



INTRODUCTION

Large-scale social network services such as Facebook and Twitter have been dominating the online social media. In this work we consider a network vector autoregression model for the social network structure and study its statistical inference using two kinds of bootstrap methods. In the network vector autoregression model with large-size social media users, a continuous response of each user at a given time point is a linear combination of momentum effect, network effect, user-specific effect and independent noise. To analyze the model, a bootstrap version of the ordinary least square estimator is developed by means of the stationary bootstrap and the classical residual bootstrap. Its asymptotic properties and multivariate normal approximation are discussed on theory as well as visualized on a simulation study. For a finite sample validity, root mean square errors of the bootstrap estimates are computed and bootstrap confident intervals are constructed along with their empirical coverage probabilities and average lengths in a Monte-Carlo experiment.

MODEL and MAIN RESULT

Network Vector Auto-Regression Model

To describe relational data in large-scale social network of N users, an adjacency matrix $A \in \mathbb{R}^{N \times N}$, whose entries represent link between pairs of users and Y_{it} which is response of user i at time t (e.g., log-transformed total tweet length) are used.

$$A = (a_{ij}) = \begin{cases} 1, & \text{if user } i \text{ follows user } j \\ 0, & \text{otherwise.} \end{cases}$$

In Zhu et al. (2017), the network vector auto-regression model $\{Y_{it}\}$ is given by

$$Y_{it} = \beta_0 + Z_i^\top \gamma + \beta_1 n_i^{-1} \sum_{j=1}^N a_{ij} Y_{j(t-1)} + \beta_2 Y_{i(t-1)} + \varepsilon_{it}, 1 \leq i \leq N, 1 \leq t \leq T \quad (1)$$

where $Z_i = (Z_{i1}, \dots, Z_{iq})^\top \in \mathbb{R}^q$: q -dimensional user specific random vector,
 $n_i = \sum_{j \neq i} a_{ij} \in \mathbb{R}^1$: total number of users that user i follows, $\varepsilon_{it} \in \mathbb{R}^1$: error with mean zero variance σ^2
 $\gamma = (\gamma_1, \dots, \gamma_q)^\top \in \mathbb{R}^q$: user specific effect associated parameters, $\beta = (\beta_0, \beta_1, \beta_2)^\top \in \mathbb{R}^3$: parameters.

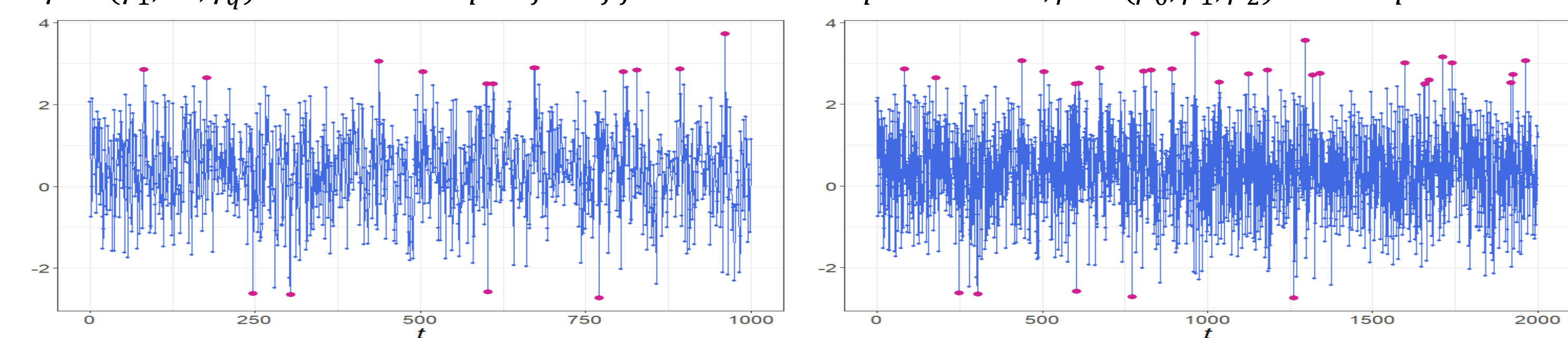


Fig. 1. Plots of Network Vector Auto-Regression model Y_{it}

Ordinary Least Square Estimator(OLSE)

We rewrite (1) as follows:

$$Y_{it} = \beta_0 + Z_i^\top \gamma + \beta_1 w_i^\top Y_{t-1} + \beta_2 Y_{i(t-1)} + \varepsilon_{it} = X_{it}^\top \theta + \varepsilon_{it} \quad (2)$$

where

$$w_i = \left(\frac{a_{i1}}{n_i}, \frac{a_{i2}}{n_i}, \dots, \frac{a_{iN}}{n_i} \right)^\top \in \mathbb{R}^N, W = (w_1, w_2, \dots, w_N), Y_{t-1} = (Y_{1(t-1)}, Y_{2(t-1)}, \dots, Y_{N(t-1)})^\top \in \mathbb{R}^N,$$

$$X_{it(t-1)} = (1, w_i^\top Y_{t-1}, Y_{i(t-1)}, Z_i^\top)^\top \in \mathbb{R}^{q+3}, \theta = (\beta_0, \beta_1, \beta_2, \gamma^\top)^\top.$$

We can finally rewrite (2) in matrix form with $i = 1, 2, \dots, N$ in rows as

$$Y_t = X_{t-1} \theta + \varepsilon_t, \text{ where } X_t = (X_{1t}, X_{2t}, \dots, X_{Nt})^\top \in \mathbb{R}^{N \times (q+3)}, \varepsilon_t = (\varepsilon_{1t}, \varepsilon_{2t}, \dots, \varepsilon_{Nt})^\top \in \mathbb{R}^N.$$

The OLSE is given by

$$\hat{\theta} = \left(\sum_{t=1}^T X_{t-1}^\top X_{t-1} \right)^{-1} \left(\sum_{t=1}^T X_{t-1}^\top Y_t \right) = \left(\sum_{t=1}^T \sum_{i=1}^N X_{i(t-1)} X_{i(t-1)}^\top \right)^{-1} \left(\sum_{t=1}^T \sum_{i=1}^N X_{i(t-1)} Y_{it} \right)$$

Monte-Carlo Simulation

To demonstrate the performance of the proposed model, Monte-Carlo simulations are used. The asymptotic properties and normal approximation of OLSE are verified through simulation. Specifically, we follow Nowicki and Snijders (2001), and randomly assign for each user label with equal probability.

User label K is simulated from discrete uniform $\{1, 2, \dots, K_{max}\}$

Adjacency matrix A is simulated as $a_{ij} = \begin{cases} 0.3 \times N^{-0.3}, & \text{if label of user } i \text{ is equal with label of user } j \\ 0.3 \times N^{-1}, & \text{otherwise.} \end{cases}$

Random error ε_{it} is simulated from $N(0, 1)$.

User specific vector $Z_i = (Z_{i1}, \dots, Z_{i3})^\top \in \mathbb{R}^3$ is simulated from $N_3(\mathbf{0}, \Sigma_z)$,
where $\Sigma_z = (\sigma_{j_1 j_2}), \sigma_{j_1 j_2} = 0.5^{|j_1 - j_2|}$.

Initial value Y_0 is simulated from $N_N(\mu, \Gamma(0))$,

where $\mu = (I - G)^{-1} \text{vec}\{\Gamma(0)\} = \sigma^2(I - G \otimes G)^{-1} \text{vec}(I)$, $G = \beta_1 W + \beta_2 I$.

Parameter θ is fixed to be $\theta = (\beta_0, \beta_1, \beta_2, \gamma^\top)^\top = (0, 0.1, -0.2, -0.5, 0.3, 0.8)$.

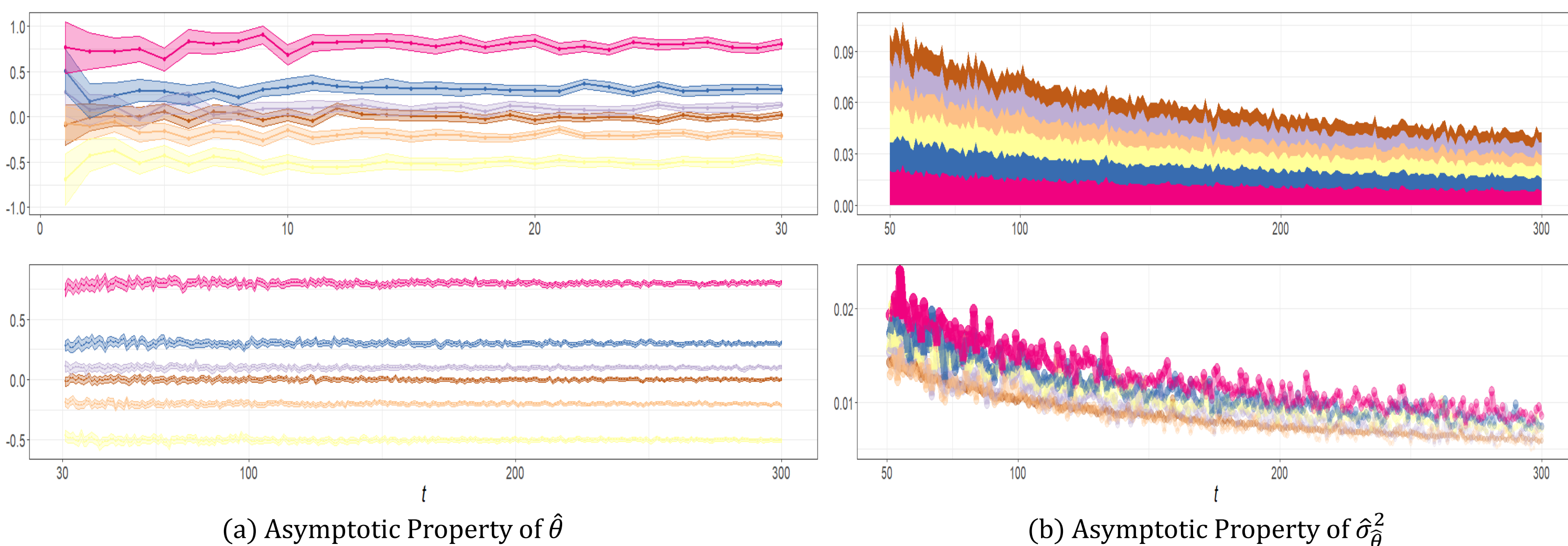


Fig. 2. Plots of Asymptotic Property

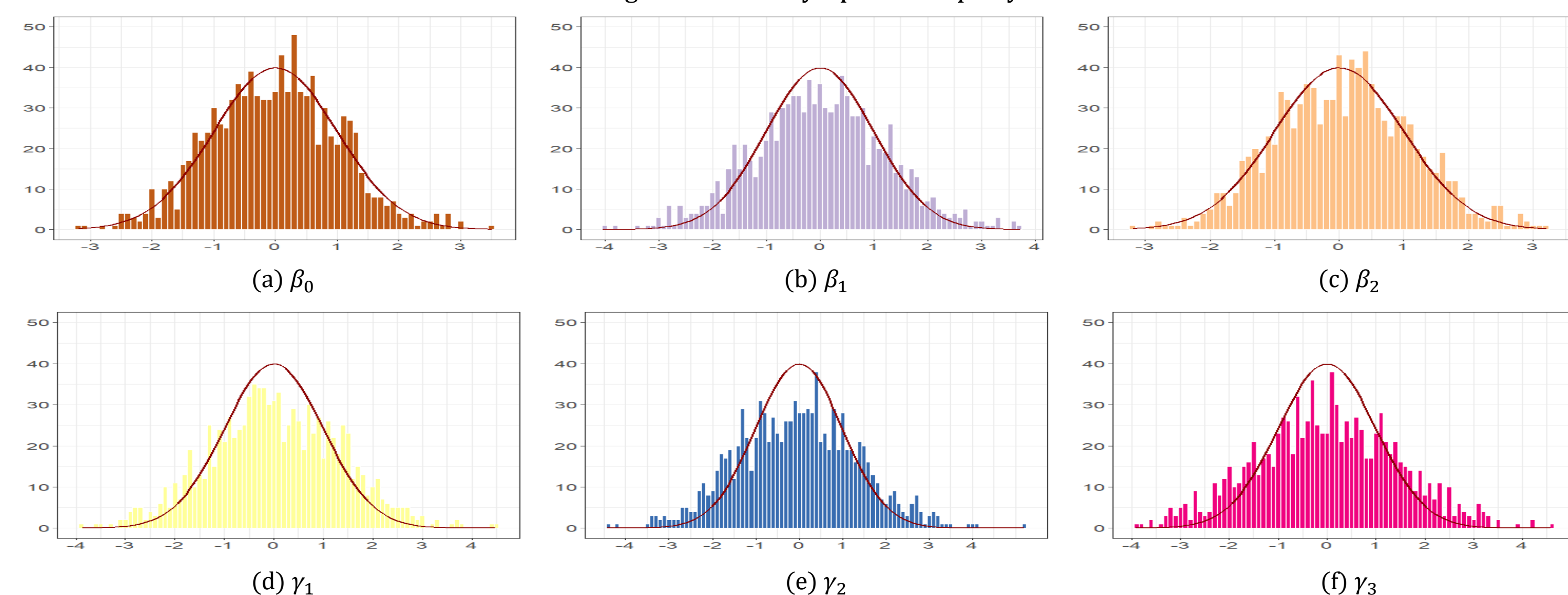


Fig. 3. Normal Approximation of $\sqrt{NT}(\hat{\theta} - \theta)$

Reference

- [1] Xuening Z, Rui P, Guodong L, Yuewen L and Hansheng W. (2017), Network Vector AutoRegression, The Annals of Statistics Volume 45, Number 3, 1096-1123.
- [2] Nowicki K. and Snijders T. A. B. (2001), Estimation and Prediction for Stochastic Block structures. Journal of the American Statistical Association. Volume 96, 1077-1087.

Bootstrap Estimator and Theorem

The OLSE residuals are computed as

$$\hat{\varepsilon}_{it} = Y_{it} - X_{it}^\top \hat{\theta}, \quad 1 \leq i \leq N, 1 \leq t \leq T.$$

Bootstrap version of residual array $\{\hat{\varepsilon}_{it}^*\}$ is generated through its empirical distribution function $\hat{F}_{\varepsilon}(\cdot)$:

$$\hat{F}_{\varepsilon}(x) = \frac{1}{NT} \sum_{t=1}^T \sum_{i=1}^N \mathbb{I}_{(-\infty, \hat{\varepsilon}_{it}]}(x).$$

Using the stationary bootstrap sample $\{Y_t^* = (Y_{1t}^*, \dots, Y_{Nt}^*)^\top : 0 \leq t \leq T-1\}$ obtained by

$$Y_{U_1}, Y_{U_1+1}, \dots, Y_{U_1+L_1-1}, \dots, Y_{U_k}, Y_{U_k+1}, \dots, Y_{U_k+L_k-1}, \dots$$

where $\{U_i : i = 1, 2, \dots\}$ are i.i.d. with discrete uniform $\{1, 2, \dots, T\}$ with probability $\frac{1}{T}$ and $\{L_i : i = 1, 2, \dots\}$

are i.i.d. with geometric $(p = 0.05 \times \frac{T-1}{100})$ the stationary bootstrap version of array X_t^* is defined as

$$X_t^* = (X_{1t}^*, \dots, X_{Nt}^*)^\top \text{ with } X_{it}^* = (1, w_i^\top Y_t^*, Y_{it}^*, Z_i^\top)^\top, 1 \leq i \leq N.$$

We define bootstrap estimator $\hat{\theta}^*$ of θ as

$$\hat{\theta}^* = \left(\sum_{t=1}^T X_{t-1}^{*\top} X_{t-1}^* \right)^{-1} \left(\sum_{t=1}^T X_{t-1}^{*\top} Y_t^* \right) = \left(\sum_{t=1}^T \sum_{i=1}^N X_{i(t-1)}^* X_{i(t-1)}^{*\top} \right)^{-1} \left(\sum_{t=1}^T \sum_{i=1}^N X_{i(t-1)}^* Y_{it}^* \right)$$

where $\tilde{Y}_{it}^* = (Y_{1t}^*, \dots, Y_{Nt}^*)^\top$, $\tilde{Y}_{it}^* = X_{t-1}^{*\top} \hat{\theta} + \hat{\varepsilon}_{it}^*$.

Theorem 1. Under the stationary condition $|\beta_1| + |\beta_2| < 1$, as $\min\{N, T\} \rightarrow \infty$,

$$\sqrt{NT}(\hat{\theta}^* - \theta) \xrightarrow{d} N(0, \sigma^2 \Sigma^{-1}), \quad \text{where } \Sigma = \begin{pmatrix} 1 & c_\beta & c_\beta & \mathbf{0}^\top \\ c_\beta & \Sigma_1 & \Sigma_2 & \kappa_8 \gamma^\top \Sigma_z \\ c_\beta & \Sigma_2 & \Sigma_3 & \kappa_3 \gamma^\top \Sigma_z \\ \mathbf{0} & \kappa_8 \Sigma_z \gamma & \kappa_3 \Sigma_z \gamma & \Sigma_z \end{pmatrix}$$

$c_\beta = \beta_0(1 - \beta_1 - \beta_2)^{-1}$, $\Sigma_1 = c_\beta^2 + \kappa_5 \gamma^\top \Sigma_z \gamma + \kappa_6$, $\Sigma_2 = c_\beta^2 + \kappa_7 \gamma^\top \Sigma_z \gamma + \kappa_2$, $\Sigma_3 = c_\beta^2 + \kappa_4 \gamma^\top \Sigma_z \gamma + \kappa_1$,
 $\kappa_1 = N^{-1} \text{tr}\{\Gamma(0)\}$, $\kappa_2 = N^{-1} \text{tr}\{W\Gamma(0)\}$, $\kappa_3 = N^{-1} \text{tr}\{(I - G)^{-1}\}$, $\kappa_4 = N^{-1} \text{tr}\{Q\}$, $\kappa_5 = N^{-1} \text{tr}\{WQW^\top\}$,
 $\kappa_6 = N^{-1} \text{tr}\{W\Gamma(0)W^\top\}$, $\kappa_7 = N^{-1} \text{tr}\{WQ\}$, $\kappa_8 = N^{-1} \text{tr}\{W(I - G)^{-1}\}$.

Theorem 2. As $\min\{N, T\} \rightarrow \infty$ and $Tp \rightarrow \infty$

$$\sup_x |P^*(\sqrt{NT}[\hat{\theta}^* - \hat{\theta}] \leq x) - P(\sqrt{NT}[\hat{\theta} - \theta] \leq x)| \xrightarrow{p} 0.$$

Bootstrap Confidence Interval

According to Theorem 2, we construct a confidence interval for each bootstrap estimator $\hat{\theta}^*$. Let $q_{0.025}^*$ and $q_{0.975}^*$ be the 0.025 and 0.975 quantiles of the bootstrap estimates. The 95% confidence interval based on the bootstrap is constructed by $0.95 = P(q_{0.025}^* - \hat{\theta} \leq \hat{\theta}^* - \hat{\theta} \leq q_{0.975}^* - \hat{\theta})$

$$\approx P(q_{0.025}^* - \hat{\theta} \leq \hat{\theta} - \theta \leq q_{0.975}^* - \hat{\theta}) = P(2\hat{\theta} - q_{0.975}^* \leq \theta \leq 2\hat{\theta} - q_{0.025}^*).$$

Thus, 95% bootstrap confidence interval = $[2\hat{\theta} - q_{0.975}^*, 2\hat{\theta} - q_{0.025}^*]$. Evaluation of bootstrap estimates via root mean square error, coverage probability and average length are presented.

K	T	$\beta_0 = 0$	$\beta_1 = 0.1$	$\beta_2 = -0.2$	$\gamma_1 = -0.5$	$\gamma_2 = 0.3$	$\gamma_3 = 0.8$	K	T	$\beta_0 = 0.1$	$\beta_1 = -0.3$	$\beta_2 = 0.2$	$\gamma_1 = 0.7$	$\gamma_2 = 0.1$	$\gamma_3 = -0.8$
5	50	(1.55)	(1.63)	(1.42)	(1.74)	(1.94)	(2.11)	5	50	(1.44)	(1.34)	(1.40)	(2.10)	(1.78)	(2.33)
		93.9%	95.9%	96.0%	94.8%	95.3%	95.5%			95.4%	95.4%	94.8%	95.6%	94.8%	96.8%
		[5.77]	[6.49]	[5.64]	[6.86]	[7.44]	[8.43]			[5.66]	[5.25]	[5.49]	[8.34]	[6.82]	[9.59]
	100	(1.00)	(1.23)	(1.01)	(1.29)	(1.41)	(1.48)		100	(1.13)	(1.11)	(1.02)	(1.62)	(1.32)	(1.93)
		95.8%	94.2%	94.8%	94.4%	95.0%	95.7%			94.0%	94.9%	94.4%	94.6%	95.8%	94.8%
		[4.08]	[4.75]	[3.99]	[4.85]	[5.62]	[5.97]			[4.15]	[4.31]	[3.88]	[6.16]	[5.42]	[7.35]
10	200	(0.72)	(0.85)	(0.74)	(0.88)	(0.93)	(1.21)	10	200	(0.72)	(0.68)	(0.69)	(1.10)	(0.96)	(1.25)
		95.5%	94.8%	95.5%	95.8%	96.1%	94.7%			96.3%	95.7%	95.7%	95.1%	94.4%	95.5%
		[2.92]	[3.40]	[2.83]	[3.58]	[3.74]	[4.68]			[2.85]	[2.62]	[2.73]	[4.29]	[3.72]	[4.80]
	50	(1.45)	(1.60)	(1.45)	(1.74)	(1.91)	(2.27)		50	(1.52)	(1.41)	(1.38)	(2.08)	(2.01)	(2.35)
		95.0%	95.5%	95.5%	94.6%	93.9%	94.1%			94.1%	94.4%	95.3%	94.1%	93.4%	95.4%
		[5.75]	[6.35]	[5.66]	[6.83]	[7.36]	[8.85]			[5.85]	[5.36]	[5.48]	[8.18]	[7.64]	[9.22]
20	100	(1.05)	(1.13)	(1.02)	(1.16)	(1.44)	(1.52)	20	100	(1.05)	(0.88)	(0.99)	(1.54)	(1.42)	(1.59)
		94.7%	96.3%	95.4%	96.6%	95.6%	94.9%			94.6%	94.2%	94.9%	95.2%	94.4%	94.2%
		[4.05]	[4.51]	[4.00]	[4.74]	[5.81]	[5.98]			[4.04]	[3.46]	[3.86]	[6.19]	[5.41]	[6.14]
	200	(0.74)	(0.93)	(0.73)	(0.91)	(0.96)	(1.01)		200	(0.75)	(0.71)	(0.69)	(0.99)	(0.95)	(1.14)
		95.5%	95.3%	95.0%	94.6%	93.8%	96.4%			94.9%	94.5%	95.1%	94.6%	94.9%	94.5%
		[2.92]	[3.64]	[2.82]	[3.43]	[3.60]	[4.07]			[2.95]	[2.78]	[2.74]	[3.93]	[3.76]	[4.47]
50	50	(1.44)	(1.81)	(1.42)	(1.60)	(1.82)	(2.08)	50	50	(1.59)	(1.58)	(1.38)	(2.05)	(1.80)	(2.37)
		95.4%	94.5%	95.7%	95.3%	96.4%	94.1%			93.7%	95.2%	95.8%	96.8%	94.3%	95.5%
		[5.78]	[6.78]	[5.67]	[6.56]	[7.38]	[7.89]			[5.92]	[6.36]	[5.53]	[8.38]	[6.94]	[9.19]
	100	(1.01)	(1.31)	(1.03)	(1.43)	(1.41)	(1.53)		100	(1.01)	(0.93)	(1.01)	(1.51)	(1.28)	(1.62)
		95.7%	94.7%	95.2%	94.3%	95.3%	94.1%			95.4%	94.2%	95.0%	94.0%	95.8%	94.4%
		[4.03]	[4.99]	[3.99]	[5.44]	[5.41]	[5.70]			[4.06]	[3.71]	[3.88]	[5.84]	[5.07]	[6.24]
200	200	(0.71)	(0.95)	(0.73)	(0.87)	(0.89)	(0.97)	200	200	(0.76)	(0.81)	(0.69)	(1.04)	(1.03)	(1.16)
		95.4%	95.3%	94.5%	95.7%	94.5%	94.6%			94.4%	94.4%	96.6%	95.1%	94.8%	96.2%
		[2.88]	[3.86]	[2.83]	[3.40]	[3.52]	[3.81]			[2.85]	[3.05]	[2.77]	[4.19]	[3.97]	[4.65]

Table 1. (RMSE $\times 10^2$) of Bootstrap Estimates, Coverage Probability(%) and [Average Length $\times 10^2$].
of Replications = 1000, # of Bootstrap Samples = 1000.

Asymptotic Normality of Bootstrap Estimates

According to Theorem 1, $\sqrt{NT}(\hat{\theta}^* - \theta)$ approximates to the multivariate normal distribution. To verify this, three pairs of bivariate normalized approximations of beta and gamma, respectively, via Monte-Carlo simulations are presented.

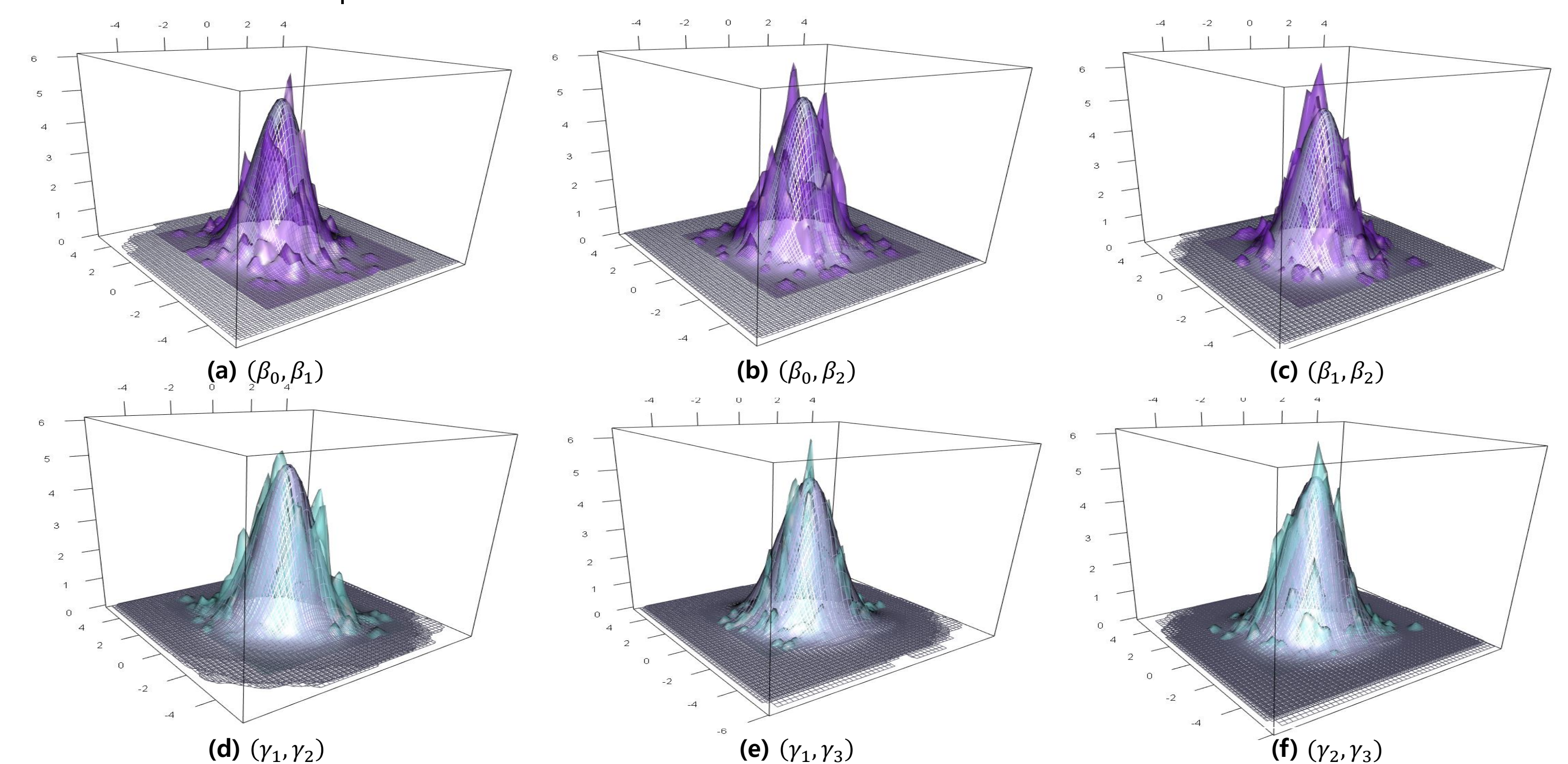


Fig. 4. Normal Approximation of $\sqrt{NT}(\hat{\theta}^* - \theta)$

저자 : 홍만호, 가천대학교 응용통계학과 학부생, ghdksgsh123@naver.com
황은주, 가천대학교 응용통계학과 조교수, ehwang@gachon.ac.kr
사사 : 본 연구는 가천대학교 캠퍼스 디자인 지원을 받아 수행되었음.