data Characteristics.

|  | HC | AD |
|---|---|---|
| Num | 1476 | 21 |
| Gender (M/F) | 647/829 | 11/10 |
| Age (mean± std) | 68.5±6.4 | 77.4±3.4 |
| ICV (mean± std) | 60960.2±10249.9 | 61172.8±10493.2 |

for automatic removal of motion artifacts.

### 2.2. Feature Generation

Table 1 shows the characteristics of healthy control (HC) and Alzheimer's disease (AD) subjects used in this study. The OASIS-3 database spans a considerable amount of collection time, and some subjects were collected from multiple time periods. Thus, we treated the data of every 100-day period as a subject and assumed no significant deviation change in the same subject during the 100-day interval. This will result in 1497 subjects with 1476 HC subjects and 21 AD patients. For each subject, the voxel-level BOLD time series were first averaged into 100 regions-of-interest (ROIs) at each time point based on the Schaefer parcellation [15], and then averaged across time points to generate the ROI features as input. In addition, we consider the age, gender, and intracranial volume (ICV) as potentially confounding covariates, which were included in our model as conditional variables. The ICV for each subject was calculated by summing together all 100 ROI values. Both age and ICV were divided into 10 quantile-based bins and featured as one-hot encoded vectors. Specifically, the covariates from each subject were represented as a 22-dimensional vector with two dimensions used to one-hot encode the gender of male and female attributes.

### 2.3. Normative Modeling

**Overview.** Fig. 1 illustrates our end-to-end architecture of ACVAE for normative modeling consisting of two components: conditional variational autoencoder (bottom) and adversarial learning (top). The former generates useful low-dimensional latent representations of ROI features in HC

group and is conditioned by confounding covariates, while the latter employs adversarial training to shape the latent code distribution so that it resembles the predetermined prior distribution. Next, we give the details of each component.

**Conditional Variational Autoencoder (CVAE).** CVAE is a variant of variational autoencoder (VAE) to conditional tasks, which allows to learn a low-dimensional latent space from the input data with multivariable control in an unsupervised manner. To eliminate the effect of confounding covariates present in the data on the latent space in the neural network, we used a CVAE model for normative modeling. Similar to VAE, the CVAE network has three main components: the encoder, the latent distribution, and the decoder. However, the encoder and decoder of the CVAE receive additional conditional variables, i.e., age, gender, and ICV. First, the encoder receives the conditional variables and input data, and then generates pre-defined distribution parameters, i.e., mean and variance. The decoder receives both the samples sampled from the latent distribution and the conditional variables to output the reconstructed input. Furthermore, the loss function can be formulated as:

$$L_{\text{CVAE}} = E[logP(X|z,c)] - D_{\text{KL}}[Q(z|X,c)||P(z|c)], \quad (1)$$

where $X$ is the input data, $c$ is the conditional variables, and $z$ is the latent code. In addtion, $Q(z|X,c)$, $P(X|z,c)$, $P(z|c)$ are represented as encoder, decoder, prior respectively. The first term of the above loss function is the reconstruction mean squared error that measures the difference between the input data and the reconstructed output. The second term refers to the Kullback-Leibler (KL) divergence, which measures how far the pre-defined distribution ($Q(z|X,c)$) is from the true distribution ($P(z|c)$).

**Adversarial Learning.** Previous study has shown that combining VAE with adversarial learning allows to complement the VAE reconstruction loss with the perceptual-level representation of the discriminator [16]. Thus we combine the benefit of adversarial learning with CVAE. Adversarial learning consists of two parts: the discriminator and the generator. The discriminator accepts two inputs, one is a random sample sampled in the prior distribution and the other is sampled from the latent distribution. The discriminator will try to identify whether the input is a sample sampled from the prior distribution or sampled from the latent distribution. Yet the generator wants to generate samples that cheat the discriminator. The objective function can be expressed as follows:

$$L_{\text{Adv}} = E[log(D(z))] + E[log(1 - D(P(X|z,c)))], \quad (2)$$

where $D(z)$ is the discriminator and $P(X|z,c))$ is the generator, in this case also the decoder. Since the discriminator will output smaller values for samples from the prior distribution, otherwise it will output higher values. As a result, the discriminator seeks to maximize the loss. However, the generator wants the discriminator to think that the generated data

**Table 2**. Performance Comparison.

| Category | Method | ROC-AUC |
|---|---|---|
| Non-Conditional Model | AE | 60.41% |
| | VAE | 62.66% |
| Conditional Model | AAE | 63.88% |
| | CVAE | 64.67% |
| | ACVAE (ours) | 66.25% |

is a sample from a priori distribution, so the generator wants to minimize the loss. This constitutes an adversarial learning process between discriminator $D$ and generator $G$ to concurrently train their respective neural networks.

**Deviation Metric.** Similar to the previous work [7], we used the standard mean square error (MSE) as a performance function to compute the deviation between the input data and the reconstructed output as follows:

$$D_{\text{MSE}} = \frac{1}{|R|} \sum_{i \in R} (x_i - \hat{x}_i)^2, \quad (3)$$

where $x_i$ is the value of the brain region $i$ after normalization, $\hat{x}_i$ is the value of the brain region $i$ reconstructed by the decoder. $R$ is a set representing all brain regions of interest and $|R|$ denotes the cardinality of $R$ (i.e., $|R| = 100$).

## 3. EXPERIMENTS

**Experimental Settings.** We divided the whole data into a training set and a test set. The training set was obtained by randomly selecting 80% of the HC subjects, and the rest was grouped into the test set along with the AD patients. Specifically, we scaled the ROI values of each subject by dividing them by the ICV (as a means of adjusting for different brain sizes). Then, we normalized the training and test sets separately using the robust scaler method from the scikit-learn library, which is robust to outliers. We first scaled the features of the training set by subtracting the median and then dividing by the interquartile range (25% value – 75% value). After this, we used the same statistics (median and interquartile range) of training set to normalize the test set. In order to ensure the robustness of the results, we used the bootstrap resampling technique to repeat the operation 10 times and reported the averaged results.

**Competing Methods.** To demonstrate the performance of the proposed ACVAE, we compared it with four other deep normative modeling methods: the vanilla AE [4], VAE [6], CVAE [6], and AAE [17], each of which represents a different normative strategy. In particular, both CVAE and AAE are conditional autoencoder methods, and we also used age, gender, and ICV as conditional variables of these models to have a fair comparison.

**Parameter Screening.** In the proposed ACVAE, both the encoder and decoder architectures are highly flexible in code length and code rate, here we used two hidden layers
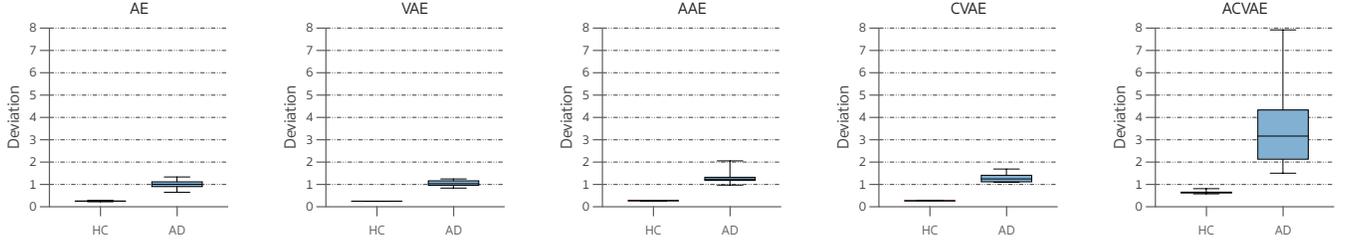
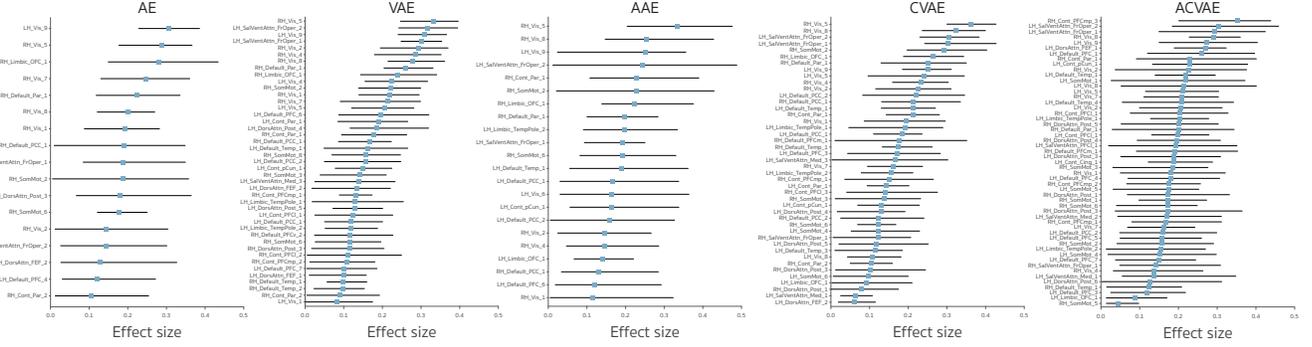**Fig. 2**. The observed mean deviation of HC and AD for each method.



**Fig. 3**. Mean effect sizes for HC and AD. The $y$-axis depicts brain regions selected for each method with significant differences.

for each module. We filtered different potential latent dimensions (10, 20) and encoder/decoder layer sizes (90, 100, and 110). We used the Adam optimizer for training with a total of 200 epochs and a cyclic learning rate containing a minimum bound and a maximum bound [18], where the minimum bound is set to 0.0001 and the maximum bound is set to 0.005. The decay parameter gamma was set to 0.98. Finally, the mini-batch method with a batch size of 256 was used for training. In addition, to avoid the gradient vanishing problem, we used LeakyReLU as our activation function. The weight coefficient of $L_{CVAE}$ loss was chosen among {0.1, 1, 5, 10, 100, 1000}, and the weight coefficient of the discriminator loss $L_{Adv}$ was chosen among {0, 2, 4, ..., 200}. For other competing methods, we use their public codes and the same parameter settings for a fair comparison.

**Results.** Table 2 shows the performance comparison of five methods in terms of the ROC-AUC score. From the experimental results, it can be observed that the proposed method outperforms the other competing methods, especially our model improves from 64.67% to 66.25% on the basis of CVAE, which is the best baseline method. Additionally, we observed that the method that takes conditional variables into account is superior to the method that does not. This demonstrated the ability of ACVAE to separate the influence of covariates from the latent vectors. Moreover, Fig. 2 shows the deviation boxplots for the HC and AD groups for each model in the test set, which is the corresponding deviation average over 10 repeated experiments. It can be seen that our ACVAE model can better distinguish HC and AD. Furthermore, we explore which brain regions had a larger effect.

To this end, we calculated 95% confidence intervals for the effect size of the difference in mean deviations for HC and AD. If the interval contains 0, it means that the difference is insignificant, otherwise, it means that the difference is significant. Fig. 3 shows the selected brain regions of each method with the most significant differences. Compared to the baselines, our model is capable of detecting more brain regions with significant effects. Taken together, this suggests that our method is more robust and sensitive.

## 4. CONCLUSION

In this paper, a new normative model is presented for quantifying deviations in Alzheimer's disease (AD) and health control (HC) at the individual level, named as adversarial conditional variational autoencoder (ACVAE). Unlike the case-control studies, ACVAE does not require training in a dataset with a reasonable balance of AD and HC groups. It is trained with only HCs, allowing it to use large cohorts of HC participants. Besides, it enables to reduce the effect of covariates in the external cohort, which highlights the potential for covariate adjustment. We validated ACVAE on OASIS-3 dataset and the experimental results showed that our approach is more effective to identify deviation values than existing models, demonstrating its potential for identifying minor pathogenic effects. Because this kind of disorder is associated with profound changes in brain morphology that are not present in the training set, this pattern is expected. This means that our approach could be applied to develop more reliable and personalized treatment plans for a variety of patients.

## 5. COMPLIANCE WITH ETHICAL STANDARDS

This research study was conducted retrospectively using human subject data made available in open access by OASIS-3 database [9]. Ethical approval was not required as confirmed by the license attached with the open access data.

## 6. REFERENCES

[1] G Li, YC Shen, YT Li, CH Chen, YW Zhau, and JM Silverman, "A case-control study of alzheimer's disease in china," *Neurology*, vol. 42, no. 8, pp. 1481–1481, 1992.

[2] AmanPreet Badhwar, G Peggy McFall, Shraddha Sapkota, Sandra E Black, Howard Chertkow, Simon Duchesne, Mario Masellis, Liang Li, Roger A Dixon, and Pierre Bellec, "A multiomics approach to heterogeneity in alzheimer's disease: focused review and roadmap," *Brain*, vol. 143, no. 5, pp. 1315–1331, 2020.

[3] Nastaran Mohammadian Rad, Twan Van Laarhoven, Cesare Furlanello, and Elena Marchiori, "Novelty detection using deep normative modeling for imu-based abnormal movement monitoring in parkinson's disease and autism spectrum disorders," *Sensors*, vol. 18, no. 10, pp. 3533, 2018.

[4] Maxime Chamberland, Sila Genc, Chantal MW Tax, Dmitri Shastin, Kristin Koller, Erika P Raven, Adam Cunningham, Joanne Doherty, Marianne van den Bree, Greg D Parker, et al., "Detecting microstructural deviations in individuals with deep diffusion mri tractometry," *Nature Computational Science*, vol. 1, no. 9, pp. 598–606, 2021.

[5] Walter HL Pinaya, Andrea Mechelli, and João R Sato, "Using deep autoencoders to identify abnormal brain structural patterns in neuropsychiatric disorders: A large-scale multi-sample study," *Human brain mapping*, vol. 40, no. 3, pp. 944–954, 2019.

[6] Ana Lawry Aguila, James Chapman, Mohammed Janahi, and Andre Altmann, "Conditional vaes for confound removal and normative modelling of neurodegenerative diseases," in *MICCAI*, 2022, pp. 430–440.

[7] Walter HL Pinaya, Cristina Scarpazza, Rafael Garcia-Dias, Sandra Vieira, Lea Baecker, Pedro F da Costa, Alberto Redolfi, Giovanni B Frisoni, Michela Pievani, Vince D Calhoun, et al., "Using normative modelling to detect disease progression in mild cognitive impairment and alzheimer's disease in a cross-sectional multi-cohort study," *Scientific Reports*, vol. 11, no. 1, pp. 1–13, 2021.

[8] Jaehyeon Kim, Jungil Kong, and Juhee Son, "Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech," in *ICML*. PMLR, 2021, pp. 5530–5540.

[9] Pamela J LaMontagne, Tammie LS Benzinger, John C Morris, Sarah Keefe, Russ Hornbeck, Chengjie Xiong, Elizabeth Grant, Jason Hassenstab, Krista Moulder, Andrei G Vlassenko, et al., "Oasis-3: longitudinal neuroimaging, clinical, and cognitive dataset for normal aging and alzheimer disease," *MedRxiv*, 2019.

[10] Oscar Esteban, Christopher J Markiewicz, Ross W Blair, Craig A Moodie, A Ilkay Isik, Asier Erramuzpe, James D Kent, Mathias Goncalves, Elizabeth DuPre, Madeleine Snyder, et al., "fmriprep: a robust preprocessing pipeline for functional mri," *Nature methods*, vol. 16, no. 1, pp. 111–116, 2019.

[11] Brian B Avants, Charles L Epstein, Murray Grossman, and James C Gee, "Symmetric diffeomorphic image registration with cross-correlation: evaluating automated labeling of elderly and neurodegenerative brain," *Medical image analysis*, vol. 12, no. 1, pp. 26–41, 2008.

[12] Douglas N Greve and Bruce Fischl, "Accurate and robust brain image alignment using boundary-based registration," *Neuroimage*, vol. 48, no. 1, pp. 63–72, 2009.

[13] Mark Jenkinson, Peter Bannister, Michael Brady, and Stephen Smith, "Improved optimization for the robust and accurate linear registration and motion correction of brain images," *Neuroimage*, vol. 17, no. 2, pp. 825–841, 2002.

[14] Raimon HR Pruim, Maarten Mennes, Daan van Rooij, Alberto Llera, Jan K Buitelaar, and Christian F Beckmann, "Ica-aroma: A robust ica-based strategy for removing motion artifacts from fmri data," *Neuroimage*, vol. 112, pp. 267–277, 2015.

[15] Alexander Schaefer, Ru Kong, Evan M Gordon, Timothy O Laumann, Xi-Nian Zuo, Avram J Holmes, Simon B Eickhoff, and BT Thomas Yeo, "Local-global parcellation of the human cerebral cortex from intrinsic functional connectivity mri," *Cerebral cortex*, vol. 28, no. 9, pp. 3095–3114, 2018.

[16] Zhihua Wang, Stefano Rosa, and Andrew Markham, "Learning the intuitive physics of non-rigid object deformations," in *NIPS Workshops*, 2018, vol. 3.

[17] Walter HL Pinaya, Cristina Scarpazza, Rafael Garcia-Dias, Sandra Vieira, Lea Baecker, Pedro F da Costa, Alberto Redolfi, Giovanni B Frisoni, Michela Pievani, et al., "Normative modelling using deep autoencoders: a multi-cohort study on mild cognitive impairment and alzheimer's disease," *bioRxiv*, 2020.

[18] Leslie N Smith, "Cyclical learning rates for training neural networks," in *WACV*. IEEE, 2017, pp. 464–472.