

# Multivariate Principal Component Analysis for Mixed-Type Functional Data with application to mHealth

Debangana Dey  
Postdoc Fellow  
National Institute of Mental Health

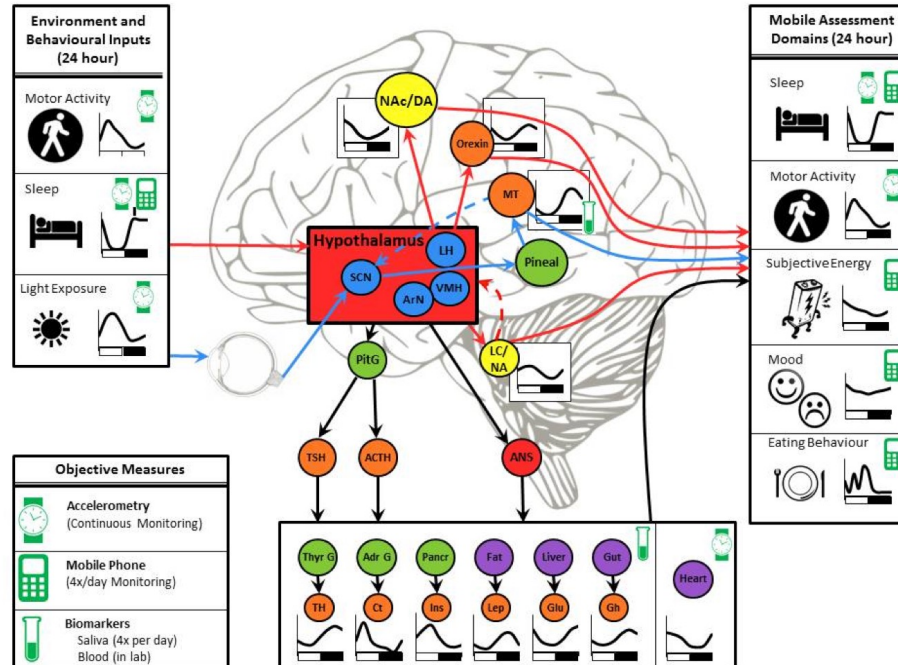
Joint work with  
Rahul Ghosal,  
Kathleen Merikangas, and  
Vadim Zipunnikov

## Disclaimer:

This work was supported by the National Institute of Mental Health Intramural Research Program.

The views and opinions expressed in this article are those of the authors and should not be construed to represent the views of any of the U.S. Government.

# Biological processes associated with regulation of homeostatic domains assessed by mHealth



Sleep



Sleep



Light



Sleep



Light



Mood



Sleep



Light

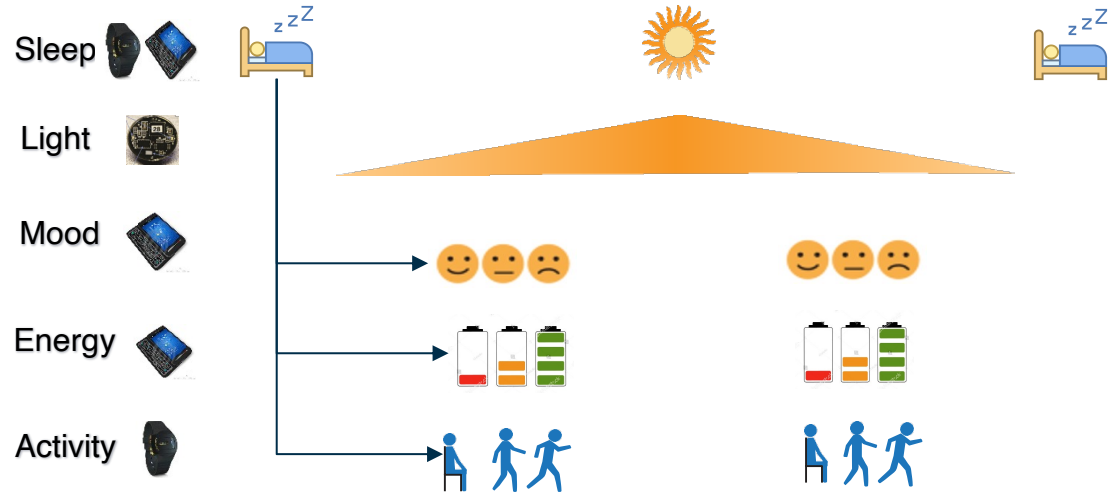


Mood

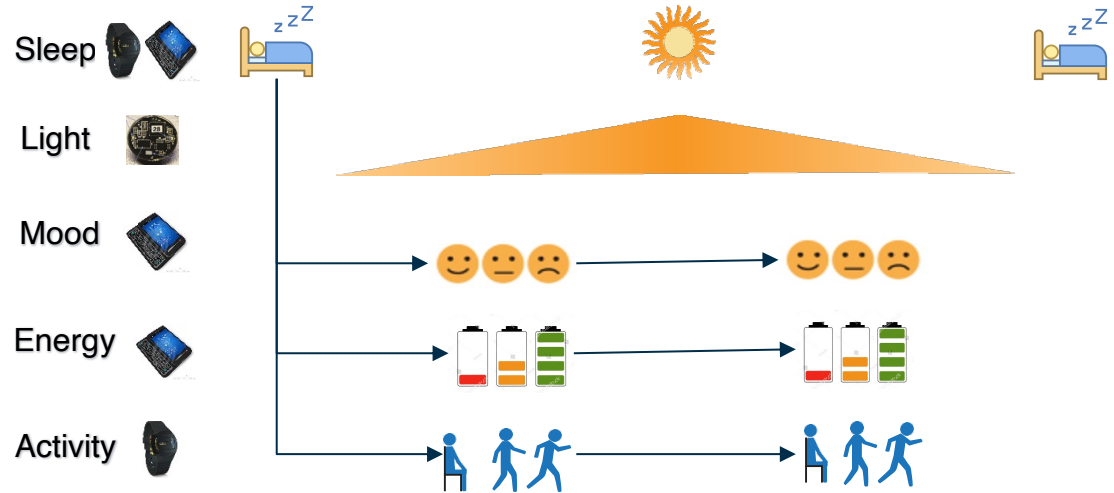


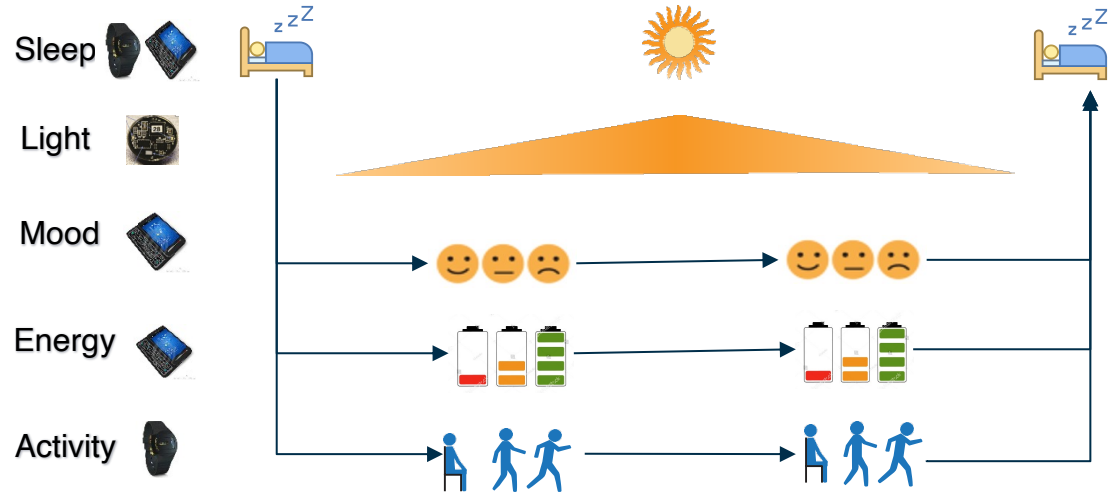
Energy

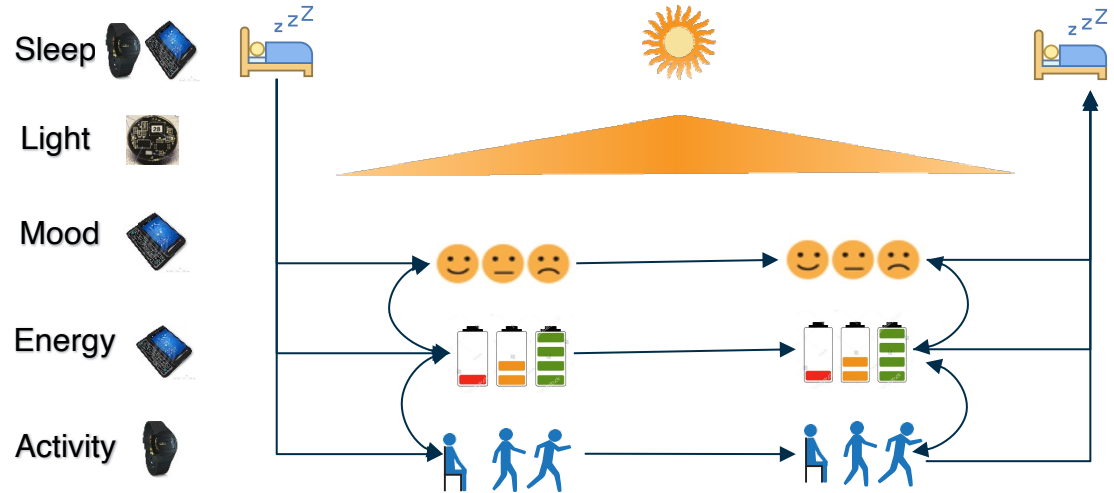


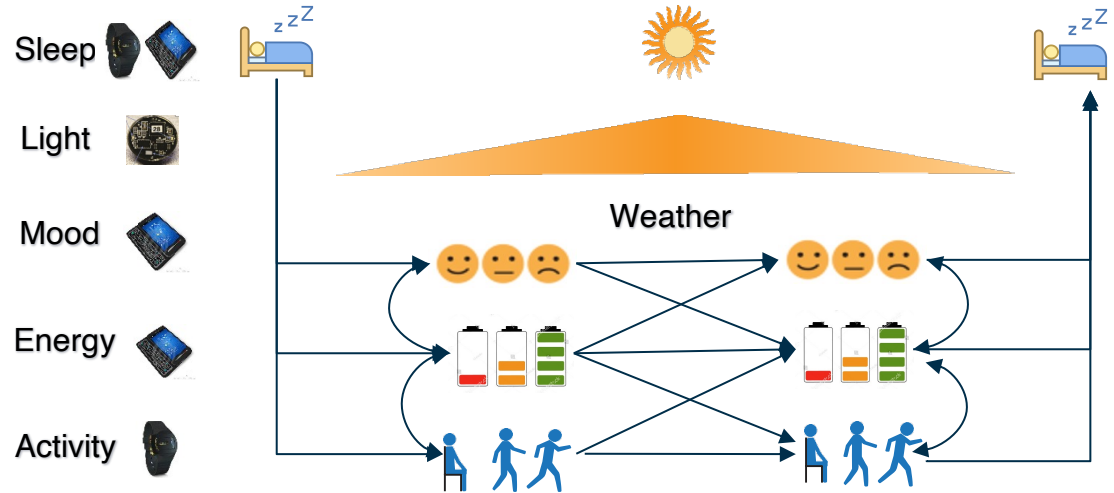


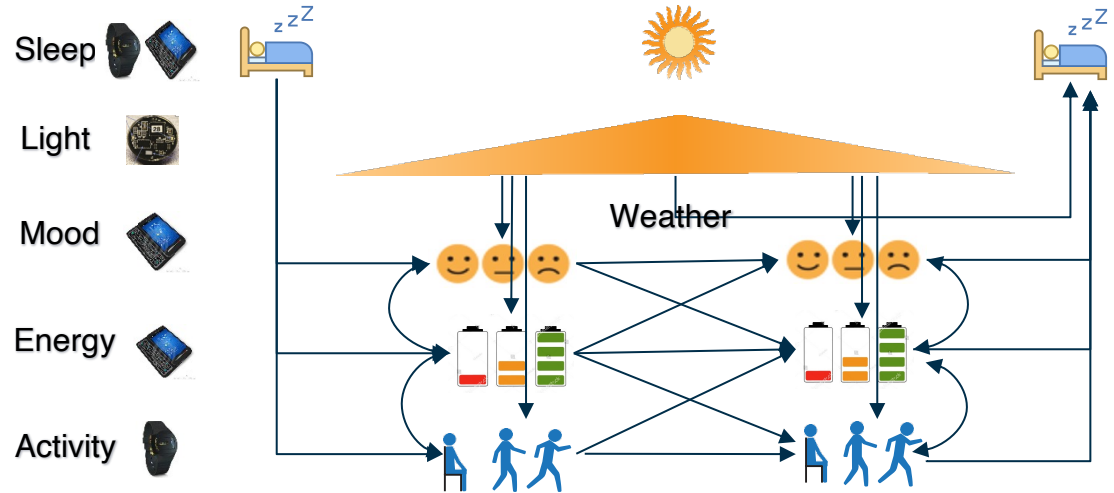


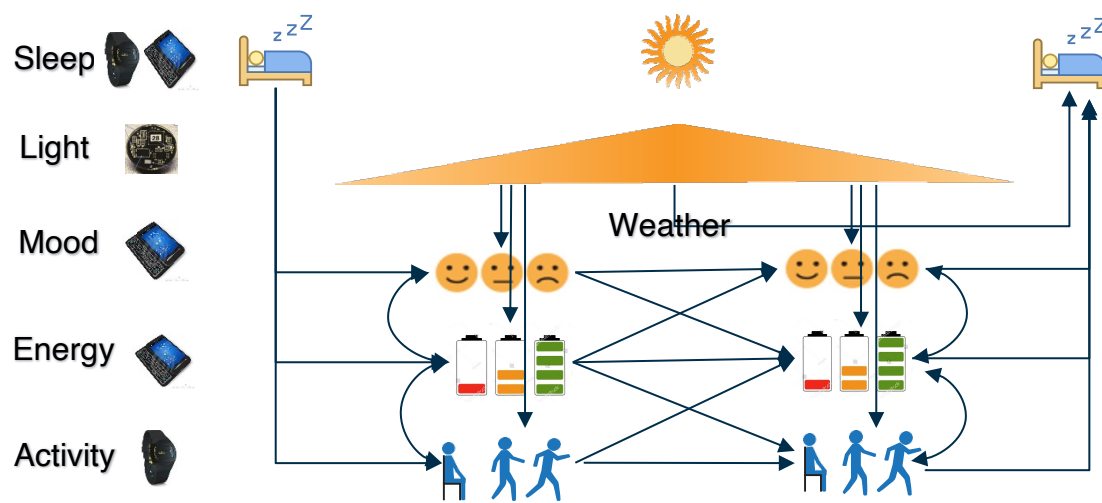












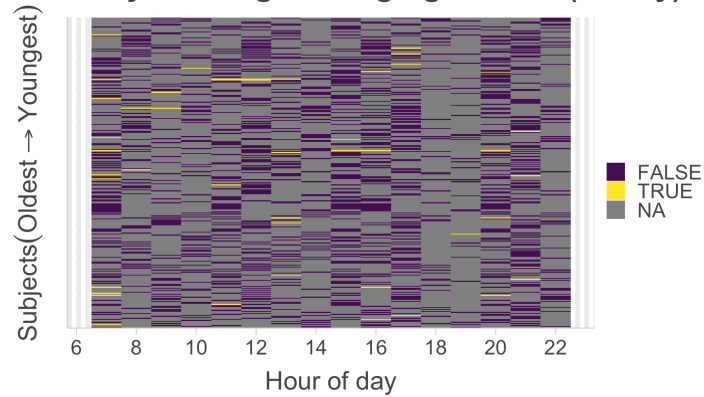
**Multivariate Principal Component Analysis for Mixed-type functional data to understand cross- and inter-correlations between processes.**

# Challenges: mHealth

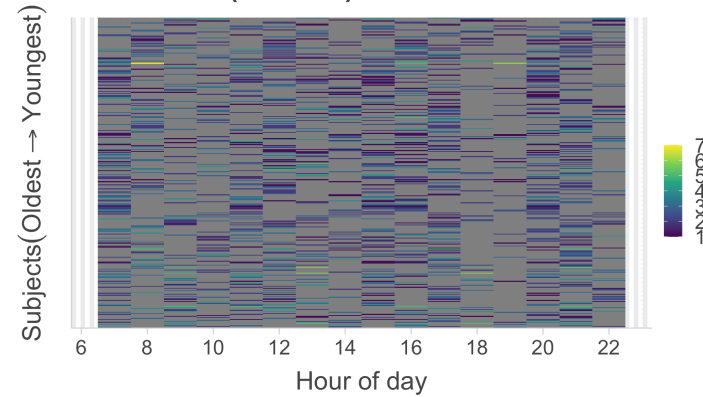


- Real-time self-reports of **mood**, **energy**, **stress**, **anxiety** (1-7), **headache** (0-1) recorded through smartphones.
- Objectively recorded **physical activity** and **sleep** through smartwatches.
- Intensive mixed-type longitudinal data.
- Mixed-type **sparse** data observed in **misaligned** time scales.
- Different measurement scales (**binary**, **ordinal**, **truncated**, **continuous**).
- Differences in subjective interpretation of **scales**.

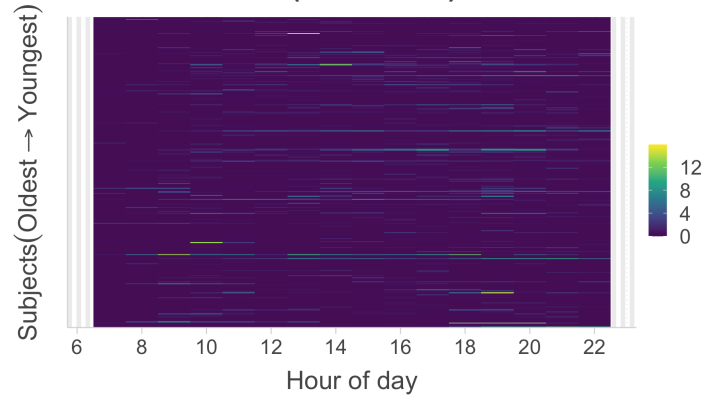
**Are you doing nothing right now? (Binary)**



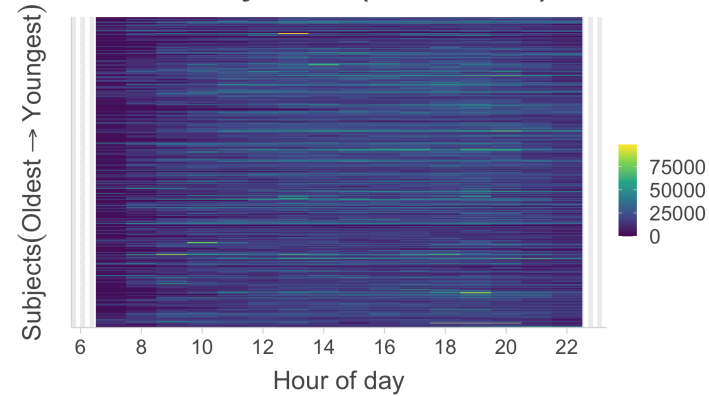
**Sad mood (Ordinal)**



**MVPA minutes (Truncated)**



**Total activity count (Continuous)**





# What we need

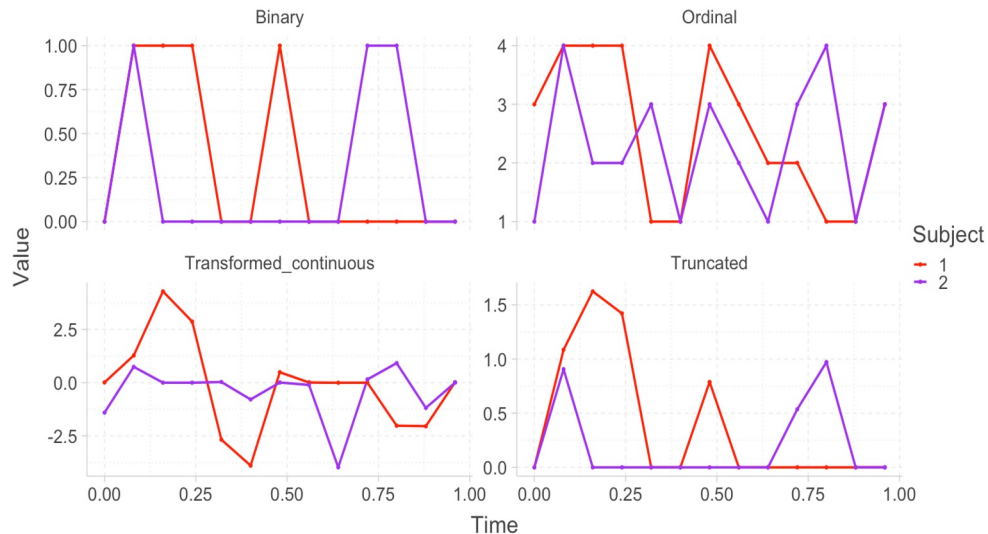
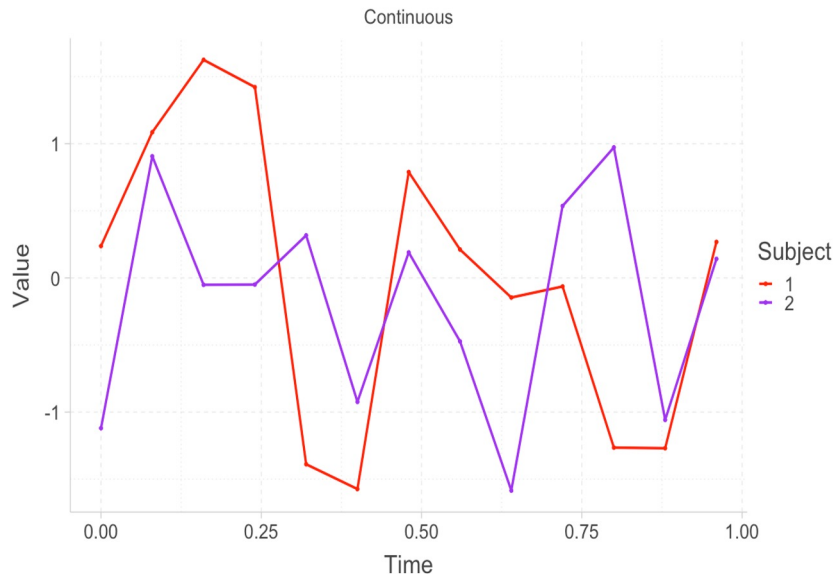
- Can we build a general modeling framework for joint modeling of **binary, ordinal, continuous and truncated type** functional data?
- Can we build **Multivariate Functional Principal Component Analysis** for such approaches?

# Generalized Latent Non-paranormal Process

Latent Gaussian Process



Observed data



# Latent Non-paranormal Distribution

**Definition 1.** (Nonparanormal distribution) A random vector  $W = (W_1, \dots, W_p)'$  follows a non-paranormal distribution denoted as  $W \sim NPN_p(0, \Sigma, f)$ , if there exists monotone transformation functions  $f = (f_1, \dots, f_p)$  such that  $L = f(W) = (f_1(W_1), \dots, f_p(W_p)) \sim N(0, \Sigma)$ , with  $\Sigma_{jj} = 1$  for  $1 \leq j \leq p$ . The constraints on diagonal elements of  $\Sigma$  are made to ensure the identifiability of the distribution as demonstrated in Liu et al.<sup>[22]</sup> and Fan et al.<sup>[26]</sup>.

# Generalized Latent Non-paranormal Process (X(t))

Let  $X_{ij}(t); t \in \mathcal{T}$  be a function for  $j$ th variable measured over a continuous variable  $t$  within subject  $i$  for  $i = 1, 2, \dots, I$  and  $j = 1, 2, \dots, J$ .

$$X_{ij}(t) = Z_{ij}(t)(\text{continuous scale}), \text{ or}$$

$$X_{ij}(t) = Z_{ij}(t)I(Z_{ij}(t) > \Delta(t)), (\text{truncated scale}), \text{ or}$$

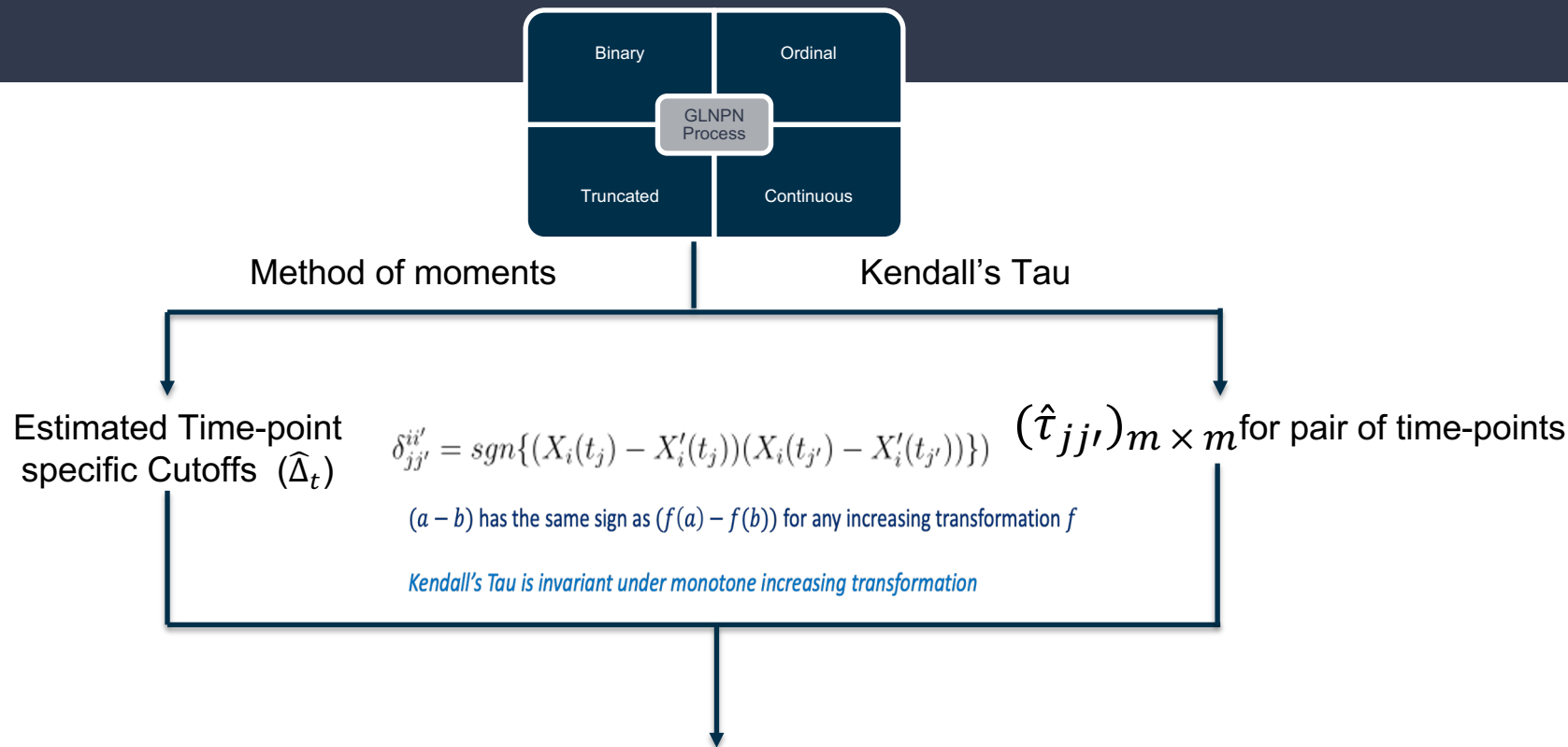
$$X_{ij}(t) = \sum_{k=0}^{l-1} kI(\Delta_{ijk}(t) \leq Z_{ij}(t) < \Delta_{ij(k+1)}(t)),$$

$$-\infty = \Delta_0(t) \leq \Delta_1(t) \leq \dots \leq \Delta_l(t) = \infty, (\text{ordinal scale}), \text{ or}$$

$$X(t) = I(Z_{ij}(t) > \Delta(t)), (\text{binary scale}).$$

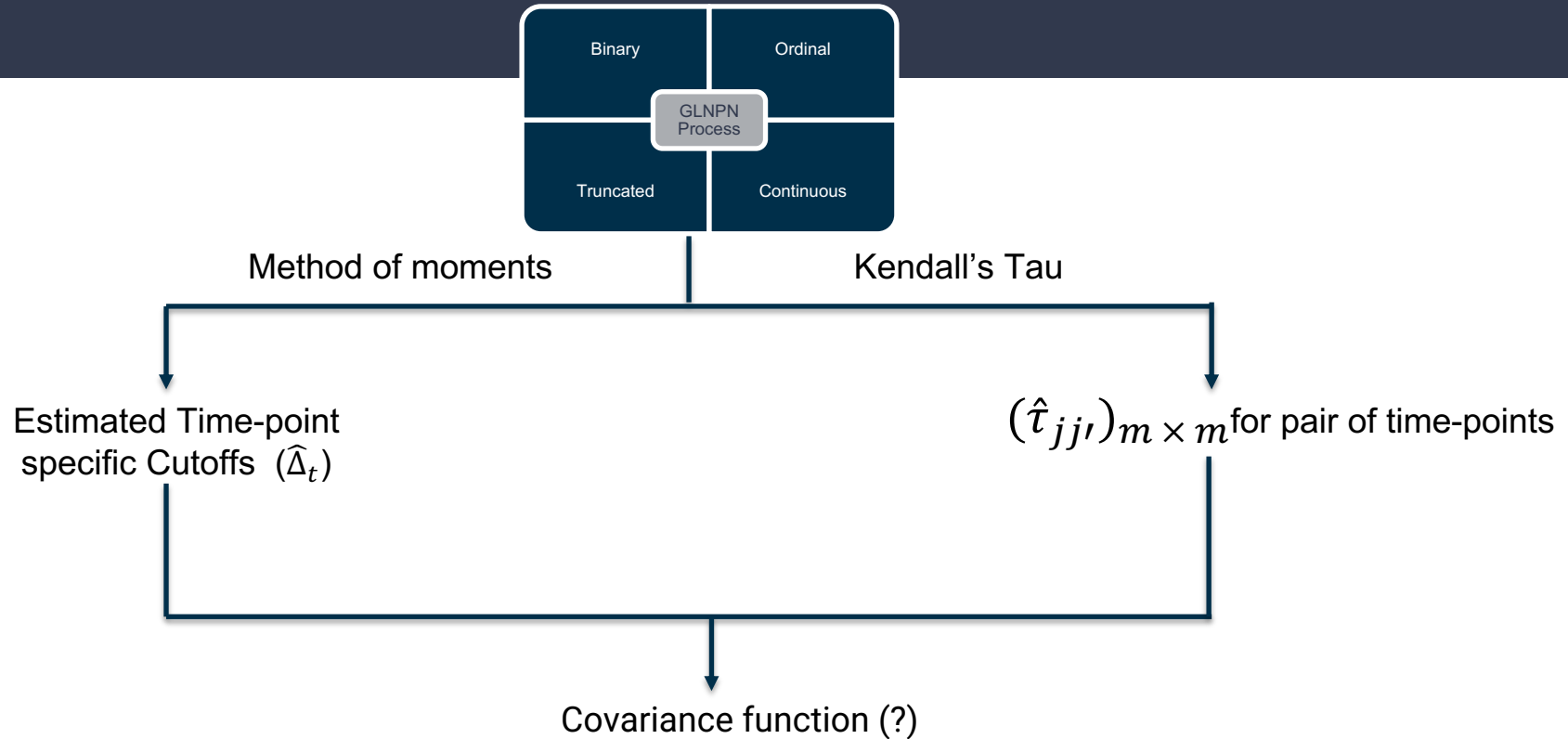
For multiple dependent outcomes  $Z(\mathcal{T}) = (Z_1(\mathcal{T}), \dots, Z_q(\mathcal{T}))^T$ , specification is  $(Z(\mathcal{T})) \sim NPN(0, C(\cdot, \cdot, f))$  where  $C = (C_{ij})$  is the  $J \times J$  multivariate covariance function.

# Univariate Marginal Covariance Estimation\*

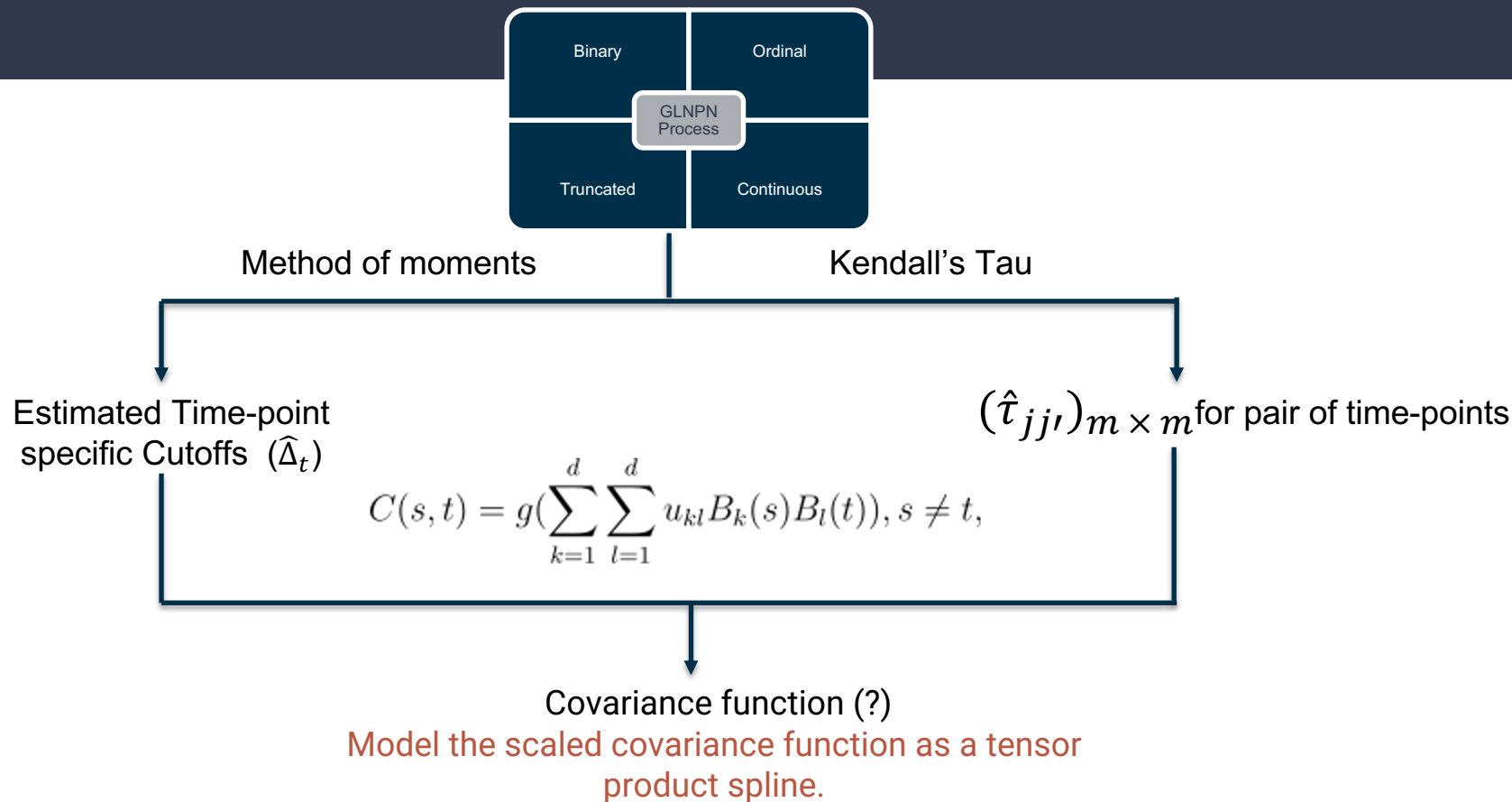


\* Functional Principal Component Analysis for Continuous non-Gaussian, Truncated, and Discrete Functional Data, Dey, Ghosal, Merikangas, Zupnik, 2024

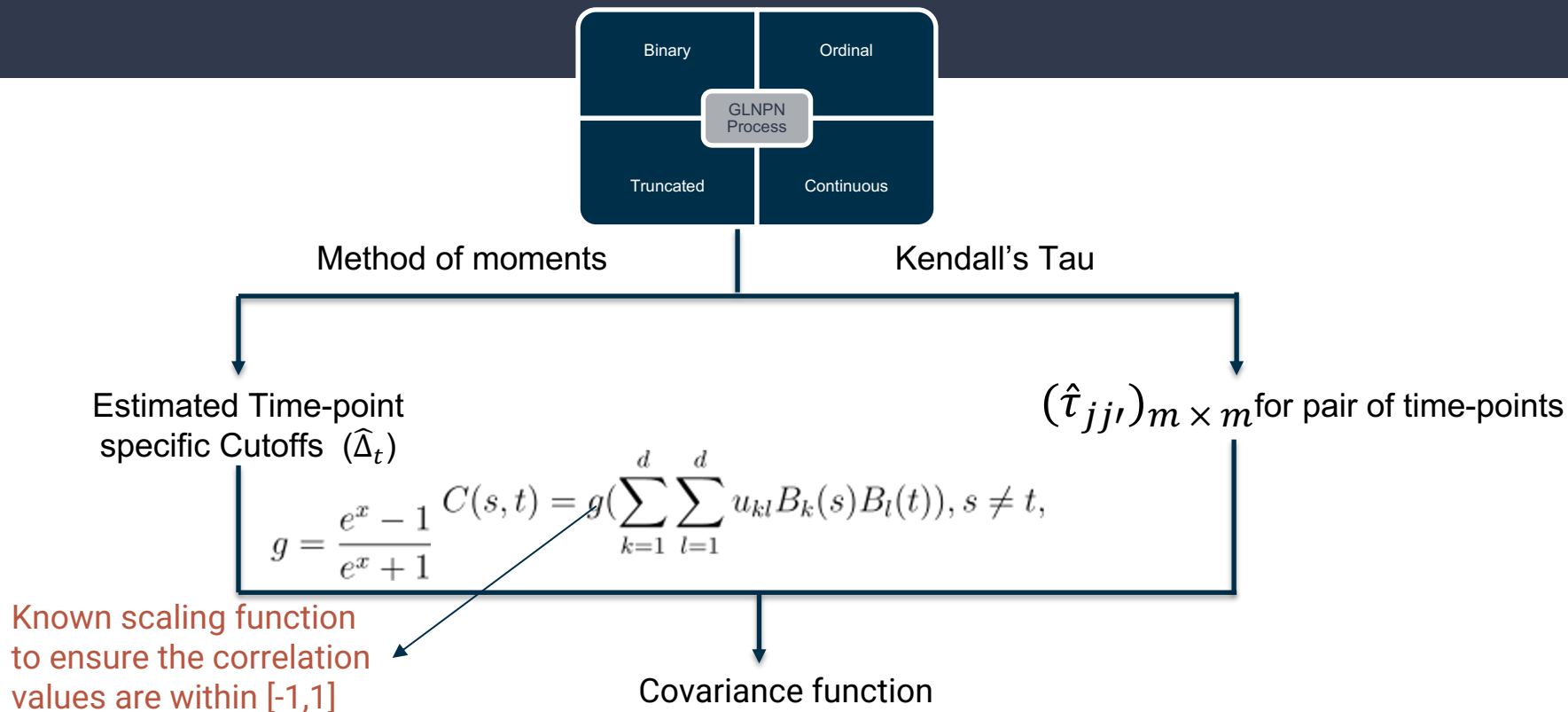
# Univariate Marginal Covariance Estimation



# Univariate Marginal Covariance Estimation

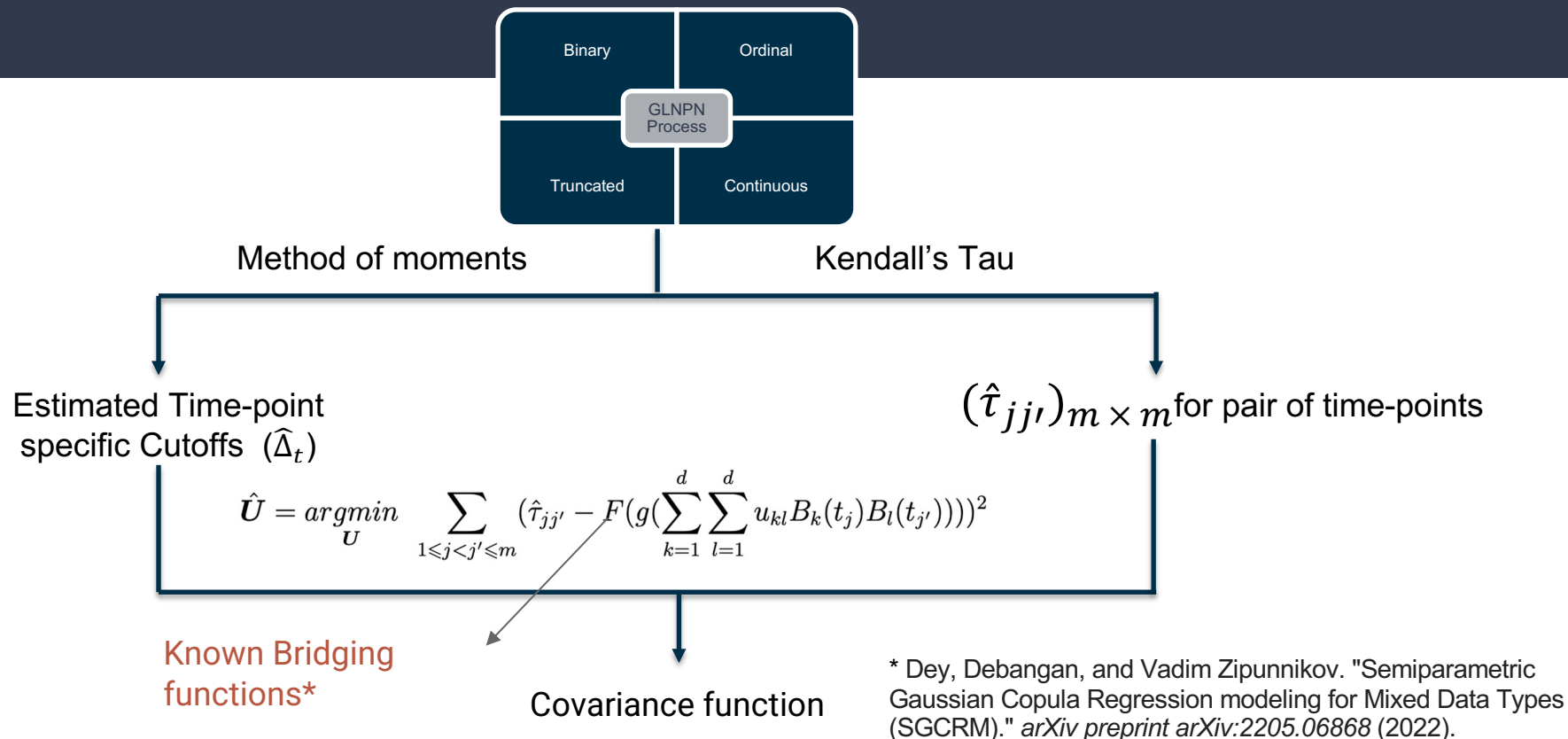


# Univariate Marginal Covariance Estimation

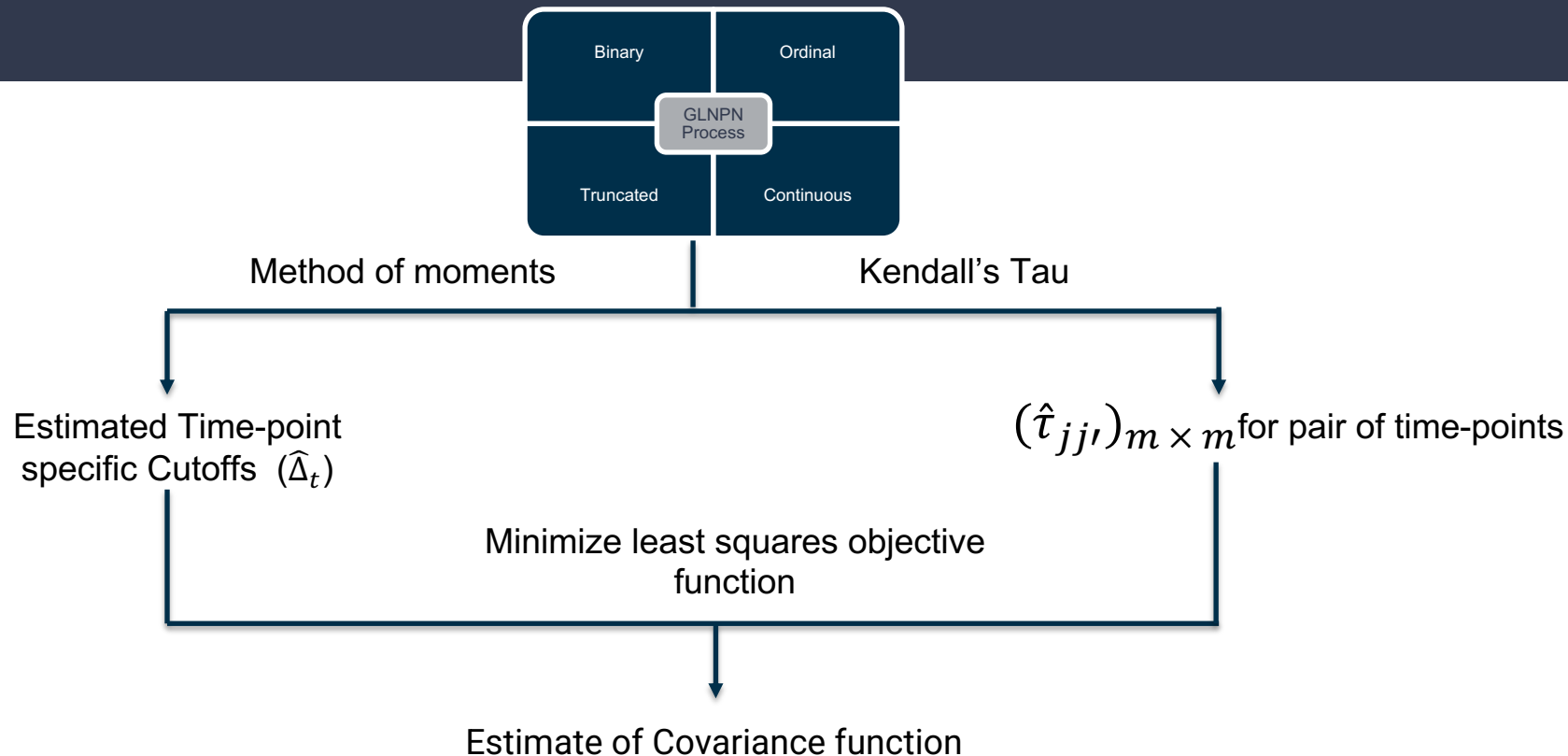




# Univariate Marginal Covariance Estimation



# Univariate Marginal Covariance Estimation



# Full Covariance Estimation

- Estimate marginal covariance matrices ( $\Sigma_{jj}$  blocks).
- Following similar steps as marginal estimation, estimate cross-covariance ( $\Sigma_{jk}$  blocks) (use cross-variable bridging function)
- Use nearPD\* function in R package Matrix to find the nearest correlation matrix made of the above blocks.
- Needs to estimate  $O(J^2)$  covariances. Could be slow.

\*Higham, Nick (2002) Computing the nearest correlation matrix - a problem from finance; IMA Journal of Numerical Analysis 22, 329--343.

# Full Covariance Estimation (Faster)

- Estimate marginal covariance matrices ( $\Sigma_{\{jj\}}$  blocks).
- We can use conditional expectation to predict **latent continuous variables** for each marginal process.
- Treat these predictions as **pseudo-observations (continuous)** to perform partially separable\* or full multivariate principal component analyses.
- Partial separable assumption can give us  $J$  vector principal component scores for desired number of levels and can be used to define graphical model between variables.

\*Partial separability and functional graphical models for multivariate Gaussian processes By J. Zapata, S. Y. Oh and A. Petersen, Biometrika, 2022

# Resulting methods

## 1. Function-on-function regression

$$E(V_1(s)|V_2(T)) = \Sigma_{12}(s, T)\Sigma_{22}(T, T)^{-1}(V_2(T))$$

In our case, we can regress mood ratings on concurrent and past physical activity from the joint covariance structure.

## 2. Multivariate functional principal component scores

# Advantages

- Can work with **sparse data** as we only need pairwise complete information for any pairs of time-points to calculate Kendall's Tau.
- We can use conditional expectation to predict **latent trajectory of a subject** at missing time-points.
- We develop approaches to calculate **latent Functional Principal Component Scores** which can be used in further predictive modeling.
- We provide **asymptotical confidence intervals** for regression coefficients and covariance values.

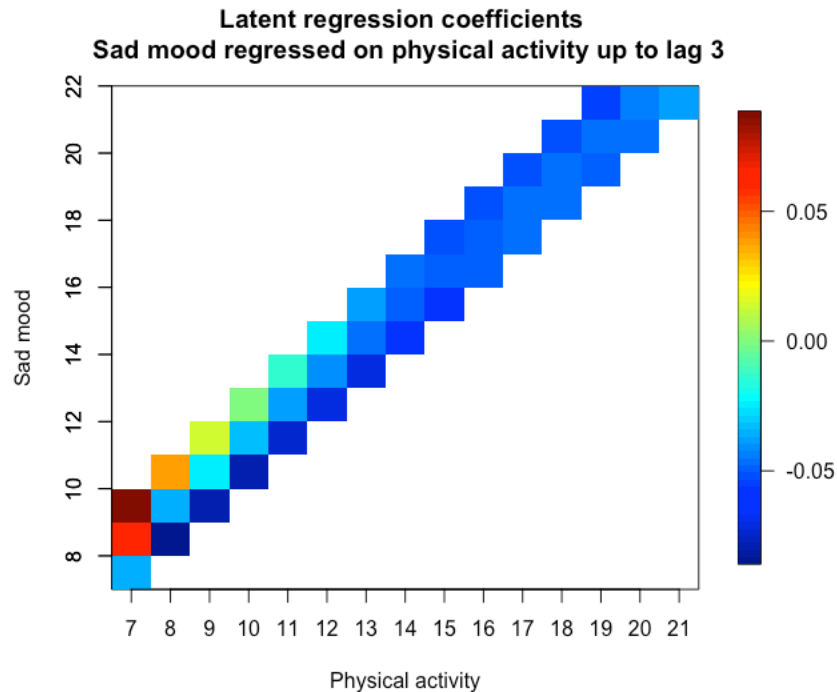
# Data analysis: NIMH Family Study of Affective Disorders

**Table 2**

*Descriptive statistics for the complete, male and female samples in the real data analysis. For continuous variable the mean and standard deviation is reported, for categorical variable the frequency in each group is mentioned. The P-values are from two-sample t-test and Chi-Square test of association with gender.*

Characteristic	Complete (n=497)	Male (n=195)	Female (n=302)	P value
	Mean(sd)	Mean(sd)	Mean(sd)	
Age	41.8 (19.5)	41.2 (21.7)	42.2(17.9)	0.56
Diagnosis: control ( <i>N</i> )	134	74	60	0.0001
Diagnosis: Anxiety ( <i>N</i> )	97	35	62	
Diagnosis: bipolar I ( <i>N</i> )	56	20	36	
Diagnosis: bipolar II ( <i>N</i> )	54	22	32	
Diagnosis: MDD ( <i>N</i> )	156	44	112	

# Mood on current and past physical activity





# Discussion

- Our approach provides a unified, general solution for modeling multivariate mixed-type functional data.
- Since we use rank-based estimators, our approach is robust.
- Next, we want to extend to handle multilevel structure of the data and take care of within-subject correlation.

Thank you.  
Questions?