

Bayesian Regularization in Multiple-Indicators Multiple-Causes Models

Lijin Zhang

Sun Yat-sen University, China, zhanglj37@mail2.sysu.edu.cn

Joint work with Dr. Xinya Liang

Department of Rehabilitation, Human Resources and Communication Disorders, University of Arkansas, USA

20 July 2021

Overview

① Introduction

MIMIC Models

② Bayesian Regularization

Ridge

Lasso

Adaptive Lasso

Spike-and-Slab

Horseshoe

③ Simulation Study

④ Empirical Illustration

⑤ Discussion

Introduction

Regularized Structural Equation Modeling (SEM)

Trade-off between model fitting and model complexity in SEM.

- complex model - overfitting, less generalizability
- simple model - omit important variables, poor model fitting

Regularization for achieving model parsimony and meaningful interpretations simultaneously (e.g., Jacobucci & Grimm, 2018a).

- shrink nuisance parameters toward zero & identify essential parameters
- retain accurate parameter estimates and improve the generalizability of estimates

Bayesian Regularization

Bayesian regularization assigns penalty priors to regularize the posterior distributions of parameters.

It is flexible in (Polson & Sokolov, 2019; van Erp et al., 2019):

- estimating the shrinkage parameters
- quantifying the uncertainty of parameter estimates
- handling number of parameters \geq sample size conditions

Different penalty priors:

- Ridge: global shrinkage
- Lasso (least absolute shrinkage and selection operator; Park & Casella, 2008; Tibshirani, 1996): global shrinkage
- Adaptive lasso (alasso; Zou, 2006): local shrinkage
- Spike-and-slab prior (SSP; Mitchell & Beauchamp, 1988): assign a discrete mixture of normal distributions on parameters
- Horseshoe (Carvalho et al., 2010): global-local shrinkage

Bayesian Regularized SEM

Methods	Measurement Models	Structural Models
Ridge	Muthén & Asparouhov (2012, 2013)	
Lasso	Chen et al. (2021), Pan et al. (2017);	-
Alasso	Chen (2021), Pan et al. (2021);	Feng et al. (2017), Jacobucci & Grimm (2018b), Brandt et al. (2018);
SSP	Lu et al. (2016)	Brandt et al. (2018)
Horseshoe	-	-

Tabel 1: Integration of Different Penalty Priors with SEM

	Comparison		Model
Chen et al. (2021)	Ridge	Lasso	Measurement
Lu et al. (2016)	Ridge	SSP	Measurement
Feng et al. (2017)	Alasso	Lasso	Structural
Brandt et al. (2018)	Alasso	Lasso	Structural

Tabel 2: Comparison between Different Penalty Priors

- Lasso and SSP have advantages in achieving parsimonious factor structures than ridge.
- Alasso has benefits in reducing appreciable bias caused by global lasso shrinkage.

Investigate the performance of different Bayesian regularization methods in parameter estimation and variable selection using MIMIC models:

- penalty priors vs non-informative prior
- global vs local vs global-local shrinkage
- under different modeling conditions (sample sizes, multicollinearity, effect sizes)

Suppose there are K latent factors ω measured by J indicators \mathbf{y} and regressed on P predictors \mathbf{X} , a MIMIC model (Jöreskog & Goldberger, 1975) can be expressed as follows:

$$\mathbf{y}_i = \boldsymbol{\mu} + \boldsymbol{\Lambda}\omega_i + \boldsymbol{\epsilon}_i, i = 1, 2, \dots, n, \quad (1)$$

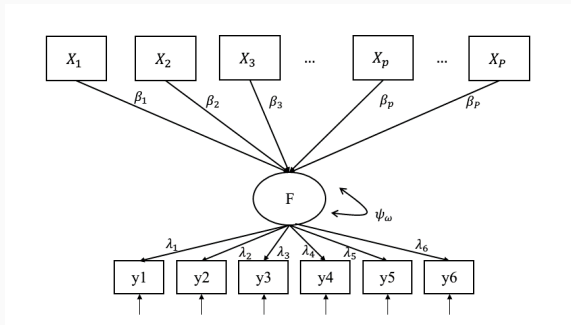
$$\omega_i = \boldsymbol{\mu}_\omega + \boldsymbol{\beta}\mathbf{X}_i + \boldsymbol{\delta}_i \quad (2)$$

- \mathbf{y}_i : observed values of J indicators for the i -th participant.
- $\boldsymbol{\mu}, \boldsymbol{\mu}_\omega$: vector of intercepts.
- $\boldsymbol{\Lambda}$: factor loading matrix.
- ω_i : latent factors.
- $\boldsymbol{\epsilon}_i$: measurement errors.
- $\boldsymbol{\beta}$: path coefficients.
- $\boldsymbol{\delta}_i$: factor disturbances.

MIMIC Models

The utility of MIMIC models is versatile (Finch & Miller, 2019)

- control the influence of covariates on latent variables
- test the measurement invariance between groups
- identify differential item functioning



MIMIC Model with One-factor and Six-indicators

Bayesian Regularization

$$\beta_p \sim N(0, \sigma^2) \quad (3)$$

- β_p : p-th parameter to be regularized.
- σ^2 : variance hyperparameter, determines the penalty strength.
- The prior variance can be fixed at a preassigned value such as .01 (Muthén & Asparouhov, 2012) or .001 (Jacobucci & Grimm, 2018a), or be estimated through a hyperprior.

Applications

- identify cross-loadings and residual correlations (e.g., Falkenström et al., 2015)
- handle small sample sizes (e.g., Crenshaw et al., 2016)
- assess measurement invariance (e.g., de Bondt & van Petegem, 2015)

$$\beta_p \sim N(0, \psi_\omega \tau_p^2), \psi_\omega^{-1} \sim \text{Gamma}(\alpha_\omega, \beta_\omega) \quad (4)$$

$$\tau_p^2 \sim \text{Gamma}(1, \frac{\gamma^2}{2}), \gamma^2 \sim \text{Gamma}(a_l, b_l) \quad (5)$$

- τ_p is included to obtain the desired Laplace distribution of the conditional prior.
- γ is the global penalty parameter.

Applications

- identify cross-loadings and residual correlations (Chen et al., 2021; Pan et al., 2017; Zhang et al., 2021)

$$\beta_p \sim N(0, \psi_\omega \tau_p^2), \psi_\omega^{-1} \sim \text{Gamma}(\alpha_\omega, \beta_\omega) \quad (6)$$

$$\tau_p^2 \sim \text{Gamma}(1, \frac{\gamma_p^2}{2}), \gamma_p^2 \sim \text{Gamma}(a_l, b_l) \quad (7)$$

- γ_p : local penalty parameter.

Bayesian adaptive lasso has been extended to

- SEMs with ordinal variables (Feng et al., 2018)
- latent change score models (Jacobucci & Grimm, 2018b)
- detect multiple linear and nonlinear effects in SEM with SSP (Brandt et al., 2018)

$$\beta_p \sim r_p N(0, c_p^2) + (1 - r_p) N(0, \sigma_p^2) \quad (8)$$

$$r_p \sim \text{Bernoulli}(.5) \quad (9)$$

- r_p : selection variable
- $N(0, \sigma_p^2)$: a point mass function (spike) commonly with a small prior variance to shrink the parameter to zero
- $N(0, c_p^2)$: the fuzzy prior (slab) that is typically assigned a large prior variance

$$\beta_p \sim N(0, \rho_p^2 v^2), \rho_p^2 \sim C^+(0, 1), v \sim C^+(0, 1) \quad (10)$$

- ρ_p, v : local and global shrinkage parameters, respectively.
- Placing the half-Cauchy distributions $C^+(0, 1)$ is similar to putting a $\text{beta}(.5, .5)$ prior on the shrinkage weight $\kappa_p = 1/(1 + \rho_p^2 v^2)$.

Global shrinkage (ridge & lasso) may induce appreciable bias, which may be solved by methods that include local shrinkage parameters.

Literature Gap

- How the global and/or local shrinkage affects the convergence rate, variable selection, and parameter estimation in SEM ?
- How do factors such as collinearity and sample sizes influence the performance of Bayesian regularization?

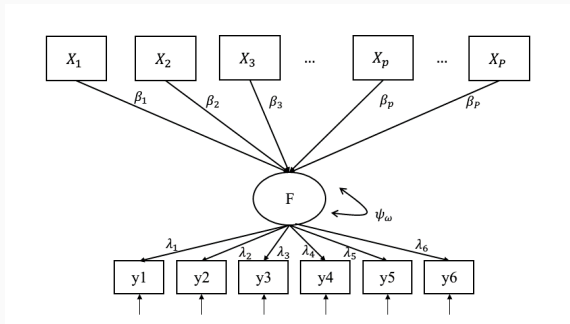
Simulation Study

Purpose and Design

We conducted a simulation study with a similar design to Jacobucci et al. (2019):

- Collinearity among covariates: 0, .2, .5, .8, .95
- Sample size: 100, 200, 300, 500, 1000

$$\frac{\text{sample size}}{\text{number of observed variables}} = .94 - 9.43$$



Other Settings

- Effect Sizes:

$$\beta_1 - \beta_{70} = 0, \beta_{71} - \beta_{80} = .2, \beta_{81} - \beta_{90} = .5, \beta_{91} - \beta_{100} = .8$$

- Factor loadings: `c(1, .8, .8, .8, .5, .5)`
- Residual variances of indicators and factor disturbance: 1
- Number of replications: 200 datasets per condition

Model Estimation

Software: R, JAGS (Plummer, 2003)

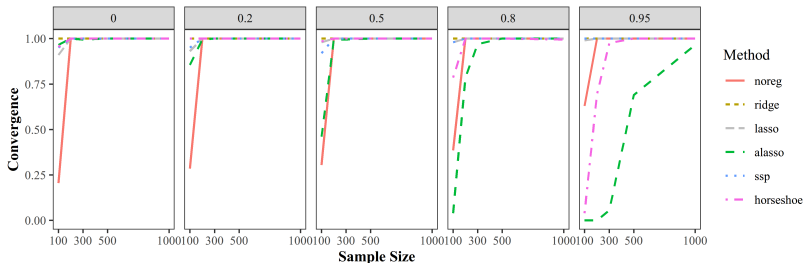
Priors and Hyperparameters for path coefficients:

- Diffuse prior: $N(0, 1000)$
- Ridge: $N(0, .01)$ (Muthén & Asparouhov, 2012)
- Lasso and Alasso: $\alpha_l = 1, \beta_l = .01$ (Chen et al., 2021)
- SSP: $\sigma_p^2 = .001, c_p^2 \sim IG(.5, .5)$ (van Erp et al., 2019)
- Horseshoe: $\rho_p \sim C^+(0, 1), v \sim C^+(0, 1)$
- Priors for other parameters (i.e., loadings): diffuse priors.
- Model convergence criteria: A less than 1.05 estimated potential scale reduction (EPSR; Gelman et al., 1996) value within 5,000 - 20,000 burn-in iterations.

- Convergence Rate
- Rejection Rate of 95% HPD interval
- Rejection Rate of Threshold: the proportion of converged replications where $|\beta_{est}| > 0.1$ (Feng et al., 2017)
- 95% Coverage Rate
- Relative Bias
- Root Mean Square Error (RMSE)

$\sqrt{\frac{1}{N} \sum_{i=1}^N (\beta_{est} - \beta_{true})^2}$ where N is the number of converged replications

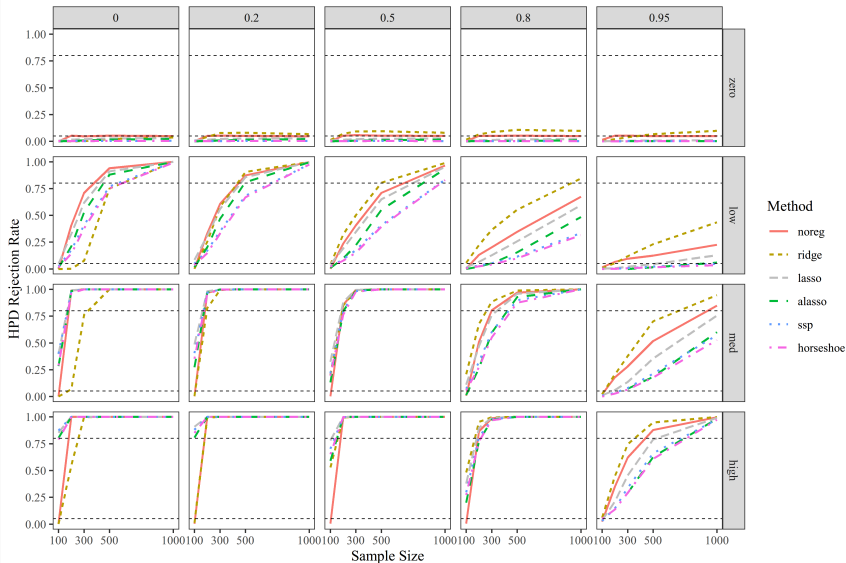
Results: Convergence Rates



noreg: diffuse prior; collinearity: 0, .2, .5, .8, and .95

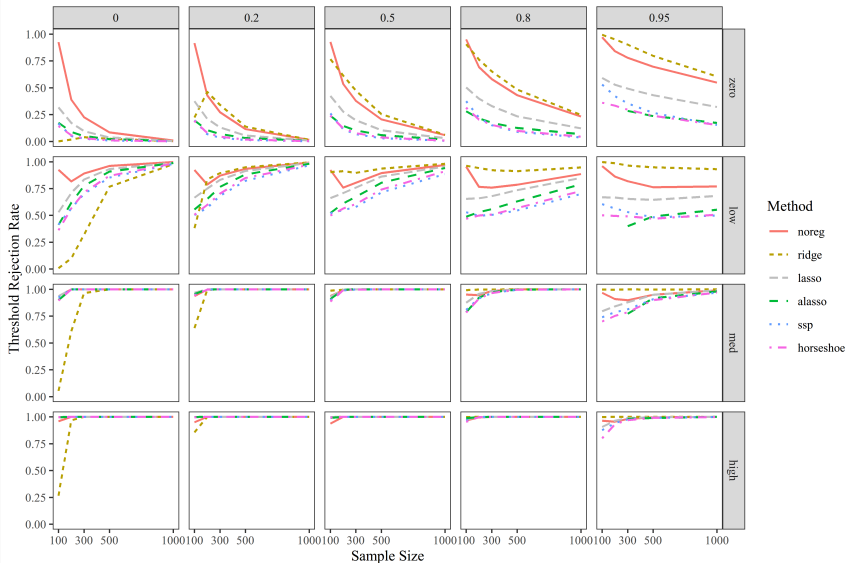
- Diffuse prior: low convergence rates (.21 - .63) with small sample size
- Ridge, lasso, and SSP: excellent convergence rates under all conditions
- The convergence rates decreased as the collinearity increased for adaptive lasso and horseshoe

Results: 95% HPD Rejection Rates

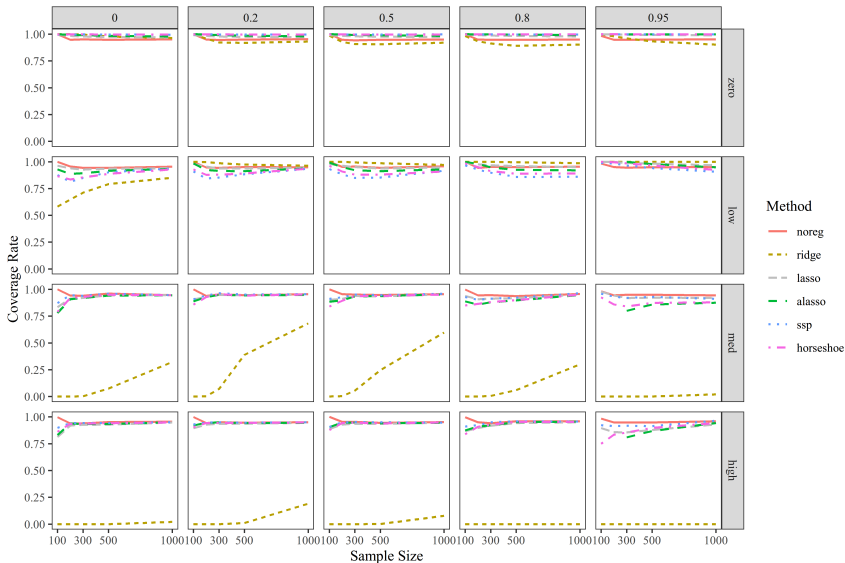


effect size: zero, low (.2), medium (.5), and high (.8)

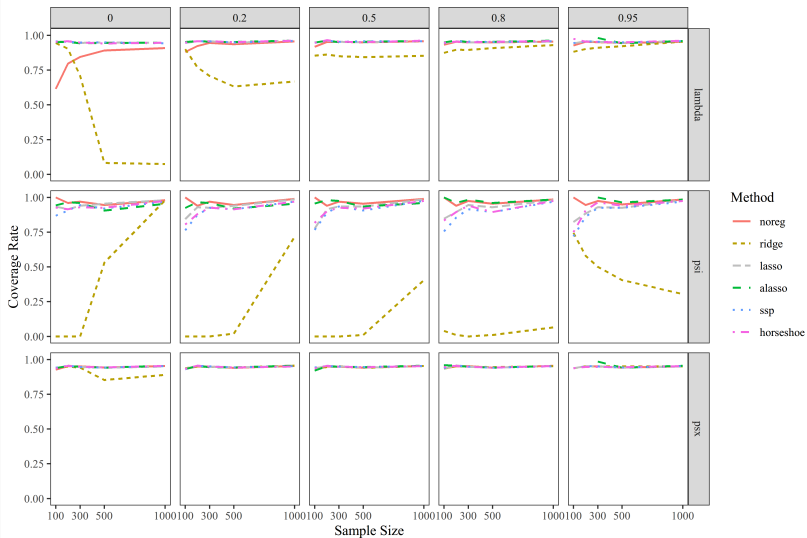
Results: Threshold Rejection Rates



Results: 95% Coverage Rate of Path Coefficients

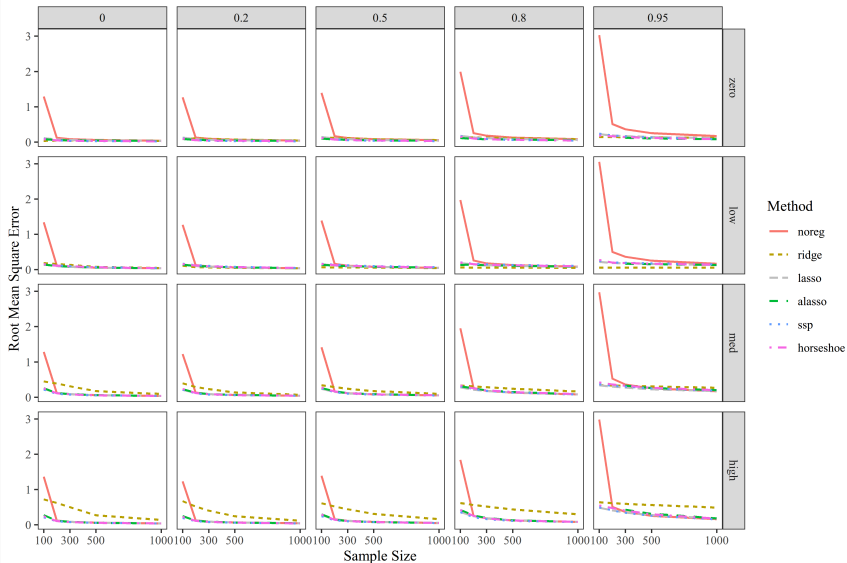


Results: 95% Coverage Rate of Other parameters

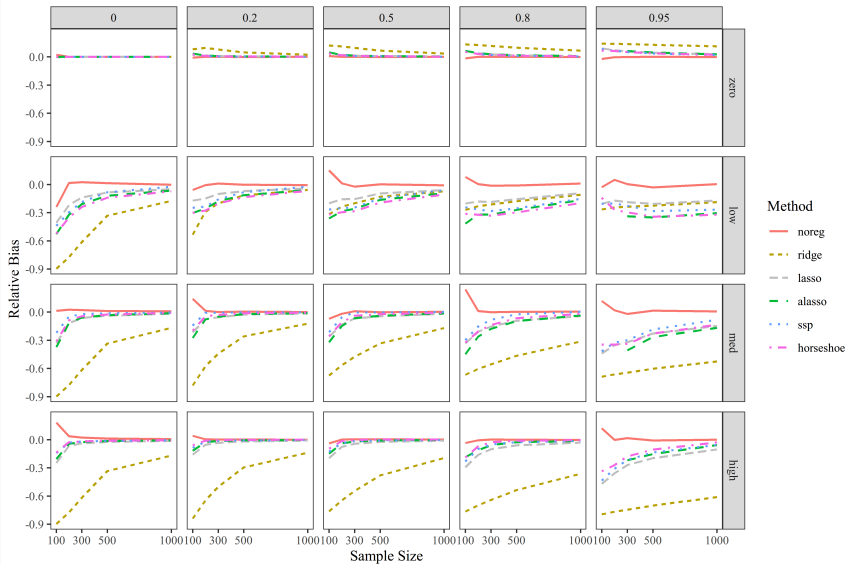


lambda = loadings; psi = factor disturbance; psx = residual
variances of indicators

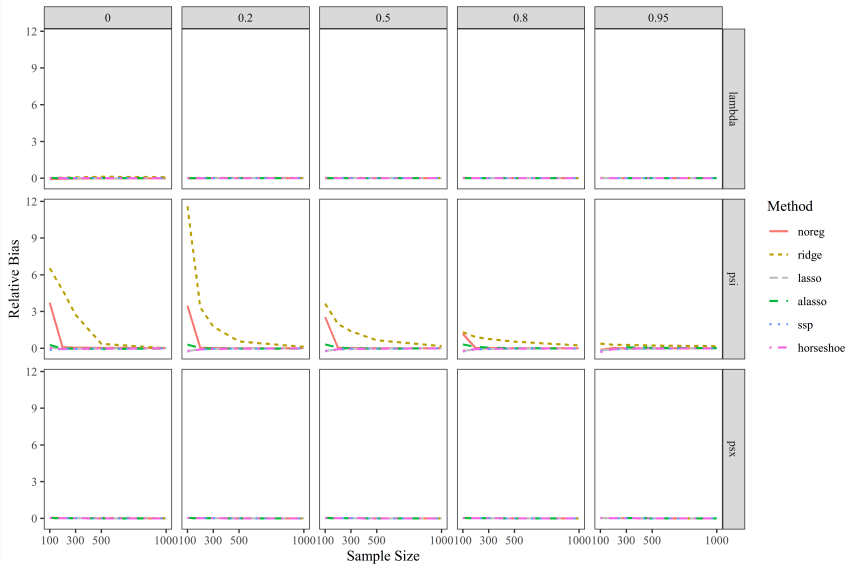
Results: RMSE of Path Coefficients



Results: Relative Bias



Results: Relative Bias



Empirical Illustration

Data

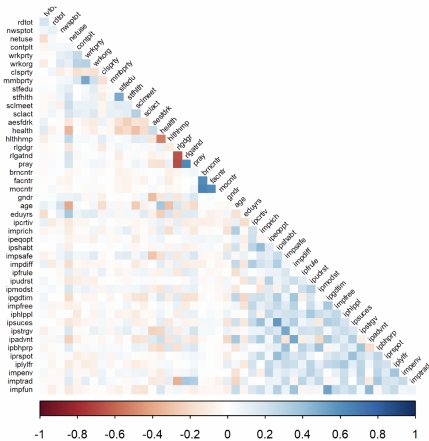
- The third round of the European Social Survey
- Randomly selected 1,000 samples (45.5% male and 54.5% female, Age: mean = 46.69, sd = 18.04)

Factor

- Center for Epidemiologic Studies - Depression Scale (CES-D, Radloff, 1977)
- e.g., "Felt depressed, how often past week"
- Eight items, 4-point Likert-type scale, treated as continuous following Van de Velde et al. (2009).

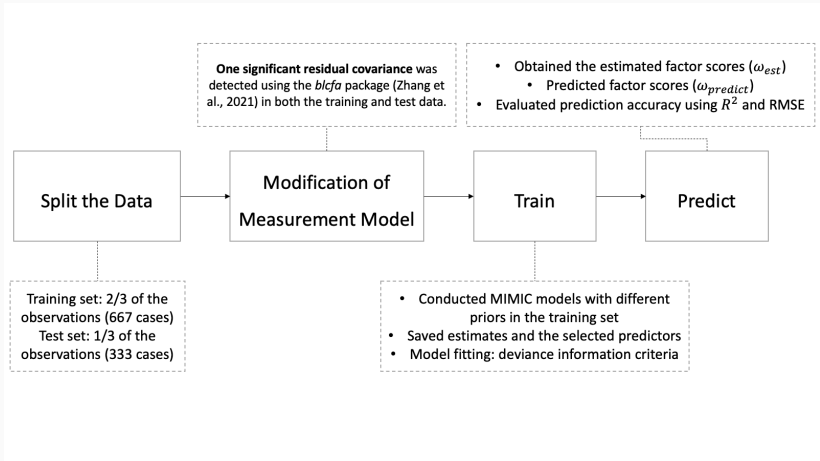
Multicollinearity

Forty-six covariates include demographic variables, health status, family status, and portrait values



Correlation Heatmap for Covariates

Hold-out Method



Results: Variable Selection and Parameter Estimation

covariate	definition	Non-informative	Ridge	Lasso	Alasso	Horseshoe	SSP
wrkprty	Worked in political party or action group last 12 months.	0.46 (0.083, 0.822) ¹	0.121 (-0.049, 0.286)	0.081 (-0.06, 0.284)	0.22 (-0.041, 0.56)	0.195 (-0.05, 0.572)	0.298 (-0.035, 0.592)
sclmeet	How often socially meet with friends, relatives or colleagues.	-0.076 (-0.14, -0.011)	-0.071 (-0.13, -0.01)	-0.057 (-0.115, 0.002)	-0.061 (-0.123, -0.004)	-0.045 (-0.108, 0.011)	-0.051 (-0.116, 0.007)
sclact	Take part in social activities compared to others of same age.	-0.133 (-0.195, -0.07)	-0.122 (-0.18, -0.064)	-0.118 (-0.178, -0.056)	-0.124 (-0.182, -0.064)	-0.127 (-0.189, -0.061)	-0.135 (-0.202, -0.066)
aesfdrk	Feeling of safety of walking alone in local area after dark.	0.07 (0.004, 0.135)	0.073 (0.012, 0.132)	0.071 (0.011, 0.133)	0.069 (0.006, 0.13)	0.074 (-0.001, 0.137)	0.067 (0.003, 0.14)
unhealth	Subjective general unhealthy.	0.22 (0.142, 0.298)	0.2 (0.132, 0.271)	0.21 (0.138, 0.283)	0.216 (0.141, 0.29)	0.235 (0.157, 0.314)	0.234 (0.158, 0.312)
hlthhmp	Hampered in daily activities by illness/disability/infirmity/mental problem.	-0.112 (-0.182, -0.043)	-0.108 (-0.171, -0.044)	-0.095 (-0.16, -0.027)	-0.099 (-0.167, -0.032)	-0.088 (-0.155, 0.001)	-0.091 (-0.167, -0.015)
gnr	gender	-0.139 (-0.268, -0.011)	-0.111 (-0.213, -0.004)	-0.086 (-0.193, 0.015)	-0.101 (-0.22, 0.011)	-0.077 (-0.21, 0.021)	-0.106 (-0.246, 0.02)
age	age	-0.099 (-0.177, -0.021)	-0.081 (-0.149, -0.011)	-0.062 (-0.127, 0.005)	-0.076 (-0.143, -0.003)	-0.061 (-0.136, 0.008)	-0.067 (-0.151, 0.001)
ipeqopt	Important that people are treated equally and have equal opportunities.	0.068 (0.006, 0.136)	0.064 (0.004, 0.123)	0.053 (-0.005, 0.11)	0.057 (-0.001, 0.117)	0.05 (-0.008, 0.111)	0.051 (-0.007, 0.115)
impsafe	Important to be humble and modest, not draw attention.	-0.08 (-0.151, -0.005)	-0.071 (-0.138, -0.004)	-0.057 (-0.122, 0.005)	-0.063 (-0.133, 0.002)	-0.048 (-0.119, 0.012)	-0.054 (-0.132, 0.01)
ipsuces	Important to be successful and that people recognize achievements.	0.114 (0.039, 0.193)	0.099 (0.029, 0.168)	0.08 (0.011, 0.149)	0.092 (0.023, 0.164)	0.076 (-0.004, 0.144)	0.084 (0.012, 0.168)
ipstrgv	Important that government is strong and ensures safety.	-0.103 (-0.176, -0.03)	-0.091 (-0.158, -0.026)	-0.077 (-0.146, -0.01)	-0.083 (-0.151, -0.014)	-0.067 (-0.14, 0.007)	-0.076 (-0.158, -0.004)
Number of Significant Covariates		12	11	6	8	2	6
DIC		14331.906	14287.801	14289.324	14331.017	14195.412	14235.981

Note. ¹ estimate (95% highest posterior density interval). Bold: significant estimates detected by the 95% HPD interval. DIC = deviance information criteria.

Results: Prediction

Correlation Among the Predicted Factor Scores and Estimated Values							R^2	RMSE
	Non-informative	Ridge	Lasso	Alasso	Horseshoe	SSP		
Non-informative	-						0.171	0.642
Ridge	0.977**						0.181	0.628
Lasso	0.921**	0.941**					0.191	0.614
Alasso	0.948**	0.969**	0.973**				0.181	0.621
Horseshoe	0.836**	0.853**	0.910**	0.904**			0.225	0.598
SSP	0.920**	0.940**	0.999**	0.974**	0.925**		0.195	0.616
Estimated Values	0.414**	0.426**	0.437**	0.425**	0.474**	0.441**	-	-

Note. ** $p < 0.01$.

Discussion

Penalty Priors vs Non-informative Priors

Variable selection

- ridge has advantages in handling high collinearity
- for low collinearity conditions, penalty priors except for ridge performed better than diffuse prior in small sample sizes

Parameter Estimation

- all regularization methods except for ridge outperformed the non-informative prior in maintaining low RMSEs in small sample size conditions.

Benefits of regularization in making predictions (empirical study) and achieving model convergence (ridge, lasso, and SSP).

Convergence

- Global shrinkage has advantages in model convergence.
- Alasso and horseshoe yielded low convergence rates with the co-existence of small sample size and high multicollinearity.

Variable Selection and Parameter Estimation

- Global shrinkage methods (ridge, lasso): variable selection
- Methods with local shrinkage parameter: parameter estimation
- SSP and Horseshoe had a similar performance in most conditions

Penalty priors compared to diffuse priors

- Robust results in small sample size conditions (simulation)
- High generalizability even with a relatively large sample size (empirical study)

Choice of different penalty priors

- For variable selection: global shrinkage (e.g., ridge in high collinearity conditions)
- For parameter estimation: penalty priors which include local shrinkage
- Model fit indexes (e.g., DIC)

- Other models and parameters (e.g., network analysis)
- Other penalty priors (e.g., Spike-and-Slab Lasso, SSP with continuous selection variables)
- Bayesian regularization vs frequentist regularization
- Model fitting assessment

Thanks for listening!