

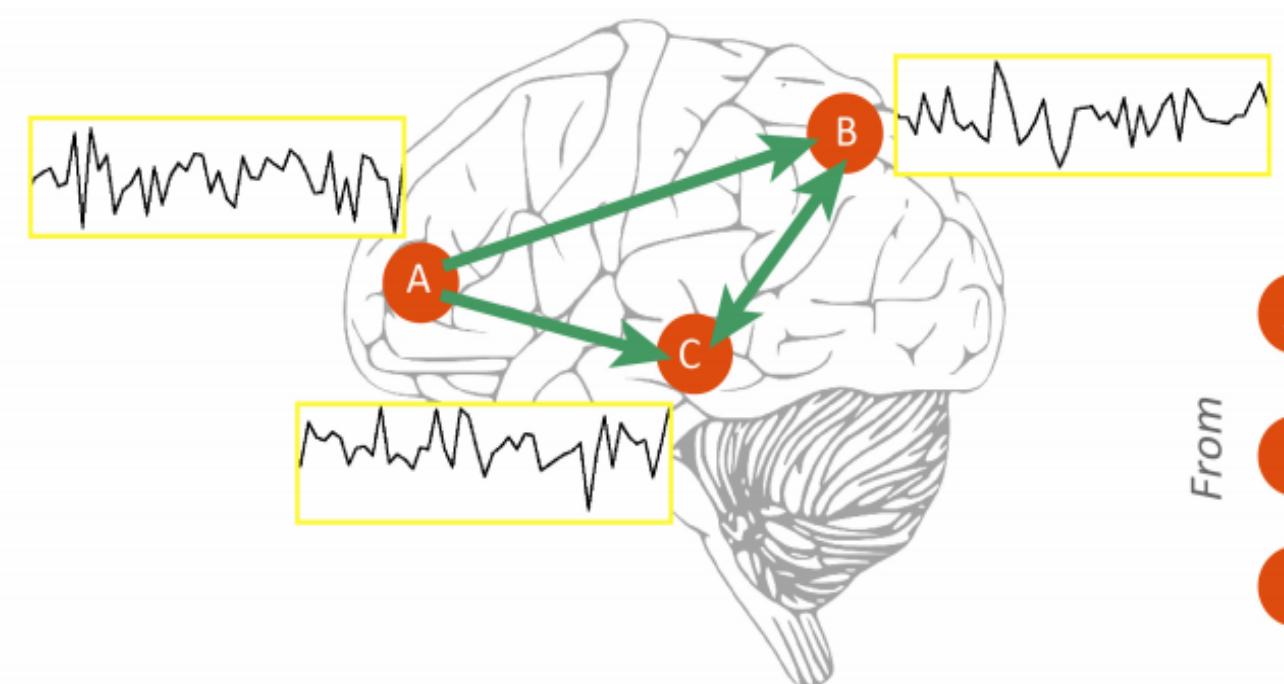
# Fingerprinting: measuring information by matching repeated measures

Charles Zheng<sup>1,2</sup>, Yuval Benjamini<sup>3</sup>, Emily S. Finn<sup>1,4</sup>, Natasha Topolski<sup>4,5</sup> and Francisco Pereira<sup>1,2</sup>

1: Functional Magnetic Resonance Imaging Facility, National Institute of Mental Health, 2: Machine Learning Team, National Institute of Mental Health, 3: Department of Statistics and Data Science, Hebrew University of Jerusalem, 4: Section of Functional Imaging Methods, National Institute of Mental Health, 5: McGovern Medical School

## Introduction

**Functional connectivity (FC)** matrices measure the degree of correlation between activation in different brain regions...



	To	A	B	C
From	A	0	0.8	0.7
B	0	0	0.5	
C	0	0.6	0	

and can be used to predict

- diagnosis (e.g. autism)
- age
- fluid intelligence
- memory
- ... and many more outcomes!

Multiple pipelines produce FC matrices, with many preprocessing choices affecting downstream tasks... creating a methodological need for quantitative evaluation of pipelines.



### MUTUAL INFORMATION FOR EVALUATING PIPELINE QUALITY

- let  $\mathbf{X}$  denote the FC matrix, and let  $\mathbf{Y}$  denote *all* the **latent subject characteristics** which include or cause the various outcomes of interest (age, diagnosis, etc.)
- mutual information  $I(\mathbf{X}; \mathbf{Y})$ : an *ideal* criterion that reflects the quality of  $\mathbf{X}$  with regard to the information it contains about  $\mathbf{Y}$ , for *many different* downstream tasks, including prediction.

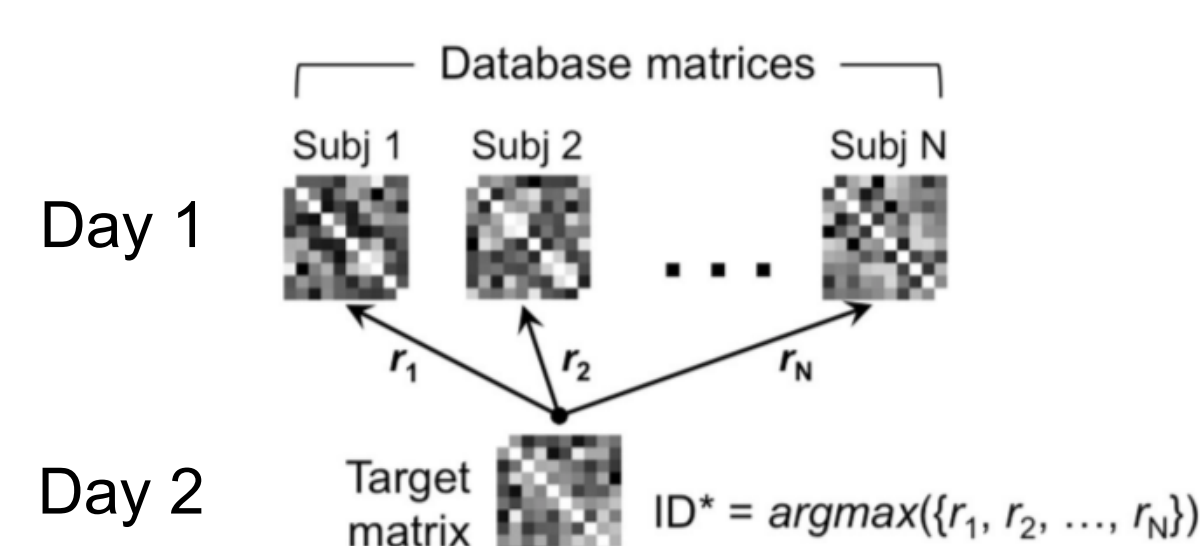
**Our contribution:** A novel approach for estimating  $I(\mathbf{X}; \mathbf{Y})$  when  $\mathbf{Y}$  is unobserved, by using a repeated measurement  $\mathbf{X}'$  as a proxy for  $\mathbf{Y}$ .

1) Fingerprinting accuracy (FA) [1] defined as accuracy of matching repeated measures from  $N$  subjects.

2) Applying results from [2] to FA yields an estimator for  $I(\mathbf{X}; \mathbf{X}')$ .

3) **Remaining challenge:** estimating  $I(\mathbf{X}; \mathbf{Y})$  from  $I(\mathbf{X}; \mathbf{X}')$

We develop a model-based estimation approach (see Theory.)



## Prior Work: estimating $I(\mathbf{X}; \mathbf{X}')$

### ESTIMATION OF MUTUAL INFORMATION

Classical nonparametric estimators of mutual information do not scale well with dimensionality of  $\mathbf{X}$  and are not practical for neuroimaging data.

Zheng and Benjamini (2016) obtain a high-dimensional estimator of mutual information based on  $k$ -class *classification accuracy*

$$\hat{I}(\mathbf{X}; \mathbf{Y}) = \frac{1}{2} \pi_k^{-1} (1 - \text{Acc})^2$$

### APPLICATION TO FINGERPRINTING

Fingerprinting accuracy, defined

$$\text{Acc} = \frac{1}{n} \sum_{i=1}^n I\{i = \arg\max_j S_{ij}\}$$

is 1-nearest neighbors classification with feature  $\mathbf{X}$  and label  $\mathbf{X}'$ .

Therefore, we have

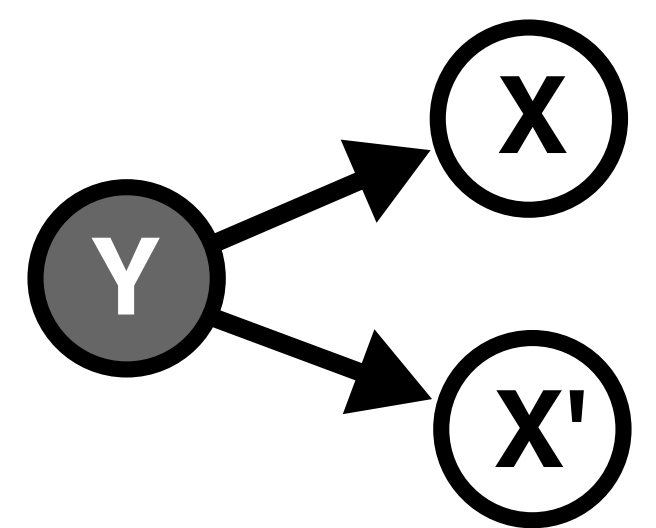
$$\hat{I}(\mathbf{X}; \mathbf{X}') = \frac{1}{2} \pi_N^{-1} (1 - \text{Acc})^2.$$

## References

- [1] Finn, Emily S., et al. "Functional connectome fingerprinting: identifying individuals using patterns of brain connectivity." *Nature neuroscience* 18.11 (2015).
- [2] Zheng, Charles Y., and Yuval Benjamini. "Estimating mutual information in high dimensions via classification error." arXiv preprint arXiv:1606.05229 (2016).
- [3] Murphy, Kevin, and Michael D. Fox. "Towards a consensus regarding global signal regression for resting state functional connectivity MRI." *Neuroimage* 154 (2017): 169-173.

## Theory: estimating $I(\mathbf{X}; \mathbf{Y})$

**Problem:** Given observed  $\mathbf{X}$  and  $\mathbf{X}'$  and unobserved  $\mathbf{Y}$ , where  $\mathbf{X}$  and  $\mathbf{X}'$  are conditionally independent given  $\mathbf{Y}$ , estimate  $I(\mathbf{X}; \mathbf{Y})$ .



**Example:**  $\mathbf{Y}$  are latent subject characteristics,  $\mathbf{X}$  and  $\mathbf{X}'$  are FC matrices measured **on the same subject** over 2 different sessions.

**Data-Processing Inequality-based estimate:** Due to conditional independence of  $\mathbf{X}'$  on  $\mathbf{X}$ , we have  $I(\mathbf{X}; \mathbf{X}') \leq I(\mathbf{X}; \mathbf{Y})$ . Hence, any estimator of  $I(\mathbf{X}; \mathbf{X}')$  serves as an estimator of  $I(\mathbf{X}; \mathbf{Y})$ . However, the Data-Processing Inequality (DPI) bound can be very loose.

**Model-based estimator:** Assume the following—

(I) **Gaussian copula:** There exists a bijection  $g(\cdot)$  such that for  $\mathbf{Z}=g(\mathbf{X})$ , we have

$$\text{Cor}(\mathbf{Y}, \mathbf{Z}) = \begin{pmatrix} \rho_1 & 0 & \cdots & 0 \\ 0 & \rho_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & \cdots & \rho_d \end{pmatrix},$$

where  $d$  is the number of independent latent components in  $\mathbf{Y}$ .

(II) **Bounded  $\rho$ :** we have  $\rho_{\max} \geq \rho_1 \geq \dots \geq \rho_d \geq \rho_{\min}$ .

**Example:**  $\mathbf{Y}$  are network activation levels; The  $\rho_d$  reflect how easy it is to infer each network activation from the measurement  $\mathbf{X}$ .

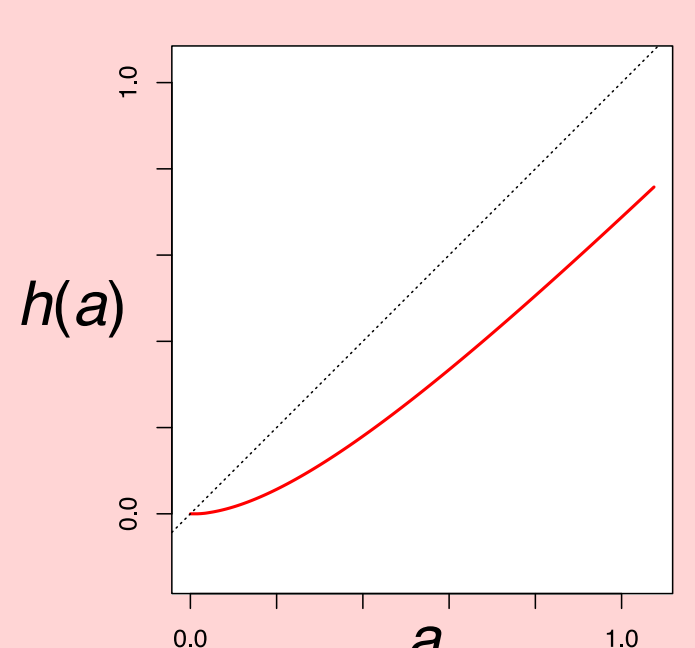
Bounded  $\rho \leftrightarrow$  components of  $\mathbf{Y}$  are neither too easy ( $\rho_{\max}$ ) nor difficult ( $\rho_{\min}$ ) to infer from  $\mathbf{X}$ .

### DERIVATION OF ESTIMATOR

1) Using assumption I:

$$I(\mathbf{X}; \mathbf{X}') = \sum_{i=1}^d I(X_i; X'_i) = \sum_{i=1}^d -\frac{1}{2} \ln(1 - \rho_i^2); \quad I(\mathbf{X}; \mathbf{Y}) = \sum_{i=1}^d I(X_i; Y_i) = \sum_{i=1}^d -\frac{1}{2} \ln(1 - \rho_i^2).$$

2) Defining  $h(a) = -\frac{1}{2} \ln(1 - \sqrt{e^{-2a}})$ , we have  $I(X_i; Y_i) = h(I(X_i; X'_i))$  and hence  $I(\mathbf{X}; \mathbf{Y}) = \sum_{i=1}^d I(X_i; Y_i) = \sum_{i=1}^d h(I(X_i; X'_i))$ .



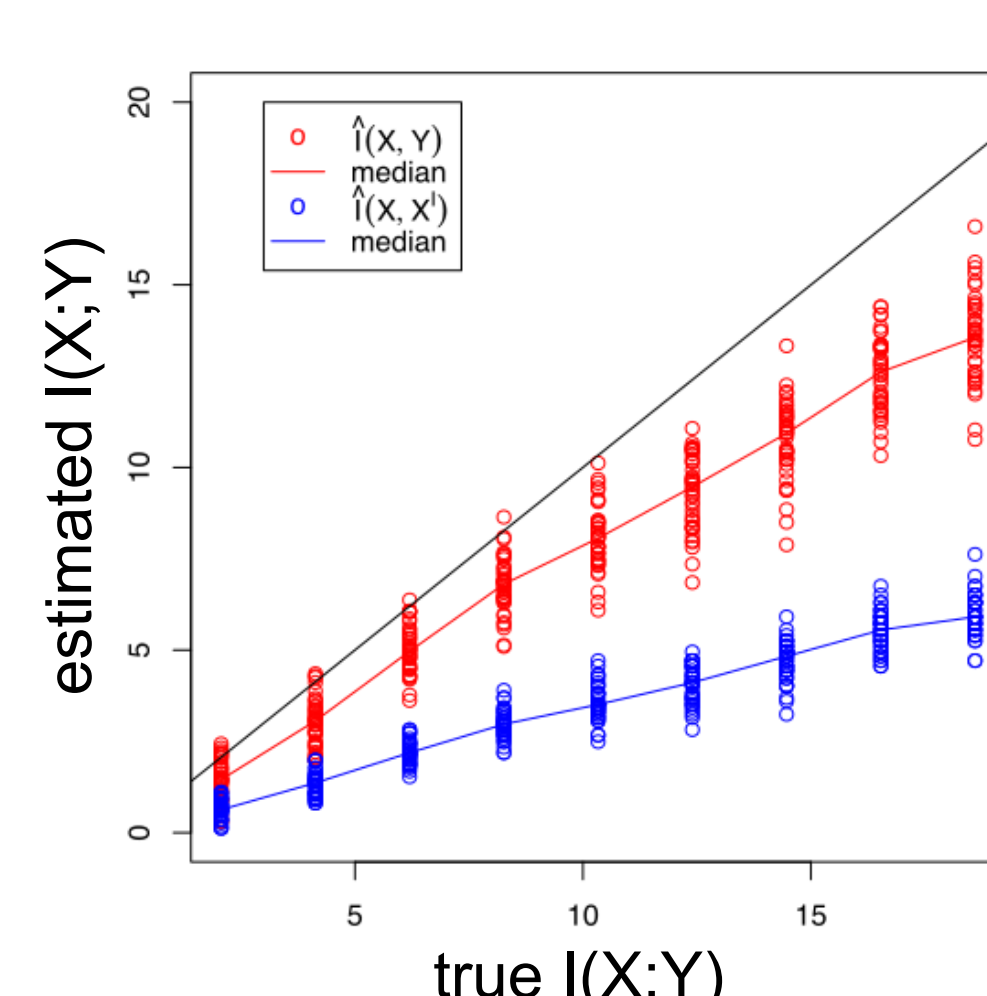
3) Using assumption II, linearly approximate  $h(a) \approx \beta_1 a + \beta_0$ , hence

$$I(\mathbf{X}; \mathbf{Y}) = \sum_{i=1}^d h(I(X_i; X'_i)) \approx \sum_{i=1}^d \beta_1 I(X_i; X'_i) + \beta_0 = \beta_1 I(\mathbf{X}; \mathbf{X}') + d\beta_0.$$

4) **Additional heuristics:** Estimate  $d$  by  $\hat{d} = \dim(\mathbf{X}) \frac{\hat{I}(\mathbf{X}; \mathbf{X}')}{\sum_{i=1}^p \hat{I}(X_i; X'_i)}$ .

Estimate  $\rho_{\text{mid}} = (\rho_{\min} + \rho_{\max})/2$  by  $\rho_{\text{mid}} = \sqrt[4]{1 - e^{-\frac{2\hat{I}(\mathbf{X}; \mathbf{X}')}{d}}}$ .

## Results: DPI vs. model-based



**Simulation:** (40 repeats per setting)

- $\mathbf{Y} \sim$  standard multivariate Gaussian
- $\dim(\mathbf{Y})$  varied from 5 to 45
- $\rho_1 = \dots = \rho_d = 0.75$
- $\dim(\mathbf{X}) = 100$ ,  $\mathbf{X} = \mathbf{AZ}$  where  $\mathbf{A}$  is a random  $d \times p$  matrix
- $N=80$  subjects

**Results:** Diagonal is ideal performance. Both **DPI (blue)** and **model-based (red)** underestimate, but model-based is less conservative.

**HCP Data:** (338 subjects)

- $\mathbf{X}$  are FC matrices,  $\dim(\mathbf{X}) \approx 36,000$
- Compared GSR[3] and non-GSR pipelines

**Results:** Non-GSR superior

with respect to fingerprinting accuracy and **DPI estimate**. GSR scores higher for **model-based  $I(\mathbf{X}; \mathbf{Y})$**  due to higher latent dimensionality.

Pipeline	Fingerprinting accuracy	DPI estimate $I(\mathbf{X}; \mathbf{Y})$	Estimated d	Model-based estimate $I(\mathbf{X}; \mathbf{Y})$
GSR	0.973	11.4	75.3	27.0
No-GSR	0.979	12.0	50.8	24.2