
Idiographic Personality Gaussian Process for Psychological Assessment

Yehu Chen, Muchen Xi, Jacob Montgomery
Joshua Jackson, Roman Garnett

Washington University in St Louis

chenyehu,m.xi,j.jackson,jacob.montgomery,garnett@wustl.edu

Abstract

Developing taxonomies for psychological assessment is crucial for understanding long-term human behaviors. However, existing psychometric methods tend to be nomothetic, lacking individualization, or rely on static cross-sectional data that overlooks the dynamic nature of psychological processes. We introduce an idiographic personality Gaussian process (IPGP) framework for time-series survey data, by leveraging Gaussian process coregionalization to conceptualize individualized taxonomies and stochastic variational inference for computational scalability. Through an extensive simulation study against benchmark methods and an exploratory factor analysis study of life outcomes of personality replication, we demonstrate that IPGP can simultaneously improve estimation of idiographic taxonomies and prediction of missing responses. We also assess IPGP using our IRB-approved data with a forecasting and a leave-one-trait-out prediction task, illustrating how IPGP identifies unique taxonomies of personality that display potential in advancing individualized approaches to psychological diagnosis.

1 Introduction

Building standard taxonomies for psychological assessment is crucial to understand long-term behaviors through repeated quantitative surveys, for instance, emotional stability after medical treatment or development of academic ability during secondary education [Molenaar, 2004, Wang et al., 2013, Dumas et al., 2020]. However, existing taxonomies face several limitations. First of all, a common group of latent concepts is usually constructed for describing a single psychological trait yet with different causes. For instance, depression is constructed to explaining behavior of people suffering from numerous syndromes that differ in etiology, symptoms, and biological processes but get grouped together and called depression [Borsboom et al., 2003, Molenaar, 2004]. Second, conceptualization of taxonomies tend not to be individualized by assuming (1) the correlation between latent traits of interest and survey responses are invariant across individuals and (2) the actual structure of the underlying latent dimensions is fixed across individuals. Lastly, current established taxonomies are usually developed from cross-sectional data that are collected only once from respondents, and might overlook the dynamic nature of psychological processes.

To address these limitations, recent work have proposed an *idiographic* approach that builds completely distinct taxonomy for everyone [Borkenau and Ostendorf, 1998, Beck and Jackson, 2020, 2021], compared to the *nomothetic* approach where everyone is described by the same taxonomy. However, complete personalization may sacrifice general interpretability to clinicians as any possible population commonality is completely ignored. Another line of research focus on building dynamic psychometric models with time-series data via item response theory [Rijmen et al., 2003, Reise and Waller, 2009, Dumas et al., 2020], longitudinal structural equation [Little, 2013, Kim and Willson, 2014, Asparouhov et al., 2018], vectorized autoregression [Lu et al., 2018, Haslbeck et al., 2021] and

Gaussian process (GP) [Wang et al., 2005, Damianou et al., 2011, Dürichen et al., 2014, Duck-Mayr et al., 2020]. Yet all these models adopt the nomothetic approach by fixing the same taxonomic structure across individuals, and hence fail to identify any deviation from subgroups. Furthermore, there has been attempt to combine approaches for creating individualization while maintaining group commonality [Beltz et al., 2016], but prioritizes predicting responses to clinical survey batteries at the expense of studying the latent structures that are the focus of domain researchers.

In this work, we propose an idiographic personality Gaussian process (IPGP) framework for assessing dynamic psychological taxonomies from time-series survey data, and combine the nomothetic and idiographic approaches by deploying a common structure for explaining the typical circumstance and individual structures for permitting deviations into distinct forms. We leverage the Gaussian process coregionalization model to conceptualize responses of grouped survey batteries, adjusted to non-Gaussian ordinal data, and utilize IPGP for hypothesis testing of domain theories. Computationally, our framework also exploits the stochastic variational inference for latent factor estimation, contrasting with other GP measurement models relying on Gibbs sampling that may not scale efficiently to intensive longitudinal setups [Dürichen et al., 2014, Duck-Mayr et al., 2020].

To our knowledge, our work is the first multi-task Gaussian process latent variable model for dynamic idiographic assessment, compared to previous literature focusing on either static setup that ignores dynamics in latent processes [Borkenau and Ostendorf, 1998, Bonilla et al., 2007, Beck and Jackson, 2021] or single-task approach that disregards inter-battery correlation [Snelson and Ghahramani, 2005, Hensman et al., 2015]. Methodologically, our approach intersects Gaussian process latent variable model (GPLVM) [Lawrence, 2003], Gaussian process dynamic system (GPDM) [Damianou et al., 2011, Dürichen et al., 2014] and GP ordinal regression for likert-type survey data [Croasmun and Ostrom, 2011, Chu and Ghahramani, 2005]. Through an extensive simulation study against benchmark methods and an exploratory factor analysis study of life outcomes of personality replication, we demonstrate that IPGP can simultaneously improve estimation of idiographic taxonomies and prediction of missing responses. We also assess IPGP using our IRB-approved data with a forecasting and a leave-one-trait-out prediction task, illustrating how IPGP identifies unique taxonomies of personality that display potential in advancing individualized approaches to psychological diagnosis.

2 Backgrounds

We start by laying out the ordinal factor model for building standard taxonomy from survey experiments [Digman, 1997, Baglin, 2014]. We then briefly discuss several existing idiographic longitudinal models in psychological assessment, and review the model of Gaussian process.

Ordinal factor analysis. Consider the following scenario in survey experiments of $i \in \{1, \dots, N\}$ units repeatedly answering the same set of $j \in \{1, \dots, J\}$ batteries over $t \in \{1, \dots, T\}$ periods with ordinal observations $y_{ijt} \in \{1, \dots, C\}$ up to C levels. For example, the responses could be Likert-typed, ranging from “strongly disagree” to “strongly agree”. The latent factor model posits that the j th underlying latent variable $f_j^{(i)}(t)$ for unit i at time t are factored as $\mathbf{w}_j^T \mathbf{x}_i(t)$, where $\mathbf{x}_i(t) \in \mathbf{R}^K$ are unit-level latent factors and $\mathbf{w}_j \in \mathbf{R}^K$ are factor loadings. The $f_j^{(i)}(t)$ s are then mapped to ordinal responses via an ordered logit model: $p(y_{ijt} = c \mid f_j^{(i)}(t) = f) = \Phi(b_c - f) - \Phi(b_{c-1} - f)$ with threshold parameters $b_0 < \dots < b_C$. Usually b_0 and b_C are fixed to $-\infty$ and $+\infty$ such that the resulted categorical probability vector sums to 1, while b_1, \dots, b_{C-1} are allowed to move freely. Stacking $\mathbf{x}_i(t)$ s, \mathbf{w}_j s and y_{ijt} ’s into matrices \mathbf{x} , \mathbf{w} and tensor \mathbf{y} , the joint likelihood can be written as $\mathcal{L}(\mathbf{y} \mid \mathbf{x}, \mathbf{w}) = \prod_i \prod_j \prod_t p(y_{ijt} \mid \mathbf{x}_i(t), \mathbf{w}_j)$, while model identification is guaranteed by the general rule of factor models with additional orthogonality and normalization constraints [Bollen, 1989]. This factor model is also known as item response model [Samejima, 1969, Van der Linden and Hambleton, 1997], which estimates parameters via maximum likelihood, weighted least squares or EM algorithm [Bock and Aitkin, 1981, Forero et al., 2009, Li, 2016].

Idiographic longitudinal assessment. In psychological assessment, the idiographic approach emphasizes *intrapersonal* variation by requiring distinct loadings $\mathbf{w}_j^{(i)}$, while the nomothetic approach identifies general *interpersonal* variation assuming shared factor loadings \mathbf{w}_j s [Salvatore and Valsiner, 2010]. In terms of data collection, the idiographic approach usually surveys each individual multiple times ($n = 1$ and large T) for learning personalized taxonomy rather than many individuals at a

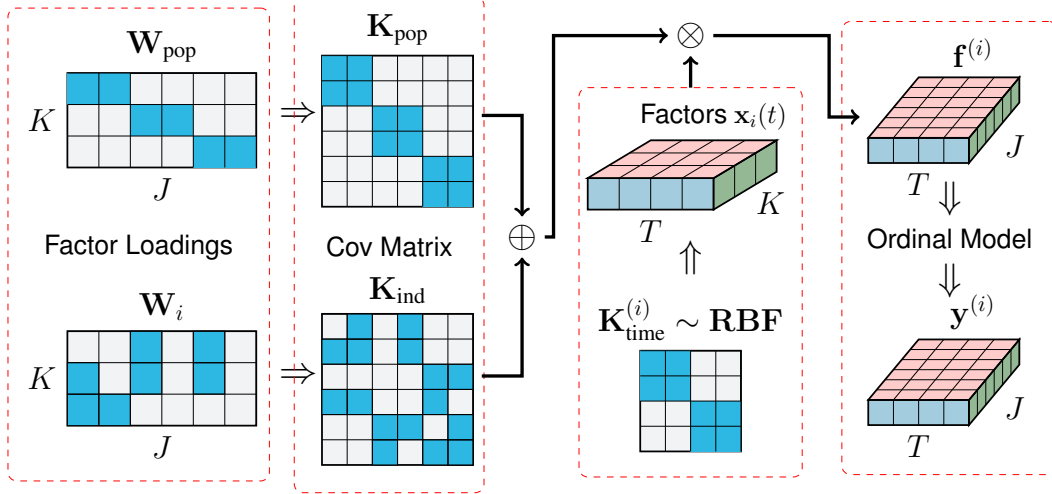


Figure 1: Proposed IPGP model for inferring latent factors and factor loadings from dynamic ordinal data. Input ordinal observations across channels are modeled as ordinal transformations of latent dynamic Gaussian processes with individualized RBF kernels and loading matrices.

single shot (large n and $T = 1$). To extract individualized dynamics from time-series data, recent psychometrics models have utilized longitudinal structural equations by explicitly specifying any intrapersonal and temporal relation yet sensitive to model mis-specification from domain theory [Little, 2013, Asparouhov et al., 2018]. Meanwhile, variants of hierarchical vector autoregression may automatically learn individual diffusion, but usually lack the ordered logit component as built in response space rather than latent space [Lu et al., 2018, Haslbeck et al., 2021].

Gaussian process. A Gaussian process (GP) can be used to define a distribution over f such that the evaluation of f at arbitrary subset of \mathcal{X} is a joint multivariate Gaussian [Rasmussen and Williams, 2005]. To determine its mean and covariance, a $\mathcal{GP}(\mu, K)$ is specified with a mean function $\mu : \mathcal{X} \rightarrow \mathbf{R}$ and a positive-definite kernel function $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbf{R}$. The most common kernel is the squared exponential (RBF) kernel $K(\mathbf{x}_1, \mathbf{x}_2) = \exp(-\frac{1}{2}\mathbf{x}_1^T \mathbf{P} \mathbf{x}_2)$ with precision matrix $\mathbf{P} = \text{diag}(1/\ell_1^2, \dots, 1/\ell_d^2)$ and $d = \text{card}(\mathcal{X})$. Posterior of a GP is usually analytical for Gaussian likelihood, but needs to be approximated in modeling latent variables. We discuss the variational approximation in Sec. (3).

3 Methodology

We propose an idiographic personality Gaussian process (IPGP) framework for assessing individualized dynamic psychological taxonomies from time-series survey data. Instead of joint estimation of latent factors and their loadings that cannot guarantee rotational and scaling invariance, we marginalize out the latent variables and focus on learning taxonomies of loadings. The overall architecture of IPGP is illustrated in Figure (1), where input ordinal observations across channels are modeled as ordinal transformations of latent dynamic GP with individualized RBF kernels and loading matrices.

3.1 Multi-task learning

Typically in psychological assessment, survey questions are meticulously grouped such that each group gauges a particular facet of personality. Hence, we conceptualize the assessment of psychological traits as a multi-task learning problem, where each question represents a distinct task but can be correlated with other tasks. A multi-task GP is an extension of the single-task GP but for vector-valued functions [Bonilla et al., 2007]. To motivate the multi-task framework, first consider the two-task scenario with two $T \times 1$ vector $\mathbf{f}_1^{(i)}$ and $\mathbf{f}_2^{(i)}$ denoting the latent temporal processes of unit i for question $j = 1, 2$. To fix the scale of latent factors, a time-level Gaussian process prior is placed on $\mathbf{x}_i(t) \sim \mathcal{GP}(\mathbf{0}, \mathbf{K}_{\text{time}}^{(i)})$. Hence, by exploiting affine property of Gaussian, the induced joint

distribution of vectorized $[\mathbf{f}_1^{(i)}, \mathbf{f}_2^{(i)}]^T$ can be written as:

$$\begin{bmatrix} \mathbf{f}_1^{(i)} \\ \mathbf{f}_2^{(i)} \end{bmatrix} \sim \mathcal{GP} \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{w}_1^T \mathbf{w}_1 \mathbf{K}_{\text{time}}^{(i)} & \mathbf{w}_1^T \mathbf{w}_2 \mathbf{K}_{\text{time}}^{(i)} \\ \mathbf{w}_2^T \mathbf{w}_1 \mathbf{K}_{\text{time}}^{(i)} & \mathbf{w}_2^T \mathbf{w}_2 \mathbf{K}_{\text{time}}^{(i)} \end{bmatrix} \right) \quad (1)$$

whose covariance of shape $2T \times 2T$ contains four block matrices $\mathbf{K}_{\text{time}}^{(i)}$ scaled by different $\mathbf{w}_j^T \mathbf{w}_{j'}$ ($j, j' \in \{1, 2\}$). Specifically, $\mathbf{w}_1^T \mathbf{w}_2$ controls the inter-task covariance between these two tasks and $\mathbf{w}_j^T \mathbf{w}_j$ s ($j \in \{1, 2\}$) control their intra-task variance. This multi-task structure is also known as the linear model of coregionalization (LMC) [Alvarez et al., 2012], where the factor structure can be recovered from the relations $\mathbf{f}_j^{(i)} = \mathbf{w}_j^T \mathbf{x}_i(t)$ of linear combinations. Compared to vector autoregression, our GP dynamic system approach can better handle intrapersonal temporal structures. Now let $\mathbf{f}^{(i)} = [\mathbf{f}_1^{(i)}, \dots, \mathbf{f}_J^{(i)}]^T$ represents the flattened $JT \times 1$ vector consisting of all J tasks. We write $\mathbf{f}^{(i)}$ in a formal multi-task GP notation using Kronecker product \otimes :

$$p(\mathbf{f}^{(i)}) \sim \mathcal{GP}(\mathbf{0}, \mathbf{K}_{\text{task}}^{(i)} \otimes \mathbf{K}_{\text{time}}^{(i)}) \quad (2)$$

$$\mathbf{K}_{\text{task}}^{(i)} = \mathbf{W}_{\text{pop}}^T \mathbf{W}_{\text{pop}} + \mathbf{w}_i^T \mathbf{w}_i + \text{diag}(\mathbf{v}) \quad (3)$$

where $\mathbf{K}_{\text{task}}^{(i)}$ denotes the unit-individualized task kernel, consisting of the self inner products of population loading $\mathbf{W}_{\text{pop}} = [\mathbf{w}_1, \dots, \mathbf{w}_J]$ for explaining the interpersonal commonality and $J \times 1$ idiographic loading \mathbf{w}_i for intrapersonal deviations, as well as a task-dependent noise component $\text{diag}(\mathbf{v}) = \text{diag}([\sigma_1^2, \dots, \sigma_J^2])$. The Kronecker product \otimes then multiplies each entry in the $J \times J$ task covariance with $\mathbf{K}_{\text{time}}^{(i)}$, and returns the stacked $JT \times JT$ covariance for $\mathbf{f}^{(i)}$. Here we use the common RBF kernel $\mathbf{K}_{\text{time}}^{(i)}(t, t') = \exp(-(t - t')^2 / \ell_i^2)$ to account for dynamic changes in the latent attributes, whose bandwidth is determined by the unit-specific length scale ℓ_i , but any other kernel can substitute RBF as practitioners see fit.

3.2 Variational learning

Due to the non-Gaussian essence of ordinal likelihood, we adopt the stochastic variational inference technique (SVI) with inducing points introduced in [Hensman et al., 2015]. Dropping superscript for demonstration, SVI utilizes a variational distribution $q(\mathbf{u}) = \mathcal{N}(\mu_{\mathbf{u}}, \Sigma_{\mathbf{u}})$ on $m \ll n$ inducing variables \mathbf{u} to approximate $p(\mathbf{f} | \mathbf{y})$ using the conditional $p(\mathbf{f} | \mathbf{u})$. Hence, the conditional log likelihood $\log p(\mathbf{y} | \mathbf{u})$ can be lower bounded by the expected log likelihood w.r.t. $p(\mathbf{f} | \mathbf{u})$, after exploiting the non-negativity of Kullback–Leibler (KL) divergence between $p(\mathbf{f} | \mathbf{u})$ and $p(\mathbf{f} | \mathbf{y})$:

$$\log p(\mathbf{y} | \mathbf{u}) \geq \mathbb{E}_{p(\mathbf{f} | \mathbf{u})} \log p(\mathbf{y} | \mathbf{f}) \quad (4)$$

Furthermore, a lower bound on model evidence (ELBO) can be obtained by combining Eq. (4) and an inequality derived by another KL divergence $\text{KL}[q(\mathbf{u}) \| p(\mathbf{u} | \mathbf{y})] \geq 0$ (see Appendix A for details):

$$\log p(\mathbf{y}) \geq \mathbb{E}_{q(\mathbf{u})} [\log p(\mathbf{y} | \mathbf{u})] - \text{KL}[q(\mathbf{u}) \| p(\mathbf{u})] \quad (5)$$

$$\geq \mathbb{E}_{q(\mathbf{f})} [\log p(\mathbf{y} | \mathbf{f})] - \text{KL}[q(\mathbf{u}) \| p(\mathbf{u})] \quad (6)$$

where the KL divergence $\text{KL}[q(\mathbf{u}) \| p(\mathbf{u})]$ between the variational $q(\mathbf{u})$ and prior $p(\mathbf{u})$ can be computed in closed form as both distributions are Gaussians. The expectation of log likelihood $\log p(\mathbf{y} | \mathbf{f})$ under the marginal distribution $q(\mathbf{f}) = \int p(\mathbf{f} | \mathbf{u}) q(\mathbf{u}) d\mathbf{u}$ is intractable but can be numerically approximated using Gauss-Hermite quadrature method. The variational parameters $\mu_{\mathbf{u}}$ and $\Sigma_{\mathbf{u}}$, individualized loadings \mathbf{w}_i and $\text{diag}(\mathbf{v})$ as well as likelihood parameters $\{b_c\}$ s are then optimized to maximize this lower bound. Finally, the predictive likelihood of new $p(\mathbf{y}^*) = \int p(\mathbf{y}^* | \mathbf{f}^*) p(\mathbf{f}^* | \mathbf{u}) q^*(\mathbf{u}) d\mathbf{u}$ is obtained by marginalizing out the optimized $q^*(\mathbf{u})$.

3.3 Theory testing

Our IPGP framework also naturally facilitates downstream tasks such as domain theory testing between models with and without shared or idiographic components. We adopt Bayes factor, the posterior $p(\mathcal{M}_i | \mathbf{y}) = \frac{p(\mathbf{y} | \mathcal{M}_i) p(\mathcal{M}_i)}{\sum_i p(\mathbf{y} | \mathcal{M}_i) p(\mathcal{M}_i)}$ over a pool of models $\{\mathcal{M}_i\}$ conditioning on observation \mathbf{y} with prior weights $p(\mathcal{M}_i)$, as the hypothesis test on whether the latent structures for each individual

are indeed distinct or are simply explainable by interpersonal commonality. Specifically, we refer the multi-task model in Eq. (3) as the *idiographic* model, and compare it with an *nomothetic* model without unit-specific components: $\mathbf{K}_{\text{task}}^{\text{pop}} = \mathbf{W}_{\text{pop}} \mathbf{W}_{\text{pop}}^T + \text{diag}(\mathbf{v})$.

Note that compared to this baseline nomothetic model, our proposed idiographic model in Eq. (3) introduces extra unit-level Jn loading parameters that enlarges the optimization space of hyperparameter. Hence, we propose to first learn the interpersonal loading matrix \mathbf{W}_{pop} using the standard cross-sectional data from a nomothetic model that focuses on learning of population taxonomy, and then use the estimated \mathbf{W}_{pop} as informative prior in the full model. We will show empirically in Sec. (4) that with this stronger prior IPGP achieves more precise estimation of individual taxonomies.

4 Experiments

We now evaluate IPGP in learning idiographic latent taxonomies and predicting actual responses against baseline methods from both psychometrics and Gaussian process literature. We then conduct an exploratory factor analysis study of life outcomes of personality replication, validating the popular Big Five personality theory standard using cross-sectional data [McCrae and John, 1992]. The inferred shared taxonomy structure is further incorporated as the informative prior for the population kernel in our case study, where we collected IRB-approved longitudinal survey data. We also highlight the predictive ability of IPGP through a forecasting and a leave-one-trait-out cross validation task, and illustrate how IPGP identifies unique taxonomies of personality that might advance individualized approaches to psychological diagnosis and inspire new theory.

4.1 Simulation and ablation

Setup. Our simulation considers longitudinal data of $n = 10$ units over $T = 30$ periods. We assume latent traits of each unit i has dimension $K = 5$, and each dimension latent vector is generated independently from a GP $\mathbf{x}_i^{(k)}(t) \sim \mathcal{GP}(\mathbf{0}, \mathbf{K}_{\text{time}}^{(i)})$ with unit-specific length scale uniformly randomly picked from $\ell_{\text{time}}^{(i)} \in [10, 20, 30]$. We split $m = 20$ batteries into K subsets of size $m/K = 4$, such that each subset dominates one dimensional in the latent traits. Specifically, we set high value of 3 in the population factor loading matrix \mathbf{W}_{pop} for entries corresponding to the k th subset for dimension k , and low values drawn from $\text{Unif}[-1, 1]$ otherwise. We also set each unit-specific loading \mathbf{w}_i from $\text{Unif}[-1, 1]$. To allow sparsity and reverse coding, we randomly drop half of loadings and flip signs of the remaining half. Finally, we generate the y_{ijts} according to the ordered logit model with $C = 5$ levels, and apply 80%/20% splitting for training and testing.

Table 1: Comparison of averaged accuracy, log lik and correlation matrix distance between IPGP and baselines and ablated models in the simulated study. The full IPGP model (indicated in bold) significantly outperforms all ablated and baseline methods in both estimated correlation matrix and either in-sample or out-of-sample prediction in paired-t tests. Results from ablations imply that IPGP succeeds in predicting the correct labels due to its idiographic components and proper likelihood, and a well-informed population kernel is crucial in recovering the factor loadings. “—” indicates baseline software that cannot handle missing values.

	TRAIN ACC \uparrow	TRAIN LL \uparrow	TEST ACC \uparrow	TEST LL \uparrow	CMD \downarrow
GRM	0.261 \pm 0.005	−3.556 \pm 0.092	0.261 \pm 0.006	−3.578 \pm 0.098	0.657 \pm 0.021
GPCM	0.562 \pm 0.017	−2.067 \pm 0.182	0.495 \pm 0.012	−2.409 \pm 0.143	0.545 \pm 0.016
SRM	0.286 \pm 0.006	−7.408 \pm 0.063	0.289 \pm 0.008	−7.341 \pm 0.084	0.300 \pm 0.024
GPDM	0.687 \pm 0.010	−4.358 \pm 0.028	0.667 \pm 0.010	−4.377 \pm 0.029	0.262 \pm 0.016
DSEM	0.539 \pm 0.021	−0.961 \pm 0.015	—	—	0.256 \pm 0.011
TVAR	0.554 \pm 0.018	−1.168 \pm 0.014	—	—	0.987 \pm 0.013
IPGP-NOM	0.807 \pm 0.007	−0.535 \pm 0.015	0.790 \pm 0.008	−0.555 \pm 0.017	0.257 \pm 0.009
IPGP-IND	0.932 \pm 0.003	−0.243 \pm 0.008	0.916 \pm 0.004	−0.267 \pm 0.009	0.530 \pm 0.005
IPGP-LOW	0.897 \pm 0.004	−0.313 \pm 0.010	0.884 \pm 0.005	−0.334 \pm 0.011	0.397 \pm 0.007
IPGP-NP	0.898 \pm 0.003	−0.318 \pm 0.009	0.883 \pm 0.005	−0.342 \pm 0.011	0.467 \pm 0.010
IPGP	0.957 \pm 0.002	−0.159 \pm 0.005	0.942 \pm 0.002	−0.184 \pm 0.006	0.128 \pm 0.006

Metrics and baselines. We consider two sets of metrics for evaluation: (1) the in-sample and out-of-sample predictive accuracy (ACC) and log likelihood (LL) of the actual responses, (2) the correlation matrix distance (CMD) between the estimated factor loading matrix and the true ones, which is defined for two covariance matrices $\mathbf{R}_1, \mathbf{R}_2$ as $d(\mathbf{R}_1, \mathbf{R}_2) = 1 - \frac{\text{tr}(\mathbf{R}_1 \mathbf{R}_2)}{\|\mathbf{R}_1\|_f \|\mathbf{R}_2\|_f}$ [Herdin et al., 2005] with l_2 Frobenius norm. Note that CMD becomes zero if $\mathbf{R}_1, \mathbf{R}_2$ are equal up to a scaling factor, and one if they are orthogonal after flattening. We compare IPGP to (1) various latent variable models for ordinal responses, including the graded response model (GRM) [Samejima, 1969], the generalized partial credit model (GPCM) [Muraki, 1992] and the sequential response model (SRM) [Tutz, 1990], (2) Gaussian process dynamic model (GPDM) [Damianou et al., 2011, Dürichen et al., 2014] where the continuous predictions are rounded to the nearest ordinal level, (3) dynamic structural equation model (DSEM) [Asparouhov et al., 2018, McNeish et al., 2023] with trait-dependent latent variables and (4) time-varying vector autoregression (TVAR) with regularized kernel smoothing [Haslbeck et al., 2021]. We also compare IPGP with several ablated models: (1) IPGP-NOM without the idiographic kernel, (2) IPGP-IND without the population kernel, (3) IPGP-LOW with lower-rank factors of 2 than actual rank of 5 in the synthetic setup and (4) IPGP-NP where the population kernel is learned from scratch rather than fixed to the informative prior. Note that \mathbf{W}_{pop} in the full IPGP model is fixed as learned from IPGP-NOM.

Results. We use 100 inducing points and ADAM optimizer of learning rate 0.05 to optimize ELBO for 10 epoches with batch size of 256. We repeat our simulation with 25 different random seeds using 300 Intel Xeon 268 CPUs. Table 1 shows comparison of averaged predictive accuracy, log likelihood and correlation matrix distance between IPGP and baselines and ablated models in the simulated study. Our IPGP model (indicated in bold) significantly outperforms all ablated models and baseline methods in estimated correlation matrix, predictive accuracy and log likelihood of both training and testing sets in paired-t tests. We found that IPGP succeeds in predicting the correct labels due to its idiographic components and proper likelihood, since IPGP-NOM and IPGP-GL are two of the worst ablations for all prediction metrics. In addition, IPGP-IND and IPGP-NP have the worst correlation matrix estimation, implying that a well-informed population kernel is crucial in recovering the underlying factor structures.

4.2 Exploratory factor analysis

We then validate the popular Big Five personality theory using standard cross-sectional data via an exploratory factor analysis, where a range of factors are tested and then determined according to model evidence rather than being fixed. We utilize the life outcomes of personality replication (LOOPR) data (see [Soto, 2019] for a full description of LOOPR), which is collected from 5,347 unique participants on the Big Five Inventory [John et al., 1999] consisting of 60 battery questions. Our validation considers a range of latent trait dimension counts from $K = 1, \dots, 5$. For each dimension count, we first apply principal component analysis (PCA) directly on the correlation matrix of the cross-sectional observations to learn a vanilla population factor loading matrix. We then initialize \mathbf{W}_{pop} in our model with this vanilla loading matrix, and optimize the loading matrix jointly with the variational parameters. Note that $T = 1$ in LOOPR, so we drop the idiographic components.

Table 2: In-sample accuracy and averaged log lik of our method and baselines for various K in LOOPR. Best model for each K is indicated in bold and the best model across different K s is further indicated in italic.

MODEL	ACC \uparrow					LL / N \uparrow				
	$K = 1$	$K = 2$	$K = 3$	$K = 4$	$K = 5$	$K = 1$	$K = 2$	$K = 3$	$K = 4$	$K = 5$
PCA	0.106	0.099	0.123	0.217	0.192	-1.957	-1.990	-2.009	-2.036	-2.051
GRM	0.238	0.107	0.178	0.113	0.146	-1.838	-1.832	-1.814	-1.838	-1.841
GPCM	0.213	0.156	0.186	0.159	0.163	-1.754	-1.761	-1.764	-1.750	-1.756
SRM	0.243	0.134	0.179	0.125	0.155	-1.784	-1.784	-1.783	-1.780	-1.767
GPDM	0.268	0.272	0.266	0.268	0.263	-2.155	-2.158	-2.158	-2.159	-2.158
DSEM	0.188	0.114	0.110	0.105	0.104	-1.997	-1.960	-1.908	-1.845	-1.775
IPGP	0.322	0.319	0.323	0.318	0.318	-1.478	-1.477	-1.477	-1.477	-1.476

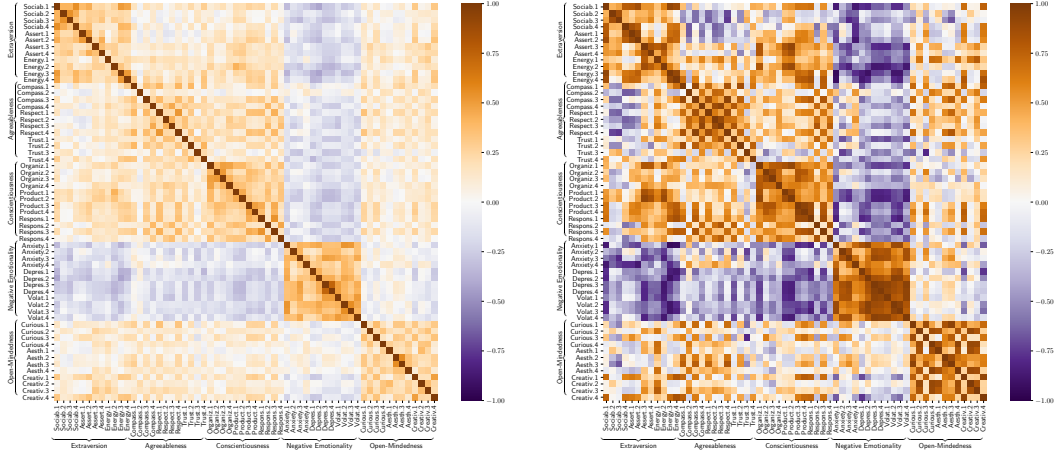


Figure 2: Illustration of raw correlation matrix (left) and our estimated Big Five loading matrix (right). Both correlation matrices display a *block* pattern, where estimated interpersonal variation show strong correlation between questions within the same factor of the Big Five personalities and weak correlation across different factors. Besides, questions corresponding negative emotionality show minor negative correlation with those corresponding to extraversion and conscientiousness, suggesting trait-by-trait interaction effects.

Validation of Big Five. Table 2 shows the predictive accuracy and averaged log likelihood of our method and baseline methods (excluding TVAR for lacking low-rank assumption) for various K in LOOPR. Best model for each K is indicated in bold numbers and the best model across different K s is further indicated in italic numbers. Despite having slightly worse in-sample predictive accuracy than factor 3 model, IPGP with factor 5 has significant higher model evidence than all the other models, with the second best model is $\exp(-79)$ more unlikely indicated by Bayes factors. Therefore, our results indicate that when psychological measurements are estimated from standard cross-sectional data, Big Five personality is necessary for learning interpersonal variation.

Estimated interpersonal variation. We also show the raw correlation and our estimated Big Five correlation in Figure 2. Both correlation matrices display a *block* pattern, where estimated interpersonal variation show strong correlation between questions within the same factor of the Big Five and weak correlation across different factors. Besides, questions corresponding negative emotionality show minor negative correlation with those corresponding to extraversion and conscientiousness, suggesting appropriate trait-by-trait interaction effects.

4.3 Case study

To further demonstrate IPGP in longitudinal setting for learning idiographic psychological taxonomies, we collected an intensive longitudinal data using experience sampling measures (ESM). We also highlight the predictive ability of IPGP through a forecasting and a leave-one-trait-out cross validation task, and illustrate how IPGP identifies unique taxonomies of personality that might advance individualized approaches to psychological diagnosis and inspire new theory.

Data collection. In ESM design, each participant was asked to complete personality assessments six times per day for three weeks, resulting maximum 126 assessments per person. With 93 valid student participants, we acquired 8,770 assessments in total with an average of 94 assessments per person. The personality assessment is derived from the BFI-2 [Soto and John, 2017] to ensure identification of latent factors and ample coverage of the late trait space. The BFI-2 includes 60 items with four unique items assessing each of the three different sub-factors for each Big-Five domains. We removed one item for each sub-factors that are not appropriate for contextualized assessments of ESM design. To mitigate the fatigue and learning effect from repeated measures, we employed a planned missing design where participants were randomly tested on only two out of three items assessing the same sub-factors, resulting only 30 items for each assessment.

Comparison between nomothetic and idiographic models. We run the full IPGP model with idiographic component and unit-specific time kernel on the collected longitudinal data. Again we set the ranks of the population and individual loading matrices to 5 and 1 respectively, and incorporate the prior knowledge of the cross-sectional data by fixing the population loadings as the Big Five loadings estimated in Sec. (4.2) and optimizing the individual loadings. We contrast our proposed idiographic model (IPGP) and baselines in Table 3, which shows the in-sample prediction, averaged log likelihood and Bayes factors. We found that IPGP outperforms IPGP-NOM with higher predictive accuracy and log likelihood, and is favored decisively by a Bayes factor of $\exp(1.06 \times 10^4)$.

Table 3: In-sample prediction and averaged log likelihood of our proposed model (IPGP) and baselines for the longitudinal data, as well as Bayes factors to IPGP. “—” indicates self comparison of Bayes factors.

	ACC	LL/N	log(BF)
GRM	0.210	-2.266	-2.32×10^5
GPCM	0.288	-1.516	-3.80×10^4
SRM	0.260	-1.927	-1.44×10^5
GPDM	0.382	-3.865	-7.80×10^5
DSEM	0.226	-1.399	-7.72×10^3
TVAR	0.382	-1.546	-4.47×10^4
IPGP-NOM	0.403	-1.410	-1.06×10^4
IPGP	0.417	-1.369	—

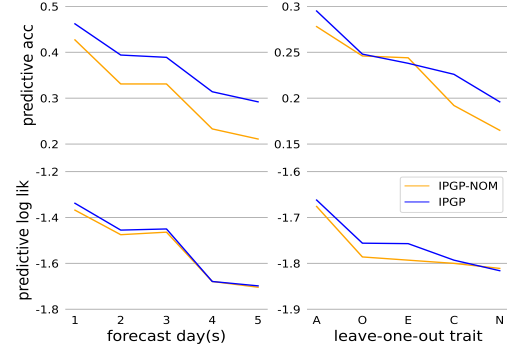


Figure 3: Predictive accuracy and log lik of IPGP and IPGP-NOM for the forecasting task and leave-one-trait-out cross-validation task.

Predictive performance of IPGP. We also evaluate the out-of-sample performance of the idiographic and nomothetic models using two prediction tasks: forecasting future responses and leave-one-trait-out cross validation. For the forecasting task, we train both models with data from the first 40 days and predict future responses for the last 5 days. For the cross validation task, we predict responses of each trait by training on data belonging to the other four traits. Figure (3) shows the predictive accuracy and log likelihood of IPGP and IPGP-NOM for the forecasting task over varying horizons and the leave-one-trait-out cross-validation task. IPGP has consistently better performance than IPGP-NOM in both tasks except for being slightly less accurate in predicting extraversion. Overall, IPGP is favored than IPGP-NOM by Bayes factors of $\exp(43)$ and $\exp(716)$ in these tasks.

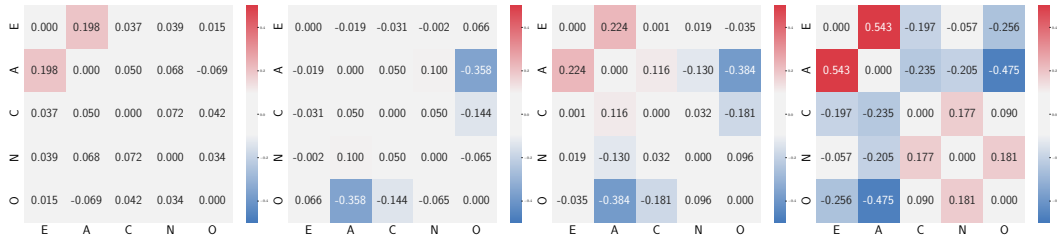


Figure 4: Four residual correlations as identified by our k-mean clustering. Each heatmap displays the trait-level residual correlation averaged across corresponding batteries for one cluster, with darker red and blue indicating larger positive and negative deviations. For instance, agreeableness (A) is more correlated to extraversion (E) than the population profile in the first profile, but less correlated to openness (O) in the second profile. Moreover, these two directions of deviations are even exacerbated in the third and fourth profiles.

Discovery of unique taxonomies. Despite our small cohort size (93 respondents), we also manage to identify distinct profiles of personality that substantially differ from the interpersonal commonality. Specifically, we first perform a k-mean clustering using all 93 estimated individual correlation matrix with CMD as the distance metric, and then compute the residual correlation between each estimated clustering centroid and the population correlation. Figure (4) illustrates four residual correlations as identified by our k -mean clustering. Each heatmap displays the trait-level residual correlation

averaged across corresponding batteries for one cluster, with darker red and blue indicating larger positive and negative deviations. For instance, agreeableness (A) is more correlated to extraversion (E) than population profile in the first profile, but less correlated to openness (O) in the second profile. Moreover, these two directions of deviations are even exacerbated in the third and fourth profiles.

5 Related work

Idiographic assessment emphasizes the important aspects of individuals otherwise missing in oversimplified taxonomies of psychological behaviors [Hamaker and Dolan, 2009]. Empirical evidence across many psychometrics fields has shown the lack of generalizability of the nomothetic models only focusing on interpersonal variation [Molenaar, 2004]. Hence, Song and Ferrer [2012] incorporated simple random effect method with dynamic factor models for analyzing psychological processes. Jongerling et al. [2015] proposed a multilevel first-order autoregressive model with random intercepts to measure daily positive effects over several weeks. Beltz et al. [2016] combined the nomothetic and idiographic approaches in analyzing clinical data by adding individual components to the group iterative multiple model (GIMME). However, all of these methods focus on modeling in response space rather than latent space for ordinal survey data.

Gaussian process latent variable model (GPLVM) is a dimensional reduction method for Gaussian data, where the latent variables are optimized after integrating out the function mappings [Lawrence, 2003, Lalchand et al., 2022]. Our proposed framework differs from GPLVM as we optimize the factor loading matrix while marginalizing the latent variables. In addition, our model contrasts GPLVM and (variational) Gaussian Process dynamical model (GPDm) [Wang et al., 2005, Damianou et al., 2011] in our non-Gaussian ordered logistic observation model. Finally, our longitudinal framework with stochastic variational inference learning differs from the static GP item response theory (GPIRT) [Duck-Mayr et al., 2020] with more computationally demanding Gibbs sampling.

Longitudinal measurement models integrate temporal dynamics into psychological theories with growing popularity of longitudinal design in survey methods [Jebb et al., 2015, Ariens et al., 2020]. For instance, families of longitudinal structural equation models (SEM) such as multiple-group longitudinal SEM and longitudinal growth curve model were developed for repeated measurement studies [Little, 2013], where *Mplus* software was developed later for dynamic SEM with Bayesian Gibbs sampling [Asparouhov et al., 2018, McNeish et al., 2023]. Dynamic item response models [Rijmen et al., 2003, Reise and Waller, 2009, Dumas et al., 2020] and time-varying vector autoregressive model [Lu et al., 2018, Haslbeck et al., 2021] were also proposed to estimate the trajectories of latent traits. Despite previous work in behavioral literature focusing on Gaussian observations [Dürichen et al., 2014], multi-task Gaussian process time series has not yet been exploited for survey experiments with non-Gaussian likelihood when exact inference is not plausible.

6 Conclusion

We propose a novel idiographic personality Gaussian process (IPGP) model for personalized psychological assessment and learning of intrapersonal taxonomy from longitudinal ordinal survey data, an under-explored setup in Gaussian process dynamic system literature. We exploit Gaussian process coregionalization for capturing between-battery structure and stochastic variational inference for scalable inference. Future directions include adaptation of IPGP to other psychological studies such as emotion, and incorporation of contextual information such as behaviors or demographics.

Our proposed IPGP framework also provides insights to domain theory testing, in this work, addressing the substantive debate in psychometrics surrounding the shared versus unique structures of psychological features. Besides predicting the actual responses, we also include learning of the true underlying taxonomy structures as evaluation metrics to minimize risks of inductive bias. Our experimental results show that IPGP is decisively favored than the nomothetic baseline, and substantive deviations from the common trend persist in considerable individuals. Hence, our framework has a great potential in advancing individualized approaches to psychological diagnosis and treatment.

Acknowledgments and Disclosure of Funding

This work was supported by the 2023 Seed Grant of Transdisciplinary Institute in Applied Data Sciences at Washington University.

References

- Peter CM Molenaar. A Manifesto on Psychology as Idiographic Science: Bringing the Person Back Into Scientific Psychology, This Time Forever. *Measurement*, 2(4):201–218, 2004.
- Xiaojing Wang, James O Berger, and Donald S Burdick. Bayesian analysis of dynamic item response models in educational testing. *The Annals of Applied Statistics*, 7(1):126–153, 2013.
- Denis Dumas, Daniel McNeish, and Jeffrey A Greene. Dynamic measurement: A theoretical–psychometric paradigm for modern educational psychology. *Educational Psychologist*, 55(2): 88–105, 2020.
- Denny Borsboom, Gideon J Mellenbergh, and Jaap Van Heerden. The theoretical status of latent variables. *Psychological Review*, 110(2):203, 2003.
- Peter Borkenau and Fritz Ostendorf. The Big Five as States: How Useful Is the Five-Factor Model to Describe Intraindividual Variations over Time? *Journal of Research in Personality*, 32(2):202–221, 1998.
- Emorie D Beck and Joshua J Jackson. Consistency and Change in Idiographic Personality: A Longitudinal ESM Network Study. *Journal of Personality and Social Psychology*, 118(5):1080, 2020.
- Emorie D Beck and Joshua J Jackson. Within-person variability. In *The Handbook of Personality Dynamics and Processes*, pages 75–100. Elsevier, 2021.
- Frank Rijmen, Francis Tuerlinckx, Paul De Boeck, and Peter Kuppens. A nonlinear mixed model framework for item response theory. *Psychological Methods*, 8(2):185, 2003.
- Steven P Reise and Niels G Waller. Item Response Theory and Clinical Measurement. *Annual Review of Clinical Psychology*, 5:27–48, 2009.
- Todd D Little. *Longitudinal Structural Equation Modeling*. Guilford Press, 2013.
- Eun Sook Kim and Victor L Willson. Testing Measurement Invariance Across Groups in Longitudinal Data: Multigroup Second-Order Latent Growth Model. *Structural Equation Modeling: A Multidisciplinary Journal*, 21(4):566–576, 2014.
- Tihomir Asparouhov, Ellen L Hamaker, and Bengt Muthén. Dynamic Structural Equation Models. *Structural Equation Modeling: a Multidisciplinary Journal*, 25(3):359–388, 2018.
- Feihan Lu, Yao Zheng, Harrington Cleveland, Chris Burton, and David Madigan. Bayesian hierarchical vector autoregressive models for patient-level predictive modeling. *PloS One*, 13(12): e0208082, 2018.
- Jonas MB Haslbeck, Laura F Bringmann, and Lourens J Waldorp. A Tutorial on Estimating Time-Varying Vector Autoregressive Models. *Multivariate Behavioral Research*, 56(1):120–149, 2021.
- Jack Wang, Aaron Hertzmann, and David J Fleet. Gaussian Process Dynamical Models. In Y. Weiss, B. Schölkopf, and J. Platt, editors, *Advances in Neural Information Processing Systems*, volume 18. MIT Press, 2005.
- Andreas Damianou, Michalis Titsias, and Neil Lawrence. Variational Gaussian Process Dynamical Systems. *Advances in Neural Information Processing systems*, 24, 2011.
- Robert Dürichen, Marco AF Pimentel, Lei Clifton, Achim Schweikard, and David A Clifton. Multi-task Gaussian Processes for Multivariate Physiological Time-Series Analysis. *IEEE Transactions on Biomedical Engineering*, 62(1):314–322, 2014.

- JBrandon Duck-Mayr, Roman Garnett, and Jacob Montgomery. GPIRT: A Gaussian Process Model for Item Response Theory. In *Conference on Uncertainty in Artificial Intelligence*, pages 520–529. PMLR, 2020.
- Adriene M. Beltz, Aidan G. C. Wright, Briana N. Sprague, and Peter C. M. Molenaar. Bridging the Nomothetic and Idiographic Approaches to the Analysis of Clinical Data. *Assessment*, 23(4): 447–458, 2016.
- Edwin V Bonilla, Kian Chai, and Christopher Williams. Multi-task Gaussian Process Prediction. *Advances in Neural Information Processing Systems*, 20, 2007.
- Edward Snelson and Zoubin Ghahramani. Sparse Gaussian Processes using Pseudo-inputs. In Y. Weiss, B. Schölkopf, and J. Platt, editors, *Advances in Neural Information Processing Systems*, volume 18. MIT Press, 2005.
- James Hensman, Alexander Matthews, and Zoubin Ghahramani. Scalable Variational Gaussian Process Classification. In *Artificial Intelligence and Statistics*, pages 351–360. PMLR, 2015.
- Neil Lawrence. Gaussian Process Latent Variable Models for Visualisation of High Dimensional Data. *Advances in Neural Information Processing Systems*, 16, 2003.
- James T Croasmun and Lee Ostrom. Using Likert-Type Scales in the Social Sciences. *Journal of Adult Education*, 40(1):19–22, 2011.
- Wei Chu and Zoubin Ghahramani. Gaussian Processes for Ordinal Regression. *Journal of Machine Learning Research*, 6(35):1019–1041, 2005.
- John M Digman. Higher-order factors of the Big Five. *Journal of Personality and Social Psychology*, 73(6):1246, 1997.
- James Baglin. Improving Your Exploratory Factor Analysis for Ordinal Data: A Demonstration Using FACTOR. *Practical Assessment, Research, and Evaluation*, 19(1):5, 2014.
- Kenneth A Bollen. *Structural Equations with Latent Variables*, volume 210. John Wiley & Sons, 1989.
- Fumiko Samejima. Estimation of latent ability using a response pattern of graded scores. *Psychometrika Monograph Supplement*, 1969.
- Wim J Van der Linden and RK Hambleton. Handbook of Item Response Theory. *Taylor & Francis Group*. Cited on page 1(7):8, 1997.
- R Darrell Bock and Murray Aitkin. Marginal maximum likelihood estimation of item parameters: Application of an EM algorithm. *Psychometrika*, 46(4):443–459, 1981.
- Carlos G Forero, Alberto Maydeu-Olivares, and David Gallardo-Pujol. Factor Analysis with Ordinal Indicators: A Monte Carlo Study Comparing DWLS and ULS Estimation. *Structural Equation Modeling: A Multidisciplinary Journal*, 16(4):625–641, 2009. doi: 10.1080/10705510903203573.
- Cheng-Hsien Li. Confirmatory factor analysis with ordinal data: Comparing robust maximum likelihood and diagonally weighted least squares. *Behavior Research Methods*, 48:936–949, 2016.
- Sergio Salvatore and Jaan Valsiner. Between the General and the Unique: Overcoming the Nomothetic versus Idiographic Opposition. *Theory & Psychology*, 20(6):817–833, 2010.
- Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine Learning*. The MIT Press, 11 2005.
- Mauricio A Alvarez, Lorenzo Rosasco, Neil D Lawrence, et al. Kernels for Vector-Valued Functions: A Review. *Foundations and Trends® in Machine Learning*, 4(3):195–266, 2012.
- Robert R McCrae and Oliver P John. An Introduction to the Five-Factor Model and Its Applications. *Journal of Personality*, 60(2):175–215, 1992.

- Markus Herdin, Nicolai Czink, Hüseyin Ozelik, and Ernst Bonek. Correlation matrix distance, a meaningful measure for evaluation of non-stationary MIMO channels. In *2005 IEEE 61st Vehicular Technology Conference*, volume 1, pages 136–140. IEEE, 2005.
- Eiji Muraki. A Generalized Partial Credit Model: Application of an EM Algorithm. *Applied Psychological Measurement*, 16(2):159–176, 1992.
- Gerhard Tutz. Sequential item response models with an ordered response. *British Journal of Mathematical and Statistical Psychology*, 43(1):39–55, 1990.
- Daniel McNeish, Jennifer A Somers, and Andrea Savord. Dynamic structural equation models with binary and ordinal outcomes in *Mplus*. *Behavior Research Methods*, pages 1–27, 2023.
- Christopher J Soto. How Replicable Are Links Between Personality Traits and Consequential Life Outcomes? The Life Outcomes of Personality Replication Project. *Psychological Science*, 30(5): 711–727, 2019.
- Oliver P John, Sanjay Srivastava, et al. The Big-Five Trait Taxonomy: History, Measurement, and Theoretical Perspectives. *Handbook of personality: Theory and research*, 2:102—138, 1999.
- Christopher J Soto and Oliver P John. The next Big Five Inventory (BFI-2): Developing and assessing a hierarchical model with 15 facets to enhance bandwidth, fidelity, and predictive power. *Journal of Personality and Social Psychology*, 113(1):117, 2017.
- Ellen L Hamaker and Conor V Dolan. Idiographic Data Analysis: Quantitative Methods—From Simple to Advanced. In *Dynamic Process Methodology in the Social and Developmental Sciences*, pages 191–216. Springer, 2009.
- Hairong Song and Emilio Ferrer. Bayesian estimation of random coefficient dynamic factor models. *Multivariate Behavioral Research*, 47(1):26–60, 2012.
- Joran Jongerling, Jean-Philippe Laurenceau, and Ellen L Hamaker. A Multilevel AR(1) Model: Allowing for Inter-Individual Differences in Trait-Scores, Inertia, and Innovation Variance. *Multivariate Behavioral Research*, 50(3):334–349, 2015.
- Vidhi Lalchand, Aditya Ravuri, and Neil D. Lawrence. Generalised GPLVM with Stochastic Variational Inference. In Gustau Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera, editors, *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, pages 7841–7864. PMLR, 28–30 Mar 2022.
- Andrew T Jebb, Louis Tay, Wei Wang, and Qiming Huang. Time series analysis for psychological research: examining and forecasting change. *Frontiers in Psychology*, 6:727, 2015.
- Sigert Ariens, Eva Ceulemans, and Janne K Adolf. Time series analysis of intensive longitudinal data in psychosomatic research: A methodological overview. *Journal of Psychosomatic Research*, 137:110191, 2020.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [\[Yes\]](#)

Justification: The main claims are matched to the experimental results in Sec 4.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: The limitation of this work is discussed in Sec 6.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[NA\]](#)

Justification: This paper does not include theoretical results.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: All information necessary for reproduction are fully discussed in Sec 4. All code and data are uploaded as supplementary materials.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: All code and data in this paper are uploaded as supplementary materials.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Experimental setup and details are fully discussed in Sec 4.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes] .

Justification: Statistical significance are fully reported in Sec 4.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.

- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Computing resources used in this paper are reported in Sec 4.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: We have reviewed NeurIPS Code of Ethics and made sure anonymity is preserved in both main paper submission and supplementary materials.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: The broad impact of this work on psychometrics is discussed in Sec 6.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA] .

Justification: This paper poses no risks for misuse.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes] .

Justification: We have provided proper citations of existing assets of data in Sec 4.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.

13. **New Assets**

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [\[Yes\]](#)

Justification: We promise to release our collected data along with the IRB number in the camera-ready version. The documentation is uploaded as supplementary materials.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. **Crowdsourcing and Research with Human Subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [\[Yes\]](#)

Justification: We have uploaded a complete documentation of experimental instructions and survey questions as supplementary materials.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [\[Yes\]](#)

Justification: Our study is approved by Institutional Review Board and we promise to report the IRB number in the camera-ready version.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

A Mathematical Details of Evidence Lower Bound

We provide the full mathematical details of the evidence lower bound defined in Eq. (6). As KL divergence is always non-negative, we first consider the KL divergence between $p(\mathbf{f} | \mathbf{u})$ and $p(\mathbf{f} | \mathbf{y})$:

$$\text{KL}[p(\mathbf{f} | \mathbf{u}) \parallel p(\mathbf{f} | \mathbf{y})] = \mathbb{E}_{p(\mathbf{f}|\mathbf{u})} \log \frac{p(\mathbf{f} | \mathbf{u})}{p(\mathbf{f} | \mathbf{y})} \quad (7)$$

$$= \mathbb{E}_{p(\mathbf{f}|\mathbf{u})} \log \frac{p(\mathbf{f} | \mathbf{u})p(\mathbf{y})}{p(\mathbf{y} | \mathbf{f})p(\mathbf{f})} \quad (8)$$

$$= \mathbb{E}_{p(\mathbf{f}|\mathbf{u})} \log \frac{p(\mathbf{f} | \mathbf{u})p(\mathbf{y} | \mathbf{u})p(\mathbf{u})}{p(\mathbf{y} | \mathbf{f})p(\mathbf{f})} \quad (9)$$

$$= \mathbb{E}_{p(\mathbf{f}|\mathbf{u})} \log \frac{p(\mathbf{y} | \mathbf{u})}{p(\mathbf{y} | \mathbf{f})} \quad (10)$$

$$= \log p(\mathbf{y} | \mathbf{u}) - \mathbb{E}_{p(\mathbf{f}|\mathbf{u})} \log p(\mathbf{y} | \mathbf{f}) \geq 0 \quad (11)$$

Moving $\mathbb{E}_{p(\mathbf{f}|\mathbf{u})} \log p(\mathbf{y} | \mathbf{f})$ to the R.H.S of the above inequality will lead to Eq. (4). We then exploit the inequality given by $\text{KL}[q(\mathbf{u}) \parallel p(\mathbf{u} | \mathbf{y})] \geq 0$:

$$\text{KL}[q(\mathbf{u}) \parallel p(\mathbf{u} | \mathbf{y})] = \mathbb{E}_{q(\mathbf{u})} \log \frac{q(\mathbf{u})}{p(\mathbf{u} | \mathbf{y})} \quad (12)$$

$$= \mathbb{E}_{q(\mathbf{u})} \log \frac{q(\mathbf{u})p(\mathbf{y})}{p(\mathbf{y} | \mathbf{u})p(\mathbf{u})} \quad (13)$$

$$= -\mathbb{E}_{q(\mathbf{u})} \log p(\mathbf{y} | \mathbf{u}) + \text{KL}[q(\mathbf{u}) \parallel p(\mathbf{u})] + \log p(\mathbf{y}) \geq 0 \quad (14)$$

Rearranging the above inequality, applying Eq. (4) and exploiting notation $q(\mathbf{f}) = \int p(\mathbf{f} | \mathbf{u})q(\mathbf{u})d\mathbf{u}$ leads to the ELBO:

$$\log p(\mathbf{y}) \geq \mathbb{E}_{q(\mathbf{u})} \log p(\mathbf{y} | \mathbf{u}) - \text{KL}[q(\mathbf{u}) \parallel p(\mathbf{u})] \quad (15)$$

$$= \mathbb{E}_{q(\mathbf{u})} [\mathbb{E}_{p(\mathbf{f}|\mathbf{u})} \log p(\mathbf{y} | \mathbf{f})] - \text{KL}[q(\mathbf{u}) \parallel p(\mathbf{u})] \quad (16)$$

$$= \mathbb{E}_{q(\mathbf{f})} \log p(\mathbf{y} | \mathbf{f}) - \text{KL}[q(\mathbf{u}) \parallel p(\mathbf{u})] \quad (17)$$

B Estimated Correlations of Selective Individuals

Figure (5) shows the estimated correlations of selective individuals for the identified four profiles in the longitudinal study.

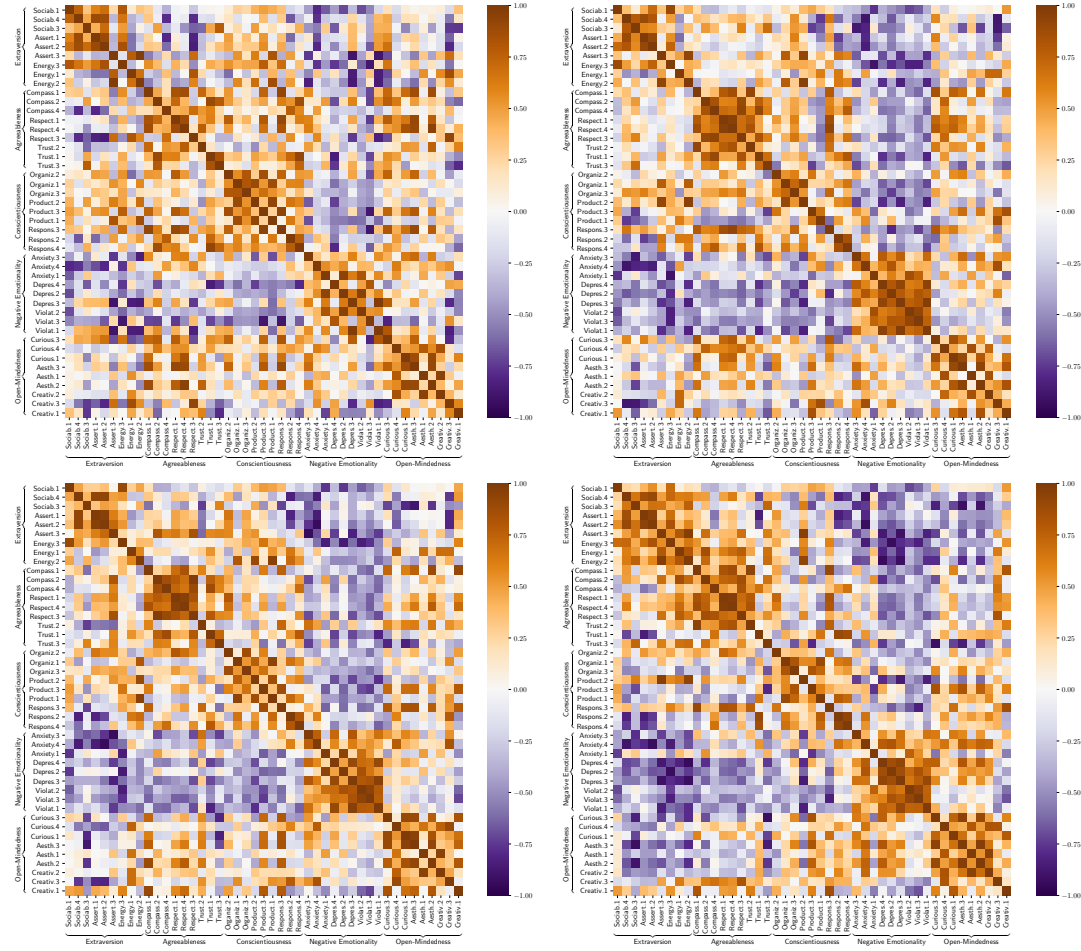


Figure 5: Estimated correlations of selective individuals for the identified four profiles in the longitudinal study.