

Attention boosted Individualized Regression

Presenter: **YANG Guang**
Department of Data Science
City University of Hong Kong

Joint with **CAO Yuan** and **FENG Long**
Department of Statistics and Actuarial Science
The University of Hong Kong

NeurIPS 2024



Department of Data Science
香港城市大學
City University of Hong Kong



SCHOOL OF
COMPUTING &
DATA SCIENCE
The University of Hong Kong

Department of 統計及精算學系
Statistics & Actuarial Science
THE UNIVERSITY OF HONG KONG

Contents

- 1 Introduction
- 2 Model
 - Attention boosted individualized regression (AIR)
 - Connection with self-attention mechanism
 - Alternating minimization algorithm
- 3 Theoretical analysis
- 4 Experiments
 - Simulation study
 - Ablation study
 - ADNI study

Contents

1 Introduction

2 Model

- Attention boosted individualized regression (AIR)
- Connection with self-attention mechanism
- Alternating minimization algorithm

3 Theoretical analysis

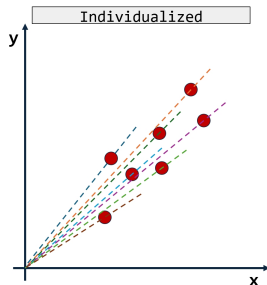
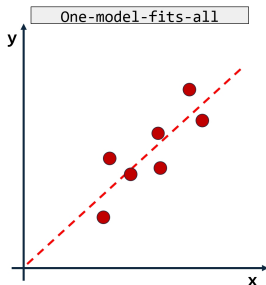
4 Experiments

- Simulation study
- Ablation study
- ADNI study

Introduction

Principle

- One-model-fits-all
- Individualized model
 - Varying-coefficient model (Time, spatially, etc.)
 - Sample-specific model (Coefficients similarity)
 - Self-attention mechanism (Contextual relationship)



Contents

1 Introduction

2 Model

- Attention boosted individualized regression (AIR)
- Connection with self-attention mechanism
- Alternating minimization algorithm

3 Theoretical analysis

4 Experiments

- Simulation study
- Ablation study
- ADNI study

Individualized model

Individualized linear regression

- Ordinary linear regression model for matrix-valued input

$$y_i = \langle \mathbf{X}_i, \mathbf{C} \rangle + \varepsilon_i, \quad i = 1, \dots, n.$$

- Individualized linear regression model for matrix-valued input

$$y_i = \langle \mathbf{X}_i, \mathbf{C}_i \rangle + \varepsilon_i, \quad i = 1, \dots, n.$$

To study internal relations among patches, we take reshaped images as input, namely $\mathbf{X}_i = \mathcal{R}_{(d_1, d_2)}(\mathbf{X}_i^{\text{ori}})$. Subsequently, $\mathbf{X}_i \mathbf{X}_i^\top$ reflects relations among patches within i -th sample, called **internal relations**.



Attention boosted individualized regression

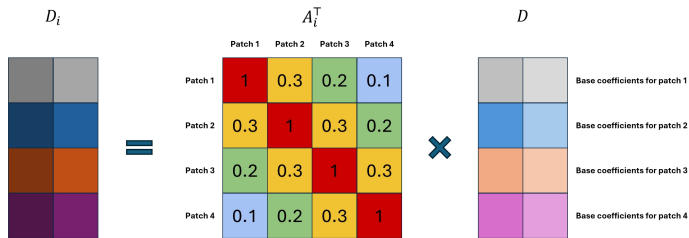
Decomposition of coefficients

I. Decomposition:

$$\begin{array}{c} C_i \\ \uparrow \\ \text{Individualized coefficient} \end{array} = \begin{array}{c} C \\ \uparrow \\ \text{homogeneous coefficient} \end{array} + \begin{array}{c} D_i \\ \uparrow \\ \text{heterogeneous coefficient} \end{array} .$$

II. Aggregation: $D_i = A_i^\top D$.

III. Internal relations: $A_i = g(X_i W X_i^\top)$.



Attention boosted individualized regression

Rotational correlation

For any two vectors \mathbf{u} and \mathbf{v} , the rotational correlation is defined as

$$\max_{\mathbf{H}} \mathbf{u}^\top \mathbf{H} \mathbf{v},$$

where the matrix \mathbf{H} is usually required to be orthogonal. This rotational correlation aims to find the maximized correlation between \mathbf{u} and \mathbf{v} with the best possible rotation.

Explanation

Note that the (j, k) -th element of the sample-specific factor can be written as

$$\{\mathbf{A}_i\}_{jk} = \{\mathbf{X}_i\}_j \cdot \mathbf{W} \{\mathbf{X}_i\}_{k^\cdot}^\top,$$

where $\{\mathbf{X}_i\}_j$ and $\{\mathbf{X}_i\}_{k^\cdot}$ are the j -th and k -th rows of \mathbf{X}_i , respectively.

To say, $\{\mathbf{A}_i\}_{jk}$ is related to the rotational correlation between $\{\mathbf{X}_i\}_j$ and $\{\mathbf{X}_i\}_{k^\cdot}$.

Our goal is not to maximize the rotational correlation between $\{\mathbf{X}_i\}_j$ and $\{\mathbf{X}_i\}_{k^\cdot}$, but to find the optimal rotation that achieves the best fitting for the responses.

Attention boosted individualized regression

Attention boosted individualized regression (AIR)

Now we obtain our individualized model in the following form:

$$y_i = \underbrace{\langle \mathbf{X}_i, \mathbf{C} \rangle}_{\text{homogeneous}} + \underbrace{\langle \mathbf{X}_i, g(\mathbf{X}_i \mathbf{W} \mathbf{X}_i^\top)^\top \mathbf{D} \rangle}_{\text{heterogeneous}} + \varepsilon_i.$$

Here, $\mathbf{C}, \mathbf{D} \in \mathbb{R}^{p \times d}$, and $\mathbf{W} \in \mathbb{R}^{d \times d}$ are the coefficient matrices to be learned.

Optimization

To learn \mathbf{C} , \mathbf{D} and \mathbf{W} , we propose the following penalized minimization problem

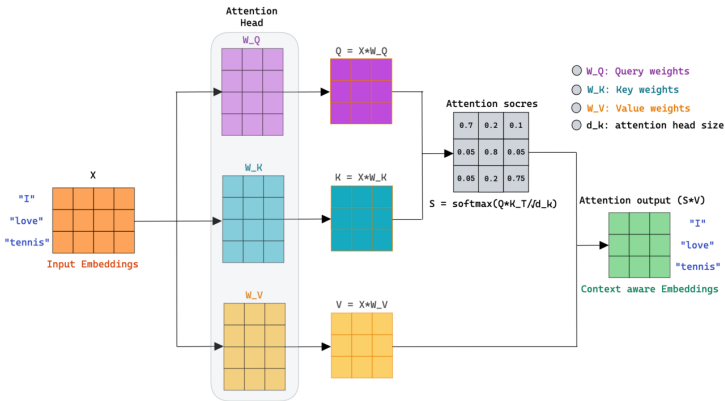
$$\begin{aligned} \min_{\mathbf{C}, \mathbf{D}, \mathbf{W}} \quad & \frac{1}{n} \sum_{i=1}^n (y_i - \langle \mathbf{X}_i, \mathbf{C}_i \rangle)^2 + \lambda_1 \|\mathbf{C}\|_F^2 + \lambda_2 \|\mathbf{D}\|_F^2, \\ \text{s.t.} \quad & \mathbf{C}_i = \mathbf{C} + g(\mathbf{X}_i \mathbf{W} \mathbf{X}_i^\top)^\top \mathbf{D}, \quad \|\mathbf{W}\|_F = 1, \end{aligned}$$

where $\|\cdot\|_F$ is the Frobenius norm and λ_1 and λ_2 are regularization parameters to balance the homogeneous and heterogeneous effects.

Self-attention mechanisms

Given the input $\mathbf{X} \in \mathbb{R}^{n \times d}$, with $\mathbf{Q} = \mathbf{XW}_Q$, $\mathbf{K} = \mathbf{XW}_K$ and $\mathbf{V} = \mathbf{XW}_V$, the Scaled Dot-Product Attention mechanism computes

$$f(\mathbf{X}) = \text{softmax}\left(\frac{\mathbf{QK}^\top}{\sqrt{d_k}}\right) \mathbf{V}.$$



Self-attention mechanisms

General self-attention

Beyond softmax function, normalization in attention could also be accomplished using a general function $g(\cdot)$. We obtain the following generalized attention:

$$f(\mathbf{X}) = g(\mathbf{Q}\mathbf{K}^\top) \mathbf{V}.$$

Linear self-attention

Shen et al. (2021) considered a linear function as scaling normalization:

$$f(\mathbf{X}) = \frac{1}{n} \mathbf{Q}\mathbf{K}^\top \mathbf{V}.$$

Linear attention mechanisms are efficient because they bypass the need to compute $n \times n$ matrices by using associative multiplication, reducing complexity from $O(n^2)$ to $O(n)$. While on the other hand, experiments show that Linear attentions does not result in a significant compromise in performance.

Individualized regression and attention mechanism

Focus on the individualized model with only heterogeneous part

$$y_i = \langle \mathbf{X}_i, g(\mathbf{X}_i \mathbf{W} \mathbf{X}_i^\top)^\top \mathbf{D} \rangle + \varepsilon_i.$$

Proposition 1

Suppose the above model holds and \mathbf{W} and \mathbf{D} could be decomposed as below

(I) $\mathbf{W} = \mathbf{W}_Q \mathbf{W}_K^\top$ for two matrices $\mathbf{W}_Q, \mathbf{W}_K \in \mathbb{R}^{d \times d_k}$ with $d_k \leq d$,

(II) $\mathbf{D} = \mathbf{B} \mathbf{W}_V^\top$ for two matrices $\mathbf{B}, \mathbf{W}_V \in \mathbb{R}^{d \times d_v}$ with $d_v \leq d$.

where $d_k, d_v \leq d$. Then, for each sample \mathbf{X}_i , the following equation holds

$$\langle \mathbf{X}_i, g(\mathbf{X}_i \mathbf{W} \mathbf{X}_i^\top)^\top \mathbf{D} \rangle = \langle g(\mathbf{Q}_i \mathbf{K}_i^\top) \mathbf{V}_i, \mathbf{B} \rangle,$$

where

$$\mathbf{Q}_i = \mathbf{X}_i \mathbf{W}_Q, \quad \mathbf{K}_i = \mathbf{X}_i \mathbf{W}_K, \quad \mathbf{V}_i = \mathbf{X}_i \mathbf{W}_V. \quad (1)$$

Computation

Equation

$$\langle \mathbf{X}_i, \mathbf{X}_i \mathbf{W}^\top \mathbf{X}_i^\top \mathbf{D} \rangle = \langle \mathbf{X}_i^\top \mathbf{D} \mathbf{X}_i^\top \mathbf{X}_i, \mathbf{W} \rangle = \langle \mathbf{X}_i \mathbf{W} \mathbf{X}_i^\top \mathbf{X}_i, \mathbf{D} \rangle.$$

Initialization

Let $\mathbf{w} = \text{vec}(\mathbf{W})$ and $\mathbf{d} = \text{vec}(\mathbf{D})$ be the vectorization of \mathbf{W} and \mathbf{D} . It holds that

$$\langle \mathbf{X}_i^\top \mathbf{D} \mathbf{X}_i^\top \mathbf{X}_i, \mathbf{W} \rangle = \langle \mathbf{Z}_i, \mathbf{w} \mathbf{d}^\top \rangle, \quad \text{where } \mathbf{Z}_i = (\mathbf{X}_i^\top \mathbf{X}_i) \otimes \mathbf{X}_i^\top$$

We start our algorithm by initializing \mathbf{w} as the top left singular vector of $\sum_{i=1}^n y_i \mathbf{Z}_i$,

$$\hat{\mathbf{w}}^{(0)} = \text{SVD}_u \left(\sum_{i=1}^n y_i \mathbf{Z}_i \right).$$

Computation

Alternating minimization algorithm

Given $\widehat{\mathbf{W}}^{(t-1)}$, denote $\mathbf{U}_i^{(t-1)} = \mathbf{X}_i \widehat{\mathbf{W}}^{(t-1)} \mathbf{X}_i^\top \mathbf{X}_i$,

$$\left(\widehat{\mathbf{C}}^{(t)}, \widehat{\mathbf{D}}^{(t)}\right) = \operatorname{argmin}_{\mathbf{C}, \mathbf{D}} \frac{1}{n} \sum_{i=1}^n \left(y_i - \left\langle \left[\mathbf{X}_i, \mathbf{U}_i^{(t-1)} \right], \left[\mathbf{C}, \mathbf{D} \right] \right\rangle \right)^2 + \lambda_1 \|\mathbf{C}\|_F^2 + \lambda_2 \|\mathbf{D}\|_F^2.$$

Given $\left(\widehat{\mathbf{C}}^{(t)}, \widehat{\mathbf{D}}^{(t)}\right)$,

$$\widehat{\mathbf{W}}^{(t)} = \operatorname{argmin}_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^n \left(y_i - \left\langle \mathbf{X}_i, \widehat{\mathbf{C}}^{(t)} \right\rangle - \left\langle \mathbf{X}_i^\top \widehat{\mathbf{D}}^{(t)} \mathbf{X}_i^\top \mathbf{X}_i, \mathbf{W} \right\rangle \right)^2,$$

$$\widehat{\mathbf{W}}^{(t)} = \widehat{\mathbf{W}}^{(t)} / \|\widehat{\mathbf{W}}^{(t)}\|_F.$$

Contents

1 Introduction

2 Model

- Attention boosted individualized regression (AIR)
- Connection with self-attention mechanism
- Alternating minimization algorithm

3 Theoretical analysis

4 Experiments

- Simulation study
- Ablation study
- ADNI study

Theoretical analysis

We focus on the heterogeneous part of the model. Let $\mathbf{w} = \text{vec}(\mathbf{W})$ and $\mathbf{d} = \text{vec}(\mathbf{D})$, the optimization problem could be written as

$$\min_{\mathbf{d}, \mathbf{w}} \frac{1}{n} \sum_{i=1}^n \left\{ y_i - \left\langle \left(\mathbf{X}_i^\top \mathbf{X}_i \right) \otimes \mathbf{X}_i^\top, \mathbf{w} \mathbf{d}^\top \right\rangle \right\}^2 + \lambda_2 \|\mathbf{d}\|_2^2.$$

For the rearranged images \mathbf{X}_i for $i = 1, \dots, n$, we define

$$\mathbf{Z} = \left(\text{vec} \left\{ \left(\mathbf{X}_1^\top \mathbf{X}_1 \right) \otimes \mathbf{X}_1^\top \right\}, \dots, \text{vec} \left\{ \left(\mathbf{X}_n^\top \mathbf{X}_n \right) \otimes \mathbf{X}_n^\top \right\} \right)^\top.$$

For the new feature matrix \mathbf{Z} , we suppose the following RIP condition.

Condition 1

(Restricted Isometry Property) For each integer $r = 1, 2, \dots$, a matrix $\mathbf{P} \in \mathbb{R}^{n \times q_1 q_2}$ is said to satisfy the r -RIP condition with constant $\delta_r \in (0, 1)$, if for all $\mathbf{M} \in \mathbb{R}^{q_1 \times q_2}$ of rank at most r , it holds that

$$(1 - \delta_r) \|\mathbf{M}\|_F^2 \leq 1/n \|\mathbf{P} \text{vec}(\mathbf{M})\|_2^2 \leq (1 + \delta_r) \|\mathbf{M}\|_F^2.$$

Theorems

Theorem 1. Estimation (Informal)

Suppose the heterogeneous model holds and solved by alternating minimization algorithm. Denote μ_0 as the initial distance, $\kappa_1, \kappa_2 < 1$ as contraction quantities and τ_1 and τ_2 as noise-related terms. Then, after t iterations we have

$$\begin{aligned} \text{dist} \left(\widehat{\mathbf{W}}^{(t)}, \mathbf{W} \right) &\leq (\kappa_1 \kappa_2)^t \mu_0 + \frac{\kappa_1 \tau_2 + \tau_1}{1 - \kappa_1 \kappa_2}, \\ \text{dist} \left(\widehat{\mathbf{D}}^{(t)}, \mathbf{D} \right) &\leq \kappa_1^{t-1} \kappa_2^t \mu_0 + \frac{\kappa_2 \tau_1 + \tau_2}{1 - \kappa_1 \kappa_2}. \end{aligned}$$

Theorem 2. Prediction (Informal)

With similar conditions in Theorem 1, after t iterations we have

$$\|\widehat{\mathbf{Y}}^{(t)} - \mathbf{Y}\|_2 \leq 3 \|\mathbf{D}\|_F \sqrt{1 + \delta_2} \left\{ (\nu_1 \nu_2)^{t-1} \mu_0 + \frac{\tau_1 + \tau_2}{1 - \nu_1 \nu_2} \right\}.$$

Contents

1 Introduction

2 Model

- Attention boosted individualized regression (AIR)
- Connection with self-attention mechanism
- Alternating minimization algorithm

3 Theoretical analysis

4 Experiments

- Simulation study
- Ablation study
- ADNI study

Simulation study

- Size of images $\mathbf{X}_i^{\text{ori}}$: 28×28 .
- Sample size: 4000 for training and 1000 for testing.
- Noise: $\varepsilon_i \sim \mathcal{N}(0, 1)$.
- Coefficients: \mathbf{C} and \mathbf{D} are generated as two circles depicted in Figure 1.
- Internal relations: assumed among blocks of size 4×4 within each image, where two blocks at random locations are correlated.
- Consider two cases: $y_i = \langle \mathbf{X}_i, \mathbf{C} \rangle + \langle \mathbf{X}_i, \mathbf{D}_i \rangle + \varepsilon_i$, subject to $\mathbf{D}_i = \mathbf{A}_i^\top \mathbf{D}$.
 - C1. With specific \mathbf{W} : $\mathbf{A}_i = \mathbf{X}_i \mathbf{W} \mathbf{X}_i^\top$ and $\mathbf{W} = 2 \cdot \mathbf{u}_1 \mathbf{v}_1^\top + 1 \cdot \mathbf{u}_2 \mathbf{v}_2^\top$.
 - C2. Without specific \mathbf{W} : $\mathbf{A}_i = \text{coef}(\mathbf{X}_i)$.
- Comparison methods:
 - Regularized matrix regression (LRMR, Zhou and Li, 2014)
 - Tensor regression with lasso penalty (TRLasso, Zhou et al., 2013)
 - Deep kronecker network (DKN, Feng and Yang, 2023)
 - Vision transformer (ViT, Dosovitskiy et al., 2020)

Simulation study

Furthermore, we consider different degrees of individuation (DI) by the relative magnitude of the heterogeneous part and homogeneous part. Specifically,

$$DI = \sqrt{\sum_{i=1}^n \langle \mathbf{X}_i, \mathbf{D}_i \rangle^2 / \sum_{i=1}^n \langle \mathbf{X}_i, \mathbf{C} \rangle^2}$$

Table: Prediction errors of different methods.

Methods	DI=0.5	DI=1.0	DI=2.0
	Case 1		
AIR	4.422 (0.130)	8.102 (0.325)	10.599 (0.816)
LRMR	6.616 (0.020)	13.101 (0.040)	26.239 (0.081)
TRLasso	8.215 (0.021)	14.465 (0.044)	27.007 (0.085)
DKN	4.886 (0.018)	7.028 (0.032)	11.741 (0.043)
ViT	18.429 (0.049)	18.351 (0.047)	24.098 (0.069)
Case 2			
AIR	3.590 (0.046)	6.632 (0.022)	13.002 (0.044)
LRMR	6.766 (0.018)	13.408 (0.037)	26.864 (0.074)
TRLasso	8.337 (0.021)	14.739 (0.039)	27.484 (0.073)
DKN	8.269 (0.018)	14.964 (0.034)	28.686 (0.060)
ViT	24.492 (0.063)	29.939 (0.084)	44.036 (0.111)

Simulation study

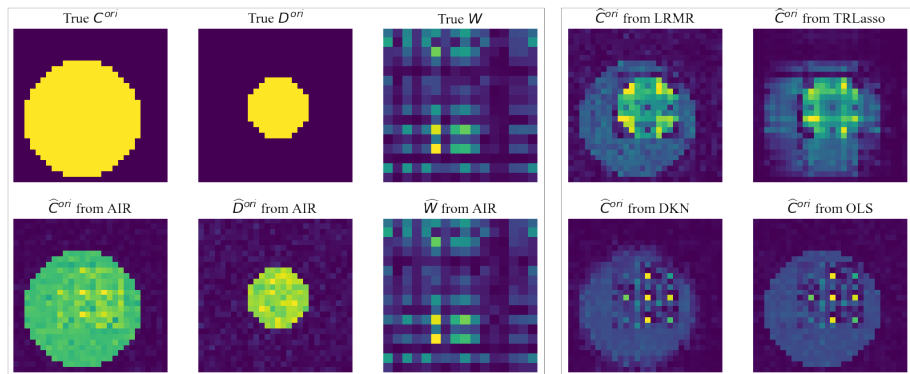


Figure: Case 1 simulation results with $DI = 1.0$. The first three columns show true parameters and estimations from AIR. The last two columns show estimations from other methods except ViT, as it has no explicit coefficient matrix. An additional OLS estimation is added for reference.

Simulation study

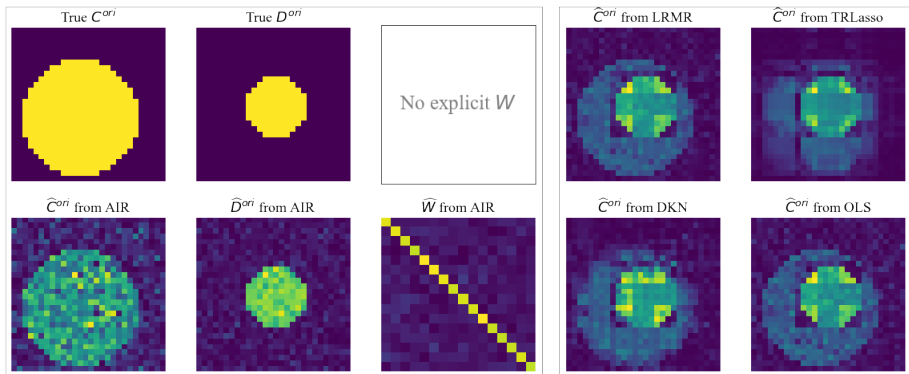
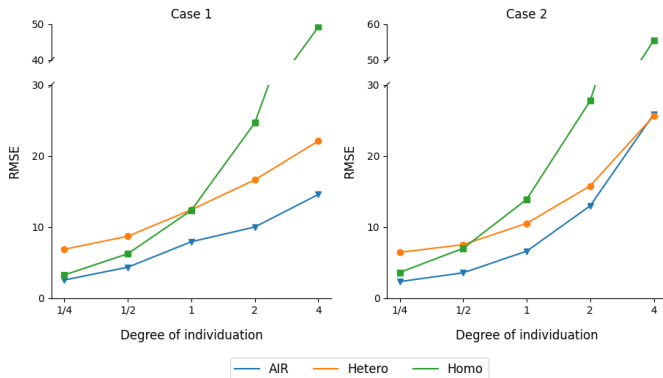


Figure: Case 2 simulation results with $DI = 1.0$. There does not exist an explicit true \mathbf{W} while the internal relation matrix \mathbf{A}_i is computed directly by patchwise Pearson correlation coefficients. Because such \mathbf{A}_i is close to a diagonal matrix, it is rational that $\widehat{\mathbf{W}}$ from AIR is close to a diagonal matrix.

Ablation study

We compare

- (1) AIR: $y_i = \langle \mathbf{X}_i, \mathbf{C} \rangle + \langle \mathbf{X}_i, \mathbf{D}_i \rangle + \varepsilon_i$, subject to $\mathbf{D}_i = \mathbf{X}_i \mathbf{W}^\top \mathbf{X}_i^\top \mathbf{D}$.
- (2) Hetero: $y_i = \langle \mathbf{X}_i, \mathbf{D}_i \rangle + \varepsilon_i$, subject to $\mathbf{D}_i = \mathbf{X}_i \mathbf{W}^\top \mathbf{X}_i^\top \mathbf{D}$.
- (3) Homo: $y_i = \langle \mathbf{X}_i, \mathbf{C} \rangle + \varepsilon_i$.

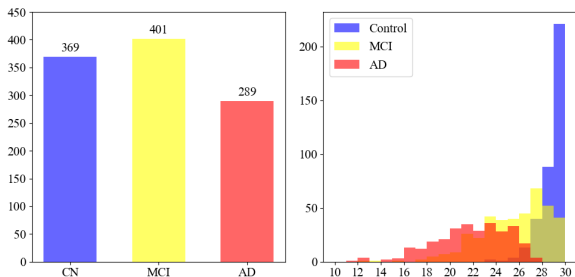


ADNI study

Alzheimer's Disease Neuroimaging Initiative (ADNI)

We collected a total of 1059 subjects from ADNI 1 and GO/2 phases with brain MRI scans, preprocessed to be of size $48 \times 60 \times 48$.

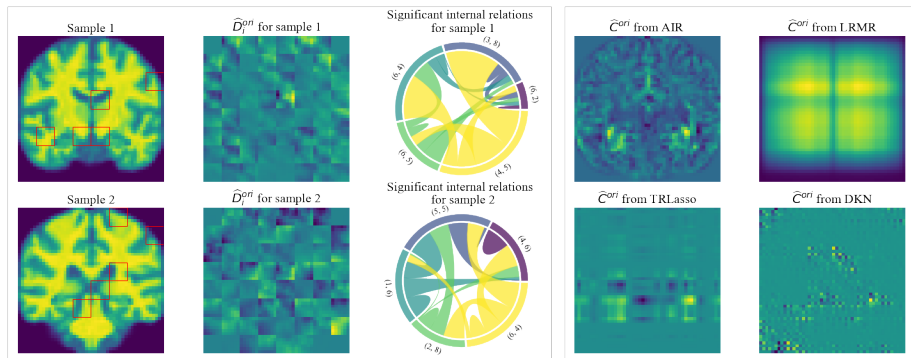
- X_i : For each subject, we extracted 10 middle coronal slices, resulting in images of size 48×48 .
- y_i : the Mini-Mental State Examination (MMSE) score.



ADNI study

Table: Prediction errors of different methods.

AIR	LRMR	TRLasso	DKN	ViT
3.145 (0.019)	3.715 (0.008)	3.292 (0.023)	3.261 (0.017)	3.282 (0.025)



Thank you!

This paper

Yang, G., Cao, Y. and Feng, L., Attention boosted Individualized Regression.
To appear in *NeurIPS 2024*.

References

- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.
- Feng, L. and Yang, G. (2023). Deep Kronecker network. *Biometrika*, 111(2):707–714.
- Shen, Z., Zhang, M., Zhao, H., Yi, S., and Li, H. (2021). Efficient attention: Attention with linear complexities. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 3531–3539.
- Zhou, H. and Li, L. (2014). Regularized matrix regression. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(2):463–483.
- Zhou, H., Li, L., and Zhu, H. (2013). Tensor regression with applications in neuroimaging data analysis. *Journal of the American Statistical Association*, 108(502):540–552.